

**Titre:** New Faithfulness-Centric Interpretability Paradigms for Natural  
Title: Language Processing

**Auteur:** Andreas Madsen  
Author:

**Date:** 2024

**Type:** Mémoire ou thèse / Dissertation or Thesis

**Référence:** Madsen, A. (2024). New Faithfulness-Centric Interpretability Paradigms for Natural  
Citation: Language Processing [Thèse de doctorat, Polytechnique Montréal]. PolyPublie.  
<https://publications.polymtl.ca/59883/>

 **Document en libre accès dans PolyPublie**  
Open Access document in PolyPublie

**URL de PolyPublie:** <https://publications.polymtl.ca/59883/>  
PolyPublie URL:

**Directeurs de  
recherche:** Sarath Chandar Anbil Parthipan  
Advisors:

**Programme:** Génie informatique  
Program:

**POLYTECHNIQUE MONTRÉAL**

affiliée à l'Université de Montréal

**New Faithfulness-Centric Interpretability Paradigms  
for Natural Language Processing**

**ANDREAS MADSEN**

Département de génie informatique et génie logiciel

Thèse présentée en vue de l'obtention du diplôme de *Philosophiæ Doctor*  
Génie informatique

Novembre 2024

**POLYTECHNIQUE MONTRÉAL**

affiliée à l'Université de Montréal

Cette thèse intitulée :

**New Faithfulness-Centric Interpretability Paradigms  
for Natural Language Processing**

présentée par **Andreas MADSEN**

en vue de l'obtention du diplôme de *Philosophiæ Doctor*  
a été dûment acceptée par le jury d'examen constitué de :

**Amal ZOUAQ**, présidente

**Sarath Chandar ANBIL PARTHIPAN**, membre et directeur de recherche

**Siva REDDY**, membre et codirecteur de recherche

**Christopher J. PAL**, membre

**Marco Tulio RIBEIRO**, membre externe

## ACKNOWLEDGEMENTS

In 2024, the Machine Learning and Artificial Intelligence field is undoubtedly highly competitive. Realistically, hard work alone won't grant anyone a Ph.D. in this field. It takes a lot of support from institutions and individuals who believe in what one can achieve, not just based on what has been achieved. Therefore, I'm sincerely grateful to everyone who has believed in me and supported me on this journey.

As a matter of personal philosophy, I try my hardest not to make comparisons about others or myself. Therefore, this acknowledgment is written in approximately chronological order.

**Early education** Somewhat surprisingly, this Ph.D. ended up involving a lot of philosophy regarding the paradigm development of science and mathematics, and the early high-level ideas were motivated by philosophy and history lessons of past paradigm shifts. Those lessons were mostly forced upon me during Gymnasium and my Bachelor's, so thanks to The H.C. Ørsted Gymnasium and the Technical University of Denmark for providing a good education. Of course, this includes more technical topics, but I consider the aspect of philosophy to be particularly noteworthy.

I'm also grateful that Denmark funded this education and granted me enough money to pay for food and a dorm room. Although I worked part-time for some of this time, the support allowed me enough spare time to take on personal projects, which became highly instructive.

**Node.js Foundation** My most influential personal project was becoming a contributor to Node.js. During this time, I established connections with some of the world's best software developers, in particular Ryan Dahl and Isaac Schlueter. Their feedback was instrumental in developing my code-design taste, which has helped me in every endeavor since then.

**Personal motivations for interpretability** In the summer of 2015, I worked with Joel Raucq at Founders on developing a prototype search engine for recruiters. This time was instrumental in helping me understand the kind of work environment I prefer. This project also made me aware of how important interpretability is, as I realized that even with a perfect model, there will be imperfect people who use it. For example, there is nothing that prevents a recruiter from selecting candidates based on their bias. After finishing the work, I thought about this problem for a long time, and eventually realized that if we model the recruiter, we can also explain them and their bias, which can lead to real change. Even long after this job,

we still occasionally talk, and Joel’s support and encouragement have always been precious.

**BSc and MSc supervisor** My supervisor for my Bachelor’s and Master’s theses at the Technical University of Denmark was Ole Winther. Although his supervision was mostly hands-off, he allowed me to work on the projects that I found interesting, which were mostly NLP, and provided critical guidance, insights, and connections to enable me to succeed.

**NearForm** After I graduated, NearForm approached me to work for them. My friend Emil Melgaard had recommended me without me knowing, and I’m incredibly grateful for Emil thinking highly about me and for the opportunities that it led to. In particular, NearForm sponsored most of my work when writing my first publication for Distill and afterward provided me the flexibility to pursue my dreams of doing more research as an independent researcher. In particular, I’m grateful for James Snell, Matteo Collina, and Conor O’Neill, who have all been tremendously supportive and encouraging.

**Distill** My work at Distill [1] was my first publication. However, at the time, I didn’t know much about how to write a publication, and with no supervisors or research colleagues, I was going in mostly blind. I submitted it to Distill, and although I was rejected 3 times, every time, Chris Olah was tremendously kind to provide clear and constructive feedback, and the work was eventually accepted. I’m very grateful for this guidance, and undoubtedly, I would not have made it as an independent researcher without it.

**Alexander R. Johansen** For my next two publications I worked with Alexander R. Johansen who was a research assistant at the time, whom I met at a local conference. As the social individual he is, we talked about research, and he suggested we collaborate on something. This eventually led to a spotlight award at ICLR 2020, and I’m incredibly grateful for his support during this journey. Alexander himself fulfilled his dream of going to Stanford University to do a Ph.D., so it was quite productive for both of us. In connection to this, I also want to thank Andrew Trasks for being a good sport when we created a lot of drama around his research and to thank the Technical University of Denmark for lending us a lot of CPU computing resources that nobody was using anymore.

**Interpretability community** For many years, I had been thinking about the interpretability problem, but most professors discouraged me from working on it, saying it was not worth spending time on. So, when I saw Been Kim at NeurIPS 2019 doing a talk on interpretability, it was hugely influential to me. I’ve had a few chances to talk and write with her since then,

and she has always been encouraging and supportive. I later got inspired by Sara Hooker’s work on faithfulness, which has influenced much of the research in this thesis. We have talked a few times, and she also invited me to do the inaugural talk at Cohere for AI, which is a huge honor. Later, in an attempt to get an internship at Google, I connected with Jasmine Bastings, who moved heaven and earth to set up an internship for me. Things got complicated, and the internship never manifested itself, but I’m truly honored to have received such kind treatment. I sincerely want to thank all of you who have been so supportive on my journey.

**My Ph.D. supervisors** In January 2020, I was accepted into the Google AI Residency. I also had some backup plans, but after facing at least 50 Ph.D. program rejections in previous years, I didn’t apply for a Ph.D. this year. Then COVID-19 happened, the Google AI Residency was canceled, and all other backup plans were canceled. Then, by some incredible luck, Prof. Sarath Chandar saw my work and suggested that I come to do a Ph.D. with him at Mila. Regardless of the circumstances, I would always have considered such an opportunity very highly. However, it is safe to say that if I had not had this chance encounter and Sarath’s belief in me, my life would have looked very different. For that, I’m sincerely grateful.

Sarath then introduced me to Prof. Siva Reddy, and I feel incredibly fortunate to have worked with both of them. When we started, neither knew much about interpretability, but both have been incredibly supportive from the beginning to the end. I’m grateful for your belief in me.

**Colleagues** As I’m not the most social creature, I’m quite fortunate to be surrounded by those who are. My fellow students and post-docs in Sarath’s and Siva’s labs have always been kind and helpful.

**Examinators** A lot of work is done during a Ph.D., and examining and reviewing that is a major task. I have good memories from my Ph.D. proposal defense, where I received extremely useful suggestions and criticisms. Thanks for all your attention and support, both then and now.

**Funding** Finally, despite applying for many of the typical scholarships, I haven’t had much luck. I’m therefore grateful to the Danish foundations, Familien Hede Nielsens Fond, Viet-Jacobsen Fonden, Thomas B. Thriges Fond, Augustinus, and Jorcks Legater Og Stipendium, who were able to provide some support, without which I would certainly have had no chance to pay my medical bills and complete my Ph.D.

## RÉSUMÉ

L'apprentissage automatique, en particulier les réseaux de neurones, est aujourd'hui fréquemment utilisé dans de nombreuses applications, telles que l'évaluation des prêts, les résultats de recherche et les recommandations d'embauche. Ces systèmes fournissent souvent des décisions automatisées qui affectent la vie de la plupart des gens. Ces dernières années, les applications de traitement automatique du langage naturel (TALN) ont connu un essor considérable, car des modèles polyvalents et très efficaces en termes de données sont devenus disponibles, en particulier les systèmes de chat, qui sont désormais largement utilisés, même directement par le grand public.

Malheureusement, ces systèmes ne sont pas exempts de défauts. Bien que dans les systèmes d'embauche, il existe des cas documentés de discrimination fondée sur le sexe, comme le fait de favoriser "Chess club member" et de défavoriser "Women's Chess Club member" dans un curriculum vitae, ou un système qui pense qu'un diplôme en informatique est une qualification nécessaire pour être dactylographe dans un hôpital. Si le premier peut peut-être être atténué en analysant le modèle pour détecter les préjugés sexistes et les corriger, le second est si spécifique qu'il ne peut être révélé qu'en expliquant la prédiction du modèle.

L'interprétabilité est le domaine qui traite de l'explication des modèles et des ensembles de données aux humains en termes compréhensibles. L'objectif est généralement d'éviter un comportement indésirable, comme dans les exemples ci-dessus. Malheureusement, ce domaine est souvent confronté avec des défis dû à la production d'explications erronées, c'est-à-dire qui ne reflètent pas le modèle. Par exemple, une explication qui indique quels mots sont importants n'est pas forcément meilleure qu'une simple indication de mots au hasard. Lorsqu'une explication reflète le modèle, cela consiste en une explication fidèle.

Malheureusement, mesurer si une explication est fidèle (métrique de fidélité) est assez difficile car la vérité est inconnue. Une raison importante est que les modèles sont souvent trop complexes pour être compris par les humains, qui ne peuvent donc pas annoter si une explication est fidèle au modèle. De plus, des mesures de fidélité mal conçues peuvent donner une fausse confiance dans les explications, et de fausses explications peuvent donner une fausse confiance dans le comportement du modèle. C'est pourquoi la question centrale de cette thèse est la suivante : ***“Comment fournir et garantir des explications fidèles pour les modèles TALN neuronaux complexes à usage général?”***

En matière d'interprétabilité, il existe actuellement deux paradigmes sur la manière de répondre à cette question : *intrinsèque* et *post-hoc*. Intrinsèque dit que "seuls les modèles avec

un architecture conçue pour être expliquée peuvent l'être" et *post-hoc* dit que "les explications peuvent et doivent être produites après que le modèle soit entraîné, pour éviter tout impact sur les performances". Après avoir analysé la littérature existante et développé de nouvelles mesures de fidélité, cette thèse estime qu'aucun des deux paradigmes n'a été productif. La fidélité des explications *post-hoc* est souvent critiquée, et les modèles *intrinsèque* ne sont pas réellement *intrinsèque* ou sont trop contraints pour être des modèles à usage général très performants.

Par conséquent, cette thèse émet l'hypothèse que le domaine devrait développer de nouveaux paradigmes pour répondre à la question centrale de la recherche en combinant le meilleur des deux paradigmes existants. En d'autres termes, elle conçoit des modèles à expliquer sans utiliser de contraintes architecturales, de manière à ce que les modèles soient polyvalents et très performants. En particulier, cette thèse présente deux paradigmes potentiels de ce type : *Modèles mesurables de fidélité* (FMM) et *auto-explications*. Les FMM répondent directement à la question centrale de la recherche, tandis que les auto-explications ne répondent pas actuellement à la question, mais pourraient le faire à l'avenir.

Les modèles mesurables de fidélité sont un nouveau paradigme proposé dans cette thèse, qui reformule le souhait intrinsèque de "concevoir des modèles qui peuvent être expliqués" en "concevoir des modèles de telle sorte que la mesure de la fidélité soit peu coûteuse et fiable". Comme le montre cette thèse, il s'agit d'un objectif beaucoup plus facile à atteindre que ce que propose le paradigme intrinsèque, car il ne nécessite pas de contraintes architecturales. La démonstration spécifique de ce paradigme applique l'idée aux mesures d'importance, qui sont des explications qui indiquent l'importance de chaque caractéristique d'entrée pour faire une prédiction. Pour ce type d'explication, les mesures d'importance peuvent être obtenues en modifiant légèrement la procédure d'apprentissage, les jetons d'entrée aléatoires étant masqués pendant l'apprentissage.

Cette procédure d'apprentissage permet d'utiliser la métrique d'effacement de la fidélité : "Si une caractéristique est vraiment importante, la prédiction du modèle devrait changer de manière significative si cette caractéristique est supprimée". Cette métrique ne peut pas être appliquée à n'importe quel modèle, car la suppression de caractéristiques (par exemple, des mots ou des jetons) entraînera des problèmes de non-distribution. Toutefois, ce modèle de métrologie prend en charge ce type de suppression et permet donc l'application de la métrique. La fidélité étant désormais une mesure peu coûteuse et fiable, il est possible d'optimiser une explication en vue d'une fidélité maximale. Ainsi, les FMM deviennent indirectement interprétables de manière intrinsèque, mais sans utiliser de contraintes architecturales, et ils répondent également à la manière de mesurer la fidélité, répondant ainsi à la question centrale

de la recherche.

L’auto-explication est un autre paradigme émergent avec des explications produites directement par le modèle. Ces types d’explications sont devenus populaires en raison de l’essor des systèmes basés sur le chat, qui articulent souvent des explications pour leurs prédictions sous la forme d’un langage naturel. En raison de leur nature libre, il est notoirement difficile d’en évaluer la fidélité. En outre, comme ces modèles présentent également des problèmes d’hallucination, il y a de bonnes raisons d’être sceptique. Malgré cela, les explications sont extrêmement courantes et souvent prises pour argent comptant, y compris par les chercheurs dans ce domaine. Pour évaluer la faisabilité de ce nouveau paradigme, cette thèse propose et évalue également des mesures de fidélité pour les auto-explications. La conclusion est que, comme les explications post-hoc, elles dépendent du modèle et de la tâche.

Cette thèse étudie également la fidélité des explications post-hoc et intrinsèques et aboutit à la même conclusion, à savoir qu’elles dépendent du modèle et de la tâche. Cependant, ce n’est pas le cas lorsque l’on utilise des modèles mesurant la fidélité, même lorsque les mêmes méthodes d’explication post-hoc sont utilisées sur les mêmes ensembles de données et en utilisant le même modèle de base.

Cela permet de conclure que la fidélité dépend par défaut du modèle et de la tâche. Cependant, même de simples modifications du modèle, telles que le masquage aléatoire de l’ensemble de données d’apprentissage, comme cela a été fait dans les modèles mesurant la fidélité, peuvent changer radicalement la situation et donner lieu à des explications toujours fidèles. Nous proposons quelques suggestions sur la manière dont cela pourrait également être possible avec les auto-explications. De plus, avec les modèles à fidélité mesurable, cette thèse démontre qu’il est possible d’identifier des nouveaux paradigmes d’interprétabilité qui peuvent surmonter les limitations du passé et répondre à la question centrale de la recherche, à savoir comment fournir et garantir des explications fidèles pour des modèles TALN neuronaux complexes à usage général.

## ABSTRACT

Machine Learning, particularly Neural Networks, is nowadays frequently used in many applications, such as loan assessment, search results, and hiring recommendations. These systems often provide automated decisions which affect most people’s lives. In recent years, Natural Language Processing (NLP) applications have, in particular, seen a great increase as very data-efficient general-purpose models have become available, especially chat systems, which are now being widely used, even by the regular public directly.

Unfortunately, these systems are not without flaws. In hiring systems alone, there are documented cases of gender discrimination, such as favoring “Chess club member” and disfavoring “Women’s Chess Club member” in a resume, or a system that thinks that a computer science degree is a necessary qualification to be a typist at a hospital. While the former can perhaps be mitigated by analyzing the model for gender bias and correcting this, the latter is so specific that it can likely only be revealed by explaining the model’s prediction.

Interpretability is the field that deals with explaining models and datasets to humans in understandable terms. The goal is typically to prevent undesired behavior, as in the above examples. Unfortunately, the field is often challenged by providing false explanations, meaning the explanations do not reflect the model. For example, an explanation that indicates which input words are important might not be better than simply pointing at random words. When an explanation does reflect the model, it is termed a faithful explanation.

Unfortunately, measuring if an explanation is faithful (faithfulness metric) is quite challenging as the ground truth is unknown. One important reason is that the models are often too complex for humans to understand, thus humans cannot annotate if an explanation is true to the model. Even worse, poorly designed faithfulness metrics may provide false confidence in explanations, and false explanations may provide false confidence in the model’s behavior. As such, the central research question of this thesis is: ***“How to provide and ensure faithful explanations for complex general-purpose neural NLP models?”***

In interpretability, there are currently two paradigms on how to answer this question: *intrinsic* and *post-hoc*. The *intrinsic* paradigm says that “Only models architecturally designed to be explained can be explained” and the *post-hoc* paradigm says that “explanations can and should be produced after the model has been trained, to avoid any performance impact”. From analyzing the existing literature and developing new faithfulness metrics, this thesis takes the position that neither paradigm has been productive. The faithfulness of *Post-hoc* explanations is often criticized, and *intrinsic* models are either not actually *intrinsic* or are too

constrained to be high-performance general-purpose models.

Therefore, this thesis hypothesizes that the field should develop new paradigms to answer the central research question by combining the best of both existing paradigms. Namely, it designs models to be explained without employing architectural constraints, such that the models are general-purpose and high-performing. In particular, this thesis presents two such potential paradigms, *Faithfulness Measurable Models* (FMMs) and *self-explanations*. FMMs directly answer the central research question, while self-explanations do not currently answer the question but may do so in the future.

Faithfulness Measurable Models is a new paradigm that this thesis proposes, which reformulates the intrinsic desirable from “design models to be explained” to “designed models such that measuring faithfulness is cheap and reliable”. As is shown in this thesis, this is a much easier objective than what the intrinsic paradigm proposes, as it does not require architectural constraints. The specific demonstration of this paradigm applies the idea to importance measures, which are explanations that indicate the importance of each input feature in making a prediction. For this kind of explanation, FMMs can be achieved using only a small alteration to the training procedure, where random input tokens are masked during training.

This training procedure allows using the faithfulness erasure-metric, “If a feature is truly important, then the model’s prediction should change significantly if that feature is removed.”. This metric cannot be applied to just any model because removing features (e.g. words or tokens) will cause out-of-distribution issues. However, this FMM supports such removal and thus enables the metric. Because faithfulness is now cheap and reliable to measure, optimizing an explanation towards maximal faithfulness is possible. As such, FMMs become indirectly intrinsically interpretable, but without employing architectural constraints, and they also answer how to measure faithfulness, thereby answering the central research question.

Self-explanation is another emerging paradigm, these explanations are produced directly as the model’s output. These types of explanations have become popular due to the rise of chat-based systems, which will often articulate explanations for their utterances in the form of natural language. However, due to the free-form nature of self-explanations, evaluating their faithfulness is notoriously difficult. Additionally, because these models also have hallucination issues, there are good reasons to be skeptical. Despite this, the explanations are extremely common and often taken at face value, including by researchers in the field. To evaluate the feasibility of this new paradigm, this thesis also proposes and evaluates faithfulness metrics for self-explanations. The finding is that similar to post-hoc explanations, they are model and task-dependent.

This thesis also investigates the faithfulness of post-hoc and intrinsic explanations and finds

the same model and task-dependent conclusion. However, this was not the case when using faithfulness measurable models, even when the same post-hoc explanation methods were used on the same datasets and using the same base model.

This leads to the overall conclusion that faithfulness is by default model and task-dependent. However, even simple modifications to the model, such as randomly masking the training dataset, as was done in faithfulness measurable models, can drastically change the situation and result in consistently faithful explanations. We provide some suggestions on how this might also be possible with self-explanations. Additionally, with faithfulness measurable models, this thesis demonstrates that it is possible to identify new interpretability paradigms that can overcome past limitations and answer the central research question of how to provide and ensure faithful explanations for complex general-purpose neural NLP models.

## TABLE OF CONTENTS

ACKNOWLEDGEMENTS . . . . .	iii
RÉSUMÉ . . . . .	vi
ABSTRACT . . . . .	ix
TABLE OF CONTENTS . . . . .	xii
LIST OF TABLES . . . . .	xvi
LIST OF FIGURES . . . . .	xix
LIST OF SYMBOLS AND ACRONYMS . . . . .	xxxvi
LIST OF APPENDICES . . . . .	xl
CHAPTER 1 INTRODUCTION . . . . .	1
1.1 Why interpretability is needed . . . . .	2
1.1.1 Limitations of bias and fairness metrics . . . . .	3
1.1.2 Interpretability for scientific discovery and understanding . . . . .	4
1.2 The current paradigms of interpretability . . . . .	4
1.2.1 Definitions . . . . .	5
1.2.2 Beliefs . . . . .	6
1.3 Why interpretability needs a new paradigm . . . . .	8
1.3.1 The case against the intrinsic paradigm . . . . .	8
1.3.2 The case against the post-hoc paradigm . . . . .	9
1.3.3 Overall trend . . . . .	10
1.4 Thesis overview . . . . .	10
1.4.1 Content overview . . . . .	11
1.4.2 Scientific contributions . . . . .	13
1.4.3 Statement of contributions in papers . . . . .	13
CHAPTER 2 BACKGROUND . . . . .	16
2.1 Organizing by method of communication . . . . .	16
2.2 Motivating Example . . . . .	18
2.3 Measures of Interpretability . . . . .	20

2.4	Methods of Interpretability . . . . .	22
2.5	Input Features . . . . .	22
2.5.1	Gradient-based . . . . .	23
2.5.2	Occlusion-based . . . . .	25
2.5.3	Attention-based . . . . .	27
2.5.4	Discussion . . . . .	29
2.6	Counterfactuals . . . . .	30
2.6.1	Polyjuice . . . . .	30
2.6.2	MiCE . . . . .	31
2.6.3	Discussion . . . . .	33
2.7	Natural Language . . . . .	33
2.7.1	Explain-then-predict . . . . .	34
2.7.2	Predict-then-explain . . . . .	35
2.7.3	Discussion . . . . .	36
CHAPTER 3 GENERAL-PURPOSE FAITHFULNESS METRIC FOR IMPOR-		
TANCE MEASURES . . . . .		37
3.1	Existing faithfulness metrics . . . . .	39
3.1.1	Correlating importance measures . . . . .	39
3.1.2	Mutate model to deceive . . . . .	39
3.1.3	Optimize model to deceive . . . . .	40
3.1.4	Known explanations in synthetic tasks . . . . .	40
3.1.5	Similar inputs, similar explanation . . . . .	40
3.2	ROAR: RemOve And Retrain . . . . .	41
3.2.1	Adaptation to NLP . . . . .	41
3.2.2	Recursive ROAR . . . . .	42
3.2.3	Validation on a synthetic problem . . . . .	43
3.3	Models . . . . .	44
3.3.1	BiLSTM-Attention . . . . .	44
3.3.2	RoBERTa . . . . .	45
3.4	Experiments . . . . .	46
3.4.1	Recursive ROAR . . . . .	48
3.4.2	Summarizing faithfulness metric . . . . .	50
3.4.3	Supporting experiments . . . . .	52
3.5	Important Findings . . . . .	54
3.6	Limitations . . . . .	55

3.7	Conclusion . . . . .	57
CHAPTER 4 FAITHFULNESS MEASURABLE MODELS . . . . .		58
4.1	Inherently faithfulness measurable models (FMMs) . . . . .	60
4.1.1	Faithfulness of importance measures . . . . .	61
4.1.2	Masked fine-tuning . . . . .	61
4.1.3	In-distribution validation . . . . .	63
4.1.4	Faithfulness metric . . . . .	65
4.1.5	Optimizing for faithfulness (Beam) . . . . .	65
4.2	Experiments . . . . .	67
4.2.1	Masked fine-tuning . . . . .	67
4.2.2	In-distribution validation . . . . .	70
4.2.3	Faithfulness metric . . . . .	72
4.3	Important Findings . . . . .	74
4.4	Limitations . . . . .	75
4.5	Conclusion . . . . .	77
CHAPTER 5 FAITHFULNESS OF SELF-EXPLANATIONS . . . . .		79
5.1	Self-explanations . . . . .	81
5.2	Faithfulness of self-explanations . . . . .	81
5.2.1	Self-consistency check . . . . .	81
5.2.2	Negative results and robustness . . . . .	82
5.3	Prompt-patterns for self-explanation and self-consistency checks . . . . .	82
5.3.1	Counterfactual explanation . . . . .	83
5.3.2	Feature attribution explanation . . . . .	83
5.3.3	Redaction explanation . . . . .	84
5.3.4	Classification . . . . .	85
5.3.5	Persona robustness . . . . .	85
5.4	Experiments . . . . .	86
5.4.1	Classification . . . . .	88
5.4.2	Faithfulness . . . . .	89
5.5	Important findings . . . . .	89
5.6	Related work . . . . .	91
5.6.1	Non-faithfulness works . . . . .	92
5.7	Limitations . . . . .	93
5.8	Conclusion . . . . .	95

CHAPTER 6	CONCLUSION . . . . .	96
6.1	Limitations . . . . .	98
6.2	Future Research . . . . .	98
6.2.1	Masked Causal Language Models . . . . .	98
6.2.2	Faithfulness Measure Models for other kinds of communication . . . . .	99
6.2.3	Preventing double-edged alignment . . . . .	100
REFERENCES	. . . . .	101
APPENDICES	. . . . .	130

# LIST OF TABLES

Table 1.1	Comparison of the definitions and underlying beliefs of the intrinsic and post-hoc paradigms. The beliefs relate to a) requirements for a faithful explanation and b) model capabilities. It should be apparent that these two views are seemingly incompatible. . . . .	5
Table 1.2	Comparison of the definitions and underlying beliefs of the new paradigms. The beliefs relate to a) explanation requirements and b) model capabilities. These new paradigms can be compared with the old paradigms in Table 1.1. . . . .	11
Table 2.1	Overview of <i>post-hoc</i> interpretability methods . . . . .	17
Table 3.1	Details on the BiLSTM-attention models' hyperparameters. Everything is exactly as done by Jain and Wallace [2]. For all datasets, ASMGrad Adam [3] is used with default hyperparameters ( $\lambda = 0.001$ , $\beta_1 = 0.9$ , $\beta_2 = 0.999$ , $\epsilon = 10^{-8}$ ) and a weight decay of $10^{-5}$ . . . . .	45
Table 3.2	Details on the RoBERTa models' hyperparameters. RoBERTa [4] is fine-tuned using the RoBERTa-base pre-trained model from HuggingFace [5] (125M parameters). The hyperparameters are those used by Liu et al. [4] on GLUE tasks [4, Appendix C]. The optimizer is AdamW [6], the learning rate has linear decay with a warmup ratio of 0.06, and there is a weight decay of 0.01. Additionally, we use a batch size of 16 and a learning rate of $2 \cdot 10^{-5}$ . . . . .	46
Table 3.3	Datasets statistics for single-sequence and paired-sequence tasks. Following Jain and Wallace [2], we use the same BiLSTM-attention model and report performance as macro-F1 for SST, IMDB, Anemia, and Diabetes; micro-F1 for SNLI; and accuracy for bAbI. The 95% confidence interval is reported in the subscript and subscript. . . . .	47
Table 3.4	Faithfulness metric defined as a <b>relative area-between-curves</b> (RACU) using Recursive ROAR, see (3.4). Higher values mean more faithful, and zero or negative values mean distinctly not faithful. . . . .	51
Table 3.5	Performance of the class-majority classifier and the BiLSTM-Attention and RoBERTa classifier on the 100% masked dataset. Performance is the standard metric for the dataset, meaning macro-F1 for SST, IMDB, Anemia, and Diabetes; micro-F1 for SNLI; and accuracy for bAbI. . . .	52

Table 4.1	Max-epoch parameters and performance statistics for each dataset. Performance metrics are the mean with a 95% confidence interval. . . .	68
Table 4.2	Datasets used, all datasets are either single-sequence or sequence-pair datasets. All datasets are sourced from GLUE [7], SuperGLUE [8], MIMIC-III [9], or bAbI [10]. The decisions regarding which metrics are used are also from these sources. The class-majority baseline is when the most frequent class is always selected. . . . .	69
Table 4.3	Faithfulness scores using Relative Area Between Curves (RACU) and the non-relative variant (ACU). The less relevant score is grayed out. Higher is better. Negative values indicate not-faithful. The comparison with Recursive-ROAR from Chapter 3 is imperfect because Recursive-ROAR has limitations. See Table C.7 for all datasets and Table C.8 for RoBERTa-large. . . . .	73
Table 5.1	List of datasets used in Chapter 5. All datasets are publicly available. .	86
Table 5.2	List of models used in Chapter 5. All models are publicly available. . .	87
Table 5.3	Generation parameters used for TGI. These are the same parameters used in HuggingFace’s online chat service ( <a href="https://huggingface.co/chat/">https://huggingface.co/chat/</a> ). . . . .	88
Table A.1	Fictive example of the top-3 words for each basis-dimension in the rotated word embeddings. . . . .	142
Table A.2	Performance on the HANS dataset provided by McCoy et al. [11]. Unfortunately, McCoy et al. [11] do not provide enough information to make a direct comparison possible. For comparison, BERT has 83% accuracy on MNLI [12], which was used for training. . . . .	147
Table B.1	Compute hardware used for each model. Note, the models were computed on a shared user system. Hence, we only report the resources allocated for our jobs. . . . .	153
Table B.2	Compute times for each model and importance measure combination. Note, there is no need to compute models for each importance measure at 0% and 100% masking. Hence, we report for 9 iterations. . . . .	154
Table C.1	The computing hardware used. Note, that a shared user system were used, only the allocated resources are reported. . . . .	163
Table C.2	Walltime for fine-tuning. Masked fine-tuning does not affect the total wall time in our setup. . . . .	164

Table C.3	Walltime for in-distribution validation. This does not include importance measure calculations. See Table C.5. . . . .	165
Table C.4	Walltime for faithfulness evaluation. This does not include importance measure calculations. See Table C.5. . . . .	166
Table C.5	Walltime for importance measures. Note that because the beam-search method (Beam) scales quadratic with the sequence-length, it is not feasible to compute for all datasets. . . . .	167
Table C.6	This table relates terminologies between the fine-tuning strategies mentioned in Section 4.2.1 and the training strategy and validation strategy terms. . . . .	168
Table C.7	Faithfulness scores for <b>RoBERTa-base</b> . Shows Relative Area Between Curves (RACU) and the non-relative variant (ACU), defined in Chapter 3. Also compares with Recursive-ROAR from Chapter 3. . . .	188
Table C.8	Faithfulness scores for <b>RoBERTa-large</b> . Shows Relative Area Between Curves (RACU) and the non-relative variant (ACU), defined in Chapter 3, note that Chapter 3 does not report results for Recursive-ROAR with RoBERTa-large. . . . .	189
Table D.1	Inference time as reported by TGI. Note that this does not correspond to wall-time. In particular, because 50 prompts are computed in parallel. Dividing the inference time by 50 is a decent approximation for wall-time.	190
Table D.2	The computing hardware used. Only the allocated resources are reported.	191

## LIST OF FIGURES

Figure 2.1	Fictive visualization of an <i>input features</i> explanation which highlights tokens and a <i>natural language</i> explanation, applied on a sentiment classification task [7]. $y = \text{pos}$ means the gold label is <i>positive</i> sentiment. . . . .	18
Figure 2.2	Three examples from the SST dataset [13]. $\mathbf{x}$ is the input, with each token denoted by an <u>underline</u> . $y$ is the gold target label, where <b>pos</b> is <i>positive</i> and <b>neg</b> is <i>negative</i> sentiment. Finally, $p(y \mathbf{x})$ is the model's estimate of $\mathbf{x}$ belonging to category $y$ . Note that the model predicts the 3rd (last) wrong, indicated with <b>red</b> . . . . .	19
Figure 2.3	Hypothetical visualization of applying $\mathbf{E}_{\text{gradient}}(\mathbf{x})$ , where $c$ is the explained class. Note that because the vocabulary dimension is reduced away, typically using the $L^2$ -norm, it is not possible to separate positive influence from negative influence. . . . .	23
Figure 2.4	A fictive visualization of LIME, where the weights of the logistic regression determine the <i>importance measure</i> . Note that for LIME, it is possible to have negative importance (indicated by blue). Furthermore, some tokens have no importance score due to the $L^1$ -regularizer. . . . .	26
Figure 2.5	Fictive visualization of <i>Kernel SHAP</i> . Note how input tokens are combined to a single feature to make <i>SHAP</i> more tractable to compute, this is the role of $h_{\mathbf{x}}(z)$ in (2.6). . . . .	27
Figure 2.6	Hypothetical results of <i>Polyjuice</i> , showing how some words were either replaced or removed to produce <i>counterfactual examples</i> . . . . .	31
Figure 2.7	Hypothetical visualization of how <i>MiCE</i> progressively creates a counterfactual $\tilde{\mathbf{x}}$ from an original sentence $\mathbf{x}$ . The highlight shows the <i>gradient</i> $\nabla_{\mathbf{x}}f(\mathbf{x}; \theta)_y$ , which <i>MiCE</i> uses to know what tokens to replace. . . . .	32
Figure 2.8	Hypothetical explanations from using <i>CAGE</i> to produce rationalizations for the prediction. . . . .	36
Figure 3.1	Example of <b>ROAR</b> . The first sentence shows the importance of various tokens. The next two sentences demonstrate the proportion of important tokens replaced by <b>[MASK]</b> . Note, the second sentence is enough to infer the sentiment. . . . .	42

Figure 3.2	Example of how a redundancy can be removed in <b>Recursive ROAR</b> by reevaluating the <i>importance measure</i> . Compare this to Figure 3.1, where redundancies are not removed and the performance can remain the same, even when the <i>importance measure</i> is faithful. . . . .	42
Figure 3.3	Using the weights of a linear model as the explanation, ROAR and Recursive ROAR are applied to the problem described in (3.1). In addition, the ground truth and worst case are shown. Recursive ROAR and the ground truth are identical. Note that the worst case does not lose performance at “13 features removed” since there are redundancies. . . . .	43
Figure 3.4	Recursive ROAR results, showing model performance at x% of tokens masked. A model performance below <i>random</i> indicates faithfulness, while above or similar to <i>random</i> indicates a non-faithful importance measure. Performance is averaged over 5 seeds with a 95% confidence interval. . . . .	49
Figure 3.5	Visualization of the faithfulness calculation done in (3.4). The <i>faithfulness</i> area is the numerator in (3.4), while the <i>normalizer</i> area is the denominator. Essentially (3.4) computes the <b>relative area-between-curves</b> (RACU) between an <i>explanation</i> curve and the <i>random</i> baseline curve. . . . .	50
Figure 3.6	ROAR and Recursive ROAR results for BiLSTM-Attention, showing model performance at x% of tokens masked. A model performance below <i>random</i> indicates faithfulness. For Recursive ROAR a curve above or similar to <i>random</i> indicates a non-faithful importance measure, while for ROAR, this case is inconclusive [14]. . . . .	53
Figure 3.7	Recursive ROAR results, showing model performance at up to 10 tokens masked. Note that because the datasets have more than 10 tokens, the conclusion one can draw from this plot may change if more tokens are considered. However, a model performance below <i>random</i> indicates faithfulness, while above or similar to <i>random</i> indicates a non-faithful importance measure. Performance is averaged over 5 seeds with a 95% confidence interval. . . . .	53
Figure 3.8	The accumulative importance score relative to the total importance score for the top-x% number of tokens. The metric is averaged over 5 seeds with a 95% confidence interval. . . . .	54

Figure 4.1	To measure faithfulness, a <i>faithfulness measurable masked language model</i> is created (a), then the model is checked for out-of-distribution issues given an explanation (b), and finally, the faithfulness is measured by masking allegedly important tokens (c). – $[\mathcal{M}]$ is the masking token.	59
Figure 4.2	Visualization of the faithfulness calculation. AUC is the <i>faithfulness</i> area, and RACU is the AUC normalized by the theoretical best explanation. See the definition for AUC and RACU in (3.4).	66
Figure 4.3	The unmasked performance for each fine-tuning strategy. <i>Plain fine-tuning</i> is the baseline (dashed line). We find that our <i>Masked fine-tuning</i> does not decrease performance. <i>All</i> is computed by taking the average of all datasets. More datasets and a more detailed ablation study can be found in Appendix C.	70
Figure 4.4	The 100% masked performance for each fine-tuning strategy. The dashed line represents the class-majority baseline. Results show that masking during training (either our <i>masked fine-tuning</i> or <i>only masking</i> ) is necessary. More datasets and a more detailed ablation study can be found in Appendix C.	70
Figure 4.5	In-distribution p-values using MaSF, for RoBERTa-base with and without masked fine-tuning. The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using <i>masked fine-tuning</i> is masking consistently not OOD. Because the results are highly consistent, the overlapping lines do not hide any important details. More datasets and models in Appendix C. Grad is “Gradient”, $x \odot grad$ is “Input times gradient”, IG is “Integrated Gradient”, LOO is “Leave-one-out”. These methods are described in Section 2.5.	71
Figure 4.6	The performance given the masked datasets, where masking is done for the x% allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “ <i>Random</i> ” baseline, this shows faithfulness. Although faithfulness is not an absolute concept, so more is better. This plot is for RoBERTa-base and separates importance measures based on their signed and absolute variants. More datasets and models in Appendix C. Legend has same meaning as in Figure 4.5 but does not separate between signed (sign) and absolute (abs) importance measures.	71

Figure 5.1	Example of an LLM providing a counterfactual self-explanation and using a self-consistency check to evaluate if it is faithful. – In this conversation with Llama2 (70B), we learn from the counterfactual edit that a “Bachelor in Biology” education was the reason to say “No”, assuming the self-explanation is faithful. Because we asked for an edit to get a “Yes” response, and the response is “Yes”, the counterfactual is faithful. Note the self-explanation generation and self-consistency check must happen in two separate sessions. . . . .	80
Figure 5.2	The explicit input-template prompt used for generating the counterfactual explanation. {opposite sentiment} is replaced with either “positive” or “negative”. {paragraph} is replaced with the content. We also consider an implicit version where “is {opposite sentiment}” is replaced with “becomes the opposite of what it currently is”. The partial output example is entirely generated by the model. . . . .	84
Figure 5.3	The input-template prompt used for generating the feature attribution explanations. The model will often generate either a bullet-point list or a comma-separated list. . . . .	84
Figure 5.4	The input-template prompt used for generating redaction explanations. We also consider a prompt where “[REMOVED]” is used instead of “[REDACTED]”. . . . .	85
Figure 5.5	Prompt-template for classification. The prompt needs to support redaction and an “unknown” class for when the classification can not be performed due to missing information. . . . .	85
Figure 5.6	The classification accuracy of each task using Llama2-70B, with different prompt variations. Performance is not affected much by the persona or redaction-instruction. . . . .	89
Figure 5.7	The classification accuracy of each task with the default prompt settings. Performance is highly dependent on model type and size. . . . .	89
Figure 5.8	Faithfulness evaluation using self-consistency checks, evaluated using Llama2-70B. Results show that Llama2-70B is not affected by prompt variations, but the faithfulness for each explanation type is task-dependent.	90

Figure 5.9	Faithfulness evaluation using self-consistency checks, evaluated with default prompt settings. Results show that faithfulness generally increases with size, the exception being Falcon, particularly in the counterfactual case. In general, the results are heavily dependent on the model type, tasks, and explanation. . . . .	90
Figure A.1	Hypothetical visualization of <i>HotFlip</i> . The highlight indicates the gradient w.r.t. the input, which HotFlip uses to select which token to change. $\mathbf{x}$ indicates the original sentence, and $\tilde{\mathbf{x}}$ indicates the adversarial sentence. . . . .	132
Figure A.2	Hypothetical results of using <i>SEA</i> [15]. Note that unlike <i>HotFlip</i> , <i>SEA</i> can change and delete multiple tokens simultaneously as it samples from a paraphrasing model. Again, $\mathbf{x}$ indicates the original sentence, $\tilde{\mathbf{x}}$ indicates the adversarial sentence, and $S(\mathbf{x}, \tilde{\mathbf{x}})$ is the semantical-equivalency-score which must be at least 0.8. . . . .	133
Figure A.3	Fictive result showing the <i>influential examples</i> $\tilde{\mathbf{x}}$ , in relation to the <i>input example</i> $\mathbf{x}$ , showing both examples with positive and negative influence. $\Delta$ is the similarity score; the scale and range may depend on the specific method. Note it is possible to measure the influence of an example on itself. This can be useful to identify mislabeled observations, as such observations will be important for their own prediction. . . . .	135
Figure A.4	Visualization of hypothetical <i>Natural Indirect Effect</i> (NIE) results, similar to Vig et al. [16]. Such visualization can reveal which attention-heads are responsible for gender bias in a small GPT-2 model. A stronger color indicates a higher NIE, meaning more responsible for the bias. . . . .	140
Figure A.5	PCA [17] and t-SNE [18] projection of GloVe [19] embeddings for the words in the semantic classification examples, as shown in Section 2.2 and elsewhere in the background chapter and appendix. . . . .	142
Figure A.6	Visualization of <i>SP-LIME</i> in a hypothetical setting. The matrix shows how each selected observation represents the different modes of the model. The left side shows two out of the four selected examples and their <i>LIME</i> explanation. . . . .	144
Figure A.7	Results by Tenney et al. [20] which shows how much each BERT [21] layer is used for each linguistic task. The $F_1$ score for each task is also presented. . . . .	148

Figure A.8	Hypothetical example showing rules which commonly break the model. The flip-rate describes how often these rules break the model. $\mathbf{x}$ represents the original sentence and $\tilde{\mathbf{x}}$ represents an adversarial example. . . . .	150
Figure A.9	Hypothetical example showing rules which activate a selected neuron. IoU is how often the rule activated the neuron, compared to cases where either the rule is true or the neuron activated (higher is better). . . . .	151
Figure B.1	Shows the accumulative importance score relative to the total importance score, for the top-k number of tokens. The metric is averaged over 5 seeds with a 95% confidence interval. Note that datasets are not equal in sequence-length, the scores are therefore hard to compare across datasets. . . . .	156
Figure B.2	The accumulative importance score relative to the total importance score for the top-x% number of tokens. The metric is averaged over 5 seeds with a 95% confidence interval. . . . .	157
Figure B.3	Recursive ROAR results, showing model performance at up to 10 tokens masked. Note that because the datasets have more than 10-tokens, the conclusion one can draw from this plot may change if more tokens were considered. However, in general, a model performance below <i>random</i> indicates faithfulness, while above or similar to <i>random</i> indicates a non-faithful importance measure. Performance is averaged over 5 seeds with a 95% confidence interval. . . . .	159
Figure B.4	ROAR and Recursive ROAR results for <b>BiLSTM-Attention</b> , showing model performance at x% of tokens masked. A model performance below <i>random</i> indicates faithfulness. For Recursive ROAR a curve above or similar to <i>random</i> indicates a non-faithful importance measure, while for ROAR by Hooker et al. [14] this case is inconclusive. Performance is averaged over 5 seeds with a 95% confidence interval. . . . .	161
Figure B.5	ROAR and Recursive ROAR results for <b>RoBERTa</b> , showing model performance at x% of tokens masked. A model performance below <i>random</i> indicates faithfulness. For Recursive ROAR a curve above or similar to <i>random</i> indicates a non-faithful importance measure, while for ROAR by Hooker et al. [14] this case is inconclusive. Performance is averaged over 5 seeds with a 95% confidence interval. . . . .	162

Figure C.1	The all aggregation for the 100% masked performance and unmasked performance. The baseline (dashed line) for 100% masked performance is the class-majority baseline. Unmasked performance is when using no masking for both validation and training. . . . .	169
Figure C.2	The unmasked performance for each validation and training strategy, using the test dataset. Not that “ <i>No masking</i> ” as a <i>training strategy</i> is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task <i>training strategy</i> “ <i>Use 50/50</i> ” works best. This plot is <b>page-1</b> . Corresponding main results in Figure 4.3. . . . .	170
Figure C.3	The unmasked performance for each validation and training strategy, using the test dataset. Not that “ <i>No masking</i> ” as a <i>training strategy</i> is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task <i>training strategy</i> “ <i>Use 50/50</i> ” works best. This plot is <b>page-2</b> . Corresponding main results in Figure 4.3. . . . .	171
Figure C.4	The 100% masked performance, using the test dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training (“ <i>Masking</i> ” or “ <i>Use 50/50</i> ”) is necessary. This plot is <b>page-1</b> . Corresponding main results in Figure 4.4. . . . .	172
Figure C.5	The 100% masked performance, using the test dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training (“ <i>Masking</i> ” or “ <i>Use 50/50</i> ”) is necessary. This plot is <b>page-2</b> . Corresponding main results in Figure 4.4. . . . .	173
Figure C.6	The unmasked performance for each validation and training strategy, using the validation dataset. Not that “ <i>No masking</i> ” as a <i>training strategy</i> is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task <i>training strategy</i> “ <i>Use 50/50</i> ” works best. This plot is <b>page-1</b> . . . . .	174
Figure C.7	The unmasked performance for each validation and training strategy, using the validation dataset. Not that “ <i>No masking</i> ” as a <i>training strategy</i> is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task <i>training strategy</i> “ <i>Use 50/50</i> ” works best. This plot is <b>page-2</b> . . . . .	175
Figure C.8	The 100% masked performance, using the validation dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training (“ <i>Masking</i> ” or “ <i>Use 50/50</i> ”) is necessary. This plot is <b>page-1</b> . . . . .	176

Figure C.9	The 100% masked performance, using the validation dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training ( <i>“Masking”</i> or <i>“Use 50/50”</i> ) is necessary. This plot is <b>page-2</b> . . . . .	177
Figure C.10	The validation performance for each epoch. Note that the max number of epochs vary depending on the dataset. This is only to limit the compute requirements when fine-tuning. The best epoch is selected by the “early-stopping” dataset, which has one copy with no masking and one copy with uniformly sampled masking ratios. This plot is <b>page-1</b> .	178
Figure C.11	The validation performance for each epoch. Note that the max number of epochs vary depending on the dataset. This is only to limit the compute requirements when fine-tuning. The best epoch is selected by the “early-stopping” dataset, which has one copy with no masking and one copy with uniformly sampled masking ratios. This plot is <b>page-2</b> .	179
Figure C.12	In-distribution p-values using MaSF, for <b>RoBERTa-base</b> with and without masked fine-tuning, <b>page-1</b> . The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5.	180
Figure C.13	In-distribution p-values using MaSF, for <b>RoBERTa-base</b> with and without masked fine-tuning, <b>page-2</b> . The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5.	181
Figure C.14	In-distribution p-values using MaSF, for <b>RoBERTa-large</b> with and without masked fine-tuning, <b>page-1</b> . The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5.	182

- Figure C.15 In-distribution p-values using MaSF, for **RoBERTa-large** with and without masked fine-tuning, **page-2**. The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5. 183
- Figure C.16 The performance given the masked datasets, where masking is done for the x% allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “*Random*” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-1** for **RoBERTa-base**. Corresponding main results in Section 4.2.3. . . . . 184
- Figure C.17 The performance given the masked datasets, where masking is done for the x% allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “*Random*” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-2** for **RoBERTa-base**. Corresponding main results in Section 4.2.3. . . . . 185
- Figure C.18 The performance given the masked datasets, where masking is done for the x% allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “*Random*” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-1** for **RoBERTa-large**. Corresponding main results in Section 4.2.3. . . . . 186
- Figure C.19 The performance given the masked datasets, where masking is done for the x% allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “*Random*” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-2** for **RoBERTa-large**. Corresponding main results in Section 4.2.3. . . . . 187
- Figure D.1 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: explicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. 192

Figure D.2	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: explicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. 193
Figure D.3	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: explicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. 194
Figure D.4	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: implicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. 195
Figure D.5	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: implicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. 196
Figure D.6	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: implicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. 197
Figure D.7	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: “redacted””. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . . 198
Figure D.8	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: “redacted””. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . . 199
Figure D.9	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: “redacted””. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . . 200

Figure D.10	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	201
Figure D.11	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	202
Figure D.12	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	203
Figure D.13	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	204
Figure D.14	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	205
Figure D.15	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	206
Figure D.16	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	207

Figure D.17	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: “removed””. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	208
Figure D.18	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: “removed””. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	209
Figure D.19	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: explicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	210
Figure D.20	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: explicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	211
Figure D.21	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: explicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . .	212
Figure D.22	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: implicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	213
Figure D.23	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: implicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “faithful”. . . . .	214

Figure D.24	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: implicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . .	215
Figure D.25	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "redacted"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	216
Figure D.26	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "redacted"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	217
Figure D.27	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "redacted"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	218
Figure D.28	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	219
Figure D.29	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "removed"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	220
Figure D.30	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "removed"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	221

Figure D.31	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: “redacted””. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	222
Figure D.32	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: “redacted””. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	223
Figure D.33	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: “redacted””. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	224
Figure D.34	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: “removed””. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . .	225
Figure D.35	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: “removed””. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . .	226
Figure D.36	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: “removed””. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . .	227
Figure D.37	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: explicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “not faithful”. . . . .	228
Figure D.38	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: explicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “faithful”. 229	

Figure D.39	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: explicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . .	230
Figure D.40	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: implicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	231
Figure D.41	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: implicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . . .	232
Figure D.42	Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: implicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”. . . .	233
Figure D.43	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: “redacted””. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	234
Figure D.44	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: “redacted””. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	235
Figure D.45	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: “redacted””. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	236

Figure D.46	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	237
Figure D.47	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "removed"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	238
Figure D.48	Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "removed"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	239
Figure D.49	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "redacted"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	240
Figure D.50	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "redacted"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	241
Figure D.51	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "redacted"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	242
Figure D.52	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	243
Figure D.53	Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "removed"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”. . . . .	244

Figure D.54      Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "removed"”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “faithful”. . . . . 245

## LIST OF SYMBOLS AND ACRONYMS

### Publication venues:

---

ACL	Association for Computational Linguistics
ACM	Association for Computing Machinery
BlackboxNLP	Workshop on analyzing and interpreting neural networks for NLP
EMNLP	Empirical Methods in Natural Language Processing
ICLR	International Conference on Learning Representations
ICML	International Conference on Machine Learning
NeurIPS	Advances in Neural Information Processing Systems

### Government related terms:

---

COVID-19	Corona Virus 2019
EU	European Union
GDPR	General Data Protection Regulation
U.S.	United States of America
UK	United Kingdom

### Field acronyms:

---

AI	Artificial Intelligence
HCI	Human Computer Interaction
ML	Machine Learning
NLP	Natural Language Processing

### Datasets:

---

bAbI	Set of synthetic toy datasets [10]
BoolQ	Boolean Question Answering dataset [22]
CB	CommitmentBank – NLI dataset [23]
CoLA	Corpus of Linguistic Acceptability [24]
CQA	multiple-choice question-answering dataset [25]
GLUE	General Language Understanding Evaluation [7]
IMDB	International Movie Database – sentiment dataset [26].
MNLI	Multi-Genre NLI [12]
MRPC	Microsoft Research Paragraph Corpus [27]

MCTest	Multi Choice dataset [28]
NLI	Natural Language Inference
RTE	Recognizing Textual Entailment [29]
SNLI	NLI dataset [30]
SST	Stanford Sentiment Treebank – sentiment dataset [13]
SST2	Two class version of the SST dataset
QA	Question and Answering
QNLI	Question NLI [31]
QQP	Quora Question Pairs – Duplicate Question Detection [32]

---

Models:

BERT	Bidirectional Encoder Representations from Transformers [21]
BiLSTM	Bi-directional Long-Short Term Memory
GPT	Generative Pre-trained Transformer
LSTM	Long-Short Term Memory [33]
Model-XB	For example, Llama2-70B refers to the 70 Billion parameter version of the Llama2 model.
RoBERTa	Robustly optimized BERT [4]
T5	Text-to-Text Transfer Transformer [34]

---

Explanation methods:

abs	Absolute importance measure, cannot separate positive and negative contributions – Section 2.5
Beam	Beam-search based explanation optimization method [35]
CAGE	Commonsense Auto-Generated Explanation [36]
CoT	Chain-of-Thought
Grad	Gradient with respect to the input [37, 38]
IG	Integrated Gradient [39]
LOO	Leave-on-out [40]
LIME	Local Interpretable Model-agnostic Explanations [41].
MiCE	Minimal Contrastive Editing [42]
NILE	Natural Language Inference with Faithful Natural Language Explanations [43]
sign	Signed importance measure, can separate positive and negative contributions – Section 2.5
SHAP	SHapley Additive exPlanation [44]

$x \odot \text{Grad}$       Input-times-gradient [45]

Categories of language models:

---

CLM	Causal Language Model
LMM	Large Language Model
MLM	Masked Language Model

Acronyms introduced in this thesis:

---

AUC	Area Between Curves – Section 3.4.1
FMM	Faithfulness Measurable model – Chapter 4
RACU	Relative Area Between Curves – Section 3.4.1

Miscellaneous:

---

CDF	Communicative Density Function
IM	Importance Measure
MaSF	Max-Simes-Fisher, out-of-distribution detection method [46]
OOD	Out of Distribution
i.i.d.	Independent and identically distributed
ROAR	Remove And Retrain [14]
TGI	Text Generation Inference – A tool for performing inference on LLMs, see <a href="https://github.com/huggingface/text-generation-inference">https://github.com/huggingface/text-generation-inference</a>
POS	Part-of-speech

Mathematical symbols:

---

$y$	Target label.
$\mathbf{x}$	Inputs, often vector of tokens.
$\tilde{\mathbf{x}}_{-i}$	Inputs, with token at position $i$ removed or masked.
$\nabla_x f(x)$	The gradient of function $f$ with respect to $x$ . If $\mathbf{x}$ is a vector, the gradient will be a vector too.
$p(y x)$	Density of label $y$ given $\mathbf{x}$ .
$\mathbb{P}(Z \leq z)$	Probability of realized value from the stochastic distribution $Z$ being less than $z$ .
$\theta$	All parameters of a model.
$\mathbf{W}$	Selected model parameters, typically in the shape of a matrix.
$\alpha$	Attention weight.
$E(x, c)$	a Local explanation of the input $\mathbf{x}$ for the predicted class $c$ .

$\sigma(z)$	The sigmoid function of $z$ .
$\odot$	Element-wise multiplication, also known as Hadamard product.
$\mathcal{O}(C)$	Expresses upper-bound complexity, with complexity $C$ .
$\mathbf{A}^\top$	Matrix $A$ transposed.
$\mathbf{A}_{r,:}$	The $r$ 'th row vector of matrix $A$ .
$\mathbf{A}_{:,c}$	The $c$ 'th columns vector of matrix $A$ .
$\mathbf{A}_{r,c}$	Scalar at row $r$ and column $c$ of matrix $A$ .

---

**LIST OF APPENDICES**

Appendix A	Literature Review, other communication methods . . . . .	131
Appendix B	General-purpose faithfulness metric for importance measures . . . . .	153
Appendix C	Faithfulness measurable models . . . . .	163
Appendix D	Faithfulness of self-explanations . . . . .	190

## CHAPTER 1 INTRODUCTION

Machine Learning (ML) is increasingly being used by the industry to perform automatic decisions that affect most lives [47]. This can have both positive and negative consequences for individuals and society. Along with this trend, the machine learning models have also become more complex and thus harder to understand [5]. Today’s neural networks utilize billions of parameters and provide no direct mechanism to ensure that they behave as intended. The consequences of this can be catastrophic, Rudin [48] explains “There have been cases of people incorrectly denied parole [49], poor bail decisions leading to the release of dangerous criminals, ML-based pollution models stating that highly polluted air was safe to breathe [50] and generally poor use of limited valuable resources in criminal justice, medicine, energy reliability, finance and in other domains [51]”.

Within Natural Language Processing (NLP), applications such as translation, dialog systems, resume screening, search, etc. [52] also suffer from ethical issues. For many of these applications, neural models have been shown to exhibit unwanted biases and other ethical issues [48, 53, 54, 55, 56, 57].

To combat these issues, there are increasing legal works, such as the GDPR in the EU, that mandate that automatic model decisions must be accompanied by “meaningful information about the logic involved” [52]. The field often responsible for achieving these explanations is called *interpretability*.

Doshi-Velez and Kim [58] define *interpretability* as the “ability to explain or to present in understandable terms to a human”. While this field has been around for a while, since statistics and decision trees, the ever-increasing demand for predictive performance and complex capabilities, such as chatting, has dramatically increased neural networks’ complexity.

While many interpretability methods have been proposed to satisfy this need, the explanations are often found not to be *faithful*, meaning they do not reflect the true reasoning process of the model they explain [59]. Unfortunately, as discussed in Section 1.3, and Chapter 2, there is a pattern in the interpretability field, where interpretability methods are proposed and then later debunked through analysis of their faithfulness. For example, an explanation algorithm indicating which input tokens are important for a prediction is later found to be less faithful than simply pointing at random tokens (Chapter 3).

This trend likely happens because interpretability methods are often proposed before their corresponding faithfulness metrics are well developed. This is a consequence of the ground-

truth explanation being inaccessible to humans, as the models are too complex to be manually analyzed. As a result, it’s often necessary to measure faithfulness using proxies [59]. For example, if this token is truly important, the prediction should change significantly when the token is removed [60]. However, even such a measurable definition has a lot of nuance to it, such as: does token removal cause out-of-distribution issues, what is a significant change, and what if multiple tokens need to be removed for the prediction to change (Chapter 3 and Chapter 4)?

Besides the challenges of the faithfulness metric, there is also the fundamental question of when an interpretability method is faithful. Currently, there is no universally recognized answer to this question; rather, there are two competing perspectives on this, namely *post-hoc* and *intrinsic* [61].

Section 1.2 properly describe these perspectives. Put briefly, the *intrinsic* perspective says that only models architecturally designed to be explained can be explained [48]. In contrast, the *post-hoc* perspective says this constraint is unnecessary and too restrictive to achieve competitive performance.

This thesis posits that these perspectives should be considered as paradigms, where a paradigm is “universally recognized scientific achievements that, for a time, provide model problems and solutions to a community of practitioners” [62].

Furthermore, it is the hypothesis of this thesis that neither paradigm has been fruitful because their underlying beliefs are problematic or unnecessary, and we should, therefore, embrace that paradigms only exist “for a time” [62] and look for new directions. Section 1.3 contains the primary support for this hypothesis.

This thesis seeks to propose and develop such new paradigms for interpretability in NLP. At the core of this introductory discussion and the thesis is how each paradigm approaches faithfulness. Faithfulness is particularly important, as false but convincing explanations can lead to unsupported confidence in models, increasing the risk of AI.

## 1.1 Why interpretability is needed

Before discussing the current paradigms and their shortcomings, it’s necessary to consider whether interpretability is needed. Many ethical motivations for interpretability are also served by bias and fairness metrics, so if the current paradigms of interpretability do not work (as we argue in Section 1.3), perhaps we should drop the idea of interpretability completely. If the models can be made accurate, unbiased, and fair enough, do we need to explain the models? In this section, we will argue that interpretability is required by examining the

limitations of bias and fairness metrics and the scientific motivations for interpretability.

### **1.1.1 Limitations of bias and fairness metrics**

There is no doubt that bias and fairness metrics present a vital role in validating models' behavior. However, a shared limitation is that they always measure known attributes [63]. For example, gender-bias metrics use gender attributes. This presents two challenges. Can we procure such attributes (known as protected attributes)? How do we prevent unanticipated biases?

#### **Protected attribute procurement**

Attributes like gender, race, age, disability, etc., are under U.S. law known as “protected attributes” [64], and collecting and using these attributes is heavily regulated in most of the world. Andrus et al. [65] writes, “In many situations, however, information about demographics can be extremely difficult for practitioners to even procure.”. Therefore, systematically measuring bias and fairness is not always practical [65].

On the other hand, explanations often don't depend on knowing these protected attributes in advance and can provide a more qualitative analysis. For example, suppose an explanation tells us that the word “Woman” from “Member of Woman's Chess Club” in a resume is important for making a hiring recommendation. In that case, there is a potential harmful bias [66]. Therefore, explanations can serve a similar practical purpose to a fairness or bias metric without performing systematical correlations.

#### **Unknown attribute bias**

Although protected attributes are important to consider and are often legally protected, many more relevant attributes are involved in ensuring a fair and unbiased system. Unfortunately, it is impossible to consider every possible bias in advance. As an alternative, interpretability offers a more qualitative and explorative validation.

Continuing the example with resumes and automated hiring recommendations, during investigations by Fuller et al. [67], the authors found that a hospital only accepted candidates with computer programming experience when they needed workers to enter patient data into a computer. Another example was a clerk position where applicants were rejected if they did not mention floor-buffing (i.e., a cleaning method for floors, which had no relevance to the position) [68].

These examples present cases of systematic unintended bias. However, they do not relate to any protected attributes, and they are so specific they can only be discovered through qualitative explanations and investigations. That said, systematic fairness/bias metrics can quantify the damage once potential biases are identified using interpretability. Afterward, those metrics can be integrated into a quality assessment system to prevent future harm.

### **1.1.2 Interpretability for scientific discovery and understanding**

Interpretability is not only used for ethics and adjacent purposes, where bias and fairness metrics have an important role. Interpretability is also used for scientific discovery and learning about what makes models work.

#### **Scientific Discovery**

An example of scientific discovery is interpretability in drug discovery [69, 70, 71]. A common approach is to use importance measures to identify regions in genomic sequences responsible for a particular behavior, such as producing a protein. While these explanations do not guarantee that such connections exist in reality, they can provide important initial hypotheses for scientists, enabling them to make more informed choices about the direction of their research.

#### **Model understanding**

An emerging field of interpretability is mechanistic interpretability, which identifies parts of a neural network that have a particular responsibility [72]. For example, identifying a collection of neurons responsible for copying content in a generative language model, etc. [73]. Such insights may not be directly relevant to downstream tasks, but they help us understand current model limitations and can lead to better model design.

## **1.2 The current paradigms of interpretability**

This thesis uses a common definition of interpretability, “the ability to explain or to present in understandable terms to a human” by [58]. However, even this definition of interpretability is not agreed upon.

Lipton says, “the term interpretability holds no agreed upon meaning, and yet machine learning conferences frequently publish papers which wield the term in a quasi-mathematical way” [61]. In 2017, a UK Government House of Lords review of AI noted after substantial

Table 1.1 Comparison of the definitions and underlying beliefs of the intrinsic and post-hoc paradigms. The beliefs relate to a) requirements for a faithful explanation and b) model capabilities. It should be apparent that these two views are seemingly incompatible.

	Intrinsic paradigm	Post-hoc paradigm
definition	The model is designed to provide explanations by making the explanation part of the model architecture.	The model is produced without regard for explanation, and the explanations are then created after model training.
underlying beliefs	<p>Only models that were designed to be explained can be explained.</p> <p>Intrinsic models can have the same performance as a black-box model.</p>	<p>Although it may be very challenging, black-box models can be explained.</p> <p>Black-box models will be more generally applicable than intrinsic models.</p>

expert evidence that “the terminology used by our witnesses varied widely. Many used the term transparency, while others used interpretability or explainability, sometimes interchangeably” [74, 91].

For this reason, there are also no clearly agreed-upon definitions of the current paradigms of interpretability [75, 76]. As such, this section defines the *intrinsic* and *post-hoc* paradigms, as well as describe their underlying beliefs, which are summarized in Table 1.1.

### 1.2.1 Definitions

Jacovi and Goldberg [59] write: “A distinction is often made between two methods of interpretability: (1) interpreting existing models via post-hoc techniques; and (2) designing inherently interpretable models. [48]”. Based on this and other sources [75, 77, 78], this thesis refers to these two ideas as 1) the *intrinsic* paradigm and 2) the *post-hoc* paradigm.

#### The intrinsic paradigm

The intrinsic paradigm works on creating so-called *inherently interpretable models*. These models are architecturally constrained, such that the explanation emerges from the architecture itself.

Classical examples are decision trees or linear regression. In the field of neural networks, some examples are 1) “Old-school” attention [2, 79], where attention points to which input tokens are important. 2) Neural Modular Networks [80, 81, 82], which produce a prediction via a

sequence of sub-models, each with known behavior. 3) Prototypical Networks [83, 84, 85], which predicts by finding similar training observations.

## The post-hoc paradigm

*Post-hoc* explanations are computed after the model has been trained. They are developed independently of the model’s architecture and how it was trained. However, some simple criteria often exist, like “the model should be differentiable”, “the training dataset is known”, or “inputs are represented as tokens”. Although general applicability is technically not a requirement, if a method is so specific that it only works on one specific model, it is likely an *intrinsic explanation*.

As an example, a common post-hoc explanation is gradient-based importance measures. Importance measures explain which input features (words, pixels, etc.) are important for making a prediction. This is achieved by differentiating the prediction with respect to the input. The idea is that if a small change in input causes a big change in the output, then that input is important [37, 86, 87].

### 1.2.2 Beliefs

As with all paradigms, there are fundamental underlying beliefs, which are why the paradigm’s followers partake in their paradigm of choice. At the core of these beliefs are two central questions. When are explanations faithful, what are the requirements for faithfulness, and how do these requirements affect the model’s general performance capabilities?

#### When are explanations faithful?

The intrinsic paradigm believes that: *only models designed to be explained, can be explained*, which their *inherently interpretable models* try to satisfy. Therefore, they argue that using black-box models is too risky, as these models can never be faithfully explained [48].

However, although their models are designed to be intrinsically explainable, this claim and their faithfulness should still be questioned [59], as many inherently interpretable model ideas are later revealed not to provide faithful explanations. For example, attention-based explanations have received notable criticism for not being faithful [2, 88, 89, 90, 91]. This is discussed more in Section 1.3.1.

The *post-hoc explanation* paradigm takes a less strict stance and believes that even models that were not designed to be explained (i.e., black-box models) can still be explained. However, as

this paradigm has no control over the model, achieving faithful explanations is very challenging; this is discussed more in Section 1.3.2.

In conclusion, the intrinsic paradigm considers explanations to be part of the model design, and post-hoc explanations are always applied after the model design and training. Hence, the two schools of thought are incompatible frameworks and can philosophically be considered paradigms [62].

### **What is the effect on the model’s general performance capabilities?**

It would seem that *intrinsic explanation* is the obvious choice. If we can control the model such that the faithfulness of explanations can be guaranteed, why consider *post-hoc explanation*?

The commonly mentioned idea is that the *post-hoc* paradigm believes that by constraining the models in the manners that the *intrinsic paradigm* requires, there is a trade-off in performance [92]. However, this trade-off does not have to be the case in practice [48, section 2].

A more accurate take, which is rarely explicitly discussed, is that the common industry prefers off-the-shelf general-purpose models and only later thinks about interpretability [47]. Additionally, most research only considers predictive performance, not interpretability. Therefore, *intrinsic* researchers are always catching up to black-box models. From the *post-hoc* perspective, it would make more sense to work on generally applicable interpretability methods for both off-the-shelf and future black-box models.

From the intrinsic perspective, while the industry might prefer off-the-shelf models now, they shouldn’t. Not validating models through intrinsic explanations can have serious consequences [48] and eventually damage their business. Additionally, with increasing legal requirements to provide explanations, the industry may have to use inherently explainable models [93].

For these reasons, the *intrinsic* paradigm believes we should not let the industry’s needs dictate our research direction, as their goals may be too short-sighted. In the long run, intrinsic models may be the only reasonable option.

In conclusion, the *post-hoc* paradigm has good intentions of providing general explanations for general-purpose models. However, from the *intrinsic* paradigm perspective, those good intentions are meaningless if it is fundamentally impossible to provide guaranteed faithful explanations without an *inherently interpretable model*.

### 1.3 Why interpretability needs a new paradigm

When there are multiple paradigms, it tends to be the case that neither of the paradigms fits the needs. However, for the case of the *post-hoc* and *intrinsic* paradigms, it could be argued that they serve different needs. For example, *intrinsic* explanations should be preferred for critical applications [48], and *post-hoc* explanations could be used for verifiable situations, such as drug discovery, where the hypothesis generated by the explanations is verified using physical experiments.

#### 1.3.1 The case against the intrinsic paradigm

The industry primarily uses post-hoc explanations, including for high-stakes applications such as insurance risk assessment and financial loan assessment [47, 94]. This is because such industries usually do not have the in-house expertise to develop custom, inherently interpretable models. They must rely on basic inherently interpretable models, like decision trees, which are not competitive or use more advanced off-the-shelf neural black-box models, like pre-trained language models, which will be competitive. In practice, the industry is thus not in a position to choose inherently interpretable models.

Another challenge with the intrinsic paradigm is that its models are often not completely interpretable because only a part of the model is architecturally constrained to be interpretable. The rest of the neural network still uses black-box components (e.g., Dense layer, Recurrent layer, etc.) which are not interpretable. As such, the intrinsic promise should not be taken at face value [59].

An example of this is classic attention-based models [2, 79]. Attention itself is interpretable, as it's a weighted sum that explains each intermediate representation's importance. However, attention is often used for token-importance. This does not work, as the intermediate representations are produced by a black-box recurrent neural network (e.g., LSTM [33]), which can mix or move the relationship between tokens and the intermediate representations. Therefore, the attention scores do not necessarily represent token-importance [95].

Likewise, Neural Modular Networks produce an executable problem composed of sub-networks, such as `find-max-num(filter(find()))`, which is interpretable [80, 81, 82]. However, each sub-networks (`find-max-num`, `filter`, `find`) is itself a black-box model with little guarantee that it operates as intended [96, 97, 98].

Overall, there are few success stories with intrinsic explanations. They are usually either not performance-wise competitive, general-purpose enough for the industry [47], or their intrinsic claims are unsupported [59].

### 1.3.2 The case against the post-hoc paradigm

Although post-hoc explanations directly address the interpretability challenge of black-box components and models, and could therefore provide more complete explanations, there are very few success stories with post-hoc, where post-hoc explanations are consistently faithful.

Most notable is perhaps post-hoc importance measure (IM) explanations, where the explanation indicates which input features are the most important for making a prediction. The pursuit of such explanations have produced countless papers [39, 40, 41, 99, 100, 101, 102, 103, 104, 105].

However, repeatedly, the faithfulness of these IM explanations is criticized [14, 45, 106, 107, 108, 109]. For example, different allegedly faithful IMs often disagree on which inputs are important, an issues known as the disagreement problem which is hard to reconcile [2, 94]. There are also theoretical works that suggest that IMs are subject to a *no free lunch theorem* [110], or it may be impossible to provide faithful post-hoc IMs [111]. This thesis will also demonstrate empirically that the faithfulness of *post-hoc* explanations is model- and task-dependent, and they, therefore, don't provide the general capabilities that the *post-hoc* paradigm desires.

Similar to the work of IM, there are visualizations of neurons in computer vision, which shows that neurons represent high-level concepts, such as nose or dog. This is done by visualizing convolutional weights or the input image that maximizes a neuron's activation [112, 113, 114], which provides very convincing evidence. However, it has been shown empirically, theoretically, and through human-computer-interaction (HCI) studies that these visualizations do not provide more useful explanations regarding the neurons' responsibility than simply using existing images [115, 116, 117]<sup>1</sup>.

Another notable example is probing explanations, where models are verified by relating the model's behavior or intermediate representation to, for example, linguistic properties (part-of-speech, etc.) [118, 119]. This idea has produced an entire subfield called BERTology [120]. BERTology, in particular, has attained substantial attention [11, 20, 120, 121, 122, 123, 124], with most of the works finding that neural networks can learn linguistic properties indirectly.

Unfortunately, like post-hoc importance measures, there are many reasons to be highly skeptical [125]. For example, using an untrained model or a randomized dataset shows an equally high correlation with linguistic properties, compared with training a regular model [126, 127]. These discoveries have put the entire methodology into question, although there is

---

<sup>1</sup>Neural networks likely do encode high-level concepts, but these visualizations are not useful for identifying the responsibility of specific neurons.

work trying to adapt to these new critics [128].

### 1.3.3 Overall trend

Post-hoc importance measures and probing explanations are just two cases where post-hoc shows initial promise through countless papers, only to be debunked repeatedly. The trend oscillates between proposing new explanation methods and debunking them. Of course, proving that there will never be a great post-hoc method is impossible. However, the lack of guarantees also makes it impossible to know when a faithful post-hoc method is proposed. Similarly, intrinsic explanations also receive criticism after a while, as has been the case with attention and Neural Modular Networks.

## 1.4 Thesis overview

Although both the intrinsic and post-hoc paradigms have significant issues, parts of their underlying beliefs have merit. The intrinsic paradigm believes that *only models designed to be explained can be explained*, while post-hoc believes that *black-box models tend to be more general purpose while providing high predictive performance*. These beliefs have merit; therefore, the first foundational idea of this thesis is to develop new paradigms that incorporate their spirit.

Secondly, there has been a trend where interpretability methods have been proposed and later debunked, as is discussed in Chapter 2. This has particularly been the case for importance measures. This indicates that there is something fundamental about faithfulness which is not understood. Therefore, the second foundational idea for this thesis is to focus on faithfulness metrics first, especially for important measures that have received the most criticism, and then take what is learned from this investigation to develop better interpretability methods.

Finally, providing interpretability methods and metrics for all neural network literature is likely too big of a scope. Instead, this thesis restricts the research to natural language tasks and NLP models. The hope is that by restricting the ambitions to natural language, it will be possible to take advantage of some of the intrinsic properties within natural language or natural language models. Although, the hope is still that the methods and metrics could be somewhat domain-agnostic.

The following research question and hypothesis frame these ideas and the thesis:

**Research question:** How to provide and ensure faithful explanations for complex general-purpose neural NLP models?

**Research hypothesis:** By developing new paradigms that design models to be explained without employing architectural constraints, by focusing on developing accurate faithfulness metrics, by focusing on importance measures that have had a notoriously troubling history regarding faithfulness, and by taking advantage of properties specific to natural language and NLP models, it is possible to learn and inform how faithful explanations for complex general-purpose neural NLP models can be produced.

Using this research question and hypothesis, this thesis presents two new paradigms, the *faithfulness measurable model* (FMM) paradigm and the *self-explanation* paradigm, which are summarized in Table 1.2 and presented in Chapter 4 and Chapter 5 respectively.

Table 1.2 Comparison of the definitions and underlying beliefs of the new paradigms. The beliefs relate to a) explanation requirements and b) model capabilities. These new paradigms can be compared with the old paradigms in Table 1.1.

	Faithfulness measurable model paradigm	Self-explanation paradigm
definition	The model is designed to measure the faithfulness of a category of explanations.	The model can produce both its prediction and an explanation for that prediction.
underlying beliefs	<p>It is computationally feasible to optimize explanations for optimal faithfulness.</p> <p>Models can be optimized to be faithfulness measurable without loss of predictive performance.</p>	<p>Models can be trained to model and articulate their own reasoning accurately and will generalize.</p> <p>Self-explanation capabilities do not negatively impact regular predictions.</p>

#### 1.4.1 Content overview

**Chapter 2** is the background chapter, which covers most types of explanations, how selected interpretability methods produce them, and the limitations or potential issues of those interpretability methods, in particular from a faithfulness perspective. As this thesis primarily focuses on importance measures and secondarily on counterfactuals, extra attention is given to these types of explanations.

**Chapter 3** develops a general purpose faithfulness metrics for importance measures; this is then applied to both common post-hoc and intrinsic methods and concludes that the faithfulness of an interpretability method is both model and task-dependent.

