



Titre: Analyse automatique des données scripturales prétraitées par des outils de visualisation

Auteurs: F. Neveu, Hélène-Sarah Bécotte-Boutin, Gilles Caporossi, Alain Hertz, Christophe Leblay, G. Bergounioux, M. H. Côté, J. M. Fournier, L. Hriba, & S. Prévost

Date: 2016

Type: Article de revue / Article


Référence: Neveu, F., Bécotte-Boutin, H.-S., Caporossi, G., Hertz, A., Leblay, C., Bergounioux, G., Côté, M. H., Fournier, J. M., Hriba, L., & Prévost, S. (2016). Analyse automatique des données scripturales prétraitées par des outils de visualisation. SHS Web of Conferences, 27, 06001 (18 pages).
Citation: <https://doi.org/10.1051/shsconf/20162706001>

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

URL de PolyPublie: <https://publications.polymtl.ca/5054/>
PolyPublie URL:

Version: Version officielle de l'éditeur / Published version
Révisé par les pairs / Refereed

Conditions d'utilisation: CC BY
Terms of Use:

 **Document publié chez l'éditeur officiel**
Document issued by the official publisher

Titre de la revue: SHS Web of Conferences (vol. 27)
Journal Title:

Maison d'édition: EDP Sciences
Publisher:

URL officiel: <https://doi.org/10.1051/shsconf/20162706001>
Official URL:

Mention légale: © The Authors, published by EDP Sciences. This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0 (<http://creativecommons.org/licenses/by/4.0/>).
Legal notice:

! "#\$%&' (#) *+, #- .) ' (/ ' & (/ + " " 0' & (&123*) 2#\$ & 320*2#-*0' &(3#2/ ' &(+) *-\$&/ ' (4-&) #5#*+ " ((

HŽl•ne-Sarah BŹcotte-Boutin^{a,1,2,3}, Gilles Caporossi^{2,3,4}, Alain Hertz^{1,3} et Christophe Leblay^{2,5}

¹ Polytechnique MontrŹal, 2900, boul. f•douard-Montpetit Campus de l'UniversitŹ de MontrŹal 2500, chemin de Polytechnique, MontrŹal (QuŹbec) H3T 1J4, Canada

² Institut des textes et manuscrits modernes (ITEM), CNRS : fcole normale supŹrieure [ENS] 45, rue d'Ulm 75005 Paris, France

³ GERAD, HEC MontrŹal 3000, chemin de la C^mte-Sainte-Catherine, MontrŹal (QuŹbec) H3T 2A7, Canada

⁴ HEC MontrŹal, 3000 chemin de la C^mte-Sainte-Catherine MontrŹal (QuŹbec) Canada H3T 2A7, Canada

⁵ University of Turku, 20014 Turun yliopisto, Finlande

60&), 0. Plusieurs mŹthodes pour analyser le processus d'Źcriture ont ŹtŹ utilisŹs afin de comprendre les stratŹgies des scripteurs. L'outil principal pour analyser le processus d'Źcriture est le fichier log, qui contient de fa•on exhaustive et dŹtaille l'ensemble des opŹrations effectuŹes par le scripteur lors de la rŹdaction d'un texte. Les donnŹes qui y sont emmagasinŹes sont de quantitŹ considŹrable et lorsqu'elles ne sont pas prŹalablement traitŹes, elles sont hostiles Źtre analysŹes par l'humain. Parmi les outils d'analyse utilisŹs, les reprŹsentations du processus d'Źcriture permettent l'agrŹgation des donnŹes grŹce Ź un prŹ-traitement. Les structures sous-jacentes des donnŹes ainsi reprŹsentŹes sont gŹnŹralement plus propices Ź l'analyse que les donnŹes brutes. Cet article vise Ź dŹmontrer diffŹrentes mŹthodes d'analyse automatique pouvant Źtre appliquŹes Ź ces structures afin de trouver ou confirmer des structures et tendances Ź travers les donnŹes.

! 7&*2#1*8 Several methods to analyze the writing process were used in order to understand the strategies of the writers. The main tool to analyze the writing process is the log file which contains all the operations performed by the writer when writing a text, in a comprehensive and detailed way. The data stored in it is of considerable amount and when not previously treated, it is not made to be analyzed by humans. Among the analytical tools used, the representations of the writing process allow aggregation of data through a pre-treatment. The underlying data structures as shown by these tools are generally conducive to analyzing the raw data afterwards. This article aims to demonstrate various automatic analysis methods that can be applied to these structures to find or confirm the structures and trends through data.

* Auteur de correspondance : helene.becotte@polymtl.ca

9(: " *2+ /) 1* -+ " (

Les techniques d'analyse des textes produits foisonnent et de nouvelles disciplines sont entièrement créées autour de ce concept. Les techniques utilisées sont variées et proviennent de plusieurs domaines [1]. Les statistiques, les tests A/B, la fusion et intégration de données, la régression, la classification, les règles d'association, le data mining, la visualisation, le text mining, l'analyse de sentiments, les réseaux de neurones, l'analyse de réseaux, le traitement du langage naturel, la simulation, l'analyse de séries chronologiques sont toutes des techniques ou groupes de techniques qui se croisent [2].

Bien que ces techniques soient utilisées sur des textes produits et des données statiques et que le processus d'écriture soit pour sa part composé de données dynamiques multidimensionnelles, il est possible de s'en inspirer pour bénéficier des possibilités de telles techniques. C'est ce que nous nous proposons de faire dans ce travail, en nous concentrant sur la visualisation.

L'étude du processus d'écriture est complexe pour plusieurs raisons. La toute première raison est due à la singularité du processus individuel qui fait en sorte qu'il n'existe aucune méthode unique pour rédiger un texte de qualité [3]. La seconde est qu'il est composé principalement de données qui varient dans le temps. Ainsi, les opérations d'écriture, telles qu'elles ont été décrites par la grammaire textuelle dans un environnement numérique [4], peuvent être listées en fonction de leur ordre chronologique d'apparition. Ces opérations et les pauses qui y sont attachées, sont difficiles à interpréter par l'humain lorsqu'elles ne sont pas combinées et contextualisées pour former des unités identifiables par l'homme : des mots, des phrases ou une section de ce qui deviendra du texte.

Parmi la panoplie disponible des outils d'analyse disponibles, les représentations du processus scriptural sont les plus utilisées grâce à leur facilité à traiter l'information.

Avec l'avènement de technologies récentes à l'analyse du processus depuis les années 90, les possibilités d'automatisation se sont multipliées. Malgré cela, peu de méthodes permettent aujourd'hui aux chercheurs d'accélérer leur analyse et de se concentrer sur l'interprétation et la compréhension des mécanismes de révisions qui pourraient permettre à un scripteur d'améliorer la qualité de ses textes.

Dans ce travail, nous nous proposons de nous intéresser, tout d'abord, aux structures sous-jacentes aux données impliquées dans le processus scriptural en distinguant nettement entre les représentations réalisées dans le cadre de la linguistique cognitive et celles réalisées dans le cadre de la linguistique grammaticale. Ce qui nous permettra d'introduire à la méthodologie suivie et d'exposer les modalités d'enregistrement. Nous finirons par illustrer notre propos à l'aide d'un exemple tiré de notre corpus afin de préciser les objectifs poursuivis.

