

Un correcteur orthographique informatisé à l'épreuve de dictées d'élèves: quelle efficacité?

Jean Ravestein

MOTS CLÉS: Orthographe, correcteurs informatisés, TICE

Les correcteurs orthographiques et grammaticaux informatisés sont aujourd'hui implantés systématiquement dans les traitements de texte. Les élèves, comme leurs professeurs, s'en servent de fait sans véritable questionnement sur les implications didactiques de cette utilisation. Cette recherche mesure l'efficacité d'un correcteur informatisé sur des dictées d'élèves de CM2 (cinquième primaire). En fonctionnement automatique, il «corrige» 30% des erreurs, résultat encourageant. Une amélioration des algorithmes et une didactique spécifique pour ce public devraient en faire un outil efficace pour l'apprentissage de l'écriture à l'école, incontournable dans l'avenir, comme l'est devenue la calculette pour les mathématiques.

KEY WORDS: Spelling, computerized correctors, spell-checker, TICE

Computerized spelling and grammatical correctors are nowadays systematically built-in tools in word processing software. The pupils and their teachers are brought to use it without a clear vision of its didactical implications. This study measures the effectiveness of a computerized corrector on pupil's dictation (5th primary class). Under automatic process; it "corrects" 30% of the errors. This is an encouraging result. An improvement of the algorithms and specific didactics for this public should make a helpful tool for the training of writing at school, impossible to circumvent in the future, like the calculator became in mathematics.

Introduction

En France et dans les pays francophones, l'orthographe connaît aujourd'hui les faveurs d'un nombre croissant de citoyens. Il suffit pour s'en convaincre de constater l'engouement populaire pour les concours de dictées, relayés par les médias et parrainés par des éminences en matière littéraire, le tout avec le concours actif des institutions éducatives¹.

Sujet de passion et de débats contradictoires dans tous les milieux dès qu'on s'avise de vouloir en modifier les règles² (Catach, 1991; Goosse, 1991; Millet, Lucci & Billiez, 1990), elle reste néanmoins rétive à une maîtrise totale par les élèves comme par leurs professeurs. Pour cette «science qu'il n'y a aucune gloire à connaître, mais qu'il y a honte à ignorer» (Thimonnier, 1976, p. 26), il faut bien admettre que chacun porte une part de cette honte.

En effet, malgré qu'elle soit massivement enseignée dans les pays francophones dès l'âge de six ans et jusqu'à seize ans, puis l'objet des attentions des enseignants jusqu'à l'université, d'aucuns peuvent se targuer de la maîtriser parfaitement. L'apparition des correcteurs orthographiques puis grammaticaux couplés aux traitements de texte informatisés contribue aujourd'hui à alléger de manière non négligeable les angoisses des scripteurs, qu'ils soient élèves, secrétaires, et même universitaires.

Vu l'accroissement du parc de machines comportant ce type de logiciel, tant dans le cadre scolaire que privé, il semblerait que les correcteurs soient de plus en plus voués à être utilisés «naturellement».

La question vient se poser, pour le système d'enseignement, de savoir s'il faut superbement les ignorer en avançant sur un mode un peu rétrograde «qu'il ne faut pas méconnaître la nécessité innéductable de l'effort ingrat dans l'apprentissage» (Brunot & Gossot, 1957, p. 213) ou au contraire en cerner d'abord les avantages et les limites puis engager les chercheurs et les professeurs dans leur mise au point, leurs évolutions, et leur utilisation didactique.

Notre travail s'inscrit dans la deuxième proposition de l'alternative et dans son premier temps, en tentant de répondre à la question: comment «réagit» un correcteur orthographique standard confronté à des productions d'élèves?

Orthographe, école et machines

Si l'enseignement de l'orthographe «est marqué par des tensions entre options didactiques divergentes qui ne partagent pas les mêmes fondements conceptuels» (Allal, Béatrix Köhler, Rieben, Rouiller, Saada-Robert & Wegmuller, 2001, p. 53), il utilise paradoxalement des méthodes souvent voisines et d'une grande longévité (dont l'emblème est la méthode d'Édouard Bled) et reste un domaine de compétence central à l'école et au collège pour trois catégories de raisons: sociohistoriques, docimologiques et didactiques.

En France, sous la monarchie de Juillet, l'orthographe académique devient la seule obligatoire et officielle dans l'enseignement (décret de 1835). Ceci manifestait notamment une volonté politique d'unité nationale face à la diversité des langues régionales. Ce poids identitaire teinté de morale est encore présent dans les esprits aujourd'hui: ne parle-t-on pas de «faute» d'orthographe plutôt que d'erreur orthographique? Une faute s'expie devant le corps

social qui a produit les règles, une erreur, simplement, se corrige. Bien orthographier est donc souvent perçu comme une marque d'appartenance ou, plus fortement, d'allégeance, à la communauté qui a établi les règles: à une dissertation «géniale» – mais – «bourrée de fautes», on associera facilement une personnalité «intéressante mais désinvolte ou dilettante» et la note s'en ressentira.

L'orthographe reste également, sous sa forme «dictée», une manière commode de contourner les biais de la mesure en évaluation dénoncés par les travaux en docimologie (Piéron, 1969). En effet, à partir du moment où le barème de correction est précisément défini et où le correcteur s'y tient, on peut espérer que la mesure soit «juste», dans la double acceptation du terme: arithmétique et juridique. Ainsi, la dictée a-t-elle servi de base (coefficient élevé) à la sélection dans des concours et examens.

D'un point de vue didactique, l'enseignement de l'orthographe est théoriquement chose aisée, même si on peut, de manière récurrente (Burney, 1955), être déçu de ses résultats³. En effet, dans le système didactique qui tient solidaires, par le biais du contrat didactique, savoir transposé, maître et élèves dans des interrelations complexes, plus le savoir enseigné est bien défini dans la sphère savante, consigné, socialement nécessaire, correctement étalé dans le curriculum, plus il est «transparent» pour tous et accepté par tous comme objet de savoir légitime (Ravestein, 1999).