Motivated by the model and task-dependent conclusion from Chapter 3, it’s clear that it will be necessary to measure the faithfulness for the specific model and task at hand. The methodology presented in Chapter 3 does not support this because it’s computationally expensive. Even if computing was not a problem, another limitation is that it measures faithfulness using different models than the one that would be in deployment, thus risking unsupported confidence in the faithfulness of an interpretability method.

**Chapter 4** solves these limitations and more by proposing a new paradigm, namely *faithfulness measurable models* (FMMs). The idea with FMMs is that the model inherently provides the means to cost-effectively and precisely measure the faithfulness of an explanation. This is different from the intrinsic paradigm, which inherently provides the means to compute an explanation. Chapter 4 shows that reframing the interpretability problem from “inherently explainable” to “inherently measurable” means that the model does not require architectural constraints. Finally, because faithfulness is now easy and cheap to measure, optimizing an explanation towards maximal faithfulness is possible, thus making FMMs indirectly designed to be explained. This property and the lack of architectural constraints means that this paradigm archives the goal of taking the best part from both paradigms.

Chapter 4 analyses explanations produced by interpretability methods and algorithms. However, with the emergence of instruction-tuned Large Language Models (LLMs), there is a new type of explanation approach where the model explains itself, so-called self-explanations. This approach presents a new interpretability paradigm. In principle, this direction could be productive because if LLMs do have reasoning capabilities, they should be in a better position to explain themselves than anything else. Mathematically, it has access to all its weights and the input, which produces the prediction. In this sense, the model is designed to explain itself and is also a very general-purpose model. However, the self-explanations are also produced by a black box, which creates a potential danger as they may be hallucinated. Finally, many of the LLMs are only accessible via APIs, meaning using most regular explanation methods or turning them into FMMs may be infeasible.

**Chapter 5** addresses these concerns by taking the lessons from the survey in Chapter 2 and the faithfulness metric in Chapter 3. It proposes a methodology for measuring the faithfulness of self-explanations produced by instruction-tuned LLMs, using only API access to the model. This methodology is then applied to variants of the importance measure idea and counterfactual explanations. The results show, similar to Chapter 3, that faithful is model and task-dependent, and additionally explanation-dependent. This conclusion means that we should not trust self-explanations in general, highlighting the risk of using LLMs for high-stakes decisions.

### 1.4.2 Scientific contributions

This thesis presents many findings and scientific contributions, which each chapter discusses in detail. The most major contributions are:

- To develop an extensive survey focusing on the faithfulness and human-groundedness of each important interpretability method. Demonstrating a critical mindset that the field should adopt.
- To communicate and support the position that interpretability should develop new paradigms regarding faithfulness rather than stick to intrinsic and post-hoc.
- To develop solid general-purpose faithfulness measures for both regular importance measures and self-explanations (herein including counterfactuals), with well-documented limitations.
- To make the observation that faithfulness tends to be model and task-dependent, both regarding regular importance measures and self-explanations. Additionally, in the context of self-explanations, faithfulness is also explanation-dependent. Thus, arriving at the conclusion, we should not trust self-explanations and importance measures in general.
- To propose and develop the *faithfulness measurable model* (FMM) paradigm, which provides a new direction for interpretability, that both designs models to be explained and also works with general purpose black-box models.
- To demonstrate the *faithfulness measurable model* paradigm with masked language models and importance measure explanations. The results show significant improvements across all tasks, and for synthetic datasets, achieves near theoretical perfect faithfulness.

### 1.4.3 Statement of contributions in papers

This thesis includes content from 5 papers, the personal contributions of Andreas Madsen and co-authors are enumerated here. Written confirmation that the following is true was obtained prior to submitting this thesis.

**Chapter 1** borrows discussion from our position paper titled “AI Interpretability Needs a New Paradigm” [129] submitted to “Communications of the ACM” and written by Andreas Madsen, Himabindu Lakkaraju, Siva Reddy, and Sarath Chandar.

- Andreas Madsen’s contributions: identified paradigms, the position, and wrote the entire paper.
- Prof. Himabindu Lakkaraju’s contributions: Provided references to theoretical post-hoc criticisms, the learn-to-faithfulness explain paradigm, and the self-explanation paradigm; identified regularization techniques as belonging to learn-to-faithfulness, provided overall feedback to the tone and message of the paper.
- Prof. Siva Reddy’s contributions: reviewed paper, provided feedback to the self-explanation paradigm.
- Prof. Sarath Chandar’s contributions: proposed to write a position paper; provided feedback on paradigms, paper structure, collaborations, and submission venues; and reviewed paper.

**Chapter 2** borrows heavily from our survey titled “Post-hoc Interpretability for Neural NLP: A Survey” [130] published in *ACM Computing Surveys* and written by Andreas Madsen, Siva Reddy, and Sarath Chandar.

- Andreas Madsen’s contributions: wrote the entire survey, proposed and developed a categorization system, and performed a literature review.
- Prof. Siva Reddy proposed papers to include and provided feedback on the writing and categorization.
- Prof. Sarath Chandar proposed papers to include and provided feedback on the writing and categorization.

**Chapter 3** contains results and discussion from the paper titled “Evaluating the Faithfulness of Importance Measures in NLP by Recursively Masking Allegedly Important Tokens and Retraining” [131] published at *Findings in EMNLP 2022* and *BlackboxNLP 2022*; and written by Andreas Madsen, Vaibhav, Adlakha, Nicholas Meade, and Siva Reddy.

- Andreas Madsen’s contributions: written the entire paper, except for the method section. Implemented SST, Anemia, and Diabetese datasets; single-sequence model; gradient, integrated-gradient, mutual-information explanations; ROAR and Recursive ROAR; faithfulness metric; HPC setup; and plots.
- Vaibhav Adlakha’s contributions: written method section and implemented SNLI, IMDB, bAbI datasets, paired-sequence model.

- Nicholas Meade’s contributions: investigated and implemented attention sparsity analysis, integrated-gradient hyper-parameters analysis, and performed code review.
- Prof. Siva Reddy’s contributions: proposed using mutual-information and to develop a faithfulness metric. Also reviewed and provided feedback on the paper.

**Chapter 4** contains discussion and results from the paper titled “Faithfulness Measurable Masked Language Models” Madsen et al. [132], published at ICML 2024 (received a spotlight award) and written by Andreas Madsen, Siva Reddy, and Sarath Chandar.

- Andreas Madsen’s contributions: proposed the idea, implemented all experiments, and wrote the entire paper.
- Prof. Siva Reddy’s contributions: helped analyze above 100% RACU scores and provided feedback to the paper.
- Prof. Sarath Chandar’s contributions, scoped the research project, identified MaSF for out-of-distribution detection, and provided feedback to the paper.

**Chapter 5** contains discussion and results from the paper titled “Are self-explanations from Large Language Models faithful? ” [133] published at ACL Findings 2024 and written by Andreas Madsen, Sarath Chandar, and Siva Reddy.

- Andreas Madsen’s contributions: developed explanation and self-consistency prompts, implemented all experiments, and wrote entire papers.
- Prof. Sarath Chandar’s contributions: provided feedback on scope, scientific methodology, structure, and reviewed the paper.
- Prof. Siva Reddy’s contributions: proposed the general research direction, and reviewed the paper.

## CHAPTER 2 BACKGROUND

Doshi-Velez and Kim [58] define *interpretability* as the “ability to explain or to present in understandable terms to a human”. However, what constitutes an “understandable” explanation is an interdisciplinary question. An important work from social science by Miller [134] argues that *effective explanations* must be selective in the sense one must select “one or two causes from a sometimes infinite number of causes”. Such observation necessitates organizing interpretability methods by how and what they selectively communicate.

This background chapter presents such an organization in Table 2.1, where each row represents a communication approach. For example, the first row describes *input feature* explanations that communicate what tokens are most relevant for a prediction. Each row is ordered by how abstract the communication approach is, although this is an approximation. Organizing by the method of communication is discussed further in Section 2.1.

The communication approaches covered in this chapter are *input feature*, *counterfactuals*, and *natural language* explanations, as these have relevance to the thesis. For other communication approaches, see Appendix A. Generally, the reader is assumed to understand neural networks and NLP, particularly typical models like LSTM-based and Transformer-based architectures. Alternatively, the reader can reference Devlin et al. [21], Goodfellow et al. [135], Vaswani et al. [136], Graves [137], Jurafsky and Martin [138] to learn about these topics.

Each interpretability method uses different kinds of information to produce its explanation; in Table 2.1, this is indicated by the columns. The columns are ordered by an increasing level of information. Again, this is an inexact ranking but serves as a useful tool to contrast the methods.

Finally, before discussing interpretability methods, Section 2.2 will provide a motivating example to better frame the categorization, and Section 2.3 will cover the general concepts regarding how well interpretability is satisfied.

### 2.1 Organizing by method of communication

As a categorization of communication strategies, it is standard in the interpretability literature to distinguish between methods that explain a single observation, called *local explanations*, and methods that explain the entire model, called *global explanations* [47, 58, 75, 139, 140, 141]. In this background chapter, we also consider an additional category of methods that explains an entire output-class, which we call *class explanations*.

		less information				more information →	
		post-hoc				intrinsic	
		black-box	dataset	gradient	embeddings	white-box	model specific
lower abstraction	local explanation						
	input features	Occlusion-based § 2.5.2		Gradient-based § 2.5.1			Attention-based § 2.5.3
	adversarial examples	SEA <sup>M</sup> § A		HotFlip § A			
	influential examples		Influence Functions <sup>H</sup> § A TracIn <sup>C</sup> § A		Representer Pointers <sup>†</sup> § A		Prototype Networks
	counter-factuals	Polyjuice <sup>M,D</sup> § 2.6.1	MiCE <sup>M</sup> § 2.6.2				
class explanation	natural language	predict-then-explain <sup>M</sup> § 2.7.2				explain-then-predict <sup>M</sup> § 2.7.1	
	concepts					NIE <sup>D</sup> § A	
global explanation							
	vocabulary				Project § A, Rotate § A		
higher abstraction ↓	ensemble	SP-LIME § A					
	linguistic information	Behavioral Probes <sup>D</sup> § A			Structural Probes <sup>D</sup> § A	Structural Probes <sup>D</sup> § A	Auxiliary Task <sup>D</sup>
	rules	SEAR <sup>M</sup> § A	Compositional Explanations of Neurons <sup>†</sup> § A				

Table 2.1 Overview of *post-hoc* interpretability methods, where § indicates the section the method is discussed. Rows describe how the explanation is communicated, while columns describe what information is used to produce the explanation. The order of both rows and columns indicates the level of abstraction and amount of information, respectively. However, this order is only approximate.

**Columns:** *Black-box*: the method only evaluates the model. *Dataset*: the method has access to all training and validation observations. *Gradient*: the gradient of the model is computed. *Embeddings*: the method uses the word embedding matrix. *White-box*: the method knows everything about the model, such as all weights and all operations. However, the method is not specific to a particular architecture. *Model specific*: the method is specific to the architecture. Note that neural models in NLP are usually differentiable and have an embedding matrix. We therefore do not consider these as architectural constraints.

**Superscript:** <sup>C</sup>: Depends on checkpoints during training. <sup>D</sup>: Depends on supplementary dataset. <sup>H</sup>: Depends on second-order derivative. <sup>M</sup>: Depends on supplementary model. <sup>†</sup>: Depends only on dataset and white-box access.

To subdivide these categories further, Table 2.1 orders each communication strategy by its abstraction level. As an example, see Figure 2.1, where an *input features* explanation highlights the input tokens that are most responsible for a prediction; because this must refer to specific tokens, its ability to provide abstract explanations is limited. For a highly abstract explanation, consider the *natural language* category, which explains a prediction using a sentence and can therefore use abstract concepts in its explanation.

explanation	$y$	communication approach
<u>the</u> <u>year</u> 's <u>best</u> <u>and</u> <u>most</u> <u>unpredictable</u> <u>comedy</u>	pos	<i>input feature</i>
<i>unpredictable comedies are funny</i>	-	<i>natural language</i>

Figure 2.1 Fictive visualization of an *input features* explanation which highlights tokens and a *natural language* explanation, applied on a sentiment classification task [7].  $y = \text{pos}$  means the gold label is *positive* sentiment.

Communication methods that have a higher abstraction level are typically easier to understand (more *human-grounded*), but the trade-off is that they may reflect the model’s behavior less (less *faithful*). Because the purpose of interpretability is to communicate the model to a human, this trade-off is necessary [48, 134]. The communication strategy must be decided by considering the applications and to whom the explanation is communicated. In Section 2.3, we discuss *human-groundedness* and *faithfulness* in-depth and how to measure them such that an informed decision can be made.

The organization in Table 2.1 does have some limitations. Firstly, ordering explanation methods by their abstraction level is an approximation, and while *global explanations* are generally more abstract than *local explanations*, this is not always true. For example, the explanation “simply print all weights” (not included in Table 2.1) is arguably the lowest possible abstraction level. However, it’s also a *global explanation*.

## 2.2 Motivating Example

To make the method sections as concrete and comparable as possible, this background chapter will show fictive examples often based on the “Stanford Sentiment Treebank” (SST) dataset [13]. The SST dataset has been modeled using LSTM [7], Self-Attention-based models [21], etc., all of which are popular examples of neural networks.

We use a sequence-to-class problem, as this is what most interpretability methods apply to. Although some are agnostic to the problem type, and others are specific to sequence-

to-sequence problems. Throughout this background chapter, we attempt to highlight which types of problems each method applies to.

		$p(y \mathbf{x}; \theta)$	$y$
$\mathbf{x}$	<u>the</u> <u>year</u> <u>'s</u> <u>best</u> <u>and</u> <u>most</u> <u>unpredictable</u> <u>comedy</u>	0.91	pos
$\mathbf{x}$	<u>we</u> <u>never</u> <u>feel</u> <u>anything</u> <u>for</u> <u>these</u> <u>characters</u>	0.95	neg
$\mathbf{x}$	<u>handsome</u> <u>but</u> <u>unfulfilling</u> <u>suspense</u> <u>drama</u>	0.18	neg

Figure 2.2 Three examples from the SST dataset [13].  $\mathbf{x}$  is the input, with each token denoted by an underline.  $y$  is the gold target label, where **pos** is *positive* and **neg** is *negative* sentiment. Finally,  $p(y|\mathbf{x})$  is the model’s estimate of  $\mathbf{x}$  belonging to category  $y$ . Note that the model predicts the 3rd (last) wrong, indicated with **red**.

The model responsible for the predictions in Figure 2.2 can be explained by asking different questions, each of which communicates a different aspect of the model covered in this chapter’s sections. Sometimes, these explanation relates to a single observation; other times, the explanation relates to the whole model.

**local explanations** explain a single observation:

Input Features	<i>Which tokens are most important for the prediction, Section 2.5.</i>
Adversarial Examples	<i>What would break the model’s prediction, Appendix A.</i>
Influential Examples	<i>What training examples influenced the prediction, Appendix A.</i>
Counterfactuals	<i>What does the model consider a valid opposite example, Section 2.6.</i>
Natural Language	<i>What would a generated natural language explanation be, Section 2.7.</i>

**Class explanations** summarize the model, but only with regard to one selected class:

Concepts	<i>What concepts (e.g. movie genre) can explain a class, Appendix A.</i>
----------	--

**Global explanations** summarize the entire model with regard to a specific aspect:

Vocabulary	<i>How does the model relate words to each other, Appendix A.</i>
Ensemble	<i>What set of local explanations are representative of the model, Appendix A.</i>
Linguistic information	<i>What linguistic information does the model use, Appendix A.</i>
Rules	<i>Which general rules can summarize an aspect of the model, Appendix A.</i>

### 2.3 Measures of Interpretability

Because interpretability is, by definition, about explaining the model to humans [58, 134], and these explanations are often qualitative, it is not clear how to quantitatively evaluate and compare interpretability methods. This ambiguity has led to much discussion. A notable is the *intrinsic* interpretability method *attention*, where different measures of interpretability have been published, resulting in conflicting findings [2, 88, 89, 91, 142].

In general, there is no consensus on how to measure interpretability. However, validation is still paramount. As such, this section attempts to cover the general categories, themes, and methods that have been proposed. Additionally, each method section, starting from *input features*, in Section 2.5, will briefly cover how the authors choose to evaluate their method.

To describe the evaluation strategies, we use the terminology defined by Doshi-Velez and Kim [58], which separates the evaluation of interpretability into three categories which are collectively called *groundedness*: *faithfulness* (also called *functionally-grounded*), *human-grounded*, and *application-grounded*. This categorization reflects the need to have explanations that are useful to humans (*human-grounded*) and accurately reflect the model (*faithfulness*).

**Application-grounded** evaluation is when the interpretability method is evaluated in the environment in which it will be deployed. For example, do the explanations result in higher survival rates in a medical setting, higher grades in a homework-hint system, or a better model in a label-correction setting [58, 143]. Importantly, this evaluation should include the baseline where the explanations are provided by humans.

Due to this approach’s application-specific and time-consuming nature, *application-grounded* evaluation is rarely done in NLP interpretability research. Instead, more synthetic and general evaluation setups are being used, which is what *faithfulness* and *human-grounded* evaluation is about. These categories each provide an important but different aspect for validating interpretability and should, therefore, be used in combination.

**Human-grounded** evaluation checks if the explanations are useful to humans. Unlike *application-grounded*, the task is often simpler, and the task itself can be evaluated immediately. Additionally, expert humans are often not required [58]. In other literature, this is known as *simulatability* [61] and *comprehensibility* [144].

Although *human-grounded* evaluation is much more efficient than *application-grounded* evaluation, the human aspect still takes time. Therefore, an unfortunate but common approach is to replace the human with a simulated user. This is unfortunate as providing explanations that

are informative to humans is a non-trivial task and often involves interdisciplinary knowledge from the human-computer interaction (HCI) and social science fields. Replacing a human with a simulated user leads to over-optimistic results.

Miller [134] provides an excellent overview of what effective explanation is from the social science perspective and criticizes current works by saying “most work in explainable artificial intelligence uses only the researchers’ intuition of what constitutes a ‘good’ explanation.”

It is therefore critical that interpretability methods are *human-grounded*. These are common strategies to measure *human-grounded*, used both in NLP and other fields:

- Humans have to choose the best model based on an explanation [41].
- Humans have to predict the model’s behavior on new data [36].
- Humans have to identify an outlier example called an intruder [145]. While it can be used in other fields, it is most common in NLP where it is used with *vocabulary* explanations [146].

**Faithfulness** evaluation checks how well the explanation reflects the model. This is also known as *functionally-groundedness* [41, 59, 142, 147] or sometimes *fidelity* [144].

It might seem surprising that an explanation that is directly produced from the model would not reflect the model. However, even intrinsically interpretable methods such as *attention* and *Neural Modular Networks* have been shown to not reflect the model [2, 96, 97, 98].

Interestingly, *human-grounded* interpretability methods cannot reflect the model perfectly because humans require explanations to be selective, meaning the explanation should select “one or two causes from a sometimes infinite number of causes” [134]. Regardless, the explanations must still reflect the model to some extent, which surprisingly is not always the case [48, 59]. Additionally, explanations that provide a similar type of explanation, with similar selectiveness, should compete to provide the explanation that best reflects the model.

For some tasks, measuring if an interpretability method is *faithful* is trivial. In the case of *adversarial examples*, it is enough to show that the prediction changed, and the adversarial example is a paraphrase. In other cases, most notably *input features*, providing a *faithfulness* metric can be very challenging [14, 45, 59, 107, 109].

In general, common evaluation strategies, both in NLP and other fields, are:

- Comparing with an intrinsically interpretable model, such as logistic regression [41].
- Comparing with other post-hoc methods [2].
- Proposing axiomatic desirables [39].
- Benchmarking against random explanations [14].

## 2.4 Methods of Interpretability

The main objective of this background chapter is to give an overview of interpretability methods and categorize them by how they communicate. Section 2.5 to Section 2.7 and Appendix A are dedicated towards this goal. Note, Appendix A exists only for completeness and is not required to understand the remainder of the thesis.

Each method section covers one communication approach, corresponding to one row in Table 2.1, and can be read somewhat independently. Each section discusses the purpose of the communication approach and covers the most relevant methods and how they are evaluated. Because interpretability is a large field, this background section chooses methods based on historical progression and diversity regarding what information they use. Finally, at the end of each method section, the general trends and issues related to that communication approach are discussed.

## 2.5 Input Features

An *Input feature* explanation is a *local explanation*, where the goal is to determine how important an *input feature*, e.g. a token, is for a given prediction. This approach is highly adaptable to different problems, as the input features are always known and are often meaningful to humans. Especially in NLP, the input features will often represent words, sub-words, or characters. Knowing which words are the most important can be a powerful explanation method. An *input feature* explanation of the input  $\mathbf{x}$ , is represented as

$$\mathbf{E}(\mathbf{x}, c) : \mathbb{I}^{T \times d} \rightarrow \mathbb{R}^T, \text{ where } \mathbb{I} \text{ is the input domain,} \quad (2.1)$$

d is the input dimensionality,

and T is the sequence length.

Note that when the output is a score of importance, the explanation is called an *importance measure*; alternatives could be a simple listing or ranking of important tokens.

Additionally, there exists a secondary categorization, which separates methods that can distinguish between positive and negative contribution (termed signed) and those that just tell if or how much something is contributing (termed absolute). In the case of importance measures, a signed importance measure can be transformed into an absolute importance measure by simply using the *abs*( $\cdot$ ) operation. However, this may not always be the case for listings or rankings of importance.

Importantly, *input feature* explanations can only explain one scalar, meaning one class at

one timestep. In a sequence-to-sequence application, the explanation is therefore repeated for each time step [1, 38] although this may not respect the combinatorial complexities [148]. Additionally, the selected class is either the most likely or true-label class. In this section, the explained class is denoted with  $c$ . For all methods in this section,  $c$  can be set as desired.

### 2.5.1 Gradient-based

The essential idea in gradient-based input feature explanations, which are typically importance measures, is that if a small change in the input affects the output a lot, then this indicates importance. Such a relationship can be estimated using the gradient with respect to the input, of which there exist many variations.

**Gradient** The simplest approach [37, 38] is the gradient with respect to the input, as defined in (2.2).

$$\mathbf{E}_{\text{gradient}}(\mathbf{x}, c) = L_p(\nabla_{\mathbf{x}} p(c|\mathbf{x}; \theta)), \text{ where } L_p \in \{L_1, L_2, L_\infty\} \quad (2.2)$$

and  $p(c|\mathbf{x}; \theta)$  is the model's probability output.

Because NLP features are often discrete, the gradient is w.r.t. the one-hot-encoding, which is done by treating it as continuous. Because the one-hot-encoding has shape  $\mathbf{x} \in \mathbb{I}^{T \times V}$ , where  $V$  is the vocabulary size, it is necessary to reduce away the vocabulary dimension (often using an  $L_p$ -norm) such that  $\mathbf{E}(\mathbf{x}, c) \in \mathbb{R}^T$ . This normalization means that (2.2) is an *absolute* importance measure.

		$p(y \mathbf{x}; \theta)$	$y$	$c$
x	the year 's <b>best</b> and most <b>unpredictable</b> comedy	0.91	pos	pos
x	<b>we</b> <b>never</b> <b>feel</b> <b>anything</b> <b>for</b> <b>these</b> <b>characters</b>	0.95	neg	neg
x	<b>handsome</b> <b>but</b> <b>unfulfilling</b> <b>suspense</b> <b>drama</b>	0.18	neg	pos

Figure 2.3 Hypothetical visualization of applying  $\mathbf{E}_{\text{gradient}}(\mathbf{x})$ , where  $c$  is the explained class. Note that because the vocabulary dimension is reduced away, typically using the  $L^2$ -norm, it is not possible to separate positive influence from negative influence.

The primary argument for (2.2) being *faithful*, is that for a linear model  $\mathbf{xW}$ , the explanation would be  $\mathbf{W}_{c,:}^\top$ , which is clearly a valid explanation [45]. However, this does not guarantee

*faithfulness* for non-linear models, although it will relate to a first-order Taylor approximation [38].

**Input times gradient** The simplest extension of (2.2) is to also consider the scale of  $\mathbf{x}$ , hence the extension  $\mathbf{x} \odot \nabla_{\mathbf{x}} p(c|\mathbf{x}; \theta)$  is sometimes preferred. Although, a counter-argument is that  $\mathbf{x}$  does not directly relate to the model, and this can therefore result in a less faithful explanation [45].

Note that because  $\mathbf{x}$  is a one-hot encoding, only one element per input word will be non-zero. Therefore, instead of using a norm to reduce away the vocabulary dimension, it's possible to just select the non-zero element. Therefore, this variation can be a *signed* importance measure.

**Integrated gradient** Parts of the input may be important but have zero gradients, for example due to the truncation in  $\text{ReLU}(\cdot)$ . In such cases, the previous gradient-based methods won't show any attribution.

Sundararajan et al. [39] call this desirable *sensitivity*. Specifically, if there exists a combination of  $\mathbf{x}$  and baseline  $\mathbf{b}$  (often an empty sequence), where the logit outputs of  $f(\mathbf{x}; \theta)$  and  $f(\mathbf{b}; \theta)$  are different, then the feature that changed should get a non-zero attribution.

Additionally, Sundararajan et al. [39] suggest the desirable *completeness*. Meaning, that the sum of importance scores assigned to each token should equal the model output relative to the baseline  $\mathbf{b}$ .

To satisfy these desirables, Sundararajan et al. [39] propose *integrated gradient* as defined in (2.3). This integrates the gradients between an uninformative baseline  $\mathbf{b}$  and the observation  $\mathbf{x}$  [39], using an approximative sampling of  $k$  steps.

$$\mathbf{E}_{\text{integrated-gradient}}(\mathbf{x}, c) = (\mathbf{x} - \mathbf{b}) \odot \frac{1}{k} \sum_{i=1}^k \nabla_{\tilde{\mathbf{x}}_i} f(\tilde{\mathbf{x}}_i; \theta)_c, \quad \tilde{\mathbf{x}}_i = \mathbf{b} + i/k(\mathbf{x} - \mathbf{b}), \quad (2.3)$$

where  $f(\mathbf{x}; \theta)$  is the model logits.

This approach has been successfully applied to NLP, where the uninformative baseline can be an empty sentence, such as padding tokens [149].

Although Integrated Gradient has become a popular approach, it has recently received criticism in computer vision (CV) community for not being *faithful* [14]. In NLP, Bastings et al. [91] use synthetic NLP tasks and conclude its *faithfulness* is task- and model-dependent. Finally, Bilodeau et al. [111] provide a theoretical framework that also says that this explanation will

always be task- and model-dependent.

### 2.5.2 Occlusion-based

With occlusion-based feature attribution explanations, the idea is to either mask or remove input tokens and then see how the model responds.

**Leave-one-out** The simplest approach is perhaps leave-one-out, which removes or masks one token at a time; the explanation is the difference in the model’s logit [40]. This is a signed importance measure because there is no need for normalization. Let  $\tilde{\mathbf{x}}_{-i}$  be the input  $\mathbf{x}$  with token at position  $i$  removed or masked. Then this explanation can be expressed as:

$$\mathbf{E}_{\text{LOO}}(\mathbf{x}, c) = [f(\mathbf{x}; \theta)_c - f(\tilde{\mathbf{x}}_{-i}; \theta)_c]_{t=1}^T \quad (2.4)$$

One concern is that leave-on-out may not be *faithful* because the input may be ungrammatical or the model is not designed to have its input masked, thus causing out-of-distribution issues [35].

**LIME** This approach samples nearby observations  $\tilde{\mathbf{x}}$  and uses the model estimate  $p(c|\tilde{\mathbf{x}})$  to fit a logistic regression. The parameters  $\mathbf{w}$  of the logistic regression then represent a signed *importance measure*.

$$\begin{aligned} \mathbf{E}_{\text{LIME}}(\mathbf{x}, c) = \underset{\mathbf{w}}{\operatorname{argmin}} \quad & \frac{1}{k} \sum_{i=1}^k (p(c|\tilde{\mathbf{x}}_i; \theta) \log(q(\tilde{\mathbf{x}}_i)) + (1 - p(c|\tilde{\mathbf{x}}_i; \theta)) \log(1 - q(\tilde{\mathbf{x}}_i)) + \lambda \|\mathbf{w}\|_1 \\ & \text{where } q(\tilde{\mathbf{x}}) = \sigma(\mathbf{w}\tilde{\mathbf{x}}) \end{aligned} \quad (2.5)$$

One complication of *LIME* is how to sample  $\tilde{\mathbf{x}}$ , representing the nearby observations. In the original paper [41], they use a Bag-Of-Words (BoW) representation with a cosine distance. While this approach remains possible with a model that works on sequential data, such distance metrics may not effectively match the model’s internal space. Recent works [150], therefore sample  $\tilde{\mathbf{x}}$  by masking words of  $\mathbf{x}$ . However, this requires a model that supports such masking.

An advantage of *LIME* is it uses a LASSO logistic regression, which is a normal logistic regression with an  $L_1$ -regularizer. This means that its explanation is selective, as in sparse, which may be essential for providing a human-friendly explanation [134].

Ribeiro et al. [41] show that LIME is *faithful* by applying LIME on *intrinsically* interpretable

		$p(y \mathbf{x}; \theta)$	$y$	$c$
x	the year 's best and most unpredictable comedy	0.91	pos	pos
x	we never feel anything for these characters	0.95	neg	neg
x	handsome but unfulfilling suspense drama	0.18	neg	pos

Figure 2.4 A fictive visualization of LIME, where the weights of the logistic regression determine the *importance measure*. Note that for LIME, it is possible to have negative importance (indicated by blue). Furthermore, some tokens have no importance score due to the  $L^1$ -regularizer.

models, such as a logistic regression model, and then compare the LIME explanation with the *intrinsic* explanation from the logistic regression. They also show *human-groundedness* by conducting a human trial experiment, where non-experts have to choose the best model, based on the provided explanation, given a “wrong classifier” trained on a biased dataset and a “correct classifier” trained on a curated dataset.

**Kernel SHAP** A limitation of *LIME* is that the weights in a linear model are not necessarily *intrinsically* interpretable. When multicollinearity exists (input features are linearly correlated with each other), the model weights can be scaled arbitrarily, resulting in misleading importance scores.

To avoid the multicollinearity issue, one approach is to compute Shapley values [151], which are derived from game theory. Shapley values not only solve the multicollinearity issue, but more broadly is an approach to assign importance to individual features even when there are co-dependencies. For example, if a model takes two features  $\{x_1, x_2\}$  but the model can perform identically correct predictions with just one feature (e.g.  $\{x_1\}$ ), most other explanations will assign their individual importance as being zero, despite the feature actually being important.

The central idea with Shapley values is to consider every permutation of features enabled. For example in the two feature case ( $\{x_1, x_2\}$ ), the Shapley values would be computed by considering the model outputs with the features  $\{\emptyset\}, \{x_1\}, \{x_2\}, \{x_1, x_2\}$ . Thus, if there are  $T$  features, this would require  $\mathcal{O}(2^T)$  permutations.

While this method works in theory, it is intractable in practice due to the exponential number of permutations. Lundberg and Lee [44] present a framework for producing Shapley values in a more tractable manner. They introduce a model-agnostic approach called *Kernel SHAP*. It combines 3 ideas: it reduces the number of features via a mapping function  $h_{\mathbf{x}}(\mathbf{z})$ , it uses squared-loss instead of cross-entropy by working on logits, and it weighs each observation by

how many features there are enabled.

$$\begin{aligned} \mathbf{E}_{\text{SHAP}}(\mathbf{x}, c) &= \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{\mathbf{z} \in \mathbb{Z}^M} \pi(\mathbf{z}) (f(h_{\mathbf{x}}(\mathbf{z}); \theta)_c - g(\mathbf{z}))^2 \\ \text{where } g(\mathbf{z}) &= \mathbf{w}\mathbf{z} \\ \pi(\mathbf{z}) &= \frac{M-1}{(M \text{ choose } |\mathbf{z}|)|\mathbf{z}|(M-|\mathbf{z}|)} \end{aligned} \quad (2.6)$$

In (2.6),  $\mathbf{z}$  is a  $\{0, 1\}^M$  vector that describes which combined features are enabled. This is then used in  $h_{\mathbf{x}}(\mathbf{z})$ , which enables those features in  $\mathbf{x}$ . Furthermore,  $\mathbb{Z}^M$  represents all permutations of enabled combined features, and  $|\mathbf{z}|$  is the number of enabled combined features. Figure 2.5, demonstrates a fictive example of how input features can be combined and visualizes their shapley values.

		$p(y \mathbf{x}; \theta)$	$y$	$c$
x	<u>the year 's</u> <u>best and most</u> <u>unpredictable comedy</u>	0.91	pos	pos
x	<u>we never feel anything</u> <u>for these characters</u>	0.95	neg	neg
x	<u>handsome</u> <u>but unfulfilling</u> <u>suspense drama</u>	0.18	neg	pos

Figure 2.5 Fictive visualization of *Kernel SHAP*. Note how input tokens are combined to a single feature to make *SHAP* more tractable to compute, this is the role of  $h_{\mathbf{x}}(\mathbf{z})$  in (2.6).

Lundberg and Lee [44] show *faithfulness* theoretically, using that Shapley values uniquely satisfy a set of axiomatic desirables (such as *sensitivity* and *completeness*, which *Integrated Gradient* also satisfy) and that because *Kernel SHAP* also satisfy these axioms, they are also Shapley values. Lundberg and Lee [44] show *human-groundedness* by asking humans to manually produce importance measures and correlate them with the *Kernel SHAP* values.

*Kernel SHAP* and Shapley values in general are heavily used in the industry [47]. This popularity is likely due to their mathematical foundation and the `shap` library, and there are many other approximations aside from *Kernel SHAP* with different tradeoffs [152]. That being said, SHAP’s linearity and completeness also mean that its *faithfulness* is theoretically proven to be model and task-dependent [111], and there will be cases where these explanations are not more faithful than a random explanation.

### 2.5.3 Attention-based

Attention is an example of an intrinsic explanation, as attention is part of the model architecture. As such, it has become a popular way of explaining models, particularly transformer

models, as they use attention. However, despite being an intrinsic explanation, the way it's interpreted is often not *faithful* [95].

**Classical attention** One of the simplest cases of attention is the BiLSTM-Attention classification model. In this case, there are two types of models: single-sequence and paired-sequence. However, they are nearly identical and only differ in how the context vector  $\mathbf{b}$  is computed [2].

In both cases, a  $d$ -dimensional word embedding followed by a bidirectional LSTM (BiLSTM) encoder is used to transform the one-hot encoding into the hidden states  $\mathbf{h}_x \in \mathbb{R}^{T \times 2d}$ . These hidden states are then aggregated using an additive attention layer  $\mathbf{h}_\alpha = \sum_{i=1}^T \alpha_i \mathbf{h}_{x,i}$ .

To compute the attention weights  $\alpha_i$  for each token:

$$\alpha_i = \frac{\exp(\mathbf{u}_i^\top \mathbf{v})}{\sum_j \exp(\mathbf{u}_j^\top \mathbf{v})}, \quad \mathbf{u}_i = \tanh(\mathbf{W}_x \mathbf{h}_{x,i} + \mathbf{b}) \quad (2.7)$$

where  $\mathbf{W}_x, \mathbf{v}$  are model parameters. Finally, the  $\mathbf{h}_\alpha$  is passed through a fully connected layer to obtain the logits  $f(\mathbf{x})$ . In the single-sequence case,  $\mathbf{b}$  is a learned parameter. While in the paired-sequence case, the second sequence is  $\mathbf{s} \in \mathbb{R}^{T_s \times V}$ , which is transformed using a separate BiLSTM encoder to get the hidden states  $\mathbf{h}_b$ . Finally,  $\mathbf{b} = \mathbf{W}_b \mathbf{h}_{b,T_s}$ , where  $\mathbf{W}_b$  is a learnable model parameters.

Because there is an attention-weight  $\alpha_i$  for each input token, which is used as a weight in the sum  $\mathbf{h}_\alpha$ , each attention-weight intrinsically explains how relevant  $\mathbf{h}_{x,i}$  is [79]. The problem is that this interpretation is then extended to how relevant the input  $x_i$  is. This is not *faithful* as the BiLSTM layer can swap, merge, or move the tokens. Therefore, there is no intrinsic way to align the relevance of  $\mathbf{h}_{x,i}$  with  $x_i$  [95]. In practice, Bastings et al. [91] show that the *faithfulness* is both task- and model-dependent.

**Transformer attention** The idea of attention has since been made popular by Vaswani et al. [136], where multiple layers and multiple “heads” of attention are used.

Each layer is composed of the multi-head self-attention mechanism, ignoring some minor details. (2.8) defines this mechanism. In (2.8),  $\mathbf{h}_\ell$  is the hidden representation for layer  $\ell$ , with  $\mathbf{h}_1 = g(\mathbf{x})$ . Each  $\alpha_{i,\ell}$  is the self-attention, of which there are multiple “heads” ( $i \in [1, H]$ ).  $\mathbf{h}_\ell = f(\alpha_{:, \ell})$  represents other neural-network components, such as normalization and activation functions, which are not relevant to this discussion on attention.

$$\mathbf{h}_\ell = f(\alpha_{1,\ell} \mathbf{W}_{K,1,\ell}, \dots, \alpha_{H,\ell} \mathbf{W}_{K,H,\ell}), \quad \alpha_{i,\ell} = \text{softmax}(\mathbf{h}_{\ell-1} \mathbf{W}_{Q,i,\ell} \mathbf{W}_{K,i,\ell}^T \mathbf{h}_{\ell-1}^T) \quad (2.8)$$

Because each row in  $\alpha_{i,\ell}$  are normalized to sum to one, and because they are used in a matrix product, they can be interpreted as importance weights similar to the classical attention case. However, this transformer attention has the additional complication of having multiple layers and multiple heads for each layer. Hence, it’s unclear which attention matrix to investigate. Additionally, each layer can move, swap, or merge the relationship between intermediate representations ( $\mathbf{h}$ ) and the input tokens  $\mathbf{x}$  [98].

To address the issue of multiple attention matrices, Abnar and Zuidema [153] take the perspective that all the attention matrices can be represented as a directional graph, and then a max-flow algorithm [154] can be used to consider the total importance for the entire model. Ethayarajh and Jurafsky [155] then show that using max-flow creates Shapley values; this gives it the same theoretical foundation as Kernel SHAP. However, it also gives it the same theoretical disadvantages, where there is no guarantee for *faithfulness* [111].

To address the issue of information mixing [156], Tutek and Snajder [157] apply regularizations such as weight tying and show that increasing the relationship between input tokens and intermediate representations is possible. However, it’s unclear how strong this relationship must be to call attention *faithful*.

#### 2.5.4 Discussion

**Groundedness** The *faithfulness* of *input feature* explanations have received a lot of attention and discussion. However, there is still little consensus on what is *functionally-grounded* or how to even measure it [2, 45, 88, 89, 91, 107, 142].

It has been suggested, that a general *functionally-grounded* post-hoc *input feature* explanation method just doesn’t exist [48], gradient-based methods have been shown to be arbitrarily manipulable [158], and many explanations methods, both post-hoc and intrinsic, have been theoretically shown to be subject to a no-free-lunch theorem [110].

**Future work** At present, it’s unclear what it would take to get *functionally-grounded input feature* explanations. Such high-level questions are likely difficult to answer without a more fundamental understanding of what the *faithfulness* desirables are. Therefore, we advocate for continuing the effort in measuring *faithfulness* but to focus more on establishing the fundamental desirables.

Additionally, the methods that had the most theoretical motivation, such as Kernel SHAP, integrated gradient, and max-flow attention, have also been shown to be theoretically questionable [111]. Indicating theoretical motivations should be taken with skepticism. We would

therefore advocate for a more empirical evidence-based approach, where it’s made clear what the *functionally-grounded* metric is before the method is proposed.

## 2.6 Counterfactuals

*Counterfactual explanations* are essentially answering the question “how would the input need to change for the prediction to be different?”. Furthermore, these *counterfactual examples* should be a minimal edit from the original example and fluent. However, all of these properties can also be said of *adversarial explanations*, and indeed some works confuse these terms. The critical difference is that *adversarial examples* should have the same gold label as the original example, while *counterfactual examples* should have a different gold label (often opposite) as the original example [42]. Because *Counterfactual explanations* are defined by the output class they are limited to sequence-to-class models.

Another common confusion is with *counterfactual datasets*, also known as *Contrast Sets*. These datasets are used in robustness research and could consist of *counterfactual examples*. However, these datasets are generated without using a model [159, 160], and can therefore not be used to explain the model. However, *Contrast Sets* are important for ensuring a robust model.

In social sciences, *counterfactual explanations* are considered highly useful for a person’s ability to understand causal connections. Miller [134] explains that “why” questions are often answered by comparing *facts* with *foils*, where the term *foils* is the social sciences term for *counterfactual examples*.

### 2.6.1 Polyjuice

*Polyjuice* by Wu et al. [150] is primarily a *counterfactual dataset* generator, and the generation is therefore detached from the model. However, by strategically filtering these generated examples such that the model’s prediction is changed the most, they condition the *counterfactual* generation on the model, thereby making a *post-hoc* explanation.

The generation is done by fine-tuning a GPT-2 model [161] on existing *counterfactual datasets* [11, 159, 160, 162, 163, 164]. For each pair of original and counterfactual example, they produce a training prompt, see (2.9) for an instantiated example of this (with structure annotated). What the conditioning code is and what is replaced in (2.9) is determined by the existing *counterfactual datasets*.

$$\begin{aligned}
\text{prompt} = & \underbrace{\text{"It is great for kids <GENERATE>"}_{\text{original sentence}} \\
& \underbrace{[\text{negation}]}_{\text{conditioning code}} \underbrace{\text{It is [BLANK] great for [BLANK]}}_{\text{masked counterfactual}} \\
& \text{<REPLACE>} \underbrace{\text{not [ANSWER] children [ANSWER] <EOS>}}_{\text{masking answers}}
\end{aligned} \tag{2.9}$$

For *counterfactual* generation, they specify the original sentence and optionally the condition code and then let the model generate the *counterfactuals*. These *counterfactuals* are independent of the model. To make them dependent on the model, they filter the *counterfactuals* and select those examples that change the prediction the most. One important detail is that they adjust the prediction change with an *importance measure* (*SHAP*), such that the *counterfactual examples* that could have been generated by an *importance measure* are valued less. An example of this explanation can be seen in Figure 2.6.

	<b>x</b>	$p(y \mathbf{x}; \theta)$	<b>y</b>
x	<u>the year 's best and most unpredictable comedy</u>	0.91	pos
$\tilde{x}$	<u>the year 's worst and least unpredictable comedy</u>	0.11	-
x	<u>we never feel anything for these characters</u>	0.95	neg
$\tilde{x}$	<u>we feel everything for these characters</u>	0.02	-

Figure 2.6 Hypothetical results of *Polyjuice*, showing how some words were either replaced or removed to produce *counterfactual examples*.

To validate *Polyjuice*, for a *human-grounded* experiment, they show that humans were unable to predict the model’s behavior for the *counterfactual examples*, thereby concluding that their method highlights potential robustness issues. Whether *Polyjuice* is *functionally-grounded* is somewhat questionable, because the model is not a part of the generation process itself, it is merely used as a filtering step.

### 2.6.2 MiCE

Like *Polyjuice* [150], *MiCE* [42] also uses an auxiliary model to generate *counterfactuals*. However, unlike *Polyjuice*, *MiCE* does not depend on auxiliary datasets and the counterfactual generation is more tied to the model being explained, rather than just using the model’s predictions to filter the *counterfactual examples*.

The counterfactual generator is a T5 model [34], a sequence-to-sequence model, which is fine-tuned by input-output-pairs, where the input consists of the gold label and the masked sentence, while the output is the masking answer, see (2.10) for an example.

$$\begin{aligned}
 \text{input} &= \text{"label: } \underbrace{\text{positive}}_{\text{gold label}}, \text{ input: } \underbrace{\text{This movie is [BLANK]!}}_{\text{masked sentence}}\text{"} \\
 \text{target} &= \text{"[CLR] } \underbrace{\text{really great}}_{\text{masking answer}} \text{ [EOS]"}
 \end{aligned} \tag{2.10}$$

The *MiCE* approach to selecting which tokens to mask is to use an *importance measure*, specifically *the gradient w.r.t. the input*, and then mask the top  $x\%$  most important consecutive tokens.

For generating counterfactuals, *MiCE* again masks tokens based on the *importance measure*, but then also inverts the gold label used for the T5-input (2.10). This way, the model will attempt to infill the mask so that the sentence will have an opposite semantic meaning. This process is then repeated via a beam-search algorithm, which stops when the model prediction changes; an example of this can be seen in Figure 2.7.

	$\mathbf{x}$	$p(y \mathbf{x};\theta)$	$y$
$\mathbf{x}$	the year 's <b>best</b> and most unpredictable comedy	0.91	pos
	the year 's <b>worst</b> and most unpredictable comedy	0.59	-
$\tilde{\mathbf{x}}$	the year 's worst and most predictable comedy	0.04	-
$\mathbf{x}$	we <b>never</b> feel anything for these characters	0.95	neg
	we <b>can</b> feel anything for these characters	0.73	-
$\tilde{\mathbf{x}}$	we can feel anything for these animals	0.01	-

Figure 2.7 Hypothetical visualization of how *MiCE* progressively creates a counterfactual  $\tilde{\mathbf{x}}$  from an original sentence  $\mathbf{x}$ . The highlight shows the *gradient*  $\nabla_{\mathbf{x}}f(\mathbf{x};\theta)_y$ , which *MiCE* uses to know what tokens to replace.

Because *MiCE* uses the model prediction to stop the beam search, it will inherently be somewhat *functionally-grounded*. However, it may be that using the *gradient* as the *importance measure*, is not *functionally-grounded*. Ross et al. [42] validate that using the *gradient* is *functionally-grounded*, by looking at the number of edits and fluency of *MiCE* and compare it to a version of *MiCE* where random tokens are masked. They find that using the *gradient* significantly improves both fluency and reduces the number of edits it takes to change a prediction.

### 2.6.3 Discussion

**Groundedness** While *counterfactual examples* are great for *human-grounded* explanation, they struggle with *faithfulness*. The challenge comes from the desirables. On one side, it is desirable to provide a counterfactual example with the opposite gold label, an objective that is independent of the model. Simultaneously, the search procedure should be directed by the model’s behavior. These objectives can, at times, appear opposite, although *MiCE* provides a great example of how it can be done.

**Future work** Because the motivation for *counterfactual examples* is often robustness, the search procedure often becomes only weakly dependent on the model such as *Polyjuice* or sometimes completely independent such as *Contrast Sets*.

While robustness is a perfectly valid research objective, we recommend being careful when using both robustness and interpretability to motivate the same method, as this often leads to *faithfulness* issues. We would, therefore, advocate for more counterfactual research, which focuses only on interpretability and *faithfulness*.

## 2.7 Natural Language

A common concern for many of the explanation methods presented in this thesis is that they are difficult for people without specialized knowledge to understand. It is, therefore, attractive to directly generate an explanation in the form of *natural language*, which can be understood by simply reading the explanation for a given example. Because these utterances explain just a single example, they are *local explanations*.

Historically, research in the area of *natural language* explanation uses the explanations to improve the model’s predictive performance. The idea is that by enforcing the model to reason about its behavior, the model can generalize better [36, 43, 165, 166, 167, 168].

However, more recently, interpretability has also become a motivation for generating these explanations. The work in this field can be categorized as *intrinsic* and *post-hoc*, like the other communication approaches. However, in this field, it is often referred to as explain-then-predict and predict-then-explain, respectively. The terms are somewhat self-explanatory, where explain-then-predict means an explanation is produced first, then used to inform the prediction, and vice versa for predict-then-explain.

Additionally, explain-then-predict is also referred to as *reasoning*, which goes back to the original motivation for *natural language* explanation and has with the popularity of instruction-

tuned models become popular as chain-of-thought explanations.

Predict-then-explain is referred to as *rationalization*, in the sense that they attempt to explain after a prediction has been made [36]. Note that the term is a misnomer, as rationalizations in the dictionary sense<sup>1</sup> can also be false, which is not a desired outcome in this case.

### 2.7.1 Explain-then-predict

Because of the original “enforcing the model to reason about its behavior” motivation, most of the work in *natural language* explanation is on explain-then-predict [36, 43, 165, 166, 167, 168].

There is also a further sub-categorization of this approach [169]. In particular, there are extractive methods where the explanation is extracted from a corpus based on relevance. For example, in fact-checking, explanations of reasoning can be extracted from abstracts of scientific papers. These are then later used to perform the prediction [170]. The alternative is an abstractive explanation, where the explanations are generated using a model. These models are often trained using a human-annotated corpus, such as e-SNLI [166]; a simple example of this is NILE [43].

**NILE** The method presented by Kumar and Talukdar [43] is to consider an NLI task (e.g., the SNLI dataset), and then for each possible label (entailment, contradiction, and neutral) generate an explanation that supports that label. As such, each observation will have 3 explanations. The explanations are generated by a GPT-2 model [161]. The soundness of an explanation is then estimated using a RoBERTa model [4], where the input and an explanation are provided as the input. This provides 3 soundness scores, which are then used to provide a final prediction. To provide just one explanation, only the explanation for the final prediction is considered, and because the prediction is conditioned on this explanation, it’s an *intrinsic* explanation.

Kumar and Talukdar [43] show that their approach doesn’t degrade the performance accuracy. They measure *human-groundedness* using human annotators on 100 observations and show NILE outperforms *post-hoc* (predict-then-explain) methods. However, this result is likely biased since all the generations are based on the same dataset (e-SNLI). To show *faithfulness*, they measure comprehensiveness (what happens when the explanation is removed) and sufficiency (what happens when only the explanation remains) [171]. They find that these removals drastically decrease the accuracy, indicating that the model does use the explanation.

---

<sup>1</sup>“the action of attempting to explain or justify behavior or an attitude with logical reasons, even if these are not appropriate.” – Oxford definition of *rationalization*.

However, it's unclear if this is just caused because the model is not trained on this input-distribution [14].

## 2.7.2 Predict-then-explain

There isn't as much work in predict-then-explain as explain-then-predict. Additionally, the methods that do exist, are usually abstractive [169]. Meaning they provide a very high-level explanation without explicitly referring to specifics of the context or the question. In this section, we consider CAGE [36] as a simple example of an abstractive predict-then-explain method.

**Rationalizing Commonsense Auto-Generated Explanations (CAGE)** Rajani et al. [36] provides explanations for the Common Sense Question Answering (CQA) dataset, which is a multiple-choice question-answering dataset [25]. The explanations are independent of the model and are provided via Amazon Mechanical Turk. They then fine-tune a GPT model [172], using the question, answers, and explanation to provide rationalization explanations. See (2.11) for an example of the exact prompt construction. To clarify, this GPT model is not the explained model but provides the explanations, which is known as an explainer model.

$$\begin{aligned}
 \textit{input} = & \text{“} \underbrace{\text{What could people do that involves talking?}}_{\text{question}} \underbrace{\text{confession}}_{\text{choice 1}}, \underbrace{\text{carnival}}_{\text{choice 2}} \\
 & \text{, or } \underbrace{\text{state park?}}_{\text{choice 3}} \underbrace{\text{confession}}_{\text{answer}} \text{ because ”} \\
 \textit{target} = & \text{“} \underbrace{\text{confession is the only vocal action.}}_{\text{rational explanation}} \text{”}
 \end{aligned} \tag{2.11}$$

For simpler tasks, such as “Stanford Sentiment Treebank” [173], the prompt could simply be “[input]. [answer] because [explanation]”, see Figure 2.8 for hypothetical explanations using such a setup. Because *CAGE* uses a generative model, where [answer] can be a sequence, it is not limited to sequence-to-class problems.

Rajani et al. [36] find that rationalization explanations provide nearly identical explanations as reasoning explanations. The method is validated to be *human-grounded*, by tasking humans to use the explanation to predict the model behavior, again they find identical performance. It is questionable if *CAGE* is *faithful*, as its only connection to the explained model is during inference, where the **answer** is produced by the explained model. Because there are no other connections to the explained model, there is little reason to think the GPT explainer model can reflect the model's behavior. If the humans who provided explanations had specialist

	$\mathbf{x}$	$p(y \mathbf{x};\theta)$	$y$
$\mathbf{x}$	<u>the year 's best and most unpredictable comedy</u>	0.91	pos
	<i>unpredictable comedies are funny</i>	-	-
$\mathbf{x}$	<u>we never feel anything for these characters</u>	0.95	neg
	<i>it is important to feel for characters</i>	-	-

Figure 2.8 Hypothetical explanations from using *CAGE* to produce rationalizations for the prediction.

insight into the model, then an argument could be made for *CAGE* to be *faithful*. However, as the humans were Mechanical Turk workers, this is unlikely.

### 2.7.3 Discussion

**Groundedness** Early work in *natural language* explanations have received criticism in NLP for not evaluating *faithfulness* [174]. This issue is even more problematic because the annotated explanations are provided by humans who have no insights into the model’s behavior [175]. The explanation model, therefore, just learns about humans’ thought processes rather than the model’s logical process. This issue is somewhat unique to the NLP literature and is better treated in other fields [176].

More recently, the work in *faithfulness* has increased, with many proposed metrics [177, 178, 179, 180, 181]. However, these are all collectively criticized for just measuring some variation of consistency and still not measuring faithfulness [182]. Parcalabescu and Frank [182] claims that there presently aren’t any general-purpose faithfulness metrics for *natural language* explanations.

**Future work** Most work on natural *natural language* explanations uses *intrinsic* methods, under the motivation that forcing the model to “reason about itself” will make it more accurate. Unfortunately, this hypothesis has received criticism because the little *post-hoc* work there exists shows that this is not the case. Additionally, there are theoretical arguments for why this would not be the case [183]. Therefore, investing more effort into the *post-hoc* baseline would be natural.

The field should also move away from human-written explanations. There is no reason to think these will ever reflect the model, as they are produced just by annotating a dataset; at no point is a model even involved [166, 184]. Finally, the field must establish some solid metric of *faithfulness* [182].

### CHAPTER 3 GENERAL-PURPOSE FAITHFULNESS METRIC FOR IMPORTANCE MEASURES

A major challenge in the field of interpretability is ensuring that an explanation is *faithful*, where “a faithful interpretation is one that accurately represents the reasoning process behind the model’s prediction” [59]. For example, as discussed in Section 2.5, *importance measures* that are claimed to have strong theoretical foundations and are widely used in practice [47] often later turn out to be unfaithful [2, 14, 45, 107, 142]. To prevent such issues, accurately measuring if an explanation is faithful is paramount.

Unfortunately, because models are too complex to know what the correct explanation is, there is no obvious measure of faithfulness [59] and it is necessary to use proxies. Doshi-Velez and Kim [58] say a *faithfulness* metric should use “some formal definition of interpretability as a proxy for explanation quality.”

In Section 3.1, we argue that current metrics have fundamental issues. In general, the cost of proxies has been some combination of incorrect assumptions, expensive computations, or using a proxy-model, like in Chapter 3 or by Jain and Wallace [2], Bastings et al. [91]. Based on previous work, we propose the following desirables:

- a) The method does not assume a known true explanation.
- b) The method measures faithfulness of an explanation w.r.t. a specific model instance and single observation. For example, it is not a proxy-model that is measured.
- c) The method uses only the original dataset, e.g. does not introduce spurious correlations.
- d) The method only uses inputs and intermediate representations that are in-distribution w.r.t. the model.
- e) The method is computationally cheap by not training/fine-tuning repeatedly and only computes explanations of the test dataset.
- f) The method can be applied to any classification task.
- g) The method can be applied to any importance measure.

In this chapter we develop *Recursive ROAR*, which satisfies (a), (c), (d), (f), and (g). In Chapter 4, we then take the insights from this chapter to develop and motivate a method that satisfies all desirables. In particular, (b) and (e) are what enable the *faithfulness measurable model* (FMM) paradigm, where it’s possible to optimize an explanation towards optimal faithfulness.

In this chapter and the next, we use the *erasure-metric* definition of *faithfulness* by Samek et al. [60] and Hooker et al. [14]: if information (input tokens) is truly important, then removing it should result in worse model performance compared to removing random information (tokens). We build upon the ROAR metric by Hooker et al. [14], which adds that it is necessary to retrain the model after information is removed, to avoid out-of-distribution issues. Finally, the model performance is compared with removing random information.

A limitation of ROAR is that it is theoretically impossible to measure the faithfulness of an *importance measure* when dataset redundancies exist. For example, if two tokens are equally relevant but only one of them is identified as important, ROAR fails to remove the second token.

We propose *Recursive ROAR*, which solves this limitation. In addition to the *Recursive ROAR* metric, we introduce a summarizing metric (RACU), which aggregates the results into a scalar metric. We hope that such a metric will make comparing importance measures across papers more feasible.

Using the proposed faithfulness metrics, we perform a comprehensive comparative study of 4 different *importance measures* and two popular architectures: BiLSTM-Attention and RoBERTa [4]. We use 8 different datasets that are commonly used in the faithfulness of *attention* literature [2].

**To summarize, the contributions of this chapter are:**

- Develop the faithfulness metric *Recursive ROAR*.
- Propose a summarizing scalar metric.
- Use these metrics to perform a comprehensive comparative study of current *importance measures*.

Our comparative study reveals that no *importance measure* is consistently better than others. Instead, we find that faithfulness is both task and model-dependent. This is valuable knowledge, as although each *importance measure* might be equal in faithfulness, they are not equal in computational requirements or human-groundedness.

In particular, we find that *attention* generally provides more sparse explanations than *gradient* or *integrated gradient*. Although their faithfulness may be the same, a sparser explanation is often easier for humans to understand [134].

Computationally speaking, *integrated gradient* is approximately 50 times more expensive than the *gradient* method. This additional complexity is usually justified by being considered more faithful than *gradient*. However, our results indicate that this is rarely a worthwhile

trade-off.

### 3.1 Existing faithfulness metrics

Much recent work in NLP has been devoted to investigating the faithfulness of *importance measures*, particularly *attention*. This section categorizes these faithfulness metrics according to their underlying principle and discusses their limitations. The limitations are annotated as (a) to (g) and refer to the desirables mentioned in this chapter’s introduction.

The works on *attention* are all based on the BiLSTM-Attention models and datasets from Jain and Wallace [2], they are therefore highly comparable. We use the same models and datasets, while also analyzing RoBERTa.

#### 3.1.1 Correlating importance measures

One early idea was to compare two importance measures. The claim is that a correlation would be a very unlikely coincidence unless both explanations are faithful. Jain and Wallace [2] specifically compare attention, the *gradient* method, and the *leave-one-out* method. Meister et al. [90] repeat this experiment in a broader context.

Both Jain and Wallace [2] and Meister et al. [90] find little correlation between attention, gradient, and leave-one-out; the explanations are therefore not faithful. Jain and Wallace [2] do acknowledge the limitations of their approach, as it assumes each importance measure is faithful to begin with (a). Therefore, a lack of correlation does not inform about unfaithfulness; rather, it just indicates a mathematical relation.

#### 3.1.2 Mutate model to deceive

Jain and Wallace [2] propose that if there exist alternative attention weights that produce the same prediction, *attention* is unfaithful.

They implement this idea by directly mutating the attention (g) such that there is no prediction change but a large change in *attention* and find that alternative attention distributions exist. Vashishth et al. [89] and Meister et al. [90] apply a similar method and achieve similar results.

Wiegrefe and Pinter [142] find this analysis problematic because the attention distribution is changed directly, thereby creating an out-of-distribution issue (d). This means that the new attention distribution may be impossible to obtain naturally from just changing the input, and it therefore says little about the faithfulness of attention.

### 3.1.3 Optimize model to deceive

Because the *mutate attention to deceive* approach has been criticized for using direct mutation, an alternative idea is to learn an adversarial *attention* (**g**).

Wiegrefe and Pinter [142] investigate maximizing the KL-divergence between normal attention and adversarial attention while minimizing the prediction difference between the two models. By varying the allowed prediction difference over several runs (**e**), they show that it is not possible to significantly change the attention weights without affecting performance. Importantly, Wiegrefe and Pinter [142] only use this experiment to invalidate the *mutate attention to deceive* experiments, not to measure faithfulness. However, Meister et al. [90] do use this experiment setup as a faithfulness metric.

Pruthi et al. [185] perform a similar analysis but report a contradictory finding. They find it is possible to significantly change the attention weights without affecting performance. They use this to show that attention is not faithful.

We find this approach problematic because by changing the optimization criteria, the analysis is no longer about the standard BiLSTM-attention model [2], which is the subject of interest (**b**). Therefore, this analysis only works as a criticism of the *mutate attention to deceive* approach, not as an evaluation of faithfulness.

### 3.1.4 Known explanations in synthetic tasks

Arras et al. [186] construct a purely synthetic task, where the true explanation is known, therefore the correlation can be applied appropriately. Unfortunately, this approach cannot be used on real datasets (**f**). Instead, Bastings et al. [91] introduce spurious correlations into real datasets, creating partially synthetic tasks. They then evaluate if importance measures can detect these correlations. It is assumed that if an explanation fails this test, it is generally unfaithful. Bastings et al. [91] conclude that faithfulness is both model and task-dependent.

Both methods are valid when measuring faithfulness on models trained on (partially) synthetic data. However, the model and task-dependent conclusion also mean that we can't generalize the faithfulness findings to the models (**b**) and datasets of interest (**c**), thus limiting the applicability of this approach.

### 3.1.5 Similar inputs, similar explanation

Jacovi and Goldberg [59] suggest that if similar inputs show similar explanations, then the explanation method is faithful. Zaman and Belinkov [187] apply this idea using a

multilingual dataset ( $\mathbf{f}$ ), where each language example is explained using an importance measure, and the importance is then aligned to the English example’s importance measure using a known alignment mapping. If the correlation between language pairs is high, this indicates faithfulness.

Besides being limited to multilingual datasets, the metric assumes the model behaves similarly among languages. However, languages may have different linguistic properties or spurious correlations. A faithful explanation would then yield different explanations for each language.

### 3.2 ROAR: RemOve And Retrain

To address the shortcomings of the current faithfulness measures as described in Section 3.1, we base our metric on ROAR [14].

ROAR has been used in computer vision to evaluate the faithfulness of *importance measures* and to a limited extent in NLP [188]. The central idea is that if information is truly important, then removing it from the dataset and retraining a model on this reduced dataset should worsen model performance. This can then be compared with an uninformative baseline, where information is removed randomly.

For example, at a step size of 10%, one can remove the top- $\{10\%, 20\%, \dots 90\%\}$  allegedly important tokens, evaluate the model performance, and compare this with removing  $\{10\%, 20\%, \dots 90\%\}$  random tokens. If the *importance measure* is faithful, the former should result in a worse model performance than the latter.

This section covers how ROAR is adapted to an NLP context. Furthermore, we explain the dataset redundancy issue, which our proposed Recursive ROAR metric solves. Finally, we show that Recursive ROAR is an improvement on ROAR using a synthetic task.

#### 3.2.1 Adaptation to NLP

ROAR was originally proposed as a faithfulness metric in computer vision. In this context, pixels measured to be important are “removed” by replacing them with an uninformative value, such as a gray pixel [14].

In this work, ROAR is applied to sequence classification tasks. Because these models use tokens, the uninformative value is a special [MASK] token (example in Figure 3.1). We choose a [MASK] token rather than removing the token to keep the sequence length, which is an information source unrelated to *importance measures*.

---

0%	The	movie	is	great	.	I	really	liked	it	.
10%	The	movie	is	[MASK]	.	I	really	liked	it	.
20%	The	[MASK]	is	[MASK]	.	I	really	liked	it	.

---

Figure 3.1 Example of **ROAR**. The first sentence shows the importance of various tokens. The next two sentences demonstrate the proportion of important tokens replaced by [MASK]. Note, the second sentence is enough to infer the sentiment.

### 3.2.2 Recursive ROAR

With ROAR, there are two conclusions: either 1) the *importance measure* is to some degree faithful, or 2) the faithfulness is unknown. The former is observed when the model’s performance is statistically significantly below the random baseline. In the latter case, Hooker et al. [14] explain that the *importance measure* can either be not faithful or there can be a dataset redundancy. Recursive ROAR solves this redundancy issue and thereby provides a more informative conclusion.

Dataset redundancies affect the conclusion because the model does not need to use the redundant information. A faithful importance measure would therefore not highlight redundancies as important. After the important information that the importance measure highlighted is removed and the model is retrained, the redundant information can still keep the model’s performance high. An example of this issue is demonstrated in Figure 3.1.

We solve this issue by recursively recomputing the *importance measure* at each iteration of information removal. This way, if the *importance measure* is faithful, it would quickly mark the redundant information as important, after which it would be removed. Note that tokens that have already been masked are kept masked. We call this Recursive ROAR and provide an example in Figure 3.2.

---

0%	The	movie	is	great	.	I	really	liked	it	.
10%	The	movie	is	[MASK]	.	I	really	liked	it	.
20%	The	movie	is	[MASK]	.	I	really	[MASK]	it	.

---

Figure 3.2 Example of how a redundancy can be removed in **Recursive ROAR** by reevaluating the *importance measure*. Compare this to Figure 3.1, where redundancies are not removed and the performance can remain the same, even when the *importance measure* is faithful.

In the example of Figure 3.2, each iteration masks one more token. However, because

ROAR requires retraining the model, for every evaluation step, this is infeasible. Instead, we approximate it by removing a relative number of tokens. The disadvantage of this approximation is that Recursive ROAR might not remove all redundancies unless the step size is one token. We discuss this more in Appendix B.

### 3.2.3 Validation on a synthetic problem

To show that Recursive ROAR provides an optimal faithfulness metric, we validate it on the same generated synthetic problem (with input  $\mathbf{x}$  and output  $y$ ) presented in the original ROAR paper [14]:

$$\mathbf{x} = \frac{\mathbf{a}z}{10} + \mathbf{d}\eta + \frac{\epsilon}{10}, \quad y = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}. \quad (3.1)$$

Quoting Hooker et al. [14] “All random variables were sampled from a standard normal distribution. The vectors  $\mathbf{a}$  and  $\mathbf{d}$  are 16-dimensional vectors that were sampled once to generate the dataset. In  $\mathbf{a}$  only the first 4 values have nonzero values to ensure that there are exactly 4 informative features. The values  $z$ ,  $\eta$ , and  $\epsilon$  are sampled independently for each example.”

The ground truth removal order is to remove the first 4 features (the specific order does not matter) followed by the remaining irrelevant features. Note that these first 4 features are mutually redundant. For comparison, we also include the *worst case*, where the first 4 features are removed last.

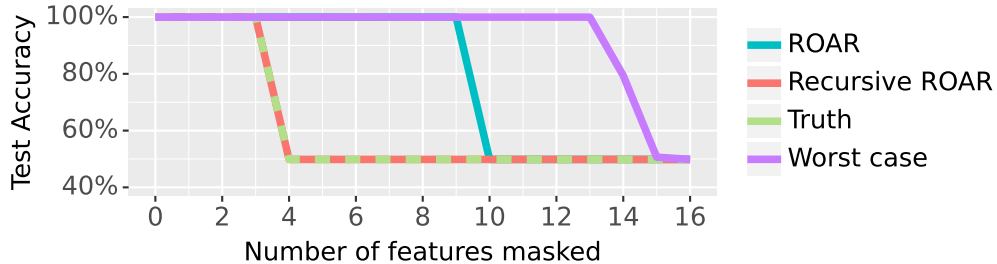


Figure 3.3 Using the weights of a linear model as the explanation, ROAR and Recursive ROAR are applied to the problem described in (3.1). In addition, the ground truth and worst case are shown. Recursive ROAR and the ground truth are identical. Note that the worst case does not lose performance at “13 features removed” since there are redundancies.

In Hooker et al. [14], they do not use a specific importance measure. Instead, they use predefined removal orders. This avoids the redundancy issue in the synthetic task, although

they do mention it as a limitation. Instead, we use the weights of a linear model, i.e. a logistic regression model that learns  $y$  given  $\mathbf{x}$ , as the importance measure and apply ROAR and Recursive ROAR using this explanation.

Figure 3.3 shows that Recursive ROAR is identical to the ground truth, while ROAR is worse, meaning that ROAR falls between the ground truth and the worst case.

### 3.3 Models

In the experimental Section 3.4, both BiLSTM-attention models and masked language models (RoBERTa) are used. In this section, we describe these models, with particular emphasis on the BiLSTM-Attention models by Jain and Wallace [2], as this is not an off-the-shelf model. Although, it is quite standard in the literature.

#### 3.3.1 BiLSTM-Attention

The BiLSTM-Attention models, hyperparameters, and pre-trained word embeddings are the same as those from Jain and Wallace [2]. We repeat the configuration details in Table 3.1, as they have not been clearly documented in prior work.

There are two types of models, single-sequence and paired-sequence, however, they are nearly identical. They only differ in how the context vector  $\mathbf{b}$  is computed. In general, we refer to  $\mathbf{x} \in \mathbb{R}^{T \times V}$  as the one-hot encoding of the primary input sequence, of length  $T$  and vocabulary size  $V$ . The logits are then  $f(\mathbf{x})$  and the target class is denoted as  $c$ .

#### Single-sequence

A  $d$ -dimensional word embedding followed by a bidirectional LSTM (BiLSTM) encoder is used to transform the one-hot encoding into the hidden states  $\mathbf{h}_x \in \mathbb{R}^{T \times 2d}$ . These hidden states are then aggregated using an additive attention layer  $\mathbf{h}_\alpha = \sum_{i=1}^T \alpha_i \mathbf{h}_{x,i}$ .

The attention weights  $\alpha_i$  for each token is then computed using:

$$\alpha_i = \frac{\exp(\mathbf{u}_i^\top \mathbf{v})}{\sum_j \exp(\mathbf{u}_j^\top \mathbf{v})}, \quad u_i = \tanh(\mathbf{W}\mathbf{h}_{x,i} + \mathbf{b}) \quad (3.2)$$

where  $\mathbf{W}, \mathbf{b}, \mathbf{v}$  are model parameters. Finally, the  $\mathbf{h}_\alpha$  is passed through a fully-connected layer to obtain the logits  $f(\mathbf{x})$ .

Dataset	Variant	Embedding initialization	Embedding size	nb. of parameters	Batch size	Max epochs
Anemia	Singe	Word2Vec trained on MIMIC	300	5 352 158	32	8
Diabetes	Single	Word2Vec trained on MIMIC	300	6 138 158	32	8
IMDB	Single	Pretrained Fast-Text	300	4 218 458	32	8
SNLI	Paired	Pretrained Glove (840B)	300	13 601 939	128	25
SST	Single	Pretrained Fast-Text	300	4 603 658	32	8
bAbI-1	Paired	Standard Normal Distribution	50	55 048	50	100
bAbI-2	Paired	Standard Normal Distribution	50	55 048	50	100
bAbI-3	Paired	Standard Normal Distribution	50	55 048	50	100

Table 3.1 Details on the BiLSTM-attention models’ hyperparameters. Everything is exactly as done by Jain and Wallace [2]. For all datasets, ASMGrad Adam [3] is used with default hyperparameters ( $\lambda = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ ) and a weight decay of  $10^{-5}$ .

### Paired-sequence

For paired-sequence problems, the two sequences are denoted as  $\mathbf{x} \in \mathbb{R}^{T_x \times V}$  and  $\mathbf{b} \in \mathbb{R}^{T_b \times V}$ . The inputs are then transformed to embeddings using the same embedding matrix, and then transformed using two separate BiLSTM encoders to get the hidden states,  $\mathbf{h}_x$  and  $\mathbf{h}_b$ . Likewise, they are aggregated using additive attention  $\mathbf{h}_\alpha = \sum_{i=1}^{T_x} \alpha_i \mathbf{h}_{x,i}$ .

The attention weights  $\alpha_i$  are computed as:

$$\alpha_i = \frac{\exp(\mathbf{u}_i^\top \mathbf{v})}{\sum_j \exp(\mathbf{u}_j^\top \mathbf{v})} \quad (3.3)$$

$$\mathbf{u}_i = \tanh(\mathbf{W}_x \mathbf{h}_{x,i} + \mathbf{W}_b \mathbf{h}_{b,T_b}),$$

where  $\mathbf{W}_x, \mathbf{W}_b, \mathbf{v}$  are model parameters. Finally,  $\mathbf{h}_\alpha$  is transformed with a dense layer.

### 3.3.2 RoBERTa

We use RoBERTa [4] as a transformer-based model due to its consistent convergence. Consistent convergence is helpful as ROAR and Recursive ROAR require the model to be trained

many times. We use the RoBERTa-base pre-trained model and only perform fine-tuning. The hyperparameters are those defined in Liu et al. [4, Appendix C] on GLUE tasks. We list the additional hyperparameters in Table 3.2.

Dataset	Variant	Max epochs
Anemia	Single	3
Diabetes	Single	3
IMDB	Single	3
SNLI	Paired	3
SST	Single	3
bAbI-1	Paired	8
bAbI-2	Paired	8
bAbI-3	Paired	8

Table 3.2 Details on the RoBERTa models’ hyperparameters. RoBERTa [4] is fine-tuned using the RoBERTa-base pre-trained model from HuggingFace [5] (125M parameters). The hyperparameters are those used by Liu et al. [4] on GLUE tasks [4, Appendix C]. The optimizer is AdamW [6], the learning rate has linear decay with a warmup ratio of 0.06, and there is a weight decay of 0.01. Additionally, we use a batch size of 16 and a learning rate of  $2 \cdot 10^{-5}$ .

RoBERTa makes use of a beginning-of-sequence [CLS] token, a end-of-sequence [EOS] token, a separation token [SEP] token, and a masking token [MASK] token. The masking token used during pre-training is the same token that we use for masking allegedly important tokens.

For the single-sequence tasks, we encode as [CLS] ... *sentence* ... [EOS]. For the paired-sequence tasks, we encode as [CLS] ... *main sentence* ... [SEP] ... *auxiliary sentence* ... [EOS]. Note that only the main sentence is considered when computing the importance measures. This is to be consistent with the BiLSTM-attention model.

### 3.4 Experiments

The datasets, performance metrics, and the BiLSTM-attention model are identical to those used in Jain and Wallace [2] and most other literature evaluating the faithfulness of *attention*. In addition, we use the RoBERTa-base model with the standard fine-tuning procedure [4]. Code is available at <https://github.com/AndreasMadsen/nlp-roar-interpretability>.

All datasets are public works. There have been no attempts to identify any individuals. The use is consistent with their intended use, and all tasks were already established by Jain and Wallace [2]. The MIMIC-III dataset [9] is an anonymized dataset of health records. I obtained

a HIPAA certification to access this. Additionally, the MIMIC-III data has not been shared with anyone else, including supervisors.

Below, we provide more details on each dataset. In Table 3.3, we provide dataset statistics.

Dataset	Size			Performance [%]		
	Train	Validation	Test	LSTM by Jain and Wallace [2]	LSTM	RoBERTa
Anemia	4262	729	1242	92	$88^{+1.1}_{-2.2}$	$86^{+0.6}_{-0.7}$
Diabetes	8066	1573	1729	79	$81^{+2.2}_{-2.9}$	$76^{+0.7}_{-0.6}$
IMDB	17212	4304	4362	78	$90^{+0.4}_{-0.7}$	$95^{+0.2}_{-0.2}$
SNLI	549367	9842	9824	88	$78^{+0.2}_{-0.3}$	$91^{+0.1}_{-0.1}$
SST	6579	848	1776	81	$82^{+0.6}_{-1.0}$	$94^{+0.3}_{-0.3}$
bAbI-1	8500	1500	1000	100	$100^{+0.0}_{-0.1}$	$100^{+0.0}_{-0.0}$
bAbI-2	8500	1500	1000	48	$68^{+9.1}_{-19.1}$	$100^{+0.1}_{-0.1}$
bAbI-3	8500	1500	1000	62	$60^{+6.5}_{-4.9}$	$81^{+6.8}_{-20.0}$

Table 3.3 Datasets statistics for single-sequence and paired-sequence tasks. Following Jain and Wallace [2], we use the same BiLSTM-attention model and report performance as macro-F1 for SST, IMDB, Anemia, and Diabetes; micro-F1 for SNLI; and accuracy for bAbI. The 95% confidence interval is reported in the subscript and subscript.