; (< 20*2#-** , ' " *(= (&*2) 1*) 2 &/ ' &/ + " " 0' &(&12-3*) 2#\$ &

Depuis les années 80, plusieurs méthodes pour analyser le processus d'écriture ont été utilisées [5]. L'outil principal pour analyser le processus d'écriture est le fichier d'enregistrement, appelé log, qui contient de façon exhaustive et détaillée l'ensemble des opérations effectuées par le scripteur lors de la rédaction d'un texte [6]. Les données qui y sont emmagasinées sont considérables, et lorsqu'elles ne sont pas préalablement traitées, elles sont hostiles à l'analyse humaine [7]. Ce traitement préalable nous semble donc déterminant dans l'accès aux données scripturales enregistrées. Autrement dit, il est impossible d'exploiter les programmes d'enregistrement de l'écriture (dits aussi de temps réel), sans procéder, au préalable, à un redéploiement visuel des données obtenues. Les structures sous-jacentes des données ainsi représentées sont généralement plus propices à l'analyse que les données brutes inexploitées. La principale raison est que lorsqu'on cherche à explorer des données, plusieurs étapes sont nécessaires avant de pouvoir utiliser les techniques d'analyse à proprement parler. Le processus

simplifié par lequel les données sont acquises, traitées et analysées se résume ainsi : 1) enregistrement, 2) partage, 3) tri, 4) recherche, 5) représentation visuelle des données et 6) analyse [2]. Plusieurs de ces actions ont déjà été étudiées de manière individuelle et ont fait l'objet de créations de logiciels visant à enregistrer le processus d'écriture et à le représenter dans l'une ou l'autre de ses dimensions [8]. Ces logiciels sont tous différents et généralement destinés à un projet de recherche particulier [6]. Certains sont conçus pour traiter l'information et la retransmettent ensuite à l'utilisateur de façon simplifiée, sous la forme d'une représentation visuelle.

L'utilisation de représentations graphiques de données, appelées de façon plus générale visualisations, consiste à explorer et essayer de comprendre les grands ensembles de données [9, p. xvi]. Elles permettent notamment d'identifier des tendances, structures, irrégularités et relations entre les données sur une certaine période temporelle [10, p. 110]. Le format compact de la visualisation agrège les données et utilise les capacités cognitives de l'humain [11]. Le but principal de l'utilisation de ces images représentant le processus scriptural est de laisser l'utilisateur trouver des structures sous-jacentes parmi les données [12].

Il est important cependant de considérer que l'utilisation de représentations visuelles des données comporte certains défauts. Puisque les conclusions d'analyse seront effectuées par l'utilisateur, les capacités de traitement de celui-ci influencent fortement les résultats obtenus [13]. Tout d'abord, l'être humain a une capacité limitée à traiter des quantités d'informations et il aurait aussi parfois tendance à percevoir des tendances et associations alors qu'il n'y en a pas [1]. Pour cette raison et pour maximiser le potentiel de conclusions possibles, il est pertinent de pouvoir confirmer mathématiquement ces conclusions grâce à des modèles d'analyse de données et de s'assurer qu'il existe effectivement une corrélation forte entre différentes variables [14].

Les techniques d'analyse utilisées sont ensuite dépendantes du type de pré-traitement effectué sur les données. Structurer les données d'une certaine manière plutôt qu'une autre influence donc nécessairement les résultats d'analyse obtenus par la suite. Conséquemment, lors du choix d'une technique particulière d'analyse, ce n'est pas la technique elle-même qui est importante mais l'utilisation finale des résultats obtenus qu'il faut considérer [15, p. 134]. Dans le but d'obtenir par exemple une compréhension du système sous-jacent aux données comme dans le cas de l'analyse du temps de l'écriture, il est important de s'assurer que le modèle utilisé et les résultats obtenus soient propices à l'interprétation [16, p. 209].

Il existe deux grands modèles de représentations qui pré-traitent les données de manière à pouvoir appliquer des techniques d'analyse linguistique de données. Le tout premier, attaché aux travaux réalisés par une linguistique cognitive, est celui qui a été développé dans le cadre des Systèmes d'Information Graphique (dorsnavant SIG). Le second, en réaction au premier, a été développé en associant, de manière pluridisciplinaire, la linguistique générale et la théorie mathématique des graphes. A ces deux modèles, nous proposerons, dans ce travail, d'introduire un tout nouveau modèle, qui est en fait le mixte des deux grands modèles : la représentation progressive.

;
E:F G

Les SIG sont un type de représentation utilisé initialement pour traiter des données géographiques. Les principaux attributs sont l'identification spatiale et temporelle d'activités, qui ont une ressemblance avec les données issues de l'écriture [17]. Ces attributs sont utilisés en tant que coordonnées. En géographie, ces coordonnées sont souvent la latitude et la longitude [18, p. 1] mais peuvent également contenir l'élévation [18, p. 73]. En analyse du processus,

dans la majorité des cas, ces coordonnées vont être composées d'une mesure de temps et d'une quantité de révisions effectuées [17].

Voici un exemple d'une représentation effectuée sur un logiciel dédié pour les SIG [19] :

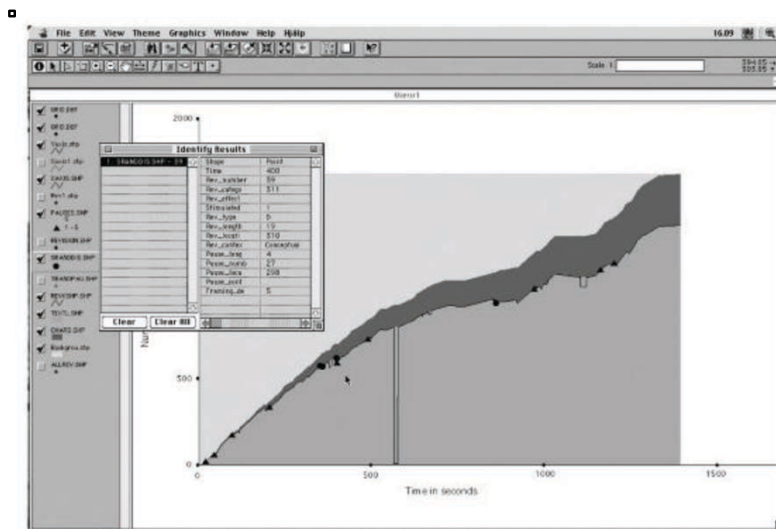


Fig. 1. Représentation SIG.

Notre propos n'étant pas de rentrer dans l'analyse de ces représentations, nous nous contenterons, dans le cadre de ce travail, d'en noter les principales caractéristiques. Ainsi, comme il est possible de le voir dans la figure 1, la contribution clé des SIG réside dans la structure qu'ils donnent aux données et aux possibilités de les analyser considérant l'emplacement physique des données dans le texte. Cette dimension du processus d'écriture se nomme spatialité. Des méthodes d'analyse spécifiques sont adaptées à ces structures de données de manière à mettre en relief certaines problématiques. Les principales techniques comprennent (1) l'analyse des chevauchements entre différentes caractéristiques, (2) l'identification de sous-groupes, (3) l'exploration de tendances spatiales selon certaines variables et (4) l'exploration de la relation entre plusieurs variables en fonction de leur emplacement spatial [20, p. 2].

Ont utilisé ce mode de représentation (SIG) principalement les études suivantes : LS Graph [21, 17], Gen se du texte [4], GIS Graph [19], Timeline [22] et Inputlog [23]. Toutes ces représentations ont en commun l'utilisation des axes X et Y, comme la figure 1 illustre, chacune exploitant des caractéristiques différentes : chronologie, topologie, avec un jeu de couleurs.

Les méthodes les plus pertinentes dans le cas de l'analyse du processus de l'écriture sont sans contredit celles qui permettent de mettre de l'avant l'aspect spatial des données. La visualisation est considérée comme étant la première étape pour n'importe quelle analyse spatiale [20, p. 19 ch.2]. Par la suite, la requête de données est la méthode la plus simple pour débiter l'analyse. Il s'agit d'une façon d'aggr ger les données à un niveau supérieur pour obtenir des informations concentrées sur une problématique particulière. Les statistiques sont l'étape suivante pour comprendre les caractéristiques des données.

Finalement, plusieurs sources d'erreur sont possibles et un choix de modélisation de données inapproprié est aussi considéré critique dans le cas des SIG [20, p. 18 ch.2].

