



HAL
open science

Variation graphique dans les documents d’Ancien Régime : Nouvelles approches scriptométriques

Simon Gabay, Philippe Gambette, Rachel Bawden, Jonathan Poinhos, Eleni Kogkitsidou, Benoît Sagot

► To cite this version:

Simon Gabay, Philippe Gambette, Rachel Bawden, Jonathan Poinhos, Eleni Kogkitsidou, et al.. Variation graphique dans les documents d’Ancien Régime : Nouvelles approches scriptométriques. Journée d’étude : “ Pour une histoire de la langue ‘par en bas’: textes privés et variation des langues dans le passé ”, Sep 2021, Paris, France. hal-03357080

HAL Id: hal-03357080

<https://inria.hal.science/hal-03357080>

Submitted on 14 Nov 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L’archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d’enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Variation graphique dans les documents d'Ancien Régime

Nouvelles approches scriptométrique

Simon Gabay¹ Philippe Gambette² Rachel Bawden³ Eleni Kogkitsidou²
Jonathan Poinhos² Benoît Sagot³

¹Université de Genève (Suisse) – prenom.nom@unige.ch

²LIGM, Université Gustave Eiffel (France) – prenom.nom@u-pem.fr

³INRIA (France) – prenom.nom@inria.fr

Au commencement était la marquise

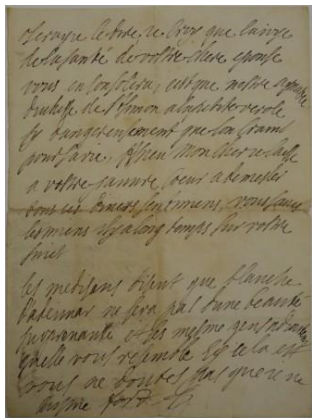


Figure 1 – Lettre de Sévigné/Grignan, MS Lowell Autograph File 282, f°3r.

- ▶ Tout a commencé avec la correspondance de la marquise
- ▶ Problème pour différencier les parties écrites par sa mère et sa fille
- ▶ Usage de l'accent, distinction <i>/<j> et <u>/<v>, usage de l'apostrophe, usage de la ponctuation...
- ▶ Très peu de travaux sur le français « classique »

Plan de la présentation

Des *scriptae* en français classique

De l'utilité de la scriptométrie

Approche par règles

Approche par réseau de neurones

Conclusion

Des *scriptae* en français classique

Des *scriptae* françaises ?

Deux grands courants orthographiques (Pellat 1991)

- Courant de l'« orthographe ancienne », apparu vers le XIII^e s., qui cherche à maintenir le lien avec le latin
- Courant de l'« orthographe nouvelle », qui cherche un système propre au français (et plus ancien que l'orthographe ancienne)

Réalisation de cette opposition en deux *scriptae* au XVII^e s. :

- Par analogie avec les *scriptae* romanes, dérivées spontanément du latin et organisées diatopiquement (wallon, picard, anglo-normand)
- Nous proposons de parler de *scriptae* françaises, dérivées de la norme médiévale et organisées diastratiquement sur le modèle de la Querelle des Anciens et des Modernes

[J.-Ch. Pellat, « Norme et variation orthographique au XVII^e s. », *Rencontres linguistiques en pays rhénan* 5/6, 1994, pp. 245-260.]

Exemple de Bossuet : théorie

Le choix du système graphique est théorisé par nombre d'auteurs, et influencé par leur position dans la Querelle.

[La compaignie] ne peut souffrir une fausse regle qu'on a uoulu introduire, d'écrire comme on prononce, parce qu'en uoulant instruire les étrangers et leur faciliter la prononciation de nostre langue, on la fait meconnoistre aux François mesmes. Si on escrivoit tans, chan, cham, emais ou émés, connoissans, anterremans, faisaict, qui reconnoistroit ces mots? [...] Il y a aussi une autre ortographe, qui s'attache scrupuleusement a toutes les lettres tirées des langues dont la nostre a pris ses mots, et qui ueut écrire nuict, ecripture, etc. Cella blesse les yeux d'une autre sorte en leur remettant en ueüe des lettres dont ils sont desaccoustumez et que l'oreille n'a iamais connus. [Bossuet, Cahiers de Mezeray]