### Single-sequence tasks

1. *Stanford Sentiment Treebank (SST)* [13] – Sentences are classified as positive or negative. The original dataset has 5 classes. Following Jain and Wallace [2], we label (1,2) as negative, (4,5) as positive, and ignore the neutral sentences.
2. *IMDB Movie Reviews* [26] – Movie reviews are classified as positive or negative.
3. *MIMIC (Diabetes)* [9] – Uses health records to detect if a patient has Diabetes.
4. *MIMIC (Chronic vs Acute Anemia)* [9] – Uses health records to detect whether a patient has chronic or acute anemia.

### Paired-sequence tasks

5. *Stanford Natural Language Inference (SNLI)* [30] – Inputs are premise and hypothesis. The hypothesis either entails, contradicts, or is neutral w.r.t. the premise.
6. *bAbI* [10] – A set of artificial text for understanding and reasoning. We use the first three tasks, which consist of questions answerable using one, two, and three sentences from a passage, respectively.

### 3.4.1 Recursive ROAR

To evaluate the faithfulness of importance measures, we apply *Recursive ROAR* to all datasets and both models, as described in Section 3.2.2, after the initial fine-tuning. The results are presented in Figure 3.4 and discussed in Section 3.5.

In Appendix B, we report the compute times. Because BiLSTM-Attention is a small model and RoBERTa-base is only fine-tuned, Recursive ROAR is feasible when *importance measure* can be evaluated on every observation. For some *importance measures*, like SHAP [44], which have exponential compute complexity, ROAR would not be feasible. Additionally, for larger language models, like T5 [34], ROAR would also be difficult to apply, as fine-tuning these models is generally challenging.

#### How to interpret

If the model performance of a given *importance measure* is below the random baseline, then this indicates the importance measure is faithful. Note that “faithful” is not absolute, rather we measure the degree of faithfulness. However, if the model performance is not statistically significant below the random baseline, then the *importance measure* is not considered to be faithful. With the (*Not Recursive*) *ROAR* measure, this latter case would be inconclusive as the faithfulness could be hidden by dataset redundancies.

Figure 3.4 also presents the model performance at 100% masking, which provides a lower bound for the model performance and is helpful as the datasets are often biased. These biases come from unbalanced classes or the secondary sequence for the paired-sequence tasks [189]. It could also come from a sequence-length. However, for these datasets sequence-length bias is not a concern (see Section 3.4.3).

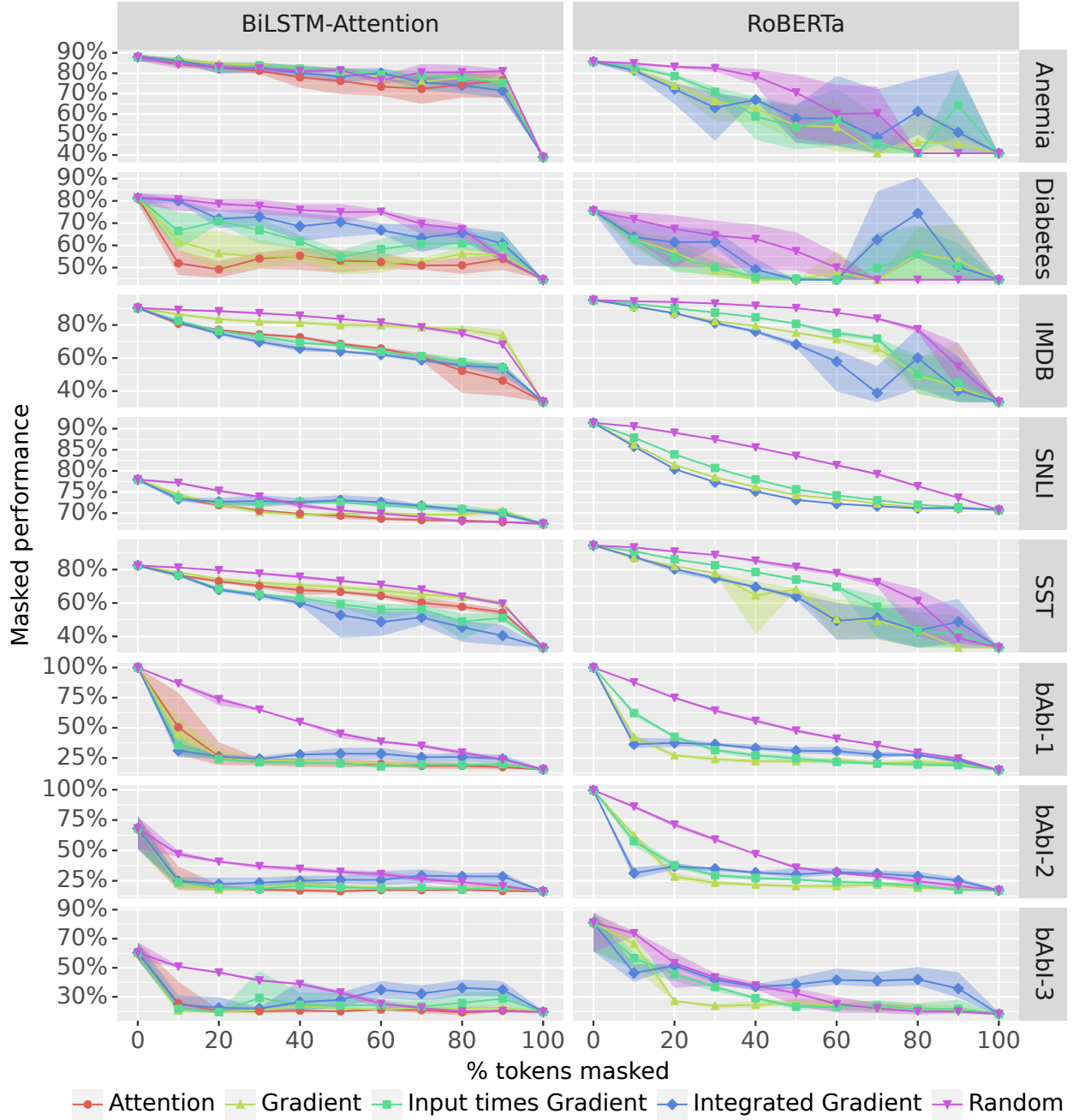


Figure 3.4 Recursive ROAR results, showing model performance at  $x\%$  of tokens masked. A model performance below *random* indicates faithfulness, while above or similar to *random* indicates a non-faithful importance measure. Performance is averaged over 5 seeds with a 95% confidence interval.

### 3.4.2 Summarizing faithfulness metric

While a ROAR plot can provide valuable insights, such as “this importance measure is only faithful for the top-20% most important tokens,” it does not summarize the faithfulness to a scalar metric. Such a metric is useful as it allows for easy comparisons between models, methods, tasks, and, in particular, papers.

To provide a scalar metric, we propose using a **relative area-between-curves** (RACU) metric. Intuitively, an importance measure is more faithful if it has a larger area between the random baseline curve and the importance measure curve. Additionally, a negative area is considered when the importance measure is above the random baseline. Finally, the metric is normalized by an upper bound, where the performance at 100% masking is achieved immediately. A visualization of this calculation can be seen in Figure 3.5.

Using an area-between-curves is useful because, unlike many other summarizing statistics, it is invariant to the step-size used in ROAR. In this case, we have a step size of 10%. Future work may choose a smaller or larger step size depending on their computational resources.

Let  $r_i$  be the masking ratio at step  $i$  out of  $I$  total steps, in our case  $r = \{0\%, 10\%, \dots, 100\%\}$ . Let  $p_i$  be the model performance for a given importance measure and  $b_i$  be the random baseline performance. With this, the metric is defined in (3.4), and we present the results in Table 3.4.

$$\begin{aligned} \text{ACU} &= \sum_{i=0}^{I-1} \frac{1}{2} \Delta x_i (\Delta p_i + \Delta p_{i+1}) \\ \text{RACU} &= \frac{\text{ACU}}{\sum_{i=0}^{I-1} \frac{1}{2} \Delta x_i (\Delta b_i + \Delta b_{i+1})} \end{aligned} \quad (3.4)$$

where  $\Delta x_i = x_{i+1} - x_i$  *step size*

$\Delta p_i = b_i - p_i$  *performance delta*

$\Delta b_i = b_i - b_I$  *baseline delta*

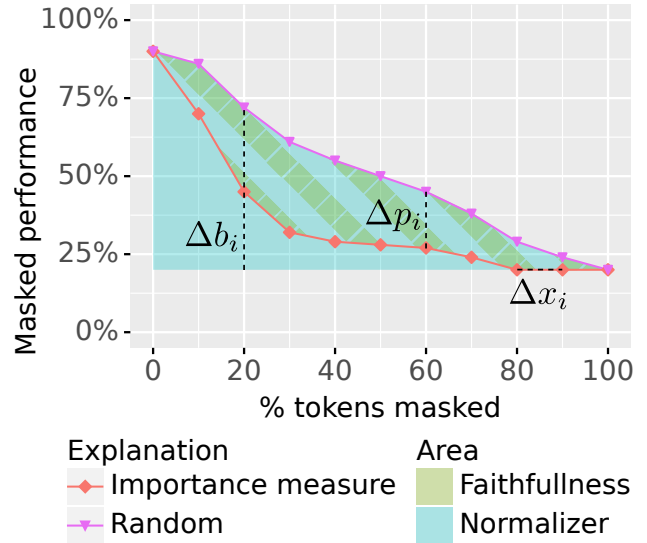


Figure 3.5 Visualization of the faithfulness calculation done in (3.4). The *faithfulness* area is the numerator in (3.4), while the *normalizer* area is the denominator. Essentially (3.4) computes the **relative area-between-curves** (RACU) between an *explanation* curve and the *random* baseline curve.

Dataset	Importance	RACU Faithfulness [%]	
	Measure	LSTM	RoBERTa
Anemia	Attention	$7.6^{+7.9}_{-6.8}$	—
	Gradient	$1.0^{+2.8}_{-4.1}$	$18.2^{+11.8}_{-13.8}$
	Input times gradient	$0.8^{+2.5}_{-3.5}$	$8.8^{+22.7}_{-22.8}$
	Integrated Gradient	$4.9^{+2.7}_{-1.8}$	$12.5^{+11.3}_{-7.0}$
Diabetes	Attention	$66.5^{+6.5}_{-13.0}$	—
	Gradient	$57.4^{+7.8}_{-7.0}$	$57.9^{+14.4}_{-19.8}$
	Input times gradient	$33.7^{+7.0}_{-15.7}$	$53.4^{+23.2}_{-29.3}$
	Integrated Gradient	$11.4^{+8.4}_{-15.0}$	$26.1^{+12.0}_{-25.1}$
IMDB	Attention	$29.8^{+5.0}_{-3.4}$	—
	Gradient	$3.1^{+2.4}_{-3.3}$	$25.4^{+3.1}_{-2.0}$
	Input times gradient	$28.4^{+1.0}_{-0.9}$	$16.9^{+1.1}_{-3.0}$
	Integrated Gradient	$32.5^{+0.9}_{-1.0}$	$35.1^{+2.4}_{-1.7}$
SNLI	Attention	$36.5^{+3.0}_{-3.5}$	—
	Gradient	$18.7^{+5.1}_{-3.5}$	$50.7^{+1.1}_{-0.8}$
	Input times gradient	$-10.7^{+6.1}_{-5.7}$	$41.0^{+0.4}_{-0.5}$
	Integrated Gradient	$-13.9^{+5.0}_{-5.0}$	$56.7^{+1.0}_{-1.1}$
SST	Attention	$15.7^{+2.4}_{-2.4}$	—
	Gradient	$7.6^{+2.3}_{-2.0}$	$26.1^{+1.6}_{-2.2}$
	Input times gradient	$28.0^{+5.6}_{-4.4}$	$18.6^{+4.1}_{-4.6}$
	Integrated Gradient	$37.8^{+4.6}_{-5.3}$	$32.9^{+1.8}_{-1.5}$
bAbI-1	Attention	$66.5^{+9.2}_{-9.2}$	—
	Gradient	$66.1^{+5.9}_{-6.5}$	$64.2^{+2.6}_{-2.6}$
	Input times gradient	$71.2^{+4.0}_{-4.2}$	$52.1^{+1.8}_{-3.7}$
	Integrated Gradient	$59.1^{+6.8}_{-7.4}$	$48.2^{+4.1}_{-5.7}$
bAbI-2	Attention	$75.4^{+4.9}_{-8.1}$	—
	Gradient	$66.3^{+4.2}_{-5.1}$	$57.8^{+2.0}_{-2.0}$
	Input times gradient	$66.7^{+8.0}_{-12.4}$	$48.1^{+3.2}_{-3.5}$
	Integrated Gradient	$34.6^{+13.4}_{-14.8}$	$42.0^{+3.8}_{-4.8}$
bAbI-3	Attention	$77.7^{+9.6}_{-8.1}$	—
	Gradient	$73.0^{+9.1}_{-7.6}$	$34.0^{+14.6}_{-15.1}$
	Input times gradient	$53.9^{+10.7}_{-24.1}$	$22.4^{+15.9}_{-12.4}$
	Integrated Gradient	$25.9^{+8.5}_{-9.1}$	$-27.9^{+18.0}_{-49.1}$

Table 3.4 Faithfulness metric defined as a **relative area-between-curves** (RACU) using Recursive ROAR, see (3.4). Higher values mean more faithful, and zero or negative values mean distinctly not faithful.

### 3.4.3 Supporting experiments

**Class bias and sequence-length bias.** Because Recursive ROAR masks tokens, the sequence length remains the same. At 100% masking, the only information the model has is the sequence length. To understand the relevance of the sequence length, we compare the 100% masking model performance with a basic class-majority classifier. The results in Table 3.5 show that the sequence-length does not have much relevance. SNLI does show a significant difference, but this relates to the secondary sequence being a very good predictor on its own, not the sequence length [189].

Dataset	Majority	LSTM	RoBERTa
Anemia	39%	39% <sup>+0.0%</sup> <sub>-0.0%</sub>	41% <sup>+0.0%</sup> <sub>-0.0%</sub>
Diabetes	45%	45% <sup>+0.0%</sup> <sub>-0.0%</sub>	45% <sup>+0.0%</sup> <sub>-0.0%</sub>
IMDB	34%	33% <sup>+0.1%</sup> <sub>-0.4%</sub>	33% <sup>+0.1%</sup> <sub>-0.3%</sub>
SNLI	34%	67% <sup>+0.3%</sup> <sub>-0.3%</sub>	71% <sup>+0.1%</sup> <sub>-0.1%</sub>
SST	33%	33% <sup>+0.0%</sup> <sub>-0.0%</sub>	33% <sup>+0.0%</sup> <sub>-0.0%</sub>
bAbI-1	15%	15% <sup>+0.8%</sup> <sub>-0.6%</sub>	15% <sup>+0.0%</sup> <sub>-0.0%</sub>
bAbI-2	19%	16% <sup>+0.3%</sup> <sub>-0.4%</sub>	17% <sup>+0.4%</sup> <sub>-0.4%</sub>
bAbI-3	19%	20% <sup>+0.8%</sup> <sub>-1.1%</sub>	18% <sup>+1.2%</sup> <sub>-0.9%</sub>

Table 3.5 Performance of the class-majority classifier and the BiLSTM-Attention and RoBERTa classifier on the 100% masked dataset. Performance is the standard metric for the dataset, meaning macro-F1 for SST, IMDB, Anemia, and Diabetes; micro-F1 for SNLI; and accuracy for bAbI.

**Effect of redundancies.** In Figure 3.6, we compare *ROAR* and *Recursive ROAR*. These results show dataset redundancies interfere with *ROAR*. For example, consider the Diabetes dataset, only when using *Recursive ROAR* is the *gradient* IM shown to be faithful. A comparison of all datasets and models, along with more detailed analysis, can be found in Appendix B.

**Effect of relative masking.** In Figure 3.7, we avoid the approximation of removing a relative number of tokens at 10% increments by instead removing exactly one token in each iteration. These results show that the approximation does affect the results, but not the conclusions that can be drawn from the results. A comparison of all datasets and models and more detailed analysis can be found in Appendix B.

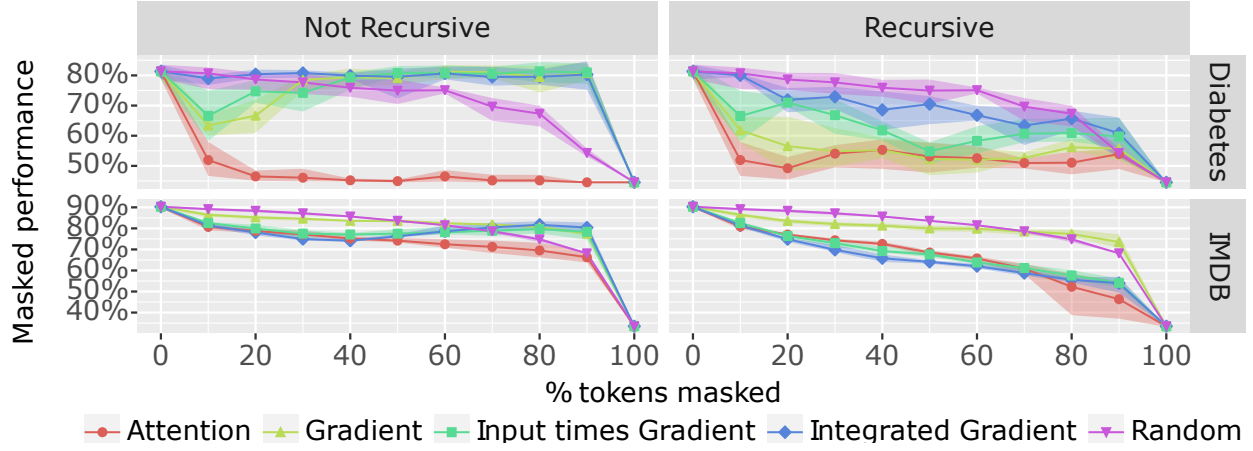


Figure 3.6 ROAR and Recursive ROAR results for BiLSTM-Attention, showing model performance at  $x\%$  of tokens masked. A model performance below *random* indicates faithfulness. For Recursive ROAR a curve above or similar to *random* indicates a non-faithful importance measure, while for ROAR, this case is inconclusive [14].

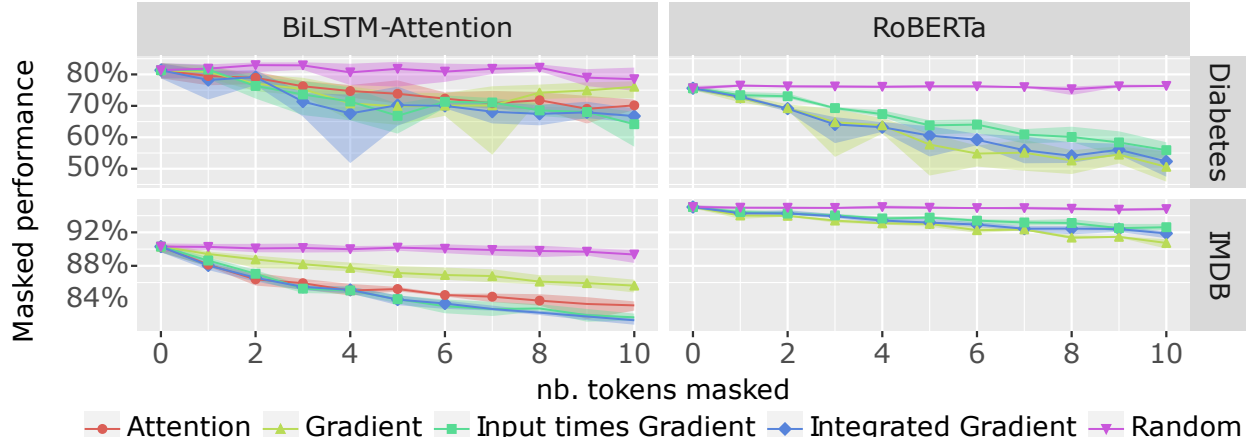


Figure 3.7 Recursive ROAR results, showing model performance at up to 10 tokens masked. Note that because the datasets have more than 10 tokens, the conclusion one can draw from this plot may change if more tokens are considered. However, a model performance below *random* indicates faithfulness, while above or similar to *random* indicates a non-faithful importance measure. Performance is averaged over 5 seeds with a 95% confidence interval.

**Sparsity of explanations.** In Figure 3.8, we report the sparsity of each importance measure and find that *attention* is significantly more sparse than other importance measures. If the faithfulness is equal, this may make it more desirable as sparse explanations are more understandable to humans [134]. A comparison of all datasets and models, along with more detailed analysis, can be found in Appendix B.

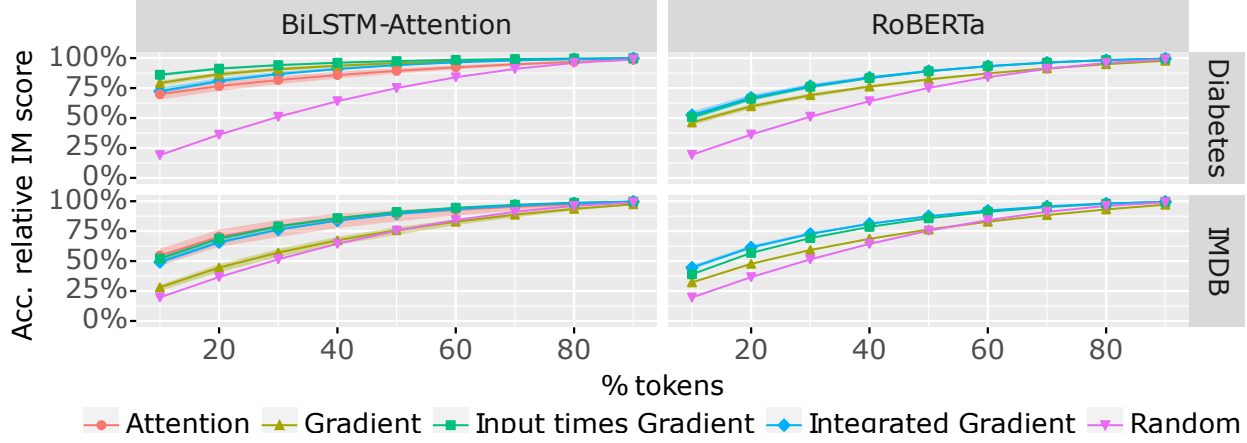


Figure 3.8 The accumulative importance score relative to the total importance score for the top- $x\%$  number of tokens. The metric is averaged over 5 seeds with a 95% confidence interval.

### 3.5 Important Findings

Based on the results in Figure 3.4 and Table 3.4, we highlight the following important findings.

#### Faithfulness is model-dependent

In particular, the faithfulness with SNLI is highly model-dependent as seen in Table 3.4. Furthermore, comparing the faithfulness between the two models, the faithfulness of *Gradient* on IMDB and *Integrated Gradient* on bAbI-3 is significantly affected by the model architecture.

#### Faithfulness is task-dependent

For BiLSTM-Attention, in Table 3.4, *Attention* is best for SNLI while *Input times Gradient* and *Integrated Gradient* is best for SST.

For RoBERTa, *Integrated Gradient* is best for IMDB and SNLI, while *Gradient* is best for bAbI-1 and bAbI-2. In fact, *Integrated Gradient* is worst in all bAbI tasks.

### Attention can be faithful

In Table 3.4, *Attention* is among the top explanations in terms of faithfulness, except for SST. This contradicts many of the previous results mentioned in Section 3.1, which found attention to be unfaithful.

Because attention is computationally free and attention is more sparse (Appendix B), which is important for human understanding [134], attention can be an attractive explanation.

### Integrated Gradient is not necessarily more faithful than Gradient or Input times Gradient

It’s commonly considered that *Integrated Gradient* is more faithful than other gradient methods. However, for BiLSTM-Attention, in Table 3.4, bAbI-2, bAbI-3, and SNLI there is at least one gradient-based importance measure which is significantly more faithful than *Integrated Gradient*. For RoBERTa, we find the same for bAbI-2 and bAbI-3. These results contradict the claim that Integrated Gradient is theoretically superior [39]. This is a valuable finding, as Integrated Gradient is significantly more computationally expensive than other gradient-based importance measures.

### Importance measures often work best for the top-20% most important tokens.

In Figure 3.4, we observe that the largest drop tends to happen at about 10% or 20% tokens masked. This indicates that importance measures are best at ranking the most important tokens, while for less important tokens, they become noisy. This is particularly observed in bAbI for both models and Diabetes with the BiLSTM-Attention model.

## 3.6 Limitations

### Measures on a retrained model

Recursive ROAR requires the model to be retrained. This means it is not possible to evaluate the faithfulness of a specific model instance, rather we evaluate the faithfulness of the model architecture. The confidence intervals we provide then inform us about what can be statistically expected in terms of the faithfulness for a model instance.

## Computationally expensive

The retraining dependence also means Recursive ROAR can only measure the faithfulness of a task-model combination that is feasible to train/fine-tune repeatedly and importance measures that are feasible to compute across the entire dataset.

## Potential class leakage

Because the importance measures explain predictions of the target label, they can leak the target label when allegedly important tokens are masked. This can make an importance measure appear less faithful than it actually is. However, this issue cannot make an importance measure appear more faithful than it is. The mechanism for this is rather complicated and non-intuitive, hence a longer explanation is provided here.

**Explanation.** When importance measures are computed, it is the prediction of the gold label that is explained. For example, for the *Gradient* method, it is  $\nabla_{\mathbf{x}}f(x)_y$  that is computed, where  $\mathbf{x}$  is the input and  $y$  is the gold label.

We want an importance measure for the correct label, as removing the tokens that are relevant for making a wrong prediction, would help the performance of the model. If the gold label was not used, the faithfulness results would be affected by the model performance. As faithfulness and model performance should be unrelated, this is not a desired outcome.

In ROAR and Recursive ROAR, this issue is expressed as an increase in the model performance. Intuitively, it should not be possible for the model performance to increase with more information removed compared to less. However, because the importance measures are w.r.t. the gold label, they can leak the gold label which can increase the model performance.

**Thought experiment.** Consider a sentiment classification task, like IMDB. Let's say that the **awful** token correlates with the negative label, but can still appear in positive sentiment sentences like **I have an awful strong crush on this actor.**

Then propose that just using the **awful** token provides an 80% accurate classification of negative labels and the model learns this. A faithful importance measure would therefore highlight the **awful** token as being important for the prediction of negative sentiment. When measuring faithfulness, the **awful** tokens are thus removed from sentences with negative sentiment as the gold label. This creates a new dataset where the existence of an **awful** token is now a perfect predictor of positive sentiment, and thus the model performance may increase.

Assuming a faithful importance measure, in the next iteration of Recursive ROAR the `awful` token would now be important for predicting positive sentiment and would be removed. In practice, this is not guaranteed since `awful` may not be among the most important tokens, and in the case where a relative number of tokens are masked, the removal of other tokens may leak the gold label.

**Implications.** This issue is particularly observed with bAbI-3 using RoBERTa in Figure 3.4, where the performance increases slightly at 60% tokens masked. This issue affects both ROAR and Recursive ROAR (Appendix B). In fact, it likely affects most faithfulness metrics. Additionally, because ROAR presents a more qualitative metric (Figure 3.4) where a curve can be observed to increase, this issue is more apparent. Had we just presented the summarizing metric (Table 3.4), as most faithfulness metrics do, the issue would have been hidden.

### 3.7 Conclusion

We show that Recursive ROAR is an improvement on ROAR. In a synthetic setting, Recursive ROAR matches the ground truth, while ROAR does not. Additionally, we argue why other faithfulness metrics may be either invalid or limited in scope.

We then use Recursive ROAR to measure the faithfulness of the most common importance measures, including attention. This is done on both recurrent and transformer-based neural models.

We provide a list of the most important findings in Section 3.5. In general, we find that the faithfulness of importance measures is both model-dependent and task-dependent. This means that no general recommendation can be made for NLP practitioners considering the current importance measures. Instead, it is necessary to measure the faithfulness of different importance measures given a task and a model.

In the next chapter we present a new category of models that makes measuring faithfulness much easier to apply, and solve the mentioned limitations of this approach. However, the Recursive ROAR metric remains the only general metric that works for general models and real-world datasets, not just synthetic problems. We therefore hope it can serve as a foundation for measuring faithfulness of importance measures in NLP.

## CHAPTER 4 FAITHFULNESS MEASURABLE MODELS

In Chapter 3, we developed the faithfulness metric *Recursive ROAR* based on the erasure-metric [60] and used retraining to solve the out-of-distribution issue that masking tokens create. The conclusion was that faithfulness is both model and task-dependent. This conclusion was also made by another simultaneous work that used a synthetic dataset with a known explanation to measure faithfulness [91].

This leaves us in the unfortunate situation where we cannot say an explanation method is generally faithful, rather it's necessary to measure it for a given task and model. This makes *Recursive ROAR* and most other faithfulness metrics covered in Section 3.1 troublesome to use.

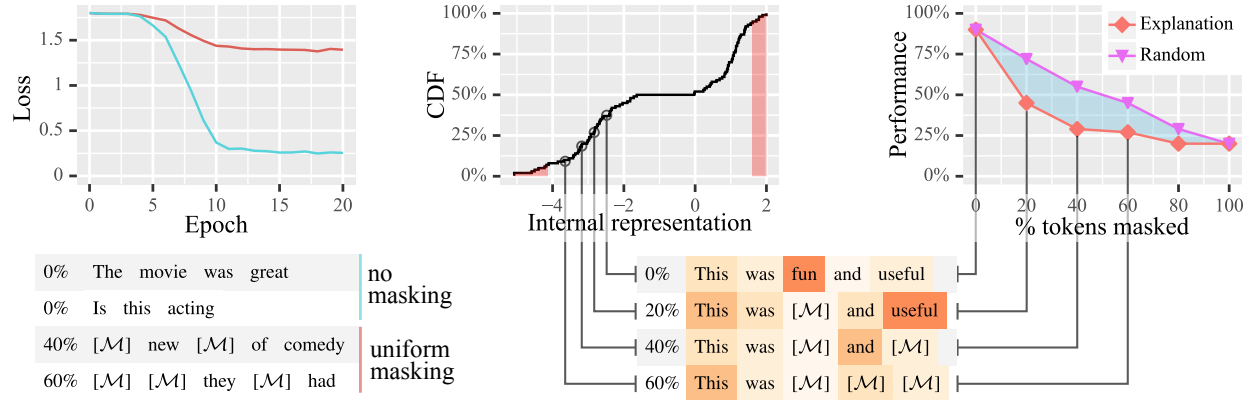
Firstly, *Recursive ROAR* is computationally expensive because it requires repeated retraining for every explanation method, model, and task. Additionally, because *Recursive ROAR* retrains the model, it's no longer measuring faithfulness on the deployed model of interest. Therefore, we may get a misleading perception of the faithfulness w.r.t. the deployed model, leading to a false confidence in an explanation. This is particularly possible due to the model-dependent conclusion. Then, *Recursive ROAR* also had an issue where the gold-class may be leaked, causing deviances in the faithfulness metric. Finally, it may also be the case that faithfulness is instance-dependent, meaning given the same model, explanation, and task, the faithfulness for each observation's explanation may be very different.

This is what motivates the desirables **(b)** and **(e)** from Chapter 3:

- b) The method measures faithfulness of an explanation w.r.t. a specific model instance and single observation. For example, it is not a proxy-model that is measured.
- e) The method is computationally cheap by not training/fine-tuning repeatedly and only computes explanations of the test dataset.

All the mentioned issues with *Recursive ROAR* are because of retraining. However, without retraining masking tokens can create out-of-distribution issues. Therefore, the key idea in this chapter is to use the erasure-metric but without using retraining to solve the out-of-distribution (OOD) issue, which avoids all the limitations.

Instead of retraining, in-distribution support for masking can be achieved by including masking in the fine-tuning procedure of masked language models as a data augmentation (Figure 4.1a). This is possible because language models are heavily over-parameterized and can thus support such additional complexity. Although our approach applies to Masked Language Models



(a) **Masked fine-tuning.** In-distribution support for masking any permutation of tokens is achieved by uniformly masking half of the mini-batch during fine-tuning. The other half is left unmasked to maintain regular unmasked performance.

(b) **In-distribution validation.** CDFs of the model’s embeddings given a masked validation dataset, provide in-distribution p-values and validate that test observations masked according to an explanation are in-distribution.

(c) **Faithfulness metric.** Observations are masked according to an explanation. A model performance lower than masking random tokens means the explanation is faithful. A larger area between the curves means more faithful.

Figure 4.1 To measure faithfulness, a *faithfulness measurable masked language model* is created (a), then the model is checked for out-of-distribution issues given an explanation (b), and finally, the faithfulness is measured by masking allegedly important tokens (c). –  $[M]$  is the masking token.

(MLMs), we suspect future work could apply this idea to any language model with sufficient capacity.

Our approach is significantly different from previous literature, which is completely model agnostic. Rather than considering the challenge of measuring faithfulness after the model has been trained, we propose to design a model using simple fine-tuning such that measuring faithfulness of explanations is cheap and precise. We call such designed models: **inherently faithfulness measurable models** (FMMs).

To validate that masking is in-distribution, we generalize previous OOD detection work from computer vision [46]. This serves as a statistically grounded meta-validation of the faithfulness measure itself (Figure 4.1b). Finally, once the model is validated, the erasure-metric can be applied (Figure 4.1c).

Note, the concept of an *inherently faithfulness measurable model* (FMM) is significantly different from inherently explainable models from the *intrinsic* paradigm, which are interpretable by design [59]. The inherently explainable models provide an explanation by design, whereas

an FMM provides a faithfulness measure by design. As such, an FMM does not guarantee that an explanation exists, and an inherently explainable model doesn’t provide a means of measuring faithfulness.

However, with an FMM, measuring faithfulness is computationally cheap. Therefore, optimizing an explanation w.r.t. faithfulness is possible, as proposed by Zhou and Shah [35]. However, they did not solve the OOD issue caused by masking, as it was “orthogonal” to their idea, but our *inherently faithfulness measurable model* fills that gap, making it indirectly inherently explainable.

Finally, for completeness, we compare a large variety of existing explanation methods and modify some existing explanations to be able to separate positive from negative contributions. In general, we find that the explanations that take advantage of masking (occlusion-based) are more faithful than gradient-based methods. However, the robustness provided by a *faithfulness measurable model* also makes some gradient-based methods more faithful.

**To summarize, this chapter’s contributions are:**

- Introducing the concept of an *inherently faithfulness measurable model* (FMM).
- Proposing *masked fine-tuning* that enables masking to be in-distribution.
- Establishing a statistically grounded meta-validation for the faithfulness measurable model, using out-of-distribution detection.
- Making existing occlusion-based explanations more faithful, as they no longer cause out-of-distribution issues.
- Introducing signed variants of existing importance measures, which can separate between positive and negative contributing tokens.

#### 4.1 Inherently faithfulness measurable models (FMMs)

As an alternative to existing faithfulness methods for importance measures, which all aim to work with any models, we propose creating *inherently faithfulness measurable models* (FMMs). These models provide the typical output (e.g., classification) for a given task and, by design, provide the means to measure the faithfulness of an explanation. Importantly, this allows measuring the faithfulness of a specific model, as there is no need for proxy models, an important property in a real deployment setting. This idea is similar to what Hase et al. [190] and Vafa et al. [191] proposed for top-k explanations.

An FMM does have the limitation that a specific model is required. However, our proposed method is very general, as it only requires a modified fine-tuning procedure applied to a

masked language model.

#### 4.1.1 Faithfulness of importance measures

In this chapter, we look at importance measures (IMs), which are explanations that either score or rank how important each input token is for making a prediction. A faithfulness metric measures how much such an explanation reflects the true reasoning process of the model [59]. Importantly, such a metric should work regardless of how the importance measure is calculated.

For IMs, there are multiple definitions of truth. In this chapter, we again use the erasure-metric definition: *if information (tokens) is truly important, then masking them should result in worse model performance compared to masking random information (tokens)* [14, 60, 171].

The challenge with an erasure-metric is that fine-tuned models do not support masking tokens. Even masked language models are usually only trained with 12% or 15% masking [4, 21, 192], and an erasure-metric use between 0% and 100% masking. Furthermore, catastrophic forgetting of the masking token is likely when fine-tuning.

In Chapter 3, we introduced Recursive ROAR, which solved this issue by retraining the model with partially masked inputs. Unfortunately, that approach has issues, as discussed in Section 3.6 and the beginning of this chapter. It is computationally expensive, leaks the gold label, and measures a proxy model instead of the true model.

We find that the core issue is the need for retraining. Instead, if the fine-tuned model supports masking any permutation of tokens, then retraining would not be required, eliminating all issues. We propose a new fine-tuning procedure called *masked fine-tuning* to achieve this.

To evaluate faithfulness of an importance measure, we propose a three-step process, as visualized in Figure 4.1:

1. Create a faithfulness measurable masked language model, using *masked fine-tuning*. – See Section 4.1.2.
2. Check for out-of-distribution (OOD) issues by using a statistical in-distribution test. – See Section 4.1.3.
3. Measure the faithfulness of an explanation. – See Section 4.1.4.

#### 4.1.2 Masked fine-tuning

To provide masking support in the fine-tuned model, we propose randomly masking the training dataset by uniformly sampling a masking rate between 0% and 100% for each

observation and then randomly masking that ratio of tokens. However, half of the mini-batch remains unmasked to maintain the regular unmasked performance. This is analogous to multi-task learning, where one task is masking support, and the other is regular performance. However, it is slightly different, as some multi-task learning methods may sample randomly between the two tasks, where we split deterministically. Other methods may also switch between the two tasks in each step, we don't do this as it can create unstable oscillations. Instead, both tasks are included in the same mini-batch.

To include masking support in early stopping, the validation dataset is duplicated, where one copy is unmasked, and one copy is randomly masked.

The high-level idea is similar to Hase et al. [190] which enabled masking support for a fixed number of tokens. However, to support a variable number of tokens we improve upon this work by sampling a masking-ratio. There are many approaches to implementing masked fine-tuning with identical results. However, Algorithm 1 presents our implementation.

---

**Algorithm 1** Creates the mini-batches used in masked fine-tuning.

---

**Precondition:**  $B$  is a mini-batch with  $N$  randomly sampled observations from the training dataset.  $[\mathcal{M}]$  is the masking token.

```

1 function MINIBATCH( $B$ )
2    $M \leftarrow \emptyset$  ▷ Stores new mini-batch
3   for  $i \leftarrow 1$  to  $N$  do
4     if  $i$  is even then
5        $r_i \leftarrow \text{SAMPLEUNIFORM}(0, 1)$ 
6        $M_i \leftarrow \text{MASKTOKENS}(B_i, r_i)$  ▷ Masks  $r\%$  randomly
       selected tokens in  $B_i$ .
7     else
8        $M_i \leftarrow B_i$ 
9   return  $M$ 

10 function MASKTOKENS( $x, r$ )
11   for  $t \leftarrow 1$  to  $T$  do
12      $s_t \leftarrow \text{SAMPLEUNIFORM}(0, 1)$ 
13     if  $s_t < r$  then
14        $\tilde{x}_t \leftarrow [\mathcal{M}]$  ▷ Masks token  $t$  with probability  $r$ .
15     else
16        $\tilde{x}_t \leftarrow x_t$ 
17   return  $\tilde{x}$ 

```

---

### 4.1.3 In-distribution validation

Erasure-based metrics are only valid when the input masked according to the importance measure is in-distribution, and previous works did not validate for this [14, 131, 190, 191]. Additionally, in-distribution is the statistical null-hypothesis and can never be proven. However, we can validate this using an out-of-distribution (OOD) test.

We use the *MaSF* method by Matan et al. [46] as the OOD test, which provides non-parametric p-values under the in-distribution global-null-hypothesis, it thus tests if all of the models embeddings given an observation are in-distribution. Then, to provide a p-value for the entire masked test dataset being in-distribution we perform another Simes [193] aggregation.

At its core, *MaSF* is an aggregation of many in-distribution p-values, where each p-value is from an in-distribution test of a latent embedding. That is, given a history of embedding observations, which presents a distribution, what is the probability of observing the new embedding or something more extreme? For example, if that probability is less than 5%, it could be classified as out-of-distribution at a 5% risk of a false-positive.

However, a model has many internal embeddings, and thus, there will be many p-values. If each were tested independently, there would be many false positives. This is known as p-hacking. To prevent this, the p-values are aggregated using the Simes and Fisher methods, which are aggregation methods for p-values that prevent this issue. Once the p-values are aggregated, it becomes a global null-test. This means the aggregated statistical test checks if any of the embeddings are out-of-distribution.

This is where the term *MaSF* comes from, it’s an acronym that stands for Max-Simes-Fisher, which is the order of aggregation functions it uses. Matan et al. [46] presented some other combinations and orders of aggregation functions but found this to be the best. However, many combinations had similar performance in their benchmark, so we do not consider this particular choice to be important.

**Empirical CDF.** Each p-value is computed using an empirical cumulative distribution function (CDF). A nice property of an empirical CDF is that it doesn’t assume any distribution, a property called non-parametric. It is however still a model, only if an infinite amount of data was available would it represent the true distribution.

A CDF measures the probability of observing  $z$  or less than  $z$ , i.e.,  $\mathbb{P}(Z \leq z)$ . The empirical version simply counts how many embeddings were historically less than the tested embedding, as shown in (4.1). However, we are also interested in cases where the embedding is abnormally large, hence we also use  $\mathbb{P}(Z > z) = 1 - \mathbb{P}(Z \leq z)$ . We are then interested in the most unlikely

case, which is known as the two-sided p-value, i.e.  $\min(\mathbb{P}(Z \leq z), 1 - \mathbb{P}(Z \leq z))$ .

$$\mathbb{P}(Z \leq z) = \frac{1}{|Z_{\text{emp}}|} \sum_{i=1}^{|Z_{\text{emp}}|} 1[Z_{\text{emp},i} < z] \quad (4.1)$$

The historical embeddings are collected by running the model on the validation dataset ( $\mathcal{D}_V$ ). Note that for this to be accurate, the validation dataset should be i.i.d. with the training dataset. This can easily be accomplished by randomly splitting the datasets, which is common practice, and transforming the validation dataset the same way as the training dataset (i.e. applying the masked fine-tuning transformation).

---

**Algorithm 2** MaSF algorithm, provides in-distribution p-values.

---

**Precondition:**  $x$  is the input.  $f_e(x) \in \mathbb{R}^{T \times H \times L}$  provides the model embeddings, where  $T$  is sequence-length,  $H$  is the hidden-size, and  $L$  is the number of layers.  $\mathbb{P}$  are the empirical CDFs; these are collected by running the MaSF algorithm on a validation dataset.

```

1 function MASF( $x, \mathbb{P}$ )
2    $e \leftarrow f_e(x)$  ▷ Get embeddings
3   for  $l \leftarrow 1$  to  $L$  do
4     for  $h \leftarrow 1$  to  $H$  do
5        $z_{l,h}^{(1)} \leftarrow \max_{t=1}^T e_{l,h,t}$  ▷ Ma-step
6        $\tilde{p}_{l,h}^{(1)} \leftarrow \mathbb{P}_{l,h}^{(1)}(Z < z_{l,h}^{(1)})$ 
7        $p_{l,h}^{(1)} \leftarrow \min(\tilde{p}_{l,h}^{(1)}, 1 - \tilde{p}_{l,h}^{(1)})$ 
8        $z_l^{(2)} \leftarrow \text{Simes}(p_{l,:}^{(1)})$  ▷ S-step
9        $\tilde{p}_l^{(2)} \leftarrow \mathbb{P}_l^{(2)}(Z < z_l^{(2)})$ 
10       $p_l^{(2)} \leftarrow \min(\tilde{p}_l^{(2)}, 1 - \tilde{p}_l^{(2)})$ 
11       $z^{(3)} \leftarrow \text{Fisher}(p^{(2)})$  ▷ F-step
12       $\tilde{p}^{(3)} \leftarrow \mathbb{P}^{(3)}(Z < z^{(3)})$ 
13       $p^{(3)} \leftarrow 1 - \tilde{p}^{(3)}$ 
14      return  $p^{(3)}$ 

15 function SIMES( $p$ )
16    $q \leftarrow \text{SortAscending}(p)$ 
17   return  $\min_{i=1}^N q_i \frac{N}{i}$ 

18 function FISHER( $p$ )
19   return  $-2 \sum_{i=1}^N \log(p_i)$ 

```

---

**Algorithm.** In the case of *MaSF*, the embeddings are first aggregated along the sequence dimension using the max operation. Matan et al. [46] only applied *MaSF* to computer vision,

in which case it was the width and height dimensions. However, we generalize this to NLP by swapping width and height with the sequence dimension. Additionally, in the case of RoBERTa we use the embeddings after the layer-normalization, which is standard practice in other embedding-based applications [194].

The max-aggregated embeddings from the validation dataset provide the historical data for the empirical CDFs. If a network has  $L$  layers, each with  $H$  latent dimensions, there will be  $L \cdot H$  CDFs. The same max-aggregated embeddings are then transformed into p-values using those CDFs. Next, the p-values are aggregated using Simes’s method [193] along the latent dimension, which provides another set of CDFs and p-values, one for each layer ( $L$  CDFs and p-values). Finally, those p-values are aggregated using Fisher’s method [195], providing one p-value for each observation and one CDF.

The algorithm for *MaSF* can be found in Algorithm 2. While this algorithm does work, practical implementation is in our experience non-trivial, as for the entire test dataset ( $\mathcal{D}_T$ ) there are  $\mathcal{O}(|\mathcal{D}_T| \cdot H \cdot L)$  CDFs evaluations, each involving  $\mathcal{O}(|\mathcal{D}_V|)$  comparisons with the validation dataset ( $\mathcal{D}_V$ ). While this is computationally trivial on a GPU, it requires a lot of memory usage when done in parallel. Therefore, we found that a practical implementation must batch in a mesh style over both the test and validation datasets.

#### 4.1.4 Faithfulness metric

To measure faithfulness on a model trained using *masked fine-tuning*: the importance measure (IM) is computed for a given input, then  $x\%$  (e.g., 10%) of the most important tokens are masked, then the IM is calculated on this masked input, finally an additional  $x\%$  of the most important tokens are masked. This is repeated until 100% of the input is masked. The importance measure is re-calculated because otherwise, dataset redundancies will interfere with the metric, as shown in Chapter 3.

At each iteration, the masked input is validated using MaSF and the performance is measured. Faithfulness is shown if and only if the performance is less than when masking random tokens. This procedure is identical to Recursive-ROAR in Chapter 3, but without retraining and with in-distribution validation, as shown in Algorithm 3 and Figure 4.2.

#### 4.1.5 Optimizing for faithfulness (Beam)

With a fast and validated faithfulness metric, optimizing the explanation for maximal faithfulness is possible. In this chapter, we use the method proposed by Zhou and Shah [35]. The central idea is to use beam-search, where the generated sequence is the optimal masking order

---

**Algorithm 3** Measures the masked model performance given an explanation.

---

**Precondition:**  $\text{IM}(f, x, y) \in \mathbb{R}^T$  explains the model  $f$  for the input  $x$  and label  $y$ , with sequence-length  $T$ .  $\delta$  is the iterations step-size (e.g. 10%) and  $[\mathcal{M}]$  is the masking token.

---

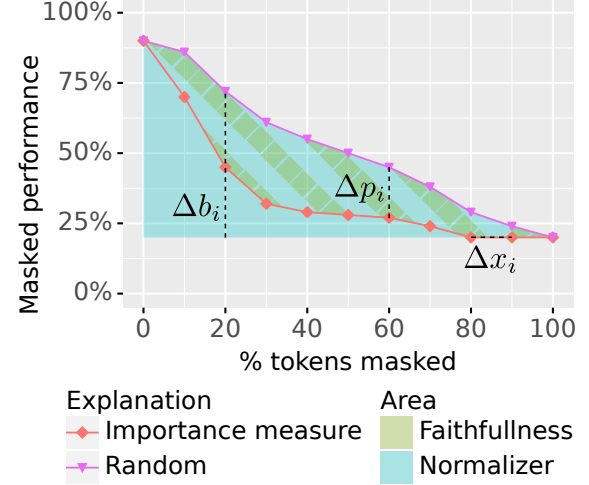
```

1 function RECURSIVEEVAL(IM,  $f$ ,  $x$ ,  $y$ ,  $\delta$ )
2    $\tilde{x}_0 \leftarrow x$ 
3    $p_0 \leftarrow \text{PerformanceMetric}(f(x), y)$ 
4   for  $i \leftarrow 1$  to  $1/\delta$  do
5      $e_i \leftarrow \text{IM}(f, \tilde{x}_{i-1}, y)$ 
6      $\tilde{x}_i \leftarrow \text{AddMask}(e_i, \tilde{x}_{i-1}, \delta) \triangleright \text{Mask}$ 
7      $\delta \cdot T$  more tokens in  $\tilde{x}_{i-1}$  using scores  $e_i$ .
8      $p_i \leftarrow \text{PerformanceMetric}(f(\tilde{x}_i), y)$ 
9   return  $p$ 

```

---

Figure 4.2 Visualization of the faithfulness calculation. AUC is the *faithfulness* area, and RACU is the AUC normalized by the theoretical best explanation. See the definition for AUC and RACU in (3.4).



of tokens. Each iteration of the beam-search masks one additional token, where the token is selected by testing every possibility and maximizing the faithfulness metric. This could be reframed as a recursive version of leave-one-out. They propose several optimization targets, but since our faithfulness metric is analog to comprehensiveness we use this variation.

Note that Zhou and Shah [35] simply ignored the out-of-distribution caused by masking. Thus causing an otherwise sound method to become unfaithful. However, this should not be a concern because the proposed *faithfulness measurable masked language model* supports masking. Importantly, this exemplifies how a *faithfulness measurable model* may also produce more faithful explanations.

Additionally, the number of forward passes is approximately  $\mathcal{O}(B \cdot T^2)$ ; ( $B$  is the beam-size and  $T$  is the sequence-length). Thus, this approach is quite computationally costly, although one advantage is that this explanation is inherently recursive, hence it is not necessary to reevaluate the importance measure in each iteration of the faithfulness metric. However, for long sequence datasets, such as IMDB, BaBi-3, Anemia, and Diabetes, it is not feasible to apply this explanation. In our experiments, we use a beam-size of  $B = 10$ .

## 4.2 Experiments

We choose the RoBERTa model because it converges consistently, and reasonable hyperparameters are well established. This should make reproducing the results easier. We use both the **base** (125M parameters) and **large** size (355M parameters). We use the default GLUE hyperparameters provided by Liu et al. [4, Appendix C, GLUE].

Although these hyperparameters are for the GLUE tasks, we use them for all tasks. The one exception is that the maximum number of epochs is higher. However, when using early stopping with the validation dataset, the optimization is not sensitive to the specific number of epochs, lower numbers are only used to reduce the compute time.

We present results on 16 classification datasets in the appendix but only include BoolQ, MRPC, and bAbI-2 in this chapter. These were chosen as they represent the general trends we observe, although we observe very consistent results across all datasets. All datasets used are public and listed below, and are all used for their intended use, which is measuring classification performance.

Note that only the first sequence is considered when computing the importance measure on paired-sequence tasks. This is to stay consistent with previous work [2] and because it does not make sense to mask the question for tasks like document-based Q&A (e.g., bAbI).

The essential statistics for each dataset, and which part is masked and auxiliary, are specified in Table 4.2. Details regarding max-epochs and reproduced baseline performance are in Table 4.1.

We use 5 seeds for each experiment and present their means with their 95% confidence interval (error-bars or ribbons). The 95% confidence interval is computed using the bias-corrected and accelerated bootstrap method [196, 197]. When relevant, each seed is presented as a plus (+).

### 4.2.1 Masked fine-tuning

There are two criteria for learning our proposed *faithfulness measurable masked language model*:

1. The usual performance metric, where no data is masked, should not decrease.
2. Masking any permutations of tokens should be in-distribution.

In Section 4.1.2, we propose *masked fine-tuning*, where one half of a mini-batch is uniformly masked between 0% and 100% and the other half is unmasked. Additionally, the validation dataset contains a masked copy and an unmasked copy.

Dataset	max epoch	Performance	
		RoBERTa- base	RoBERTa- large
BoolQ	15	80% <sup>+0.2</sup> <sub>-0.2</sub>	85% <sup>+0.3</sup> <sub>-0.3</sub>
CB	50	65% <sup>+17.6</sup> <sub>-47.9</sub>	87% <sup>+3.1</sup> <sub>-8.2</sub>
CoLA	15	59% <sup>+1.3</sup> <sub>-1.1</sub>	66% <sup>+0.8</sup> <sub>-0.8</sub>
IMDB	10	95% <sup>+0.2</sup> <sub>-0.2</sub>	96% <sup>+0.2</sup> <sub>-0.4</sub>
Anemia	20	84% <sup>+0.8</sup> <sub>-0.7</sub>	84% <sup>+0.5</sup> <sub>-0.8</sub>
Diabetes	20	76% <sup>+0.9</sup> <sub>-0.9</sub>	77% <sup>+0.6</sup> <sub>-1.6</sub>
MNLI	10	87% <sup>+0.4</sup> <sub>-0.2</sub>	90% <sup>+0.3</sup> <sub>-0.2</sub>
MRPC	20	86% <sup>+0.8</sup> <sub>-0.7</sub>	87% <sup>+0.5</sup> <sub>-1.2</sub>
QNLI	20	92% <sup>+0.1</sup> <sub>-0.1</sub>	94% <sup>+0.1</sup> <sub>-0.2</sub>
QQP	10	90% <sup>+0.1</sup> <sub>-0.1</sub>	91% <sup>+0.0</sup> <sub>-0.1</sub>
RTE	30	75% <sup>+1.4</sup> <sub>-2.9</sub>	83% <sup>+1.3</sup> <sub>-1.4</sub>
SNLI	10	91% <sup>+0.1</sup> <sub>-0.2</sub>	92% <sup>+0.1</sup> <sub>-0.2</sub>
SST2	10	94% <sup>+0.2</sup> <sub>-0.2</sub>	96% <sup>+0.2</sup> <sub>-0.2</sub>
bAbI-1	20	100% <sup>+0.0</sup> <sub>-0.1</sub>	100% <sup>+0.0</sup> <sub>-0.0</sub>
bAbI-2	20	99% <sup>+0.1</sup> <sub>-0.1</sub>	100% <sup>+0.1</sup> <sub>-0.1</sub>
bAbI-3	20	90% <sup>+0.2</sup> <sub>-0.3</sub>	90% <sup>+0.5</sup> <sub>-0.5</sub>

Table 4.1 Max-epoch parameters and performance statistics for each dataset. Performance metrics are the mean with a 95% confidence interval.

**Unmasked performance.** To validate the first goal, Figure 4.3 presents an ablation study. It compares *masked fine-tuning* with using only unmasked data (*plain fine-tuning*), as is traditionally done, and using only uniformly masked data (*only masking*). The unmasked performance is then measured (the usual benchmark).

We observe that no performance is lost when using our *masked fine-tuning*; some tasks even perform better, likely because masking has a regularizing effect. However, when using *only masking* performance is lost and unstable convergence is frequent. For bAbI-2&3, we do observe unstable convergence when using *masked fine-tuning*. However, this is less frequent (worst case: 3/5) and only for RoBERTA-large (see also Appendix C). Note the default RoBERTa hyperparameters are not meant for synthetic datasets like bAbI. Therefore, optimizing the hyperparameter would likely solve the stability issues with *masked fine-tuning*. Finally, when using *masked fine-tuning*, the models do need to be trained for slightly more epochs (twice more or less); see Appendix C. Again, tuning the hyperparameters would likely help.

**100% Masked performance.** Measuring in-distribution support for masked data is challenging, as there is generally no known performance baseline. However, for 100% masked data,

Type	Dataset	Size			Inputs		Performance		Citation
		Train	Validation	Test	masked	auxiliary	metric	class-majority	
NLI	RTE	1992	498	277	sentence1	sentence2	Accuracy	47%	[29]
	SNLI	549367	9842	9824	premise	hypothesis	Macro F1	34%	[30]
	MNLI	314162	78540	9815	premise	hypothesis	Accuracy	35%	[12]
	QNLI	83794	20949	5463	sentence	question	Accuracy	51%	[31]
	CB	200	50	56	premise	hypothesis	Macro F1	22%	[23]
Paraphrase	MRPC	2934	734	408	sentence1	sentence2	Macro F1	41%	[27]
	QQP	291077	72769	40430	question1	question2	Macro F1	39%	[32]
Sentiment	SST2	53879	13470	872	sentence	–	Accuracy	51%	[13]
	IMDB	20000	5000	25000	text	–	Macro F1	33%	[26]
Diagnosis	Anemia	4262	729	1243	text	–	Macro F1	39%	[9]
	Diabetese	8066	1573	1729	text	–	Macro F1	45%	[9]
Acceptability	CoLA	6841	1710	1043	sentence	–	Matthew	0%	[24]
QA	BoolQ	7542	1885	3270	passage	question	Accuracy	62%	[22]
	bAbI-1	8000	2000	1000	paragraph	question	Micro F1	15%	[10]
	bAbI-2	8000	2000	1000	paragraph	question	Micro F1	19%	[10]
	bAbI-3	8000	2000	1000	paragraph	question	Micro F1	18%	[10]

Table 4.2 Datasets used, all datasets are either single-sequence or sequence-pair datasets. All datasets are sourced from GLUE [7], SuperGLUE [8], MIMIC-III [9], or bAbI [10]. The decisions regarding which metrics are used are also from these sources. The class-majority baseline is when the most frequent class is always selected.

only the sequence length is left as information. Therefore, the performance of a model should be at least that of the class-majority baseline, where the most frequent class is “predicted” for all observations. We present an ablation study using this baseline in Figure 4.4. In Section 4.2.2, we perform a more in-depth validation.

From Figure 4.4, we observe that training with unmasked data (*Plain fine-tuning*) performs worse than the class-majority baseline, clearly showing an out-of-distribution issue. However, when using masked data, either *only masking* or *masked fine-tuning*, both effectively achieve in-distribution results for 100% masked data.

**The best approach used in the following experiments.** Appendix C contains a more detailed ablation study separation of the training and validation strategy. We find that the choice of validation dataset is not very significant. However, we consider it the most principled approach to using both unmasked and masked data, i.e. *Masked fine-tuning*. Besides this, the conclusion is the same. *Masked fine-tuning* is the only method that achieves good results for both the unmasked and 100% masked cases.

For the following experiments in Section 4.2.2 and Section 4.2.3, the *masked fine-tuning* method is used. Additionally, we will only present results for RoBERTa-base for brevity.

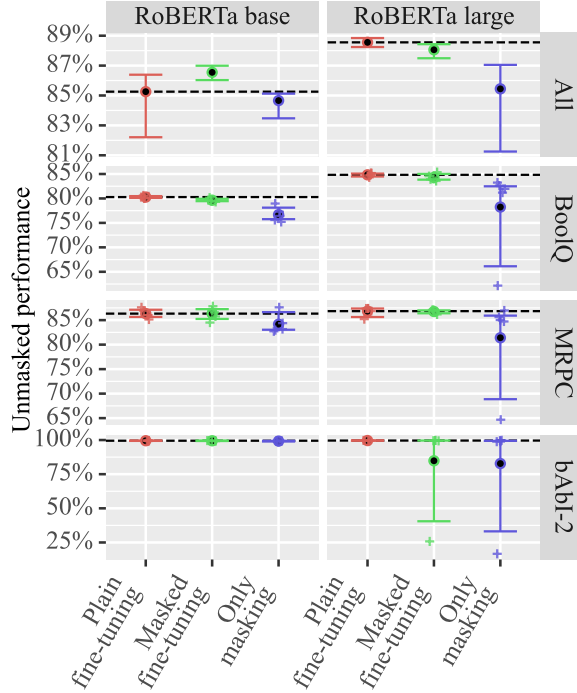


Figure 4.3 The unmasked performance for each fine-tuning strategy. *Plain fine-tuning* is the baseline (dashed line). We find that our *Masked fine-tuning* does not decrease performance. *All* is computed by taking the average of all datasets. More datasets and a more detailed ablation study can be found in Appendix C.

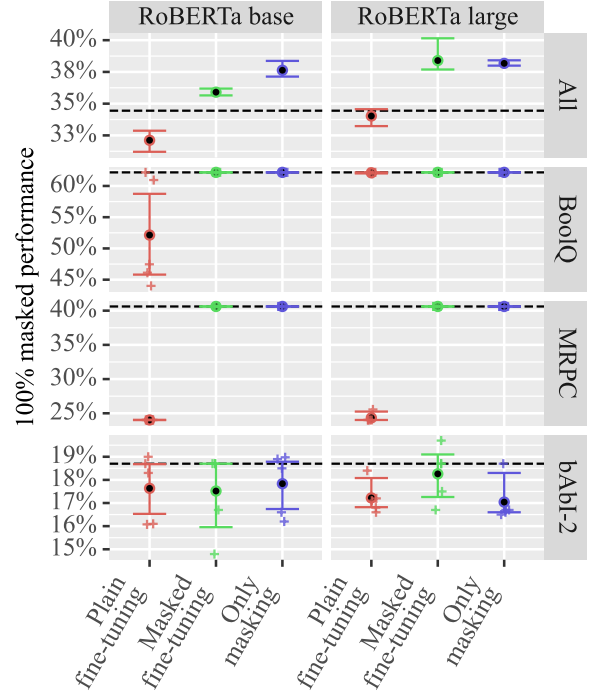


Figure 4.4 The 100% masked performance for each fine-tuning strategy. The dashed line represents the class-majority baseline. Results show that masking during training (either our *masked fine-tuning* or *only masking*) is necessary. More datasets and a more detailed ablation study can be found in Appendix C.

RoBERTa-large results are included in the appendix.

#### 4.2.2 In-distribution validation

Because the expected performance for masked data is generally unknown, a statistical in-distribution test called *MaSF* [46] is used instead, as was explained in Section 4.1.3.

MaSF provides an in-distribution p-value for each observation. To test if all masked test observations are in-distribution, the p-values are aggregated using Simes’s method [193]. Because in-distribution is the null-hypothesis, we can never confirm in-distribution; we can only validate it. Rejecting the null hypothesis would mean that some observation is out-of-distribution.

Because random uniform masking is not the same as strategically masking tokens, we validate in-distribution for each importance measure, where the masking is done according to the

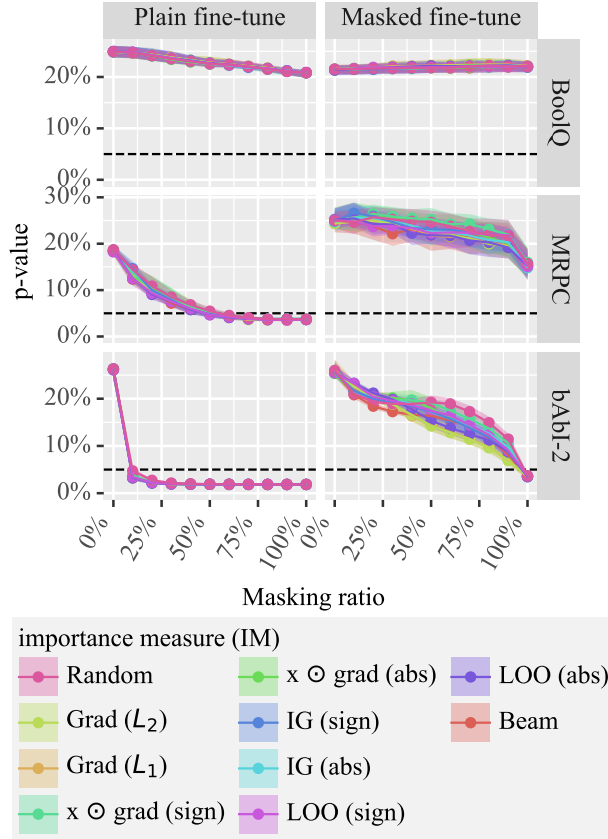


Figure 4.5 In-distribution p-values using MaSF, for RoBERTa-base with and without masked fine-tuning. The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using *masked fine-tuning* is masking consistently not OOD. Because the results are highly consistent, the overlapping lines do not hide any important details. More datasets and models in Appendix C.

Grad is “Gradient”,  $x \odot grad$  is “Input times gradient”, IG is “Integrated Gradient”, LOO is “Leave-one-out”. These methods are described in Section 2.5.

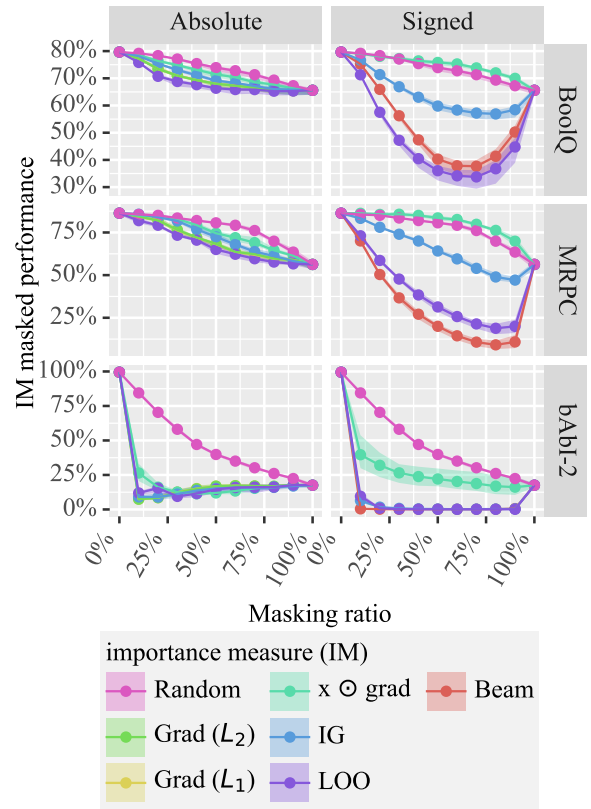


Figure 4.6 The performance given the masked datasets, where masking is done for the  $x\%$  allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “Random” baseline, this shows faithfulness. Although faithfulness is not an absolute concept, so more is better. This plot is for RoBERTa-base and separates importance measures based on their signed and absolute variants. More datasets and models in Appendix C.

Legend has same meaning as in Figure 4.5 but does not separate between signed (sign) and absolute (abs) importance measures.

importance measure, identically to how the faithfulness metric is computed (Section 4.1.4).

Additionally, because MaSF does not consider the model’s performance, it is necessary to consider these results in combination with regular performance metrics, see Section 4.2.1.

The results for when using *masked fine-tuning* and *plain fine-tuning* (no masking) are presented in Figure 4.5. The results show that masked datasets are consistently in-distribution only when using masked fine-tuning.

In the case of BoolQ, we suspect that because the training dataset is fairly small (7542 observations), the model does not completely forget the mask token. Additionally, a few datasets, such as bAbI-2, become out-of-distribution at 100% masking when using masked fine-tuning (see also Appendix C). This contradicts the performance results for 100% masked data (Figure 4.4 and Appendix C), which clearly show in-distribution performance. This is likely a limitation of MaSF because the empirical CDF in MaSF has very little 100% masked data, as the masking ratio is uniform between 0% and 100%. Fortunately, this is not a concern because Figure 4.4 shows in-distribution results for 100% masked data.

### 4.2.3 Faithfulness metric

Based on previous experiments, we can conclude that *masked fine-tuning* achieves both objectives: unaffected regular performance and support for masked inputs. Therefore, it is safe to apply the faithfulness metric to these models.

Briefly, the faithfulness metric works by showing that masking using an importance measure (IM) is more effective at removing important tokens than using a known false explanation, such as a random explanation. Therefore, if a curve is below the random baseline, the IM is faithful. Although faithfulness is not an absolute [59], so further below would indicate more faithful.

Because signed importance measures can differentiate between positive and negative contributing tokens, while absolute tokens are not, it is to be expected that signed importance measures are more faithful. However, comparing them might not be fair because of this difference in capability. We let the reader decide this for themselves.

**Relative Area Between Curves (RACU)** In Chapter 3, we propose to compute the area between the random curve and an explanation curve (RACU). This is then normalized by the theoretical optimal explanation, which would achieve the performance of 100% masking immediately. However, the normalization is only theoretically optimal for an absolute importance measure (IM). Signed IMs can trick the model into predicting the opposite label,

Dataset	IM	Faithfulness [%]		
		Our		Recursive ROAR
		ACU	RACU	RACU
Anemia	Gradient ( $L_2$ )	23.8 <sup>+0.6</sup> <sub>-0.5</sub>	62.1 <sup>+1.4</sup> <sub>-1.7</sub>	18.2 <sup>+11.8</sup> <sub>-13.8</sub>
	Gradient ( $L_1$ )	23.8 <sup>+0.6</sup> <sub>-0.6</sub>	62.2 <sup>+1.4</sup> <sub>-2.1</sub>	—
	Input times gradient (sign)	9.7 <sup>+2.6</sup> <sub>-2.5</sub>	25.1 <sup>+6.5</sup> <sub>-6.2</sub>	—
	Input times gradient (abs)	16.6 <sup>+1.3</sup> <sub>-1.3</sub>	43.2 <sup>+3.0</sup> <sub>-3.7</sub>	8.8 <sup>+22.7</sup> <sub>-22.8</sub>
	Integrated gradient (sign)	62.0 <sup>+1.6</sup> <sub>-1.5</sub>	161.8 <sup>+2.7</sup> <sub>-2.3</sub>	—
	Integrated gradient (abs)	20.0 <sup>+0.9</sup> <sub>-1.6</sub>	52.1 <sup>+2.5</sup> <sub>-4.5</sub>	12.5 <sup>+11.3</sup> <sub>-7.0</sub>
	Leave-on-out (sign)	63.3 <sup>+1.4</sup> <sub>-1.6</sub>	165.2 <sup>+2.9</sup> <sub>-3.2</sub>	—
	Leave-on-out (abs)	18.9 <sup>+1.0</sup> <sub>-1.3</sub>	49.2 <sup>+2.5</sup> <sub>-3.8</sub>	—
	Beam	—	—	—
Diabetes	Gradient ( $L_2$ )	19.7 <sup>+1.2</sup> <sub>-0.7</sub>	91.8 <sup>+0.6</sup> <sub>-0.9</sub>	57.9 <sup>+14.4</sup> <sub>-19.8</sub>
	Gradient ( $L_1$ )	19.6 <sup>+1.0</sup> <sub>-0.7</sub>	91.6 <sup>+0.5</sup> <sub>-0.9</sub>	—
	Input times gradient (sign)	10.9 <sup>+1.9</sup> <sub>-1.3</sub>	51.1 <sup>+9.1</sup> <sub>-7.8</sub>	—
	Input times gradient (abs)	18.8 <sup>+1.3</sup> <sub>-0.7</sub>	87.9 <sup>+1.4</sup> <sub>-2.0</sub>	53.4 <sup>+23.2</sup> <sub>-29.3</sub>
	Integrated gradient (sign)	24.8 <sup>+1.5</sup> <sub>-2.1</sub>	115.8 <sup>+2.8</sup> <sub>-7.8</sub>	—
	Integrated gradient (abs)	19.4 <sup>+1.0</sup> <sub>-0.6</sub>	90.5 <sup>+0.6</sup> <sub>-1.2</sub>	26.1 <sup>+12.0</sup> <sub>-25.1</sub>
	Leave-on-out (sign)	41.5 <sup>+2.3</sup> <sub>-5.9</sub>	193.4 <sup>+8.9</sup> <sub>-17.4</sub>	—
	Leave-on-out (abs)	19.1 <sup>+1.1</sup> <sub>-0.6</sub>	89.0 <sup>+0.4</sup> <sub>-0.6</sub>	—
	Beam	—	—	—
SST2	Gradient ( $L_2$ )	12.2 <sup>+0.6</sup> <sub>-0.7</sub>	40.4 <sup>+3.0</sup> <sub>-1.7</sub>	26.1 <sup>+1.6</sup> <sub>-2.2</sub>
	Gradient ( $L_1$ )	12.1 <sup>+0.7</sup> <sub>-0.7</sub>	40.3 <sup>+3.3</sup> <sub>-1.8</sub>	—
	Input times gradient (sign)	-3.7 <sup>+1.5</sup> <sub>-1.6</sub>	-12.2 <sup>+4.5</sup> <sub>-6.0</sub>	—
	Input times gradient (abs)	7.1 <sup>+0.2</sup> <sub>-0.2</sub>	23.5 <sup>+1.9</sup> <sub>-1.1</sub>	18.6 <sup>+4.1</sup> <sub>-4.6</sub>
	Integrated gradient (sign)	31.8 <sup>+2.8</sup> <sub>-2.2</sub>	105.6 <sup>+7.7</sup> <sub>-7.7</sub>	—
	Integrated gradient (abs)	13.7 <sup>+0.8</sup> <sub>-0.8</sub>	45.3 <sup>+4.1</sup> <sub>-2.8</sub>	32.9 <sup>+1.8</sup> <sub>-1.5</sub>
	Leave-on-out (sign)	51.6 <sup>+1.4</sup> <sub>-0.9</sub>	171.3 <sup>+5.8</sup> <sub>-6.2</sub>	—
	Leave-on-out (abs)	16.6 <sup>+1.2</sup> <sub>-1.0</sub>	54.9 <sup>+2.1</sup> <sub>-1.5</sub>	—
	Beam	56.4 <sup>+0.5</sup> <sub>-0.7</sub>	187.3 <sup>+8.1</sup> <sub>-7.1</sub>	—
bAbI-2	Gradient ( $L_2$ )	28.5 <sup>+0.8</sup> <sub>-0.8</sub>	96.3 <sup>+6.8</sup> <sub>-2.8</sub>	57.8 <sup>+2.0</sup> <sub>-2.0</sub>
	Gradient ( $L_1$ )	28.5 <sup>+0.9</sup> <sub>-0.8</sub>	96.3 <sup>+6.8</sup> <sub>-2.7</sub>	—
	Input times gradient (sign)	19.7 <sup>+6.6</sup> <sub>-8.1</sub>	65.7 <sup>+24.1</sup> <sub>-26.3</sub>	—
	Input times gradient (abs)	27.3 <sup>+1.7</sup> <sub>-1.5</sub>	92.0 <sup>+2.5</sup> <sub>-3.1</sub>	48.1 <sup>+3.2</sup> <sub>-3.5</sub>
	Integrated gradient (sign)	40.3 <sup>+0.9</sup> <sub>-0.8</sub>	136.3 <sup>+4.4</sup> <sub>-6.4</sub>	—
	Integrated gradient (abs)	29.1 <sup>+1.0</sup> <sub>-1.3</sub>	98.3 <sup>+5.5</sup> <sub>-3.9</sub>	42.0 <sup>+3.8</sup> <sub>-4.8</sub>
	Leave-on-out (sign)	40.2 <sup>+1.2</sup> <sub>-0.8</sub>	136.0 <sup>+4.1</sup> <sub>-6.5</sub>	—
	Leave-on-out (abs)	28.5 <sup>+0.9</sup> <sub>-1.4</sub>	96.3 <sup>+9.2</sup> <sub>-3.6</sub>	—
	Beam	41.1 <sup>+1.0</sup> <sub>-0.7</sub>	139.2 <sup>+5.0</sup> <sub>-7.3</sub>	—

Table 4.3 Faithfulness scores using Relative Area Between Curves (RACU) and the non-relative variant (ACU). The less relevant score is grayed out. Higher is better. Negative values indicate not-faithful. The comparison with Recursive-ROAR from Chapter 3 is imperfect because Recursive-ROAR has limitations. See Table C.7 for all datasets and Table C.8 for RoBERTa-large.

thus achieving even lower performance. For this reason, we also show the un-normalized metric (ACU) in Table 4.3, as defined in Section 3.4.2.