; 8 (>-" ?) -&*-.) ' (?0" 0*-.) ' (*2 320& " *#*+" (3#2?2#3C' &(

La représentation par les graphes permet de mettre davantage en relief l'aspect dynamique de l'écriture. Celle-ci est orientée sur la chronologie du processus d'écriture [24]. L'une des particularités est de bien gérer la transformation et les mouvements du texte grâce aux nœuds représentant des parties de texte (ajouts, suppressions) reliés ensemble par des arcs définissant leur relation, soit chronologique ou spatiale [25]. Il est possible de voir le contenu des nœuds. Comme mentionné par Caporossi et Leblay [8], cette représentation montre néanmoins l'aspect temporel de la rédaction, soit le moment exact où le scripteur a effectué chacune des opérations représentées.

Voici l'exemple d'une représentation par graphes produite par Caporossi et Leblay [8] :

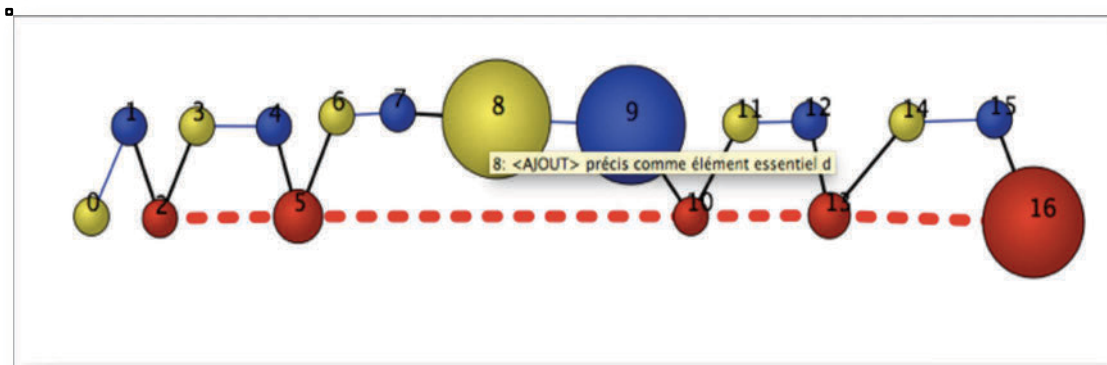


Fig. 2. Représentation graphique par graphes.

Cette représentation qui reprend une description graphique à partir des opérations élémentaires d'écriture (Ajout versus Suppression) doit être lue de la manière suivante:

- (0): Il me
- (1): <SUPPR> em Il
- (2): <AJOUT> <CR><CR>Je ne vo
- (3): <AJOUT> lis
- (4): <SUPPR> sil
- (5): <AJOUT> is pas un lieu
- (6): <AJOUT> q
- (7): <SUPPR> q
- (8): <AJOUT> précis comme élément essentiel d
- (9): <SUPPR> d leitesse tnement emmoc sicrp
- (10): <AJOUT> ou
- (11): <AJOUT> en end
- (12): <SUPPR> dne ne
- (13): <AJOUT> un endroit pr
- (14): <AJOUT> +ec

(15): <SUPPR> ce+

(16): <AJOUT> Žcis comme ŽŽment essentiel d

Ce mode de reprŽsentation n'oublie pas l'importance du texte produit : il est possible de le suivre au moyen des traits pointillŽs rouges. Il est facile de noter tout le travail d'Žcriture rŽalisŽ par le scripteur, soit les 2/3 de ce qui est Žcrit (11 n'uds sur 16) n'apparaıtra pas dans le texte produit.

Mais, au delˆ de la simple apprŽhension de cette reprŽsentation mathŽmatique, et afin de bien en comprendre les subtilitŽs, nous proposons un succinct rappel des ressources thŽoriques des graphes.

Un graphe est une reprŽsentation visuelle simple constituŽe de points reliŽs ensemble par des lignes [26, p. 1]. Leur structure est utilisŽe dans des domaines variŽs puisqu'elle reprŽsente simplement des relations entre objets [27, p. 1]. La terminologie est plutˆt intuitive. Les points sont appelŽs n'uds tandis que les lignes sont appelŽes arˆtes. Les n'uds peuvent reprŽsenter des villes sur une carte, des individus dans un rŽseau social ou encore simplement des informations reliŽes entre elles par les arˆtes [28, p. 2]. Deux n'uds sont dits adjacents ou voisins s'ils sont reliŽs par une arˆte. Finalement, le degrŽ d'un n'ud est le nombre de voisins que possŽde celui-ci [29, p. 4]. De maniŽre plus formelle, un graphe $G = (V, E)$ est constituŽ d'un ensemble fini V de n'uds et d'un ensemble E d'arˆtes. Un graphe peut ˆtre orientŽ, dans lequel cas $G = (V, A)$ a un ensemble A d'arcs. Un arc qui relie un sommet u ˆ un sommet v sera notŽ $[u, v]$ tandis qu'un arc qui relie un sommet u ˆ un sommet v sera notŽ $[u, v]$.

De maniŽre appliquŽe, en prenant l'exemple des rŽseaux sociaux, un graphe non orientŽ serait utilisŽ pour dŽcrire une plateforme telle que Facebook, pour dŽcrire la relation mutuelle que deux individus partagent en Žtant amis. Un graphe orientŽ plutˆt serait utile dans un cas comme Twitter oˆ la relation n'est pas nŽcessairement rŽciproque. Un usager peut trˆs bien suivre un autre usager, sans que l'inverse soit vrai. S'il y a une relation de prŽcŽdence entre des actions, l'orientation des arcs peut ˆtre pertinente [2, p. 29].

L'analyse de rŽseaux, que certains auteurs appellent plutˆt graph mining [30], est relative ˆ la thŽorie des graphes et est une technique d'analyse de donnŽes utilisŽe en big data [2, p. 29]. Cette technique permet en effet d'utiliser de bons algorithmes pour l'analyse et le traitement des donnŽes complexes pouvant ˆtre reprŽsentŽes sous forme de graphes [31]. Depuis une dizaine d'annŽes, le dŽveloppement d'algorithmes de data mining dŽdiŽs aux graphes a ŽtŽ source d'une augmentation d'intŽrˆt de la part des milieux scientifiques et industriels [32]. Parmi les techniques adaptŽes les plus frŽquentes, on retrouve la classification et l'identification de tendances [30]. Nous nous sommes donc posŽs la question de l'adaptabilitŽ de cette thŽorie mathŽmatique aux sciences humaines, et plus particuliŁrement ˆ la linguistique de l'Žcrit.

Ainsi, les possibilitŽs d'analyse de ces structures ont d'jˆ ŽtŽ ŽtudiŽes dans le contexte du processus de l'Žcriture. La mise en Žvidence de sous-graphes particuliers reprŽsant les patterns des opŽrations les plus frŽquentes (ajout vs insertion, suppressions immŽdiate et diffŽrŽe, remplacement), a d'jˆ ŽtŽ rŽalisŽe. L'identification de ceux-ci est utile dans l'analyse du graphe reprŽsant le temps de l'Žcriture [7]. Il serait aussi possible d'avoir une visŽe plus exploratoire pour l'analyse et chercher par exemple des tendances ou autres sous-structures dans les donnŽes.

La recherche de ce type de structures peut se faire sur un seul large graphe pour y trouver des structures qui sont rŽcurrentes un nombre minimal de fois. Elle peut Žgalement se faire sur un groupe de graphes, pour lequel l'objectif est de trouver des structures qui sont rŽcurrentes dans une certaine fraction des graphes faisant partie de ce groupe [30]. En raison de la diffŽrence fondamentale entre les caractŽristiques sous-jacentes aux donnŽes et ˆ la formulation du problŁme, les algorithmes qui sont dŽveloppŽs pour l'analyse d'un groupe de graphes ne

peuvent être appliqués à l'analyse des caractéristiques d'un seul graphe. Inversement, les algorithmes développés pour l'analyse d'intérieur d'un large graphe peuvent être facilement adaptés pour l'analyse d'un ensemble de graphes [32].