[L. Biederman-Pasques, *Les Grands Courants orthographiques au XVII^e s. et la formation de l'orthographe moderne*, Tübingen : Max Niemeyer Verlag, 1992]

scripta des Anciens vs *scripta* des Modernes

-	Anciens	Modernes
Lettres ramistes	Position (<i>vniuers</i>)	Dissimilation (<i>univ<u>er</u>s</i>)
Ancien hiatus	<i>eu</i> (<i>veu</i>)	<i>u</i> (<i>vu</i>) ou <i>û</i> (<i>vû</i>)
Lettre calligraphique	maintenue (<i>ay</i>)	supprimée (<i>ai</i>)
Pluriel nominal	-s, -z, -x	-s
Consonne muette (diac., étym, hist.)	maintenue (<i>doubte</i>)	supprimée (<i>doute</i>)
Voyelles longues	diacritique (<i>teste</i>) dédoublément (<i>aage</i>)	circonflexe (<i>tête</i>)
...

Table 1 – Exemples d'opposition entre graphie ancienne et graphie moderne

- le système alphabétique, avec les couples <i>/<j> et <u>/<v> étant conçu comme des variantes graphiques ou comme renvoyant à des phonèmes différents
- les logogrammes lexicaux, avec l'utilisation ou la suppression de lettres en surcharge, qu'elles soient historiques (*huile*) ou étymologiques (*doubter*)
- les signes auxiliaires, avec l'utilisation de lettres diacritiques (*hospital*) plutôt que d'accents ou de trémas.

Exemple de Bossuet : pratique

Opposition relativement nette entre :

- Bossuet moderne, avant la fin de ses études (doctorat en 1652)
- Bossuet ancien, après avoir rejoint les ordres en Lorraine (1654)

Manuscrit	tant	être	cette	même	paraître	avec	a-t-il
BNF Fr. 12822, f°370 <i>Brièveté de la vie</i> , 1648	tans	estre	ceste	mesme	parêtre	auèque	a t'il
BNF Fr. 12823, f°130 <i>Fête du Rosaire</i> , oct. 1651	tans	être	cête	même	parêtra	auéc	a t'il
BNF Fr. 12824, f°119 <i>Sur la Prov.</i> , mai 1656	temps	estre	ceste	mesme	paroist	auéc/ auèque	a til
BNF Fr. 12822, f°61 <i>Sur les démons</i> , fév. 1660	temps	estre	cette	mesme	paroistre	-	atil

Table 2 – Exemples de l'évolution du système graphique de Bossuet

→ Grande richesse d'une analyse graphématique du point de vue linguistique, mais aussi littéraire et historique

Plan de la présentation

Des *scriptae* en français classique

De l'utilité de la scriptométrie

Approche par règles

Approche par réseau de neurones

Conclusion

De l'utilité de la scriptométrie

Pourquoi une approche computationnelle ?

Quelle est l'utilité d'une approche computationnelle ?

- ▶ L'apparition et l'accélération des besoins numériques nécessite un outillage de qualité
 - La qualité des outils (lemmatisation, annotations diverses...) dépend de notre connaissance de la langue
 - Mais, à rebours, les outils peuvent aussi servir à améliorer notre connaissance de la langue
- ▶ Il existe un besoin de produire des textes normalisés (lecture, requête dans un corpus...)
 - La normalisation repose sur des règles (computationnelles) de transformation
 - Il devrait être possible d'étudier l'application de ces règles pour analyser la langue des textes

Ces approches computationnelles « *data-driven* » nécessitent cependant la création d'au moins deux corpus :

- Un corpus « gold » d'entraînement et d'évaluation des modèles de normalisation, proposant le même corpus dans une version normalisée ou non. → PARALLEL17
- Un corpus de recherche, sur lequel appliquer les outils développés avec le corpus « gold ». → D'ALEMBERT

Parallel17

L'approche la plus logique est celle d'un apprentissage supervisé : l'ordinateur doit déduire des règles à partir d'exemples. À cette fin, nous avons conçu le corpus PARALLEL17 (plus de 650 000 tokens)

- ▶ On trouve une cinquantaine de textes du XVII^e s., parfois complets (*Andromaque* de Racine), parfois non (*Astrée* d'Urfé).
- ▶ le corpus se veut représentatif, et contient du vers comme de la prose, différents genres (surtout littéraires), est raisonnablement distribué diachroniquement par décennie...
- ▶ Chaque texte est découpé en segments constitués de (sous-)phrases (délimiteurs retenus : < ; >, < : >, < ? >, < ! >, < . >)
- ▶ Chaque segment offre une transcription diplomatique et sa version normalisée, alignée sur le français contemporain à quelques exceptions près (*i.e.* licence poétique comme l'élision pour raison métrique).