C'est le cas de l'orthographe: les règles peuvent être énoncées, reconnues facilement à l'extérieur de l'institution qui enseigne, servir de vecteur de communication aux familles sur l'apprentissage des enfants. Par exemple, si on s'entend dire: votre enfant n'est pas «bon» en orthographe, cela fait sens sans beaucoup d'explications supplémentaires et la famille peut se sentir à même de participer à l'amélioration de la situation, car elle partage souvent un savoir orthographique avec l'institution. Maîtriser l'orthographe est aussi pour l'enseignant un moyen sûr de positionnement par rapport à l'élève, qui, lui, ne se prive pas de créer des ruptures de contrat didactique (de la simple remarque à la vive protestation) lorsqu'il détecte une erreur commise publiquement par le maître (au tableau, par exemple) et ainsi le déstabiliser provisoirement.

Plusieurs tentatives institutionnelles ont fait l'assaut de la forteresse orthographe dans le milieu enseignant, de l'application de la réforme officielle de 1990 qui visait à la simplifier et à introduire des tolérances. Il s'ensuivit des instructions accompagnant la «réforme des cycles» de 1989 qui dédramatisaient, en considérant l'orthographe comme une «activité automatique» (CNDP, 1992, p. 55) qui devait bien finir par se mettre en place par imprégnation. Toutes semblent avoir fait les frais d'une résistance passive du corps enseignant relayé par le corps social⁴ qui jalouse presque ses enfants pour des facilités qu'on leur accorderait dans cet apprentissage où il a tant souffert.

De même, l'introduction des machines ou artefacts dans l'enseignement a été diversement accueillie (Baron et Bruillard, 1996, en ont décrit précisément les différentes phases): avec défiance par ceux qui l'ont vue comme une disqualification de l'utilisateur, avec enthousiasme par ceux qui ont pris un point de vue anthropocentrique (Clegg, 1988, Norman & Draper, 1986, comme initiateurs). Ce point de vue s'appuie sur les compétences existantes des utilisateurs et cherche à les développer: les artefacts sont alors considérés comme des outils modulables.

L'exemple le plus emblématique est celui de l'introduction des machines à calculer dès le cursus élémentaire, où l'on voit que vingt ans après, certains les refusent toujours à leurs élèves

alors que d'autres prévoient un enseignement spécifique de leur utilisation, se conformant ainsi aux instructions officielles.

En ce qui concerne les ordinateurs, leur coût, leur difficulté de maintenance, leur ergonomie défaillante du départ (pas d'interface de type «souris», sous DOS), la confusion qui a régné sur la formation des maîtres (enseignerait-on l'algorithmique ou le traitement de texte?), ont considérablement freiné leur implantation dans les établissements scolaires.

Il semble bien aujourd'hui en France que l'école est à la remorque de la société, comme pour l'audiovisuel en son temps, en ce qui concerne l'introduction raisonnée des machines dans les classes⁵. Comme «forcés de faire avec» (Le Brevet informatique et internet, B2i, qui doit être mis en place à l'école depuis la fin de 2000, en est un repère institutionnel⁶), émus, et parfois effrayés de l'habileté dont font preuve les plus jeunes élèves face à l'ordinateur, les enseignants semblent devoir faire leur deuil d'une part de leur territoire dans le domaine des procédures de transmission des savoirs⁷.

L'orthographe, on l'a vu, est encore une de leur «chasse gardée». On arrive donc à poser une question circulaire: pour convaincre, les correcteurs orthographiques devront être efficaces⁸, pour le devenir, il faudra l'engagement des maîtres, si, bien sûr, on suit un paradigme anthropocentrique, comme a pu le décrire Rabardel (1995):

- les artefacts sont transformés par l'activité des utilisateurs;
- les artefacts ne doivent pas être analysés comme des «choses» mais dans la façon dont ils médiatisent l'usage;
- les artefacts sont porteurs de partage et sont incorporés dans une pratique sociale.

Dans un premier temps, il convient de regarder comment un correcteur orthographique se comporte sur un corpus ordinaire de production d'élèves et, donc, dans quelle mesure il pourrait devenir ou non une aide efficace quant au produit rendu. Un travail de recherche sur des sujets apprenant le français comme langue étrangère (Charnet & Panckhurst, 1998), dont l'objectif était de faire un examen de type «banc d'essai» partant de situations didactiques réelles et de les confronter avec un logiciel de correction orthographique et syntaxique (*Correcteur 101* de la société Machina Sapiens), avait déjà conclu que le correcteur pouvait être un bon auxiliaire et permettait d'avoir une idée globale de la compétence linguistique du sujet.

Ainsi, nous avons soumis au travail d'un correcteur orthographique «courant⁹» (Cordial 6, intégré aujourd'hui à Microsoft Word sous forme simplifiée) un corpus de dictées de 100 élèves de CM2 (cinquième primaire)¹⁰.

Un correcteur orthographique à l'épreuve de productions d'élèves¹¹

Un même texte a été dicté aux 100 élèves, en fin d'année scolaire, par leur enseignant, à titre d'exercice ordinaire. Les 100 copies ont été ensuite tapées telles quelles (sous format.txt) et soumises au travail du correcteur orthographique. Enfin, à titre exploratoire, les 100 textes ont été soumis au correcteur grammatical. Voici le texte de la dictée:

La première hirondelle

Depuis plusieurs jours déjà, j'étais prêt: je guettais le retour des hirondelles. Un matin, je me mis à travailler à mon bureau, près de la fenêtre, et, tout à coup, elle a crié. La première était arrivée. Bientôt, avec un léger bruit d'aile(s), elle a filé comme une flèche puis elle est revenue se poser sous la gouttière. Je me suis alors approché sans bruit pour l'observer. Elle tendait son cou et bougeait

ses yeux sans jamais les fixer. Soudain, elle m'aperçoit, me salue d'un petit chant et s'enfuit bien vite. Je pensais aussitôt: «tu reviendras demain ma petite amie, tu sais bien que je ne te veux pas de mal». En effet, dès le soir, elles sont apparues, innombrables, en un vol désordonné au-dessus du pré.

Il s'agit d'une dictée réelle (relevée dans un manuel) de l'année scolaire considérée, légèrement modifiée. Elle présente l'intérêt de contenir plusieurs difficultés, notamment sur le plan des homophones. Le temps des verbes est varié mais sans forme rare. Le vocabulaire employé est relativement usuel. Elle présente aussi l'avantage d'être assez courte, donc facilement traitable (126 mots, sans les doublons).