Dans l'analyse du temps de l'écriture, les deux approches sont également pertinentes, pour d'abord repérer les stratégies courantes d'un scripteur, et pour ensuite comparer les stratégies des différents scripteurs entre eux. Les méthodes d'analyse seront appliquées sur les différents niveaux de graphes pour tenter de détecter des stratégies ou opérations nouvelles.

L'utilisation de ces propriétés en analyse de l'écriture permettrait ainsi de comparer plusieurs textes d'un même auteur pour consolider et mieux comprendre ses propres stratégies de rédactions sous forme de sous-ensembles de données. Il serait également possible de prendre plusieurs textes d'un ensemble de scripteurs ayant un niveau de maîtrise d'écriture semblable et ainsi de tenter de découvrir les recurrences entre des stratégies individuelles de rédaction.

H (I +23) &' *(, 0*C+ / +\$?-' (

Nous nous proposons maintenant de préciser le cadre méthodologique dans lequel a été enregistré notre corpus.

H B ((((J 0*C+ / +\$?-' (

Ce projet est centré sur l'écriture, comme processus descriptible linguistiquement, s'inscrivant dans le temps et l'espace, c'est-à-dire dans l'avant d'un texte finalement stabilisé. Ce choix s'appuie sur le constat méthodologique, selon lequel le produit n'est pas l'image fidèle de la production. En adoptant une approche génétique du texte, associée à un logiciel d'enregistrement de l'écriture, nous cherchons à décrire comment divers scripteurs donnent lieu à diverses rédactions et révisions, invisibles dans le texte achevé. Trois orientations émergent alors: (1) le rôle des opérations d'écriture-écriture qui sous-tendent toute activité scripturale, (2) la chronologie qui caractérise les écritures experte, novice, universitaire et scolaire, et (3) l'impact de la visualisation de cette chronologie sur la description et la reproduction de ces écritures. L'objectif est de démontrer que la temporalité est une dimension essentielle de l'écriture, tout autant que la spatialité. L'étude, bien trop souvent négligée de cette temporalité, est intimement liée aux modes de représentations et de visualisations numériques. Ce que nous nous proposons d'étudier, afin de répondre à une interrogation contemporaine fondamentale: comment donner à voir ce qui se passe précisément sur un écran, pendant le temps de l'écriture. L'équipe² a développé une méthodologie propre, à partir d'un outil mathématique (graphes) adéquat, qui permet, en parallèle aux modèles de représentations (SIG) déjà consacrés par la linguistique cognitive, une vraie cartographie génétique de la dynamique scripturale. Ce qui a nécessité la mise en place d'un logiciel dédié proche des attentes de la linguistique génétique.

H B (((((K " (" +) 4' #) (32+ ?2#, , ' (/ A " 2 ?-&*2 , ' " *(

Le programme GenographiX (GGX), a été mis en place spécifiquement pour mettre en évidence le travail génétique de l'écriture des scripteurs. Voici l'environnement graphique offert par le programme :

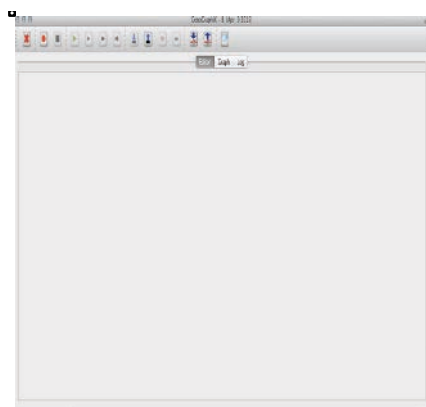


Fig. 3. Le programme GenographiX.

Comme la figure 3 le montre, le programme GenoGraphiX est expérimental, fonctionnant uniquement, pour l'instant du moins, dans un environnement fermé. Le logiciel permet d'enregistrer le détail de l'activité du scripteur, et de représenter cette activité par l'entremise d'un graphe mathématique tel qu'il a été décrit précédemment. Les données brutes (log) sont enregistrées dans un format similaire à un programme de type Scriptlog. La liste des événements est alors analysée par le logiciel qui regroupe ces événements en cellules (ou nœuds, ou arêtes) et identifie les liens (ou arcs) qui peuvent exister entre les cellules. Dans la version courante du logiciel, chaque cellule est représentée par un rectangle dans lequel on peut lire le texte correspondant (les versions antérieures représentaient chaque cellule par un cercle, mais le texte correspondant n'était pas visible). Ces cellules peuvent être de trois natures : 1) du texte ajouté que le scripteur a supprimé avant la fin de la rédaction (en jaune), 2) du texte ajouté qui restera dans le document final (en rouge) et 3) du texte supprimé (en bleu). Des liens entre cellules représentent une relation temporelle ou spatiale. Comme nous l'avons déjà signalé, pour des raisons méthodologiques, une attention particulière est portée sur le texte final qui peut toujours être lu en suivant les cellules et les liens rouges. Il est toutefois aussi possible de suivre le travail de manière chronologique en suivant les liens en trait continu (les autres liens sont en pointillés). Afin de marquer la relation particulière qu'il y a entre l'ajout d'une portion de texte et sa suppression ultérieure, la cellule d'ajout et la cellule de suppression correspondant sont reliées par un lien vert. Un exemple de graphe est représenté par la figure 5.

HE(K'' (/ +) 7\$ (*L1C' (

Il s'agit d'un corpus numérique d'écritures enregistrées, composé de deux tâches successives. La première, de nature narrative, se présente comme suit :

Ç Racontez une histoire. Vous devez intégrer dans l'ordre et sans les modifier les trois phrases :

-Elle habitait dans cette maison depuis longtemps.

-Il se retourna en entendant ce grand bruit.

-Depuis cette aventure, les enfants ne sortent plus la nuit. "

La seconde tâche, de nature argumentative, est la suivante :

Ç Qu'est-ce qui est important pour vivre ensemble ? É. Un sujet qui se prête à tous les âges, sans reposer sur une structure de type pour/contre représentative d'un mode dominant des textes argumentatifs.

Précisons que la première reprend un protocole qui a déjà fait l'objet de recueil de corpus de textes auprès d'enfants et de futurs enseignants en formation, proposé au départ par M. Charolles dans le cadre d'études de la cohérence textuelle et repris notamment dans Garcia-

Debanc et Bonnemaïson [33]. L'analyse des processus d'écriture trouvera ainsi des points d'appui dans les études existantes des textes produits, étant donné que ces études, comme nous l'avons déjà souligné dans notre introduction, sont déjà très installées. Ce corpus est déjà partiellement enregistré ; il se compose comme suit :

Tableau 1. Corpus d'écritures enregistrées.

PUBLICS			SCRIPTEURS	TÂCHES (20 minutes/tâche)	
				Narration	Argument
Enfants	-	-	-	-	-
Adultes	Etudiants (Langue Étrangère)	Licence	3	3	3
		Master	1	1	1
		Doctorat	-	-	-
	Etudiants (Langue maternelle)	Licence	6	6	1
		Master	3	3	3
		Doctorat	2	2	2
Enseignants-chercheurs	Doctorat	3	3	3	
TOTAL en heures et minutes				= 6h	= 4h20

Il a été fait le choix de consignes même de fonctionner auprès de publics variés (âge, langues maternelle vs étrangère et niveau d'expertise), la distinction enfant vs adulte recoupant principalement des lieux d'enregistrement. Le corpus est donc composé, pour l'heure actuelle, de 14 scripteurs en langue maternelle et de 4 scripteurs en langue étrangère, pour un total de 10h 20 minutes d'enregistrement. Ce qui semble déjà caractéristique est le fait que ce soit le nombre d'heures d'écriture qui importe, et non plus, le nombre de pages rédigées. Cette distinction reste éminente entre des corpus de produits et des corpus de production.