En plus du corpus d'entraînement, il nous faut un corpus de recherche

- ▶ Nous sommes confrontés à un problème important : celui de la normalisation quasi intégrale des éditions de textes classiques
- ▶ Il existe des corpus plus ou moins accessibles (et plutôt moins que plus), du type du RCFC, pas forcément utilisables pour des analyses graphématiques (type corpus Vachon, amputé pour le XVII^e s.)
- ▶ L'apparition de nouvelles technologies de numérisation et d'analyse amène l'utilisation de corpus de très grande taille

Nous nous sommes donc lancés dans la constitution d'une super corpus du français moderne

- ▶ Couvre la période XVI^e s.-XVIII^e s.
- ▶ Avec des données riches et harmonisées (auteur, date, genre, normalisation...)
- ▶ Nous avoisinons les 300 Mo de texte, mais il devrait continuer de grandir
- ▶ Problème de distribution des données et de publication sous forme de ressource

Ce corpus devrait permettre l'entraînement de modèles de langue dont l'objectif est de venir en soutien aux tâches de TAL, notamment la normalisation

Plan de la présentation

Des *scriptae* en français classique

De l'utilité de la scriptométrie

Approche par règles

Approche par réseau de neurones

Conclusion

Approche par règles

Détection automatique de *scripte*

Approche fondée sur l'alignement de corpus parallèles

- ▶ Analyse scriptologique intégrée à la tâche de normalisation linguistique
- ▶ Utilisation du corpus parallèle PARALLEL17 comme base d'apprentissage
- ▶ Développement de la méthode ABA (*Alignment-Based Approach*) :
 - alignement au niveau des mots pour chaque (sous-)phrase de PARALLEL17
 - alignement au niveau des caractères pour chaque mot
 - pour chaque caractère modifié entre la version originale et normalisée, détection de la règle appliquée
 - Utilisation de ces règles pour normaliser un texte
 - La règle utilisée est sauvegardée pour une analyse linguistique

Alignement avec l'algorithme de Needleman-Wunsch

L'algorithme de Needleman-Wunsch est un algorithme qui effectue un alignement global maximal de deux chaînes de caractères. On utilise pour cela une matrice de similarité

	A	p	o	f	t	r	e
A	■						
p		■					
ô			■	■			
t					■		
r						■	
e							■

- similarité
- substitution
- suppression
- insertion

Table 3 – Matrice de similarité pour la version originale et normalisée d'*Apofre*.

Les règles dans la méthode ABA

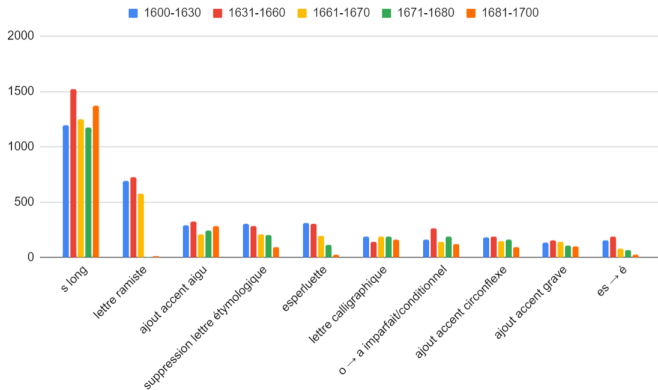
Les transformations (substitution, suppression, insertion) sont regroupées selon des grands types en s'appuyant sur la littérature existante (par ex. Vachon 2010).

Exemples :

- ▶ lettre ramiste : un *i*, un *j*, un *u*, un *v* dans le mot en version originale associé respectivement à un *j*, un *i*, un *v*, un *u* en version normalisée (*vniuers*→*univers*) ;
- ▶ lettre calligraphique : un *y* en fin de mot en version originale est associé à un *i* en version normalisée (*roy*→*roi*) ;
- ▶ lettre diacritique : *s* dans *hospital estoit* remplacé par un accent circonflexe (*hôpital*) ou aigu (*était*).