Élaboration d'une grille de correction

Pour montrer l'intérêt et les limites du travail du correcteur, il a fallu dans un premier temps élaborer une grille de correction des erreurs pertinente au corpus. La documentation est abondante dans ce domaine et nous nous en sommes largement inspiré. On citera Gey (1987) dont la grille a le mérite de distinguer les différents homophones, Chervel et Manesse (1989) pour qui la division des erreurs orthographiques a été guidée par l'utilisation effective de copies d'élèves, ce qui permet de dégager des catégories comme: *mauvais découpage, mot sauté, mot tronqué*, ainsi que *substitution de mot*. On peut dès à présent y ajouter: *mot ajouté ou déplacé*; en effet, dans un corpus constitué de dictées réelles, ces catégories sont très utiles.

C'est aussi la seule grille qui différencie la catégorie *cumul des erreurs* et *simple faute* lexicale ou grammaticale. Il est vrai qu'on ne peut pas mettre sur le même plan les erreurs suivantes: **goutière* et **coutelière*.

On retiendra également de l'essai de typologie graphique de Millet (1989) une catégorie: *la variation scripto-graphique faisant appel à d'autres signes* (autres caractères, n'appartenant pas à l'alphabet). Dans notre corpus, on trouve en effet l'utilisation du chiffre arabe 1 qui remplace le pronom indéfini «un».

La grille orthographique de Catach (2003) est assez complexe car certaines catégories sont très difficiles à distinguer: par exemple, les erreurs à dominante idéogramme et les erreurs à dominante logogramme. En revanche, il est très intéressant de rassembler les erreurs morphogrammiques (lexicaux et grammaticaux), selon l'utilisation prévue de cette grille.

Plus pragmatique, la grille orthographique selon Burfin (1991) est celle communément utilisée par les professeurs. Il s'agit d'un découpage sans affinement et applicable pour la correction des copies, mais elle paraît trop rudimentaire pour notre travail.

La grille établie par Branca-Rosoff (1979) est très intéressante, mais nous ne l'avons pas utilisée car la distinction qu'elle opère entre la catégorie morphologie lexicale relevant d'un enseignement régulier et morphologie relevant d'un enseignement non régulier est très fragile. L'exercice de la dictée met en œuvre différentes règles ou exceptions vues pendant les cours. Il serait donc peu acceptable de demander une connaissance élaborée dans un cadre extrascolaire.

Une deuxième remarque porte sur la distinction entre morphologie grammaticale et chaîne d'accord. Dans ses exemples, elle cite: **il conclue* -> morphologie grammaticale (méconnaissance des flexions), et **il a chanter* -> chaîne d'accord grammatical (déclencheur non activé). La distinction de ces deux catégories ne semble pas être nécessaire en ce qui concerne notre corpus, étant donné que toutes ces catégories rentreront dans le domaine de la correction syntaxique (qui se trouve hors champs de notre étude).

Ces différentes grilles varient selon l'utilisation prévue. Dans le corpus qui nous intéresse et son application avec le logiciel de correction, une typologie légèrement différente s'impose:

- il nous faut tenir compte de tous les cas de figure du corpus brut (c'est-à-dire non saisi), pour éviter d'ajouter ou d'oublier certaines erreurs, que l'on pourra comparer avec la saisie informatique pour respecter au mieux le corpus brut;
- il nous faut distinguer les erreurs sur les plans morphologique et syntagmatique, qui relèvent de deux corrections différentes.

Nous proposerons la typologie suivante, avec, pour chaque type d'erreur retenue pertinente à notre corpus, une «étiquette», qui facilitera le traitement.

- *Plan calligraphique*
 - Erreurs sur tous les signes diacritiques:
 - **aîle*, pour aile(s) sera étiqueté + diac, montrant l'ajout du diacritique.
 - **a*, pour à sera étiqueté - diac, montrant la suppression du diacritique.
 - **á*, pour à sera étiqueté diac, montrant que le signe diacritique est erroné.
 - Erreurs sur les majuscules et minuscules (non comptabilisées).
 - Mauvaise troncation des mots en fin de ligne (non comptabilisées).
 - Mots ajoutés, sautés, déplacés (non comptabilisées).
 - Graphème non conventionnel: **en l vol*, pour l'article indéfini 'un' étiqueté GNC.
- *Plan morphologique*
 - Quant au phonème de base:
 - **envuit*, pour enfuit qui sera étiqueté: Pb (Aph) pour noter l'erreur du phonème de base qui entraîne une modification de sa valeur phonique (souvent la confusion se produit entre les sourdes et les sonores).
 - **ficsér*, pour fixer sera étiqueté: Pb (Sph) erreur quant au phonème de base sans changement phonique.
 - * *apprues*, pour apparues, noté: - Pb (Aph) absence d'un phonème de base avec une altération phonique.
 - **aussistôt*, pour aussitôt, noté: + Pb (Aph) ajout d'un phonème de base entraînant un changement phonique.
 - **elles sont apparues*, pour apparues, noté: + Pb (Sph) ajout d'un phonème de base, sans changement phonique.
 - Quant au graphème auxiliaire:
 - **travaller*, pour travailler sera étiqueté: - Gaux, notant l'absence du graphème auxiliaire.
 - Quant au graphème de position:
 - **au desus*, pour au-dessus: - Gp.
 - Quant au graphème à valeur phonique zéro:
 - **demein*, pour demain: Gz, erreur sur le graphème sans changement phonique.
 - **peti* pour petit: - Gz, absence du graphème zéro.
 - **demien* pour demain: Gz (Aph), interversion du graphème zéro entraînant une altération phonique.
 - **allors* pour alors: + Gz, ajout d'un graphème à valeur zéro sans changement phonique.
 - Quant aux digrammes:

- **revidras* pour *reviendras*: - dig, absence totale du digramme.
- **hirodelle* pour *hirondelle*: dig, digramme mal réalisé (absence d'un graphème).
- **ele* pour *aile(s)*: dig (Aph), erreur sur le digramme avec une altération phonique.
- **eile* pour *aile(s)*: dig (Sph), erreur sur le digramme sans altération phonique.
- **étaint*, pour *était*, noté: + dig (Aph) ajout d'un digramme, avec changement phonique.
- **prait*, pour *prêt*, noté: + dig (Sph) ajout d'un digramme, sans changement phonique.
- Intersion de deux graphèmes:
 - **geuttais* pour *guettais*, **premeire* pour *première*, noté IG: erreur sur la position de deux caractères d'un mot.
- Séquence exclue:
 - **innonbrables* pour *innombrables*, **penssais* pour *pensais*: SE, erreur nommée séquence exclue, quelquefois erreur dans le code.
- Substitution:
 - **ail*, **air* pour *ailles*: Sub, substitution d'un mot pour un autre.
- Mot à erreurs multiples:
 - **alloir* pour *alors*: Mmult, qui regroupe deux ou trois erreurs sur un seul et même mot: le détail de cette erreur sera noté ainsi (+ Gz; + Pb (Aph); - Gz).
- Mot trop éloigné:
 - **plers-soi* pour *aperçoit*: MTE, mot trop éloigné de la forme originale.
- Plan syntagmatique
 - Morphèmes grammaticaux:
 - **aperçois* pour *aperçoit*: Mgram (+), erreurs graphiques: confusion de la flexion du verbe avec une autre flexion existante.
 - **étai* pour *étais*: Mgram (-), erreur sur la flexion verbale.
 - Les homophones:
 - **champ* pour *chant*: Hom.
 - Erreur quant au genre:
 - **petit* pour *petite*: genre (Aph), erreur sur le genre avec un changement phonique.
 - **ami* pour *amie*: genre (Sph), erreur sur le genre sans changement phonique.
 - Erreur quant au nombre:
 - **elle sont* pour *elles sont*: nombre, erreur sans changement phonique.
 - La segmentation:
 - **aussi tôt* pour *aussitôt*: S, segmentation du mot.
 - L'agglutination:
 - **toutacoup* pour *tout à coup*: A, agglutination de plusieurs mot.

Le problème des homophones est le suivant: le principe de base des correcteurs est la comparaison avec un dictionnaire de base. Les homophones sont existants pour le correcteur puisqu'ils figurent dans ses listes. L'erreur est à cet endroit une erreur d'interprétation de sens, et les correcteurs, quel qu'ils soient, ne peuvent pas appliquer un traitement sémantique.

Cette grille va nous permettre un étiquetage des différentes erreurs relevées dans notre corpus de copies. Il sera donc possible de procéder à une confrontation plus fine avec le correcteur Cordial.

Les étiquettes sont les suivantes:

Catégorie globale *Étiquettes*

Diacritique	- diac
Diacritique	+ diac
Diacritique	diac
Digamme	- dig
Digamme	+ dig (Aph)
Digamme	+ dig (Sph)
Digamme	dig
Digamme	dig (Aph)
Digamme	dig (Sph)
Erreur simple - Gaux	
Erreur simple - Gp	
Erreur simple GNC	
Erreur simple IG	
Erreur simple SE	
Erreurs multiples	Mmult
Erreurs multiples	MTE
Flexion	genre (Aph)
Flexion	genre (Sph)
Flexion	Mgram (+)
Flexion	Mgram (-)
Flexion	nombre
Graphème zéro	- Gz
Graphème zéro	+ Gz
Graphème zéro	Gz
Graphème zéro	Gz (Aph)
Phonème de base	- Pb (Aph)
Phonème de base	+ Pb (Aph)
Phonème de base	+ Pb (Sph)
Phonème de base	Pb (Aph)
Phonème de base	Pb (Sph)
Remplacement du mot	Hom
Remplacement du mot	Sub
Unité du mot	A
Unité du mot	S

Toutes ces typologies sont conçues avec des objectifs différents, et sont donc très diverses. L'objectif définit la typologie. Notre objectif était de trouver un étiquetage d'erreur le plus précis possible afin de bien dégager les endroits où le correcteur avait des difficultés.

Fonctionnement du correcteur informatisé

Les correcteurs informatisés utilisent différentes techniques pour rechercher les erreurs potentielles. Ces différentes techniques sont bien synthétisées par Véronis (1988a):

- «Wagner et Fisher utilisent quatre "opérations d'édition": la substitution, l'insertion, la suppression et l'interversion d'une lettre, qui servent de reconstitution du mot originel.
- «Une autre technique utilisée est la comparaison dite séquentielle; il s'agit de comparer le mot saisi avec d'autres mots de noyaux communs. Pour cela, Damereau a conçu son algorithme de la façon suivante: deux lettres interverties, la lettre suivante est omise, ajoutée, ou substituée.

- «Un gain de temps s'opère lorsque le dictionnaire enregistré servant de référent n'est pas parcouru dans sa totalité; en effet, on peut prédire la longueur minimum et maximum de la chaîne graphique à comparer: sa longueur comprendra une omission Ix-1I, une substitution ou inversion IxI et au maximum une insertion Ix+1I par rapport au mot saisi.
- «Pollock et Zamora utilisent une méthode nommée "les clés de similarité"; il s'agit de retenir les cinquante mots qui diffèrent du mot saisi que d'une seule lettre. Une variante appelée "les clés insensibles" proposent de ne pas tenir compte de l'ordre des lettres.»

Il faut, dans toutes ses méthodes, se rendre compte que la qualité résulte de deux principaux facteurs: le pourcentage de réussite, bien sûr, mais aussi du temps qu'il faut pour mettre en œuvre ses différentes recherches; c'est pourquoi, lorsqu'il s'agit d'opérer une correction sur l'écrit d'un apprenant, tous ces paramètres gênent la détection, et ceci montre qu'ils ne sont pas conçus pour ce public. En effet, la typologie des erreurs chez un élève ne sera pas seulement d'ordre typographique.

Une expérience rend compte de la difficulté de cette tâche (Véronis, 1988b): il s'agit de transcrire phonétiquement plus de trois mille mots fréquents automatiquement. Cette expérience a montré qu'un mot sur deux possédait une graphie complexe.

D'autres méthodes ont vu le jour, acceptant dans un même mot plusieurs fautes phonographiques et éliminant l'erreur en initiale; par exemple «Soundex Code», créé par Odell et Russel; mais un inconvénient apparaît: la base de référence est réduite, pour une efficacité optimum, entre 500 et 1 000 mots¹².