REK" (" *2 M ') OF\$(2 320& " *#*+ "(, -O" (

Pour illustrer notre propos, nous proposons de suivre l'écriture d'une adulte de 30 ans. La représentation progressive est un hybride entre la représentation graphique par les graphes et les SIG [34]. Elle a été créée dans le but de combiner les avantages de ces deux représentations et de pouvoir représenter plus de variables en s'inspirant de bonnes pratiques en matière de visualisations. Voici, avant de donner l'exemple d'une représentation mixte, nommée progressive par Becotte-Boutin, Caporossi et Hertz [34], le texte produit (Lil-Gon-30-F-T1-Fi-Fr.txt). Il s'agit d'un enregistrement de 20 minutes :

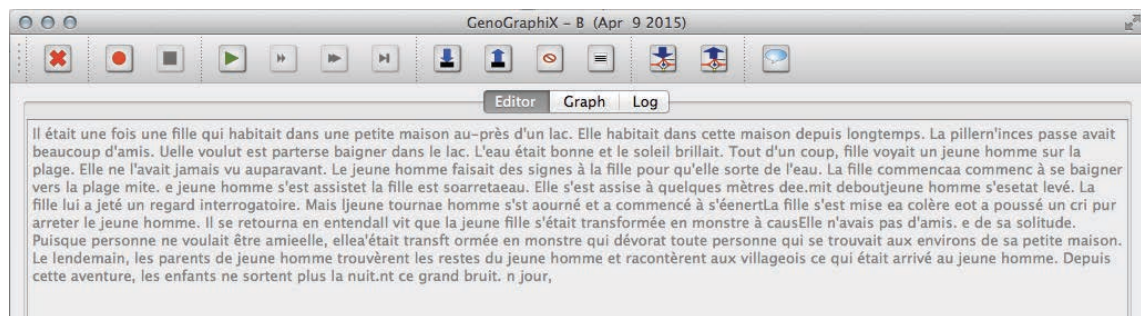


Fig. 4. Texte produit.

Comme il est possible de le voir dans l'éditeur de GGX, le texte produit apparaît sous le format txt. Voyons maintenant les 20 premières secondes de cet enregistrement, sous la forme d'un graphe :

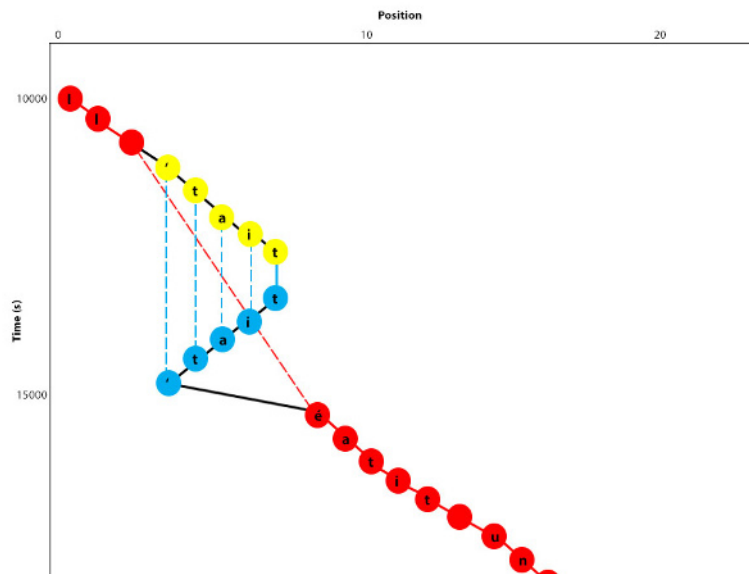


Fig. 5. Représentation mixte, dite progressive.

Ce graphe est mixte parce qu'il reprend une représentation inspirée des axes X (position spatiale) et Y (chronologie). Dans ce modèle mixte, il faut noter que les trois premières opérations effectuées, soit l'insertion d'un L, e et d'un espace sont en rouge. Ceci signifie qu'elles sont présentes dans la version finale du texte. Par la suite, la présence des opérations successives Ç Ò tait È en jaune et leur double en bleu qui se situe au même emplacement sur l'axe des abscisses signifie que ces opérations ont d'abord été effectuées (en jaune) mais par la suite supprimées (en bleu). Le lien bleu pointillé qui relie l'opération jaune avec son double bleu signifie qu'elles occupent le même emplacement spatial dans le texte. Au bas de la représentation il est aisé de remarquer les opérations successives Ç tait une È, qui sont en rouge tout comme le début et sont d'ailleurs liées à l'espace en position 3, ce qui signifie que ces opérations sont successives dans la version finale du texte. En lisant successivement les opérations en rouge, liées chronologiquement par des liens de couleur rouge également, il est possible de constater que le texte final de cette portion représentée, peut se lire Ç Il tait une È, conformément à ce qu'il est indiqué à la figure 4.

Puisque les SIG ont cette particularité de bien mettre en relation l'aspect temporel et spatial du temps de l'écriture [17], ils ont en contrepartie le désavantage de nous donner une vue d'ensemble du processus. La conséquence de cet attribut est de rendre difficile pour l'analyste de le lier avec opérations et le contexte d'écriture [8]. Pour sa part, la représentation par les graphes met davantage l'accent sur le lien entre les opérations et leur séquence [7]. Il est cependant difficile de connaître le contenu exact de chaque nœud et la position temporelle des opérations.

Cette représentation reprend d'abord l'idée des graphes génétiques mais est incluse dans un graphique ayant comme axe des abscisses la position et comme axe des ordonnées le temps. Contrairement aux SIG dont le déroulement du temps se lit de bas en haut, cette représentation se lit de gauche à droite et de haut en bas, en suivant le modèle logique d'un texte [34]. Par

rapport à la représentation par les graphes, dans ce cas-ci, chaque nœud représente une opération individuelle.

La sous-structure des données de cette représentation possède entièrement les attributs des deux représentations. Dans un graphe, il est possible d'avoir des attributs complémentaires, souvent considérés comme des poids [26, p. 15]. Chaque sommet représente ainsi une opération effectuée. L'opération peut être soit une insertion ou une suppression. De façon similaire aux fichiers log, chaque opération possède les attributs temps, position et type. Les arcs possèdent également des attributs. Nécessairement, ils ont une origine et une destination mais dans ce cas-ci ils ont aussi un type [34]. Cette structure permet ainsi d'appliquer les notions mathématiques et modèles propres à la théorie des graphes pour tenter de déceler des tendances et sous-structures.

Les attributs des sommets peuvent dans un deuxième temps être utilisés plus simplement en tant que coordonnées et reprendre les méthodes propres aux SIG.

Conclusion

Au sein des approches interactives par les processus d'écriture-lecture sur écran, les outils numériques d'évaluation sont en train de devenir un véritable objet de recherche, tout autant que les modèles de productions textuelles. Parmi ces outils d'évaluation, la visualisation de données occupe une place de choix. Se situant directement dans la tradition des transcriptions réalisées par la graphologie (schématiquement diplomatiques versus linéaires), la visualisation des données scripturales prolonge ces paradigmes en proposant des données dynamiques, et non plus statiques. Il n'est alors plus seulement question d'évaluer par l'image ce qui vient d'être écrit, mais bien d'évaluer par l'image ce qui est en train d'être écrit, d'être écrit sur écran, dans des sequentialités singulières.

Le principal avantage d'utiliser des méthodes dérivées d'une structure de données propre à un type de représentation du processus d'écriture est de pouvoir à la fois bénéficier des capacités cognitives humaines à détecter des tendances et récurrences dans des données visualisées et de confirmer ces hypothèses avec des méthodes mathématiques. Parmi les différentes représentations explorées, il semble que l'utilisation de la représentation mixte, dite progressive, permette de structurer les données en ayant accès à un bassin plus large de méthodes. Il serait tout de même intéressant de comparer ces représentations et certaines de leurs méthodes d'analyse respectives afin de confirmer leurs avantages et élaborer un cadre méthodologique pour permettre aux chercheurs d'avoir un accès facile à ces compléments d'analyse.