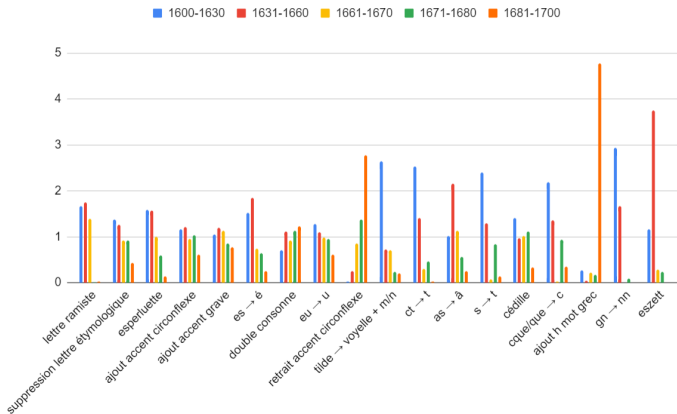
Résultats sur le corpus découpé en 5 périodes

Nombre d'occurrences de chaque règle pour 10 000 mots



Évolutions significatives au cours du 17^e siècle

Règles présentant de fortes évolutions sur la période (1 = fréquence de référence pour le 17^e siècle)



Une valeur de 1 correspond à la valeur moyenne de la fréquence de la règle (nombre d'occurrences divisé par nombre total de mots de l'ensemble du corpus).

Plan de la présentation

Des *scriptae* en français classique

De l'utilité de la scriptométrie

Approche par règles

Approche par réseau de neurones

Conclusion



Achevez, Seigneur, votre ambassade.

- Un modèle de normalisation automatique (un modèle à base de réseau de neurones) entraîné sur PARALLEL17

Approche par réseau de neurones

Achevez, Seigneur, votre ambassade.



Achevez, Seigneur, votre ambaffade.

- Un modèle de normalisation automatique (un modèle à base de réseau de neurones) entraîné sur PARALLEL17

Approche par réseau de neurones

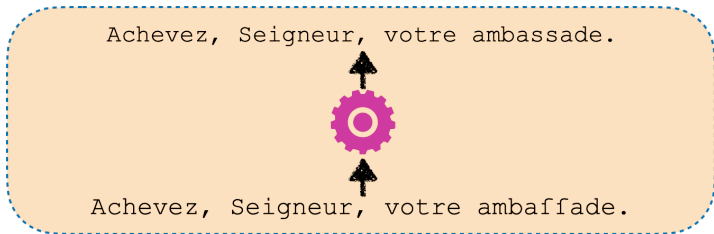
Achevez, Seigneur, votre ambassade.



Achevez, Seigneur, votre ambaffade.

- Un modèle de normalisation automatique (un modèle à base de réseau de neurones) entraîné sur PARALLEL17
- Que pouvons-nous apprendre d'un tel modèle de normalisation, concernant la variation graphique ?

Approche par réseau de neurones



- Un modèle de normalisation automatique (un modèle à base de réseau de neurones) entraîné sur PARALLEL17
- Que pouvons-nous apprendre d'un tel modèle de normalisation, concernant la variation graphique ?
- Deux expériences préliminaires :
 1. Le modèle est-il sensible au changement linguistique ?
 2. Comment extraire les informations apprises par le modèle ?

1. Sensibilité au changement linguistique

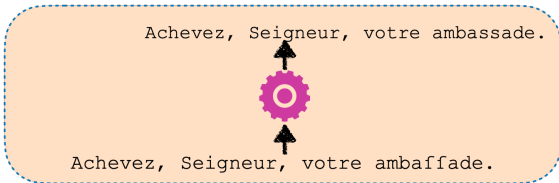
- Y a-t-il suffisamment d'informations dans les phrases originales pour prédire le texte normalisé **et la décennie de rédaction** ?



Achevez, Seigneur, votre ambassade.

1. Sensibilité au changement linguistique

- Y a-t-il suffisamment d'informations dans les phrases originales pour prédire le texte normalisé **et la décennie de rédaction** ?



1. Sensibilité au changement linguistique

- Y a-t-il suffisamment d'informations dans les phrases originales pour prédire le texte normalisé **et la décennie de rédaction** ?