Le logiciel qui nous concerne (Cordial 6) fonctionne de la manière suivante¹³:

- Le redoublement superflu d'une lettre: **allors*, le correcteur Cordial va proposer trois propositions: «alors» où il traite en priorité l'ajout d'une lettre, «aller» qui suggère que le redoublement de la consonne est correct, et que la faute n'est qu'une substitution de la voyelle «e» avec un ajout de la lettre «s»; et «allers» où il estime que le nombre de lettres est conforme, et qu'il n'existe qu'une substitution. Dans cet exemple, le correcteur admet la possibilité de plusieurs erreurs différentes. Mais ce n'est pas toujours le cas.
- L'omission d'une lettre: **alor*, Cordial proposera *alors* en restituant la lettre manquante. Pourquoi ne suggère-t-il pas aller? Le correcteur compare d'abord les mots de la même taille, et ne passe à une comparaison de la base de donnée avec des mots de la taille X + 1, que lorsqu'il n'a pas trouvé un mot approchant.
- La substitution d'une lettre: **alorr*, le correcteur proposera une seule proposition *alors*; tandis que pour **alore* c'est *allure* qui est proposée. Pourquoi? Dans l'exemple **alorr* la lettre substituée est une consonne, le correcteur cherchera à remplacer la lettre «r» par une autre consonne, étant donné que le redoublement d'un «r» final n'existe pas en français. Or, dans le second exemple **alore*, la syllabe «re» n'est pas irréaliste. Le correcteur examine la possibilité que le «o» soit la lettre erronée.

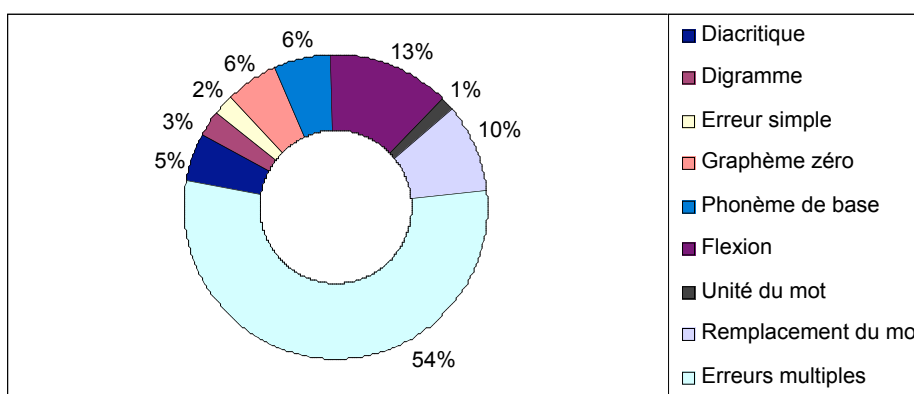
Résultats

Performance des élèves

Comme on pouvait s'y attendre, les élèves ont commis des fautes sur la dictée présentée. Celle-ci contient un total de 126 mots (sans les doublons). Pour les 100 élèves de CM2, le total d'erreurs est de 3 421, soit un pourcentage de 20%; car les erreurs n'ont pas fait l'objet d'une hiérarchie (comme le ferait un professeur, en distinguant par exemple, faute d'accord ou faute d'usage). Chaque variante est comptabilisée, et considérée sur le même plan. La première constatation est d'observer le petit nombre d'erreur. Ceci est dû au caractère uniforme des valeurs imputées aux variantes. De plus, certains mots¹⁴ ont bénéficié d'un pourcentage nul d'erreur, qui a eu pour effet de rehausser le pourcentage total.

Ces mots sont les suivants: *matin, ne*. Beaucoup de mots ont un pourcentage faible d'erreur (souvent inférieur à 3 ou 4%), qui a aussi contribué à ce faible taux. Mais, la majorité des mots erronés obtient le «soutien» de la quasi-totalité des élèves: *innombrables* remporte 96% des erreurs.

Voici la répartition des différents types de fautes commises par les élèves, en grande catégorie.



Cette division nous montre bien toute la difficulté de corriger un tel corpus. La majorité des erreurs est constituée par trois catégories d'erreurs pour les élèves de CM2:

- erreurs multiples (*Mmult*: deux ou trois erreurs cumulées sur un même mot, et *MTE*: mot trop éloigné avec plus de trois erreurs);
- la catégorie *Flexion*, constituée des erreurs de flexion verbale existante ou non (*Mgram (+)* et *Mgram (-)*), des erreurs sur l'accord en genre et en nombre, avec ou sans changement phonique (*genre (Aph)*, *genre (Sph)* et *nombre*);
- et *Remplacement de mot*, constituée des erreurs sur les homophones: *Hom*, et les erreurs de substitutions: *Sub*.

Ces trois erreurs qui représentent la majorité sont des erreurs indétectables pour Cordial. D'une part, Cordial considère la possibilité d'une seule erreur dans un mot (première catégorie); d'autre part, les mots existant sur la base de données de Cordial sont considérés automatiquement comme juste (catégories 2 et 3).

Les performances brutes du correcteur Cordial 6

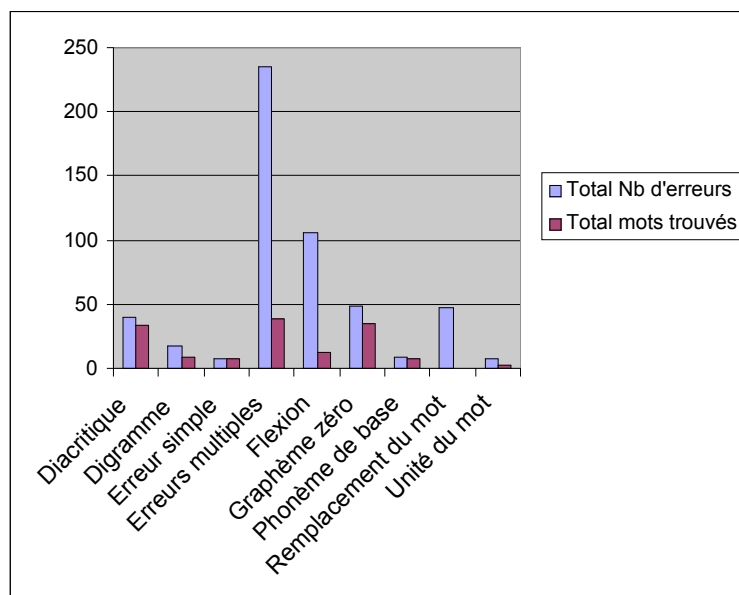
Le correcteur souligne des mots qu'il considère comme erronés, puis propose plusieurs autres orthographes avec une hiérarchie de son point de vue.