Dans l'état actuel de notre projet, les attentes sont les suivantes: 1) Affiner notre outil d'enregistrement (GGX) afin de pouvoir aussi enregistrer dans des environnements logiciels comme, par exemple, *Word* ou *Open Office*. Pour l'instant, cet outil reste confidentiel. 2) Mettre en place une base de données, tendue, composée d'heures d'enregistrement. Une telle base de données n'existe pas actuellement en langue française. 3) Pouvoir rendre accessible cette base de données sur la plate-forme d'Idée (TEMPUS), l'idée étant bien de mettre à disposition des chercheurs des données temporelles variées. 4) Développer les concepts théoriques nécessaires. Et 5) de procéder à des descriptions affinées en littérature d'habitué et experte.

(

6002 "1' &((

(

- [1] D. Boyd et K. Crawford, *ÇCritical questions for big data,È Information, Communication & Society, vol. 15, n^o 15, pp. 662-679, 2012.*
- [2] J. Manyika, M. Chui, B. Brown, J. Bughin, R. Dobbs, C. Roxburgh et A. Hung Byers, *ÇBig data: the next frontier for innovation, competition, and productivity,È McKinsey Global Institute, 2011.*
- [3] S. Plane, D. Alamargot et J.-L. Lebrave, *ÇTemporalitŽ de l'Žcriture et r^mle du texte produit dans l'activitŽ rŽdactionnelle,È Langages, vol. 177, n^o 11, pp. 11-32, 2010.*
- [4] C. Doquet-Lacoste, *ftude GŽnŽtique de l'fcriture sur Traitement de Texte d'fl•ves de Cours Moyen 2, AnnŽe 1995-1996, Paris: UniversitŽ Sorbonne nouvelle, 2003, p. 633.*
- [5] K. S. Miller et K. P. Sullivan, *ÇKeystroke Logging: An introduction,È chez Computer keystroke logging and writing, Oxford, Elsevier, 2006, pp. 1-10.*
- [6] K. P. Sullivan et E. Lindgren, *ÇLa rŽvision en production Žcrite enregistrŽe,È chez Temps de l'criture: enregistrements et repr sentations, C. Leblay et G. Caporossi, fds., Louvain-la-Neuve, Academia, 2014, pp. 71-92.*
- [7] G. Caporossi et C. Leblay, *ÇOutils de visualisation de donnŽes enregistrŽes,È chez Temps de l'criture: enregistrements et repr sentations, C. Leblay et G. Caporossi, fds., Louvain-la-Neuve, Academia, 2014, pp. 147-166.*
- [8] G. Caporossi et C. Leblay, *ÇOnline Writing Data Representation : A Graph Theory Approach,È chez Lecture Notes in Computer Sciences 7014, 2011, pp. 80-89.*
- [9] N. Yau, *Visualize this: the flowing data guide to design, visualization and statistics, Indianapolis: Wiley Publishing, 2011, p. 358.*
- [10] M. Minelli, M. Chambers et A. Dhiraj, *Big data, big analytics : Emerging business intelligence and analytic trends for today's businesses, Wiley, 2013, p. 224.*
- [11] F. Blanchard, *Visualisation et classification de donn es multidimensionnelles application aux images multicomposantes, UniversitŽ de Reims Champagne Ardenne , 2005.*
- [12] M. Tory et T. Moller, *ÇHuman factors in visualization research,È IEEE Transactions on visualization and computer graphics, vol. 10, n^o 11, pp. 72-84, 2004.*
- [13] G. Dzemyda, O. Kurasova et J. Zilinskas, *ÇMultidimensional Data Visualization: Methods and Applications,È Springer, Vilnius, 2013.*
- [14] T. Moller, B. Hamann et R. Russell, *ÇPreface,È chez Mathematical foundation of scientific visualization, computer graphics, and massive data exploration, Berlin, Springer, 2009, pp. V-VII.*
- [15] H. A. Karimi, *ÇBig data :Techniques and technologies in geoinformatics,È Boca Taton: CRC Press, 2014.*

- [16] S. Tuffřry, *Data mining et statistique dřcisionnelle: l'intelligence des donnřes*, fditions Technip, 2010, p. 705.
- [17] E. Lindgren et K. P. Sullivan, řThe LS Graph : A Methodology for Visualizing Writing Revision,ř Language Learning, vol. 52, nř 13, pp. 565-595, 2002.
- [18] J. Van Sickle, *Basic GIS coordinates*, Second Edition, Boca Raton: CRC Press, 2010.
- [19] E. Lindgren, K. P. H. Sullivan, U. Lindgren et K. Spelman Miller, řGIS for Writing: Applying Geographical Information Systems Techniques to Data Mine Writings' Cognitive Processes,ř chez Writing and Cognition, Amsterdam, Elsevier, 2007, pp. 83-96.
- [20] C. D. Lloyd, *Spatial data analysys: an introduction for GIS users*, New York: Oxford University Press, 2010.
- [21] M. Leijten et L. Van Waes, řInputlog : New Perspectives on the Logging of On-Line Writing Processes,ř chez Computer Keystroke Logging and Writing, K. P. S. a. E. Lindgren, f.d., Elsevier, 2006, pp. 73-94.
- [22] A. Wengelin, M. Torrance, K. Holmwvist, S. Simpson, D. Galbraith, V. Johansson et R. Johansson, řCombined eyetracking and keystroke-logging methods for studying cognitive processes in text production,ř Behavior Research Methods, vol. 41, nř 12, pp. 337-351, 2009.
- [23] M. Leijten et L. Van Waes, řKeystroke Logging in Writing Research : Using Inputlog to Analyze and Visualize Writing Processes,ř vol. 30, nř 13, pp. 358-392, 2013.
- [24] C. Leblay, řLe Temps de l'řcriture. Gen•se, durře, reprřsentations,ř 2011. [En ligne]. Available: <https://www.jyu.fi/ajankohtaista/arkisto/2011/11/tiedote-2011-11-04-10-14-59-722468>. [Acc•s le 15 12 2013].
- [25] V. Southavilay, K. Yacef, P. Reimann et R. A. Calvo, řAnalysis of Collaborative Writing Processes Using Revision Maps and Probabilistic Topic Models,ř Proceedings of the Third International Conference on Learning Analytics and Knowledge, pp. 38-47, 2013.
- [26] A. J. Bondy et U. S. R. Murty, *Graph theory with applications*, New York: Elsevier Science Publishing Co. Inc., 1982.
- [27] R. K. Ahuja, T. L. Magnanti et J. B. Orlin, *Network flows: Algorithms and applications*, Upper Saddle River: Prentice Hall, 1993.
- [28] E. J. Henley et R. A. Williams, *Graph theory in modern engineering*, Houston: Academic Press, 1973.
- [29] S. Saha Ray, *Graph theory with algorithms and its applications in applied science and technology*, Rourkela, India: Springer India, 2013.
- [30] C. C. Aggarwal et H. Wang, řGraph data management and mining: a survey of algorithms and applications,ř chez Managing and mining graph data, C. C. Aggarwal et H. Wang, f.ds., New York, Springer, 2010, pp. 13-68.
- [31] L. Takac et M. Zabovsky, řData analysis in public social networks,ř Lomza, Poland, 2012.