Note that comparing with Recursive ROAR is troublesome because Recursive ROAR has issues, such as leaking the gold label. Additionally, while Chapter 3 also uses RoBERTa-base, it’s not the same model because we use masked fine-tuning.

### 4.3 Important Findings

Based on the faithfulness scores as seen in Table 4.3, and also Table C.7 and Table C.8, we here discuss the most important findings.

#### Consistently faithfulness importance measures

from the RACU results, in Table 4.3 and Appendix C, it becomes apparent that there exist importance measures which are consistently faithful across all 16 datasets and both RoBERTa-base and RoBERTa-large. In fact of all the explanations, only the signed variant of *input-times-gradient* is not consistently faithful. This is a drastic improvement compared to Chapter 3 where the conclusion was that faithfulness is model- and task-dependent.

#### Major improvements in faithfulness

While comparing with Recursive-ROAR is not exact, we can observe roughly a 2 to 5 times improvement for most importance measures.

#### Near perfect explanations

For absolute importance measures it’s possible to consider a theoretical perfect explanation. At 100% masking the performance is at it’s lowest. Thus, the best possible explanation is one where just masking the most important token, attains the lowest possible performance level (i.e. same as 100% masking). This is not practically possible for most datasets, because most datasets have some redundancies. However, for synthetic datasets like bAbI, it’s only necessary to mask 2 words to remove all relevant information. Thus, for such datasets getting close to 100% relative faithfulness is realistic. As seen in Table 4.3, we do indeed get near theoretically perfect explanations.

### Occlusion-based importance measures are the best

Leave-one-out and Beam are consistently among the best explanations. This is reasonable as they directly take advantage of the masking support that masked fine-tuning provides. Thus there is a synergy between the explanation and the model. Furthermore, Beam is also directly optimizing for faithfulness using a beam-search method. Seeing that it’s always among the top explanations, therefore, validates that optimizing for faithfulness is possible.

### Gradient-based importance measures also improve

We also observe that gradient-based explanations are more faithful when using our model. We suspect this is partially also because there is no leakage issue. However, previous work has also shown that gradient-based methods behave more favorably on robust models in computer vision [198]. Using masked fine-tuning can be seen as a robustness objective, as the model becomes robust to missing information.

### Faithfulness measures have lower variance

Table 4.3 show that the RACU scores have a lower variance (confidence interval) using our methodology compared to Recursive ROAR. This is likely because Recursive ROAR from Chapter 3 leaks the gold label, which causes oscillation in the faithfulness curve.

## 4.4 Limitations

### No faithfulness ablation with regular fine-tuning

We claim *masked fine-tuning* makes importance measures (IMs) more faithful. However, there is no ablation study where we measure faithfulness without *masked fine-tuning*. This is because, without *masked fine-tuning*, masking is out-of-distribution which makes the faithfulness measure invalid.

However, our argument for occlusion-based IMs has a theoretical foundation, as occlusion (i.e., masking) is only in-distribution because of *masked fine-tuning*. We also observe that occlusion-based IMs are consistently more faithful than gradient-based IMs. Finally, for gradient-based IMs, we compare with Recursive ROAR from Chapter 3, and this approach provides more faithful explanations, although this comparison is imperfect as discussed in Section 4.2.3.

## Not a post-hoc method

While this work solves existing limitations with previous methods, it introduces the significant limitation that it, by definition, requires a *faithfulness measurable model*. As such, the question of faithfulness needs to be considered ahead of time when developing a model. It can not be an afterthought, which is often how interpretability is approached [47].

While this is a significant limitation, considering explanation ahead of deployment is increasingly becoming a legal requirement [52]. Currently, the European Union provides a “right to explanation” regarding automatic decisions, which includes NLP models [93].

## In-distribution is impossible to prove

Because in-distribution is always the null hypothesis, it is impossible to statistically show that inputs are truly in-distribution. The typical approach to similar statistical questions<sup>1</sup> is to keep validating in-distribution using various methods. Unfortunately, the literature on this topic in deep learning is extremely limited [46, 199, 200, 201, 202].

Therefore, we would advocate for more work on identifying out-of-distribution inputs using non-parametric methods that primarily consider the model’s internal state. Using parametric methods or works that use axillary models is more well-explored but not useful for our purpose.

## Requires repeated measures on the test dataset

Because datasets have redundancies, it is necessary to reevaluate the importance measures, as discussed in Chapter 3. This leads to an increased computational cost.

However, unlike Recursive ROAR from Chapter 3, this method only requires reevaluation of the test dataset, which is often quite small. Additionally, some IMs, such as the beam-search method [35], take dataset redundancies into account and therefore do not require reevaluation. Reevaluation could be done if desired but would result in the exact same results.

## Uses masked language models (MLMs)

Masked fine-tuning leverages pre-trained MLMs’ partial support for token masking. Therefore, our approach does not immediately generalize to casual language models (CLM). However, despite CLMs’ popularity for generative tasks, MLMs are still very relevant for classification

---

<sup>1</sup>A similar well-explored statistical question is how to show that the error in a linear model is normally distributed.

tasks [203] and for non-NLP tasks, such as analyzing biological sequences (genomes, proteins, etc.) [204].

Additionally, it is possible to introduce the mask tokens to CLMs by masking random tokens in the input sequence while keeping the generation objective the same, similar to how unknown-word tokens are used. This approach could also be done in an additional pre-training step using existing pre-trained models. Regardless, masking support for CLMs is likely a more complex task and is left for future work.

Another direction useful for classification tasks, is to transform CLMs into MLMs, which has been shown to be quite straightforward [205]. It may also be possible to simply prompt an instruction-tuned CLM, such that it understands what masking means, for example in Chapter 5 we prompt with “The following content may contain redacted information marked with [REDACTED]”.

In terms of supporting sequential outputs rather than just classification outputs, our methodology only requires a performance metric. Using sequential performance metrics such as ROUGE [206] or BLEU should therefore work perfectly well.

## 4.5 Conclusion

Using only a simple modified fine-tuning method, called *masked fine-tuning*, we are able to turn a typical general-purpose masked language model (RoBERTa) into an *inherently faithfulness measurable model* (FMM). Meaning that the model, by design, inherently provides a way to measure the faithfulness of importance measure (IM) explanations.

To the best of our knowledge, this is the first work that proposes creating a model designed to be faithfulness measurable. Arguably, previous work in top-k explanations [190, 191] and counterfactual explanations [150, 160] have indirectly achieved something similar. However, their motivation was to provide explanations or robustness, not measuring faithfulness.

Importantly, our approach is very general, simple to apply, and satisfies critical desirables that previous faithfulness measures didn’t. The *masked fine-tuning* method does not decrease performance on all 16 tested datasets while also adding in-distribution support for token masking, which we are able to verify down to fundamental statistics using an out-of-distribution test.

We find that when using masked fine-tuning there exist consistently faithful importance measures. This is a significantly different finding than in Chapter 3, where faithfulness was found to always be model and task-dependent. However, either due to the robustness produced by masked fine-tuning or the correctness of the faithfulness metric in FMMs (likely

a combination of both), this model and task-dependent conclusion no longer hold.

In particular, we find that occlusion-based IMs are among the most faithful. This is to be expected, as they take advantage of the masking support. Additionally, Beam uses beam-search to optimize towards faithfulness [35], which our proposed faithfulness measurable masked language model makes computationally efficient to evaluate.

It is worth considering the significance of this. While our proposed model is not an *inherently explainable model* [59], it is *indirectly* inherently explainable because it provides a built-in way to measure faithfulness, which can then be optimized for. It does this without sacrificing the generality of the model, as it is still a RoBERTa model. As such, FMMs provide a new direction for interpretability, which bridges the gap between *post-hoc* [130] and *inherent* interpretability [48]. It does so by prioritizing faithfulness measures first and then the explanation, while previous directions have worked on explanation first and then measure faithfulness [129].

However, beam-search is just an approximative optimizer, which only achieves perfect explanations at infinite beam-width, and Leave-one-out does occasionally outperform Beam for that reason. Future work could look at better optimization methods to improve the faithfulness of importance measures.

## CHAPTER 5 FAITHFULNESS OF SELF-EXPLANATIONS

The previous chapters have looked at recurrent neural networks and masked language models (e.g., RoBERTa [4]). However, during the development of this Ph.D., instruction-tuned large language models (LLMs), such as Llama2 [207], Falcon [208], Mistral [209], or GPT4 [210], have increasingly become popular. They are even becoming mainstream among the general population due to their capabilities and availability.

These models can also provide very convincing explanations for their utterances and will often do so unprompted. Because LLMs produce these explanations themselves and they provide justification for their own behavior, we term them *self-explanations*. It’s also well established that LLMs hallucinate [211, 212]; therefore, it’s possible these self-explanations are unfaithful, which would create unsupported confidence in the model’s capabilities [213, 214].

Unfortunately, because of self-explanations’ free-formed nature, they are extra hard to measure the faithfulness of [182]. Additionally, many contemporary LLMs only provide an inference API and often refuse to make discrete predictions when critical information is missing; these properties make previous faithfulness metrics that depend on confidence scores impractical to apply [215].

**To solve these challenges, we propose** a faithfulness metric that only uses an inference API and takes advantage of the model’s reluctance to answer when critical information is missing. We achieve this by limiting the scope of self-explanations to those verifiable using self-consistency checks and by carefully prompting the model regarding both the explanation and classification generation.

A *self-consistency check* is when re-evaluation is used to check if the explanation is consistent with the model’s behavior. For example, consider a hiring recommendation system. In 2018, Amazon found that “Women’s chess club membership” was a negative signal while “chess club membership” was a positive signal, indicating a harmful gender bias [66]. Such a contrastive example is known as a counterfactual (see Section 2.6). With LLMs, one can ask the model itself to edit the resume so that it would make the opposite prediction. We can then learn about the model behavior by comparing the two resumes, assuming the counterfactual is faithful. To check faithfulness, we ask the model for its hiring recommendation using the counterfactual resume. If the recommendation changed, as requested, it’s a faithful explanation. Figure 5.1 shows the prediction, explanation, and self-consistency check workflow.

In this chapter, we evaluate the faithfulness of the following types of self-explanations:

- **Counterfactual** explanations replace content, to get the opposite prediction. For example, replace “Women’s chess club” with “Chess club”.
- **Feature attribution** explanations list the necessary words for making a prediction, such that without these words a prediction can not be made. For example, 1) “Women” 2) “BSc degree”.
- **Redaction** explanations remove all relevant words for making a prediction, such that a prediction can no longer be made. For example, “[REDACTED] chess club”.

We demonstrate our approach on four datasets with varying tasks: sentiment classification (IMDB [26]), multi-choice classification (bAbI and MCTest [10, 28]), and two-paragraph classification (RTE [29]). Additionally, we apply the approach to Llama2 (70B, 7B), Falcon (40B, 7B), and Mistral (7B). The variability of these tasks and models shows the generality of our approach.

We find that the faithfulness of instruction-tuned LLMs depends on the model, explanation, and task. For example, regarding Llama2 (70B), counterfactuals only work with IMDB, and feature attribution only works with RTE and bAbI, clearly showing a task dependence. Additionally, we show our findings are robust to prompt variations.

Because faithfulness is explanation and task-dependent, we suggest self-explanations cannot generally be trusted and propose how future work might address this challenge.

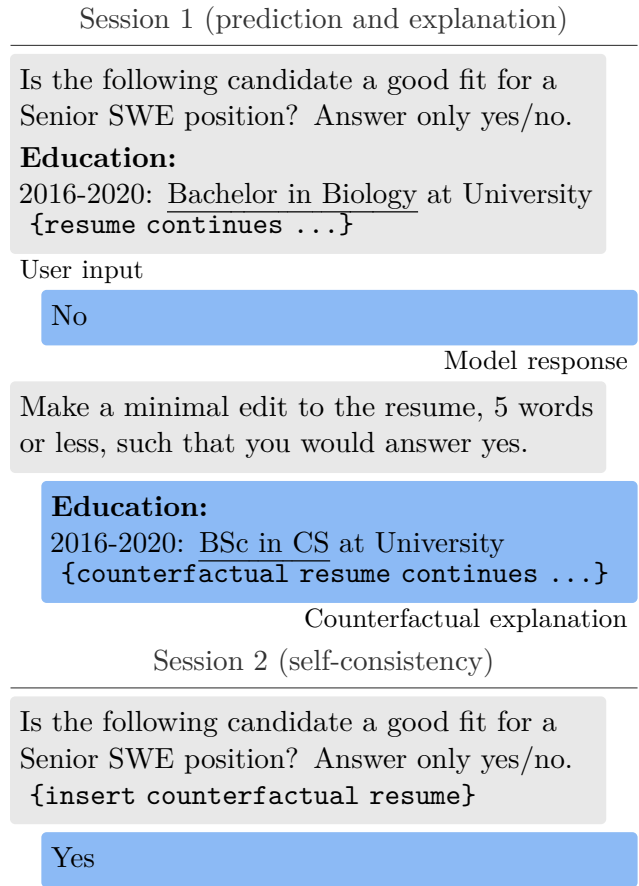


Figure 5.1 Example of an LLM providing a counterfactual self-explanation and using a self-consistency check to evaluate if it is faithful. – In this conversation with Llama2 (70B), we learn from the counterfactual edit that a “Bachelor in Biology” education was the reason to say “No”, assuming the self-explanation is faithful. Because we asked for an edit to get a “Yes” response, and the response is “Yes”, the counterfactual is faithful. Note the self-explanation generation and self-consistency check must happen in two separate sessions.

## 5.1 Self-explanations

LLMs’ ability to produce highly convincing self-explanations is a new development in the field of interpretability. As shown in Chapter 2, previously a separate model or algorithm generated the explanation, not the predictive model itself. This development creates new challenges and opportunities [216].

For example, counterfactuals are often generated using an explanation model which is almost completely independent of the classification model they explain [42, 160]. The explanations may be generated by fine-tuning on a dataset with human-annotated explanations, where the classification model is only used to filter the generated explanations [150, 217]. This is a problem because humans have no insight into the classification model’s behavior [59, 175]. Therefore, there is little reason to think that the counterfactuals represent the classification model [174].

Instead, we let the language model generate both the classification and the counterfactual explanation. In principle, because the weights are the same, the explanation generation can access or simulate the classification behavior; a property known as self-model capabilities [218]. Anthropic Team [218] explores the self-model capabilities of LLMs using a self-consistency check and find that LLMs have poor self-modeling capabilities.

That being said, it is possible that self-model capabilities are helpful but not necessary for generating self-explanations. If the LLM reproduces human behavior sufficiently well in both the classification and counterfactual case, then the explanations would be self-consistent and thus faithful.

Because Anthropic Team [218] have already explored self-model capabilities, and it might not be necessary for faithful self-explanations, we limit the scope in this chapter to faithfulness evaluation.

## 5.2 Faithfulness of self-explanations

This section describes the general methodology we propose and its principles, using counterfactuals as an example. Section 5.3 then goes into detail on the different kinds of explanations and datasets where this approach is possible and how we precisely formulate the prompts.

### 5.2.1 Self-consistency check

Counterfactuals are explanations where the input is modified such that the model predicts the opposite label. The intention is to provide a contrastive example to explain the model’s

behavior [42, 150]. For example, in the hiring-recommendation example (Figure 5.1), instead of asking for a general explanation (e.g. “Why are you recommending not to hire”), we directly asked the language model to modify the resume such that it would yield a positive hiring recommendation.

Because the goal is to produce the opposite label, we can re-evaluate the prediction with the counterfactual and check that this goal is satisfied. If satisfied, the counterfactual is faithful. We refer to such re-evaluation as a **“self-consistency check”**.

Note that when using an instruction-tuned LLM in a conversational setting, it’s important to perform the re-evaluation using a new chat session. Otherwise, the chat model may predict the opposite label only because it was previously prompted to do so within the chat history.

### 5.2.2 Negative results and robustness

If the self-consistency check does not pass, we don’t know if the model cannot generate faithful counterfactuals in general or if a different prompt or generation sample would generate faithful explanations. Therefore, negative results are less informative than positive results.

However, it’s important to consider the practical application where an end user asks a chat model for an explanation. Regular users may not be critical regarding the explanation. As such, we desire models that are robust to natural prompt variations and minor sampling differences, by providing faithful explanations in all cases.

Therefore, when measuring faithfulness, we use six different prompt patterns in order to assess the LLMs’ robustness. Furthermore, we consider finding just one non-faithful prompt pattern significantly problematic.

## 5.3 Prompt-patterns for self-explanation and self-consistency checks

This section provides the specific details and considerations behind the prompts used in the experiment section. Specifically, the prompts that generate the counterfactual, feature attribution, and redaction explanations. As well as the prompt for generating the classification.

To facilitate the “specific details,” we use sentiment classification as an example. However, the methodology generalizes to multi-choice Q&A datasets (e.g., bAbI) and multi-paragraph classification (e.g., NLI). Prompts for those tasks are in Appendix D and Appendix D, respectively.

All prompts are used with chat-tuned models, as these models are typically deployed and

provide a well-defined framing mechanism between input and output<sup>1</sup>. Such framing helps to avoid some out-of-distribution issues compared to purely instruction-tuned models. However, the prompts only convey a one-time input-output relation and thus could be used with purely instruction-tuned models given the proper framing.

It should be noted that it’s impossible to prove that the model understands the user’s intent. However, it’s worth prioritizing the user’s intent rather than the model’s understanding when evaluating the potential impact of deployed models.

### 5.3.1 Counterfactual explanation

We generate counterfactuals by asking for an edit that generates the opposite sentiment. In Figure 5.2 we explicitly express the target sentiment in the prompt. To evaluate robustness to prompt variations, we also consider an implicit version by replacing “is **{opposite sentiment}**” with “becomes the opposite of what it currently is”.

The “Do not explain the answer.” part of the prompt template in Figure 5.2 (and the other prompts) is there to prevent the model from providing additional details about why it produced the counterfactual. These details prolong inference time and cannot be validated using our self-consistency framework. Despite this instruction, some models will occasionally provide them regardless; if this happens, the extra explanation is removed.

The model output is fairly systematic, often prefixing the counterfactual paragraph with “Paragraph:” or providing a clear separation, making it reasonably easy to extract the counterfactual paragraph.

### 5.3.2 Feature attribution explanation

A common alternative to counterfactual explanations is feature attribution. These explanations highlight which input words are important for making a prediction. The faithfulness of these explanations can be evaluated using a self-consistency check, where the important words are redacted/masked [39, 41, 87]. Given a faithful explanation, it will be impossible for the model to perform the classification task [60, 171], assuming that the model understands the meaning of redaction/masking during classification. We discuss this and the classification setup in Section 5.3.4.

The model is only provided with the input prompt shown in Figure 5.3. The model response is a list of important words, and the matching words in the paragraph are replaced with “[REDACTED]”.

---

<sup>1</sup>For example, Llama2 frames the user message with special [INST] and [/INST] tokens [207].

Edit the following paragraph such that the sentiment is "{opposite sentiment}". Make as few edits as possible. Do not explain the answer.

Paragraph: {paragraph}

input prompt template

Paragraph: The movie was excellent ...

partial output example

Figure 5.2 The explicit input-template prompt used for generating the counterfactual explanation. `{opposite sentiment}` is replaced with either “positive” or “negative”. `{paragraph}` is replaced with the content. We also consider an implicit version where “is `{opposite sentiment}`” is replaced with “becomes the opposite of what it currently is”. The partial output example is entirely generated by the model.

List the most important words for determining the sentiment of the following paragraph, such that without these words the sentiment can not be determined. Do not explain the answer.

Paragraph: {paragraph}

input prompt template

Important words: "great," "amazing," ...

partial output example

Figure 5.3 The input-template prompt used for generating the feature attribution explanations. The model will often generate either a bullet-point list or a comma-separated list.

### 5.3.3 Redaction explanation

Redaction explanations are a less common variation of feature attribution explanations. Instead of asking the model to list the most important words and then algorithmically replacing those words with “[REDACTED]”, we ask the model to perform the replacement directly. The prompt-template is shown in Figure 5.4. As a prompt variation, we also use “[REMOVED]” instead of “[REDACTED]”.

Redaction may be easier for the model because the LLMs likely have built-in mechanisms for copying content [219]. In principle, the model could redact the entire paragraph, as we don’t constrain the redaction amount. This would be a faithful explanation but not a very useful explanation to humans [58], we also don’t observe such behavior in practice (see for example Appendix D).

Besides the different replacement mechanisms, the faithfulness metric works the same. The explanation is faithful if the model can not classify the redacted paragraph.

Redact the most important words for determining the sentiment of the following paragraph, by replacing important words with [REDACTED], such that without these words the sentiment can not be determined. Do not explain the answer.

Paragraph: {paragraph}

Paragraph: The movie was [REDACTED] ...

Figure 5.4 The input-template prompt used for generating redaction explanations. We also consider a prompt where “[REMOVED]” is used instead of “[REDACTED]”.

What is the sentiment of the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: {paragraph}

Positive

Figure 5.5 Prompt-template for classification. The prompt needs to support redaction and an “unknown” class for when the classification can not be performed due to missing information.

### 5.3.4 Classification

So far, we have discussed how to generate explanations. However, the self-consistency evaluation depends on a classification of the original paragraph and the explanation paragraph (or, in the case of feature attribution, it’s the paragraph modified using the explanation).

In Section 5.3.2 and Section 5.3.3, we use a “[REDACTED]” string to indicate that content is missing. We do this rather than removing content, as removing them creates ungrammatical content issues, which the models are not designed to support. This is similar to the out-of-distribution issue discussed in Chapter 3.

Finally, in the case of feature attribution (Section 5.3.2) and redaction explanations (Section 5.3.3), the paragraph of faithful explanation can not be classified. Therefore, the classification should allow for an “unknown” class prediction.

Importantly, the same prompt template is used in all cases (Figure 5.5); as in, for all explanations and both before and after the explanation step. The “unknown” and “[REDACTED]” support is not required for the counterfactual case but is kept for consistency and comparability.

### 5.3.5 Persona robustness

In the past sections, we presented some prompt modifications specific to each explanation. Inspired by Deshpande et al. [220], we propose the idea of using “persona” as a prompt

modification that can be applied to any prompt. In our setup, a “persona” means that the subject of the explanation or classification request is either “you” or “a human”.

For example, the previously presented prompts (e.g., Figure 5.5) asked the question (e.g., “What is the sentiment ...”) in an objective manner. Instead, it’s possible to ask “What would you classify the sentiment as” or “What would a human classify the sentiment as” – the exact prompts are provided in Appendix D. We hypothesize that personas could be relevant for the model’s classification and explanation [220]. For example, the “you” persona may be significant if the model has self-modeling capabilities [218].

## 5.4 Experiments

We perform all experiments with sentiment classification (IMDB), multi-choice Q&A tasks (bAbI-1 and MCTest), and an entailment/NLI task (RTE). These are all publicly available datasets, see Table 5.1. We chose these datasets to have diversity regarding how the tasks, inputs, and targets are represented. For example, sentiment classification and multi-choice Q&A are quite different in structure.

Type	Name	Test observations	explained content	reference	example
Sentiment	IMDB	25000	<b>text</b>	[26]	Appendix D
QA Multi-Choice	bAbI-1	1000	<b>paragraph</b>	[10]	Appendix D
	MCTest	600	<b>story</b>	[28]	
NLI	RTE	277	<b>sentence1</b>	[29]	Appendix D

Table 5.1 List of datasets used in Chapter 5. All datasets are publicly available.

Although the methods presented in Chapter 5 can be applied to any instruction-tuned generative language model, including API-only models like ChatGPT, we have limited the scope to only publicly available models without an indemnity clause. The motivations for this are to provide an impartial judgment and ensure reproducibility. We also did not analyze derived models that are fine-tuned versions of existing models; such analysis would add extra computing costs and is unlikely to provide valuable insights. The models analyzed (Llama 2, Falcon, and Mistral) are listed in Table 5.2 and are all intended for public consumption using a chat interface. Many of the models are or have been publically and freely available at <https://huggingface.co/chat> and have also provided their own web interface, which can be used to interact with these models (e.g., <https://huggingface.co/spaces/HuggingFaceH4/falcon-chat>). Due to the availability and accessibility of these models, analyzing the faithfulness of their self-explanations is paramount.

Name	size	HuggingFace repo	license	reference
Llama 2	70B	<code>meta-llama/Llama-2-70b-chat-hf</code>	Llama2 License	[207]
	7B	<code>meta-llama/Llama-2-7b-chat-hf</code>		
Falcon	40B	<code>tiiuae/falcon-40b-instruct</code>	Apache 2.0	[208]
	7B	<code>tiiuae/falcon-7b-instruct</code>		
Mistral	7B	<code>mistralai/Mistral-7B-Instruct-v0.1</code>	Apache 2.0	[209]

Table 5.2 List of models used in Chapter 5. All models are publicly available.

The prompts used for the experiments were developed using the training splits. The results shown in this section are all for the test splits. Furthermore, we have no reason to suspect the results are affected by the split.

The prompts for each dataset are slightly different; see Appendix D. The variability among the datasets demonstrates that our methodology generalizes to both paragraph-based multi-choice questions and two-paragraph tasks. In all cases, the main paragraph is modified through the explanation. The questions, choices, or hypotheses are not modified.

To evaluate prompt sensitivity, all the prompt variations from Section 5.3 (details in Appendix D) are evaluated using the Llama2-70B model fine-tuned for chatting [207]. In addition, we use a default case to show differences across model types and sizes as shown in Table 5.1.

**Default case:** In the experimental results and discussion, there are sometimes references to “default” prompt parameters. This is when the “[REDACTED]” token is used, counterfactuals use explicit targets, and all prompts use an objective persona.

**Inference generation details.** All generation inferences were made using Text Generation Inference (TGI) version 1.1.0 by HuggingFace (<https://github.com/huggingface/text-generation-inference>). The generation parameters are the same as those used for HuggingFace’s online chat service (<https://huggingface.co/chat/>) and thus accurately represent the conditions that a regular user would face. The only difference is that the generation has a pre-determined seed of 0 to allow for reproduction studies. The parameters are in Table 5.3. As the results are seeded, the responses are not affected much by the seed, and performing the inferences is expensive, we only perform a single inference sample.

Regarding the system prompt, as per current recommendations<sup>2</sup>, Llama2 and Mistral do not

<sup>2</sup><https://github.com/facebookresearch/llama/pull/626/commits/a971c41bde81d74f98bc2c2c451da235f1f1d37c>

use one. Falcon uses the default system prompt <sup>3</sup>. We have made the code used for generating all the prompts and results publicly available at <https://github.com/AndreasMadsen/11m-introspection>.

Parameter	Value
temperature	0.1
top_p	0.95
repetition_penalty	1.2
top_k	50
seed	0

Table 5.3 Generation parameters used for TGI. These are the same parameters used in HuggingFace’s online chat service (<https://huggingface.co/chat/>).

### 5.4.1 Classification

Before evaluating the faithfulness, it’s worth first investigating the classification accuracy for each task. In addition, this experiment also investigates the effect of persona (you, human, and objective) and the redaction-instruction ([REDACTED] or [REMOVED]). To validate that the redaction instruction does not cause issues, an ablation study with no redaction instruction (None) is also performed.

Figure 5.6 shows that neither the redaction-instruction nor the persona affects the results much. This is the desired result, as the redaction instruction should not affect the classification and supports just analyzing the default case will be sufficient.

Using the default case, Figure 5.7 shows the effect of different model types and sizes. There is quite a lot of variation between the different models. The most surprising result is that Falcon 40B performs slightly worse than Falcon 7B. Unfortunately, due to the complexity and lack of documentation regarding these models, it’s hard to make an educated guess on why this is.

Most task and model combinations do not perform well compared to regular fine-tuned models (e.g. when comparing with FMMs from Section 4.2.1). This could be problematic if, for

<sup>3</sup>“The following is a conversation between a highly knowledgeable and intelligent AI assistant, called Falcon, and a human user, called User. In the following interactions, User and Falcon will converse in natural language, and Falcon will answer User’s questions. Falcon was built to be respectful, polite and inclusive. Falcon was built by the Technology Innovation Institute in Abu Dhabi. Falcon will never decline to answer a question, and always attempts to give an answer that User would be satisfied with. It knows a lot, and always tells the truth. The conversation begins.” <https://huggingface.co/spaces/HuggingFaceH4/falcon-chat/blob/b20d83ddac4f79e772e3395621089d78804c166c/app.py>

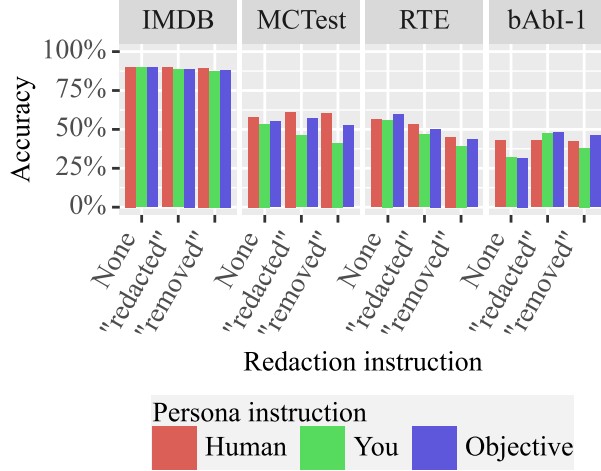


Figure 5.6 The classification accuracy of each task using Llama2-70B, with different prompt variations. Performance is not affected much by the persona or redaction-instruction.

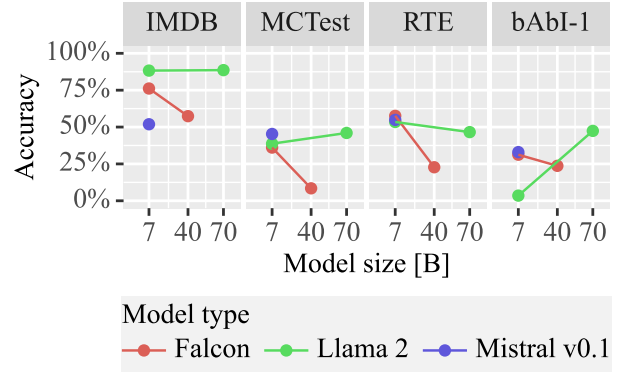


Figure 5.7 The classification accuracy of each task with the default prompt settings. Performance is highly dependent on model type and size.

example, an LLM classifies a positive-sentiment input as negative, and the counterfactual explanation asks for it to become negative, then it would appear that the explanation is faithful despite the explanation having made potentially no changes to the input. As the scope of this chapter (and the thesis) is faithfulness evaluation, not classification performance, we do not attempt to improve the classification performance. Instead, only the correctly predicted observations are used to evaluate faithfulness; the rest are discarded.

#### 5.4.2 Faithfulness

Using only the observations that are correctly classified and produce meaningful results (e.g., discarding “As an AI model I cannot do that.”), we evaluate the faithfulness of each observation. Because our self-consistency method determines whether or not an observation is faithful, faithfulness in this chapter refers to the ratio of faithful observations.

Figure 5.8 shows the faithfulness, for each prompt-variation for Llama2-70B. Figure 5.9 shows faithfulness with the default prompt settings for each model type and size.

### 5.5 Important findings

The overall conclusion from the faithfulness results in Figure 5.8 and Figure 5.9 is that the faithfulness is model-dependent, task-dependent, and explanation-dependent.

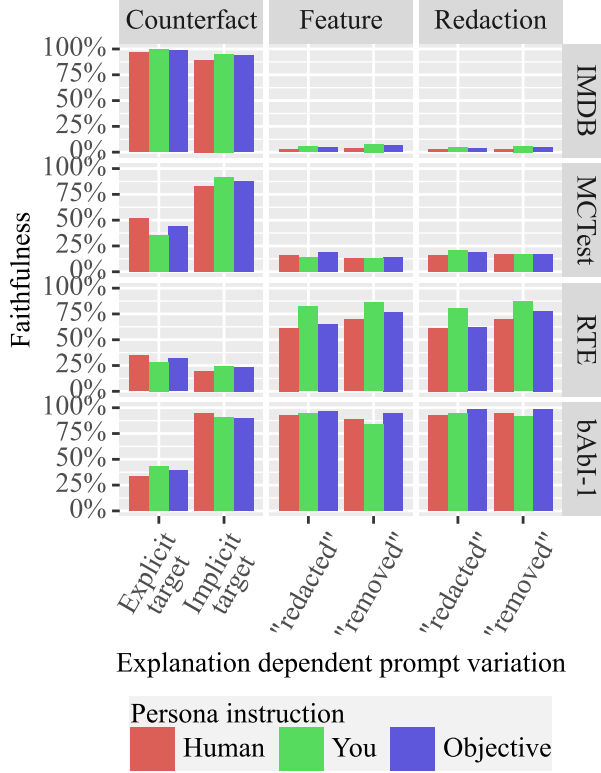


Figure 5.8 Faithfulness evaluation using self-consistency checks, evaluated using Llama2-70B. Results show that Llama2-70B is not affected by prompt variations, but the faithfulness for each explanation type is task-dependent.

## Counterfactual

From the prompt-variation results in Figure 5.8, we find that the persona has little effect. Making the counterfactual target implicit or explicit also does not affect faithfulness much. The exception here is for MCTest and bAbI-1, which is to be expected as these are multi-choice datasets thus for an implicit-target there are multiple correct answers, while there is only one correct answer using an explicit-target, thus a difference in performance is observed. As such, this is all positive, as the goal is to have models that are robust to prompt variations.

From the model-variation results in Figure 5.9, we find that Llama2 and Mistral are only consistently faithful for IMDB sentiment classification, while Falcon never performs well.

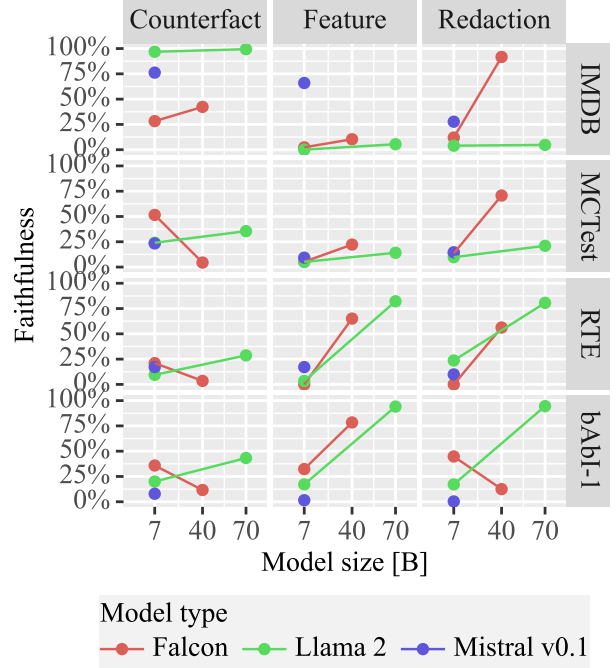


Figure 5.9 Faithfulness evaluation using self-consistency checks, evaluated with default prompt settings. Results show that faithfulness generally increases with size, the exception being Falcon, particularly in the counterfactual case. In general, the results are heavily dependent on the model type, tasks, and explanation.

## Feature attribution

The feature attribution experiments show again that Llama2-70B is robust to prompt variations (Figure 5.8), which is the desired outcome. However, we find that only for RTE and bAbI-1 is Llama2-70B faithful. The size particularly affects faithfulness, where Llama2-7B and Falcon-7B perform very poorly. Despite this, Mistral (7B) is surprisingly good at identifying important words for sentiment classification, drastically outperforming every other model (Figure 5.9).

While it might be surprising that the models generally perform poorly, it makes sense that bAbI-1 is an easy task, as there are, from a human perspective, only two relevant words for answering the questions in bAbI-1 (a location and individual).

## Redaction

The redaction experiments show results similar to the feature attribution experiment. This makes sense as the two explanations are similar in nature. This is particularly true for Llama2-70B, which have nearly identical results for the different prompt variations (Figure 5.8).

This similarity also addresses a potential concern with feature attribution explanations, which is that masking (i.e. the [REDACTED] token) inputs might create out-of-distribution results. However, for the redaction explanations, the masking is generated by the model, and should therefore be in-distribution. The observation that redaction and feature attribution have similar results, validates that masking does not cause any issues for the feature attribution explanations.

Finally, we do observe that Falcon-40B generally performs much better for this task compared to other models, given the IMDB and MCTest datasets. Also, surprisingly Falcon-40B performs worse than Falcon-7B for bAbI-1. This is particularly surprising as bAbI-1 is a synthetic dataset with few important words, hence it should be an easy task.

## 5.6 Related work

The self-consistency idea has previously been adopted for self-explanations. Huang et al. [215], applied the same idea to self-explanation feature attribution from ChatGPT on a sentiment classification task (a 100 observation subset of SST [173]). Because ChatGPT does not provide logits, they instruct ChatGPT to provide a confidence score between 1 and 0 for either positive or negative sentiment. Similarly, they instruct ChatGPT to provide scores for the importance of each word. Using this approach, previous methods can in principle be easily

applied. Unfortunately, they find their approach ineffective as the confidence and importance scores are not calibrated, and ChatGPT resists classifying as positive or negative sentiment when information is missing. For this reason, they call for better faithfulness metrics targeted LLM self-explanations.

In contrast, our approach does not depend on any scores (confidence or importance). Instead, the models perform discrete classification and can predict “unknown”. This solves the aforementioned issues identified by Huang et al. [215].

Self-consistency checks have also been used to measure the faithfulness of mathematical explanations. In particular, Lanham et al. [181] apply self-consistency checks to Chain-of-Thought (CoT) self-explanations from LLMs, where the CoTs provide mathematical context (e.g. “ $2 + 3 = 5$ ”) to a mathematical question. Lanham et al. [181] then edit the CoT to contain false information (e.g. “ $2 + 3 = 6$ ”) and check that the prediction follows. The issue here is that injecting false facts may create out-of-distribution results or be interpreted as typos by the LLM, thus it’s unclear if this method is completely valid. Regardless, they find similar to this chapter, that faithfulness is model and task-dependent.

### 5.6.1 Non-faithfulness works

Self-consistency checks also have been used to analyze other LLMs’ capabilities. For example, Anthropic Team [218] analyzed LLM self-modeling capabilities by comparing asking if the LLM knows the answer (Yes/No) and asking the LLM to answer directly. Li et al. [221] perform a similar task with mathematical questions (e.g., comparing “What is  $7 + 8$ ?” with “ $7+8=15$ , True or False”). Finally, Hu and Levy [222] evaluate meta-linguistic capabilities using self-consistency checks.

Outside of self-consistency checks, Chen et al. [214] measure the simulatability of counterfactual self-explanations. Simulatability means that humans can predict the model’s behavior on input  $\tilde{x}$ , given an explanation for a similar input  $x$ . Chen et al. [214] find that self-explanations are very convincing to humans (plausibility) but there is no correlation with simulatability. This shows that self-explanations from LLMs can be misleading [213]. These findings highlight the importance of providing faithful explanations. However, a challenge with simulatability is it’s difficult to ensure that humans use the explanations and not their world-knowledge to predict the model on input  $\tilde{x}$ . Therefore, it’s important to also measure faithfulness, as this is never evaluated by humans.

Finally, it’s worth briefly mentioning works on rationalizations. These are explanations by humans (e.g., e-SNLI [166]) or by a model that has learned from human explanations.

Rationalizations should not be used to explain the model [59] but may help convince humans of a prediction. As such, most works are on plausibility, but there are also simulatability studies [174] and faithfulness evaluations [177]. Although Parcalabescu and Frank [182] argue that the works on faithfulness only measure consistency, not faithfulness.

## 5.7 Limitations

### Absolute faithfulness

Jacovi and Goldberg [59] argue that the faithfulness field should move towards a more nuanced view on faithfulness, where the metric is not if an explanation is faithful or not, but how much faithfulness it has. We agree with this notion, but because the LLMs are designed to provide discrete outputs (a sequence of tokens), we have been unable to do that.

That being said, our absolute faithfulness is only for a single observation. We still provide an aggregation average, which provides a nuanced view of the model’s faithfulness as a whole.

### Assuming the prompt is interpreted correctly

Our work is based on a fundamental assumption that the model interprets the prompts correctly. This is not possible to verify completely. However, we attempt to ensure this by providing natural-sounding prompts. Hence, if the model doesn’t understand these prompts and generates unfaithful explanations or inaccurate classification, it’s a limitation of the model’s comprehension capabilities.

Additionally, we find that we get high classification accuracy or faithfulness for at least one dataset. This indicates that the model interprets these prompts sufficiently; it’s just not able to perform the requested task.

### No in-context learning

Because in-context learning has shown good results on many downstream tasks [223], it’s tempting to think the same is possible for generating explanations. Unfortunately, this is problematic as it would require known true explanations, which are not possible for humans to provide. Therefore, this chapter only uses zero-shot explanations.

## **Poor classification performance creates challenging comparison**

For most tasks we use, the LLMs perform poorly at the classification task and achieve low accuracy. As the goal of this chapter is to measure faithfulness, not accuracy, we simply discard the incorrect observations. However, this does mean that comparisons will be more challenging for future work. There may also be a class imbalance in what observations were selected.

## **Impossible to show lack of capability**

Our results demonstrate that the LLMs will, in some cases, consistently not provide faithful explanations. However, this does not show that the model is fundamentally unable to provide explanations. It only shows that the explanations are not faithful given the specific prompt templates. Because our work shows that generating faithful explanations is very challenging, users should not trust the explanations. Unfortunately, users often find these explanations very convincing [214]; we thus think this is problematic enough to warrant concern about AI Safety. However, no definitive evidence exists that the models can never explain themselves.

## **Limited number of explanation tasks**

We provide faithfulness results on 3 kinds of explanations: counterfactuals, importance measures, and redaction. Importance measures and redaction are also quite similar in the explanation they provide. As such, there is not a great variety of explanations in this chapter. However, these are currently the only explanations for which self-consistency can be used to measure faithfulness. We hope that future work can identify new ways of using self-consistency checks.

## **Not measuring on ChatGPT, GPT-4, etc.**

We do not evaluate ChatGPT, GPT-4, and other popular closed models. This is because such models are not static models. They are frequently updated, and we don't have the means to choose a specific model version. As such, it would not be possible to reproduce our results.

Additionally, these models use a license agreement with an unrestricted indemnity clause. As such, if we were to show that the explanations are not faithful, this could damage the companies behind these products, and we may become financially liable for damages. Due to this personal risk, there is a conflict of interest that may prevent a genuinely unbiased analysis.

## 5.8 Conclusion

Our investigation reveals that self-explanations' faithfulness is highly model-dependent, explanation-dependent, and dataset-dependent. This conclusion is similar to previous works [91, 181] and our conclusion in Chapter 3.

This chapter's primary contribution is developing the ability to measure the faithfulness of LLMs' self-explanations, specifically counterfactuals, feature attribution, and redaction explanations.

The task dependence is concerning as it means LLM self-explanations cannot generally be trusted. There is also no reason to trust more free-form explanations where faithfulness can not be evaluated using self-consistency checks. This increases the risk with LLMs, as individuals may have the misconception that LLMs can explain themselves [214].

## CHAPTER 6 CONCLUSION

A very essential but easily overlooked observation that becomes apparent from all the experiments and some existing literature is that faithfulness is, by default, model- and task-dependent. However, with just a small modification, faithfulness can become consistent.

The model- and task-dependent conclusion was observed in Chapter 3, where masking was used as an intervention to measure faithfulness of common post-hoc methods and attention. Bastings et al. [91] produced a simultaneous work that used a partially synthetic dataset to provide a known true explanation and found the same conclusion.

Chapter 5 then analyzed self-explanations, and again arrived at this model- and task-dependent conclusion, and even revealed that it holds across different types of explanations. Lanham et al. [181], another simultaneous work, looked at Chain-of-Thought (CoT) self-explanations by making interventions on the CoT and also arrived at a model- and task-dependent conclusion.

However, the conclusion was different in Chapter 4. Here, the results showed consistency in the explanations’ faithfulness. These differences could be caused by using a more accurate faithfulness metric. However, given that other works have also shown that faithfulness is model- and task-dependent, it’s more likely because masked fine-tuning regularizes the model so that those explanations become faithful. It’s understandable that the occlusion-based importance measures are consistently faithful as they use masking, which synergizes with the built-in masking support that masked fine-tuning provides. However, the same gradient-based importance measures used in Chapter 3 also became consistent.

The idea that some regularization can cause more faithful explanations is not new, as it has been proposed in other works [115, 198, 224, 225]. However, this hypothesis has mostly been hypothetical or anecdotal, without any systematic analysis. What is new is that data augmentation, such as masked fine-tuning, doesn’t just improve faithfulness but entirely solves the model- and task-dependent faithfulness issues.

The repeated observation that faithfulness is by default model- and task-dependent and the importance of regularization to combat this, means that post-hoc is unlikely to be a viable interpretability paradigm for developing general-purpose explanation methods. The work on faithfulness measurable models in Chapter 4 also shows that, at least for importance measures, the post-hoc worry about not affecting performance is not sufficient motivation to justify post-hoc as a paradigm, as there were no performance penalties on any of the 16 datasets.

For intrinsic methods, there is not as much analysis of task- and model-dependent faithfulness.

This is in part because, by definition, the explanation is tied to the model. Therefore, it’s impossible to compare the same explanation on different models. However, in Chapter 5, we also investigated attention and found its faithfulness to be task- and model-dependent, and the simultaneous work by Bastings et al. [91] produced the same conclusion.

This by-default model and task-dependent conclusion is likely what has behind the scenes caused much confusion and debate around faithfulness in the interpretability literature. It’s only through a very comprehensive analysis and carefully designed faithfulness metrics that this conclusion becomes apparent. Suppose a work is not extremely principled about the faithfulness metrics or only evaluates a short set of datasets or models, then it’s easy to conclude that an explanation method consistently either works or does not work. Thus, a comparative survey of the literature only reveals confusion and debate.

The idea that we can optimize or regularize models such that existing explanations become more faithful is already an emerging paradigm known as “Learn-to-faithfully-explain” [129]. However, the faithfulness measurable model paradigm takes the desirables much further by also answering how faithfulness is measured while making it cheap and precise by design. As discussed in Chapter 4, the properties of cheap and precise means that optimizing explanations towards maximum faithfulness is feasible. Although the beam-search-based optimization procedure used in Chapter 4 is not perfect, this possibility for this direction also answers how best to explain a model in terms of optimal faithfulness.

As such, faithfulness measurable models answer how to satisfy both of the beliefs from the intrinsic and post-hoc paradigms, *“only models that were designed to be explained can be explained”* and *“black-box models will be more generally applicable than intrinsic models”*. The first is because it’s possible to optimize an explanation, and the second is because it can be done without architectural constraints and provide the same general-purpose high performance as regular black-box models.

As such, faithfulness measurable models (FMMs) answer the research question of this thesis, “How to provide and ensure faithful explanations for complex general-purpose neural NLP models?” and did so by following the stated research hypothesis in Section 1.4: “by developing new paradigms that design models to be explained without employing architectural constraints, by focusing on developing accurate faithfulness metrics, by focusing on importance measures that have had a notoriously troubling history regarding faithfulness, and by taking advantage of properties specific to natural language and NLP models”.

Regarding self-explanation, as explored in Chapter 5, it may provide another new paradigm for providing faithful explanations. However, based on the analysis in Chapter 5 and other works [180, 181, 215, 226] that is not currently the case. However, we discuss in future works

Section 6.2.3 how this could become realized. As well as how the faithfulness measurable models paradigm could be applied to other models and explanations.

## 6.1 Limitations

In each chapter, careful attention has been paid to describing the limitations of the methodology. However, there is an overarching limitation to each chapter and the thesis as a whole, which has not been discussed. Namely, that only faithfulness is discussed. As interpretability is about explaining models to humans, it’s not only important that the explanation reflects the model’s behavior (i.e., faithfulness), but the explanation must also be useful to humans [58]. This desirable was mentioned in Section 2.3, but it’s worth mentioning again as this thesis does not address it.

The choice of restricting the considerations to faithfulness comes from the observation that the two most discussed paradigms (post-hoc and intrinsic) are about “what makes explanations faithful?”. Therefore, new alternative paradigms should also answer this question. How to best communicate said explanations and measure human-groundedness is an orthogonal matter [227, 228, 229, 230, 231, 232, 233]. Faithfulness is also a quasi-requisite to discuss human-groundedness because if the explanation is not true, to begin with, it’s hard to say if the way it is communicated is productive. Although human-groundedness could be discussed with simple intrinsically explainable models, it is unknown if lessons from these setups would generalize to complex models.

Finally, it’s worth pointing out that there are also attempts at developing new paradigms in the matter of human-groundedness. For example, Schut et al. [234] and Kim [235] propose the new idea that it is not enough to frame explanations in terms that humans already understand. We should also develop new language and mental abstractions for humans to understand machines.

## 6.2 Future Research

### 6.2.1 Masked Causal Language Models

In Chapter 4, it was shown how enabling masking of any permutation of tokens enabled a masked language model (MLM) to become faithfulness measurable and also made explanations more faithful. Given the popularity of causal language models (CLMs), especially chat systems, a natural next step would be to extend this capability to LLMs. In Section 4.4, it was briefly mentioned that this might be possible by injecting random masking tokens during pre-training

or potentially as an additional pre-training on an already pre-trained model while keeping the next-token-prediction objective the same. Additionally, to avoid catastrophic forgetting of the mask token while fine-tuning (e.g., instruction-tuning), the masking token would also need to be included here, similar to the proposed masked fine-tuning.

Pre-training large language models (LLMs) is quite expensive, so it’s worth considering what else we might gain from such masking support. One potential advantage is the ability to infer multiple tokens in parallel. Current CLMs need to infer one token at a time [207, 209], which makes inference slow. With masking support, it would be possible to append masking tokens to the input sequence, allowing for parallel inference. For example, the input tokens “The quick brown [MASK] [MASK] [MASK]” could map to “quick brown fox jumps over the”. The parallelly inferred tokens can then be validated in the next iteration, with the input “The quick brown fox jumps over the [MASK] [MASK]”. Ideas like this are already being applied but using additional models or layers, thus adding compute cost [236, 237, 238, 239]. In this case, the compute cost would be unchanged.

It would also be possible to apply standard techniques from masked language models to causal language models, such as using masking to investigate bias [240, 241]. For example, to understand the relationship between the occupation and pronoun in  $p(\textit{pronoun}|\text{The doctor washed})$ , one can perform an intervention on the occupation by replacing “doctor” with a mask token.

### 6.2.2 Faithfulness Measure Models for other kinds of communication

The faithfulness measurable model presented in Chapter 4 used masking support to provide an intrinsic faithfulness metric for importance measures and then optimize towards maximal faithfulness. While more could be done to improve the optimization or the types of models it is applied to; it’s worth thinking about the high-level idea of the paradigm and if it can be applied to other types of explanations, not just importance measures.

At the core, the high-level idea is that the question “how can models be designed to be explained” can be reformulated to “how can models be designed to be faithfulness measurable” while still allowing the model to be indirectly intrinsically interpretable. This reformulation may be more achievable, not require architectural constraints, and also answer how faithfulness is measured.

A relevant area could be concepts (Appendix A), where it’s explained to which extent concepts like “gender” are relevant for the prediction of a class (e.g., occupation) or the prediction related to a specific observation. Interventions of the intermediate representations are often used to measure the faithfulness of concept explanations. However, this has the same out-of-

distribution challenge as masking inputs. Thus, it may be productive to train a model to be robust against such interventions on the intermediate representations. Similar to how in Chapter 4 the model is robust against masking interventions on the input.

### 6.2.3 Preventing double-edged alignment

In Chapter 5, the results showed that self-explanations’ faithfulness is model- and task-dependent. It’s worth considering why this might be and how it can be solved.

One reasonable take is that because humans select the preferred utterances during alignment and humans can’t know how the model works [59], this leads to plausible self-explanations rather than faithful self-explanations [213]. This may become even more dangerous, as humans prefer fair-sounding explanations. As such, the self-explanations may hide potentially unfair model behaviors, a property known as fairwashing [242, 243].

Ensuring faithful self-explanations would prevent these risks. The model might still operate unfairly, but at least it will be transparent about it, leading to further investigation and research. Future work that develops instruction-tuned LLMs should thus also evaluate the faithfulness of self-explanation.

To align the model to provide faithful self-explanations, treating a faithfulness metric, such as the self-consistency check proposed in Chapter 5, as a reward function may be possible. Some work in this direction already exists, such as Anthropic Team [218], which showed that it’s possible to fine-tune a model to improve its self-modeling capabilities. As such, it stands to reason that improving self-explanation capabilities is also possible.

Finally, due to the difficulty in evaluating self-explanations, we suggest that self-explanation faithfulness should be treated as an out-of-domain evaluation problem. For example, one could fine-tune only for counterfactual faithfulness and show that feature attribution and redaction explanations also improve. This would give credit to the model’s explanation capabilities in general, including more free-formed explanations, which cannot be as easily evaluated using self-consistency checks.

Works in this direction may also help meditate other issues caused by alignment. For example, “The doctor washed his hands.” is a perfectly fair and unbiased sentence at an instance level, and thus a preference annotator would score it highly. However, aligning towards this increases the group-level bias that associates doctors with the male gender [244]. Work on broadening how reward functions in alignment are constructed may also help with such issues.

## REFERENCES

- [1] A. Madsen, “Visualizing memorization in RNNs,” *Distill*, vol. 4, no. 3, 3 2019. [Online]. Available: <https://distill.pub/2019/memorization-in-rnns>
- [2] S. Jain and B. C. Wallace, “Attention is not Explanation,” in *Proceedings of the 2019 Conference of the North*, vol. 1. Stroudsburg, PA, USA: Association for Computational Linguistics, 2 2019, pp. 3543–3556. [Online]. Available: <http://aclweb.org/anthology/N19-1357>
- [3] S. J. Reddi, S. Kale, and S. Kumar, “On the convergence of Adam and beyond,” *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, pp. 1–23, 4 2018. [Online]. Available: <http://arxiv.org/abs/1904.09237>
- [4] Y. Liu *et al.*, “RoBERTa: A Robustly Optimized BERT Pretraining Approach,” *arXiv*, 7 2019. [Online]. Available: <http://arxiv.org/abs/1907.11692>
- [5] T. Wolf *et al.*, “Transformers: State-of-the-Art Natural Language Processing,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Stroudsburg, PA, USA: Association for Computational Linguistics, 10 2020, pp. 38–45. [Online]. Available: <http://arxiv.org/abs/1910.03771><https://www.aclweb.org/anthology/2020.emnlp-demos.6>
- [6] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *7th International Conference on Learning Representations, ICLR 2019*, 2019.
- [7] A. Wang *et al.*, “GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding,” in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=rJ4km2R5t7>
- [8] —, “SuperGLUE: A stickier benchmark for general-purpose language understanding systems,” *Advances in Neural Information Processing Systems*, vol. 32, no. July, pp. 1–30, 2019.
- [9] A. E. Johnson *et al.*, “MIMIC-III, a freely accessible critical care database,” *Scientific Data*, vol. 3, no. 1, p. 160035, 12 2016. [Online]. Available: <http://www.nature.com/articles/sdata201635>
- [10] J. Weston *et al.*, “Towards AI-complete question answering: A set of prerequisite toy tasks,” *4th International Conference on Learning Representations*,

- ICLR 2016 - Conference Track Proceedings*, 2 2016. [Online]. Available: <http://arxiv.org/abs/1502.05698>
- [11] T. McCoy, E. Pavlick, and T. Linzen, “Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 3428–3448. [Online]. Available: <https://www.aclweb.org/anthology/P19-1334>
- [12] A. Williams, N. Nangia, and S. Bowman, “A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, vol. 1. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 1112–1122. [Online]. Available: <http://aclweb.org/anthology/N18-1101>
- [13] R. Socher *et al.*, “Parsing with compositional vector grammars,” in *ACL 2013 - 51st Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference*. Association for Computational Linguistics, 2013, vol. 1, pp. 455–465. [Online]. Available: <https://aclanthology.org/P13-1045/>
- [14] S. Hooker *et al.*, “A benchmark for interpretability methods in deep neural networks,” in *Advances in Neural Information Processing Systems*, vol. 32, 6 2019. [Online]. Available: <http://arxiv.org/abs/1806.10758>
- [15] M. T. Ribeiro, S. Singh, and C. Guestrin, “Semantically Equivalent Adversarial Rules for Debugging NLP models,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 856–865. [Online]. Available: <http://aclweb.org/anthology/P18-1079>
- [16] J. Vig *et al.*, “Investigating Gender Bias in Language Models Using Causal Mediation Analysis,” in *Advances in Neural Information Processing Systems*, H. Larochelle *et al.*, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 12 388–12 401. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/92650b2e92217715fe312e6fa7b90d82-Paper.pdf>
- [17] K. Pearson, “LIII. On lines and planes of closest fit to systems of points in space,” *The London, Edinburgh, and Dublin Philosophical Magazine and*

- Journal of Science*, vol. 2, no. 11, pp. 559–572, 11 1901. [Online]. Available: <https://www.tandfonline.com/doi/full/10.1080/14786440109462720>
- [18] L. Van Der Maaten and G. Hinton, “Visualizing data using t-SNE,” *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008. [Online]. Available: <https://www.jmlr.org/papers/v9/vandermaaten08a.html>
- [19] J. Pennington, R. Socher, and C. Manning, “Glove: Global Vectors for Word Representation,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2014, pp. 1532–1543. [Online]. Available: <http://aclweb.org/anthology/D14-1162>
- [20] I. Tenney, D. Das, and E. Pavlick, “BERT Rediscovered the Classical NLP Pipeline,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 5 2019, pp. 4593–4601. [Online]. Available: <https://www.aclweb.org/anthology/P19-1452>
- [21] J. Devlin *et al.*, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, vol. 1. Association for Computational Linguistics (ACL), 10 2019, pp. 4171–4186. [Online]. Available: <http://arxiv.org/abs/1810.04805>
- [22] C. Clark *et al.*, “Boolq: Exploring the surprising difficulty of natural yes/no questions,” in *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, vol. 1, 2019, pp. 2924–2936.
- [23] M.-C. d. Marneffe, M. Simons, and J. Tonhauser, “The CommitmentBank: Investigating projection in naturally occurring discourse,” *Proceedings of Sinn und Bedeutung*, vol. 23, no. 2, pp. 107–124, 2019. [Online]. Available: <https://ojs.ub.uni-konstanz.de/sub/index.php/sub/article/view/601>
- [24] A. Warstadt, A. Singh, and S. R. Bowman, “Neural Network Acceptability Judgments,” *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 625–641, 11 2019. [Online]. Available: <https://direct.mit.edu/tacl/article/43528>
- [25] A. Talmor *et al.*, “CommonsenseQA: A Question Answering Challenge Targeting Commonsense Knowledge,” in *Proceedings of the 2019 Conference of the North*.

- Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 4149–4158. [Online]. Available: <http://aclweb.org/anthology/N19-1421>
- [26] A. L. Maas *et al.*, “Learning word vectors for sentiment analysis,” in *ACL-HLT 2011 - Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, vol. 1. Portland, Oregon, USA: Association for Computational Linguistics, 6 2011, pp. 142–150. [Online]. Available: <https://www.aclweb.org/anthology/P11-1015>
- [27] W. B. Dolan and C. Brockett, “Automatically Constructing a Corpus of Sentential Paraphrases,” in *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005, pp. 9–16. [Online]. Available: <https://research.microsoft.com/apps/pubs/default.aspx?id=101076>
- [28] M. Richardson, C. J. Burges, and E. Renshaw, “MCTest: A challenge dataset for the open-domain machine comprehension of text,” *EMNLP 2013 - 2013 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*, vol. D13-1020, no. October, pp. 193–203, 2013.
- [29] I. Dagan, O. Glickman, and B. Magnini, “The PASCAL Recognising Textual Entailment Challenge,” in *Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Textual Entailment*, J. Quiñero-Candela *et al.*, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 177–190. [Online]. Available: [http://link.springer.com/10.1007/11736790\\_9](http://link.springer.com/10.1007/11736790_9)
- [30] S. R. Bowman *et al.*, “A large annotated corpus for learning natural language inference,” in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2015, pp. 632–642. [Online]. Available: <http://aclweb.org/anthology/D15-1075>
- [31] P. Rajpurkar *et al.*, “SQuad: 100,000+ questions for machine comprehension of text,” *EMNLP 2016 - Conference on Empirical Methods in Natural Language Processing, Proceedings*, pp. 2383–2392, 2016.
- [32] S. Iyer, N. Dandekar, and K. Csernai, “First Quora Dataset Release: Question Pairs,” 2017. [Online]. Available: <https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs>
- [33] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 11 1997. [Online]. Available: <https://www.mitpressjournals.org/doi/abs/10.1162/neco.1997.9.8.1735>

- [34] C. Raffel *et al.*, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research*, vol. 21, pp. 1–67, 2020. [Online]. Available: <https://jmlr.org/papers/v21/20-074.html>
- [35] Y. Zhou and J. Shah, “The Solvability of Interpretability Evaluation Metrics,” in *Findings of the Association for Computational Linguistics: EACL*, 2023. [Online]. Available: <http://arxiv.org/abs/2205.08696>
- [36] N. F. Rajani *et al.*, “Explain Yourself! Leveraging Language Models for Commonsense Reasoning,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 4932–4942. [Online]. Available: <https://www.aclweb.org/anthology/P19-1487>
- [37] D. Baehrens *et al.*, “How to explain individual classification decisions,” *Journal of Machine Learning Research*, vol. 11, pp. 1803–1831, 12 2010. [Online]. Available: <http://arxiv.org/abs/0912.1128>
- [38] J. Li *et al.*, “Visualizing and Understanding Neural Models in NLP,” in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2016, pp. 681–691. [Online]. Available: <http://aclweb.org/anthology/N16-1082>
- [39] M. Sundararajan, A. Taly, and Q. Yan, “Axiomatic attribution for deep networks,” in *34th International Conference on Machine Learning, ICML 2017*, vol. 7, 3 2017, pp. 5109–5118. [Online]. Available: <http://arxiv.org/abs/1703.01365>
- [40] J. Li, W. Monroe, and D. Jurafsky, “Understanding Neural Networks through Representation Erasure,” *arXiv*, 2016. [Online]. Available: <http://arxiv.org/abs/1612.08220>
- [41] M. T. Ribeiro, S. Singh, and C. Guestrin, ““Why should i trust you?” Explaining the predictions of any classifier,” in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, vol. 13-17-Aug. New York, NY, USA: ACM, 8 2016, pp. 1135–1144. [Online]. Available: <https://dl.acm.org/doi/10.1145/2939672.2939778>
- [42] A. Ross, A. Marasović, and M. Peters, “Explaining NLP Models via Minimal Contrastive Editing (MiCE),” in *Findings of the Association for Computational Linguistics: ACL-*

- IJCNLP 2021*. Stroudsburg, PA, USA: Association for Computational Linguistics, 12 2021, pp. 3840–3852. [Online]. Available: <https://aclanthology.org/2021.findings-acl.336>
- [43] S. Kumar and P. Talukdar, “NILE : Natural Language Inference with Faithful Natural Language Explanations,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 5 2020, pp. 8730–8742. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.771>
- [44] S. Lundberg and S.-I. Lee, “A Unified Approach to Interpreting Model Predictions,” in *Advances in Neural Information Processing Systems*, 5 2017, pp. 4766–4775. [Online]. Available: <http://arxiv.org/abs/1705.07874>
- [45] J. Adebayo *et al.*, “Sanity checks for saliency maps,” in *Advances in Neural Information Processing Systems*, vol. 2018-Decem. Curran Associates, Inc., 10 2018, pp. 9505–9515. [Online]. Available: <http://arxiv.org/abs/1810.03292>
- [46] H. Matan *et al.*, “A Statistical Framework for Efficient Out of Distribution Detection in Deep Neural Networks,” *International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=Oy9WeuZD51>
- [47] U. Bhatt *et al.*, “Explainable Machine Learning in Deployment,” *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 648–657, 9 2019. [Online]. Available: <https://dl.acm.org/doi/10.1145/3351095.3375624>
- [48] C. Rudin, “Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead,” *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019. [Online]. Available: <http://www.nature.com/articles/s42256-019-0048-x>
- [49] R. Wexler, “When a computer program keeps you in jail: How computers are harming criminal justice,” 2017. [Online]. Available: <https://www.nytimes.com/2017/06/13/opinion/how-computers-are-harming-criminal-justice.html>
- [50] M. McGough, “How bad is Sacramento’s air, exactly? Google results appear at odds with reality, some say,” 2018. [Online]. Available: <https://www.sacbee.com/news/california/fires/article216227775.html>
- [51] K. R. Varshney and H. Alemzadeh, “On the Safety of Machine Learning: Cyber-Physical Systems, Decision Sciences, and Data Products,” *Big Data*, vol. 5, no. 3, pp. 246–255, 9 2017. [Online]. Available: <http://www.ncbi.nlm.nih.gov/pubmed/28933947>

- [52] F. Doshi-Velez *et al.*, “Accountability of AI Under the Law: The Role of Explanation,” *SSRN Electronic Journal*, vol. Online, 11 2017. [Online]. Available: <https://www.ssrn.com/abstract=3064761>
- [53] Z. Obermeyer *et al.*, “Dissecting racial bias in an algorithm used to manage the health of populations,” *Science*, vol. 366, no. 6464, pp. 447–453, 10 2019. [Online]. Available: <https://science.sciencemag.org/content/366/6464/447>
- [54] T. B. Brown *et al.*, “Language Models are Few-Shot Learners,” in *Advances in Neural Information Processing Systems*, H. Larochelle *et al.*, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf>
- [55] E. M. Bender *et al.*, “On the Dangers of Stochastic Parrots,” in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. New York, NY, USA: ACM, 3 2021, pp. 610–623. [Online]. Available: <https://dl.acm.org/doi/10.1145/3442188.3445922>
- [56] N. Mehrabi *et al.*, “A Survey on Bias and Fairness in Machine Learning,” *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–35, 2021.
- [57] I. Garrido-Muñoz *et al.*, “A Survey on Bias in Deep NLP,” *Applied Sciences*, vol. 11, no. 7, p. 3184, 4 2021. [Online]. Available: <https://www.mdpi.com/2076-3417/11/7/3184>
- [58] F. Doshi-Velez and B. Kim, “Towards A Rigorous Science of Interpretable Machine Learning,” *arXiv*, 2 2017. [Online]. Available: <http://arxiv.org/abs/1702.08608>
- [59] A. Jacovi and Y. Goldberg, “Towards Faithfully Interpretable NLP Systems: How Should We Define and Evaluate Faithfulness?” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 4 2020, pp. 4198–4205. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.386>
- [60] W. Samek *et al.*, “Evaluating the Visualization of What a Deep Neural Network Has Learned,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 11, pp. 2660–2673, 11 2017. [Online]. Available: <https://ieeexplore.ieee.org/document/7552539/>
- [61] Z. C. Lipton, “The mythos of model interpretability,” *Communications of the ACM*, vol. 61, no. 10, pp. 36–43, 9 2018. [Online]. Available: <https://dl.acm.org/doi/10.1145/3233231>

- [62] T. S. Kuhn, *The Structure of Scientific Revolutions*, 3rd ed. University of Chicago Press, 1996.
- [63] S. Barocas, M. Hardt, and A. Narayanan, *Fairness and Machine Learning: Limitations and Opportunities*. fairmlbook.org, 2019. [Online]. Available: <https://fairmlbook.org/>
- [64] A. Xiang and I. D. Raji, “On the Legal Compatibility of Fairness Definitions,” *Workshop on Human-Centric Machine Learning at the 33rd Conference on Neural Information Processing Systems*, 2019. [Online]. Available: <http://arxiv.org/abs/1912.00761>
- [65] M. Andrus *et al.*, “What we can’t measure, We can’t understand: Challenges to demographic data procurement in the pursuit of fairness,” *FAccT 2021 - Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pp. 249–260, 2021.
- [66] A. A. Kodiyan, “An overview of ethical issues in using AI systems in hiring with a case study of Amazon’s AI based hiring tool,” *Researchgate Preprint*, pp. 1–19, 2019.
- [67] J. B. Fuller *et al.*, “Hidden Workers: Untapped Talent,” *Harvard Business School Project on Managing the Future of Work and Accenture*, 2021. [Online]. Available: <https://www.pw.hks.harvard.edu/post/hidden-workers-untapped-talent>
- [68] J. Fuller, “Companies Need More Workers. Why Do They Reject Millions of Résumés?” *The project on workforce*, 2021. [Online]. Available: <https://www.pw.hks.harvard.edu/post/companies-need-more-workers-wsj>
- [69] K. Preuer *et al.*, “Interpretable deep learning in drug discovery,” *Explainable AI: interpreting, explaining and visualizing deep learning*, pp. 331–345, 2019.
- [70] J. Jiménez-Luna, F. Grisoni, and G. Schneider, “Drug discovery with explainable artificial intelligence,” *Nature Machine Intelligence*, vol. 2, no. 10, pp. 573–584, 2020.
- [71] S. Dara *et al.*, “Machine learning in drug discovery: a review,” *Artificial Intelligence Review*, vol. 55, no. 3, pp. 1947–1999, 2022.
- [72] N. Cammarata *et al.*, “Thread: Circuits,” *Distill*, vol. 5, no. 3, 3 2020. [Online]. Available: <https://distill.pub/2020/circuits>
- [73] N. Elhage *et al.*, “A Mathematical Framework for Transformer Circuits,” *Anthropic*, 2021. [Online]. Available: <https://transformer-circuits.pub/2021/framework/index.html>

- [74] U. G. House of Lords, “AI in the UK: Ready, Willing and Able?” 2017. [Online]. Available: <https://publications.parliament.uk/pa/ld201719/ldselect/ldai/100/10007.htm>
- [75] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso, “Machine Learning Interpretability: A Survey on Methods and Metrics,” *Electronics*, vol. 8, no. 8, p. 832, 7 2019. [Online]. Available: <https://www.mdpi.com/2079-9292/8/8/832>
- [76] M. Flora *et al.*, “Comparing Explanation Methods for Traditional Machine Learning Models Part 1: An Overview of Current Methods and Quantifying Their Disagreement,” *arXiv*, pp. 1–22, 2022. [Online]. Available: <http://arxiv.org/abs/2211.08943>
- [77] V. Arya *et al.*, “One Explanation Does Not Fit All: A Toolkit and Taxonomy of AI Explainability Techniques,” *arXiv*, 9 2019. [Online]. Available: <http://arxiv.org/abs/1909.03012>
- [78] W. J. Murdoch *et al.*, “Definitions, methods, and applications in interpretable machine learning,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 116, no. 44, pp. 22 071–22 080, 10 2019. [Online]. Available: <http://www.pnas.org/lookup/doi/10.1073/pnas.1900654116>
- [79] D. Bahdanau, K. H. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” in *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*. International Conference on Learning Representations, ICLR, 9 2015, pp. 1–15. [Online]. Available: <https://arxiv.org/abs/1409.0473>
- [80] J. Andreas *et al.*, “Neural Module Networks,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 6 2016, pp. 39–48. [Online]. Available: <http://ieeexplore.ieee.org/document/7780381/>
- [81] N. Gupta *et al.*, “Neural Module Networks for Reasoning over Text,” in *International Conference on Learning Representations (ICLR)*, 12 2020. [Online]. Available: <https://openreview.net/forum?id=SygWvAVFPr>
- [82] H. Fashandi, “Neural module networks: A review,” *Neurocomputing*, vol. 552, p. 126518, 2023. [Online]. Available: <https://doi.org/10.1016/j.neucom.2023.126518>
- [83] J. Bien and R. Tibshirani, “Classification by Set Cover: The Prototype Vector Machine,” *arXiv*, pp. 1–24, 2009. [Online]. Available: <http://arxiv.org/abs/0908.2284>

- [84] B. Kim, C. Rudin, and J. Shah, “The Bayesian case model: A generative approach for case-based reasoning and prototype classification,” *Advances in Neural Information Processing Systems*, vol. 3, no. January, pp. 1952–1960, 2014.
- [85] C. Chen *et al.*, “This looks like that: Deep learning for interpretable image recognition,” *Advances in Neural Information Processing Systems*, vol. 32, 6 2019. [Online]. Available: <http://arxiv.org/abs/1806.10574>
- [86] J. Seo *et al.*, “Noise-adding Methods of Saliency Map as Series of Higher Order Partial Derivative,” in *2018 ICML Workshop on Human Interpretability in Machine Learning*, 6 2018. [Online]. Available: <http://arxiv.org/abs/1806.03000>
- [87] A. Karpathy, J. Johnson, and L. Fei-Fei, “Visualizing and Understanding Recurrent Networks,” *arXiv*, pp. 1–12, 6 2015. [Online]. Available: <http://arxiv.org/abs/1506.02078>
- [88] S. Serrano and N. A. Smith, “Is Attention Interpretable?” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 6 2019, pp. 2931–2951. [Online]. Available: <https://www.aclweb.org/anthology/P19-1282>
- [89] S. Vashishth *et al.*, “Attention Interpretability Across NLP Tasks,” *arXiv*, 9 2019. [Online]. Available: <http://arxiv.org/abs/1909.11218>
- [90] C. Meister *et al.*, “Is Sparse Attention more Interpretable?” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 8 2021, pp. 122–129. [Online]. Available: <http://arxiv.org/abs/2106.01087><https://aclanthology.org/2021.acl-short.17>
- [91] J. Bastings *et al.*, ““Will You Find These Shortcuts?” A Protocol for Evaluating the Faithfulness of Input Saliency Methods for Text Classification,” in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2022, pp. 976–991. [Online]. Available: <https://aclanthology.org/2022.emnlp-main.64>
- [92] DARPA, “Explainable Artificial Intelligence (XAI) DARPA-BAA-16-53,” *Defense Advanced Research Projects Agency (DARPA)*, pp. 1–52, 2016. [Online]. Available: <https://www.darpa.mil/attachments/DARPA-BAA-16-53.pdf>

- [93] B. Goodman and S. Flaxman, "European union regulations on algorithmic decision making and a "right to explanation",” *AI Magazine*, vol. 38, no. 3, pp. 50–57, 2017.
- [94] S. Krishna *et al.*, "The Disagreement Problem in Explainable Machine Learning: A Practitioner’s Perspective,” *arXiv*, 2022. [Online]. Available: <http://arxiv.org/abs/2202.01602>
- [95] J. Bastings and K. Filippova, "The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?” in *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 149–155. [Online]. Available: <https://www.aclweb.org/anthology/2020.blackboxnlp-1.14>
- [96] M. Amer and T. Maul, "A review of modularization techniques in artificial neural networks,” *Artificial Intelligence Review*, vol. 52, no. 1, pp. 527–561, 6 2019. [Online]. Available: <http://link.springer.com/10.1007/s10462-019-09706-7>
- [97] S. Subramanian *et al.*, "Obtaining faithful interpretations from compositional neural networks,” *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pp. 5594–5608, 2020. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.495>
- [98] Q. Lyu, M. Apidianaki, and C. Callison-Burch, "Towards Faithful Model Explanation in NLP: A Survey,” *Computational Linguistics*, vol. 50, no. 2, pp. 657–723, 6 2024. [Online]. Available: <http://arxiv.org/abs/2209.11326><https://direct.mit.edu/coli/article/50/2/657/119158/Towards-Faithful-Model-Explanation-in-NLP-A-Survey>
- [99] A. Binder *et al.*, "Layer-Wise Relevance Propagation for Neural Networks with Local Renormalization Layers,” in *Artificial Neural Networks and Machine Learning – ICANN 2016*, vol. 9887 LNCS, 2016, pp. 63–71. [Online]. Available: [http://link.springer.com/10.1007/978-3-319-44781-0\\_8](http://link.springer.com/10.1007/978-3-319-44781-0_8)
- [100] A. Shrikumar, P. Greenside, and A. Kundaje, "Learning important features through propagating activation differences,” in *34th International Conference on Machine Learning, ICML 2017*, vol. 7, 2017, pp. 4844–4866. [Online]. Available: <https://arxiv.org/>
- [101] D. Smilkov *et al.*, "SmoothGrad: removing noise by adding noise,” *ICML workshop on visualization for deep learning*, 2017. [Online]. Available: <https://goo.gl/EfVzEE>.