Cordial trouve un total de 1 235 mots erronés (sans prendre en considération sa position parmi ses suggestions), et 690 mots sont corrigés, soit un total de 20%; mais à quelle position se trouvent les bonnes propositions de Cordial? La réponse à cette question est intéressante; en effet, si l'on se place dans une situation didactique, la position peut influencer le choix de l'élève.

Parmi toutes les propositions de Cordial, les bonnes graphies sont situées en majorité de la première à la troisième position. Cette répartition est rassurante. Les positions ne varient pas selon la variation graphique, sauf une légère augmentation pour la septième position.

Les moyennes globales tenant compte des doublons nous montrent que Cordial propose en moyenne 2,47 mots dont la bonne graphie se trouve en 0,6^e position.

Examinons ci-dessous la proportion entre le nombre de mots corrigés par Cordial.



Les erreurs les plus fréquentes sont les erreurs multiples, qui sont aussi les moins corrigées.

En revanche, les erreurs sur les diacritiques, les digrammes, les graphèmes à valeur zéro, les erreurs sur les phonèmes de base et les erreurs simples ont une majorité de mots corrigés. L'absence totale de correction est flagrante pour les erreurs consistant à remplacer le mot.

Le tableau ci-dessous reprend l'ensemble des catégories et synthétise les résultats bruts de Cordial.

Étiquettes	Nb d'erreurs	Nb de corrections	%
- diac	14		857,14
- dig	4		00
- Gaux	1		1100
- Gp	4		4100
- Gz	28		22 78,57
- Pb (Aph)	8		688,8
+ diac	9		888,88
+ dig (Aph)	4		125
+ dig (Sph)	1		1100
+ Gz	16		956,25
+ Pb (Aph)	22		13 59,09
+Pb (Sph)	/	/	/
A	5		5100
diac	23		20 86,69
dig	8		787,5
dig (Aph)	4		125
dig (Sph)	8		562,5
genre (Aph)	3		00
genre (Sph)	1		00
GNC	2		00
Gz	18		11 61,11
Gz (Aph)	6		350
Hom	40		00
IG	3		3100
Mgram (+)	85		33,52
Mgram (-)	21		523,8
Mmult	348		69 5,48
MTE	164		95,48
nombre	11		00
Pb (Aph)	14		750
Pb (Sph)	7		457,14
S	8		112,5
SE	5		5100
Sub	50		00

Les erreurs les mieux corrigées automatiquement correspondent aux types d'erreur pris en compte dans les algorithmes de Cordial, à savoir les erreurs sur les diacritiques, sur un graphème, sur les séquences exclues et les agglutinations simples (c'est-à-dire sans erreur à l'intérieur de l'agglutination) et, enfin, les segmentations.

En revanche, les erreurs qui ne sont pas corrigées sont des erreurs où la détection se fait par référence au contexte: le genre et le nombre, les flexions verbales, les homophones, les substitutions. Le domaine des erreurs non corrigées, les mots à graphie trop éloignés et les mots contenant deux ou trois erreurs rentrent aussi dans cette catégorie. Mais la raison est différente: le cumul d'erreurs sur un seul et même mot.

Pour les autres erreurs, le pourcentage d'erreur varie de 45,5% à 66,6%. Ces proportions sont plutôt encourageantes.

Ce corpus met l'accent sur des erreurs non connues par le correcteur. Les élèves commettent des fautes de types très divers. Le palmarès des erreurs revient aux erreurs cumulées dans un seul et même mot. Erreur que Cordial ne peut corriger à l'heure actuelle.

Pour pallier ce manque, il faudrait que Cordial se penche sur ces différentes erreurs et résolve le problème lié à l'accumulation des erreurs sur un seul mot.

Regard sur le correcteur syntaxique

Le correcteur syntaxique de Cordial réagit de trois façons:

- il propose une correction qui se trouve être juste;
- il propose une correction, mais erronée;
- il n'avance qu'une suggestion sans correction (qui peut être juste ou fausse).

Pour la correction des 100 copies des élèves de CM2, on constate une correction (proposition juste) de 340 mots. Cette correction mériterait une étude plus minutieuse, comme celle que nous avons menée sur le correcteur orthographique. Nous nous limiterons à l'évaluation de la bonne correction de certaines expressions. Cette deuxième correction permet pour les élèves de CM2 d'enlever exactement 9,93% de leurs erreurs. Une recherche plus précise sur la réaction des élèves aux propositions du correcteur pourrait être entreprise pour mesurer l'impact des propositions fausses choisies par les élèves sur le total de fautes commises après vérification à l'aide de Cordial.

Pistes d'amélioration du correcteur

La solution d'amélioration du correcteur pour une situation d'enseignement serait de prendre en considération, dans la conception du programme, les erreurs dites multiples, c'est-à-dire l'accumulation de plusieurs erreurs par mot. Pour cela, l'utilisation des valeurs des graphèmes peut être une bonne piste de recherche pour un éventuel travail de standardisation à des fins de traitement automatique.

Partant des observations de notre corpus, la question est: comment formaliser la base d'un mot? La solution réside peut-être dans la combinaison des différents graphèmes.

Un mot est constitué d'une succession de graphèmes divisibles en deux ensembles finis:

1. les voyelles (a, e, i, o, u, y);
2. les consonnes (le reste de l'alphabet + y).

Plusieurs combinaisons sont acceptées en français:

1. C + V + C + V = P + A + P + A = "papa"
2. CC + V + C = C + H + A + T = "chat"
3. C + VV = F + E + U = "feu"

Les consonnes peuvent se grouper au maximum en trois graphèmes: -thr-, -chr-, etc.

Les voyelles forment aussi des groupes de trois voyelles au maximum: -eau-, -oua-, etc.

Ces deux types de formes se combinent entre eux, et forment tout le lexique français.