- [32] M. Kuramochi et G. Karypis, ÇFinding frequent patterns in a large sparse graph,È Data mining and knowledge discovery, vol. 11, n^o 13, pp. 243-271, 2005.
- [33] C. Garcia-Debanco et K. Bonnemaïson, ÇLa gestion de la cohésion textuelle par des Žlives de 11-12 ans: rŽussites et difficultŽs,È chez CMLF 2014, Toulouse.
- [34] H.-S. Becotte-Boutin, G. Caporossi et A. Hertz, ÇThe progressive visualization, a new tool for analyzing the writing process,È Cahiers du GERAD, vol. 131, 2015.
- [35] H. Zha, Y. Yang, J. Wang et L. Wen, ÇTransforming XPDL to Petri Nets,È chez Business Process Management Workshops.
- [36] Y. Wu, L. Wang, J. Ren et W. Ding, ÇMining sequential patterns eith periodic wildcard gaps,È Applied intelligence, vol. 41, n^o 11, pp. 99-116, 2014.
- [37] WritingPro, ÇWriting Pro,È 2014. [En ligne]. Available: <http://www.writingpro.eu/>. [Acc’s le 02 08 2014].
- [38] I. Witten, ÇText mining,È Hamilton, New Zealand, 2004.
- [39] G. Weikum, J. Hoffart, N. Nakashole, M. Spaniol, F. Suchanek et M. A. Yosef, ÇBig data methods for computational linguistics,È IEEE Data Engineering Bulletin, vol. 35, pp. 46-55, 2012.
- [40] J. Veronis, ÇComputerized Correction of Phonographic Errors,È Computers and the Humanities, vol. 22, n^o 11, pp. 43-56, 1988.
- [41] L. Van Waes et P. J. Schellens, ÇWriting Profiles : The Effect of the Writing Mode on Pausing and Revision Patterns of Experienced Writers,È Journal of Pragmatics, vol. 35, pp. 829-853, 2003.
- [42] L. Van Waes et M. Leijten, ÇInputlog 6.0: State of the art.,È chez Paper presented at the Keystroke logging training school, Antwerp, 2014.
- [43] L. Van Waes et M. Leijten, ÇInputlog 6.0: Pause and fluency analysis.,È chez Paper presented at the keystroke logging training school, Antwerp, 2014.
- [44] R. E. Tarjan, ÇProblems in data structures and algorithms,È chez Graph theory, combinatorics and algorithms, M. C. Golumbic et I. B. Hartman, eds., Springer, 2005, pp. 17-39.
- [45] S. Stromqvist, K. Holmqvist, V. Johansson, H. Karlsson et A. Wengelin, ÇWhat Keystroke Logging can Reveal about Writing,È chez Computer Keystroke Logging and Writing, K. P. S. & E. Lindgren, eds., Elsevier, 2006, pp. 45-71.
- [46] K. Spelman Miller et K. P. Sullivan, ÇKeystroke Logging : An Introduction,È chez Computer Keystroke Logging and Writing, Elsevier, 2006, pp. 1-10.
- [47] K. Severinson Eklundh et P. Kollberg, ÇA Computer Tool and Franework for Analysing On-Line Revisions,È chez The Science of Writing : Theories, Methods, Individual Differences, and Applications, Mahwah, NJ, Lawrence Erlbaum, 1996, pp. 163-188.
- [48] D. H. Roen et R. Willey, ÇThe Effects of Audience Awareness on Drafting and Revising,È Research in the Teaching of English, vol. 22, n^o 11, pp. 75-88, 1988.

- [49] S. Plane, D. Alamargot et J.-L. Levrabe, *ÇTemporalitŽ de l'écriture et R^mle du Texte Produit dans l'ActivitŽ RŽdactionnelle*, *È Langages*, vol. 177, nⁱ %11, pp. 11-32, 2010.
- [50] D. Perrin, *ÇProgression Analysis (PA): Investigating Writing Strategies at the Workplace*, *È Journal of Pragmatics*, vol. 35, pp. 907-921, 2003.
- [51] R. D. Owston, S. Murphy et H. H. Wideman, *ÇThe Effects of Word Processing on Student's Writing Quality and Revision Strategies*, *È Research in the Teaching of English*, vol. 26, nⁱ %13, pp. 249-276, 1992.
- [52] E. New, *ÇComputer-Aided Writing in French as a Foreign Language : A Qualitative and Quantitative Look at the Process of Revision*, *È The Modern Language Journal*, vol. 83, nⁱ %11, pp. 80-97, 1999.
- [53] G. Miner, J. Elder, B. Nibset, D. Delen, A. Fast et T. Hill, *Practical text mining and statistical analysis for non-structured text data applications*, Academic Press Žd., Saint-Louis, 2012, p. 1094.
- [54] Merriam-Webster, *ÇApopheia*, Merriam-Webster, È 2014. [En ligne]. Available: <http://www.merriam-webster.com/dictionary/apopheia>. [Acc•s le 03 08 2014].
- [55] J. Leskovec, J. Kleinberg et C. Faloutsos, *ÇGraph evolution: Densification and shrinking diameters*, *È ACM Transactions on Knowledge Discovery from Data*, vol. 1, nⁱ %11, p. Article 2, 2007.
- [56] M. Leijten et L. Van Waes, *ÇInputlog features*, È 2014. [En ligne]. Available: http://www.inputlog.net/description_features.html. [Acc•s le 25 07 2014].
- [57] J.-L. Lebrave et A. GrŽsillon, *ÇLinguistique et gŽnŽtique des textes : un dŽcalogue*, È 23 03 2009. [En ligne]. Available: <http://item.ens.fr/index.php?id=434571>. [Acc•s le 14 12 2013].
- [58] J.-L. Lebrave, *ÇComment Žcriront-ils?*, È *Diog ne*, vol. 196, nⁱ %14, pp. 163-171, 2001.
- [59] C. Leblay et G. Caporossi, *ÇIntroduction aux donnŽes temporelles de l'Žcriture*, È *chez Temps de l'criture: enregistrements et repr•sentations*, C. Leblay et G. Caporossi, *fds.*, Louvain-la-Neuve, Academia, 2014, pp. 5-15.
- [60] C. Leblay, *ÇEn de• du bien et du mal Žcrire*, È 14 06 2012. [En ligne]. Available: <http://www.item.ens.fr/index.php?id=578258>. [Acc•s le 07 07 2014].
- [61] M. M. A. Latif, *ÇA State-of-the-Art Review of the Real-Time Computer-Aided Study of the Writing Process*, È *International Journal of English Studies*, vol. 8, nⁱ %11, pp. 29-50, 2008.
- [62] P. Kollberg, *ÇS-notation as a tool for analysing the episodic nature of revisions*, È Barcelona, 1996.
- [63] R. Karp, *ÇOptimization problems related to internet congestion control*, È *chez Graph theory, combinatorics and algorithms*, M. C. Golubic et I. B. Hartman, *fds.*, Springer, 2005, pp. 1-16.
- [64] G. Kanawaty, *Introduction ^ l'Žtude du travail 3e Ždition*, Bureau International du Travail Žd., Gen•ve, 1996.