- [102] I. Ahern *et al.*, “Normlime: A new feature importance metric for explaining deep neural networks,” *arXiv*, 9 2019. [Online]. Available: <http://arxiv.org/abs/1909.04200>
- [103] J. Thorne *et al.*, “Generating Token-Level Explanations for Natural Language Inference,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, vol. 1. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 963–969. [Online]. Available: <http://aclweb.org/anthology/N19-1101>
- [104] R. ElShawi *et al.*, “ILIME: Local and Global Interpretable Model-Agnostic Explainer of Black-Box Decision,” in *Advances in Databases and Information Systems*, T. Welzer *et al.*, Eds. Cham: Springer International Publishing, 2019, pp. 53–68. [Online]. Available: [http://link.springer.com/10.1007/978-3-030-28730-6\\_4](http://link.springer.com/10.1007/978-3-030-28730-6_4)
- [105] A. Sangroya *et al.*, “Guided-LIME: Structured sampling based hybrid approach towards explaining blackbox machine learning models,” in *CEUR Workshop Proceedings*, vol. 2699, 2020.
- [106] J. Adebayo *et al.*, “Post hoc Explanations may be Ineffective for Detecting Unknown Spurious Correlation,” in *International Conference on Learning Representations*, 2021, pp. 1–13. [Online]. Available: <https://openreview.net/forum?id=xNOVfCCvDpM>
- [107] P.-J. Kindermans *et al.*, “The (Un)reliability of Saliency Methods,” in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. Springer, 11 2019, vol. 11700 LNCS, pp. 267–280. [Online]. Available: [http://link.springer.com/10.1007/978-3-030-28954-6\\_14](http://link.springer.com/10.1007/978-3-030-28954-6_14)
- [108] D. Slack *et al.*, “Fooling LIME and SHAP,” in *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*. New York, NY, USA: ACM, 2 2020, pp. 180–186. [Online]. Available: <https://dl.acm.org/doi/10.1145/3375627.3375830>
- [109] C.-K. Yeh *et al.*, “On the (In)fidelity and Sensitivity of Explanations,” in *Advances in Neural Information Processing Systems 32*, H. Wallach *et al.*, Eds. Vancouver, Canada: Curran Associates, Inc., 2019, pp. 10 967–10 978. [Online]. Available: <https://arxiv.org/abs/1901.09392>
- [110] T. Han, S. Srinivas, and H. Lakkaraju, “Which Explanation Should I Choose? A Function Approximation Perspective to Characterizing Post Hoc Explanations,” *Advances in Neural Information Processing Systems*, vol. 35, no. NeurIPS, 2022. [Online]. Available: <http://arxiv.org/abs/2206.01254>

- [111] B. Bilodeau *et al.*, “Impossibility theorems for feature attribution,” *Proceedings of the National Academy of Sciences*, vol. 121, no. 2, pp. 1–38, 1 2024. [Online]. Available: <https://pnas.org/doi/10.1073/pnas.2304406120><http://arxiv.org/abs/2212.11870>
- [112] C. Olah, A. Mordvintsev, and L. Schubert, “Feature Visualization,” *Distill*, vol. 2, no. 11, 11 2017. [Online]. Available: <https://distill.pub/2017/feature-visualization>
- [113] A. Nguyen, J. Yosinski, and J. Clune, “Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks,” *Visualization for Deep Learning workshop at ICML*, 2016. [Online]. Available: <http://arxiv.org/abs/1602.03616>
- [114] J. Yosinski *et al.*, “Understanding Neural Networks Through Deep Visualization,” in *Deep Learning Workshop at 31st International Conference on Machine Learning*, 2015. [Online]. Available: <http://arxiv.org/abs/1506.06579>
- [115] R. Geirhos *et al.*, “Don’t trust your eyes: on the (un)reliability of feature visualizations,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2306.04719>
- [116] J. Borowski *et al.*, “Exemplary Natural Images Explain Cnn Activations Better Than State-of-the-Art Feature Visualization,” *ICLR 2021 - 9th International Conference on Learning Representations*, pp. 1–41, 2021.
- [117] R. S. Zimmermann *et al.*, “How Well do Feature Visualizations Support Causal Understanding of CNN Activations?” *Advances in Neural Information Processing Systems*, vol. 14, no. NeurIPS, pp. 11 730–11 744, 2021.
- [118] Y. Belinkov and J. Glass, “Analysis Methods in Neural Language Processing: A Survey,” *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 49–72, 4 2019. [Online]. Available: [https://doi.org/10.1162/tacl\\_a\\_00254](https://doi.org/10.1162/tacl_a_00254)
- [119] Y. Belinkov, S. Gehrmann, and E. Pavlick, “Interpretability and Analysis in Neural NLP,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 1–5. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-tutorials.1>
- [120] A. Rogers, O. Kovaleva, and A. Rumshisky, “A Primer in BERTology: What We Know About How BERT Works,” *Transactions of the Association for Computational Linguistics*, vol. 8, pp. 842–866, 12 2020. [Online]. Available: <https://direct.mit.edu/tacl/article/96482>

- [121] A. Coenen *et al.*, “Visualizing and Measuring the Geometry of BERT,” in *Advances in Neural Information Processing Systems*, H. Wallach *et al.*, Eds., vol. 32. Curran Associates, Inc., 6 2019, pp. 8594–8603. [Online]. Available: <https://proceedings.neurips.cc/paper/2019/file/159c1ffe5b61b41b3c4d8f4c2150f6c4-Paper.pdf>
- [122] K. Clark *et al.*, “What Does BERT Look at? An Analysis of BERT’s Attention,” in *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 276–286. [Online]. Available: <https://www.aclweb.org/anthology/W19-4828>
- [123] L. Clouatre *et al.*, “Local Structure Matters Most: Perturbation Study in NLU,” in *Findings of the Association for Computational Linguistics: ACL 2022*. Stroudsburg, PA, USA: Association for Computational Linguistics, 7 2022, pp. 3712–3731. [Online]. Available: <https://aclanthology.org/2022.findings-acl.293>
- [124] A. Conneau *et al.*, “What you can cram into a single  $\$ \& ! \# *$  vector: Probing sentence embeddings for linguistic properties,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 2126–2136. [Online]. Available: <http://aclweb.org/anthology/P18-1198>
- [125] Y. Belinkov, “Probing Classifiers: Promises, Shortcomings, and Advances,” *arXiv*, pp. 1–12, 2 2021. [Online]. Available: <http://arxiv.org/abs/2102.12452>
- [126] K. Zhang and S. Bowman, “Language Modeling Teaches You More than Translation Does: Lessons Learned Through Auxiliary Syntactic Task Analysis,” in *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 359–361. [Online]. Available: <http://aclweb.org/anthology/W18-5448>
- [127] J. Hewitt and P. Liang, “Designing and Interpreting Probes with Control Tasks,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 2733–2743. [Online]. Available: <https://www.aclweb.org/anthology/D19-1275>
- [128] E. Voita and I. Titov, “Information-Theoretic Probing with Minimum Description Length,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Stroudsburg, PA, USA: Association

- for Computational Linguistics, 3 2020, pp. 183–196. [Online]. Available: <https://www.aclweb.org/anthology/2020.emnlp-main.14>
- [129] A. Madsen *et al.*, “Interpretability Needs a New Paradigm,” *arXiv*, 5 2024. [Online]. Available: <http://arxiv.org/abs/2405.05386>
- [130] A. Madsen, S. Reddy, and S. Chandar, “Post-hoc Interpretability for Neural NLP: A Survey,” *ACM Computing Surveys*, vol. 55, no. 8, pp. 1–42, 8 2022. [Online]. Available: <https://dl.acm.org/doi/10.1145/3546577>
- [131] A. Madsen *et al.*, “Evaluating the Faithfulness of Importance Measures in NLP by Recursively Masking Allegedly Important Tokens and Retraining,” in *Findings of the Association for Computational Linguistics: EMNLP 2022*. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, 12 2022, pp. 1731–1751. [Online]. Available: <https://aclanthology.org/2022.findings-emnlp.125>
- [132] —, “Faithfulness Measurable Masked Language Models,” in *Forty-first International Conference on Machine Learning*, 2024. [Online]. Available: <https://openreview.net/forum?id=tw1PwpuAuNhttp://arxiv.org/abs/2310.07819>
- [133] A. Madsen, S. Chandar, and S. Reddy, “Are self-explanations from Large Language Models faithful?” *The 62nd Annual Meeting of the Association for Computational Linguistics*, 1 2024. [Online]. Available: <https://openreview.net/forum?id=0fB5OROAiq>
- [134] T. Miller, “Explanation in artificial intelligence: Insights from the social sciences,” *Artificial Intelligence*, vol. 267, pp. 1–38, 2 2019. [Online]. Available: <http://arxiv.org/abs/1706.07269https://linkinghub.elsevier.com/retrieve/pii/S0004370218305988>
- [135] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [136] A. Vaswani *et al.*, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, vol. 2017-Decem. Association for Computational Linguistics (ACL), 6 2017, pp. 5999–6009. [Online]. Available: <http://arxiv.org/abs/1706.03762>
- [137] A. Graves, *Supervised Sequence Labelling with Recurrent Neural Networks*, ser. Studies in Computational Intelligence. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, vol. 385. [Online]. Available: <https://link.springer.com/10.1007/978-3-642-24797-2>
- [138] D. Jurafsky and J. Martin, “Speech and Language Processing,” *Speech and Language Processing.*, vol. 3, pp. 441–458, 2014.

- [139] A. Adadi and M. Berrada, “Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI),” *IEEE Access*, vol. 6, pp. 52 138–52 160, 2018. [Online]. Available: <https://ieeexplore.ieee.org/document/8466590/>
- [140] C. Molnar, *Interpretable Machine Learning*. Independent, 2019. [Online]. Available: <https://christophm.github.io/interpretable-ml-book/>
- [141] A. Chatzimparmpas *et al.*, “The State of the Art in Enhancing Trust in Machine Learning Models with the Use of Visualizations,” *Computer Graphics Forum*, vol. 39, no. 3, pp. 713–756, 6 2020. [Online]. Available: <https://onlinelibrary.wiley.com/doi/10.1111/cgf.14034>
- [142] S. Wiegrefe and Y. Pinter, “Attention is not not Explanation,” *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 11–20, 8 2019. [Online]. Available: <https://www.aclweb.org/anthology/D19-1002>
- [143] J. J. Williams *et al.*, “AXIS: Generating Explanations at Scale with Learnersourcing and Machine Learning,” in *Proceedings of the Third (2016) ACM Conference on Learning @ Scale*. New York, NY, USA: ACM, 4 2016, pp. 379–388. [Online]. Available: <https://dl.acm.org/doi/10.1145/2876034.2876042>
- [144] M. Robnik-Šikonja and M. Bohanec, *Perturbation-Based Explanations of Prediction Models*. Springer International Publishing, 2018. [Online]. Available: [http://dx.doi.org/10.1007/978-3-319-90403-0\\_9](http://dx.doi.org/10.1007/978-3-319-90403-0_9)
- [145] J. Chang *et al.*, “Reading Tea Leaves: How Humans Interpret Topic Models,” in *Advances in Neural Information Processing Systems*, Y. Bengio *et al.*, Eds., vol. 22. Curran Associates, Inc., 2009, pp. 288–296. [Online]. Available: <https://proceedings.neurips.cc/paper/2009/file/f92586a25bb3145facd64ab20fd554ff-Paper.pdf>
- [146] S. Park, J. Bak, and A. Oh, “Rotated Word Vector Representations and their Interpretability,” in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2017, pp. 401–411. [Online]. Available: <http://aclweb.org/anthology/D17-1041>
- [147] M. Du, N. Liu, and X. Hu, “Techniques for interpretable machine learning,” *Communications of the ACM*, vol. 63, no. 1, pp. 68–77, 12 2019. [Online]. Available: <https://dl.acm.org/doi/10.1145/3359786>

- [148] D. Alvarez-Melis and T. Jaakkola, “A causal framework for explaining the predictions of black-box sequence-to-sequence models,” in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2017, pp. 412–421. [Online]. Available: <http://aclweb.org/anthology/D17-1042>
- [149] P. K. Mudrakarta *et al.*, “Did the model understand the question?” in *ACL 2018 - 56th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (Long Papers)*, vol. 1, 5 2018, pp. 1896–1906. [Online]. Available: <https://www.aclweb.org/anthology/P18-1176/>
- [150] T. Wu *et al.*, “Polyjuice: Generating Counterfactuals for Explaining, Evaluating, and Improving Models,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 1 2021, pp. 6707–6723. [Online]. Available: <https://aclanthology.org/2021.acl-long.523>
- [151] Shapley, “A value for N-Person Games,” *Contributions to the Theory of Games (AM-28), Volume II*, pp. 307–317, 1953. [Online]. Available: <https://apps.dtic.mil/dtic/tr/fulltext/u2/604084.pdf>
- [152] C. Molnar, *Interpreting Machine Learning Models With SHAP*, 2023.
- [153] S. Abnar and W. Zuidema, “Quantifying Attention Flow in Transformers,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 4190–4197. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.385>
- [154] T. H. Cormen *et al.*, *Introduction to Algorithms, Third Edition*, 3rd ed. The MIT Press, 2009.
- [155] K. Ethayarajh and D. Jurafsky, “Attention Flows are Shapley Value Explanations,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2021, pp. 49–54. [Online]. Available: <https://aclanthology.org/2021.acl-short.8>
- [156] G. Brunner *et al.*, “On Identifiability in Transformers,” in *International Conference on Learning Representations (ICLR 2020)*, 8 2020. [Online]. Available: <https://openreview.net/forum?id=BJg1f6EFDB>

- [157] M. Tutek and J. Snajder, “Staying True to Your Word: (How) Can Attention Become Explanation?” in *Proceedings of the 5th Workshop on Representation Learning for NLP*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 131–142. [Online]. Available: <https://www.aclweb.org/anthology/2020.repl4nlp-1.17>
- [158] S. Srinivas and F. Fleuret, “Rethinking the Role of Gradient-Based Attribution Methods for Model Interpretability,” *ICLR 2021 - 9th International Conference on Learning Representations*, 2021.
- [159] M. Gardner *et al.*, “Evaluating Models’ Local Decision Boundaries via Contrast Sets,” in *Findings of the Association for Computational Linguistics: EMNLP 2020*. Stroudsburg, PA, USA: Association for Computational Linguistics, 4 2020, pp. 1307–1323. [Online]. Available: <https://www.aclweb.org/anthology/2020.findings-emnlp.117>
- [160] D. Kaushik, E. Hovy, and Z. C. Lipton, “Learning The Difference That Makes A Difference With Counterfactually-Augmented Data,” in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=SkLgs0NFvr>
- [161] A. Radford *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019. [Online]. Available: <https://openai.com/blog/better-language-models/>
- [162] Y. Zhang, J. Baldridge, and L. He, “PAWS: Paraphrase adversaries from word scrambling,” in *Proceedings of the 2019 Conference of the North*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 1298–1308. [Online]. Available: <http://aclweb.org/anthology/N19-1131>
- [163] K. Sakaguchi *et al.*, “WinoGrande: An Adversarial Winograd Schema Challenge at Scale,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 05, pp. 8732–8740, 4 2020. [Online]. Available: <https://aaai.org/ojs/index.php/AAAI/article/view/6399>
- [164] J. Wieting and K. Gimpel, “ParaNMT-50M: Pushing the Limits of Paraphrastic Sentence Embeddings with Millions of Machine Translations,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 451–462. [Online]. Available: <http://aclweb.org/anthology/P18-1042>

- [165] T. Lei, R. Barzilay, and T. Jaakkola, “Rationalizing Neural Predictions,” in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2016, pp. 107–117. [Online]. Available: <http://aclweb.org/anthology/D16-1011>
- [166] O.-M. Camburu *et al.*, “e-SNLI: Natural Language Inference with Natural Language Explanations,” in *Advances in Neural Information Processing Systems*, vol. 2018-Decem, 12 2018, pp. 9539–9549. [Online]. Available: <http://arxiv.org/abs/1812.01193>
- [167] H. Liu, Q. Yin, and W. Y. Wang, “Towards Explainable NLP: A Generative Explanation Framework for Text Classification,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 5570–5581. [Online]. Available: <https://www.aclweb.org/anthology/P19-1560>
- [168] V. Lattinik and J. Berant, “Explaining Question Answering Models through Text Generation,” *arXiv*, 4 2020. [Online]. Available: <http://arxiv.org/abs/2004.05569>
- [169] S. Gurrapu *et al.*, “Rationalization for explainable NLP: a survey,” *Frontiers in Artificial Intelligence*, vol. 6, 2023.
- [170] A. Rana *et al.*, “RerrFact: Reduced Evidence Retrieval Representations for Scientific Claim Verification,” in *CEUR Workshop Proceedings*, vol. 3164, 2 2022, pp. 3–7. [Online]. Available: <http://arxiv.org/abs/2202.02646>
- [171] J. DeYoung *et al.*, “ERASER: A Benchmark to Evaluate Rationalized NLP Models,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 11 2020, pp. 4443–4458. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.408>
- [172] A. Radford *et al.*, “Improving Language Understanding by Generative Pre-Training,” *OpenAI*, 2018. [Online]. Available: <https://openai.com/blog/language-unsupervised/>
- [173] R. Socher *et al.*, “Recursive deep models for semantic compositionality over a sentiment treebank,” *EMNLP 2013 - 2013 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*, pp. 1631–1642, 2013.
- [174] P. Hase *et al.*, “Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language?” *Findings of the Association for Computational Linguistics Findings of ACL: EMNLP 2020*, pp. 4351–4367, 2020. [Online]. Available: <https://www.aclweb.org/anthology/2020.findings-emnlp.390>

- [175] S. Wiegreffe and A. Marasović, “Teach Me to Explain: A Review of Datasets for Explainable Natural Language Processing,” in *35th Conference on Neural Information Processing Systems (NeurIPS 2021)*, 2 2021. [Online]. Available: <http://arxiv.org/abs/2102.12060>
- [176] J. Andreas, A. Dragan, and D. Klein, “Translating neuralese,” in *ACL 2017 - 55th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (Long Papers)*, 2017, pp. 232–242. [Online]. Available: <http://github>.
- [177] P. Atanasova *et al.*, “Faithfulness Tests for Natural Language Explanations,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, vol. 2. Stroudsburg, PA, USA: Association for Computational Linguistics, 5 2023, pp. 283–294. [Online]. Available: <https://aclanthology.org/2023.acl-short.25>
- [178] S. Wiegreffe, A. Marasović, and N. A. Smith, “Measuring Association Between Labels and Free-Text Rationales,” *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 10 266–10 284, 2020. [Online]. Available: <http://arxiv.org/abs/2010.12762><https://aclanthology.org/2021.emnlp-main.804>
- [179] P. Atanasova *et al.*, “Fact Checking with Insufficient Evidence,” *Transactions of the Association for Computational Linguistics*, vol. 10, pp. 746–763, 7 2022. [Online]. Available: [https://direct.mit.edu/tacl/article/doi/10.1162/tacl\\_a\\_00486/112498/Fact-Checking-with-Insufficient-Evidence](https://direct.mit.edu/tacl/article/doi/10.1162/tacl_a_00486/112498/Fact-Checking-with-Insufficient-Evidence)
- [180] M. Turpin *et al.*, “Language Models Don’t Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting,” in *Thirty-seventh Conference on Neural Information Processing Systems*, 5 2023, pp. 1–32. [Online]. Available: <http://arxiv.org/abs/2305.04388><https://openreview.net/forum?id=bzs4uPLXvi>
- [181] T. Lanham *et al.*, “Measuring Faithfulness in Chain-of-Thought Reasoning,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2307.13702>
- [182] L. Parcalabescu and A. Frank, “On Measuring Faithfulness of Natural Language Explanations,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2311.07466>
- [183] M. Jang and T. Lukasiewicz, “Are Training Resources Insufficient? Predict First Then Explain!” *arXiv*, 8 2021. [Online]. Available: <http://arxiv.org/abs/2110.02056>

- [184] V. Do *et al.*, “e-SNLI-VE-2.0: Corrected Visual-Textual Entailment with Natural Language Explanations,” *IEEE CVPR Workshop on Fair, Data Efficient and Trusted Computer Vision, 2020*, 2020. [Online]. Available: <https://github.com/>
- [185] D. Pruthi *et al.*, “Learning to Deceive with Attention-Based Explanations,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 4782–4793. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.432>
- [186] L. Arras, A. Osman, and W. Samek, “CLEVR-XAI: A benchmark dataset for the ground truth evaluation of neural network explanations,” *Information Fusion*, vol. 81, pp. 14–40, 5 2022. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S1566253521002335>
- [187] K. Zaman and Y. Belinkov, “A Multilingual Perspective Towards the Evaluation of Attribution Methods in Natural Language Inference,” in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*, 4 2022, pp. 1556–1576. [Online]. Available: <http://arxiv.org/abs/2204.05428>
- [188] T. M. Pham *et al.*, “Double Trouble: How to not Explain a Text Classifier’s Decisions Using Counterfactuals Synthesized by Masked Language Models?” in *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Y. He *et al.*, Eds. Online only: Association for Computational Linguistics, 11 2022, pp. 12–31. [Online]. Available: <https://aclanthology.org/2022.aacl-main.2>
- [189] S. Gururangan *et al.*, “Annotation Artifacts in Natural Language Inference Data,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, vol. 2. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 107–112. [Online]. Available: <http://aclweb.org/anthology/N18-2017>
- [190] P. Hase, H. Xie, and M. Bansal, “The Out-of-Distribution Problem in Explainability and Search Methods for Feature Importance Explanations,” *Advances in Neural Information Processing Systems*, vol. 5, no. NeurIPS, pp. 3650–3666, 2021.
- [191] K. Vafa *et al.*, “Rationales for Sequential Predictions,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA,

- USA: Association for Computational Linguistics, 2021, pp. 10 314–10 332. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.807>
- [192] A. Wettig *et al.*, “Should You Mask 15% in Masked Language Modeling?” *EACL 2023 - 17th Conference of the European Chapter of the Association for Computational Linguistics, Proceedings of the Conference*, pp. 2977–2992, 2 2023. [Online]. Available: <http://arxiv.org/abs/2202.08005>
- [193] R. J. Simes, “An Improved Bonferroni Procedure for Multiple Tests of Significance,” *Biometrika*, vol. 73, no. 3, p. 751, 12 1986. [Online]. Available: <https://www.jstor.org/stable/2336545?origin=crossref>
- [194] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer Normalization,” *Arxiv*, 2016. [Online]. Available: <http://arxiv.org/abs/1607.06450>
- [195] R. A. Fisher, “Statistical Methods for Research Workers,” in *Breakthroughs in Statistics: Methodology and Distribution*, S. Kotz and N. L. Johnson, Eds. New York, NY: Springer New York, 1992, pp. 66–70. [Online]. Available: [http://link.springer.com/10.1007/978-1-4612-4380-9\\_6](http://link.springer.com/10.1007/978-1-4612-4380-9_6)
- [196] S. T. Buckland, A. C. Davison, and D. V. Hinkley, “Bootstrap Methods and Their Application,” *Biometrics*, vol. 54, no. 2, p. 795, 6 1998.
- [197] Michael R. Chernick and R. A. LaBudde, *An introduction to bootstrap methods with applications to R*. John Wiley & Sons, 2011.
- [198] N. Bansal, C. Agarwal, and A. Nguyen, “SAM: The Sensitivity of Attribution Methods to Hyperparameters,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, 6 2020, pp. 11–21. [Online]. Available: <https://ieeexplore.ieee.org/document/9150607/>
- [199] J. Yang *et al.*, “Generalized Out-of-Distribution Detection: A Survey,” *arXiv*, 10 2021. [Online]. Available: <http://arxiv.org/abs/2110.11334>
- [200] R. Sun and C. H. Lampert, “KS(conf): A Light-Weight Test if a Multiclass Classifier Operates Outside of Its Specifications,” *International Journal of Computer Vision*, vol. 128, no. 4, pp. 970–995, 4 2020. [Online]. Available: <http://link.springer.com/10.1007/s11263-019-01232-x>
- [201] A. Dziedzic *et al.*, “ $\mathcal{P}$ -DkNN: Out-of-Distribution Detection Through Statistical Testing of Deep Representations,” *arXiv*, 7 2022. [Online]. Available: <http://arxiv.org/abs/2207.12545>

- [202] G. Kwon *et al.*, “Backpropagated Gradient Representations for Anomaly Detection,” *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 12366 LNCS, pp. 206–226, 7 2020. [Online]. Available: <http://arxiv.org/abs/2007.09507>
- [203] B. Min *et al.*, “Recent Advances in Natural Language Processing via Large Pre-trained Language Models: A Survey,” *ACM Computing Surveys*, vol. 56, no. 2, pp. 1–40, 2 2024. [Online]. Available: <https://dl.acm.org/doi/10.1145/3605943>
- [204] S. Zhang *et al.*, “Applications of transformer-based language models in bioinformatics: a survey,” *Bioinformatics Advances*, vol. 3, no. 1, 1 2023. [Online]. Available: <https://academic.oup.com/bioinformaticsadvances/article/doi/10.1093/bioadv/vbad001/6984737>
- [205] N. Muennighoff *et al.*, “Generative Representational Instruction Tuning,” *arXiv*, 2024. [Online]. Available: <http://arxiv.org/abs/2402.09906>
- [206] C.-Y. Lin, “ROUGE: A Package for Automatic Evaluation of Summaries,” in *Text Summarization Branches Out*. Barcelona, Spain: Association for Computational Linguistics, 7 2004, pp. 74–81. [Online]. Available: <https://aclanthology.org/W04-1013>
- [207] Meta, “Llama 2: Open Foundation and Fine-Tuned Chat Models,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2307.09288>
- [208] G. Penedo *et al.*, “The RefinedWeb Dataset for Falcon LLM: Outperforming Curated Corpora with Web Data, and Web Data Only,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2306.01116>
- [209] A. Q. Jiang *et al.*, “Mistral 7B,” *arXiv*, pp. 1–9, 2023. [Online]. Available: <http://arxiv.org/abs/2310.06825>
- [210] OpenAI, “GPT-4 Technical Report,” *OpenAI*, vol. 4, pp. 1–100, 3 2023. [Online]. Available: <http://arxiv.org/abs/2303.08774>
- [211] Y. Bang *et al.*, “A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2302.04023>
- [212] J.-Y. Yao *et al.*, “LLM Lies: Hallucinations are not Bugs, but Features as Adversarial Examples,” *arXiv*, pp. 1–13, 2023. [Online]. Available: <http://arxiv.org/abs/2310.01469>

- [213] C. Agarwal, S. H. Tanneru, and H. Lakkaraju, “Faithfulness vs. Plausibility: On the (Un)Reliability of Explanations from Large Language Models,” *arXiv*, 2024. [Online]. Available: <http://arxiv.org/abs/2402.04614>
- [214] Y. Chen *et al.*, “Do Models Explain Themselves? Counterfactual Simulatability of Natural Language Explanations,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2307.08678>
- [215] S. Huang *et al.*, “Can Large Language Models Explain Themselves? A Study of LLM-Generated Self-Explanations,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2310.11207>
- [216] C. Singh *et al.*, “Rethinking Interpretability in the Era of Large Language Models,” *arXiv*, 2024. [Online]. Available: <http://arxiv.org/abs/2402.01761>
- [217] Z. Li *et al.*, “Using Interactive Feedback to Improve the Accuracy and Explainability of Question Answering Systems Post-Deployment,” in *Findings of the Association for Computational Linguistics: ACL 2022*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2022, pp. 926–937. [Online]. Available: <https://aclanthology.org/2022.findings-acl.75>
- [218] Anthropic Team, “Language Models (Mostly) Know What They Know,” *Anthropic*, 7 2022. [Online]. Available: <http://arxiv.org/abs/2207.05221>
- [219] C. McDougall *et al.*, “Copy Suppression: Comprehensively Understanding an Attention Head,” in *NeurIPS 2023 Workshop on Attributing Model Behavior at Scale*, 2023. [Online]. Available: <http://arxiv.org/abs/2310.04625>
- [220] A. Deshpande *et al.*, “Toxicity in chatgpt: Analyzing persona-assigned language models,” in *Findings of the Association for Computational Linguistics: EMNLP 2023*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2023, pp. 1236–1270. [Online]. Available: <https://aclanthology.org/2023.findings-emnlp.88>
- [221] X. L. Li *et al.*, “Benchmarking and Improving Generator-Validator Consistency of Language Models,” *arXiv*, pp. 1–15, 2023. [Online]. Available: <http://arxiv.org/abs/2310.01846>
- [222] J. Hu and R. Levy, “Prompt-based methods may underestimate large language models’ linguistic generalizations,” *arXiv*, 2023. [Online]. Available: <http://arxiv.org/abs/2305.13264>

- [223] Q. Dong *et al.*, “A Survey on In-context Learning,” *arXiv*, 12 2022. [Online]. Available: <http://arxiv.org/abs/2301.00234>
- [224] A. Ross and F. Doshi-Velez, “Improving the Adversarial Robustness and Interpretability of Deep Neural Networks by Regularizing Their Input Gradients,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, pp. 1660–1669, 4 2018. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/11504>
- [225] S. Srinivas *et al.*, “Efficiently Training Low-Curvature Neural Networks,” *Advances in Neural Information Processing Systems*, vol. 35, no. NeurIPS, pp. 1–21, 6 2022. [Online]. Available: <http://arxiv.org/abs/2206.07144>
- [226] W. J. Yeo *et al.*, “How Interpretable are Reasoning Explanations from Prompting Large Language Models?” *arXiv*, 2 2024. [Online]. Available: <http://arxiv.org/abs/2402.11863>
- [227] C. Sen *et al.*, “Human Attention Maps for Text Classification: Do Humans and Neural Networks Focus on the Same Words?” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 4596–4608. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.419>
- [228] P. Hase and M. Bansal, “Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 5540–5552. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.491>
- [229] G. Prasad *et al.*, “To what extent do human explanations of model behavior align with actual model behavior?” in *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2021, pp. 1–14. [Online]. Available: <https://aclanthology.org/2021.blackboxnlp-1.1>
- [230] A. V. González, A. Rogers, and A. Søgaard, “On the Interaction of Belief Bias and Explanations,” in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2021, pp. 2930–2942. [Online]. Available: <https://aclanthology.org/2021.findings-acl.259>
- [231] H. Schuff *et al.*, “Human Interpretation of Saliency-based Explanation Over Text,” in *2022 ACM Conference on Fairness, Accountability, and Transparency*.

- New York, NY, USA: ACM, 6 2022, pp. 611–636. [Online]. Available: <https://dl.acm.org/doi/10.1145/3531146.3533127>
- [232] P. Lertvittayakumjorn and F. Toni, “Human-grounded Evaluations of Explanation Methods for Text Classification,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, no. Section 3. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 5194–5204. [Online]. Available: <https://www.aclweb.org/anthology/D19-1523>
- [233] D. Nguyen, “Comparing Automatic and Human Evaluation of Local Explanations for Text Classification,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, vol. 1. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 1069–1078. [Online]. Available: <http://aclweb.org/anthology/N18-1097>
- [234] L. Schut *et al.*, “Bridging the Human-AI Knowledge Gap: Concept Discovery and Transfer in AlphaZero,” *arXiv*, pp. 1–61, 10 2023. [Online]. Available: <http://arxiv.org/abs/2310.16410>
- [235] B. Kim, “Beyond interpretability: developing a language to shape our relationships with AI,” in *The International Conference on Learning Representations*, 2022. [Online]. Available: <https://iclr.cc/Conferences/2022/Schedule?showEvent=7237>
- [236] M. Stern, N. Shazeer, and J. Uszkoreit, “Blockwise parallel decoding for deep autoregressive models,” *Advances in Neural Information Processing Systems*, vol. 2018-Decem, no. Nips, pp. 10 086–10 095, 2018.
- [237] B. Spector and C. Re, “Accelerating LLM Inference with Staged Speculative Decoding,” *arXiv*, no. Llm, 2023. [Online]. Available: <http://arxiv.org/abs/2308.04623>
- [238] Y. Fu *et al.*, “Break the Sequential Dependency of LLM Inference Using Lookahead Decoding,” *arXiv*, 2024. [Online]. Available: <http://arxiv.org/abs/2402.02057>
- [239] T. Cai *et al.*, “Medusa: Simple LLM Inference Acceleration Framework with Multiple Decoding Heads,” *arXiv*, 2024. [Online]. Available: <http://arxiv.org/abs/2401.10774>
- [240] M. Nadeem, A. Bethke, and S. Reddy, “StereoSet: Measuring stereotypical bias in pre-trained language models,” *ACL-IJCNLP 2021 - 59th Annual Meeting of the Association*

- for *Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, Proceedings of the Conference*, vol. 2, pp. 5356–5371, 2021.
- [241] N. Nangia *et al.*, “CrowS-Pairs: A challenge dataset for measuring social biases in masked language models,” *EMNLP 2020 - 2020 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*, pp. 1953–1967, 2020.
- [242] U. Aïvodji *et al.*, “Fairwashing: The risk of rationalization,” *36th International Conference on Machine Learning, ICML 2019*, vol. 2019-June, pp. 240–252, 2019.
- [243] —, “Characterizing the risk of fairwashing,” *Advances in Neural Information Processing Systems*, vol. 18, no. NeurIPS, pp. 14 822–14 834, 2021.
- [244] A. Glaese *et al.*, “Improving alignment of dialogue agents via targeted human judgements,” *arXiv*, pp. 1–77, 2022. [Online]. Available: <http://arxiv.org/abs/2209.14375>
- [245] W. Wang *et al.*, “Towards a Robust Deep Neural Network against Adversarial Texts: A Survey,” *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2021. [Online]. Available: <https://ieeexplore.ieee.org/document/9557814/>
- [246] J. Ebrahimi *et al.*, “HotFlip: White-Box Adversarial Examples for Text Classification,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, vol. 2. Stroudsburg, PA, USA: Association for Computational Linguistics, 2018, pp. 31–36. [Online]. Available: <http://aclweb.org/anthology/P18-2006>
- [247] E. Wallace *et al.*, “Universal Adversarial Triggers for Attacking and Analyzing NLP,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 2153–2162. [Online]. Available: <https://www.aclweb.org/anthology/D19-1221>
- [248] K. Papineni *et al.*, “BLEU,” in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02*. Morristown, NJ, USA: Association for Computational Linguistics, 2001. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1073083.1073135>
- [249] R. D. Cook and S. Weisberg, “Characterizations of an Empirical Influence Function for Detecting Influential Cases in Regression,” *Technometrics*, vol. 22, no. 4, pp. 495–508, 11 1980. [Online]. Available: <http://www.tandfonline.com/doi/abs/10.1080/00401706.1980.10486199>

- [250] P. W. Koh and P. Liang, “Understanding Black-box Predictions via Influence Functions,” *34th International Conference on Machine Learning, ICML 2017*, vol. 4, pp. 2976–2987, 3 2017. [Online]. Available: <http://arxiv.org/abs/1703.04730>
- [251] C.-K. Yeh *et al.*, “Representer Point Selection for Explaining Deep Neural Networks,” in *Advances in Neural Information Processing Systems*, 11 2018, pp. 9291–9301. [Online]. Available: <http://arxiv.org/abs/1811.09720>
- [252] X. Han, B. C. Wallace, and Y. Tsvetkov, “Explaining Black Box Predictions and Unveiling Data Artifacts through Influence Functions,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 5553–5563. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.492>
- [253] H. Guo *et al.*, “FastIF: Scalable Influence Functions for Efficient Model Interpretation and Debugging,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 12 2021, pp. 10 333–10 350. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.808>
- [254] B. Schölkopf, R. Herbrich, and A. J. Smola, “A Generalized Representer Theorem,” in *International Conference on Computational Learning Theory*. Springer, 2001, pp. 416–426. [Online]. Available: [http://link.springer.com/10.1007/3-540-44581-1\\_27](http://link.springer.com/10.1007/3-540-44581-1_27)
- [255] G. Pruthi *et al.*, “Estimating Training Data Influence by Tracing Gradient Descent,” in *Advances in Neural Information Processing Systems*, 2 2020. [Online]. Available: <http://arxiv.org/abs/2002.08484>
- [256] Y. Goyal, U. Shalit, and B. Kim, “Explaining classifiers with causal concept effect (CaCE),” *arXiv*, 7 2019. [Online]. Available: <http://arxiv.org/abs/1907.07165>
- [257] B. Kim *et al.*, “Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV),” *35th International Conference on Machine Learning, ICML 2018*, vol. 6, pp. 4186–4195, 11 2018. [Online]. Available: <http://arxiv.org/abs/1711.11279>
- [258] J. Mu and J. Andreas, “Compositional Explanations of Neurons,” in *Advances in Neural Information Processing Systems*, 6 2020. [Online]. Available: <http://arxiv.org/abs/2006.14032>

- [259] J. Pearl, “Direct and Indirect Effects,” in *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, ser. UAI’01. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2001, p. 411–420. [Online]. Available: <https://dl.acm.org/doi/10.5555/2074022.2074073>
- [260] A. Ghorbani *et al.*, “Towards automatic concept-based explanations,” in *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [261] T. Mikolov *et al.*, “Efficient estimation of word representations in vector space,” in *1st International Conference on Learning Representations, ICLR 2013 - Workshop Track Proceedings*, 2013. [Online]. Available: <http://ronan.collobert.com/senna/>
- [262] T. Bolukbasi *et al.*, “Man is to computer programmer as woman is to homemaker? Debiasing word embeddings,” in *Advances in Neural Information Processing Systems*, 2016, pp. 4356–4364.
- [263] A. B. Costello and J. W. Osborne, “Best practices in exploratory factor analysis: Four recommendations for getting the most from your analysis,” *Practical Assessment, Research and Evaluation*, vol. 10, no. 7, pp. 1–9, 2005.
- [264] C. B. Crawford and G. A. Ferguson, “A general rotation criterion and its use in orthogonal rotation,” *Psychometrika*, vol. 35, no. 3, pp. 321–332, 9 1970. [Online]. Available: <http://link.springer.com/10.1007/BF02310792>
- [265] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent Dirichlet allocation,” *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003. [Online]. Available: <https://jmlr.org/papers/v3/blei03a.html>
- [266] M. Ibrahim *et al.*, “Global Explanations of Neural Networks,” in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. New York, NY, USA: ACM, 1 2019, pp. 279–287. [Online]. Available: <https://dl.acm.org/doi/10.1145/3306618.3314230>
- [267] K. N. Ramamurthy *et al.*, “Model Agnostic Multilevel Explanations,” *arXiv*, 3 2020. [Online]. Available: <http://arxiv.org/abs/2003.06005>
- [268] P. Michel, O. Levy, and G. Neubig, “Are Sixteen Heads Really Better than One?” *Advances in Neural Information Processing Systems*, vol. 32, pp. 1–13, 5 2019. [Online]. Available: <http://arxiv.org/abs/1905.10650>
- [269] T. Linzen, E. Dupoux, and Y. Goldberg, “Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies,” *Transactions of the Association for*

- Computational Linguistics*, vol. 4, no. 1990, pp. 521–535, 12 2016. [Online]. Available: <https://direct.mit.edu/tacl/article/43378>
- [270] K. Sinha *et al.*, “UnNatural Language Inference,” in *ACL 2021 - 59th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference*, 2021. [Online]. Available: <http://arxiv.org/abs/2101.00010>
- [271] X. Shi, I. Padhi, and K. Knight, “Does string-based neural MT learn source syntax?” in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2016, pp. 1526–1534. [Online]. Available: <http://aclweb.org/anthology/D16-1159>
- [272] Y. Adi *et al.*, “Fine-grained Analysis of Sentence Embeddings Using Auxiliary Prediction Tasks,” in *International Conference on Learning Representations (ICLR)*, 8 2017, pp. 1–12. [Online]. Available: <http://arxiv.org/abs/1608.04207>
- [273] G. Brunner *et al.*, “Natural language multitasking analyzing and improving syntactic saliency of latent representations,” in *31st Conference on Neural Information Processing Systems (NIPS 2017)*, 1 2017. [Online]. Available: <http://arxiv.org/abs/1801.06024>
- [274] A. Köhn, “What’s in an Embedding? Analyzing Word Embeddings through Multilingual Evaluation,” in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2015, pp. 2067–2073. [Online]. Available: <http://aclweb.org/anthology/D15-1246>
- [275] I. Tenney *et al.*, “What do you learn from context? Probing for sentence structure in contextualized word representations,” in *7th International Conference on Learning Representations, ICLR 2019*, 2019, pp. 1–17. [Online]. Available: <https://openreview.net/forum?id=SJzSgnRcKX>
- [276] M. T. Ribeiro, S. Singh, and C. Guestrin, “Anchors: High-precision model-agnostic explanations,” in *32nd AAAI Conference on Artificial Intelligence, AAAI 2018*, 2018, pp. 1527–1535. [Online]. Available: <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16982>

## APPENDIX A LITERATURE REVIEW, OTHER COMMUNICATION METHODS

### A.1 Adversarial Examples

An *adversarial example* is an input that causes a model to produce a wrong prediction due to the limitations of the model. The adversarial example is often produced from an existing example for which the model produces a correct prediction. Because the *adversarial example* serves as an explanation, in the context of an existing example, it is a *local explanation*.

Wang et al. [245] provide a thorough survey on *adversarial example* explanations, and also goes in-depth regarding taxonomy, using *adversarial examples* for robustness, and similarity scores between the existing example and the *adversarial example*. Additionally, the survey by Belinkov and Glass [118] also has a section on adversarial examples.

In this chapter, we therefore focus on just two explanation methods. These *adversarial example* methods inform us about the support boundaries of a given example, which then informs us about the logic involved and, therefore, provides interpretability. In fact, this explanation can be similar to the *input feature* methods, discussed in Section 2.5. Many of those methods also indicate what words should be changed to alter the prediction. An important difference is that *adversarial* explanations are contrastive, meaning they explain by comparing with another example, while *input features* explain only concerning the original example. Contrastive explanations are, from a social science perspective, generally considered more *human-grounded* [134].

In the following discussions, we refer to the original example as  $\mathbf{x}$  and the adversarial example as  $\tilde{\mathbf{x}}$ . The goal is to develop an adversarial method  $A$ , that maps from  $\mathbf{x}$  to  $\tilde{\mathbf{x}}$ :

$$A(\mathbf{x}) \rightarrow \tilde{\mathbf{x}} \tag{A.1}$$

Importantly, to ensure that an *adversarial example* method is *faithful*, one only needs to assert that the predicted label changes while the gold label remains the same. Additionally, it's desirable to have the original and adversarial example to be similar, in many applications this can be framed as paraphrasing. Compared to other explanation types, these properties are reasonably trivial to measure. See Section 2.3 for a general discussion on measures of interpretability.

Finally, because *adversarial example* explanations are framed by the output class, these

explanations do not generalize easily to sequence-to-sequence problems. One could imagine, for example, an offensive-text classifier, which reduces the sequence-to-sequence model back to a sequence-to-class model.

### A.1.1 HotFlip

A great example of the relation between *input feature* explanations and *adversarial examples* is *HotFlip* [246]. Here the effect of changing token  $v$  to another token  $\tilde{v}$  at position  $t$ , on the model loss  $\mathcal{L}$ , is estimated via using gradients

$$\mathcal{L}(y, \tilde{\mathbf{x}}_{t:v \rightarrow \tilde{v}}) - \mathcal{L}(y, \mathbf{x}; \theta) \approx \frac{\partial \mathcal{L}(y, \mathbf{x}; \theta)}{\partial x_{t, \tilde{v}}} - \frac{\partial \mathcal{L}(y, \mathbf{x}; \theta)}{\partial x_{t, v}}, \quad (\text{A.2})$$

where  $\tilde{\mathbf{x}}_{t:v \rightarrow \tilde{v}}$  is the one-hot-encoded input  $\mathbf{x}$ , with the token  $v$  at position  $t$  changed to  $\tilde{v}$ . Additionally,  $x_{t, \tilde{v}}$  and  $x_{t, v}$  are the scalar components of the one-hot-encoded input  $\mathbf{x}$ .

Had a gradient approximation not been used, the alternative would be to compute a forward pass for every possible token swap exactly. Instead, this approximation only requires one backward pass. The authors use a beam-search approach to produce an adversarial sentence with multiple tokens changed. A visualization of *HotFlip* can be seen in Figure A.1.

$$A_{\text{HotFlip}}(\mathbf{x}) = \underset{\tilde{\mathbf{x}}_{t:v \rightarrow \tilde{v}}}{\operatorname{argmax}} \frac{\partial \mathcal{L}(y, \mathbf{x}; \theta)}{\partial x_{t, \tilde{v}}} - \frac{\partial \mathcal{L}(y, \mathbf{x}; \theta)}{\partial x_{t, v}} \quad (\text{A.3})$$

			$p(y \mathbf{x}; \theta)$	$y$
$\mathbf{x}$	the year 's <b>best</b> and most unpredictable comedy		0.91	pos
	the year 's finest and most unpredictable comedy	▼	0.30	-
$\tilde{\mathbf{x}}$	the year 's finest and most unforeseeable comedy		0.08	-
$\mathbf{x}$	we <b>never</b> feel anything for these <b>characters</b>		0.95	neg
$\tilde{\mathbf{x}}$	we never feel anything for these people	▼	0.03	-

Figure A.1 Hypothetical visualization of *HotFlip*. The highlight indicates the gradient w.r.t. the input, which HotFlip uses to select which token to change.  $\mathbf{x}$  indicates the original sentence, and  $\tilde{\mathbf{x}}$  indicates the adversarial sentence.

The *HotFlip* paper [246] primarily investigates character-level models, for which the desire is to build a model that is robust against typos. However, in terms of word-level models, it is necessary to constrain the possible changes such that the adversarial sentence is a

paraphrase. They do this via word embeddings, such that the adversarial and original words are constrained to have a cosine similarity of at least 0.8.

The *HotFlip* approach has proven effective for other adversarial explanation methods, such as the aforementioned Universal Adversarial Triggers [247].

### A.1.2 Semantically Equivalent Adversaries (SEA)

An alternative approach to produce adversarial examples that are ensured to be paraphrased is to sample from a paraphrasing model  $q(\tilde{\mathbf{x}}|\mathbf{x})$ . Ribeiro et al. [15] do this by measuring a semantical-equivalency-score  $S(\mathbf{x}, \tilde{\mathbf{x}})$ , as the relative likelihood of  $q(\tilde{\mathbf{x}}|\mathbf{x})$  compared to  $q(\mathbf{x}|\mathbf{x})$ . It is then possible to maximize the similarity while still having a different model prediction. The exact method is defined in (A.4), which also constrains the optimization with a minimum semantical-equivalency-score and ensures the predicted label is different.

$$\begin{aligned}
 A_{\text{SEA}}(\mathbf{x}) &= \underset{\tilde{\mathbf{x}} \sim q(\tilde{\mathbf{x}}|\mathbf{x})}{\operatorname{argmax}} S(\mathbf{x}, \tilde{\mathbf{x}}) \\
 &\text{s.t. } S(\mathbf{x}, \tilde{\mathbf{x}}) \geq 0.8 \\
 &\quad \underset{i}{\operatorname{argmax}} p(i|\mathbf{x}; \theta) \neq \underset{i}{\operatorname{argmax}} p(i|\tilde{\mathbf{x}}; \theta) \\
 &\text{where } S(\mathbf{x}, \tilde{\mathbf{x}}) = \min \left( 1, \frac{q(\tilde{\mathbf{x}}|\mathbf{x})}{q(\mathbf{x}|\mathbf{x})} \right)
 \end{aligned} \tag{A.4}$$

The reason why a relative score is necessary, as opposed to just using  $S(\mathbf{x}, \tilde{\mathbf{x}}) = q(\tilde{\mathbf{x}}|\mathbf{x})$ , is that for two normal sentences  $\mathbf{x}_1$  and  $\mathbf{x}_2$  of different length, longer sentences are just inherently less likely. Therefore, to maintain a comparative semantical-equivalency-score normalizing by  $q(\mathbf{x}|\mathbf{x})$  is necessary [15].

	$\mathbf{x}$	$p(y \mathbf{x}; \theta)$	$y$	$S(\mathbf{x}, \tilde{\mathbf{x}})$
$\mathbf{x}$	<u>the</u> <u>year</u> <u>'s</u> <u>best</u> <u>and</u> <u>most</u> <u>unpredictable</u> <u>comedy</u>	0.91	pos	-
$\tilde{\mathbf{x}}$	<u>the</u> <u>best</u> <u>and</u> <u>most</u> <u>unpredictable</u> <u>comedy</u> <u>this</u> <u>year</u>	0.13	-	0.87
$\mathbf{x}$	<u>we</u> <u>never</u> <u>feel</u> <u>anything</u> <u>for</u> <u>these</u> <u>characters</u>	0.95	neg	-
$\tilde{\mathbf{x}}$	<u>we</u> <u>never</u> <u>empathize</u> <u>for</u> <u>these</u> <u>characters</u>	0.11	-	0.93

Figure A.2 Hypothetical results of using *SEA* [15]. Note that unlike *HotFlip*, *SEA* can change and delete multiple tokens simultaneously as it samples from a paraphrasing model. Again,  $\mathbf{x}$  indicates the original sentence,  $\tilde{\mathbf{x}}$  indicates the adversarial sentence, and  $S(\mathbf{x}, \tilde{\mathbf{x}})$  is the semantical-equivalency-score which must be at least 0.8.

### A.1.3 Discussion

**Groundedness** *Adversarial example* are as mentioned, easy to measure *faithfulness* on and should be *human-grounded* due to their contrastive nature [134]. However, we are not aware of any work that explicitly tests for *human-groundedness*. This is likely because it is considered to be a given, but we advocate for testing such a hypothesis anyway.

**Future work** The difficulty with *adversarial example* explanations lies in the search procedure. For example, *HotFlip* [246] uses a greedy sequential search algorithm and would therefore not be able to identify combinatorial effects like a double-negative. While *SEA* Ribeiro et al. [15] depends on an expensive paraphrase generation model.

One typical limitation of *adversarial example* methods is that they provide no control of the search direction. Hypothetically, while changing “unpredictable” to “unforeseeable” could provide the largest source of error due to a robustness issue, it might be more interesting to discover that changing “women’s chess club” to “men’s chess club” also flips the label. Unfortunately, this aspect is usually not considered because the motivation for *adversarial example* generation is often robustness and debasing.

## A.2 Influential examples

For a given *input example*, an *influential examples* explanation finds examples from the training dataset that, in terms of the model’s understanding, look like the *input example*. Because this explanation method centers around a specific *input example* it is a *local explanation*. Note that it is different from just a distance metric on the inputs, such as BLEU [248], as this does not depend on the model.

*Influential examples* explanations can be quite useful. For example, for discovering dataset artifacts as some of the *influential examples* may have nothing to do with the *input example*, except for the artifacts. Additionally, they are commonly used to discover mislabeled observations.

The *influential examples* can always be presented as just the examples and a similarity score, see Figure A.3. Because the only presentation difference is the similarity score, this chapter does not include example figures for each method.

### A.2.1 Influence functions

*Influence functions* is a classical technique from robust statistics [249]. However, robust statistics have strong assumptions regarding convexity, low-dimensionality, and differentiability.

	$\mathbf{x}$	$p(y \mathbf{x};\theta)$	$y$	$\Delta$
$\mathbf{x}$	<u>the year 's best and most unpredictable comedy</u>	0.91	pos	0.21
$\tilde{\mathbf{x}}$	<u>a delightfully unpredictable , hilarious comedy</u>	0.95	pos	3.82
$\tilde{\mathbf{x}}$	<u>loud and thoroughly obnoxious comedy</u>	0.98	neg	-1.51

Figure A.3 Fictive result showing the *influential examples*  $\tilde{\mathbf{x}}$ , in relation to the *input example*  $\mathbf{x}$ , showing both examples with positive and negative influence.  $\Delta$  is the similarity score; the scale and range may depend on the specific method. Note it is possible to measure the influence of an example on itself. This can be useful to identify mislabeled observations, as such observations will be important for their own prediction.

Recent efforts in deep learning remove the low-dimensionality constraint and to some extent the convexity constraint [250].

The central idea in *influence functions*, is to estimate the effect on the loss  $\mathcal{L}$ , of removing the observation  $\tilde{\mathbf{x}}$  from the dataset. The most influential examples are those where the loss changes the most. Let  $\tilde{\theta}$  be the model parameters if  $\tilde{\mathbf{x}}$  had not been included in the training dataset, then the loss difference can be estimated using

$$\mathcal{L}(y, \mathbf{x}; \tilde{\theta}) - \mathcal{L}(y, \mathbf{x}; \theta) \approx \frac{1}{n} \nabla_{\theta} \mathcal{L}(y, \mathbf{x}; \theta)^{\top} H_{\theta}^{-1} \nabla_{\theta} \mathcal{L}(\tilde{y}, \tilde{\mathbf{x}}; \theta). \quad (\text{A.5})$$

Importantly, the Hessian  $H_{\theta}$  needs to be positive-definite, which can only be guaranteed for convex models. The authors Koh and Liang [250] avoid this issue by adding a diagonal to the Hessian until it is positive-definite. Additionally, they solve the computational issue of computing an inverse Hessian by formulating (A.5) as an inverse Hessian-vector product. Such formulation can be approximated in  $\mathcal{O}(np)$  time, where  $n$  is the number of observations and  $p$  is the number of parameters, hence a computational complexity identical to one training epoch. Note, however, that the inverse Hessian-vector product needs to be computed for every explained test observation  $\mathbf{x}$ .

One limitation of *influence functions* is that computing the *influence functions* is not always numerically stable [251], because (A.5) uses the gradient  $\nabla_{\theta} \mathcal{L}(\tilde{y}, \tilde{\mathbf{x}}; \theta)$  which is optimized to be close to zero.

Koh and Liang [250] looked at support-vector-machines, which are known to be convex, and convolutional neural networks, which are generally non-convex. Han et al. [252] then extended the analysis of *influence functions* to BERT [21]. This is a crucial step, as BERT may be much further from convexity than CNNs, thus causing the *influence functions* to be less *faithful*.

Han et al. [252] validates for *faithfulness* by removing the 10% most influential training examples from the dataset and then retraining the model. The results show a significant decrease in the model's performance on the test split, compared to removing the 10% least influential examples and 10% random examples, validating that the influential examples are important.

Additionally, Koh and Liang [250] measures *faithfulness* by setting 10% of training observations to a wrong label. *Influence functions* is then used to select a fraction of the dataset for which labels are corrected. The metric is then how many mislabeled observations were identified and the performance difference. The idea being, wrongly labeled observations should affect the loss more than correctly labeled observations, hence *influence functions* will tend to find wrongly labeled observations. Han et al. [252] perform a similar experiment, but instead, remove observations based on importance and then measure the performance difference. Both experiments validate that *influence functions* are *faithful*.

**Performance considerations.** A criticism of influence functions has been that it is computationally expensive. Although  $\nabla_{\theta} \mathcal{L}(y, \mathbf{x}; \theta)^{\top} H_{\theta}^{-1}$  can be cached for each test example, it is still too computationally intensive for real-time inspection of the model. Additionally, having to compute the weight-gradient  $\nabla_{\theta} \mathcal{L}(\tilde{y}, \tilde{\mathbf{x}}; \theta)$  and inner-product for every training observation, does not scale sufficiently. To this end, Guo et al. [253] propose to only use a subset of training data, using a KNN clustering. Additionally, they show that the hyperparameters when computing  $\nabla_{\theta} \mathcal{L}(y, \mathbf{x}; \theta)^{\top} H_{\theta}^{-1}$  can be tuned to reduce the computation to less than half.

### A.2.2 Representer Point Selection

An alternative to *influence functions* is the Representer theorem [254]. The central idea is that the logits of a test example  $\mathbf{x}$ , can be expressed as a decomposition of all training samples  $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i \kappa(\mathbf{x}, \tilde{\mathbf{x}}_i)$ . The original Representer theorem [254] works on *reproducing kernel Hilbert spaces*, which is not applicable to deep learning. However, recent work has applied the idea to neural networks [251].

Let  $\boldsymbol{\theta}_L$  be the weight matrix of the final layer, such that the logits  $f(\mathbf{x}; \theta) = \boldsymbol{\theta}_L \mathbf{z}_{L-1}(\mathbf{x}; \theta)$ , then if the regularized loss  $\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\tilde{y}_i, \tilde{\mathbf{x}}_i; \theta) + \lambda \|\boldsymbol{\theta}_L\|^2$ , is a stationary point and  $\lambda > 0$ , then

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i \mathbf{z}_{L-1}(\tilde{\mathbf{x}}_i; \theta)^{\top} \mathbf{z}_{L-1}(\mathbf{x}; \theta), \text{ where } \alpha_i = \frac{1}{2\lambda \cdot n} \frac{\partial \mathcal{L}(\tilde{y}_i, \tilde{\mathbf{x}}_i; \theta)}{\partial \mathbf{z}_{L-1}(\mathbf{x}_i; \theta)}. \quad (\text{A.6})$$

To understand the importance of each training observation  $\tilde{\mathbf{x}}_i$ , regarding the prediction of class  $c$  for the test example  $\mathbf{x}$ , one just looks at the  $c$ 'th element of each term  $\alpha_i \mathbf{z}_{L-1}(\tilde{\mathbf{x}}_i; \theta)^{\top} \mathbf{z}_{L-1}(\mathbf{x}; \theta)$ .

This approach is more numerically stable than *influence functions* [251], but has the downside of only depending on the intermediate representation of the final layer, while *influence functions* employs the entire model.

Because *Representer Point Selection* does depend on a specific model setup, where the last layer is regularized, this could be considered an *intrinsic* method. However, Yeh et al. [251] show that the stationary solution can be achieved *post-hoc*, meaning after learning, with minimal impact on the model predictions. They do this via the optimization problem

$$\boldsymbol{\theta}_L = \underset{\mathbf{W}}{\operatorname{argmin}} \left( \frac{1}{n} \sum_{i=1}^n \mathcal{L}(p(\cdot|\tilde{\mathbf{x}}_i; \theta), \mathbf{W} \mathbf{z}_{L-1}(\tilde{\mathbf{x}}_i; \theta)) + \lambda \|\mathbf{W}\|^2 \right), \quad (\text{A.7})$$

where  $\theta$  is the original model parameters,  $\boldsymbol{\theta}_L$  are the new parameters for the last layer, and  $\mathcal{L}$  is the full cross-entropy loss. Because this is a fairly low-dimensional problem, fine-tuning this can be done with an L-BFGS optimizer or similar [251].

Yeh et al. [251] show this method is *faithful* on a computer vision task, using a label-correction experiment similar to that in *influence functions*. In this case,  $|\boldsymbol{\alpha}_{i,c}|$  is used to select the observations to perform label correction on. Their results show that *Representer Point Selection* and *influence functions* can identify wrong labels equally well, but that the observations which *Representer Point Selection* selects affects the models performance more. Unfortunately, Yeh et al. [251] only show anecdotal results on an NLP task.

### A.2.3 TracIn

The idea behind *TracIn* by Pruthi et al. [255] is to accumulate loss changes during training. Specifically, the loss change on the test observation  $\mathbf{x}$  when optimizing  $\tilde{\mathbf{x}}$ . Pruthi et al. [255] first introduce an idealized version of this, which assumes optimization is done on one observation at a time (for example, SGD):

$$\text{TracInIdeal}(\tilde{\mathbf{x}}, \mathbf{x}) = \sum_{t \in \mathcal{T}_{\tilde{\mathbf{x}}}} \mathcal{L}(y, \mathbf{x}, \theta_t) - \mathcal{L}(y, \mathbf{x}, \theta_{t+1}), \text{ where } \mathcal{T}_{\tilde{\mathbf{x}}} \text{ is timestep which optimized } \tilde{\mathbf{x}} \quad (\text{A.8})$$

*TracIn* *TracIn* is then a relaxation of this idealized version. Rather than using a direct loss difference, gradients are used. Rather than assuming stochastic gradient descent (or similar), mini-batches can be used. Rather than checking every time step, checkpoints collected during training can be used.

$$\text{TracIn}(\tilde{\mathbf{x}}, \mathbf{x}) = \frac{1}{b} \sum_{t \in \mathcal{C}} \eta_t \nabla_{\theta_t} \mathcal{L}(y, \mathbf{x}, \theta_t) \cdot \nabla_{\theta_t} \mathcal{L}(\tilde{y}, \tilde{\mathbf{x}}, \theta_t), \quad (\text{A.9})$$

where  $\mathcal{C}$  are checkpoints,  $b$  is batch-size, and  $\eta_t$  is learning-rate.

Note, that the (A.9) formulation is still based on plain gradient descent. However, Pruthi et al. [255] instruct how to adapt this to most learning algorithms (AdaGrad, Adam, Newton, etc).

As a *faithfulness* evaluation, Pruthi et al. [255] repeat the label-correction experiment of *influence functions* and *Representer Point Selection*, and find that their method can better select mislabeled observations. Note that this was evaluated on CIFAR-10 and MNIST. Unfortunately, Pruthi et al. [255] does not do any evaluation on NLP tasks, but they do anecdotally show it works on an NLP application.

#### A.2.4 Discussion

**Groundedness** *Influential example* explanations, is one of the few categories with a non-trivial but appropriate *faithfulness* metric, namely the label-correction experiment, which is used somewhat consistently across papers. Unfortunately, this experiment has not been used on NLP tasks, and in general, very little *faithfulness* validation has been done in NLP.

Additionally, the label-correction experiment is somewhat limited, as it evaluates the influence of a training observation on itself. This is not how a *Influential examples* explanation would be used in most applications, for example, dataset artifact discovery. Therefore, we suggest future work also include the experiment from Guo et al. [253], which uses information removal.

**Future work** A natural question when asking what training observations are influential is to also what part of them is important. *Influence functions* can answer this, although at an increased computational cost. However, *TracIn* can not. For sequential outputs, it is interesting to also be able to select parts of the output and ask what influenced this. Both of these questions are becoming increasingly relevant with large-scale language models, where there is a large interest in understanding what caused a particular generation.

### A.3 Concepts

A *concept explanation* attempts to explain the model in terms of an abstraction of the input, called a *concept*. A classic example in computer vision is how the concept of stripes affects the classification of a zebra. Understanding this relationship is important, as a computer vision model could classify a zebra based on a horse-like shape and a savanna background.

Such a relation may yield a high accuracy score but is logically wrong.

The term *concept* is much more common in computer vision [256, 257, 258] than in NLP. Instead, the subject is often framed more concretely as bias-detection in NLP. For example, Vig et al. [16] uses the concept of occupation-words like *nurse*, and relates it to the classification of the words *he* and *she*.

Regardless of the field, in both NLP and CV, only a single class or a small subset of classes is analyzed. For this reason, *concept explanation* belongs in its own category of *class explanations*. However, we will likely see more types of *class explanations* in the future.

### A.3.1 Natural Indirect Effect (NIE)

Consider a language model with the prompt  $\mathbf{x}$  = “The nurse said that”. To measure if the gender-stereotype of “nurse” is female, it is natural to compare  $p(\text{she}|\mathbf{x};\theta)$  with  $p(\text{he}|\mathbf{x};\theta)$ , or alternatively  $p(\text{they}|\mathbf{x};\theta)$ . Generalized, Vig et al. [16] express this as

$$\text{bias-effect}(\mathbf{x};\theta) = \frac{p(\text{anti-stereotypical}|\mathbf{x};\theta)}{p(\text{stereotypical}|\mathbf{x};\theta)}. \quad (\text{A.10})$$

Vig et al. [16] then provide insight into which parts of the model are responsible for the bias. They do this by measuring the *Natural Indirect Effect* (NIE) from causal mediation analysis. Although this approach applies to a sequence-to-sequence model, only one token is considered at a time. It is, therefore, possible to also apply it to purely sequence-to-class models.

Given a model  $f(\mathbf{x};\theta)$ , mediation analysis is used to understand how a latent representation  $z(\mathbf{x};\theta)$  (called the mediator) affects the final model output. This latent representation can either be a single neuron or several neurons, like an attention head. The *Natural Indirect Effect* measures the effect that goes through this mediator.

To measure causality, an *intervention* on the concept measured must be made. As an intervention, Vig et al. [16] replace “nurse” with “the man” or “woman” for oppositely biased occupations. They call this replace operation **set-gender**.

Then to measure the effect of the mediator Vig et al. [16] introduce

$$\text{mediation-effect}_{m_1,z,m_2}(\mathbf{x};\theta) = \frac{\text{bias-effect}_{z(m_2(\mathbf{x});\theta)}(m_1(\mathbf{x});\theta)}{\text{bias-effect}(\mathbf{x};\theta)}, \quad (\text{A.11})$$

where  $m \in \{\text{identity}, \text{set-gender}\}$  and  $\text{bias-effect}_{z(m_2(\mathbf{x}))}(\cdot)$  is  $\text{bias-effect}(\cdot)$  but uses a modified model with the mediator values for  $z(m_1(\mathbf{x}))$  fixed. With this definition, the *Natural Indirect Effect* follows from causal mediation analysis literature [259].

$$\text{NIE}_z = \mathbb{E}_{\mathbf{x} \in \mathcal{D}} [\text{mediation-effect}_{\text{identity}, z, \text{set-gender}}(\mathbf{x}; \theta) - \text{mediation-effect}_{\text{identity}, z, \text{identity}}(\mathbf{x}; \theta)] \quad (\text{A.12})$$

Vig et al. [16] apply *Natural Indirect Effect* to a small GPT-2 model, where the mediator is an attention head. By doing this, Vig et al. [16] can identify which attention heads are most responsible for the gender bias, when considering the occupation concept. Hypothetical results, but results similar to those presented in Vig et al. [16], are presented in Figure A.4.

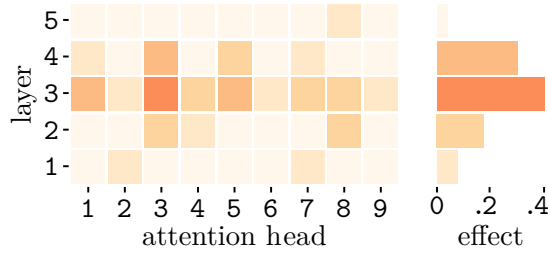


Figure A.4 Visualization of hypothetical *Natural Indirect Effect* (NIE) results, similar to Vig et al. [16]. Such visualization can reveal which attention-heads are responsible for gender bias in a small GPT-2 model. A stronger color indicates a higher NIE, meaning more responsible for the bias.

### A.3.2 Discussion

**Groundedness** As a new field, there is not much work on *groundedness*. Vig et al. [16] do not measure either *faithfulness* nor *human-groundedness* on *Natural Indirect Effect*. It is also not obvious how *faithfulness* could be measured. Note that this situation is not unique to *concept explanation*, as many other communication approaches also don't have an established measure of *faithfulness*.

**Future work** *Concept explanation* requires either a new dataset or annotation of an existing dataset. This can be quite expensive and impractical, especially when there is no concrete concept in mind, and the user wants a more exploratory explanation. However, new research is on discovering concepts automatically [260].

## A.4 Vocabulary

For this category, we define the term *vocabulary explanation* as methods that explain the whole model in relation to each word in the vocabulary and is, therefore, a *global explanation*.

In the sentiment classification context, a useful insight could be if positive and negative words are clustered together. Furthermore, perhaps there are words in those clusters that can not be considered of either positive or negative sentiment. Such a finding could indicate a bias in the dataset.

Because *vocabulary explanations* explain using the model’s vocabulary, they can often be applied to both sequence-to-class and sequence-to-sequence models. This is especially true for explanations based on the embedding matrix, which so is almost exclusively the case.

Because an embedding matrix is often used and because neural NLP models often use pre-trained word embeddings, most research on *vocabulary explanations* is applied to the pre-trained word embeddings [19, 261]. However, these explanation methods can also be applied to the word embeddings after training.

#### A.4.1 Projection

A common visual explanation is to project embeddings to two or three dimensions. This is particularly attractive, as word embeddings are of a fixed number of dimensions, and can therefore draw from the very rich literature on projection visualizations of tabular data, most notable is perhaps Principal Component Analysis (PCA) [17].

**t-SNE** Another popular and more recent method is t-SNE [18], which has been applied to word embeddings [38]. This method has, in particular, been attractive as it allows for non-linear transformations while still keeping points close to the word embedding and visualization spaces. t-SNE does this by representing the two spaces with two distance distributions; it then minimizes the KL-divergence by moving the points in the visualization space.

Note that Li et al. [38] does not go further to validate t-SNE in the context of word embeddings, except to highlight that words of similar semantic meaning are close together; we provide a similar example in Figure A.5.

**Supervised projection** A problem with using PCA and t-SNE is that they are unsupervised. Hence, while they might find a projection that offers high contrast, this projection might not correlate with what is of interest. An attractive alternative is, therefore, to define the projection, such that it reveals the subject of interest.

Bolukbasi et al. [262] are interested in how gender-biased a word is. They explore gender-bias by projecting each word onto a gender-specific vector and a gender-neutral vector. Such vectors can either be defined as the directional vector between “he” and “she”, or alternative. Bolukbasi et al. [262] also use multiple gender-specific pairs such as “daughter-son” and “herself-himself”, and then use their first Principal Component as a common projection vector.

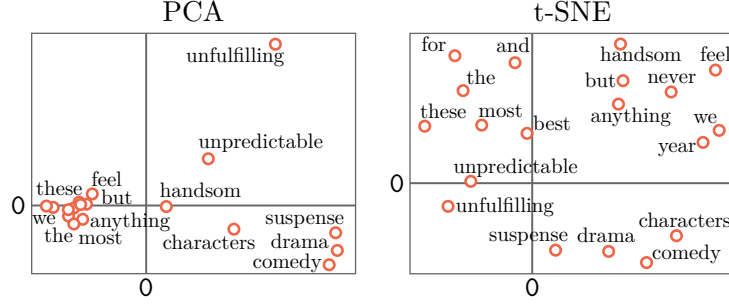


Figure A.5 PCA [17] and t-SNE [18] projection of GloVe [19] embeddings for the words in the semantic classification examples, as shown in Section 2.2 and elsewhere in the background chapter and appendix.

#### A.4.2 Rotation

The category of, for example, all positive sentiment words may have similar word embeddings. However, a particular basis dimension is unlikely to describe positive sentiment itself. Therefore, a useful interpretability method is to rotate the embedding space such that the basis dimensions in the new rotated embedding space represent significant concepts. This is distinct from *projection* methods because there is no loss of information as only a rotation is applied. Park et al. [146] perform such rotation using *Exploratory Factor Analysis (EFA)* [263]. The idea is to formalize a class of rotation matrices called the *Crawford-Ferguson Rotation Family* [264]. The parameters of this rotation formulation are then optimized to make the rotated embedding matrix only have a few large values in each row or column. As a hypothetical example, see Table A.1.

Basis-dimension	top-3 words
1	handsome, feel, unpredictable
2	most, best, anything
3	suspense, drama, comedy

Table A.1 Fictive example of the top-3 words for each basis-dimension in the rotated word embeddings.

Park et al. [146] validate this method to be *human-grounded* by using the *word intrusion* test. The classical word Intrusion test [145] provides 6 words to a human annotator, 5 of which should be semantically related, and the 6th is the intruder, which is semantically different. The human annotator then has to identify the intruder word. Importantly, semantic relatedness is, in this case, defined as the top 5 words of a given basis dimension in the rotated embedding matrix.

Unfortunately, rather than having humans detect the intruder, Park et al. [146] use a distance ratio, related to the cosine-distance, as the detector. This is problematic, as distance is directly related to how the semantically related words were chosen. In this case, the intruder should have been identified either by a human or an oracle model.

### A.4.3 Discussion

**Groundedness** In terms of *human-grounded*, *vocabulary explanation* are one of the few sub-fields that have a well-established test, namely the *word intrusion* test [145]. It is, therefore, hard to justify when methods in this category replace humans with an algorithm, as this largely invalidates the test.

**Future work** While past work, such as Latent Dirichlet Allocation (LDA) [265], have provided great *vocabulary explanations*, contemporary work using neural networks is quite limited and is mostly based on the embedding matrix. This is a pity, as the embedding matrix only provides a limited picture, and it is not hard to imagine using other information sources to create *vocabulary explanations*. For example, one could aggregate the word-contributions provided by *input feature* explanations.

## A.5 Ensemble

*Ensemble* explanations attempts to provide a *global explanation* by collecting multiple *local explanations*. This is done such that each *local explanation* represents the different modes of the model.

The extreme of this idea would be to provide a *local explanation* for every possible input, thereby providing a *global explanation*. Unfortunately, such an explanation is too much information for a human to understand and would not be *human-grounded*. As Miller [134] state, an explanation should be selective. The task of *ensemble* explanations is, therefore, to strategically select representative examples and their corresponding *local explanations*.

The assumption is that the model operates within different modes. Furthermore, that one example, or a few examples, from each mode can sufficiently represent the model’s entire behavior. For example, in the sentiment classification of movie reviews, a model may have one behavior for comments about the acting, another behavior for comments about the music score, etc.

*Ensemble* explanations is a very broad category of explanations, as for every type of *local explanation* method, there is an *ensemble* explanation could in principle be constructed. As

such, whether it can be applied to sequence-to-class or sequence-to-sequence models depends on the specific method. However, in practice, very few *ensemble* methods have been proposed, and most of them apply only to tabular data [105, 266, 267].

### A.5.1 Submodular Pick LIME (SP-LIME)

*SP-LIME* by Ribeiro et al. [41] attempts to select  $B$  observations (a budget), such that they represent the most important features based on their *LIME* explanation. Note that, while *LIME* explanations can be made for each output token and can therefore be used in a sequence-to-sequence context, *SP-LIME* do assume a sequence-to-class model.