Si ces deux ensembles finis sont doublés d'une valeur possible, comme la typologie d'erreur tentait de faire, il est possible que la correction automatique des erreurs orthographiques soit plus adaptée.

Par exemple, pour les erreurs suivantes (qui n'ont pas été trouvées par Cordial), plusieurs automatismes sont mis en œuvre:

**irrodèle*: (pour hirondelle)

1. la possibilité de l'absence du graphème muet "h" devant les mots commençant par une voyelle;
2. doublement superflu d'une lettre;

3. possibilité d'un digramme mal réalisé (en ce qui concerne les nasales, il s'agit d'une liste finie: AN, EN, IN, ON, UN);
4. diacritique masquant la possibilité d'un doublement de consonne.
**desordones*: (pour désordonné)
1. possibilité de l'absence d'un signe diacritique;
2. absence d'un graphème à valeur nulle ou zéro;
3. possibilité d'une erreur sur la flexion verbale.

Pour chacune des fautes, il est possible de disposer d'un ensemble fini de propositions qui demanderait plus de temps, mais serait plus efficace pour ce public.

Conclusion

En laissant agir Cordial de manière automatique sur le corpus, en cumulant correction orthographique et syntaxique, 30% des erreurs disparaîtraient. Ce résultat est de notre point de vue déjà très encourageant. En effet, si on imagine qu'une didactique de l'utilisation du correcteur pourrait être mise en place (comme c'est le cas aujourd'hui pour la calculette), l'interaction intelligente de l'élève avec les propositions de la machine devrait sensiblement augmenter le nombre d'erreurs corrigées. Cette interaction pourrait même être conçue comme le noyau d'un apprentissage de l'orthographe. En effet, dans les traitements de texte, et singulièrement dans celui qui est majoritairement implanté sur les machines et qui fonctionne avec le noyau de Cordial (Microsoft Word), le paramétrage du correcteur permet de renvoyer à la règle qui précise la défaillance détectée. Par ailleurs, Cordial n'est pas conçu pour travailler sur des corpus comportant des erreurs qui manifestent des incompétences orthographiques majeures, comme c'est le cas pour nombre d'élèves en fin de cursus de l'école élémentaire; on l'a vu en particulier pour le cumul d'erreurs sur un même mot. Des aménagements importants, tant sur le plan des algorithmes que sur le plan de l'interface utilisateur, devraient être entrepris pour rendre les correcteurs à la fois plus efficaces et plus «didactiques» pour un public d'apprenants. Comme bien d'autres artefacts (calculettes, messagerie électronique, «texto» des téléphones, etc.), les traitements de texte assortis de correcteurs entrent dans le quotidien des élèves avec un temps d'avance sur l'école. Les modifications qu'ils entraînent dans la pratique de l'écriture devront être tôt ou tard considérées avec toute l'attention qu'elles méritent, en particulier par les didacticiens qui pourraient, de notre point de vue, prendre une part active dans leur amélioration¹⁵... à fins pédagogiques.

NOTES

1. Par exemple, le concours «les dicos d'or», organisés dans l'éducation nationale en partenariat avec les médias et des sponsors privés, qui connaît un succès sans cesse croissant.
2. Les réticences à appliquer les recommandations orthographiques issues de la réforme de 1990 sont telles qu'une association s'est formée pour promouvoir son application.
APARO: <http://www.filtr.ucl.ac.be/FLTR/TOM/ess.html>
On trouve même sur le sujet des ouvrages savants mais polémiques (Traimond, 2001) ou, à la limite, xénophobes (Caffot, 2001).
3. En fait, la manière dont on apprend l'orthographe a fait l'objet de nombreuses modélisations, parfois concurrentes, parfois complémentaires, mais qui n'ont pas été unifiées pour l'instant, même si elles mettent toutes en évidence les difficultés importantes qui résultent de la non-correspondance phonie-graphie en français. Une analyse synthétique des recherches est proposée par Fayol et Jaffré (1999).
4. On notera avec Jacquet-Pfau (2000) que les correcteurs – eux – ont évolué avec la réforme.

-
5. On peut d'autant s'en étonner que la documentation à ce propos ne manque pas et propose de nombreuses pistes: Depover, Giardina & Marton (1998), Linard (1997), Tardif (1998), par exemple.
 6. Le texte de son institution précise: «Garantir à tous une maîtrise raisonnée de ces outils» – «Intégrer les TICE au service de l'ensemble des activités des élèves et des enseignants». Bulletin officiel de l'Éducation nationale n° 42 du 23.11.2000. Le ministre Lang précisera en 2002: «Les technologies de l'information et de la communication ne s'organisent pas en une discipline autonome. Ce sont des outils au service des diverses activités scolaires, dont l'appropriation active conduit au premier niveau du Brevet informatique et internet (B2i)». Les nouveaux programmes, *Qu'apprend-on à l'école élémentaire?*, février 2002 (p. 49).
 7. Sur cette question des rapports entre TIC et école, on peut utilement se référer à Papadoudi (2000), au rapport de l'inspection générale de l'Éducation nationale, *L'école et les réseaux numériques*, n° 2002-035, juillet 2002, mais aussi à Charlier et Peraya (2003) et à la thèse très complète d'Audran (2001).
 8. Cette préoccupation est déjà bien explicitée par Désilets (1997 et 1998). Bertin (2003) pointe avec méthode les obstacles à l'utilisation des correcteurs par des élèves, mais se montre trop pessimiste sur l'avenir, de notre point de vue.
 9. D'autres correcteurs indépendants existent sur le marché, mais le choix de Cordial s'est imposé en raison de sa large diffusion sous Word. On peut citer: Correcteur 101 (de Machina Sapiens), Antidote (de Druide Informatique), Hugo plus (de Logidisque), Prolexis (de Diagonal).
On trouvera des comparatifs sur:
<http://www.osil.ch/eval> et sur <http://www.univ-tlse2.fr/gril/TAL/Chk>
 10. La dictée a été proposée également à 94 élèves de sixième (première secondaire). Même si le traitement des résultats a permis de faire un certain nombre de remarques sur les différences entre CM2 et sixième tant sur le plan des performances des élèves que sur le plan des types d'erreurs commises, celles-ci ne feront pas partie des préoccupations centrales de cet article.
 11. Cette recherche a donné lieu à la rédaction d'un mémoire de DEA par Marie-Pierre Barthel à l'Université de Provence, sous la direction du Pr Véronis, et pour laquelle nous avons fourni toutes les données.
 12. Cet algorithme breveté est décrit dans Knuth (1998).
 13. Page Web pour plus de détails: <http://www.synapse-fr.com/CP03.htm>
 14. La notion de mot est prise dans son sens le plus large: graphie séparée d'espaces.
 15. Citons la tentative de Desmarais (1994), qui semble avoir été peu suivie.