- [65] S. Helbing et S. Baliotti, *ÇFrom social data mining to forecasting socio-economic crises,* *È The European Physical Journal Special Topics*, vol. 195, pp. 3-68, 2011.
- [66] C. Haas, *ÇHow the Writing Medium Shapes the Writing Process : Effects of Word Processing on Planning,* *È Research in the Teaching of English*, vol. 23, n^o 12, pp. 181-207, 1989.
- [67] D. Foucambert et J. Foucambert, *ÇGestes d'écriture et caractéristiques linguistiques des textes achevés,* *È chez Temps de l'écriture: enregistrements et représentations*, C. Leblay et G. Caporossi, *fds.*, Louvain-la-Neuve, *Academia*, 2014, pp. 43-70.
- [68] C. Doquet, *ÇPour une approche linguistique de l'écriture enregistrée,* *È chez Temps de l'écriture: enregistrements et représentations*, *Academia*, *fd.*, Louvain-la-Neuve, 2014, pp. 21-42.
- [69] M. Dehmer et S. C. Basak, *Statistical and machine learning approaches for network analysis*, Somerset: John Wiley & Sons, 2012.
- [70] C. De Looze, *Analyse et Interprétation de l'Empan Temporel des Variations Prosodiques en Français et en Anglais*, Marseille: Université de Provence, 2010.
- [71] M. Cox, C. Ortmeier-Hopper et K. E. Tirabassi, *ÇTeaching Writing for the "Real World": Community and Workplace Writing,* *È The English Journal*, vol. 98, n^o 15, pp. 72-80, 2009.
- [72] J. Cohen, B. Dolan, M. Dunlap, J. M. Hellerstein et C. Welton, *ÇMAD skills: new analysis practices for big data,* *È Proceedings of the VLDB Endowment 2.2*, pp. 1481-1492, 2009.
- [73] D. Chakrabarti et C. Faloutsos, *Graph mining: laws, tools, and case studies*, Morgan & Claypool Publishers series, 2012.
- [74] A. Carbone et M. Gromov, *ÇMathematical slices of molecular biology,* *È Gazette des mathématiciens*, édition spéciale, vol. 88, pp. 11-80, 2001.
- [75] I. Breetvelt, H. Van Den Bergh et G. Rijlaarsdam, *ÇRelations between Writing Processes and Text Quality : When and How,* *È Cognition and Instruction*, vol. 12, n^o 12, pp. 103-123, 1994.
- [76] M. Bramer, *Principles of data mining*, 2nd éd., Portsmouth: Springer, 2013.
- [77] M. R. Berthold, *ÇBisociative knowledge discovery,* *È chez Advances in intelligent data analysis X*, J. Gama, E. Bradley et J. Hollmen, *fds.*, Porto, Springer, 2011, pp. 1-7.
- [78] A.-L. Barabasi, R. Albert et H. Jeong, *ÇScale-free characteristics of random networks: the topology of the world-wide web,* *È Physica A*, vol. 281, pp. 69-77, 2000.
- [79] D. Alamargot, G. Caporossi, D. Chesnet et C. Ros, *ÇWhat Makes a Skilled Writer? Working Memory and Audience Awareness During Text Composition,* *È Learning and Individual Differences*, vol. 21, pp. 505-516, 2011.
- [80] R. Agrawal, T. Imielinski et A. Swami, *ÇMining association rules between sets of items in large databases,* *È New York*, 1993.
- [81] L. Flower et J. R. Hayes, *ÇA cognitive process theory of writing,* *È College Composition*

- and Communication, vol. 32, n^o 14, pp. 365-387, 1981.
- [82] P. Kollberg, *Rules for the S-notation: a computer-based method for representing revisions*, IPLab, Royal Institute of Technology (KTH), Stockholm, Sweden, 1996.
- [83] D. Alamargot et J.-L. Lebrave, *The study of professional writing: A joint contribution from cognitive psychology and genetic criticism*, European Psychologist, n^o 14, doi:10.1027/1016-9040/a000001, 2009.
- [84] T. Olive, J.-L. Lebrave, J.-M. Passerault et N. Le Bigot, *La dimension visuo-spatiale de la production de textes: approches de psychologie cognitive et de critique g \acute{e} ographique*, Langages, vol. 177, n^o 11, pp. 29-55, 2010.
- [85] W. Daelemans, P. Berck et S. Gillis, *Data mining as a Method for Linguistic Analysis: Dutch Diminutives*, Dutch Diminutives, Folia Linguistica, XXXI/I -2, pp. 57-75, 1997.
- [86] P. A. Bath, C. Craigs, R. Maheswaran, J. Raymond et P. Willett, *Use of graph theory to identify patterns of deprivation and high morbidity and mortality in public health data sets*, Journal of the American Medical Informatics Association, vol. 12, n^o 16, pp. 630-641, Nov-Dec 2005.
- [87] A. Markowetz, K. Blasiewicz, C. Montag, C. Switala et E. S. Thomas, *Psycho-Informatics: Big Data shaping modern psychometrics*, Medical Hypotheses, vol. 82, n^o 14, pp. 405-411, April 2014.
- [88] F. J. Ohlhorst, *Big Data Analytics : Turning Big Data into Big Money*, Wiley, 2013, p. 176.
- [89] E. Lindgren et K. P. Sullivan, *Writing and the Analysis of Revision: An Overview*, chez Computer Keystroke Logging and Writing, Elsevier, 2006, pp. 31-44.
- [90] ITEM, *Enjeux de recherche*, 7 mai 2014. [En ligne]. Available: <http://www.item.ens.fr/index.php?identifiant=l-item>. [Acc \acute{e} s le 14 juin 2014].
- [91] I. Breetvelt, H. v. d. Bergh et G. Rijlaarsdam, *Relations between writing processes and text quality: When and how?*, Cognition and Instruction, vol. 12, n^o 12, pp. 103-123, 1994.
- [92] C. Lynch, *How do your data grow?*, Nature, vol. 455, n^o 14, pp. 28-29, 2008.
- [93] E. Midgette, P. Haria et C. MacArthur, *The effects of content and audience awareness goals for revision on the persuasive essays of fifth- and eighth-grade students*, Reading and Writing, vol. 21, pp. 131-151, 2008.
- [94] E. W. Bethel, P. Prabhat, S. Byna, O. Rubel, K. J. Wu et M. Wehner, *Why high performance visual data analytics is both relevant and difficult*, Burlingame, California, USA, 2013.
- [95] F. A. Abukhodair, B. E. Riecke, H. I. Erhan et C. D. Shaw, *Does interactive animation control improve exploratory data analysis of animated trend visualization?*, Burlingame, California, USA, 2013.
- [96] L. Bartram, *Perceptual and interpretative properties of motion for information visualization*, Proceedings of the 1997 workshop on new paradigms in information visualization and manipulation, pp. 3-7, 1997.

- [97] D. Alamargot, G. Caporossi, D. Chesnet et C. Ros, *What makes a skilled writer? Working memory and audience awareness during text composition*, *Learning and Individual Differences*, vol. 21, n°15, pp. 505-516, 2011.
- [98] V. M. Baaijen, D. Galbraith et K. d. Glopper, *Keystroke analysis: reflections on procedures and measures*, *Written Communication*, vol. 29, n°13, pp. 246-277, 2012.
- [99] A. Kirk, *Data visualization: a successful design process* [electronic book], Packt Pub., 2012.
- [100] A. Wengelin, *Temps et pauses dans l'écriture au clavier*, *chez Temps de l'écriture: enregistrements et représentations*, Louvain-la-Neuve, Academia, 2014, pp. 97-124.
- [101] W. Aigner, S. Miksch, H. Schumann et C. Tominski, *Visualization of Time-Oriented Data*, *Springer*, London, 2011.
- [102] D. Perrin et S. Laemmel, *Application à l'écriture journalistique*, *chez Temps de l'écriture, enregistrements et représentations*, Louvain-la-Neuve, Academia-L'Harmattan s.a., 2014, pp. 171-192.
- [103] A. Unwin, C.-h. Chen et W. K. Hardle, *Introduction*, *chez Handbook of Data Visualization*, Berlin, Springer, 2008, pp. 3-14.
- [104] A. Vathy-Fogarassy et J. Abonyi, *Graph-Based clustering and data visualization*, Springer, 2013.
- [105] D. Chesnet et D. Alamargot, *Eye and Pen 2 manuel de l'utilisateur*, *Poitiers*, 2011.
- [106] E. R. Tufte, *The visual display of quantitative information*, Second *éd.*, Cheshire, Connecticut: Graphics Press, 2001.

¹ Deux graphes, ou sous-graphes, sont dits isomorphiques lorsqu'ils représentent clairement des diagrammes semblables [26, p. 4]. De manière plus formelle, un graphe $G_1 = (V_1, E_1)$ est isomorphe à un graphe $G_2 = (V_2, E_2)$ s'il y a une correspondance entre les ensembles de sommets V_1 et V_2 (et entre les ensembles d'arêtes E_1 et E_2) de telle manière que si e_1 est une arête joignant les sommets u_1 et v_1 dans le graphe G_1 , alors l'arête correspondante e_2 (dans G_2) joins les sommets u_2 et v_2 , ce qui correspond à u_1 et v_1 respectivement [29]. Les algorithmes pour la recherche de structures semblables dans les graphes se basent donc sur cette propriété [73].

² L'équipe TERS (Temps de l'Écriture et ses Représentations) est une équipe de l'Institut des Textes et Manuscrits Modernes (ITEM/ENS-CNRS).