*SP-LIME* calculates the importance of each feature  $v$ , by summing the absolute importance for all observations in the dataset; this total importance is  $\mathbf{I}_v$  in (A.13). The objective is then to maximize the sum of  $\mathbf{I}_v$  given a subset of features by strategically selecting  $B$  observations. Note that selecting multiple observations that represent the same features will not improve the objective. The specific objective is formalized in (A.13), which Ribeiro et al. [41] optimize greedily.

$$\mathbf{G}_{\text{SP-LIME}} = \underset{\tilde{\mathcal{D}} \text{ s.t. } |\tilde{\mathcal{D}}| \leq B}{\operatorname{argmax}} \sum_{v=1}^V \mathbb{1}_{[\exists \tilde{\mathbf{x}}_i \in \tilde{\mathcal{D}}: |\mathbf{E}_{\text{LIME}}(\tilde{\mathbf{x}}_i, \operatorname{argmax}_i p(i|\tilde{\mathbf{x}}_i; \theta))_v| > 0]} \mathbf{I}_v$$

where  $\tilde{\mathcal{D}} \subseteq \mathcal{D}$

$$\mathbf{I}_v = \sum_{\tilde{\mathbf{x}}_i \in \mathcal{D}} \left| \mathbf{E}_{\text{LIME}} \left( \tilde{\mathbf{x}}_i, \operatorname{argmax}_i p(i|\tilde{\mathbf{x}}_i; \theta) \right)_v \right| \quad (\text{A.13})$$

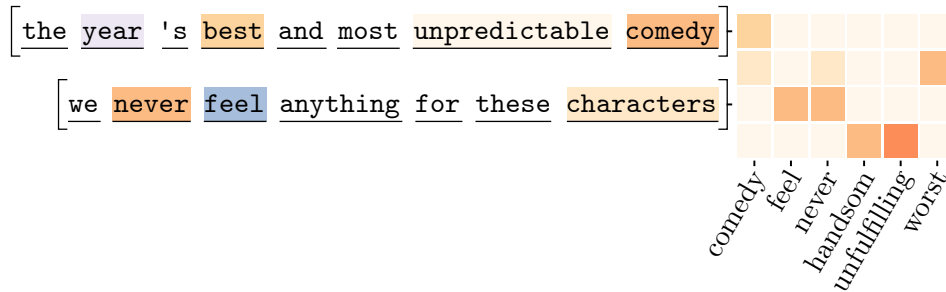


Figure A.6 Visualization of *SP-LIME* in a hypothetical setting. The matrix shows how each selected observation represents the different modes of the model. The left side shows two out of the four selected examples and their *LIME* explanation.

A major challenge with *SP-LIME* is that it requires computing a *LIME* explanation for every observation. Because each *LIME* explanation involves optimizing a logistic regression, this can be quite expensive. To reduce the number of observations that need to be explained,

Sangroya et al. [105] proposed using *Formal Concept Analysis* to strategically select which observations to explain. However, this approach has not yet been applied to NLP.

Ribeiro et al. [41] validate *SP-LIME* to be *human-grounded* by asking humans to select the best classifier, where a “wrong classifier” is trained on a biased dataset and a “correct classifier” is trained on a curated dataset. Ribeiro et al. [41] then compare *SP-LIME* with a random baseline, which simply selects random observations. From this experiment, they find that 89% of humans can select the best classifier using *SP-LIME*, where as only 75% can select the best classifier based on the random baseline.

### A.5.2 Discussion

**Groundedness** The *faithfulness* of *ensemble* explanations is very much dependent on the *faithfulness* of the *local explanation*. It is, therefore, difficult to imagine a general evaluation approach for *ensemble* explanations. However, even for *local explanations* with established validation *faithfulness* does not come for free, as also the selection algorithm also needs to be validated.

**Future work** As mentioned there is not much work using *ensemble* explanations. This is because when non-tabular data is used, comparing the selected explanations to ensure they represent different modes is more challenging. Even *SP-LIME* [41], which does apply to NLP tasks, uses a Bag-of-Word representation as a tabular proxy. Additionally, we can imagine that *ensemble* explanations are hard to scale as datasets increase and models get more complex with more modes.

That being said, we would be curious to see more work in this category. For example, an *ensemble* explanation which used a *influential example* method to show the overall most relevant observations.

## A.6 Linguistic Information

To validate that a natural language model does something reasonable, a popular approach is to attempt to align the model with the large body of linguistic theory that has been developed for hundreds of years. Because these methods summarize the model, they are a case of *global explanation*.

Methods in this category either probe by strategically modifying the input to observe the model’s reaction or show alignment between a latent representation and some linguistic representation. The former is called *behavioral probes* or *behavioral analysis*, the latter is

called *structural probes* or *structural analysis*. Which type of models these strategies apply to depends on the specific method. However, in general, *behavioral probes* applies primarily to sequence-to-class models, and *structural probes* applies to both sequence-to-class and sequence-to-sequence models.

One especially noteworthy subcategory of *Structural Probes* is *BERTology*, which specifically focuses on explaining the BERT-like models [4, 21, 54]. BERT’s popularity and effectiveness have resulted in countless papers in this category [20, 120, 121, 122, 268], hence the name *BERTology*. Some of the works use the attention of BERT and are therefore *intrinsic* explanations, while others simply probe the intermediate representations and are therefore *post-hoc* explanations.

Well-written survey papers already exist on *Linguistic Information* explanations. In particular, Belinkov et al. [119] cover *behavioral probes* and *structural probes*, Rogers et al. [120] discuss *BERTology*, and Belinkov and Glass [118] cover *structural probing* in detail. In this section, we will therefore not go in-depth but simply provide enough context to understand the field and, importantly, mention some of the criticisms that we believe have not been sufficiently highlighted by those surveys.

### A.6.1 Behavioral Probes

The research being done in *behavioral probes*, also called *behavioral analysis*, is not just for interpretability but also to measure the robustness and generalization ability of the model. For this reason, many *challenge datasets* are in the category of *behavioral analysis*. These datasets are meant to test the model’s generalization capabilities, often by containing many observations of underrepresented modes in the training datasets. However, the model’s performance on *challenge datasets* does not necessarily provide interpretability.

One of the initial papers providing interpretability via *behavioral probes* is that by Linzen et al. [269]. They probe a language model’s ability to reason about subject-verb agreement correctly. A recent work, by Clouatre et al. [123], Sinha et al. [270], finds that destroying syntax by shuffling words does not significantly affect a model trained on an NLI task, indicating that the model does not achieve natural language understanding.

As mentioned, this area of research is quite large and Belinkov et al. [119] cover *behavioral probes* in detail. Therefore, we just briefly discuss the work by McCoy et al. [11], which provides a particularly useful example of how *behavioral probes* can be used to provide interpretability.

McCoy et al. [11] look at Natural Language Inference (NLI), a task where a premise (for example, “The judge was paid by the actor”) and a hypothesis (for example, “The actor paid

the judge”) are provided, and the model should inform if these sentences are in agreement (called *entailment*). The other options are *contradiction* and *neutral*. McCoy et al. [11] hypothesise that models may not actually learn to understand the sentences but merely use heuristics to identify *entailment*.

They propose 3 heuristics based on the linguistic properties: lexical overlap, subsequence, and constituent. An example of lexical overlap is the premise “**The doctor** was **paid** by **the actor**” and hypothesis “The doctor paid the actor”. The proposed heuristic is that the model would classify this observation as *entailment* due to lexical overlap, even though this is not the correct classification.

To test for these heuristics, McCoy et al. [11] developed a dataset called HANS, which contains examples with these linguistic properties but do not have *entailment*. The results (Table A.2) validate the hypothesis that the model relies on these heuristics rather than a true understanding of the content. Had just an average score across all heuristics been provided, this would just be a robustness measure. However, by providing meta-information on which pattern each observation follows, the accuracy scores provide interpretability on where the model fails.

	Lexical Overlap	Subsequence	Constituent	Average
BERT [21]	17%	5%	17%	–
Human (Mechanical Turk)	–	–	–	77%

Table A.2 Performance on the HANS dataset provided by McCoy et al. [11]. Unfortunately, McCoy et al. [11] do not provide enough information to make a direct comparison possible. For comparison, BERT has 83% accuracy on MNLI [12], which was used for training.

In terms of *faithfulness*, McCoy et al. [11] perform no explicit evaluation. However, given that *behavioral probes* merely evaluate the model, *faithfulness* is generally not a concern. Furthermore, while McCoy et al. [11] do evaluate with humans, this is not a *human-grounded* evaluation. Because they only use humans to evaluate the dataset, not if the explanation itself is suitable to humans.

### A.6.2 Structural Probes

Probing methods primarily use a simple neural network, often just a logistic regression, to learn a mapping from an intermediate representation to a linguistic representation, such as the Part-Of-Speech (POS).

One of the early papers, by Shi et al. [271], analyzed the sentence-embeddings of a sequence-to-sequence LSTM, by looking at POS (part-of-speech), TSS (top-level syntactic sequence), SPC (the smallest phrase constituent for each word), tense (past or non-past), and voice (active or passive). Similarly, Adi et al. [272] used a multi-layer-perceptron (MLP) to analyze sentence embeddings for sentence length, word presence, and word order. More recently, Conneau et al. [124] have been using similar linguistic tasks and MLP probes but have extended previous analyses to multiple models and training methods.

Analog to these papers, a few methods use cluster algorithms instead of logistic regression [273]. Additionally, some methods only look at *word embeddings* [274]. The list of papers is very long, we suggest looking at the survey paper by Belinkov and Glass [118].

**BERTology** As an instructive example of probing in BERTology, the paper by Tenney et al. [20] is briefly described. Note that this is just one example of a vast number of papers. Rogers et al. [120] offer a much more comprehensive survey on BERTology.

Tenney et al. [20] probe a BERT model [21] by computing a learned weighted-sum  $\mathbf{z}_i(\mathbf{x}; \theta)$  for each intermediate representation  $\mathbf{h}_{l,i}(\mathbf{x}; \theta)$  of the token  $i$ , as described in (A.14).

$$\mathbf{z}_i(\mathbf{x}; \theta) = \gamma \sum_{l=1}^L s_l \mathbf{h}_{l,i}(\mathbf{x}; \theta) \quad (\text{A.14})$$

where  $\mathbf{s} = \text{softmax}(\mathbf{w})$

The weighted-sum  $\mathbf{z}_i(\mathbf{x})$  is then used by a classifier [275], and the weights  $s_l$ , parameterized by  $\mathbf{w}$ , describe how important each layer  $l$  is. The results can be seen in Figure A.7.

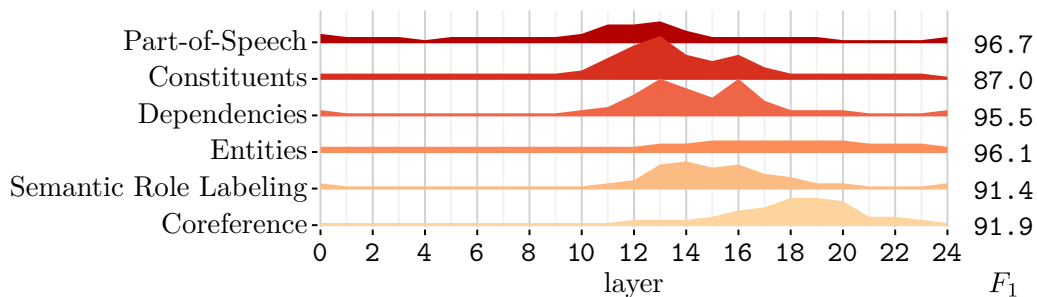


Figure A.7 Results by Tenney et al. [20] which shows how much each BERT [21] layer is used for each linguistic task. The  $F_1$  score for each task is also presented.

**Criticisms** A growing concern in probing methods is that given a sufficiently high-dimensional embedding, complex probe, and large auxiliary dataset, the probe can learn anything from

anything. If this concern is valid, it would mean that the probing methods do not provide *faithful* explanations [125].

Recent work attempts to overcome this concern by developing baselines. Zhang and Bowman [126] suggest learning a probe from an untrained model as a baseline. In that paper, they find probes can achieve high accuracy from an untrained model unless the auxiliary dataset size is dramatically decreased. Similarly, Hewitt and Liang [127] use randomized datasets as a baseline, called a control task. For example, with part-of-speech (POS), they assign a random POS-tag to each word, following the same empirical distribution of the non-randomized dataset. They find that equally high accuracy can be achieved on the randomized dataset unless the probe is made extraordinarily small.

**Information-Theoretic Probing** The solutions presented by Zhang and Bowman [126] and Hewitt and Liang [127] are useful. However, limiting the probe and dataset size could make finding complex hidden structures in the embeddings impossible.

Voita and Titov [128] attempt to overcome the criticism by a more principled approach, using information theory. More specifically, they measure the required complexity of the probe as a communication effort, called *Minimum Description Length* (MDL), and compare the MDL with a control task similar to Hewitt and Liang [127]. They find, similar to Hewitt and Liang [127], that the probes achieve similar accuracy on the probe dataset as on the control task. However, the control task is much harder to communicate (the MDL is higher), indicating that the probe is much more complex compared to training on the probe dataset.

### A.6.3 Discussion

**Groundedness** Considering the vast amount of research on *linguistic information* explanations, we find it worrying that there isn't more work on evaluating if these explanations are actually useful in terms of the *human-groundedness* and *faithfulness*. Without such evaluation, ensuring that the field of *linguistic information* explanations move in a productive direction is difficult.

**Future work** Considering the *groundedness* issues in *linguistic information* explanations, we advocate for more focus on *groundedness*. Voita and Titov [128] provide a great solution to how the *faithfulness* issues can be overcome. However, the field still lacks independent study on *human-groundedness* and *faithfulness*.

## A.7 Rules

*Rule* explanations attempt to explain the model by a simple set of rules. Therefore, they are an example of *global explanations*.

Reducing highly complex models like neural networks to a simple set of rules is likely impossible. Therefore, methods that attempt this simplify the objectivity by only explaining one particular aspect of the model.

Due to the challenges of producing rules, there is little research attempting it. We will present *Compositional Explanations of Neurons* [258] and *SEAR* [15].

### A.7.1 Semantically Equivalent Adversaries Rules (SEAR)

*SEAR* is an extension of the *Semantically Equivalent Adversaries* (SEA) method [15], where they developed a sampling algorithm for finding adversarial examples. Hence, the rule-generation objective is simplified, as only rules that describe what breaks the model need to be generated. Additionally, because *SEAR* uses an *adversarial examples* explanation, it only applies to sequence-to-class models.

		$p(y \mathbf{x}; \theta)$	$y$	Flips
$\mathbf{x}$	<u>the year 's best and most unpredictable comedy</u>	0.91	pos	-
$\tilde{\mathbf{x}}$	<u>the best and most unpredictable comedy this year</u>	0.13	-	-
rule	DET year 's $\rightarrow$ this year	-	-	1%
$\mathbf{x}$	<u>we never feel anything for these characters</u>	0.95	neg	-
$\tilde{\mathbf{x}}$	<u>we never empathize for these characters</u>	0.11	-	-
rule	<u>feel</u> $\rightarrow$ <u>empatize</u>	-	-	4%

Figure A.8 Hypothetical example showing rules which commonly break the model. The flip-rate describes how often these rules break the model.  $\mathbf{x}$  represents the original sentence and  $\tilde{\mathbf{x}}$  represents an adversarial example.

Ribeiro et al. [15] propose rules by simply observing individual word changes found by the *SEA* method discussed earlier, and then compute statistics on the bi-grams of the changed word and the Part of Speech of the adjacent word, Figure A.8 shows examples of this. If the proposed rule has a high success rate (called flip-rate) in terms of providing a semantically equivalent adversarial sample, it is considered a rule.

The authors validate this approach by asking experts to produce rules and then compare the success rate of human-generated rules and *SEAR*-generated rules. They find that the rules

generated by *SEAR* have a higher success rate.

### A.7.2 Compositional Explanations of Neurons

In *Compositional Explanations of Neurons* by Mu and Andreas [258], the rule generation problem is simplified by only relating the presence of input words to the activation of a single neuron.

The rules typically have the form of logical rules, meaning **not**, **and**, and **or**, where the booleans indicate a word is present, although Mu and Andreas [258] do not make any hard constraints here. For example, in an NLI task, there are indicators for POS presence and word overlap between the hypothesis and premise. If these rules are satisfied it means the neuron activation is above a defined threshold. For example, in a  $\text{ReLU}(\cdot)$  unit, one can threshold if its post-activation is above 0.

	IoU
$((\underline{\text{moving}} \text{ AND NOT } \underline{\text{house}}) \text{ OR } \underline{\text{feel}}) \text{ OR } \underline{\text{emotional}}$	13%
$((\underline{\text{best}} \text{ OR } \underline{\text{greatest}}) \text{ OR } \underline{\text{most}}) \text{ AND NOT } \underline{\text{bad}}$	21%

Figure A.9 Hypothetical example showing rules which activate a selected neuron. IoU is how often the rule activated the neuron, compared to cases where either the rule is true or the neuron activated (higher is better).

Given a dataset  $\mathcal{D}$ , a neuron activation  $z_n(\mathbf{x})$ , a threshold  $\tau$ , and an indicator function for the rule  $R(\mathbf{x})$ , the agreement between the rule and the neuron activation can be measured with the *Intersection over Union score*:

$$\text{IoU}(n, R) = \frac{\sum_{x \in \mathcal{D}} \mathbb{1}(z_n(\mathbf{x}) > \tau \wedge R(\mathbf{x}))}{\sum_{x \in \mathcal{D}} \mathbb{1}(z_n(\mathbf{x}) > \tau \vee R(\mathbf{x}))} \quad (\text{A.15})$$

For one particular neuron  $n$ , the combinatorial rule  $R$  is then constructed using beam-search, which stops at a pre-defined number of iterations. At each iteration, all feature indicator functions (e.g. word in  $\mathbf{x}$ ) and their negative, combined with the logical operators **and** and **or**, are scored using  $\text{IoU}(n, R)$ .

Unfortunately, Mu and Andreas [258] do not perform any *groundedness* validation of this approach. Furthermore, as the method only looks at the relation between the input and the neuron, it is unclear how much the selected neuron affects the output.

### A.7.3 Discussion

**Future work** As mentioned, there is little work on *rule* explanations. While this is definitely due to the inherent challenge, it is not too hard to imagine something like the Anchor method [276] be modified towards *global explanation*, in which case it would be a *rule* explanation.

**Groundedness** Because the category of *rule* explanations can be very diverse, *groundedness* evaluation would likely depend on the specific explanation method. However, generally *faithfulness* can be measured by asserting if the rule holds true by evaluating it on the dataset and compare with the model response. Additionally, *human-groundedness* can be evaluated by asking humans to predict the model’s output or choose the better model.

## APPENDIX B GENERAL-PURPOSE FAITHFULNESS METRIC FOR IMPORTANCE MEASURES

### B.1 Compute

In this section, we document the compute times and resources used for computing the results. Unfortunately, our compute infrastructure changed during the making of these results. The BiLSTM-attention results were computed on V100 GPUs while the RoBERTa results were computed on A100 GPUs. The A100 GPU is significantly faster than the V100 GPU, hence the compute times are not comparable across models. We could have recomputed the BiLSTM-attention results, but doing so would be a waste of resources. We report the machine specifications in Table B.1.

	BiLSTM-attention
CPU	4 cores, Intel Gold 6148 Sky-lake @ 2.4 GHz
GPU	1x NVidia V100 SXM2 (16 GB)
Memory	24 GB
	RoBERTa
CPU	6 cores, AMD Milan 7413 @ 2.65 GHz 128M cache L3
GPU	1x NVidia A100 (40 GB)
Memory	24 GB

Table B.1 Compute hardware used for each model. Note, the models were computed on a shared user system. Hence, we only report the resources allocated for our jobs.

The compute times are reported in Table B.2. All compute was done using 99% hydroelectric energy.

While the totals in Table B.2 may be large, in partial situations only one dataset is usually considered. Additionally, the variance in Figure 3.4 is quite low, making less seeds an option. Finally, the compute time of *integrated gradient* is approximately 2/3 of the total. As discussed in Section 3.5, this is rarely worth it. Practical settings may want to not consider *integrated gradient* at all for this reason.

Dataset	Importance	Walltime [hh:mm]	
	Measure	LSTM	RoBERTa
Anemia	Random	00:09	00:03
	Attention	00:09	–
	Gradient	00:11	00:04
	Input times Gradient	00:11	00:04
	Integrated Gradient	00:44	00:27
Diabetes	Random	00:17	00:05
	Attention	00:17	–
	Gradient	00:23	00:07
	Input times Gradient	00:23	00:07
	Integrated Gradient	01:46	01:09
IMDB	Random	00:05	00:08
	Attention	00:05	–
	Gradient	00:05	00:10
	Input times Gradient	00:05	00:10
	Integrated Gradient	00:20	02:10
SNLI	Random	00:49	01:03
	Attention	00:46	–
	Gradient	00:48	01:28
	Input times Gradient	00:48	01:10
	Integrated Gradient	01:09	05:41
SST	Random	00:02	00:02
	Attention	00:02	–
	Gradient	00:02	00:02
	Input times Gradient	00:02	00:02
	Integrated Gradient	00:03	00:06
bAbI-1	Random	00:08	00:04
	Attention	00:09	–
	Gradient	00:08	00:04
	Input times Gradient	00:08	00:04
	Integrated Gradient	00:10	00:11
bAbI-2	Random	00:12	00:06
	Attention	00:12	–
	Gradient	00:12	00:06
	Input times Gradient	00:12	00:06
	Integrated Gradient	00:15	00:32
bAbI-3	Random	00:24	00:11
	Attention	00:25	–
	Gradient	00:25	00:13
	Input times Gradient	00:25	00:13
	Integrated Gradient	00:32	01:12
<b>Total</b>	sum	13:38	17:20
	x9 iterations (approx.)	5 days	6.5 days
	x5 seeds (approx.)	25.5 days	32.5 days

Table B.2 Compute times for each model and importance measure combination. Note, there is no need to compute models for each importance measure at 0% and 100% masking. Hence, we report for 9 iterations.

## B.2 Sparsity

In this section, we analyse the sparsity of each importance measure. While none of the importance measures produce an actual importance for any token, they may have most of the importance assigned to just a few tokens.

This analysis serves two purposes, to show that masking a relative number of tokens is justified and to test if any importance measure are more sparse than others.

**Masking a relative number of tokens is justified.** If the majority of the importance is assigned to just a few tokens (e.g. 10 tokens have 99% of the total importance scores), then it would make more sense to perform the non-approximate version of Recursive ROAR where exactly one token is masked in each iteration.

In Figure B.1, we look at the sparsity considering the top-10 tokens. We find that the sparsity is not sufficiently high to justify masking exactly one token in each iteration. For completeness, we include this analysis in Appendix B.

There are cases where masking exactly one token in each iteration could make sense, for example, for *attention* in bAbI. However, as this is a comparative study among several importance measures and datasets, this is not enough.

**Attention is more sparse than others importance measures** If a particular importance measure is more sparse than others, while having a similar faithfulness, then the more sparse importance measure would be preferable. This is because it is more likely to be understandable to humans [134].

In Figure B.2, we look at the sparsity considering a relative number of tokens. We find that for some datasets, in particular bAbI, attention is the most sparse importance measure. Besides this, integrated gradient is usually the most sparse in nearly all cases. However, while the difference in sparsity is often statistically significant we speculate that the difference is not large enough to cause a difference in practical settings.

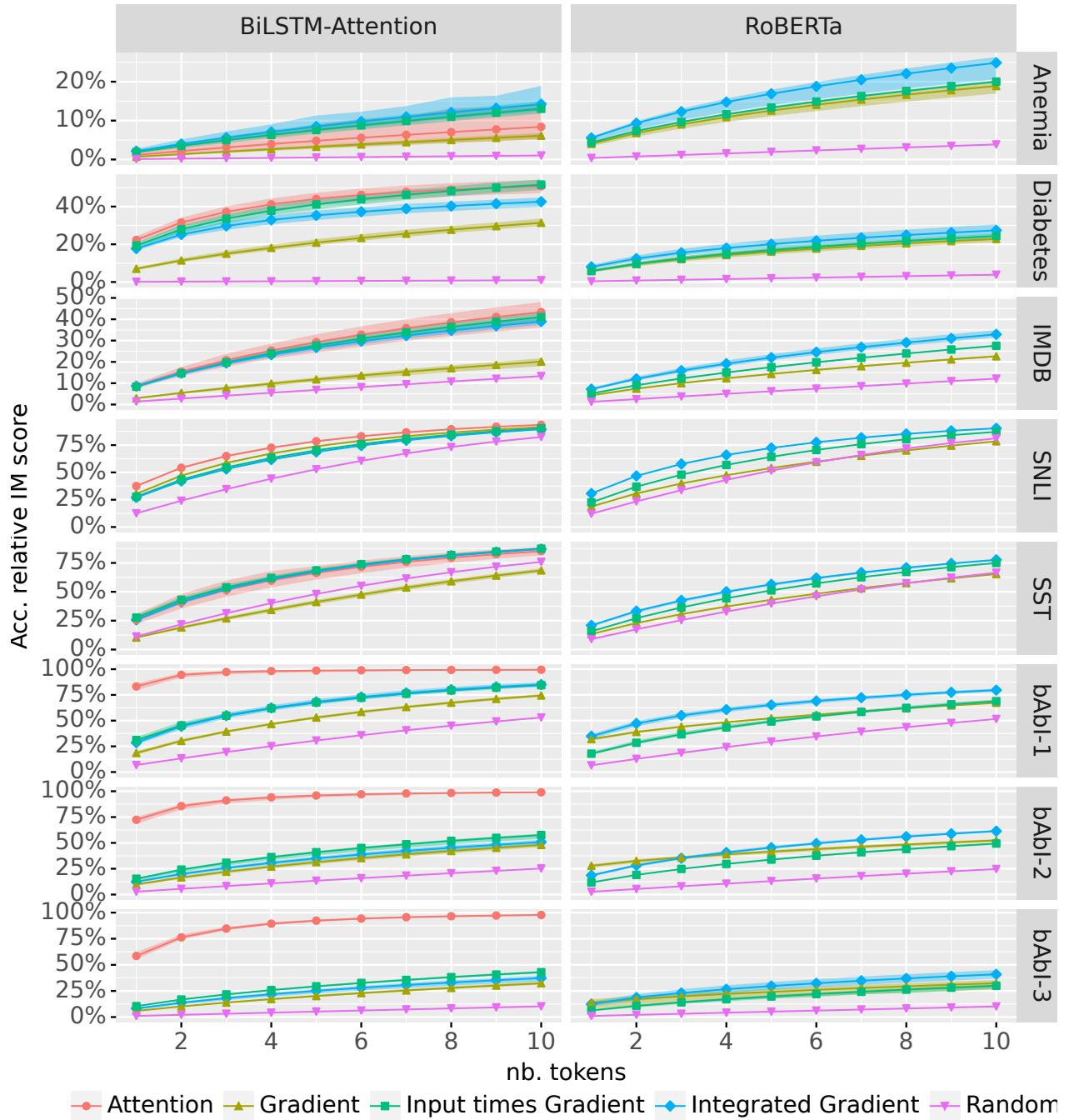


Figure B.1 Shows the accumulative importance score relative to the total importance score, for the top-k number of tokens. The metric is averaged over 5 seeds with a 95% confidence interval. Note that datasets are not equal in sequence-length, the scores are therefore hard to compare across datasets.

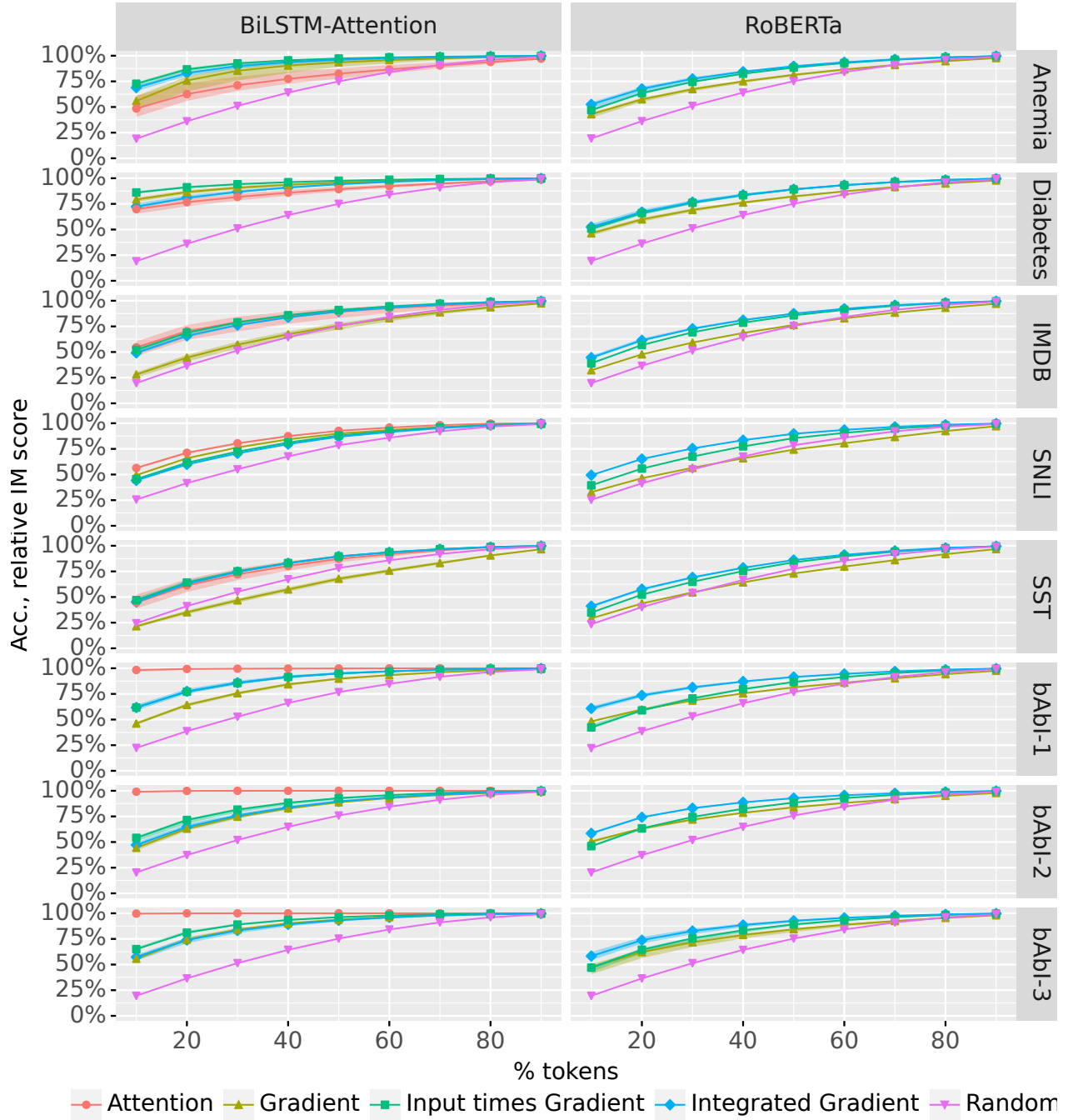


Figure B.2 The accumulative importance score relative to the total importance score for the top-x% number of tokens. The metric is averaged over 5 seeds with a 95% confidence interval.

### B.3 Recursive ROAR with a stepsize of one token

To analyze the effect of masking 10%, as opposed to masking exactly one token in each iteration, we perform the Recursive ROAR experiment with exactly one token masked. The results are in Figure B.3. Because this is computationally expensive, we only do this for up to 10 tokens. This makes it harder to draw clear conclusions from this experiment, in particular, because not all redundancies are removed when only masking 10 tokens.

In general, the results in Figure B.3 show that the approximation of masking 10% in each iteration does affect the results. However, we can draw the same conclusions. That being said, some of the conclusions are less obvious because we only look at 10 tokens.

#### B.3.1 The results are affected by the approximation

Looking just at RoBERTa, for Diabetes, *Integrated Gradient* yields 65% performance at 10% masking (approximately 51 tokens), while *Integrated Gradient* yields 55% performance at 10 tokens. Similarly for bAbI-3, *Gradient* yields 65% at 10% masking (approximately 30 tokens), while *Gradient* yields 30% at 10 tokens. Both of these cases show that a lower performance is achieved earlier when masking one token in each iteration.

This is to be expected, as masking one token in each iteration is more effective for removing redundancies. Were we to complete the experiment to eventually mask all tokens, the faithfulness scores can therefore be expected to be higher.

#### B.3.2 The conclusions are the same

In Section 3.5, we present 5 findings. Here, we briefly show that the same conclusions can be drawn from Figure B.3. However, as only 10 tokens are masked they may be less obvious and there may be less evidence.

**Faithfulness is model-dependent.** Yes, this is most clearly seen for IMDB, where BiLSTM-Attention archives significantly lower performance (higher faithfulness) compared to RoBERTa.

**Faithfulness is task-dependent.** Yes, looking at BiLSTM-Attention, for IMDB *Integrated Gradient* is the worst importance measure. However, for the bAbI tasks *Integrated Gradient* is among the best importance measures.

**Attention can be faithful.** Yes, particularly for bAbI, IMDB, and Diabetes attention is faithful.

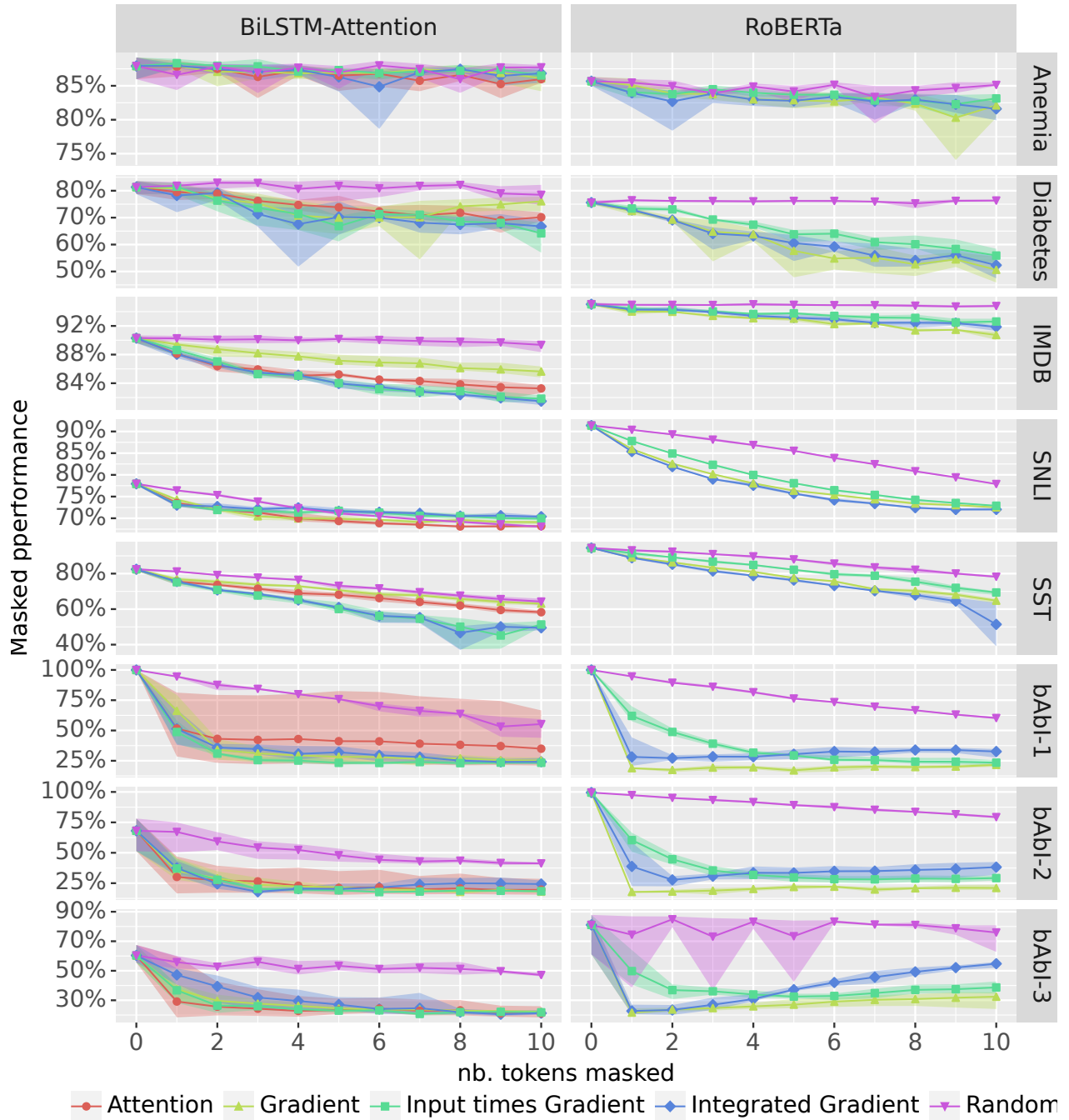


Figure B.3 Recursive ROAR results, showing model performance at up to 10 tokens masked. Note that because the datasets have more than 10-tokens, the conclusion one can draw from this plot may change if more tokens were considered. However, in general, a model performance below *random* indicates faithfulness, while above or similar to *random* indicates a non-faithful importance measure. Performance is averaged over 5 seeds with a 95% confidence interval.

**Integrated Gradient is not necessarily more faithful than Gradient or Input times Gradient.** Yes, considering BiLSTM-Attention, IMDB *Integrated Gradient* is significantly worse than other explanations. For most datasets, *Integrated Gradient* has similar faithfulness as other importance measures.

**Importance measures often work best for the top-20% most important tokens.** As Figure B.3 only shows 10 tokens, which is usually below top-20% this is hard to comment on.

**Class leakage can cause the model performance to increase.** For RoBERTa, in bAbI-3, the *Integrated Gradient* importance measure can be seen to increase performance after 2 tokens are masked.

#### B.4 ROAR vs Recursive ROAR

As an ablation study, we compare ROAR by Hooker et al. [14] with our Recursive ROAR. Figure B.4 shows the comparison for BiLSTM-Attention and Figure B.5 shows the comparison for RoBERTa. Recall that for ROAR by Hooker et al. [14] it is not possible to say that an importance measure is not faithful.

**Some datasets have redundancies that affect ROAR.** In particular, we find that Diabetes shows a significant difference when comparing ROAR with Recursive ROAR. This is both for BiLSTM-Attention (Figure B.4) and RoBERTa (Figure B.5). For both models, *Gradient* and *Input times Gradient* become faithful with Recursive ROAR. Additionally, for RoBERTa the same is the case for *Integrated Gradient*. This is not surprising, as Diabetes contains incredibly long sequences and contains redundancies.

Also, for IMDB, and to a lesser extent SST, there is a clear difference between BiLSTM-Attention and RoBERTa. This too is not surprising, as sentiment can often be inferred from just a single word. However, there are likely to be many positive or negative words in each observation.

**Class leakage affects both ROAR and Recursive ROAR.** We observe the class leakage issue for ROAR in SNLI with BiLSTM-Attention and for the bAbI tasks with RoBERTa. We observe the issue for Recursive ROAR in bAbI with BiLSTM-Attention. The fact that the issue mostly exists with bAbI is somewhat encouraging, as the bAbI datasets are synthetic. The class leakage issue appears to affect real datasets less.

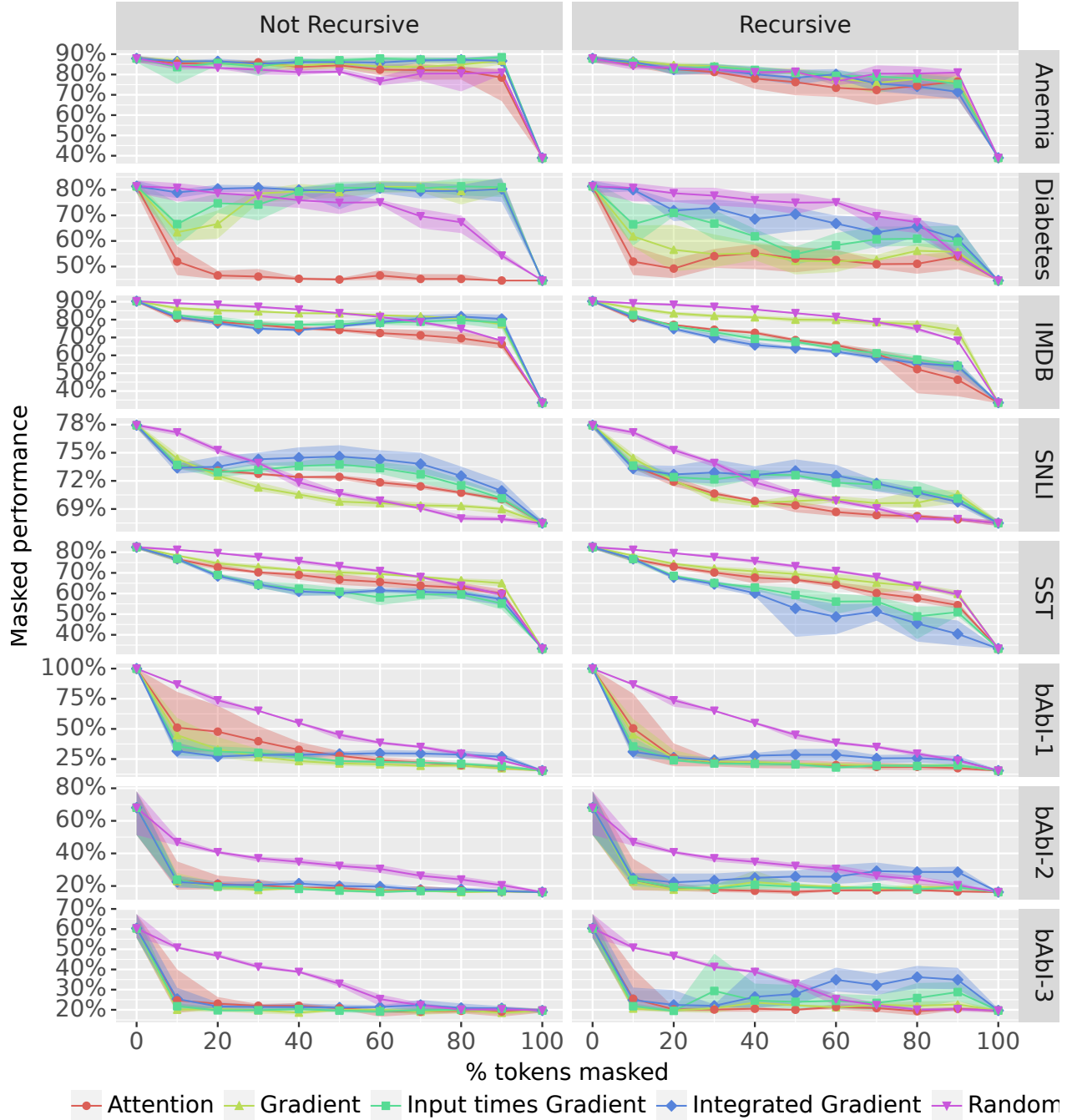


Figure B.4 ROAR and Recursive ROAR results for **BiLSTM-Attention**, showing model performance at x% of tokens masked. A model performance below *random* indicates faithfulness. For Recursive ROAR a curve above or similar to *random* indicates a non-faithful importance measure, while for ROAR by Hooker et al. [14] this case is inconclusive. Performance is averaged over 5 seeds with a 95% confidence interval.

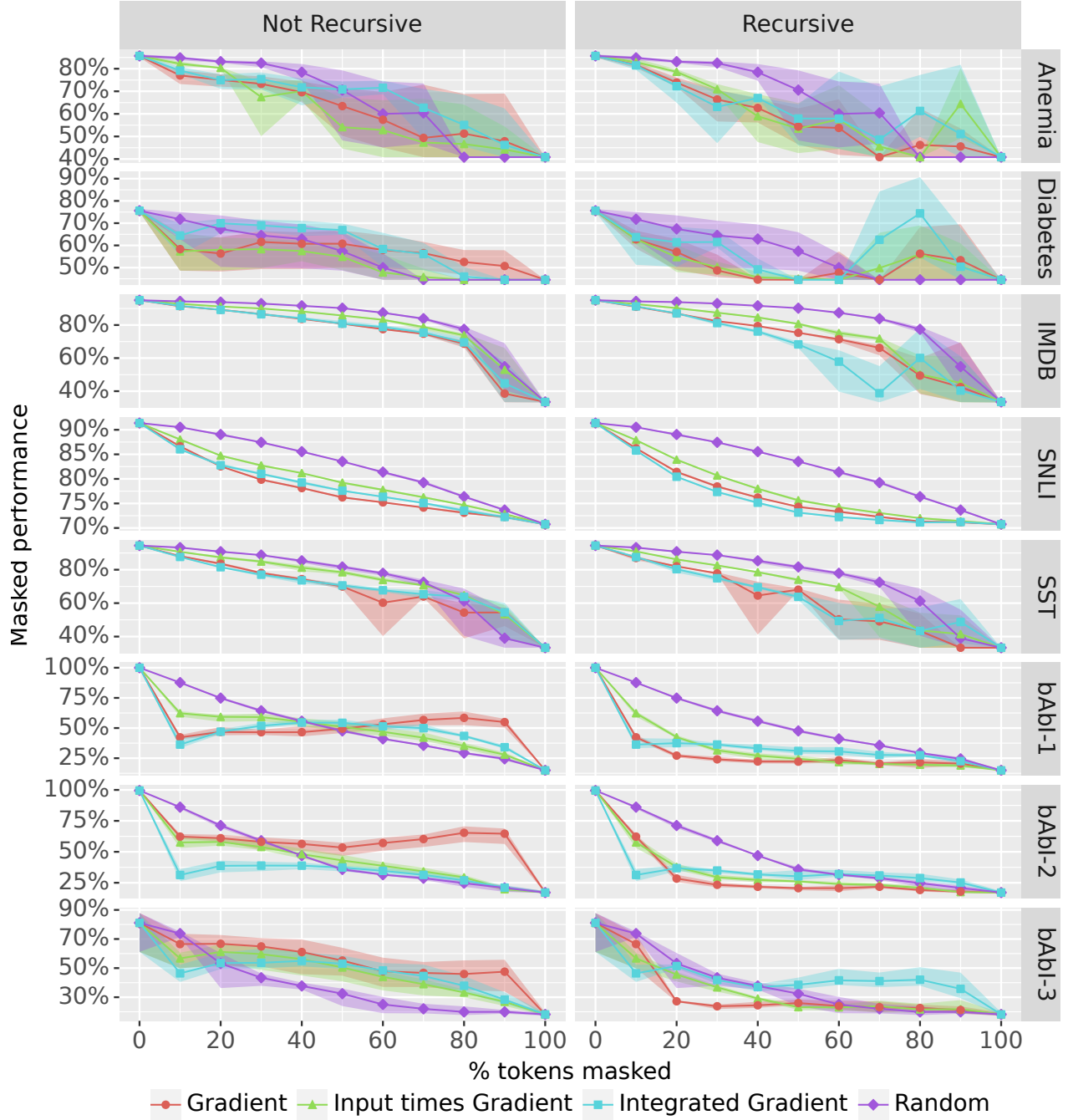


Figure B.5 ROAR and Recursive ROAR results for **RoBERTa**, showing model performance at  $x\%$  of tokens masked. A model performance below *random* indicates faithfulness. For Recursive ROAR a curve above or similar to *random* indicates a non-faithful importance measure, while for ROAR by Hooker et al. [14] this case is inconclusive. Performance is averaged over 5 seeds with a 95% confidence interval.

## APPENDIX C FAITHFULNESS MEASURABLE MODELS

### C.1 Compute

This section reports the compute resources and requirements. The compute hardware specifications are in Table C.1 and were the same for all experiments. All computing was performed using 99% hydroelectric power.

Table C.1 The computing hardware used. Note, that a shared user system were used, only the allocated resources are reported.

CPU	12 cores, Intel Silver 4216 Cascade Lake @ 2.1GHz
GPU	1x NVidia V100 (32G HBM2 memory)
Memory	24 GB

Note that the importance measures computed are the same for both the faithfulness results and the out-of-distribution results. Hence, these do not need to be computed twice. Additionally, the beam-search method is in itself recursive, so this was only computed for 0% masking.

#### C.1.1 Implementation

We use the HuggingFace implementation of RoBERTa and the TensorFlow framework. The code is available at <https://github.com/AndreasMadsen/faithfulness-measurable-models>.

#### C.1.2 Walltimes

We here include the walltimes for all experiments.

- Table C.2 shows wall-times for the masked fine-tuning.
- Table C.3 shows wall times for the in-distribution validation, not including importance measures.
- Table C.4 shows wall-times for the faithfulness evaluation, not including importance measures.
- Table C.5 shows wall-times for the importance measures.

Dataset	Walltime [hh:mm]	
	RoBERTa- base	RoBERTa- large
BoolQ	00:51	02:03
CB	00:06	00:12
CoLA	00:17	00:33
IMDB	01:44	04:02
Anemia	00:48	02:04
Diabetes	01:34	04:04
MNLI	06:39	14:47
MRPC	00:12	00:27
QNLI	04:03	09:12
QQP	05:13	11:52
RTE	00:18	00:43
SNLI	04:57	10:38
SST2	01:19	02:44
bAbI-1	00:27	01:01
bAbI-2	00:50	02:05
bAbI-3	01:43	04:28
sum	01:43	04:28
x5 seeds	08:37	22:21

Table C.2 Walltime for fine-tuning. Masked fine-tuning does not affect the total wall time in our setup.

Dataset	Walltime [hh:mm]	
	RoBERTa- base	RoBERTa- large
BoolQ	00:04	00:09
CB	00:01	00:02
CoLA	00:02	00:04
IMDB	00:15	00:44
Anemia	00:01	00:03
Diabetes	00:01	00:04
MNLI	00:20	00:57
MRPC	00:02	00:03
QNLI	00:05	00:12
QQP	00:47	02:13
RTE	00:02	00:04
SNLI	00:04	00:09
SST2	00:09	00:25
bAbI-1	00:01	00:03
bAbI-2	00:02	00:05
bAbI-3	00:01	00:03
sum	02:06	05:25
x5 seeds	10:30	27:09

Table C.3 Walltime for in-distribution validation. This does not include importance measure calculations. See Table C.5.

Dataset	Walltime [hh:mm]	
	RoBERTa- base	RoBERTa- large
BoolQ	00:02	00:05
CB	00:00	00:01
CoLA	00:01	00:02
IMDB	00:13	00:39
Anemia	00:01	00:02
Diabetes	00:01	00:03
MNLI	00:02	00:05
MRPC	00:01	00:02
QNLI	00:01	00:04
QQP	00:05	00:13
RTE	00:01	00:02
SNLI	00:01	00:03
SST2	00:01	00:01
bAbI-1	00:00	00:01
bAbI-2	00:01	00:02
bAbI-3	00:00	00:02
sum	00:38	01:34
x5 seeds	03:13	07:51

Table C.4 Walltime for faithfulness evaluation. This does not include importance measure calculations. See Table C.5.

Dataset	IM	Walltime [hh:mm]	
		RoBERTa-base	RoBERTa-large
bAbI-1	Beam	00:54	02:24
	Grad ( $L_1$ )	00:01	00:04
	Grad ( $L_2$ )	00:02	00:04
	x $\odot$ grad (abs)	00:01	00:04
	x $\odot$ grad (sign)	00:01	00:04
	IG (abs)	00:04	00:12
	IG (sign)	00:04	00:11
	LOO (abs)	00:24	00:49
	LOO (sign)	00:24	00:49
	Random	00:00	00:00
bAbI-2	Beam	20:56	61:24
	Grad ( $L_1$ )	00:02	00:06
	Grad ( $L_2$ )	00:02	00:05
	x $\odot$ grad (abs)	00:02	00:05
	x $\odot$ grad (sign)	00:02	00:05
	IG (abs)	00:10	00:29
	IG (sign)	00:10	00:29
	LOO (abs)	00:39	01:26
	LOO (sign)	00:39	01:26
	Random	00:00	00:00
bAbI-3	Beam	–	–
	Grad ( $L_1$ )	00:02	00:05
	Grad ( $L_2$ )	00:02	00:05
	x $\odot$ grad (abs)	00:01	00:04
	x $\odot$ grad (sign)	00:01	00:04
	IG (abs)	00:18	00:54
	IG (sign)	00:18	00:53
	LOO (abs)	01:09	03:18
	LOO (sign)	01:09	03:18
	Random	00:00	00:00
BoolQ	Beam	00:33	01:22
	Grad ( $L_1$ )	00:05	00:11
	Grad ( $L_2$ )	00:05	00:11
	x $\odot$ grad (abs)	00:04	00:10
	x $\odot$ grad (sign)	00:04	00:10
	IG (abs)	00:39	01:48
	IG (sign)	00:39	01:49
	LOO (abs)	00:16	00:38
	LOO (sign)	00:16	00:38
	Random	00:00	00:00
CB	Beam	00:45	02:09
	Grad ( $L_1$ )	00:01	00:03
	Grad ( $L_2$ )	00:01	00:03
	x $\odot$ grad (abs)	00:01	00:03
	x $\odot$ grad (sign)	00:01	00:03
	IG (abs)	00:01	00:04
	IG (sign)	00:01	00:04
	LOO (abs)	00:09	00:19
	LOO (sign)	00:09	00:19
	Random	00:00	00:00
CoLA	Beam	00:11	00:19
	Grad ( $L_1$ )	00:02	00:05
	Grad ( $L_2$ )	00:02	00:04
	x $\odot$ grad (abs)	00:02	00:04
	x $\odot$ grad (sign)	00:02	00:04
	IG (abs)	00:04	00:09
	IG (sign)	00:04	00:09
	LOO (abs)	00:09	00:18
	LOO (sign)	00:09	00:18
	Random	00:00	00:00
Anemia	Beam	–	–
	Grad ( $L_1$ )	00:02	00:06
	Grad ( $L_2$ )	00:02	00:06
	x $\odot$ grad (abs)	00:01	00:05
	x $\odot$ grad (sign)	00:01	00:05
	IG (abs)	00:23	01:08
	IG (sign)	00:23	01:08
	LOO (abs)	02:23	06:58
	LOO (sign)	02:23	07:01
	Random	00:00	00:00
Diabetes	Beam	–	–
	Grad ( $L_1$ )	00:03	00:07
	Grad ( $L_2$ )	00:03	00:07
	x $\odot$ grad (abs)	00:02	00:06
	x $\odot$ grad (sign)	00:02	00:06
	IG (abs)	00:32	01:34
	IG (sign)	00:32	01:34
	LOO (abs)	03:19	09:44
	LOO (sign)	03:17	09:45
	Random	00:00	00:00
MRPC	Beam	00:14	00:33
	Grad ( $L_1$ )	00:02	00:04
	Grad ( $L_2$ )	00:02	00:04
	x $\odot$ grad (abs)	00:02	00:04
	x $\odot$ grad (sign)	00:02	00:04
	IG (abs)	00:03	00:07
	IG (sign)	00:03	00:07
	LOO (abs)	00:08	00:17
	LOO (sign)	00:08	00:17
	Random	00:00	00:00
RTE	Beam	01:32	04:26
	Grad ( $L_1$ )	00:02	00:04
	Grad ( $L_2$ )	00:02	00:04
	x $\odot$ grad (abs)	00:02	00:04
	x $\odot$ grad (sign)	00:02	00:04
	IG (abs)	00:04	00:09
	IG (sign)	00:04	00:09
	LOO (abs)	00:10	00:22
	LOO (sign)	00:10	00:22
	Random	00:00	00:00
SST2	Beam	00:18	00:43
	Grad ( $L_1$ )	00:02	00:04
	Grad ( $L_2$ )	00:02	00:04
	x $\odot$ grad (abs)	00:02	00:04
	x $\odot$ grad (sign)	00:02	00:04
	IG (abs)	00:04	00:09
	IG (sign)	00:04	00:09
	LOO (abs)	00:09	00:19
	LOO (sign)	00:10	00:19
	Random	00:00	00:00
SNLI	Beam	01:10	02:38
	Grad ( $L_1$ )	00:05	00:07
	Grad ( $L_2$ )	00:06	00:07
	x $\odot$ grad (abs)	00:05	00:06
	x $\odot$ grad (sign)	00:04	00:06
	IG (abs)	00:21	00:57
	IG (sign)	00:21	00:56
	LOO (abs)	00:12	00:26
	LOO (sign)	00:12	00:26
	Random	00:01	00:00
IMDB	Beam	–	–
	Grad ( $L_1$ )	00:34	01:18
	Grad ( $L_2$ )	00:34	01:17
	x $\odot$ grad (abs)	00:22	01:03
	x $\odot$ grad (sign)	00:22	01:03
	IG (abs)	06:49	20:08
	IG (sign)	06:54	20:09
	LOO (abs)	25:02	73:17
	LOO (sign)	24:48	72:55
	Random	00:01	00:01
MNLI	Beam	05:44	15:34
	Grad ( $L_1$ )	00:05	00:11
	Grad ( $L_2$ )	00:05	00:11
	x $\odot$ grad (abs)	00:04	00:09
	x $\odot$ grad (sign)	00:04	00:09
	IG (abs)	00:35	01:35
	IG (sign)	00:35	01:34
	LOO (abs)	00:19	00:46
	LOO (sign)	00:19	00:46
	Random	00:00	00:00
QNLI	Beam	06:39	18:51
	Grad ( $L_1$ )	00:04	00:08
	Grad ( $L_2$ )	00:04	00:08
	x $\odot$ grad (abs)	00:03	00:07
	x $\odot$ grad (sign)	00:03	00:08
	IG (abs)	00:23	01:03
	IG (sign)	00:23	01:04
	LOO (abs)	00:17	00:43
	LOO (sign)	00:17	00:43
	Random	00:00	00:00
QQP	Beam	04:44	11:12
	Grad ( $L_1$ )	00:12	00:26
	Grad ( $L_2$ )	00:12	00:26
	x $\odot$ grad (abs)	00:10	00:22
	x $\odot$ grad (sign)	00:10	00:22
	IG (abs)	01:48	04:57
	IG (sign)	01:48	04:59
	LOO (abs)	00:36	01:24
	LOO (sign)	00:36	01:23
	Random	00:01	00:01
sum		145:00	406:58
x5 seeds		725:02	2034:53

Table C.5 Walltime for importance measures. Note that because the beam-search method (Beam) scales quadratic with the sequence-length, it is not feasible to compute for all datasets.

## C.2 Masked fine-tuning

In Section 4.2.1, we show selected results for unmasked performance and 100% masked performance. In this appendix, we extend those results to all 16 datasets. In addition to this, this appendix contains a more detailed ablation study, where the training strategy and validation strategy are considered separate. As such, the results in Section 4.2.1 are a strict subset of these detailed results. In Table C.6 we show how the terminologies relate.

Section 4.2.1	Training strategy	Validation strategy
Masked fine-tuning	Use 50/50	Use both
Plain fine-tuning	No masking	No masking
Only masking	Masking	Masking

Table C.6 This table relates terminologies between the fine-tuning strategies mentioned in Section 4.2.1 and the training strategy and validation strategy terms.

**Training strategy** The training strategy applies to the training dataset during fine-tuning.

**No masking** No masking is applied to the training dataset. This is what is ordinarily done in the literature.

**Masking** Masking is applied to every observation. The masking is uniformly sampled, at a masking rate between 0% and 100%.

**Use 50/50** Half of the mini-batch using the *No masking* strategy and the other half use the *Masking* strategy.

**Validation strategy** The validation dataset is used to select the optional epoch. This is similar to early stopping, but rather than stopping immediately. The training continues, and the best epoch is chosen at the end of the training.

The validation strategy applies to the validation dataset during fine-tuning.

**No masking** No masking is applied to the validation dataset. This is what is ordinarily done in the literature.

**Masking** Masking is applied to every observation. The masking is uniformly sampled, at a masking rate between 0% and 100%.

**Use both** A copy of the validation dataset has the *No masking* strategy applied to it. Another copy of the validation dataset has the *Masking* strategy applied to it. As such, the validation dataset is twice as long, but it does not add additional observations or information.

### C.2.1 Findings

We generally find that the choice of validation strategy when using the *Use 50/50* training strategy is not important. Interestingly, *Masking* for the validation dataset and *No masking* for the training dataset often works too.

However, because *Use 50/50* for training strategy and *Use both* for validation strategy, i.e. masked fine-tuning, work well in all cases and is theoretically sound, this is the approach we recommend and use throughout Chapter 4.

### C.2.2 All datasets aggregation

In Section 4.2.1 we also include an *All* “dataset”. This is a simple arithmetic mean over all the performance of all 16 datasets. This is similar to how the GLUE benchmark [7] works. To compute the confidence interval, a dataset-aggregation is done for each seed, such that the all-observation are i.i.d..

Because some seeds do not converge for some datasets, such as bAbI-2 and bAbI-3 (as mentioned in Section 4.2.1), those outliers and not included in the aggregation, also hyperparameter optimization will likely help. For complete transparency, we do include them in the statistics for the individual datasets and show all individual performances with a (+) symbol.

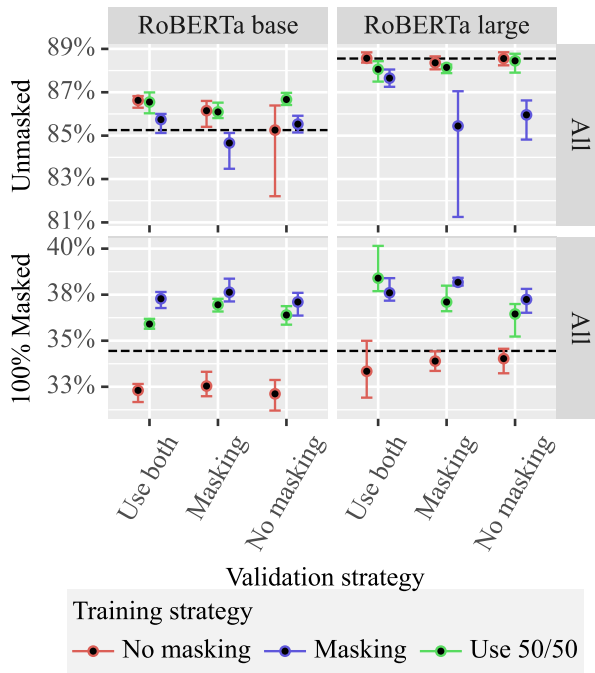


Figure C.1 The all aggregation for the 100% masked performance and unmasked performance. The baseline (dashed line) for 100% masked performance is the class-majority baseline. Unmasked performance is when using no masking for both validation and training.

### C.2.3 Test dataset

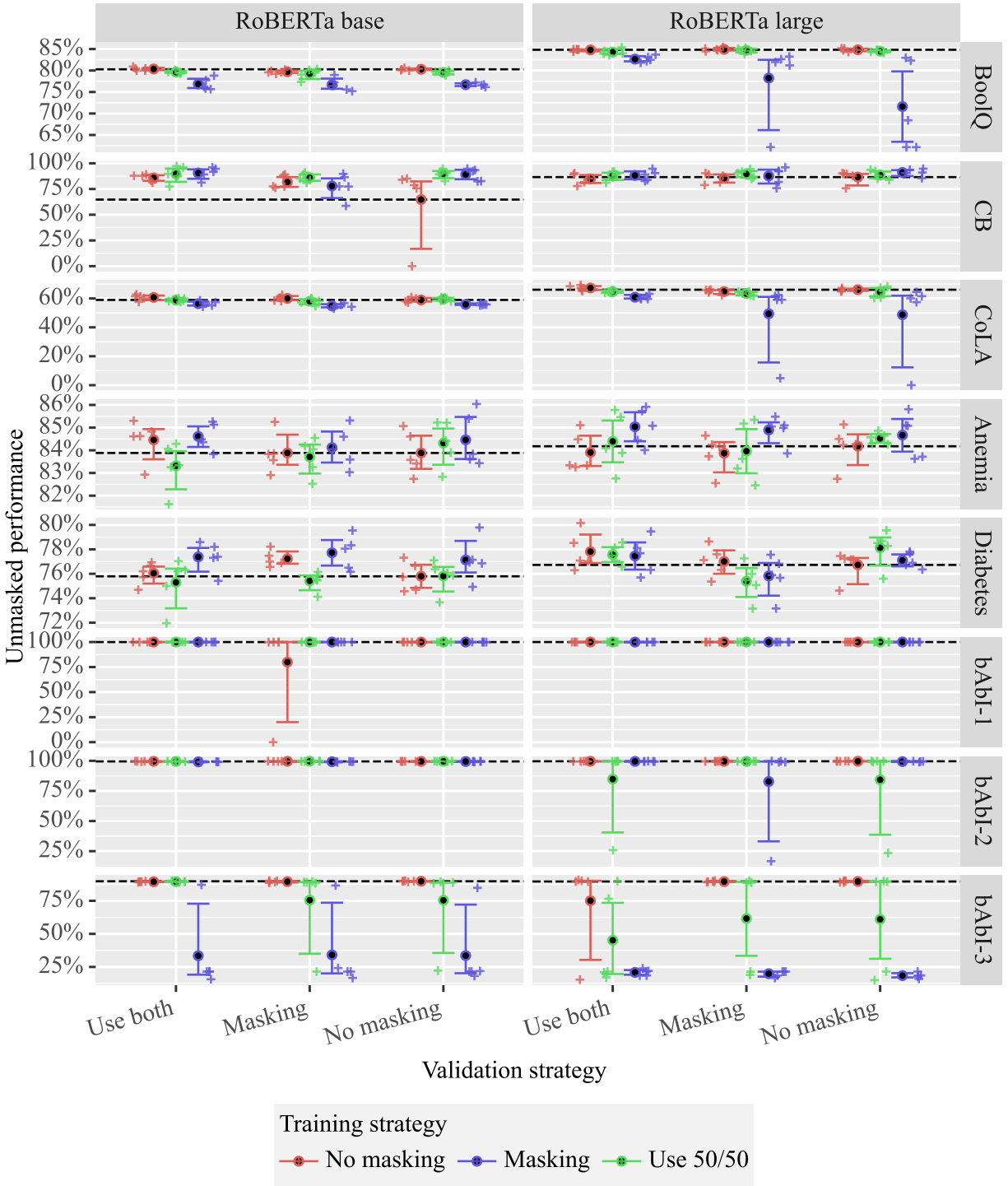


Figure C.2 The unmasked performance for each validation and training strategy, using the test dataset. Not that “No masking” as a *training strategy* is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task *training strategy* “Use 50/50” works best. This plot is **page-1**. Corresponding main results in Figure 4.3.

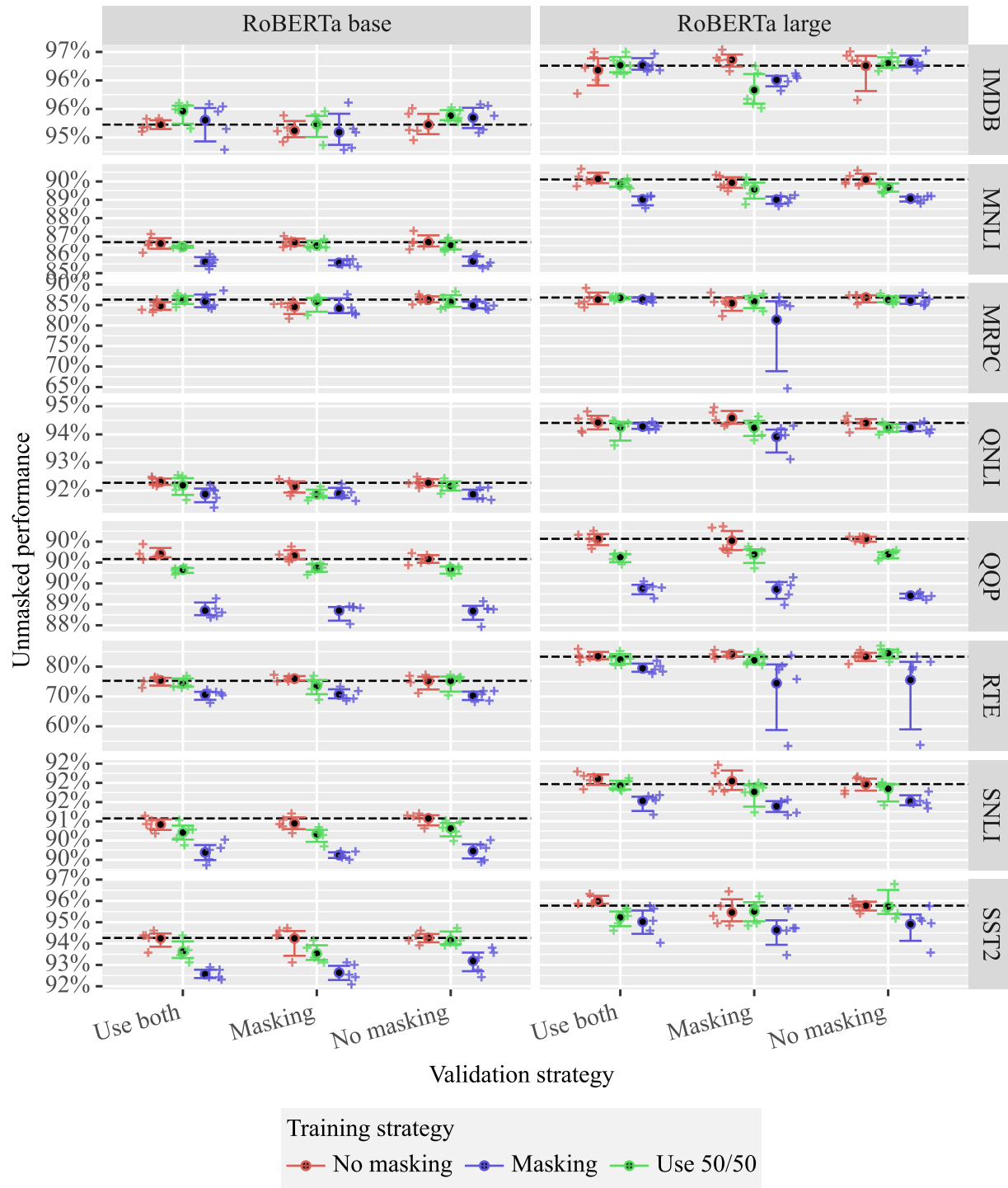


Figure C.3 The unmasked performance for each validation and training strategy, using the test dataset. Not that “*No masking*” as a *training strategy* is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task *training strategy* “*Use 50/50*” works best. This plot is **page-2**. Corresponding main results in Figure 4.3.

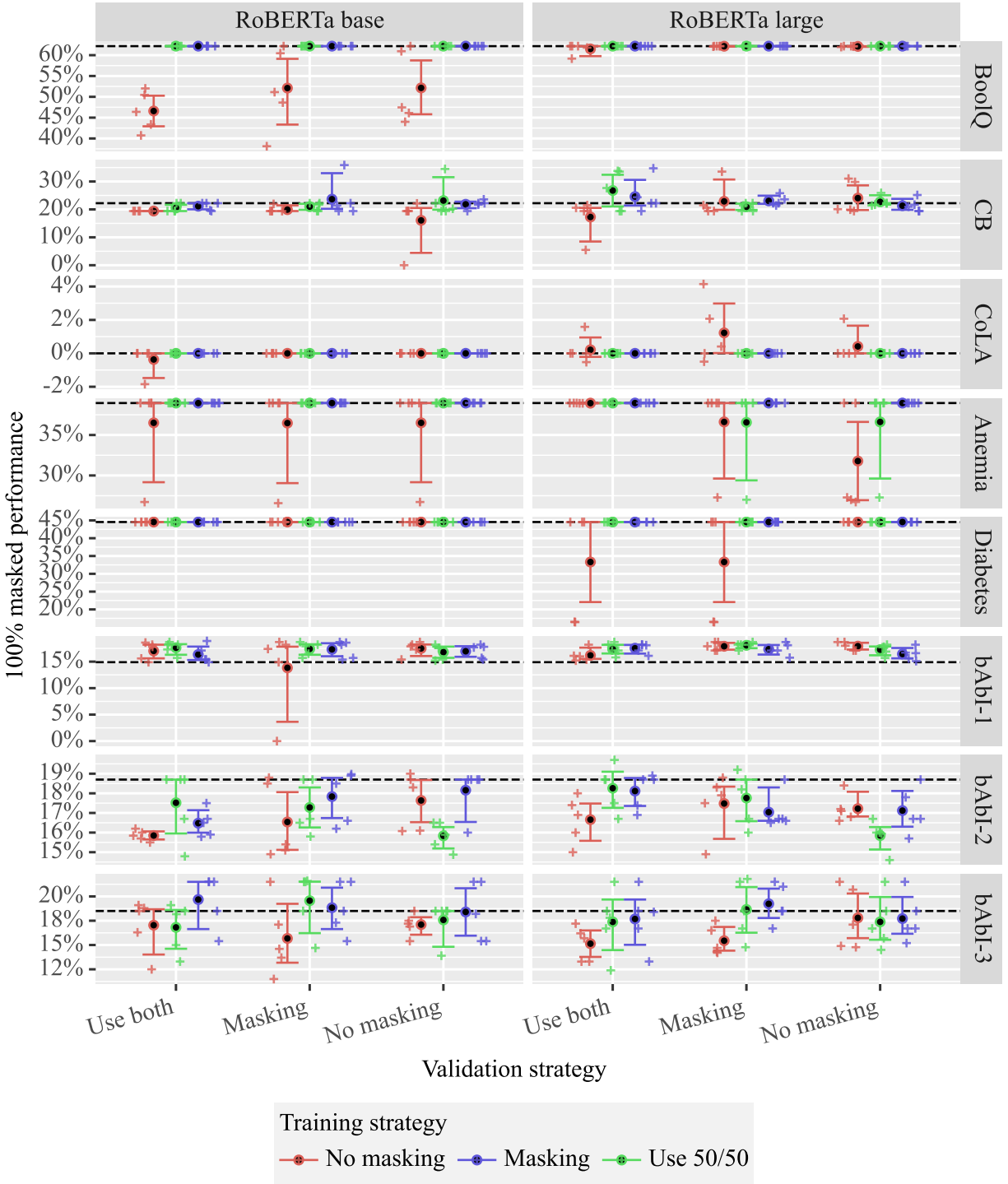


Figure C.4 The 100% masked performance, using the test dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training (“*Masking*” or “*Use 50/50*”) is necessary. This plot is **page-1**. Corresponding main results in Figure 4.4.