RÉFÉRENCES

- Allal, L., Bétrix Köhler, D., Rieben, L., Rouiller, Y., Saada-Robert, M., & Wegmuller E. (2001). *Apprendre l'orthographe en produisant des textes*. Programme national de recherche 33. Efficacité de nos systèmes de formation. Fribourg: Éditions universitaires.
- Audran, J. (2001). *Influences réciproques relatives à l'usage des nouvelles technologies de l'information et de la communication par les acteurs de l'école. Le cas des sites Web des écoles primaires françaises*. Thèse NR, Université Aix-Marseille 1.
- Baron, G.L., & Bruillard, E. (1996). *L'informatique et ses usagers dans l'éducation*. Paris: PUF.
- Berten, F. (2003). *Rapport de la Commission «français et informatique»*. Bruxelles: Secrétariat général de l'enseignement catholique. <http://users.skynet.be/ameurant/francinfo/correcteur/correcteur.html>
- Branca-Rosoff, S., & Piolat, M. (1979). L'acquisition de l'orthographe, état d'une recherche. *Recherches sur l'acquisition de l'orthographe* (pp. 3-29). Montréal: Université du Québec à Montréal.
- Brunot, F., & Gossot, H. (1957). *De la théorie... à la pratique*. Paris: Istra.
- Burfín, R. (1991). *Le français pour tous: écrire sans fautes*. Lyon: Chroniques sociales.
- Burney, P. (1955). *L'orthographe*. Paris: PUF.
- Caffot, P. (2001). *Défense et illustration de l'orthographe française*. Paris: Godefroy de Bouillon.
- Catach, N. (1991). *L'orthographe en débat: dossier pour un changement*. Paris: Nathan.
- Catach, N. (2003). *L'orthographe* (9^e édition). Paris: PUF.
- Charlier, B., & Peraya, D. (2003). *Technologie et innovation en pédagogie. Dispositifs innovants de formation pour l'enseignement supérieur*. Perspectives en éducation et formation. Bruxelles: De Boeck.
- Charnet C., & Panckhurst, R. (1998). Le correcteur grammatical: un auxiliaire efficace pour l'enseignant? Quelques éléments de réflexion. *Revue Alsic*, 1(2).
<http://alsic.u-strasbg.fr>.
- Chervel, A., & Manesse, D. (1989). *La dictée: les français et l'orthographe*. Paris: INRP, Calman-Levy.
- Clegg, C. (1988). The human side of advanced manufacturing technology. In D. Toby, C. Clegg, & J. Nigel (éds), *Chichester*. New York: Wiley.
- CNDP (1992). *Les cycles à l'école primaire*. Paris: Hachette.
- Depover, C., Giardina, M., & Marton, P. (1998). *Les environnements d'apprentissage multimédia*. Paris: L'Harmattan.
- Désilets, M. (1997). Que penser de l'utilisation des logiciels correcteurs à l'école? *Vie pédagogique*, 107.
- Désilets, M. (1998). Les correcteurs d'orthographe et l'apprentissage. *Les Cahiers pédagogiques*, 362.
- Desmarais, L. (1994). *Proposition d'une didactique de l'orthographe ayant recours au correcteur orthographique*. Mémoire de recherche, Université Laval.
- Fayol, M., & Jaffré, J.P. (éds) (1999). L'acquisition/apprentissage de l'orthographe: *Revue française de pédagogie*, 126, Paris: INRP.
- Gey, M. (1987). *Didactique de l'orthographe française*. Paris: Nathan.
- Goosse, A. (1991). *La nouvelle orthographe*. Paris: Ducolot.
- Jacquet-Pfau, C. (2000). Les correcteurs orthographiques et grammaticaux: fonctionnement et typologie. In *Les dictionnaires de langue française au sein de la francophonie: normes et orthographe d'hier à aujourd'hui*. La Journée des Dictionnaires, Colloque international, Université de Cergy-Pontoise, 17 mars 1999, Klincksieck.
- Knuth, D. (1998). *The art of computer programming* (vol. 3). Reading (Mass.): Addison-Wesley.
- Linard, C. (1997). *Des machines et des hommes*. Paris: L'Harmattan.
- Millet, A. (1989). *Quelques aspects sociolinguistiques de l'orthographe française*. Thèse de doctorat NR, Université de Grenoble 3.
- Millet, A., Lucci, V., & Billiez, J. (1990). *Orthographe mon amour!* Grenoble: PUG.
- Norman, D., & Draper, S. (1986). *User centered system design: new perspectives on human-computer interaction*. New York: Hillsdale ou Erlbaum Associates.
- Papadoudi, H. (2000). *Technologies et éducation. Contribution à l'analyse des politiques publiques*. Paris: PUF.
- Piéron, H. (1969). *Examens et docimologie*. Paris: PUF.
- Rabardel, P. (1995). *Les hommes et les technologies: approche cognitive des instruments contemporains*. Paris: Colin.
- Ravestain, J. (1999). *Autonomie de l'élève et régulation du système didactique*. Bruxelles: De Boeck.
- Tardif, J. (1998). *Intégrer les nouvelles technologies de l'information, quel cadre pédagogique?* Paris: ESF.

- Thimonnier, R. (1976). *Le système graphique du français: introduction à une pédagogie rationnelle de l'orthographe*. Paris: Plon.
- Traimond, B. (2001). *Une cause nationale: l'orthographe française: éloge de l'inconstance*. Paris: PUF.
- Véronis, J. (1988a). *Contribution à l'étude de l'erreur homme-machine en langage naturel*. Thèse soutenue le 5 octobre 1988, Université Aix-Marseille II.
- Véronis, J. (1988b). From sound to spelling in French: Simulation on a computer. *Cahiers de psychologie cognitives, European Bulletin of Cognitive Psychology*, 8(4).