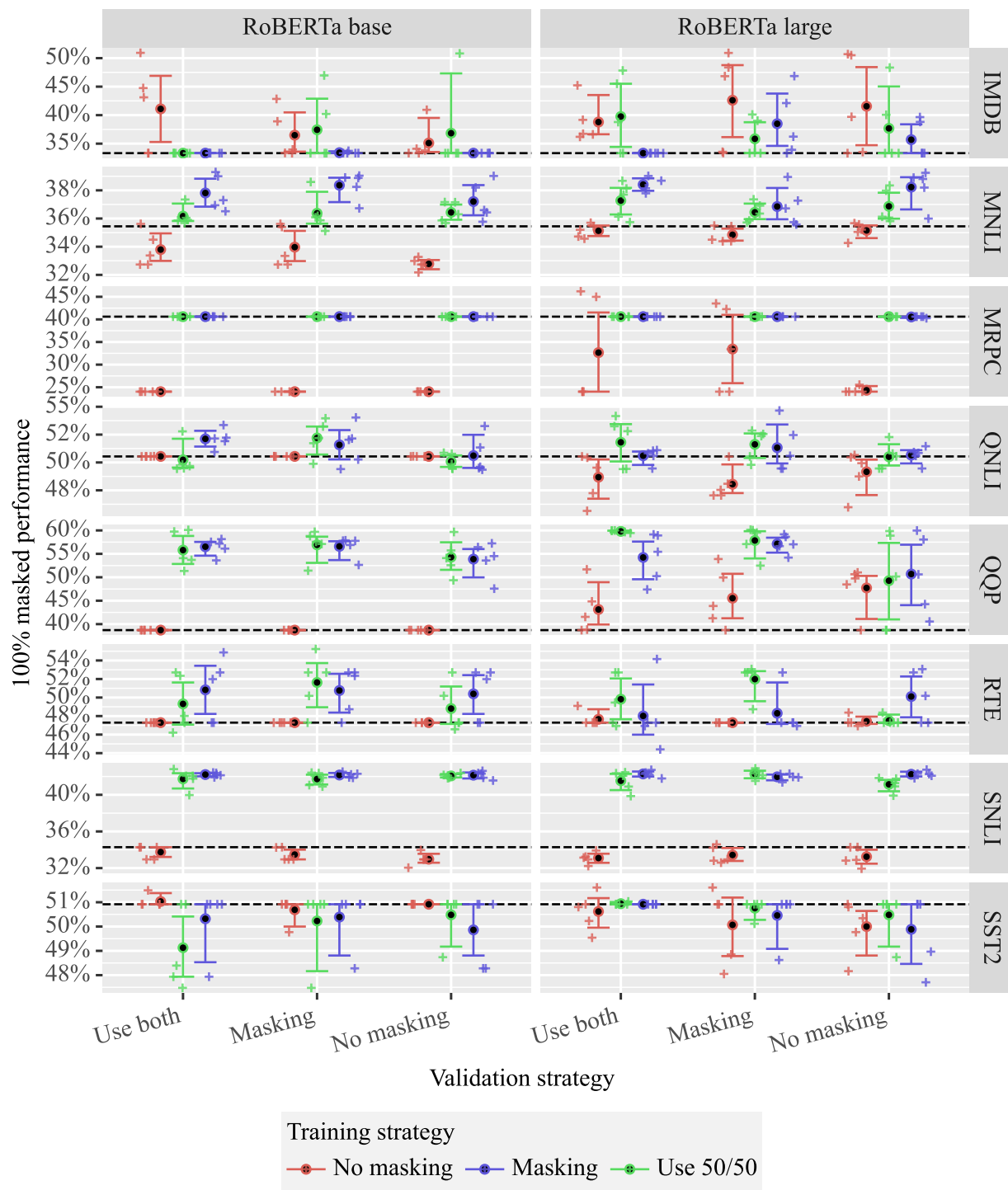


Figure C.5 The 100% masked performance, using the test dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training (“*Masking*” or “*Use 50/50*”) is necessary. This plot is **page-2**. Corresponding main results in Figure 4.4.

### C.2.4 Validation dataset

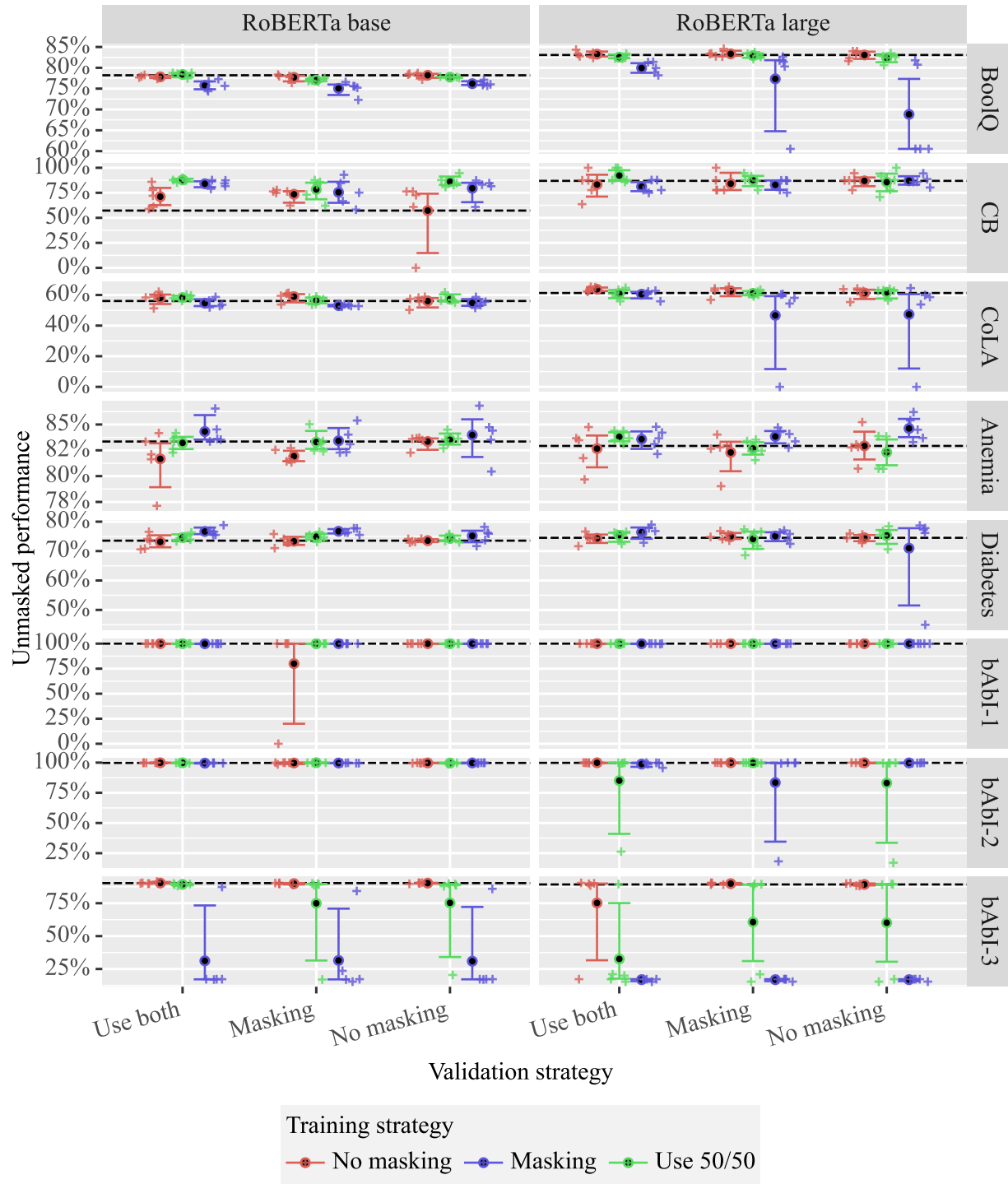


Figure C.6 The unmasked performance for each validation and training strategy, using the validation dataset. Not that “No masking” as a *training strategy* is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task *training strategy* “Use 50/50” works best. This plot is **page-1**.

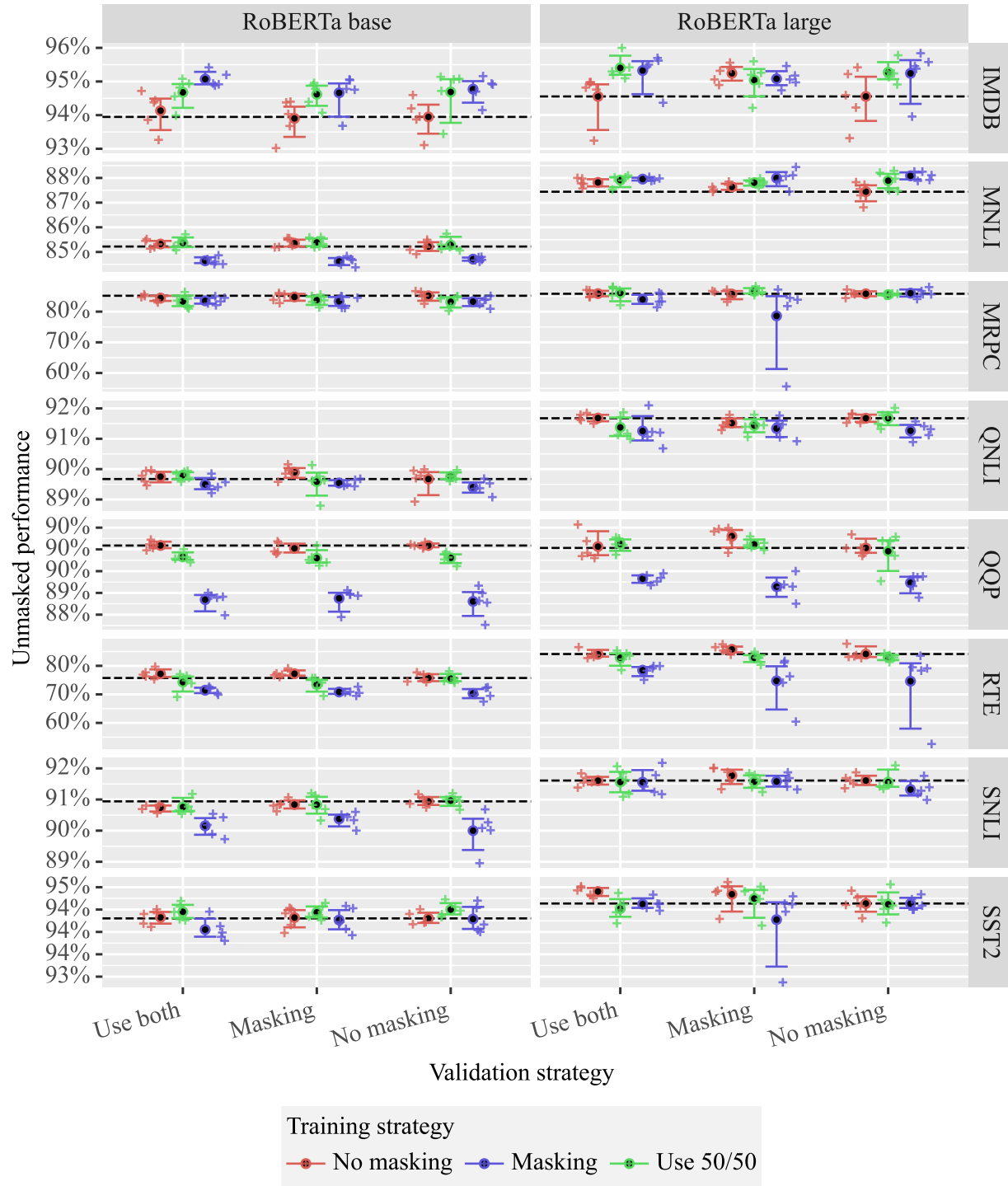


Figure C.7 The unmasked performance for each validation and training strategy, using the validation dataset. Not that “No masking” as a *training strategy* is not a valid option only a baseline, as it creates OOD issues. We find that the multi-task *training strategy* “Use 50/50” works best. This plot is **page-2**.

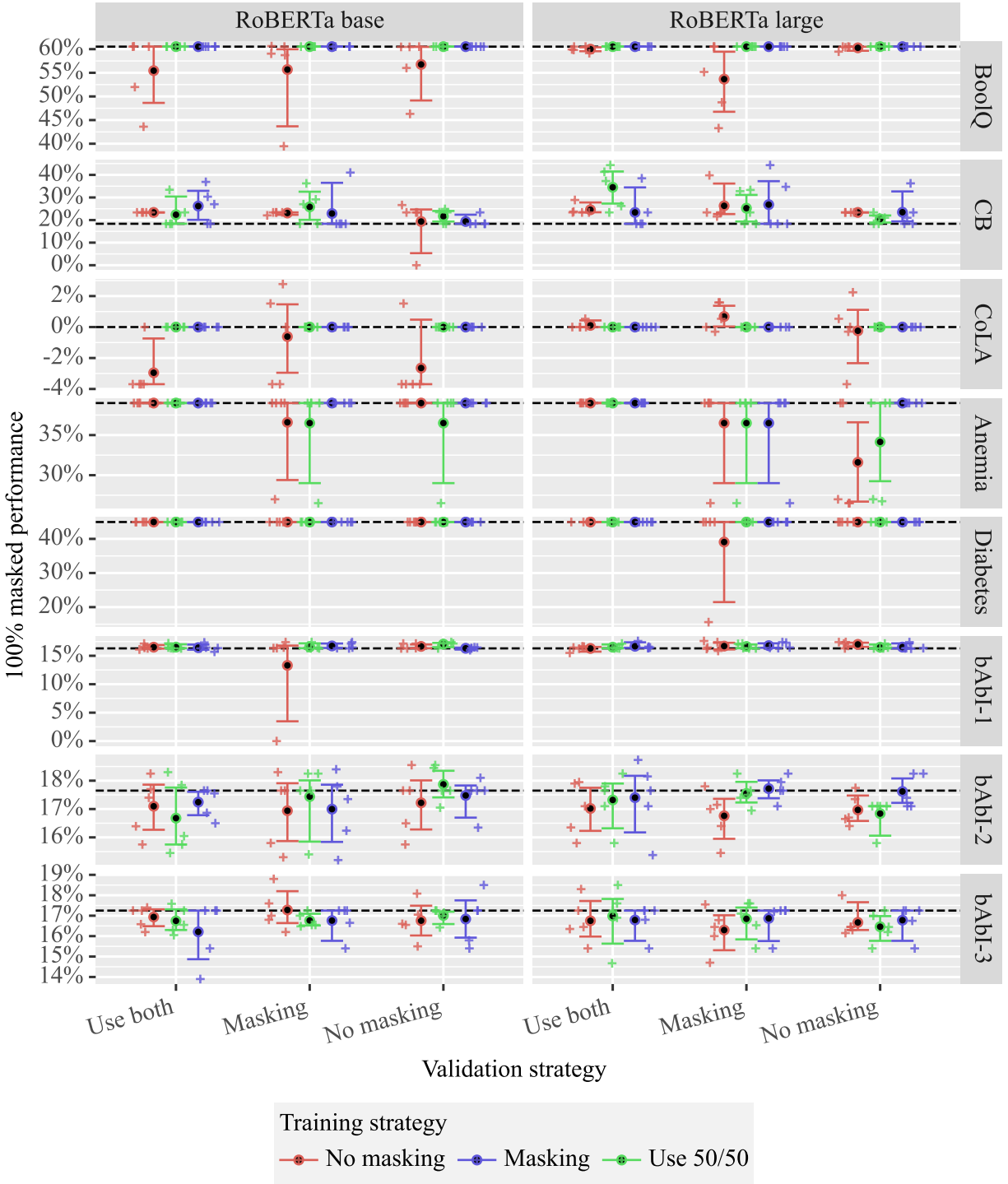


Figure C.8 The 100% masked performance, using the validation dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training (“Masking” or “Use 50/50”) is necessary. This plot is **page-1**.

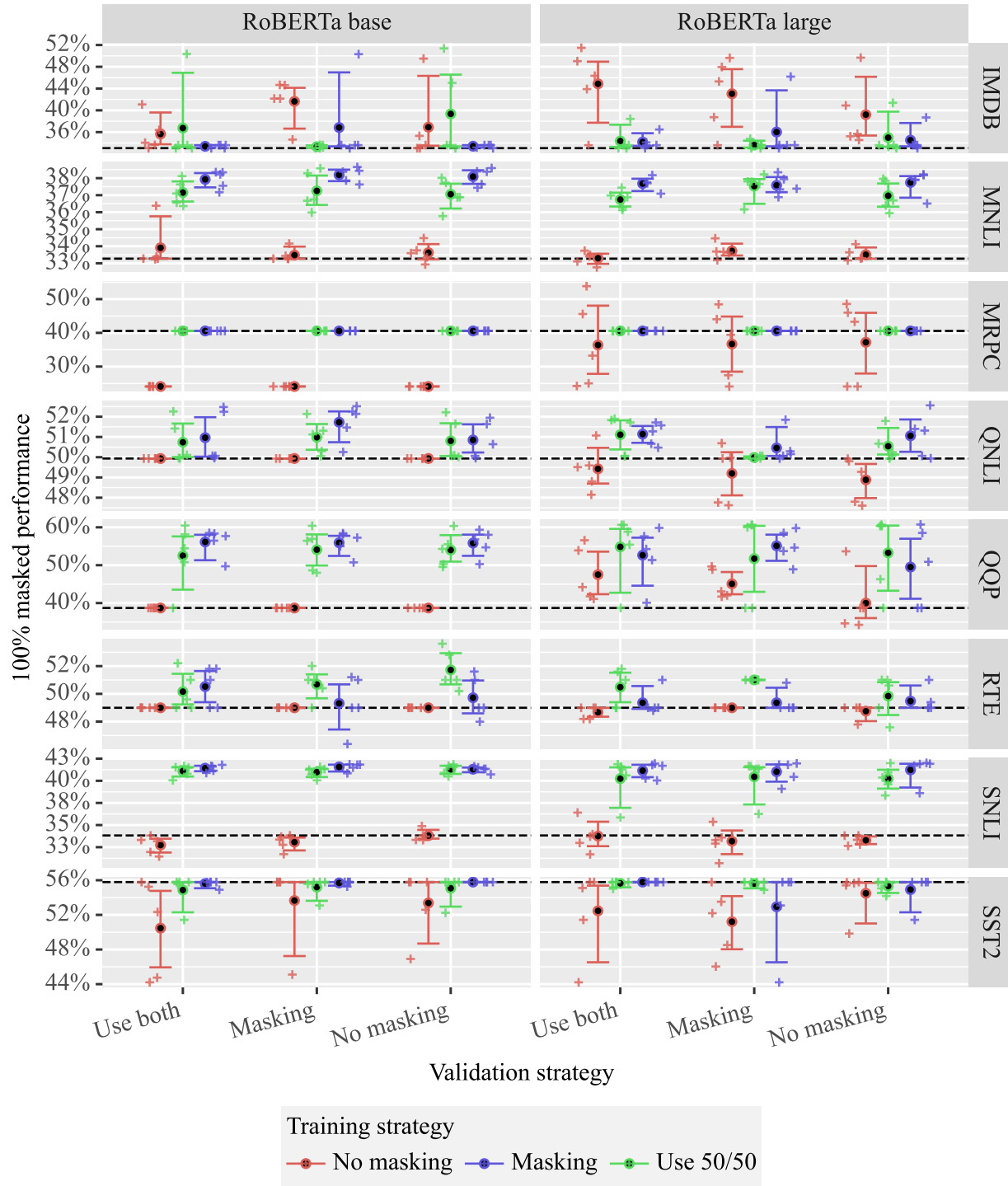


Figure C.9 The 100% masked performance, using the validation dataset. The dashed line represents the class-majority classifier baseline. Results show that masking during training (“Masking” or “Use 50/50”) is necessary. This plot is **page-2**.

### C.3 Convergence speed

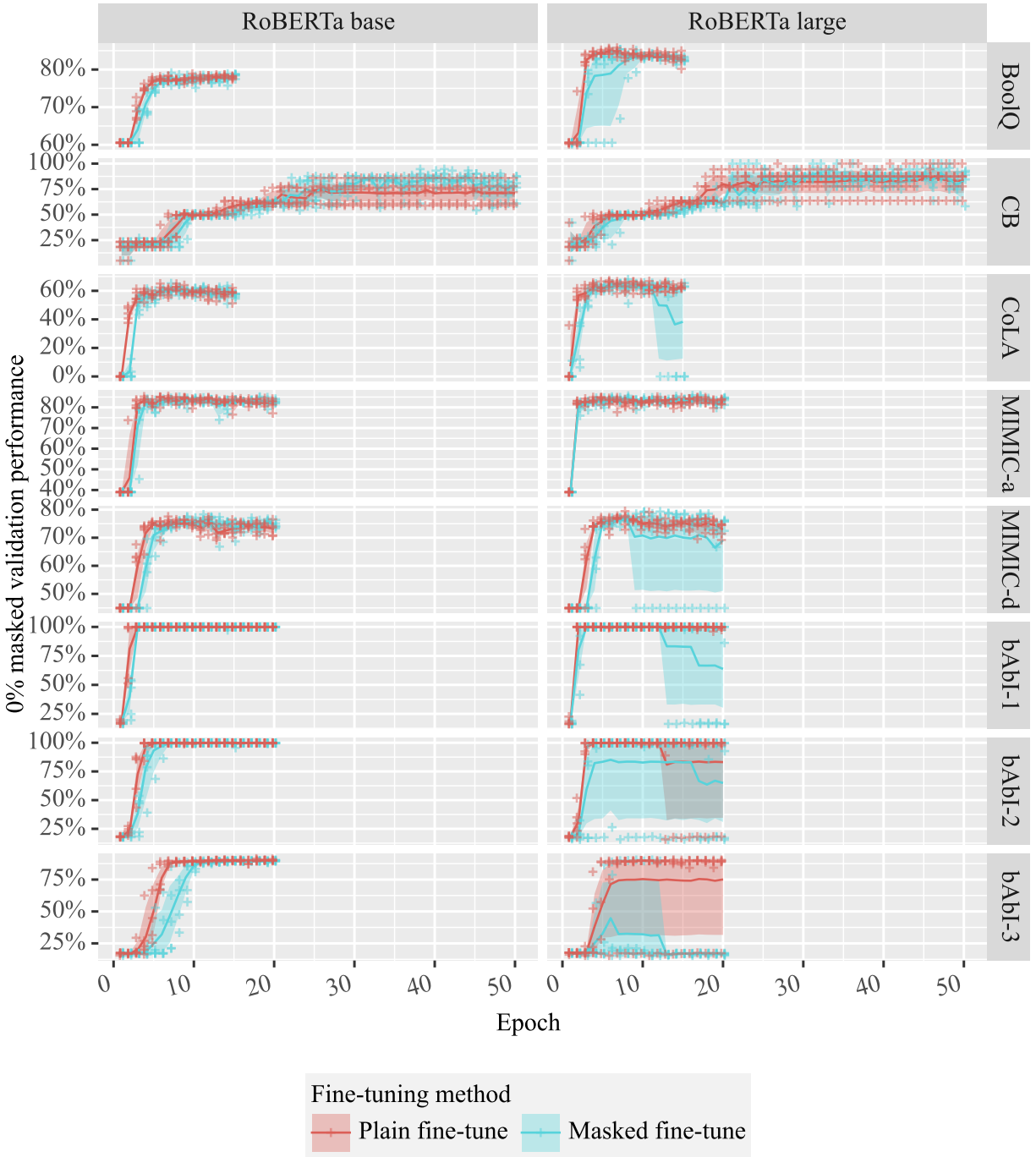


Figure C.10 The validation performance for each epoch. Note that the max number of epochs vary depending on the dataset. This is only to limit the compute requirements when fine-tuning. The best epoch is selected by the “early-stopping” dataset, which has one copy with no masking and one copy with uniformly sampled masking ratios. This plot is **page-1**.

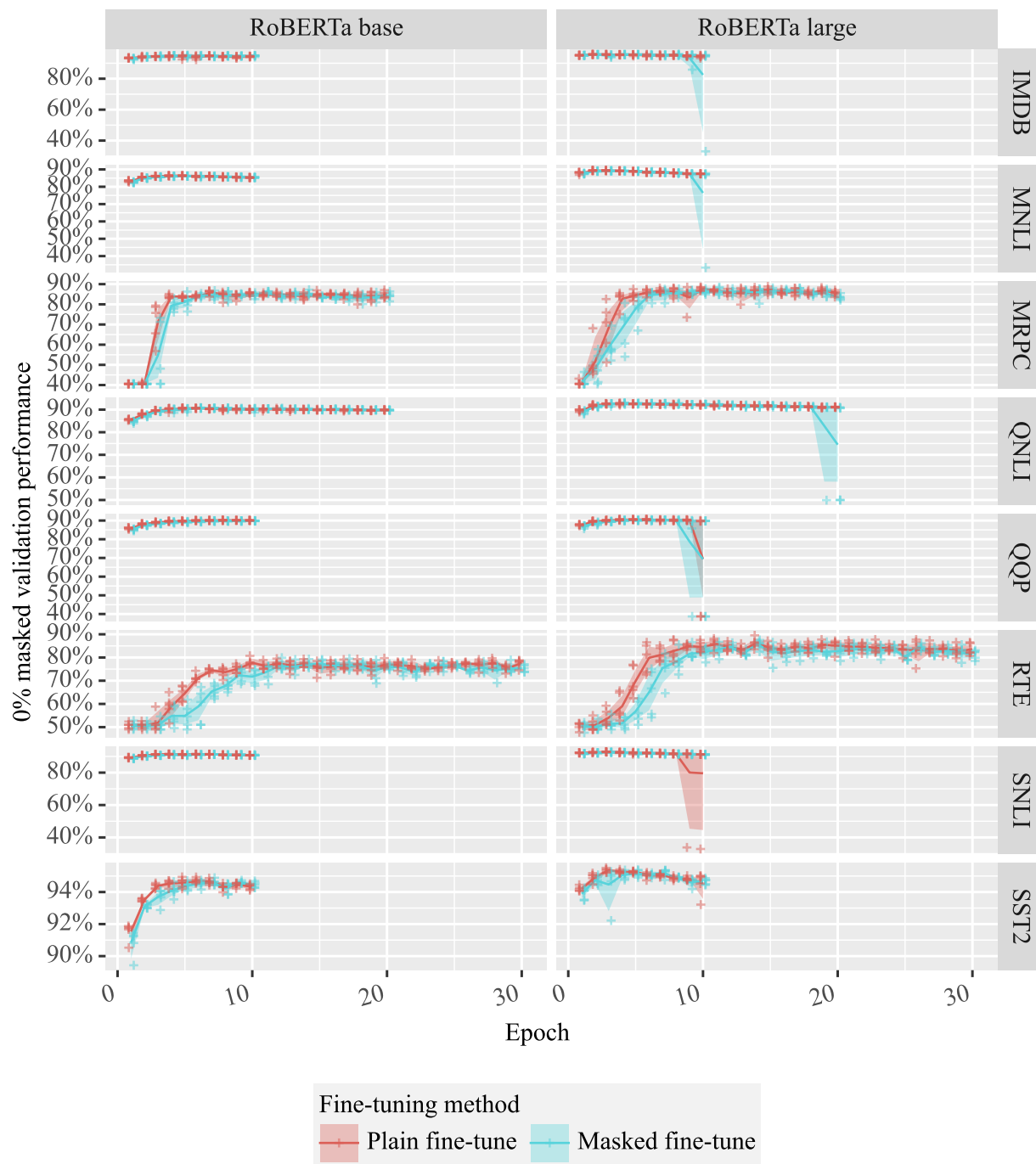


Figure C.11 The validation performance for each epoch. Note that the max number of epochs vary depending on the dataset. This is only to limit the compute requirements when fine-tuning. The best epoch is selected by the “early-stopping” dataset, which has one copy with no masking and one copy with uniformly sampled masking ratios. This plot is **page-2**.

## C.4 In-distribution validation

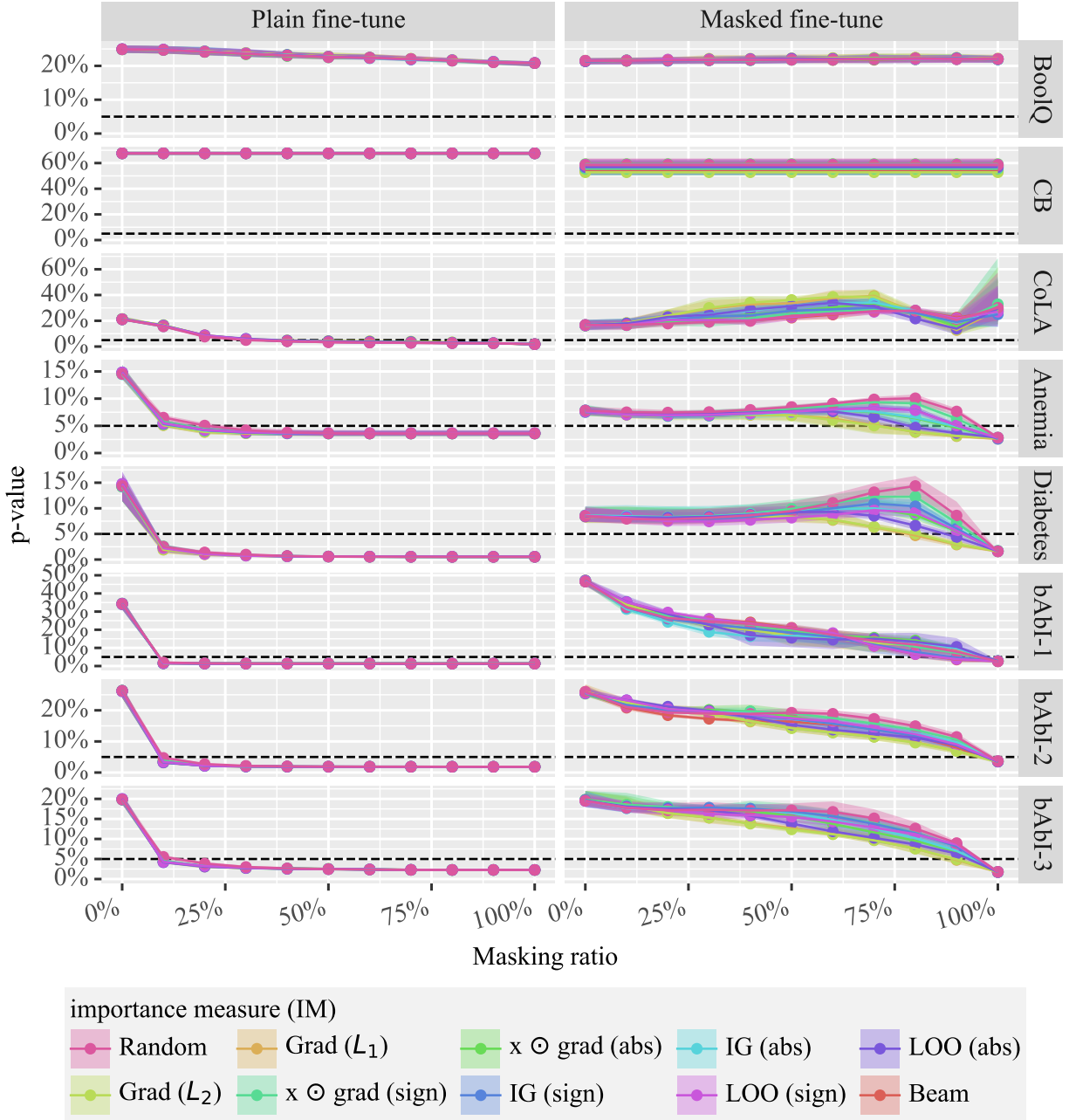


Figure C.12 In-distribution p-values using MaSF, for **RoBERTa-base** with and without masked fine-tuning, **page-1**. The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5.

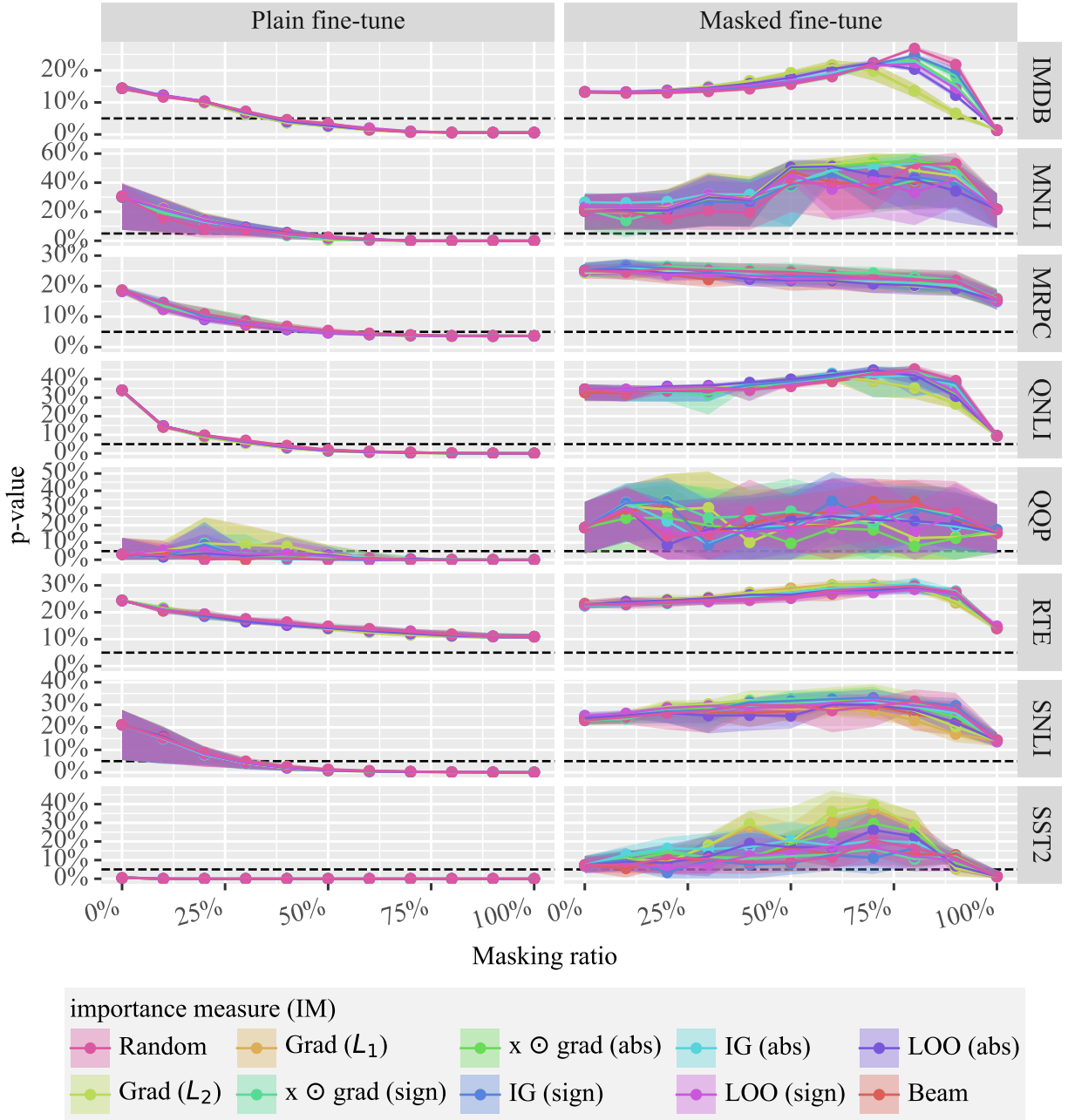


Figure C.13 In-distribution p-values using MaSF, for **RoBERTa-base** with and without masked fine-tuning, **page-2**. The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5.

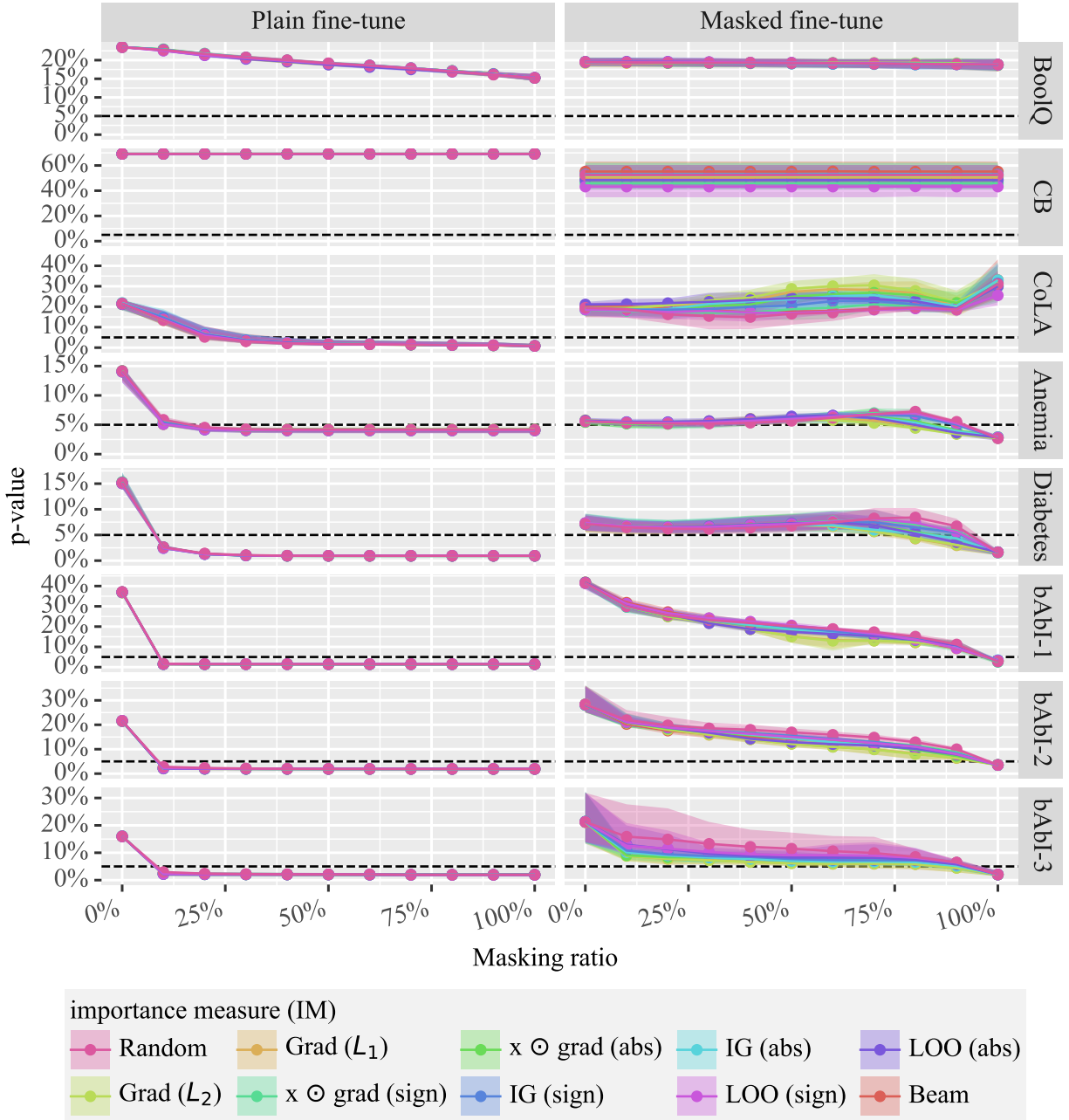


Figure C.14 In-distribution p-values using MaSF, for **RoBERTa-large** with and without masked fine-tuning, **page-1**. The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5.

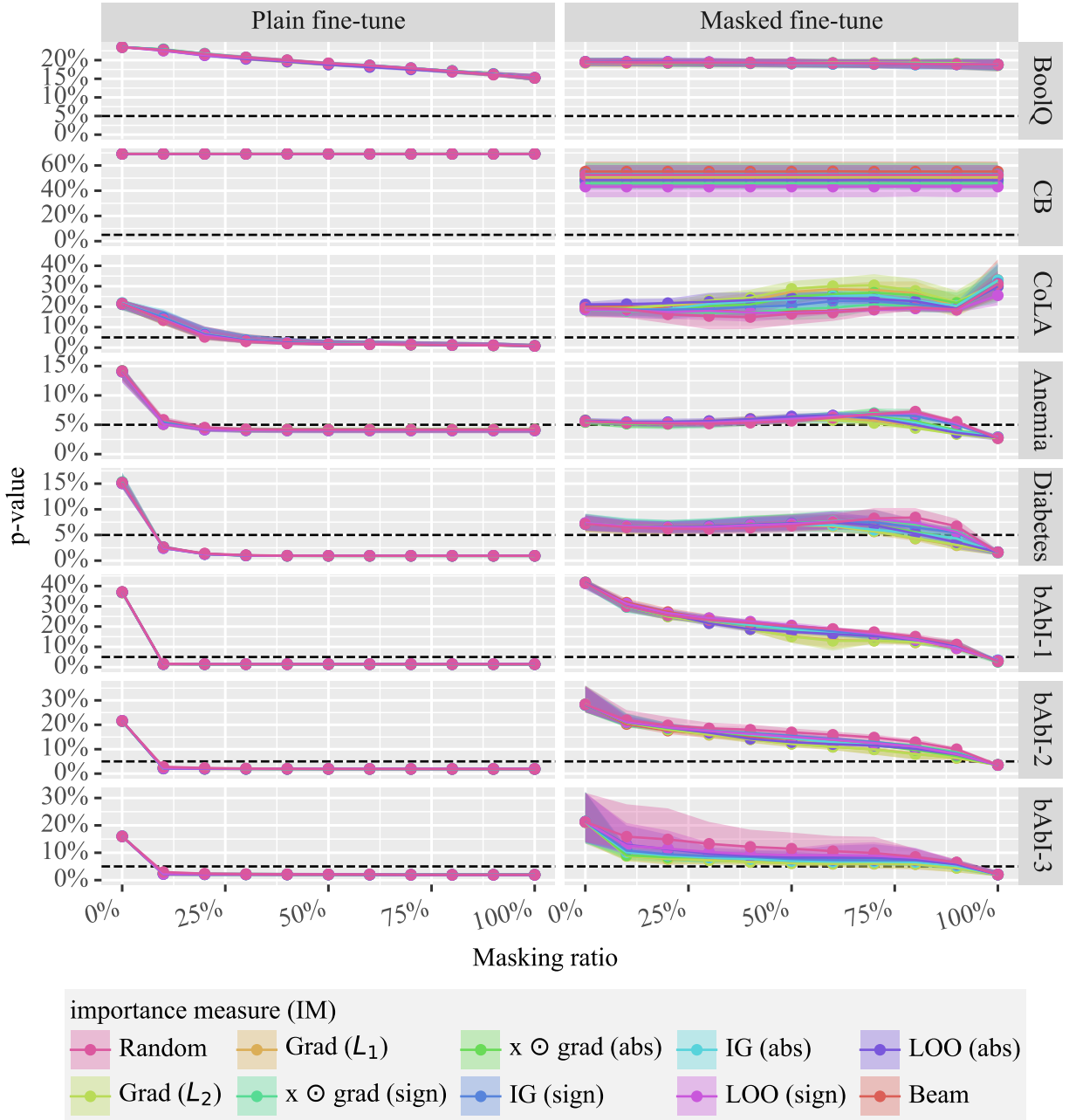


Figure C.15 In-distribution p-values using MaSF, for **RoBERTa-large** with and without masked fine-tuning, **page-2**. The masked tokens are chosen according to an importance measure. P-values below the dashed line show out-of-distribution (OOD) results, given a 5% risk of a false positive. Results show that only when using masked fine-tuning, masked data is consistently not OOD. Corresponding main results are in Figure 4.5.

## C.5 Faithfulness metrics

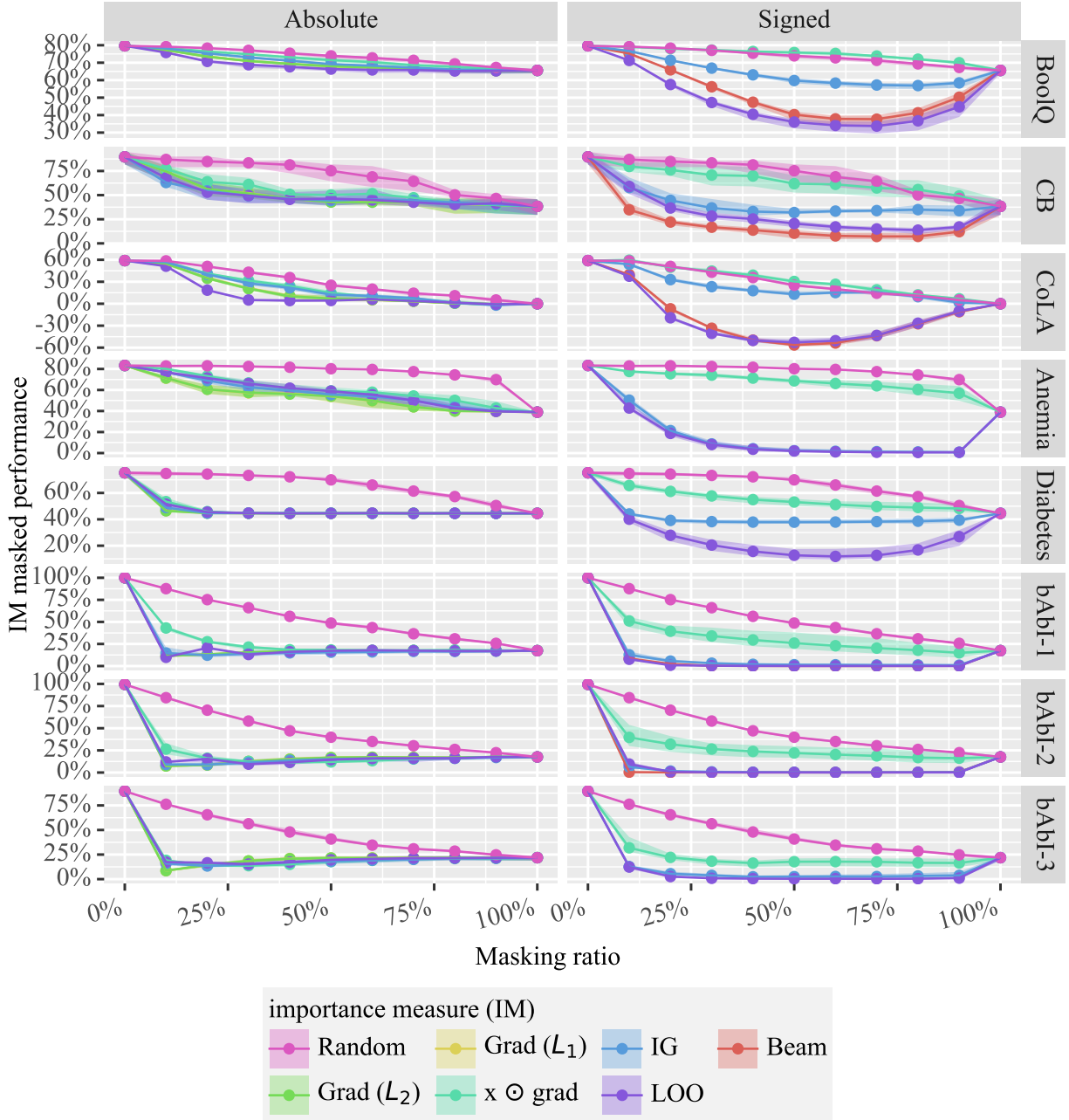


Figure C.16 The performance given the masked datasets, where masking is done for the x% allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “Random” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-1** for **RoBERTa-base**. Corresponding main results in Section 4.2.3.

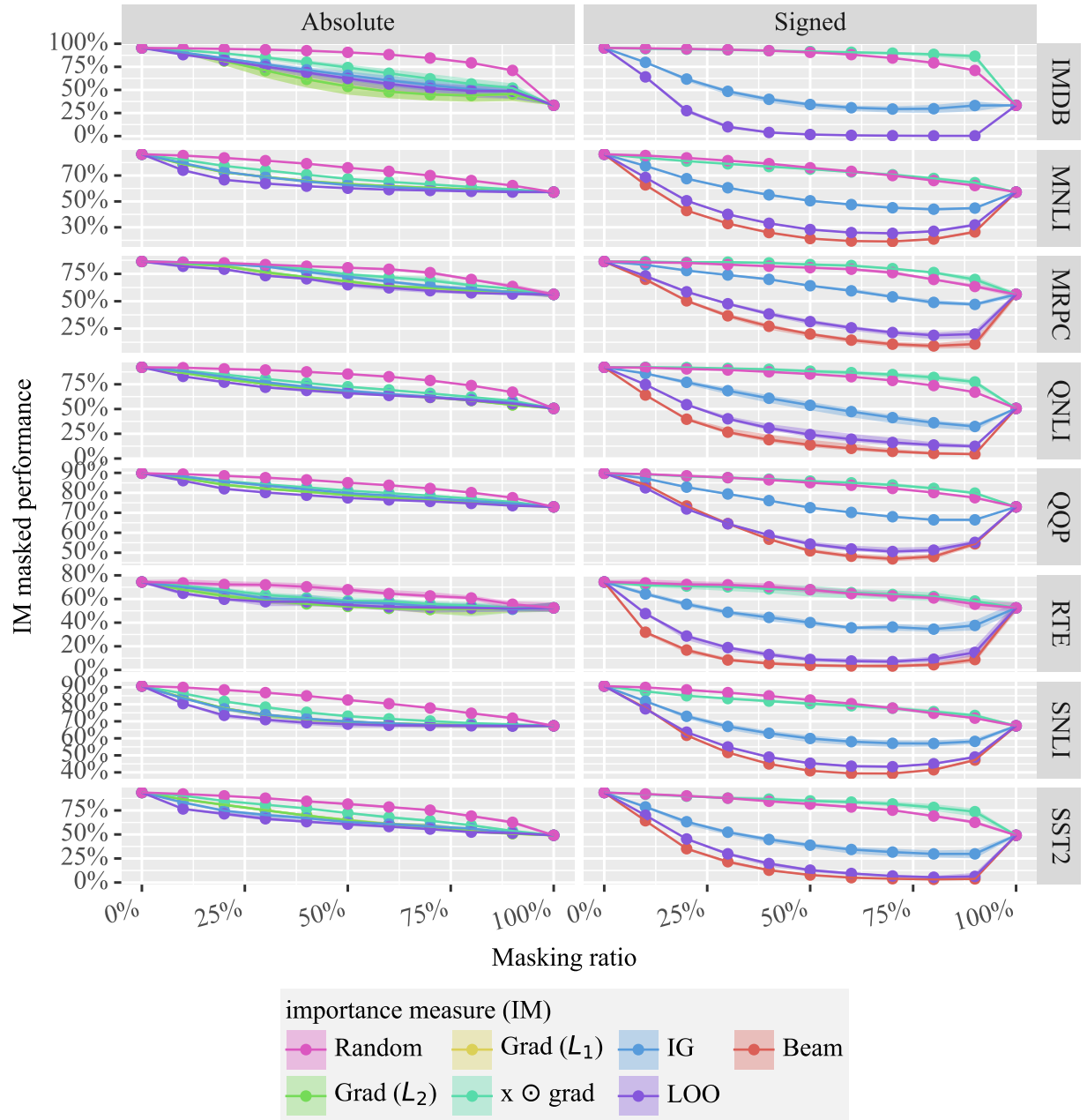


Figure C.17 The performance given the masked datasets, where masking is done for the  $x\%$  allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “Random” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-2** for **RoBERTa-base**. Corresponding main results in Section 4.2.3.

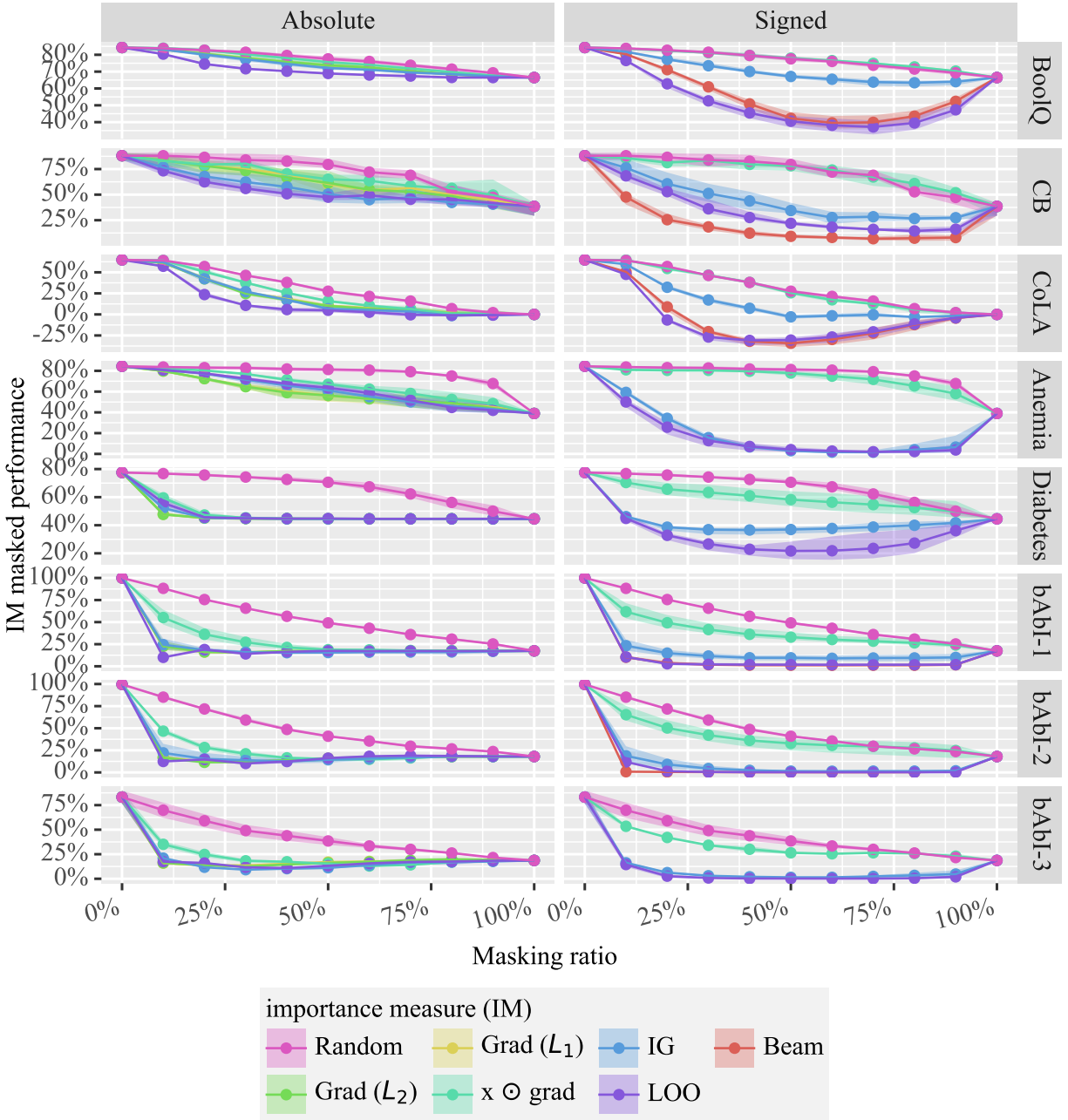


Figure C.18 The performance given the masked datasets, where masking is done for the  $x\%$  allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “Random” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-1** for **RoBERTa-large**. Corresponding main results in Section 4.2.3.

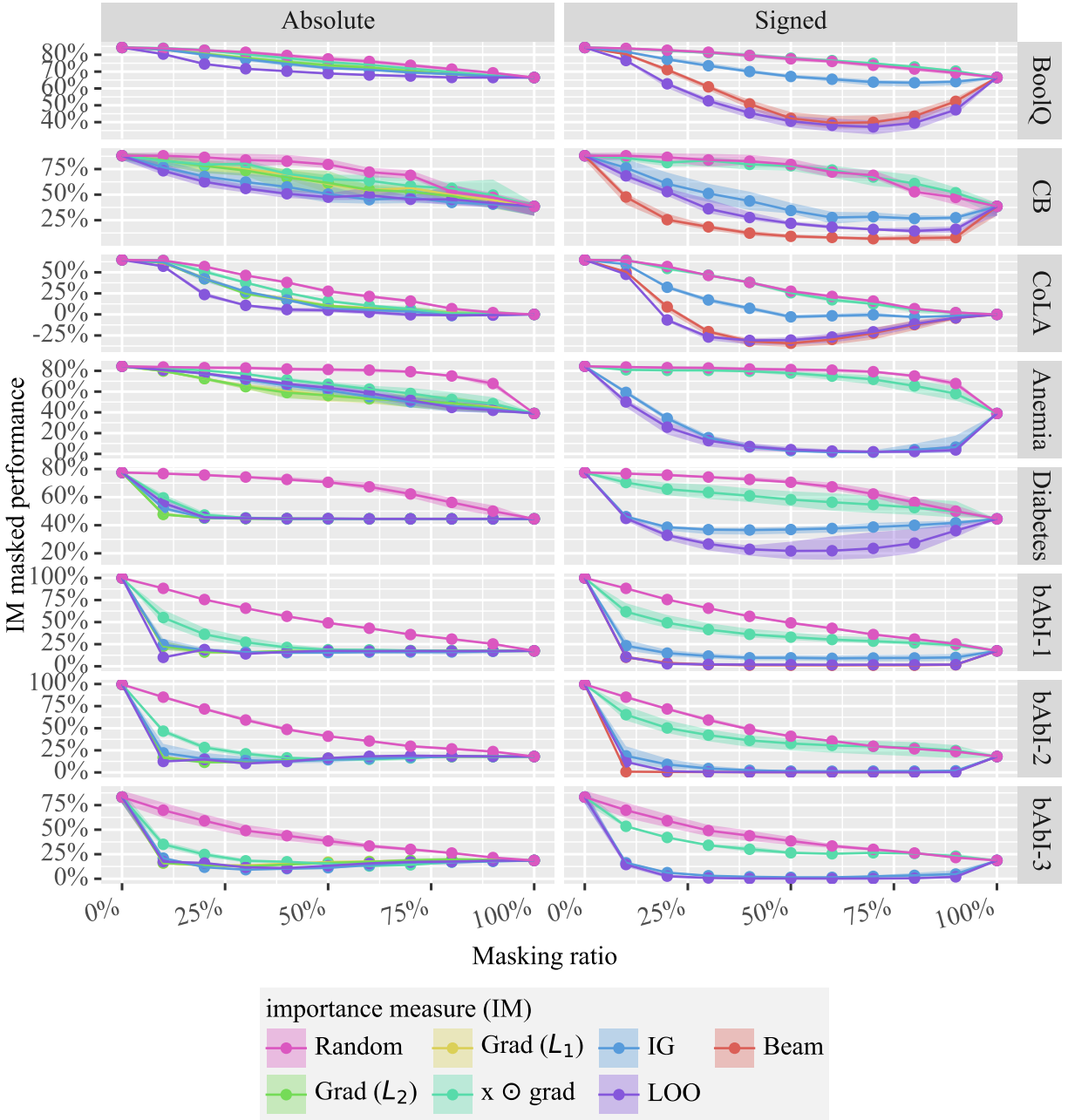


Figure C.19 The performance given the masked datasets, where masking is done for the  $x\%$  allegedly most important tokens according to the importance measure. If the performance for a given explanation is below the “*Random*” baseline, this shows faithfulness. Although, faithfulness is not an absolute concept, so more is better. This plot is **page-2** for **RoBERTa-large**. Corresponding main results in Section 4.2.3.

## C.5.1 Relative Area between Curves (RACU)

Table C.7 Faithfulness scores for **RoBERTa-base**. Shows Relative Area Between Curves (RACU) and the non-relative variant (ACU), defined in Chapter 3. Also compares with Recursive-ROAR from Chapter 3.

Faithfulness [%]					Faithfulness [%]					Faithfulness [%]				
Dataset	IM	Our		R-ROAR	Dataset	IM	Our		R-ROAR	Dataset	IM	Our		R-ROAR
		ACU	RACU	RACU			ACU	RACU	RACU			ACU	RACU	RACU
bAbI-1	Grad ( $L_2$ )	32.6 $^{+1.2}_{-1.1}$	91.9 $^{+2.0}_{-3.2}$	64.2 $^{+2.6}_{-2.6}$	CoLA	Grad ( $L_2$ )	12.7 $^{+0.9}_{-1.1}$	43.4 $^{+4.0}_{-3.2}$	—	SST2	Grad ( $L_2$ )	12.2 $^{+0.6}_{-0.7}$	40.4 $^{+3.0}_{-1.7}$	26.1 $^{+1.6}_{-2.2}$
	Grad ( $L_1$ )	32.7 $^{+1.3}_{-1.0}$	92.1 $^{+2.0}_{-3.1}$	—		Grad ( $L_1$ )	12.7 $^{+0.8}_{-1.2}$	43.2 $^{+3.7}_{-3.4}$	—		Grad ( $L_1$ )	12.1 $^{+0.7}_{-0.7}$	40.3 $^{+3.3}_{-1.8}$	—
	x $\odot$ grad (sign)	21.4 $^{+4.9}_{-7.3}$	60.5 $^{+13.7}_{-20.7}$	—		x $\odot$ grad (sign)	-2.4 $^{+1.9}_{-1.1}$	-8.3 $^{+6.0}_{-3.7}$	—		x $\odot$ grad (sign)	-3.7 $^{+1.5}_{-1.6}$	-12.2 $^{+4.5}_{-6.0}$	—
	x $\odot$ grad (abs)	27.0 $^{+1.1}_{-1.4}$	76.0 $^{+3.0}_{-1.6}$	52.1 $^{+1.8}_{-3.7}$		x $\odot$ grad (abs)	8.0 $^{+1.0}_{-0.8}$	27.1 $^{+2.9}_{-2.9}$	—		x $\odot$ grad (abs)	7.1 $^{+0.2}_{-0.2}$	23.5 $^{+1.9}_{-1.1}$	18.6 $^{+4.1}_{-4.6}$
	IG (sign)	44.1 $^{+0.8}_{-0.8}$	124.2 $^{+2.3}_{-4.5}$	—		IG (sign)	8.1 $^{+1.1}_{-0.8}$	27.5 $^{+2.8}_{-2.6}$	—		IG (sign)	31.8 $^{+2.8}_{-2.2}$	105.6 $^{+2.7}_{-2.7}$	—
	IG (abs)	33.3 $^{+2.4}_{-2.4}$	93.7 $^{+5.5}_{-4.1}$	48.2 $^{+4.1}_{-3.7}$		IG (abs)	8.7 $^{+1.7}_{-0.8}$	29.8 $^{+5.6}_{-4.0}$	—		IG (abs)	13.7 $^{+0.8}_{-0.8}$	45.3 $^{+4.1}_{-2.8}$	32.9 $^{+1.8}_{-1.5}$
	LOO (sign)	46.1 $^{+0.7}_{-0.7}$	129.9 $^{+3.4}_{-3.5}$	—		LOO (sign)	52.0 $^{+1.7}_{-2.9}$	177.2 $^{+7.2}_{-7.8}$	—		LOO (sign)	51.6 $^{+1.4}_{-0.9}$	171.3 $^{+5.8}_{-6.2}$	—
	LOO (abs)	32.4 $^{+1.2}_{-1.2}$	91.2 $^{+1.2}_{-0.7}$	—		LOO (abs)	17.0 $^{+0.4}_{-0.4}$	57.9 $^{+3.2}_{-2.9}$	—		LOO (abs)	16.6 $^{+1.2}_{-1.0}$	54.9 $^{+2.1}_{-1.5}$	—
Beam	45.9 $^{+0.7}_{-0.7}$	129.2 $^{+3.1}_{-3.2}$	—	Beam	50.6 $^{+1.3}_{-2.3}$	172.7 $^{+5.9}_{-7.2}$	—	Beam	56.4 $^{+0.5}_{-0.7}$	187.3 $^{+8.1}_{-7.1}$	—			
bAbI-2	Grad ( $L_2$ )	28.5 $^{+0.8}_{-0.8}$	96.3 $^{+6.7}_{-2.8}$	57.8 $^{+2.0}_{-2.0}$	Anemia	Grad ( $L_2$ )	23.8 $^{+0.6}_{-0.5}$	62.1 $^{+1.4}_{-1.7}$	18.2 $^{+1.8}_{-13.8}$	SNLI	Grad ( $L_2$ )	8.9 $^{+1.0}_{-0.4}$	62.2 $^{+2.7}_{-1.7}$	50.7 $^{+1.1}_{-0.8}$
	Grad ( $L_1$ )	28.5 $^{+0.9}_{-0.8}$	96.3 $^{+6.8}_{-2.7}$	—		Grad ( $L_1$ )	23.8 $^{+0.6}_{-0.5}$	62.2 $^{+1.4}_{-1.4}$	—		Grad ( $L_1$ )	8.9 $^{+1.0}_{-0.4}$	62.2 $^{+2.6}_{-1.8}$	—
	x $\odot$ grad (sign)	19.7 $^{+6.6}_{-8.1}$	65.7 $^{+24.1}_{-26.3}$	—		x $\odot$ grad (sign)	9.7 $^{+2.6}_{-2.5}$	25.1 $^{+6.5}_{-6.2}$	—		x $\odot$ grad (sign)	1.3 $^{+0.7}_{-0.8}$	9.2 $^{+5.2}_{-4.9}$	—
	x $\odot$ grad (abs)	27.3 $^{+1.7}_{-1.5}$	92.0 $^{+2.5}_{-3.1}$	48.1 $^{+3.2}_{-3.5}$		x $\odot$ grad (abs)	16.6 $^{+1.3}_{-1.3}$	43.2 $^{+3.0}_{-3.7}$	8.8 $^{+22.7}_{-22.8}$		x $\odot$ grad (abs)	6.4 $^{+0.5}_{-0.3}$	44.8 $^{+1.6}_{-1.6}$	41.0 $^{+0.4}_{-0.5}$
	IG (sign)	40.3 $^{+0.9}_{-0.8}$	136.3 $^{+4.4}_{-4.4}$	—		IG (sign)	62.0 $^{+1.6}_{-1.5}$	161.8 $^{+2.7}_{-2.3}$	—		IG (sign)	16.3 $^{+1.3}_{-1.3}$	113.5 $^{+2.6}_{-2.6}$	—
	IG (abs)	29.1 $^{+1.0}_{-1.0}$	98.3 $^{+3.5}_{-3.5}$	42.0 $^{+3.8}_{-1.8}$		IG (abs)	20.0 $^{+0.9}_{-1.6}$	52.1 $^{+2.5}_{-3.5}$	12.5 $^{+11.3}_{-7.0}$		IG (abs)	8.9 $^{+0.8}_{-0.5}$	62.3 $^{+2.5}_{-2.1}$	56.7 $^{+1.0}_{-1.1}$
	LOO (sign)	40.2 $^{+1.2}_{-0.8}$	136.0 $^{+4.1}_{-6.5}$	—		LOO (sign)	63.3 $^{+1.4}_{-1.6}$	165.2 $^{+2.9}_{-3.2}$	—		LOO (sign)	26.6 $^{+0.4}_{-0.4}$	186.6 $^{+9.1}_{-12.4}$	—
	LOO (abs)	28.5 $^{+1.3}_{-0.4}$	96.3 $^{+9.2}_{-9.2}$	—		LOO (abs)	18.9 $^{+1.3}_{-1.3}$	49.2 $^{+3.8}_{-3.8}$	—		LOO (abs)	10.5 $^{+0.5}_{-0.5}$	73.6 $^{+2.0}_{-1.9}$	—
Beam	41.1 $^{+1.0}_{-0.7}$	139.2 $^{+5.0}_{-7.3}$	—	Beam	—	—	—	Beam	29.3 $^{+0.2}_{-0.2}$	205.2 $^{+8.1}_{-8.7}$	—			
bAbI-3	Grad ( $L_2$ )	23.5 $^{+0.9}_{-1.5}$	97.3 $^{+3.5}_{-4.2}$	34.0 $^{+14.6}_{-5.1}$	Diabetes	Grad ( $L_2$ )	19.7 $^{+1.2}_{-1.0}$	91.8 $^{+0.6}_{-0.9}$	57.9 $^{+14.4}_{-19.8}$	IMDB	Grad ( $L_2$ )	24.9 $^{+5.9}_{-4.3}$	47.8 $^{+10.9}_{-8.5}$	25.4 $^{+3.1}_{-2.0}$
	Grad ( $L_1$ )	23.5 $^{+0.9}_{-1.3}$	97.2 $^{+3.6}_{-4.3}$	—		Grad ( $L_1$ )	19.6 $^{+1.0}_{-0.7}$	91.6 $^{+0.5}_{-0.9}$	—		Grad ( $L_1$ )	24.9 $^{+5.7}_{-4.4}$	47.8 $^{+10.9}_{-8.3}$	—
	x $\odot$ grad (sign)	23.1 $^{+2.4}_{-2.5}$	96.7 $^{+21.8}_{-14.9}$	—		x $\odot$ grad (sign)	10.9 $^{+1.3}_{-1.3}$	51.1 $^{+9.1}_{-7.8}$	—		x $\odot$ grad (sign)	-3.2 $^{+1.3}_{-1.0}$	-6.2 $^{+2.5}_{-1.9}$	—
	x $\odot$ grad (abs)	24.4 $^{+0.7}_{-1.2}$	101.2 $^{+4.3}_{-5.4}$	22.4 $^{+15.9}_{-12.4}$		x $\odot$ grad (abs)	18.8 $^{+1.3}_{-0.7}$	87.9 $^{+1.4}_{-2.0}$	53.4 $^{+23.2}_{-29.3}$		x $\odot$ grad (abs)	12.9 $^{+4.5}_{-2.6}$	24.7 $^{+8.7}_{-5.1}$	16.9 $^{+1.1}_{-0.5}$
	IG (sign)	36.6 $^{+1.3}_{-1.4}$	152.5 $^{+14.7}_{-17.4}$	—		IG (sign)	24.8 $^{+1.5}_{-2.1}$	115.8 $^{+2.8}_{-7.8}$	—		IG (sign)	40.3 $^{+3.3}_{-1.9}$	77.3 $^{+6.3}_{-3.6}$	—
	IG (abs)	24.4 $^{+1.0}_{-1.3}$	100.9 $^{+4.0}_{-3.9}$	-27.9 $^{+18.0}_{-49.1}$		IG (abs)	19.4 $^{+1.0}_{-0.6}$	90.5 $^{+0.6}_{-1.2}$	26.1 $^{+12.0}_{-25.1}$		IG (abs)	18.4 $^{+4.2}_{-2.8}$	35.4 $^{+8.0}_{-5.9}$	35.1 $^{+2.4}_{-1.7}$
	LOO (sign)	38.7 $^{+0.9}_{-0.7}$	160.5 $^{+13.8}_{-13.2}$	—		LOO (sign)	41.5 $^{+2.3}_{-5.9}$	193.4 $^{+3.9}_{-17.4}$	—		LOO (sign)	68.0 $^{+0.9}_{-1.3}$	130.7 $^{+1.3}_{-1.3}$	—
	LOO (abs)	23.4 $^{+0.7}_{-1.1}$	97.0 $^{+4.5}_{-4.6}$	—		LOO (abs)	19.1 $^{+1.1}_{-0.6}$	89.0 $^{+0.4}_{-0.6}$	—		LOO (abs)	20.7 $^{+4.2}_{-3.2}$	39.8 $^{+9.5}_{-6.0}$	—
Beam	—	—	—	Beam	—	—	—	Beam	—	—	—			
BoolQ	Grad ( $L_2$ )	4.1 $^{+0.2}_{-0.3}$	50.1 $^{+3.6}_{-3.2}$	—	MRPC	Grad ( $L_2$ )	7.8 $^{+0.5}_{-0.6}$	36.8 $^{+4.3}_{-6.1}$	—	MNLI	Grad ( $L_2$ )	8.8 $^{+0.3}_{-0.2}$	49.6 $^{+1.1}_{-1.2}$	—
	Grad ( $L_1$ )	4.1 $^{+0.2}_{-0.3}$	50.1 $^{+4.0}_{-3.1}$	—		Grad ( $L_1$ )	7.9 $^{+0.5}_{-0.7}$	37.4 $^{+4.2}_{-6.3}$	—		Grad ( $L_1$ )	8.8 $^{+0.3}_{-0.2}$	49.7 $^{+1.1}_{-1.2}$	—
	x $\odot$ grad (sign)	-1.3 $^{+0.7}_{-0.4}$	-16.0 $^{+8.2}_{-4.5}$	—		x $\odot$ grad (sign)	-2.9 $^{+0.5}_{-0.6}$	-13.7 $^{+2.9}_{-2.9}$	—		x $\odot$ grad (sign)	0.6 $^{+0.6}_{-0.7}$	3.5 $^{+3.7}_{-3.8}$	—
	x $\odot$ grad (abs)	1.7 $^{+0.2}_{-0.2}$	21.3 $^{+2.7}_{-3.5}$	—		x $\odot$ grad (abs)	3.3 $^{+0.6}_{-0.5}$	15.6 $^{+3.4}_{-4.6}$	—		x $\odot$ grad (abs)	5.7 $^{+0.2}_{-0.2}$	32.0 $^{+0.7}_{-0.8}$	—
	IG (sign)	9.6 $^{+0.6}_{-0.6}$	118.3 $^{+10.3}_{-7.9}$	—		IG (sign)	12.7 $^{+0.8}_{-1.0}$	59.9 $^{+6.2}_{-10.6}$	—		IG (sign)	18.5 $^{+0.6}_{-0.7}$	104.3 $^{+5.9}_{-5.6}$	—
	IG (abs)	3.0 $^{+0.2}_{-0.4}$	37.8 $^{+4.7}_{-5.3}$	—		IG (abs)	5.4 $^{+1.0}_{-0.6}$	25.4 $^{+2.0}_{-3.9}$	—		IG (abs)	9.1 $^{+0.2}_{-0.2}$	51.3 $^{+2.0}_{-2.0}$	—
	LOO (sign)	26.3 $^{+1.9}_{-1.4}$	323.1 $^{+14.4}_{-15.2}$	—		LOO (sign)	37.1 $^{+2.4}_{-1.3}$	174.4 $^{+12.4}_{-18.9}$	—		LOO (sign)	34.7 $^{+0.7}_{-0.5}$	195.7 $^{+3.2}_{-2.8}$	—
	LOO (abs)	5.3 $^{+0.3}_{-0.2}$	65.2 $^{+3.3}_{-2.2}$	—		LOO (abs)	9.9 $^{+0.8}_{-0.9}$	47.0 $^{+5.5}_{-3.9}$	—		LOO (abs)	11.8 $^{+0.3}_{-0.3}$	66.7 $^{+0.9}_{-1.1}$	—
Beam	21.2 $^{+1.1}_{-1.1}$	261.4 $^{+10.7}_{-14.5}$	—	Beam	45.6 $^{+2.6}_{-1.5}$	214.8 $^{+13.3}_{-25.1}$	—	Beam	40.5 $^{+0.8}_{-0.4}$	228.9 $^{+4.5}_{-1.5}$	—			
CB	Grad ( $L_2$ )	20.3 $^{+1.3}_{-1.3}$	64.7 $^{+9.9}_{-9.9}$	—	RTE	Grad ( $L_2$ )	9.4 $^{+0.8}_{-1.0}$	73.9 $^{+23.6}_{-19.3}$	—	QNLI	Grad ( $L_2$ )	12.7 $^{+0.6}_{-0.7}$	40.8 $^{+2.4}_{-2.0}$	—
	Grad ( $L_1$ )	20.2 $^{+1.4}_{-1.3}$	64.7 $^{+10.7}_{-10.7}$	—		Grad ( $L_1$ )	9.4 $^{+0.8}_{-1.0}$	73.5 $^{+23.3}_{-19.1}$	—		Grad ( $L_1$ )	12.7 $^{+0.6}_{-0.7}$	40.8 $^{+2.3}_{-1.9}$	—
	x $\odot$ grad (sign)	6.0 $^{+4.0}_{-3.6}$	17.8 $^{+7.0}_{-11.0}$	—		x $\odot$ grad (sign)	0.1 $^{+2.3}_{-1.7}$	-3.5 $^{+17.3}_{-13.9}$	—		x $\odot$ grad (sign)	-3.7 $^{+1.9}_{-1.1}$	-11.9 $^{+6.0}_{-3.4}$	—
	x $\odot$ grad (abs)	15.9 $^{+2.2}_{-3.2}$	52.0 $^{+15.7}_{-14.3}$	—		x $\odot$ grad (abs)	5.7 $^{+0.5}_{-0.5}$	44.3 $^{+11.8}_{-8.5}$	—		x $\odot$ grad (abs)	8.8 $^{+0.6}_{-0.6}$	28.3 $^{+2.1}_{-1.7}$	—
	IG (sign)	30.1 $^{+5.7}_{-5.7}$	97.6 $^{+22.7}_{-28.3}$	—		IG (sign)	20.2 $^{+1.3}_{-2.2}$	156.6 $^{+45.5}_{-28.4}$	—		IG (sign)	24.4 $^{+4.3}_{-2.0}$	78.0 $^{+10.4}_{-7.4}$	—
	IG (abs)	21.1 $^{+2.3}_{-2.3}$	67.5 $^{+10.8}_{-13.2}$	—		IG (abs)	7.3 $^{+1.6}_{-0.9}$	55.9 $^{+14.8}_{-13.7}$	—		IG (abs)	11.4 $^{+0.7}_{-0.6}$	36.5 $^{+2.0}_{-2.6}$	—
	LOO (sign)	40.9 $^{+5.8}_{-5.5}$	131.4 $^{+26.7}_{-27.3}$	—		LOO (sign)	44.4 $^{+1.6}_{-1.9}$	344.8 $^{+92.4}_{-75.3}$	—		LOO (sign)	46.1 $^{+1.8}_{-2.6}$	147.7 $^{+4.7}_{-8.7}$	—
	LOO (abs)	21.0 $^{+1.7}_{-2.4}$	67.4 $^{+12.8}_{-12.8}$	—		LOO (abs)	9.4 $^{+0.9}_{-0.6}$	73.2 $^{+21.1}_{-15.2}$	—		LOO (abs)	14.0 $^{+0.6}_{-0.6}$	44.9 $^{+3.2}_{-1.3}$	—
Beam	51.0 $^{+6.4}_{-6.8}$	163.5 $^{+32.7}_{-33.8}$	—	Beam	51.3 $^{+2.0}_{-1.7}$	399.7 $^{+109.9}_{-89.9}$	—	Beam	55.7 $^{+1.5}_{-2.2}$	178.4 $^{+4.8}_{-8.6}$	—			
QQP	Grad ( $L_2$ )	4.5 $^{+0.1}_{-0.1}$	40.0 $^{+0.4}_{-0.8}$	—	QQP	Grad ( $L_2$ )	4.5 $^{+0.1}_{-0.1}$	40.0 $^{+0.4}_{-0.8}$	—	QQP	Grad ( $L_2$ )	4.5 $^{+0.1}_{-0.1}$	40.0 $^{+0.4}_{-0.8}$	—
	Grad ( $L_1$ )	4.5 $^{+0.1}_{-0.1}$	40.0 $^{+0.4}_{-0.8}$	—		Grad ( $L_1$ )	4.5 $^{+0.1}_{-0.1}$	40.0 $^{+0.4}_{-0.8}$	—		Grad ( $L_1$ )	4.5 $^{+0.1}_{-0.1}$	40.0 $^{+0.4}_{-0.8}$	—
	x $\odot$ grad (sign)	-0.8 $^{+0.3}_{-0.6}$	-7.3 $^{+3.0}_{-3.5}$	—		x $\odot$ grad (sign)	-0.8 $^{+0.3}_{-0.6}$	-7.3 $^{+3.0}_{-3.5}$	—		x $\odot$ grad (sign)	-0.8 $^{+0.3}_{-0.6}$	-7.3 $^{+3.0}_{-3.5}$	—
	x $\odot$ grad (abs)	2.8 $^{+0.2}_{-0.2}$	24.5 $^{+1.4}_{-1.5}$	—		x $\odot$ grad (abs)	2.8 $^{+0.2}_{-0.2}$	24.5 $^{+1.4}_{-1.5}$	—		x $\odot$ grad (abs)	2.8 $^{+0.2}_{-0.2}$	24.5 $^{+1.4}_{-1.5}$	—
	IG (sign)	9.1 $^{+0.4}_{-0.3}$	81.3 $^{+2.6}_{-2.3}$	—		IG (sign)	9.1 $^{+0.4}_{-0.3}$	81.3 $^{+2.6}_{-2.3}$	—		IG (sign)	9.1 $^{+0.4}_{-0.3}$	81.3 $^{+2.6}_{-2.3}$	—
	IG (abs)	3.5 $^{+0.2}_{-0.2}$	31.1 $^{+1.3}_{-1.3}$	—		IG (abs)	3.5 $^{+0.2}_{-0.2}$	31.1 $^{+1.3}_{-1.3}$	—		IG (abs)	3.5 $^{+0.2}_{-0.2}$	31.1 $^{+1.3}_{-1.3}$	—
LOO (sign)	22.0 $^{+0.6}_{-0.6}$	195.4 $^{+6.1}_{-7.7}$	—	LOO (sign)	22.0 $^{+0.6}_{-0.6}$	195.4 $^{+6.1}_{-7.7}$	—	LOO (sign)	22.0 $^{+0.6}_{-0.6}$	195.4 $^{+6.1}_{-7.7}$	—			
LOO (abs)	5.6 $^{+0.1}_{-0.2}$	49.6 $^{+0.8}_{-1.8}$	—	LOO (abs)	5.6 $^{+0.1}_{-0.2}$	49.6 $^{+0.8}_{-1.8}$	—	LOO (abs)	5.6 $^{+0.1}_{-0.2}$	49.6 $^{+0.8}_{-1.8}$	—			
Beam	23.3 $^{+0.5}_{-0.6}$	207.3 $^{+6.6}_{-5.7}$	—	Beam	23.3 $^{+0.5}_{-0.6}$	207.3 $^{+6.6}_{-5.7}$	—	Beam	23.3 $^{+0.5}_{-0.6}$	207.3 $^{+6.6}_{-5.7}$	—			

Table C.8 Faithfulness scores for **RoBERTa-large**. Shows Relative Area Between Curves (RACU) and the non-relative variant (ACU), defined in Chapter 3, note that Chapter 3 does not report results for Recursive-ROAR with RoBERTa-large.

Dataset	IM	Faithfulness [%]			Dataset	IM	Faithfulness [%]			Dataset	IM	Faithfulness [%]		
		Our		R-ROAR			Our		R-ROAR			Our		R-ROAR
		ACU	RACU	RACU			ACU	RACU	RACU			ACU	RACU	RACU
bAbI-1	Grad ( $L_2$ )	31.1 <sup>+1.0</sup> <sub>-1.5</sub>	87.6 <sup>+1.6</sup> <sub>-2.3</sub>	–	CoLA	Grad ( $L_2$ )	10.9 <sup>+0.9</sup> <sub>-0.7</sub>	35.0 <sup>+2.9</sup> <sub>-2.5</sub>	–	SST2	Grad ( $L_2$ )	9.8 <sup>+1.1</sup> <sub>-1.0</sub>	32.4 <sup>+3.5</sup> <sub>-3.0</sub>	–
	Grad ( $L_1$ )	31.1 <sup>+1.1</sup> <sub>-1.6</sub>	87.7 <sup>+1.6</sup> <sub>-2.5</sub>	–		Grad ( $L_1$ )	10.4 <sup>+0.9</sup> <sub>-0.6</sub>	33.3 <sup>+3.1</sup> <sub>-2.2</sub>	–		Grad ( $L_1$ )	9.7 <sup>+1.1</sup> <sub>-0.9</sub>	32.0 <sup>+3.2</sup> <sub>-2.9</sub>	–
	x ⊙ grad (sign)	13.9 <sup>+3.1</sup> <sub>-6.5</sub>	39.1 <sup>+8.8</sup> <sub>-16.0</sub>	–		x ⊙ grad (sign)	1.2 <sup>+0.2</sup> <sub>-0.3</sub>	3.8 <sup>+0.5</sup> <sub>-0.9</sub>	–		x ⊙ grad (sign)	−3.4 <sup>+0.7</sup> <sub>-0.6</sub>	−11.4 <sup>+2.4</sup> <sub>-2.2</sub>	–
	x ⊙ grad (abs)	24.0 <sup>+2.8</sup> <sub>-3.2</sub>	67.5 <sup>+5.7</sup> <sub>-8.3</sub>	–		x ⊙ grad (abs)	6.6 <sup>+0.9</sup> <sub>-0.9</sub>	21.3 <sup>+2.9</sup> <sub>-2.9</sub>	–		x ⊙ grad (abs)	5.4 <sup>+1.6</sup> <sub>-1.1</sub>	18.0 <sup>+3.5</sup> <sub>-3.5</sub>	–
	IG (sign)	36.3 <sup>+2.9</sup> <sub>-3.1</sub>	102.4 <sup>+2.7</sup> <sub>-6.9</sub>	–		IG (sign)	17.3 <sup>+1.1</sup> <sub>-1.1</sub>	55.4 <sup>+3.6</sup> <sub>-3.5</sub>	–		IG (sign)	40.1 <sup>+3.2</sup> <sub>-1.9</sub>	133.1 <sup>+10.6</sup> <sub>-6.5</sub>	–
	IG (abs)	31.6 <sup>+1.3</sup> <sub>-1.5</sub>	88.9 <sup>+2.5</sup> <sub>-2.7</sub>	–		IG (abs)	11.7 <sup>+0.4</sup> <sub>-0.3</sub>	37.5 <sup>+1.0</sup> <sub>-0.9</sub>	–		IG (abs)	15.6 <sup>+0.9</sup> <sub>-0.7</sub>	51.6 <sup>+3.4</sup> <sub>-1.5</sub>	–
	LOO (sign)	44.4 <sup>+0.9</sup> <sub>-1.0</sub>	125.2 <sup>+1.7</sup> <sub>-1.7</sub>	–		LOO (sign)	38.9 <sup>+1.1</sup> <sub>-2.0</sub>	124.9 <sup>+11.4</sup> <sub>-4.3</sub>	–		LOO (sign)	49.4 <sup>+1.3</sup> <sub>-1.4</sub>	164.0 <sup>+1.7</sup> <sub>-2.9</sub>	–
	LOO (abs)	32.4 <sup>+0.4</sup> <sub>-0.6</sub>	91.4 <sup>+0.4</sup> <sub>-0.9</sub>	–		LOO (abs)	17.6 <sup>+0.9</sup> <sub>-1.3</sub>	56.7 <sup>+3.2</sup> <sub>-3.4</sub>	–		LOO (abs)	17.2 <sup>+1.0</sup> <sub>-0.7</sub>	57.1 <sup>+3.6</sup> <sub>-2.2</sub>	–
Beam	44.7 <sup>+0.8</sup> <sub>-1.0</sub>	126.0 <sup>+1.1</sup> <sub>-1.1</sub>	–	Beam	37.5 <sup>+4.7</sup> <sub>-2.7</sub>	120.3 <sup>+13.4</sup> <sub>-8.2</sub>	–	Beam	55.6 <sup>+1.0</sup> <sub>-0.5</sub>	184.5 <sup>+2.1</sup> <sub>-2.6</sub>	–			
bAbI-2	Grad ( $L_2$ )	28.2 <sup>+0.9</sup> <sub>-1.5</sub>	94.0 <sup>+6.3</sup> <sub>-3.4</sub>	–	Anemia	Grad ( $L_2$ )	18.6 <sup>+1.7</sup> <sub>-1.5</sub>	47.9 <sup>+3.7</sup> <sub>-3.6</sub>	–	SNLI	Grad ( $L_2$ )	8.2 <sup>+0.3</sup> <sub>-0.5</sub>	53.8 <sup>+1.3</sup> <sub>-2.0</sub>	–
	Grad ( $L_1$ )	28.0 <sup>+0.9</sup> <sub>-1.4</sub>	93.5 <sup>+5.3</sup> <sub>-2.8</sub>	–		Grad ( $L_1$ )	18.4 <sup>+2.0</sup> <sub>-1.8</sub>	47.4 <sup>+4.3</sup> <sub>-4.3</sub>	–		Grad ( $L_1$ )	8.2 <sup>+0.4</sup> <sub>-0.5</sub>	53.3 <sup>+1.6</sup> <sub>-2.1</sub>	–
	x ⊙ grad (sign)	8.2 <sup>+7.0</sup> <sub>-4.9</sub>	26.9 <sup>+20.4</sup> <sub>-19.4</sub>	–		x ⊙ grad (sign)	4.6 <sup>+1.8</sup> <sub>-1.6</sub>	11.9 <sup>+4.6</sup> <sub>-4.2</sub>	–		x ⊙ grad (sign)	−0.3 <sup>+0.3</sup> <sub>-0.3</sub>	−2.2 <sup>+1.7</sup> <sub>-2.2</sub>	–
	x ⊙ grad (abs)	22.6 <sup>+1.5</sup> <sub>-1.5</sub>	75.5 <sup>+7.7</sup> <sub>-5.2</sub>	–		x ⊙ grad (abs)	11.6 <sup>+2.1</sup> <sub>-2.1</sub>	29.8 <sup>+5.5</sup> <sub>-5.5</sub>	–		x ⊙ grad (abs)	5.6 <sup>+0.3</sup> <sub>-0.4</sub>	36.5 <sup>+1.7</sup> <sub>-2.0</sub>	–
	IG (sign)	37.9 <sup>+1.3</sup> <sub>-1.3</sub>	126.6 <sup>+5.8</sup> <sub>-9.7</sub>	–		IG (sign)	58.2 <sup>+2.9</sup> <sub>-4.0</sub>	150.1 <sup>+4.2</sup> <sub>-5.0</sub>	–		IG (sign)	14.0 <sup>+0.3</sup> <sub>-0.3</sub>	91.2 <sup>+1.0</sup> <sub>-1.0</sub>	–
	IG (abs)	27.4 <sup>+1.7</sup> <sub>-1.7</sub>	91.7 <sup>+6.8</sup> <sub>-9.5</sub>	–		IG (abs)	16.4 <sup>+1.2</sup> <sub>-1.9</sub>	42.3 <sup>+2.4</sup> <sub>-3.7</sub>	–		IG (abs)	8.0 <sup>+0.4</sup> <sub>-0.5</sub>	52.6 <sup>+1.4</sup> <sub>-1.7</sub>	–
	LOO (sign)	40.6 <sup>+1.9</sup> <sub>-0.7</sub>	135.4 <sup>+4.0</sup> <sub>-4.0</sub>	–		LOO (sign)	60.5 <sup>+3.3</sup> <sub>-3.3</sub>	156.0 <sup>+7.4</sup> <sub>-7.7</sub>	–		LOO (sign)	26.3 <sup>+0.4</sup> <sub>-0.5</sub>	172.3 <sup>+1.3</sup> <sub>-1.3</sub>	–
	LOO (abs)	28.1 <sup>+0.9</sup> <sub>-0.9</sub>	93.9 <sup>+2.8</sup> <sub>-1.8</sub>	–		LOO (abs)	15.6 <sup>+1.6</sup> <sub>-1.7</sub>	40.2 <sup>+3.5</sup> <sub>-3.7</sub>	–		LOO (abs)	11.0 <sup>+0.5</sup> <sub>-0.3</sub>	71.8 <sup>+1.7</sup> <sub>-0.9</sub>	–
Beam	41.7 <sup>+1.8</sup> <sub>-0.7</sub>	139.2 <sup>+4.3</sup> <sub>-4.3</sub>	–	Beam	–	–	–	Beam	29.6 <sup>+0.3</sup> <sub>-0.3</sub>	193.7 <sup>+5.2</sup> <sub>-4.5</sub>	–			
bAbI-3	Grad ( $L_2$ )	22.4 <sup>+3.8</sup> <sub>-3.8</sub>	94.7 <sup>+0.2</sup> <sub>-0.2</sub>	–	Diabetes	Grad ( $L_2$ )	20.1 <sup>+1.8</sup> <sub>-1.1</sub>	90.7 <sup>+0.5</sup> <sub>-0.4</sub>	–	IMDB	Grad ( $L_2$ )	13.9 <sup>+3.5</sup> <sub>-1.9</sub>	29.4 <sup>+6.0</sup> <sub>-2.4</sub>	–
	Grad ( $L_1$ )	22.2 <sup>+3.6</sup> <sub>-3.6</sub>	94.1 <sup>+0.4</sup> <sub>-0.4</sub>	–		Grad ( $L_1$ )	20.1 <sup>+1.9</sup> <sub>-1.0</sub>	90.7 <sup>+0.6</sup> <sub>-0.6</sub>	–		Grad ( $L_1$ )	13.7 <sup>+4.1</sup> <sub>-1.9</sub>	28.9 <sup>+6.0</sup> <sub>-2.5</sub>	–
	x ⊙ grad (sign)	8.5 <sup>+3.5</sup> <sub>-3.5</sub>	34.4 <sup>+9.0</sup> <sub>-9.0</sub>	–		x ⊙ grad (sign)	7.4 <sup>+3.1</sup> <sub>-4.0</sub>	34.2 <sup>+16.6</sup> <sub>-19.0</sub>	–		x ⊙ grad (sign)	−2.9 <sup>+0.4</sup> <sub>-0.4</sub>	−6.4 <sup>+1.3</sup> <sub>-1.2</sub>	–
	x ⊙ grad (abs)	19.9 <sup>+4.6</sup> <sub>-4.6</sub>	83.2 <sup>+5.5</sup> <sub>-5.5</sub>	–		x ⊙ grad (abs)	18.6 <sup>+5.0</sup> <sub>-4.9</sub>	84.2 <sup>+0.8</sup> <sub>-0.8</sub>	–		x ⊙ grad (abs)	7.7 <sup>+1.2</sup> <sub>-1.2</sub>	16.3 <sup>+3.9</sup> <sub>-2.8</sub>	–
	IG (sign)	33.0 <sup>+3.8</sup> <sub>-3.8</sub>	141.3 <sup>+2.7</sup> <sub>-2.7</sub>	–		IG (sign)	25.3 <sup>+4.0</sup> <sub>-4.0</sub>	113.5 <sup>+7.9</sup> <sub>-7.9</sub>	–		IG (sign)	53.2 <sup>+3.4</sup> <sub>-4.1</sub>	114.2 <sup>+12.8</sup> <sub>-12.8</sub>	–
	IG (abs)	24.3 <sup>+3.7</sup> <sub>-3.7</sub>	103.2 <sup>+1.5</sup> <sub>-1.5</sub>	–		IG (abs)	19.6 <sup>+1.4</sup> <sub>-1.0</sub>	88.7 <sup>+0.5</sup> <sub>-1.2</sub>	–		IG (abs)	18.9 <sup>+3.7</sup> <sub>-4.2</sub>	40.3 <sup>+6.2</sup> <sub>-9.8</sub>	–
	LOO (sign)	35.0 <sup>+3.0</sup> <sub>-3.0</sub>	150.5 <sup>+12.4</sup> <sub>-12.4</sub>	–		LOO (sign)	34.8 <sup>+6.7</sup> <sub>-4.9</sub>	156.0 <sup>+12.2</sup> <sub>-18.3</sub>	–		LOO (sign)	60.5 <sup>+1.0</sup> <sub>-1.0</sub>	130.1 <sup>+13.0</sup> <sub>-13.0</sub>	–
	LOO (abs)	23.3 <sup>+4.2</sup> <sub>-4.2</sub>	98.7 <sup>+1.2</sup> <sub>-1.2</sub>	–		LOO (abs)	19.2 <sup>+1.4</sup> <sub>-1.0</sub>	86.6 <sup>+0.6</sup> <sub>-0.9</sub>	–		LOO (abs)	16.7 <sup>+5.0</sup> <sub>-1.9</sub>	35.5 <sup>+5.6</sup> <sub>-5.6</sub>	–
Beam	–	–	–	Beam	–	–	–	Beam	–	–	–			
BoolQ	Grad ( $L_2$ )	2.6 <sup>+0.1</sup> <sub>-0.3</sub>	24.8 <sup>+1.7</sup> <sub>-1.9</sub>	–	MRPC	Grad ( $L_2$ )	6.6 <sup>+2.2</sup> <sub>-1.1</sub>	22.9 <sup>+3.1</sup> <sub>-4.3</sub>	–	MNLI	Grad ( $L_2$ )	7.9 <sup>+0.1</sup> <sub>-0.2</sub>	38.7 <sup>+0.9</sup> <sub>-1.3</sub>	–
	Grad ( $L_1$ )	2.7 <sup>+0.2</sup> <sub>-0.3</sub>	25.3 <sup>+2.3</sup> <sub>-1.9</sub>	–		Grad ( $L_1$ )	6.6 <sup>+1.9</sup> <sub>-1.1</sub>	22.7 <sup>+1.5</sup> <sub>-3.7</sub>	–		Grad ( $L_1$ )	7.8 <sup>+0.2</sup> <sub>-0.2</sub>	38.3 <sup>+1.0</sup> <sub>-1.5</sub>	–
	x ⊙ grad (sign)	−0.4 <sup>+0.2</sup> <sub>-0.1</sub>	−3.6 <sup>+1.6</sup> <sub>-1.5</sub>	–		x ⊙ grad (sign)	−0.9 <sup>+1.0</sup> <sub>-1.1</sub>	−3.4 <sup>+3.2</sup> <sub>-2.2</sub>	–		x ⊙ grad (sign)	−0.5 <sup>+0.5</sup> <sub>-0.2</sub>	−2.3 <sup>+0.6</sup> <sub>-1.1</sub>	–
	x ⊙ grad (abs)	1.2 <sup>+0.3</sup> <sub>-0.3</sub>	10.8 <sup>+2.0</sup> <sub>-2.0</sub>	–		x ⊙ grad (abs)	4.4 <sup>+1.2</sup> <sub>-1.0</sub>	15.7 <sup>+2.6</sup> <sub>-4.8</sub>	–		x ⊙ grad (abs)	5.2 <sup>+0.1</sup> <sub>-0.1</sub>	25.4 <sup>+0.8</sup> <sub>-1.1</sub>	–
	IG (sign)	6.9 <sup>+1.2</sup> <sub>-0.8</sub>	65.6 <sup>+15.7</sup> <sub>-6.9</sub>	–		IG (sign)	15.7 <sup>+1.6</sup> <sub>-2.0</sub>	55.9 <sup>+9.5</sup> <sub>-8.8</sub>	–		IG (sign)	18.6 <sup>+0.8</sup> <sub>-1.1</sub>	91.1 <sup>+4.3</sup> <sub>-4.6</sub>	–
	IG (abs)	3.2 <sup>+0.4</sup> <sub>-0.8</sub>	30.3 <sup>+4.6</sup> <sub>-5.0</sub>	–		IG (abs)	8.1 <sup>+1.3</sup> <sub>-1.1</sub>	28.8 <sup>+3.7</sup> <sub>-3.7</sub>	–		IG (abs)	9.0 <sup>+0.3</sup> <sub>-0.2</sub>	44.2 <sup>+1.5</sup> <sub>-0.7</sub>	–
	LOO (sign)	25.6 <sup>+1.7</sup> <sub>-2.7</sub>	242.9 <sup>+32.2</sup> <sub>-26.4</sub>	–		LOO (sign)	29.9 <sup>+0.8</sup> <sub>-1.3</sub>	110.0 <sup>+30.9</sup> <sub>-24.2</sub>	–		LOO (sign)	33.0 <sup>+1.0</sup> <sub>-1.0</sub>	161.9 <sup>+6.1</sup> <sub>-8.4</sub>	–
	LOO (abs)	6.2 <sup>+0.6</sup> <sub>-0.5</sub>	58.0 <sup>+3.3</sup> <sub>-1.2</sub>	–		LOO (abs)	10.2 <sup>+1.6</sup> <sub>-1.6</sub>	36.1 <sup>+8.8</sup> <sub>-5.2</sub>	–		LOO (abs)	12.4 <sup>+0.1</sup> <sub>-0.2</sub>	60.6 <sup>+1.4</sup> <sub>-0.9</sub>	–
Beam	21.5 <sup>+1.3</sup> <sub>-2.4</sub>	204.0 <sup>+20.8</sup> <sub>-21.7</sub>	–	Beam	40.1 <sup>+1.4</sup> <sub>-2.3</sub>	146.9 <sup>+41.5</sup> <sub>-30.6</sub>	–	Beam	41.2 <sup>+0.9</sup> <sub>-1.2</sub>	201.7 <sup>+6.5</sup> <sub>-9.1</sub>	–			
CB	Grad ( $L_2$ )	10.1 <sup>+4.4</sup> <sub>-4.2</sub>	30.8 <sup>+15.2</sup> <sub>-12.9</sub>	–	RTE	Grad ( $L_2$ )	7.1 <sup>+0.8</sup> <sub>-1.4</sub>	38.8 <sup>+5.9</sup> <sub>-5.9</sub>	–	QNLI	Grad ( $L_2$ )	9.5 <sup>+0.5</sup> <sub>-0.4</sub>	28.3 <sup>+1.4</sup> <sub>-0.9</sub>	–
	Grad ( $L_1$ )	8.9 <sup>+3.9</sup> <sub>-4.3</sub>	27.3 <sup>+12.6</sup> <sub>-18.2</sub>	–		Grad ( $L_1$ )	7.4 <sup>+0.9</sup> <sub>-1.2</sub>	40.4 <sup>+5.1</sup> <sub>-7.6</sub>	–		Grad ( $L_1$ )	9.4 <sup>+0.3</sup> <sub>-0.3</sub>	28.0 <sup>+1.3</sup> <sub>-0.5</sub>	–
	x ⊙ grad (sign)	−0.1 <sup>+2.7</sup> <sub>-3.3</sub>	0.6 <sup>+9.5</sup> <sub>-8.1</sub>	–		x ⊙ grad (sign)	−0.1 <sup>+1.0</sup> <sub>-1.3</sub>	−1.7 <sup>+4.8</sup> <sub>-7.2</sub>	–		x ⊙ grad (sign)	−1.4 <sup>+0.5</sup> <sub>-0.5</sub>	−4.1 <sup>+1.4</sup> <sub>-1.3</sub>	–
	x ⊙ grad (abs)	5.6 <sup>+1.8</sup> <sub>-2.9</sub>	17.4 <sup>+6.1</sup> <sub>-9.8</sub>	–		x ⊙ grad (abs)	5.0 <sup>+1.3</sup> <sub>-1.0</sub>	27.1 <sup>+6.8</sup> <sub>-3.0</sub>	–		x ⊙ grad (abs)	6.5 <sup>+0.3</sup> <sub>-0.2</sub>	19.1 <sup>+0.7</sup> <sub>-0.6</sub>	–
	IG (sign)	28.5 <sup>+3.9</sup> <sub>-3.4</sub>	85.3 <sup>+19.0</sup> <sub>-14.0</sub>	–		IG (sign)	22.6 <sup>+2.3</sup> <sub>-2.0</sub>	127.1 <sup>+28.9</sup> <sub>-31.6</sub>	–		IG (sign)	25.0 <sup>+3.3</sup> <sub>-5.1</sub>	74.0 <sup>+9.6</sup> <sub>-13.7</sub>	–
	IG (abs)	17.1 <sup>+3.2</sup> <sub>-2.2</sub>	51.3 <sup>+10.5</sup> <sub>-9.0</sub>	–		IG (abs)	7.7 <sup>+1.2</sup> <sub>-1.3</sub>	43.0 <sup>+8.8</sup> <sub>-10.7</sub>	–		IG (abs)	9.7 <sup>+1.2</sup> <sub>-1.8</sub>	28.7 <sup>+3.2</sup> <sub>-5.0</sub>	–
	LOO (sign)	39.0 <sup>+2.9</sup> <sub>-2.4</sub>	116.1 <sup>+9.5</sup> <sub>-15.8</sub>	–		LOO (sign)	38.5 <sup>+1.9</sup> <sub>-3.6</sub>	213.6 <sup>+37.1</sup> <sub>-40.5</sub>	–		LOO (sign)	40.4 <sup>+1.5</sup> <sub>-2.1</sub>	119.8 <sup>+4.0</sup> <sub>-5.0</sub>	–
	LOO (abs)	19.0 <sup>+4.0</sup> <sub>-2.4</sub>	56.6 <sup>+8.8</sup> <sub>-12.8</sub>	–		LOO (abs)	9.9 <sup>+0.4</sup> <sub>-0.4</sub>	55.4 <sup>+11.9</sup> <sub>-30.9</sub>	–		LOO (abs)	12.5 <sup>+0.2</sup> <sub>-0.2</sub>	37.0 <sup>+0.3</sup> <sub>-0.4</sub>	–
Beam	51.8 <sup>+3.0</sup> <sub>-3.0</sub>	154.6 <sup>+14.0</sup> <sub>-22.6</sub>	–	Beam	50.3 <sup>+1.2</sup> <sub>-2.9</sub>	280.1 <sup>+57.9</sup> <sub>-45.7</sub>	–	Beam	53.8 <sup>+1.7</sup> <sub>-2.1</sub>	159.5 <sup>+4.6</sup> <sub>-4.4</sub>	–			
QQP	Grad ( $L_2$ )	–	–	–	QQP	Grad ( $L_2$ )	–	–	–	QQP	Grad ( $L_2$ )	4.0 <sup>+0.3</sup> <sub>-0.2</sub>	33.5 <sup>+1.4</sup> <sub>-1.9</sub>	–
	Grad ( $L_1$ )	–	–	–		Grad ( $L_1$ )	–	–	–		Grad ( $L_1$ )	4.0 <sup>+0.3</sup> <sub>-0.2</sub>	33.0 <sup>+1.3</sup> <sub>-1.9</sub>	–
	x ⊙ grad (sign)	–	–	–		x ⊙ grad (sign)	–	–	–		x ⊙ grad (sign)	−0.4 <sup>+0.3</sup> <sub>-0.2</sub>	−3.3 <sup>+2.3</sup> <sub>-2.0</sub>	–
	x ⊙ grad (abs)	–	–	–		x ⊙ grad (abs)	–	–	–		x ⊙ grad (abs)	2.5 <sup>+0.3</sup> <sub>-0.2</sub>	20.7 <sup>+1.3</sup> <sub>-2.2</sub>	–
	IG (sign)	–	–	–		IG (sign)	–	–	–		IG (sign)	8.9 <sup>+1.0</sup> <sub>-0.6</sub>	73.7 <sup>+4.8</sup> <sub>-3.6</sub>	–
	IG (abs)	–	–	–		IG (abs)	–	–	–		IG (abs)	3.8 <sup>+0.7</sup> <sub>-0.3</sub>	31.6 <sup>+2.5</sup> <sub>-1.6</sub>	–
	LOO (sign)	–	–	–		LOO (sign)	–	–	–		LOO (sign)	20.4 <sup>+0.7</sup> <sub>-0.4</sub>	169.8 <sup>+8.6</sup> <sub>-11.1</sub>	–
	LOO (abs)	–	–	–		LOO (abs)	–	–	–		LOO (abs)	5.7 <sup>+0.2</sup> <sub>-0.3</sub>	47.3 <sup>+2.3</sup> <sub>-2.8</sub>	–
Beam	–	–	–	Beam	–	–	–	Beam	22.5 <sup>+0.7</sup> <sub>-0.8</sub>	187.0 <sup>+10.4</sup> <sub>-10.3</sub>	–			

Dataset	Model	Size	Inference time [hh:mm]			
			Classify	Counterfactual	Redacted	Feature
IMDB	Llama 2	70B	10:14	128:52	78:21	275:29
		7B	03:03	62:56	22:59	102:26
	Falcon	40B	09:53	55:34	93:18	34:26
		7B	06:28	80:38	446:40	112:25
	Mistral v0.1	7B	02:15	61:17	39:09	110:36
MCTest	Llama 2	70B	00:31	04:52	01:53	05:46
		7B	00:11	02:26	00:40	02:35
	Falcon	40B	00:14	01:02	01:22	00:44
		7B	00:10	00:36	01:05	01:27
	Mistral v0.1	7B	00:06	02:13	00:25	01:24
RTE	Llama 2	70B	00:08	00:47	00:51	00:55
		7B	00:01	00:18	00:11	00:18
	Falcon	40B	00:09	00:31	00:34	00:26
		7B	00:02	00:09	00:21	00:08
	Mistral v0.1	7B	00:01	00:13	00:11	00:14
bAbI-1	Llama 2	70B	00:53	03:03	02:35	03:19
		7B	00:10	01:01	00:41	00:56
	Falcon	40B	00:31	01:25	01:24	01:22
		7B	00:10	00:31	00:29	00:37
	Mistral v0.1	7B	00:06	00:32	00:30	00:37

Table D.1 Inference time as reported by TGI. Note that this does not correspond to wall-time. In particular, because 50 prompts are computed in parallel. Dividing the inference time by 50 is a decent approximation for wall-time.

## APPENDIX D FAITHFULNESS OF SELF-EXPLANATIONS

### D.1 Compute

The specifications for the compute hardware are provided in Table D.2. The electricity is from 99% hydroelectric power.

40B and 70B models	
CPU	6 cores, AMD Milan 7413 @ 2.65 GHz
GPU	4x NVidia A100SXM4 (40 GB), NVLink
Memory	128 GB
7B models	
CPU	6 cores, AMD Milan 7413 @ 2.65 GHz
GPU	1x NVidia A100SXM4 (40 GB)
Memory	32 GB

Table D.2 The computing hardware used. Only the allocated resources are reported.

The inference time is stated in Table D.1. Note that this is the inference time as reported by TGI and does not correspond to the wall-time. In particular, because prompt responses are generated in parallel (50 prompts simultaneously). Additionally, prompts that have been computed before are only computed once; for example, the response for the classification prompt is reused in each explanation variation.

## D.2 Prompt examples

This section provides output examples for all the prompts used in this thesis. The outputs are from the Llama2-70B chat-tuned model. Simple heuristics extract the content, such as the predicted label or the explanation. Correct extraction has been manually verified for every observation.

### D.2.1 Sentiment classification with IMDB

#### Counterfactual

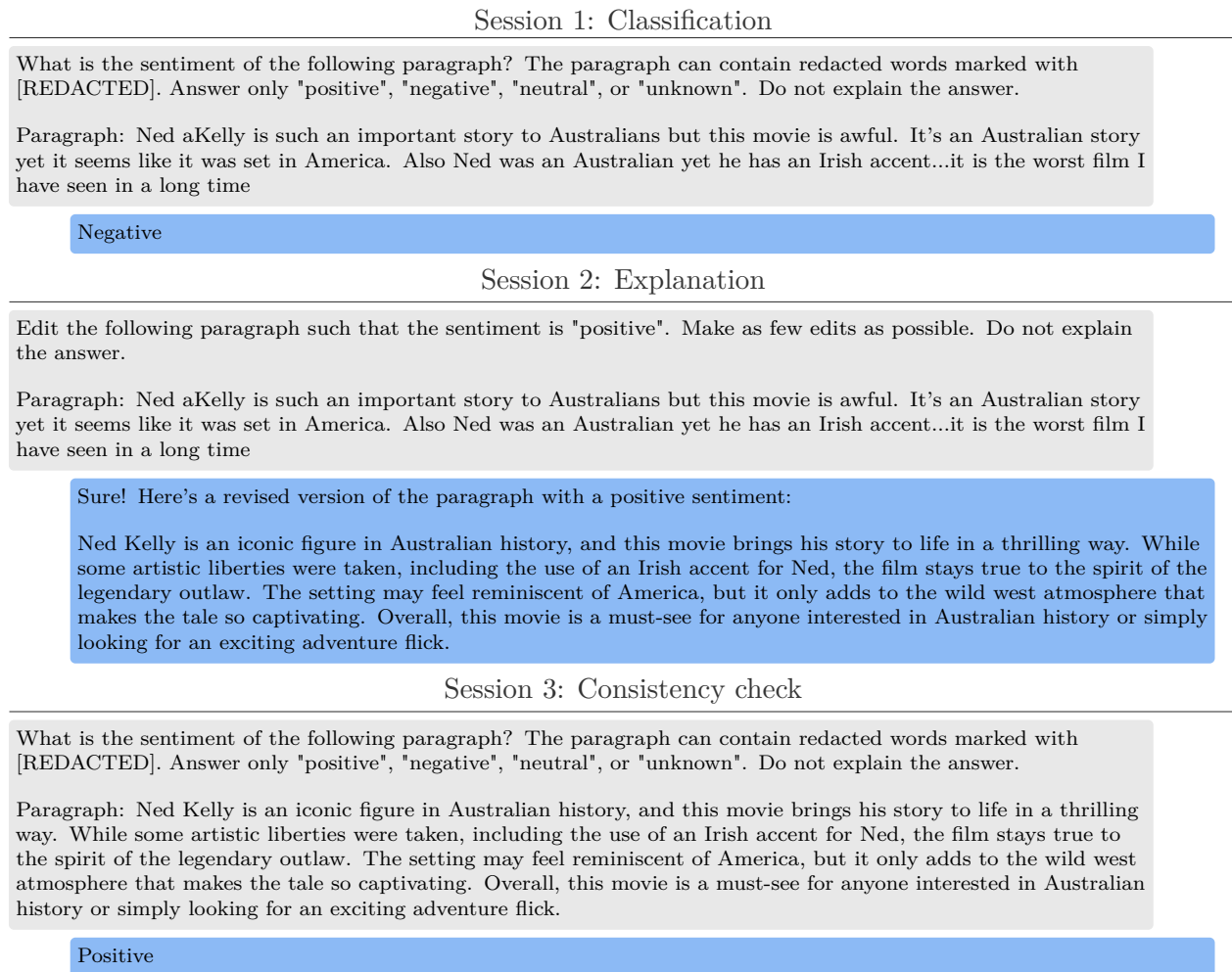


Figure D.1 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: explicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

---

### Session 1: Classification

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

---

Edit the following paragraph such a human would classify the sentiment is "positive". Make as few edits as possible. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's a revised version of the paragraph that a human would likely classify as having a positive sentiment:

Ned Kelly is such an iconic figure in Australian history, and this movie does his story justice. While some artistic liberties were taken, such as setting the film in America and giving Ned an Irish accent, the spirit of the story remains true to its roots. Overall, it's a powerful and moving film that I highly recommend.

### Session 3: Consistency check

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned Kelly is such an iconic figure in Australian history, and this movie does his story justice. While some artistic liberties were taken, such as setting the film in America and giving Ned an Irish accent, the spirit of the story remains true to its roots. Overall, it's a powerful and moving film that I highly recommend.

Positive

Figure D.2 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: explicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

Edit the following paragraph such you would classify the sentiment is "positive". Make as few edits as possible. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's a revised version of the paragraph with a positive sentiment:

Ned Kelly is an iconic figure in Australian history, and this movie brings his story to life in a thrilling way. While some artistic liberties were taken, such as the Irish accent, the film stays true to the spirit of Ned's legacy. The cinematography is stunning, capturing the beauty of the Australian outback, and the acting performances are top-notch. Overall, it's a must-see for anyone interested in Australian history or looking for a compelling drama.

### Session 3: Consistency check

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned Kelly is an iconic figure in Australian history, and this movie brings his story to life in a thrilling way. While some artistic liberties were taken, such as the Irish accent, the film stays true to the spirit of Ned's legacy. The cinematography is stunning, capturing the beauty of the Australian outback, and the acting performances are top-notch. Overall, it's a must-see for anyone interested in Australian history or looking for a compelling drama.

Positive

Figure D.3 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: explicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

What is the sentiment of the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

Edit the following paragraph such that the sentiment becomes the opposite of what it currently is. Make as few edits as possible. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Paragraph: Ned Kelly is such an unimportant story to Australians and this movie is fantastic. It's an American story yet it seems like it was set in Australia. Also, Ned was Irish yet he has an Australian accent...it is the best film I have seen in a long time.

### Session 3: Consistency check

What is the sentiment of the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned Kelly is such an unimportant story to Australians and this movie is fantastic. It's an American story yet it seems like it was set in Australia. Also, Ned was Irish yet he has an Australian accent...it is the best film I have seen in a long time.

Positive

Figure D.4 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: implicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

Edit the following paragraph such a human would classify the as the opposite of what it currently is. Make as few edits as possible. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Paragraph: Ned Kelly is such an unimportant story to Americans, but this movie is amazing. It's an American story yet it seems like it was set in Australia. Also, Ned was an American, yet he has an Australian accent...it is the best film I have seen in a long time.

### Session 3: Consistency check

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned Kelly is such an unimportant story to Americans, but this movie is amazing. It's an American story yet it seems like it was set in Australia. Also, Ned was an American, yet he has an Australian accent...it is the best film I have seen in a long time.

Positive

Figure D.5 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: implicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

Edit the following paragraph such you would classify the sentiment as the opposite of what it currently is. Make as few edits as possible. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's a revised version of the paragraph with a opposite sentiment:

Ned Kelly is such an iconic figure in Australian history, and this movie does his story justice. It's a powerful tale that highlights the struggles and triumphs of the Australian people. The setting is authentic and captures the spirit of the era, and Ned's Irish accent adds a unique touch to his character. Overall, it's a fantastic film that I highly recommend.

### Session 3: Consistency check

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned Kelly is such an iconic figure in Australian history, and this movie does his story justice. It's a powerful tale that highlights the struggles and triumphs of the Australian people. The setting is authentic and captures the spirit of the era, and Ned's Irish accent adds a unique touch to his character. Overall, it's a fantastic film that I highly recommend.

Positive

Figure D.6 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: implicit”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluted to be “faithful”.

## Feature attribution

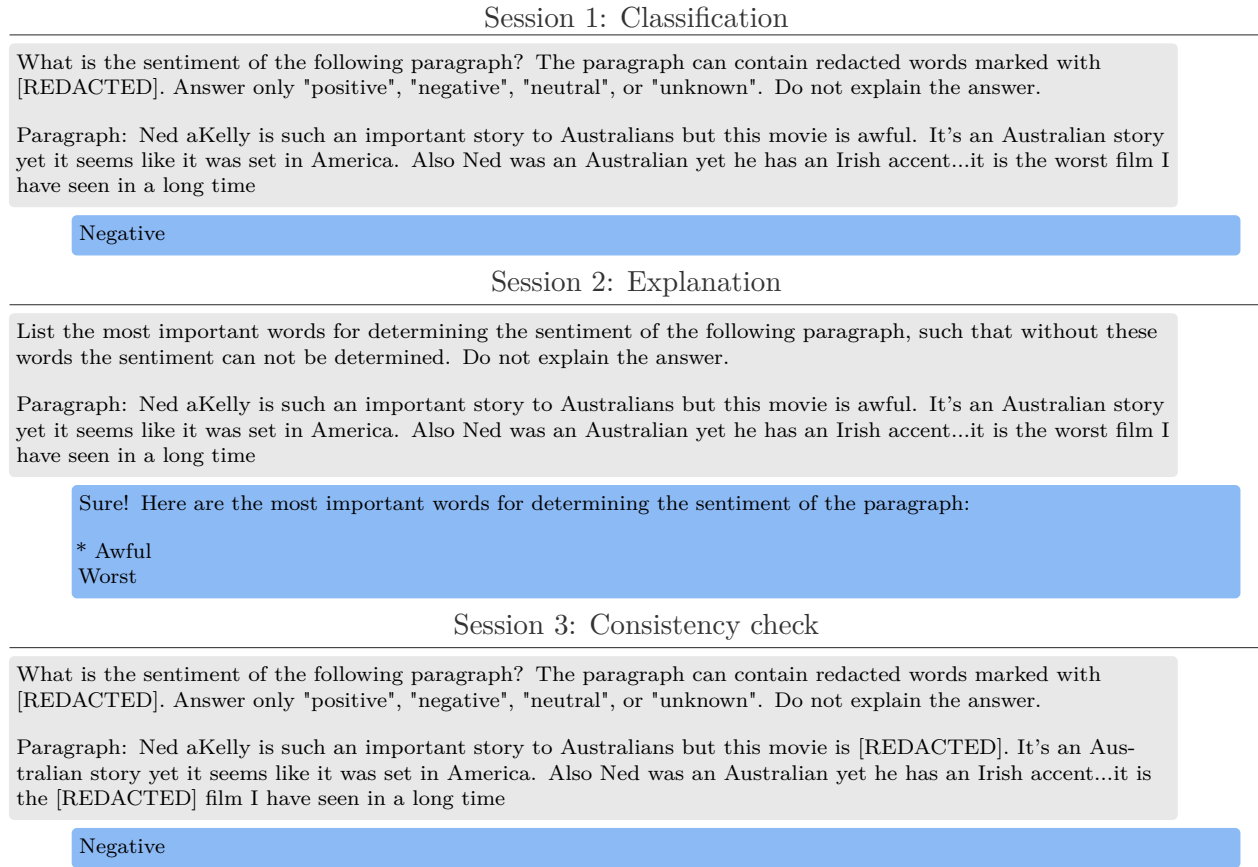


Figure D.7 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

List the most important words for determining the sentiment of the following paragraph, such that without these words a human can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here are the most important words for determining the sentiment of the given paragraph:

1. Awful
2. Worst

### Session 3: Consistency check

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is [REDACTED]. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the [REDACTED] film I have seen in a long time

Negative

Figure D.8 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

List the most important words for determining the sentiment of the following paragraph, such that without these words you can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here are the most important words for determining the sentiment of the paragraph:

\* Awful  
Worst

### Session 3: Consistency check

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is [REDACTED]. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the [REDACTED] film I have seen in a long time

Negative

Figure D.9 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

What is the sentiment of the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

List the most important words for determining the sentiment of the following paragraph, such that without these words the sentiment can not be determined. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here are the most important words for determining the sentiment of the paragraph:

\* Awful  
Worst

### Session 3: Consistency check

What is the sentiment of the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is [REMOVED]. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the [REMOVED] film I have seen in a long time

Negative

Figure D.10 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

Session 1: Classification

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

---

Session 2: Explanation

---

List the most important words for determining the sentiment of the following paragraph, such that without these words a human can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here are the most important words for determining the sentiment of the given paragraph:

1. Awful
2. Worst

---

Session 3: Consistency check

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is [REMOVED]. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the [REMOVED] film I have seen in a long time

Negative

Figure D.11 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

What would you classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

List the most important words for determining the sentiment of the following paragraph, such that without these words you can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here are the most important words for determining the sentiment of the paragraph:

\* Awful  
Worst

### Session 3: Consistency check

What would you classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is [REMOVED]. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the [REMOVED] film I have seen in a long time

Negative

Figure D.12 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

## Redacted

### Session 1: Classification

What is the sentiment of the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

Redact the most important words for determining the sentiment of the following paragraph, by replacing important words with [REDACTED], such that without these words the sentiment can not be determined. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's the redacted paragraph:

Ned [REDACTED] is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in [REDACTED]. Also Ned was an Australian yet he has a [REDACTED] accent...it is the worst film I have seen in a long time

### Session 3: Consistency check

What is the sentiment of the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned [REDACTED] is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in [REDACTED]. Also Ned was an Australian yet he has a [REDACTED] accent...it is the worst film I have seen in a long time

Negative

Figure D.13 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

Session 1: Classification

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

---

Session 2: Explanation

---

Redact the most important words for determining the sentiment of the following paragraph, by replacing important words with [REDACTED], such that without these words a human can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's the redacted version of the paragraph:

Ned [REDACTED] is such an important story to [REDACTED] but this movie is awful. It's a [REDACTED] story yet it seems like it was set in [REDACTED]. Also [REDACTED] was a [REDACTED] yet he has an [REDACTED] accent...it is the worst film I have seen in a long time

---

Session 3: Consistency check

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned [REDACTED] is such an important story to [REDACTED] but this movie is awful. It's a [REDACTED] story yet it seems like it was set in [REDACTED]. Also [REDACTED] was a [REDACTED] yet he has an [REDACTED] accent...it is the worst film I have seen in a long time

Negative

Figure D.14 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

Session 1: Classification

---

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

---

Session 2: Explanation

---

Redact the most important words for determining the sentiment of the following paragraph, by replacing important words with [REDACTED], such that without these words you can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's the redacted version of the paragraph:

Ned [REDACTED] is such an important story to [REDACTED] but this movie is awful. It's a [REDACTED] story yet it seems like it was set in [REDACTED]. Also [REDACTED] was a [REDACTED] yet he has an [REDACTED] accent...it is the worst film I have seen in a long time

---

Session 3: Consistency check

---

What would you classify the sentiment of the following paragraph as? The paragraph can contain redacted words marked with [REDACTED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned [REDACTED] is such an important story to [REDACTED] but this movie is awful. It's a [REDACTED] story yet it seems like it was set in [REDACTED]. Also [REDACTED] was a [REDACTED] yet he has an [REDACTED] accent...it is the worst film I have seen in a long time

Negative

Figure D.15 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "redacted"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

### Session 1: Classification

---

What is the sentiment of the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

### Session 2: Explanation

---

Redact the most important words for determining the sentiment of the following paragraph, by replacing important words with [REMOVED], such that without these words the sentiment can not be determined. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's the redacted paragraph:

Ned [REMOVED] is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in [REMOVED]. Also Ned was an Australian yet he has a [REMOVED] accent...it is the worst film I have seen in a long time

### Session 3: Consistency check

---

What is the sentiment of the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned [REMOVED] is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in [REMOVED]. Also Ned was an Australian yet he has a [REMOVED] accent...it is the worst film I have seen in a long time

Negative

Figure D.16 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

Session 1: Classification

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

---

Session 2: Explanation

---

Redact the most important words for determining the sentiment of the following paragraph, by replacing important words with [REMOVED], such that without these words a human can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's the redacted version of the paragraph:

Ned [REMOVED] is such an important story to [REMOVED] but this movie is awful. It's a [REMOVED] story yet it seems like it was set in [REMOVED]. Also [REMOVED] was a [REMOVED] yet he has a [REMOVED] accent...it is the worst film I have seen in a long time

---

Session 3: Consistency check

---

What would a human classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned [REMOVED] is such an important story to [REMOVED] but this movie is awful. It's a [REMOVED] story yet it seems like it was set in [REMOVED]. Also [REMOVED] was a [REMOVED] yet he has a [REMOVED] accent...it is the worst film I have seen in a long time

Negative

Figure D.17 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

Session 1: Classification

---

What would you classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Negative

---

Session 2: Explanation

---

Redact the most important words for determining the sentiment of the following paragraph, by replacing important words with [REMOVED], such that without these words you can not determine the sentiment. Do not explain the answer.

Paragraph: Ned aKelly is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in America. Also Ned was an Australian yet he has an Irish accent...it is the worst film I have seen in a long time

Sure! Here's the redacted paragraph:

Ned [REMOVED] is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in [REMOVED]. Also Ned was an Australian yet he has a [REMOVED] accent...it is the worst film I have seen in a long time

---

Session 3: Consistency check

---

What would you classify the sentiment of the following paragraph as? The paragraph can contain removed words marked with [REMOVED]. Answer only "positive", "negative", "neutral", or "unknown". Do not explain the answer.

Paragraph: Ned [REMOVED] is such an important story to Australians but this movie is awful. It's an Australian story yet it seems like it was set in [REMOVED]. Also Ned was an Australian yet he has a [REMOVED] accent...it is the worst film I have seen in a long time

Negative

Figure D.18 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "removed"”. The true label is “negative”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

## D.2.2 Multi Choice classification with bAbI-1

### Counterfactual

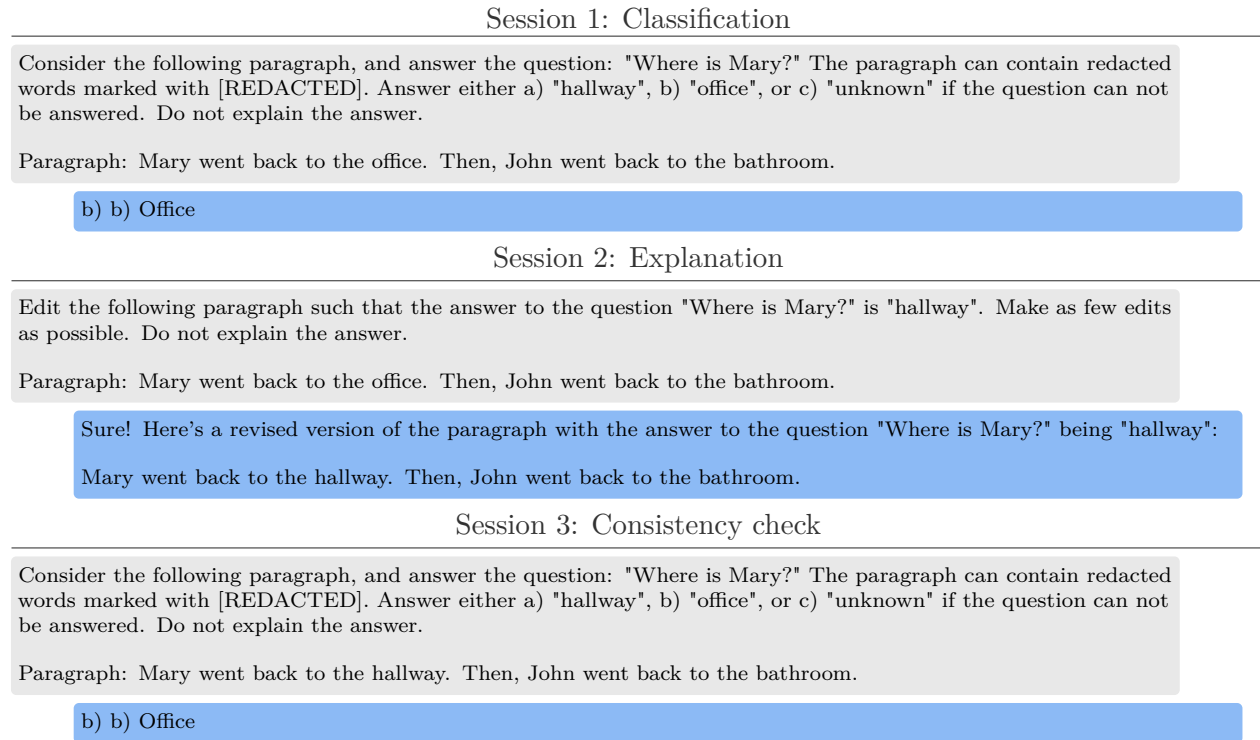


Figure D.19 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: explicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

### Session 1: Classification

---

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) The human would answer: b) Office

### Session 2: Explanation

---

Edit the following paragraph such a human would answer the question "Where is Mary?" with "hallway". Make as few edits as possible. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's an edited version of the paragraph that would elicit the response "hallway" when asked "Where is Mary?":

Mary went back to the hallway. Then, John went back to the bathroom.

### Session 3: Consistency check

---

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the hallway. Then, John went back to the bathroom.

b) The human would answer: b) Office

Figure D.20 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: explicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

Edit the following paragraph such you would answer the question "Where is Mary?" with "hallway". Make as few edits as possible. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's a revised version of the paragraph that answers the question "Where is Mary?" with "hallway":

Mary went back to the hallway. Then, John went back to the bathroom.

### Session 3: Consistency check

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the hallway. Then, John went back to the bathroom.

b) b) Office

Figure D.21 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Counterfactual target: explicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

Edit the following paragraph such that the answer to the question "Where is Mary?" is different from how it should currently be answered. The possible answer options, including the current answer, are a) "hallway", and b) "office". Make as few edits as possible. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Here's a revised version of the paragraph with a different answer to the question "Where is Mary?":

Mary went back to the hallway. Then, John went back to the bathroom.

In this version, the answer to the question "Where is Mary?" would be (b) "hallway", whereas in the original paragraph, the answer would have been (a) "office".

### Session 3: Consistency check

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the hallway. Then, John went back to the bathroom.

In this version, the answer to the question "Where is Mary?" would be (b) "hallway", whereas in the original paragraph, the answer would have been (a) "office".

Sure! Here's a new paragraph with a redacted word:

Paragraph: Mary went back to the [REDACTED]. Then, John went back to the bathroom.

The answer to the question "Where is Mary?" would be (c) "unknown".

Figure D.22 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: implicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) The human would answer: b) Office

### Session 2: Explanation

Edit the following paragraph such a human would answer the question "Where is Mary?" differently from how it should currently be answered. The possible answer options, including the current answer, are a) "hallway", and b) "office". Make as few edits as possible. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's an edited version of the paragraph that would elicit a different response to the question "Where is Mary?":

Mary returned to the hallway. Then, John went back to the bathroom.

In this revised version, a human would likely respond to the question "Where is Mary?" with option (a) "hallway", whereas before they would have said option (b) "office".

### Session 3: Consistency check

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary returned to the hallway. Then, John went back to the bathroom.

In this revised version, a human would likely respond to the question "Where is Mary?" with option (a) "hallway", whereas before they would have said option (b) "office".

Sure! A human would answer the question "Where is Mary?" with option (a) "hallway".

Figure D.23 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: implicit”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

Edit the following paragraph such you would answer the question "Where is Mary?" differently from how it should currently be answered. The possible answer options, including the current answer, are a) "hallway", and b) "office". Make as few edits as possible. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Here's an edited version of the paragraph that would allow for a different answer to the question "Where is Mary?":

Mary went back to the hallway. Then, John went back to the bathroom.

### Session 3: Consistency check

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the hallway. Then, John went back to the bathroom.

b) b) Office

Figure D.24 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Counterfactual target: implicit". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

## Feature attribution

### Session 1: Classification

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

List the most important words for answering "Where is Mary?" given the following paragraph, such that without these words the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here are the most important words for answering "Where is Mary?" based on the given paragraph:

1. Back
2. Office

### Session 3: Consistency check

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went [REDACTED] to the [REDACTED]. Then, John went [REDACTED] to the bathroom.

Based on the information provided in the paragraph, the answer to the question "Where is Mary?" would be (c) unknown.

Figure D.25 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "redacted"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) The human would answer: b) Office

### Session 2: Explanation

List the most important words for answering "Where is Mary?" given the following paragraph, such that without these words a human can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here are the most important words for answering "Where is Mary?" based on the given paragraph:

1. Back
2. Office

### Session 3: Consistency check

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went [REDACTED] to the [REDACTED]. Then, John went [REDACTED] to the bathroom.

Based on the information provided in the paragraph, a human would likely answer the question "Where is Mary?" as (c) "unknown." The paragraph does not provide enough information to determine Mary's current location, as it only mentions that she went [REDACTED] and doesn't specify where that is.

Figure D.26 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Redaction instruction: "redacted"". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

List the most important words for answering "Where is Mary?" given the following paragraph, such that without these words you can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here are the most important words for answering "Where is Mary?" based on the given paragraph:

1. Back
2. Office

### Session 3: Consistency check

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went [REDACTED] to the [REDACTED]. Then, John went [REDACTED] to the bathroom.

Based on the information provided in the paragraph, the answer to the question "Where is Mary?" would be (c) unknown.

Figure D.27 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: you, Redaction instruction: "redacted"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

List the most important words for answering "Where is Mary?" given the following paragraph, such that without these words the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here are the most important words for answering "Where is Mary?" based on the given paragraph:

1. Back
2. Office

### Session 3: Consistency check

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went [REMOVED] to the [REMOVED]. Then, John went [REMOVED] to the bathroom.

Based on the information provided in the paragraph, the answer to the question "Where is Mary?" would be (c) unknown.

Figure D.28 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) The human would answer: b) Office

### Session 2: Explanation

List the most important words for answering "Where is Mary?" given the following paragraph, such that without these words a human can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here are the most important words for answering "Where is Mary?" based on the given paragraph:

1. Back
2. Office

### Session 3: Consistency check

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went [REMOVED] to the [REMOVED]. Then, John went [REMOVED] to the bathroom.

Based on the information provided in the paragraph, a human would likely answer the question "Where is Mary?" as (c) "unknown." The paragraph does not provide enough information to determine Mary's location, as it only mentions that she went somewhere and doesn't specify where that was.

Figure D.29 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Redaction instruction: "removed"". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

List the most important words for answering "Where is Mary?" given the following paragraph, such that without these words you can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here are the most important words for answering "Where is Mary?" based on the given paragraph:

1. Back
2. Office

### Session 3: Consistency check

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went [REMOVED] to the [REMOVED]. Then, John went [REMOVED] to the bathroom.

Based on the information provided in the paragraph, the best answer to the question "Where is Mary?" would be (c) "unknown." The paragraph does not provide enough information to determine Mary's location.

Figure D.30 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Redaction instruction: "removed"". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

## Redacted

### Session 1: Classification

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

Redact the most important words for answering "Where is Mary?" given the following paragraph, by replacing important words with [REDACTED], such that without these words the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's a redacted version of the paragraph:

Mary went back to the [REDACTED]. Then, John went back to the [REDACTED].

Without the redacted words, it would be difficult to determine where Mary is.

### Session 3: Consistency check

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the [REDACTED]. Then, John went back to the [REDACTED].

Without the redacted words, it would be difficult to determine where Mary is.

b) b) Office

Figure D.31 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: objective, Redaction instruction: "redacted"". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

---

### Session 1: Classification

---

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) The human would answer: b) Office

### Session 2: Explanation

---

Redact the most important words for answering "Where is Mary?" given the following paragraph, by replacing important words with [REDACTED], such that without these words a human can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's the redacted paragraph:

Mary went back to the [REDACTED]. Then, John went back to the [REDACTED].

Without the redacted words, it would be difficult or impossible for a human to accurately answer the question "Where is Mary?"

### Session 3: Consistency check

---

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the [REDACTED]. Then, John went back to the [REDACTED].

Without the redacted words, it would be difficult or impossible for a human to accurately answer the question "Where is Mary?"

b) b) Office

Figure D.32 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Redaction instruction: "redacted"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

Redact the most important words for answering "Where is Mary?" given the following paragraph, by replacing important words with [REDACTED], such that without these words you can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's a redacted version of the paragraph:

Mary went back to the [REDACTED]. Then, John went back to the [REDACTED].

Without the redacted words, it would be difficult to determine where Mary is.

### Session 3: Consistency check

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain redacted words marked with [REDACTED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the [REDACTED]. Then, John went back to the [REDACTED].

Without the redacted words, it would be difficult to determine where Mary is.

b) b) Office

Figure D.33 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Redaction instruction: "redacted"". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

---

### Session 1: Classification

---

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

---

Redact the most important words for answering "Where is Mary?" given the following paragraph, by replacing important words with [REMOVED], such that without these words the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's a redacted version of the paragraph with the most important words removed:

Mary went back to the [REMOVED]. Then, John went back to the [REMOVED].

### Session 3: Consistency check

---

Consider the following paragraph, and answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the [REMOVED]. Then, John went back to the [REMOVED].

b) b) Office

Figure D.34 Redaction explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Redaction instruction: "removed"”. The true label is “office”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

---

Session 1: Classification

---

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) The human would answer: b) Office

---

Session 2: Explanation

---

Redact the most important words for answering "Where is Mary?" given the following paragraph, by replacing important words with [REMOVED], such that without these words a human can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's the redacted paragraph:

Mary went back to the [REMOVED]. Then, John went back to the [REMOVED].

---

Session 3: Consistency check

---

Considering the following paragraph, how would a human answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the [REMOVED]. Then, John went back to the [REMOVED].

b) b) Office

Figure D.35 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Redaction instruction: "removed"". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

### Session 1: Classification

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

b) b) Office

### Session 2: Explanation

Redact the most important words for answering "Where is Mary?" given the following paragraph, by replacing important words with [REMOVED], such that without these words you can not answer the question. Do not explain the answer.

Paragraph: Mary went back to the office. Then, John went back to the bathroom.

Sure! Here's a redacted version of the paragraph:

Mary went back to the [REMOVED]. Then, John went back to the [REMOVED].

Without the redacted words, it would be difficult to determine where Mary is.

### Session 3: Consistency check

Considering the following paragraph, how would you answer the question: "Where is Mary?" The paragraph can contain removed words marked with [REMOVED]. Answer either a) "hallway", b) "office", or c) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Mary went back to the [REMOVED]. Then, John went back to the [REMOVED].

Without the redacted words, it would be difficult to determine where Mary is.

b) b) Office

Figure D.36 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Redaction instruction: "removed"". The true label is "office". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

### D.2.3 NLI classification with RTE

#### Counterfactual

Session 1: Classification

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

Session 2: Explanation

Edit the following paragraph such that the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." does not entails from it. Make as few edits as possible. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here is a revised version of the paragraph that does not entail the statement "Microsoft's Hotmail has raised its storage capacity to 250MB":

Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced a change in its storage policy, following similar moves by rivals such as Google, Yahoo, and Lycos.

Session 3: Consistency check

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced a change in its storage policy, following similar moves by rivals such as Google, Yahoo, and Lycos.

3) Unknown

Figure D.37 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: objective, Counterfactual target: explicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “not faithful”.

### Session 1: Classification

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Edit the following paragraph such that a human would say the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." does not entail from it. Make as few edits as possible. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here is a revised version of the paragraph that does not entail the statement "Microsoft's Hotmail has raised its storage capacity to 250MB":

Only a week after declining to comment on rumors of an increase in storage capacity for its Hotmail email service, Microsoft made an announcement early Thursday regarding changes to their allowance. The company revealed that they would be following in the footsteps of competitors like Google, Yahoo, and Lycos by increasing their storage capacity.

### Session 3: Consistency check

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after declining to comment on rumors of an increase in storage capacity for its Hotmail email service, Microsoft made an announcement early Thursday regarding changes to their allowance. The company revealed that they would be following in the footsteps of competitors like Google, Yahoo, and Lycos by increasing their storage capacity.

2) No

Figure D.38 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration “Persona instruction: human, Counterfactual target: explicit”. The true label is “yes”. The initial prediction was “correct”. The interpretability-faithfulness was evaluated to be “faithful”.

### Session 1: Classification

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Edit the following paragraph such that you would say the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." does not entails from it. Make as few edits as possible. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here is a revised version of the paragraph that does not entail the statement "Microsoft's Hotmail has raised its storage capacity to 250MB":

Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced a change in its storage policy for the service, following similar moves by rivals such as Google, Yahoo, and Lycos.

### Session 3: Consistency check

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced a change in its storage policy for the service, following similar moves by rivals such as Google, Yahoo, and Lycos.

3) Unknown

Figure D.39 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Counterfactual target: explicit". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

### Session 1: Classification

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Edit the following paragraph, such that given the statement "Microsoft's Hotmail has raised its storage capacity to 250MB.", the entailment becomes the opposite of what it currently is. Make as few edits as possible. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's an edited version of the paragraph with the entailment you requested:

Only a week after it had no comment on lowering the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was reducing the allowance to 250MB, diverging from the trend set by rivals such as Google, Yahoo, and Lycos.

### Session 3: Consistency check

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on lowering the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was reducing the allowance to 250MB, diverging from the trend set by rivals such as Google, Yahoo, and Lycos.

3) Unknown. The statement mentions that Microsoft reduced the storage capacity of its Hotmail email service to 250MB, but does not mention anything about raising it to that amount. Therefore, we cannot determine whether the statement entails the given sentence.

Figure D.40 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: objective, Counterfactual target: implicit". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

### Session 1: Classification

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Edit the following paragraph, such that given the statement "Microsoft's Hotmail has raised its storage capacity to 250MB.", a human would say the entailment is the opposite of what it currently is. Make as few edits as possible. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's an edited version of the paragraph with the desired entailment:

Only a week after it had no comment on lowering the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was reducing the allowance to 250MB, following suit with its competitors Google, Yahoo, and Lycos who have also downsized their email storage capacities.

### Session 3: Consistency check

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on lowering the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was reducing the allowance to 250MB, following suit with its competitors Google, Yahoo, and Lycos who have also downsized their email storage capacities.

3) Unknown. The statement mentions that Microsoft reduced the storage capacity of its Hotmail email service to 250MB, but does not mention anything about raising it to that amount. Therefore, we cannot determine whether the statement entails the given sentence.

Figure D.41 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Counterfactual target: implicit". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

### Session 1: Classification

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Edit the following paragraph, such that given the statement "Microsoft's Hotmail has raised its storage capacity to 250MB.", you would say the entailment is the opposite of what it currently is. Make as few edits as possible. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's an edited version of the paragraph with the entailment flipped:

Only a week after it had no comment on lowering the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was reducing the allowance to 250MB in response to pressure from competitors like Google, Yahoo, and Lycos.

### Session 3: Consistency check

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on lowering the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was reducing the allowance to 250MB in response to pressure from competitors like Google, Yahoo, and Lycos.

3) Unknown

Figure D.42 Counterfactual explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Counterfactual target: implicit". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "not faithful".

## Feature attribution

### Session 1: Classification

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

List the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, such that without these words it can not be determined if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here are the most important words in the paragraph for determining entailment:

1. Microsoft
2. Hotmail
3. storage capacity
4. 250MB
5. rival(s)
6. Google
7. Yahoo
8. Lycos

### Session 3: Consistency check

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the [REDACTED] of its [REDACTED] e-mail service, [REDACTED] early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

3) Unknown.

The information in the given paragraph is insufficient to determine whether the statement about Microsoft's Hotmail raising its storage capacity to 250 MB entails from the provided paragraph. The paragraph mentions a company increasing the allowance of its email service but does not specify the name of the company or the exact amount of the increase, making it impossible to determine if it refers to the same event mentioned in the statement. Therefore, the answer is unknown.

Figure D.43 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: objective, Redaction instruction: "redacted"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

List the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, such that without these words you a human not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here are the most important words in the paragraph for determining entailment:

1. Microsoft
2. Hotmail
3. storage capacity
4. 250MB
5. rival(s)
6. Google
7. Yahoo
8. Lycos

### Session 3: Consistency check

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the [REDACTED] of its [REDACTED] e-mail service, [REDACTED] early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

3) Unknown.

The paragraph does not mention Microsoft or Hotmail, so it is unclear whether the statement about raising storage capacity refers to those entities or another email service. Additionally, without knowing what the redacted words represent, it is difficult to determine the context and meaning of the sentence. Therefore, it cannot be determined whether the given statement entails from the provided paragraph.

Figure D.44 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Redaction instruction: "redacted"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

List the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, such that without these words you can not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here are the most important words in the paragraph for determining entailment:

1. Microsoft
2. Hotmail
3. storage capacity
4. 250MB
5. rival(s)
6. Google
7. Yahoo
8. Lycos

### Session 3: Consistency check

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the [REDACTED] of its [REDACTED] e-mail service, [REDACTED] early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

3) Unknown.

The information in the provided paragraph is insufficient to determine whether the statement "Microsoft's Hotmail has raised its storage capacity to 250MB" entails from it because the name "Hotmail" is not mentioned in the paragraph.

Figure D.45 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Redaction instruction: "redacted"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

List the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, such that without these words it can not be determined if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here are the most important words in the paragraph for determining entailment:

1. Microsoft
2. Hotmail
3. storage capacity
4. 250MB
5. rival(s)
6. Google
7. Yahoo
8. Lycos

### Session 3: Consistency check

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the [REMOVED] of its [REMOVED] e-mail service, [REMOVED] early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

3) Unknown.

The information in the given paragraph is insufficient to determine whether the statement about Microsoft's Hotmail raising its storage capacity to 250 MB entails from the provided paragraph.

Figure D.46 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: objective, Redaction instruction: 'removed'". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

List the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, such that without these words you a human not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here are the most important words in the paragraph for determining entailment:

1. Microsoft
2. Hotmail
3. storage capacity
4. 250MB
5. rival(s)
6. Google
7. Yahoo
8. Lycos

### Session 3: Consistency check

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the [REMOVED] of its [REMOVED] e-mail service, [REMOVED] early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

3) Unknown.

The paragraph does not mention Microsoft's Hotmail nor the specific amount of 250 MB, so it cannot be determined whether the statement entails the given information without additional context.

Figure D.47 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Redaction instruction: "removed"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

List the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, such that without these words you can not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here are the most important words in the paragraph for determining entailment:

1. Microsoft
2. Hotmail
3. storage capacity
4. 250MB
5. rival(s)
6. Google
7. Yahoo
8. Lycos

### Session 3: Consistency check

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the [REMOVED] of its [REMOVED] e-mail service, [REMOVED] early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

3) Unknown.

The paragraph does not mention anything about Microsoft's Hotmail specifically, so it cannot be determined whether the statement about Hotmail's storage capacity entails from the given information.

Figure D.48 Feature attribution explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Redaction instruction: "removed"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

## Redacted

### Session 1: Classification

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Redact the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, by replacing important words with [REDACTED], such that without these words it can not be determined if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's the redacted paragraph:

Only a week after it had no comment on upping the storage capacity of its [REDACTED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

### Session 3: Consistency check

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its [REDACTED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

3) Unknown. The paragraph does not explicitly state that Hotmail is the email service being referred to, nor does it mention the specific amount of increased storage capacity, so it cannot be determined whether the statement "Microsoft's Hotmail has raised its storage capacity to 250MB" logically follows from the given information.

Figure D.49 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: objective, Redaction instruction: "redacted"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Redact the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, by replacing important words with [REDACTED], such that without these words a human can not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's the redacted paragraph:

Only a week after it had no comment on upping the storage capacity of its [REDACTED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

### Session 3: Consistency check

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its [REDACTED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

3) Unknown.

The paragraph does not explicitly state that Hotmail, specifically, is having its storage capacity increased, so we cannot say with certainty that the statement about Hotmail's storage capacity entails from the given information.

Figure D.50 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Redaction instruction: "redacted"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Redact the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, by replacing important words with [REDACTED], such that without these words you can not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's the redacted paragraph:

Only a week after it had no comment on upping the storage capacity of its [REDACTED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

### Session 3: Consistency check

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain redacted words marked with [REDACTED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its [REDACTED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REDACTED] to follow similar moves by rivals such as [REDACTED], [REDACTED], and [REDACTED].

3) Unknown.

The paragraph does not explicitly state that Hotmail has raised its storage capacity to 250 MB, nor does it provide enough information to deduce this fact. Therefore, we cannot say whether the given statement entails from the provided paragraph without additional context.

Figure D.51 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Redaction instruction: "redacted"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Redact the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, by replacing important words with [REMOVED], such that without these words it can not be determined if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's the redacted paragraph:

Only a week after it had no comment on upping the storage capacity of its [REMOVED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

### Session 3: Consistency check

Does the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its [REMOVED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

3) Unknown. The paragraph does not mention anything about Hotmail specifically, so we cannot infer that Hotmail's storage capacity has been increased to 250 MB based on this information alone. Additionally, there are several blank spaces in the paragraph where important information might have been included, making it difficult to determine whether the statement entails from the given text. Therefore, without additional context or complete information, we cannot confidently say yes or no.

Figure D.52 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: objective, Redaction instruction: "removed"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Redact the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, by replacing important words with [REMOVED], such that without these words a human can not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's the redacted paragraph:

Only a week after it had no comment on upping the storage capacity of its [REMOVED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

### Session 3: Consistency check

Does a human think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its [REMOVED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

3) Unknown. The paragraph does not mention anything about Hotmail specifically, so it is unclear whether the increased storage capacity applies to that particular email service or another one offered by Microsoft. Additionally, the paragraph mentions rival companies but doesn't specify which ones they are, making it difficult to determine what their role in this situation might be.

Figure D.53 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: human, Redaction instruction: "removed"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".

### Session 1: Classification

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

1) Yes

### Session 2: Explanation

Redact the most important words in the following paragraph, for determining if the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entails from it, by replacing important words with [REMOVED], such that without these words you can not determine if there is entailment or no entailment. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its Hotmail e-mail service, Microsoft early Thursday announced it was boosting the allowance to 250MB to follow similar moves by rivals such as Google, Yahoo, and Lycos.

Sure! Here's the redacted paragraph:

Only a week after it had no comment on upping the storage capacity of its [REMOVED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

### Session 3: Consistency check

Do you think the statement "Microsoft's Hotmail has raised its storage capacity to 250MB." entail from the following paragraph? The paragraph can contain removed words marked with [REMOVED]. Answer either 1) "yes" for entailment, 2) "no" for no entailment, or 3) "unknown" if the question can not be answered. Do not explain the answer.

Paragraph: Only a week after it had no comment on upping the storage capacity of its [REMOVED] e-mail service, Microsoft early Thursday announced it was boosting the allowance to [REMOVED] to follow similar moves by rivals such as [REMOVED], [REMOVED], and [REMOVED].

3) Unknown. The paragraph does not mention anything about Hotmail specifically, so we cannot infer that Hotmail's storage capacity has been raised to 250 MB based on this information alone. Additionally, there are several blank spaces in the paragraph where important information might have been included, making it difficult to determine whether the statement entails from the given text. Therefore, without additional context or information, it is impossible to say whether the statement entails from the provided paragraph.

Figure D.54 Redaction explanation and interpretability-faithfulness evaluation, with the configuration "Persona instruction: you, Redaction instruction: "removed"". The true label is "yes". The initial prediction was "correct". The interpretability-faithfulness was evaluated to be "faithful".