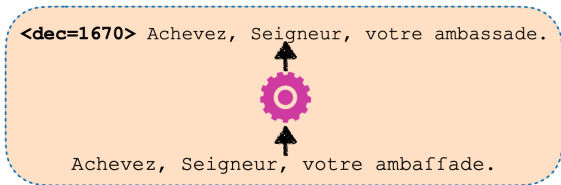
`<dec=1670>` Achevez, Seigneur, votre ambassade.



Achevez, Seigneur, votre ambaffade.

1. Sensibilité au changement linguistique

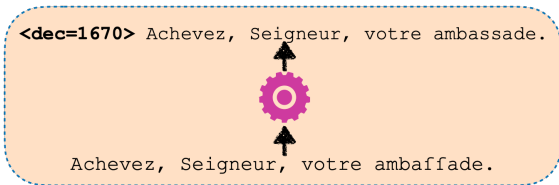
- Y a-t-il suffisamment d'informations dans les phrases originales pour prédire le texte normalisé **et la décennie de rédaction** ?



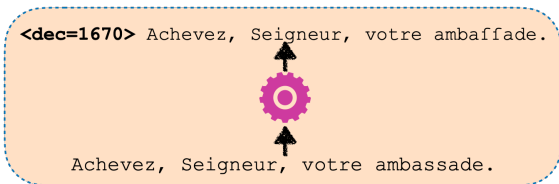
- Mais les informations apprises sont-elles de nature graphique ou lexicale ?

1. Sensibilité au changement linguistique

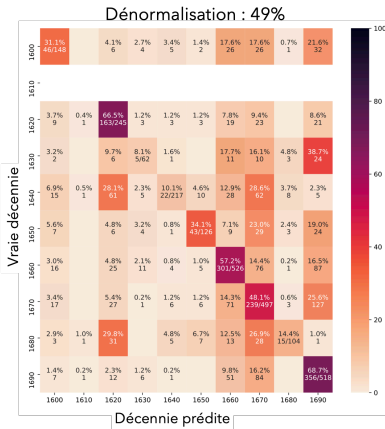
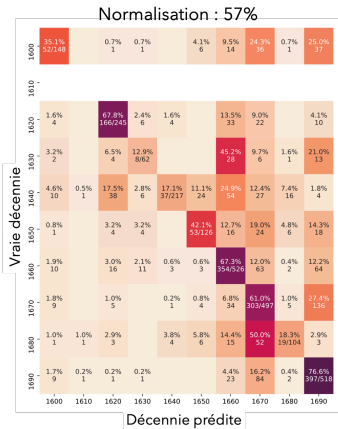
- Y a-t-il suffisamment d'informations dans les phrases originales pour prédire le texte normalisé **et la décennie de rédaction** ?



- Mais les informations apprises sont-elles de nature graphique ou lexicale ?
Expérience de contrôle : dénormalisation (réduction de l'information au simple lexique)

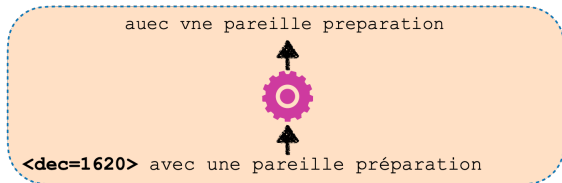


1. Sensibilité au changement linguistique - résultats



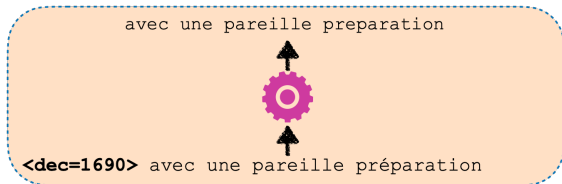
- Prédications assez bonnes pour les deux modèles
- Les scores sont plus élevés quand le modèle a accès à la variation graphique
- Le modèle ne s'appuie donc pas uniquement sur le lexique pour prédire la décennie

2. Quelles informations apprises ?



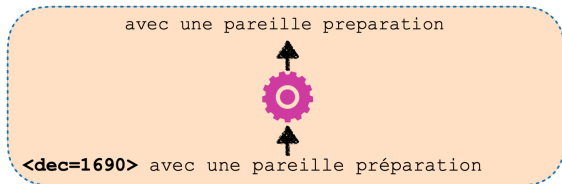
- Modèle de dénormalisation, conditionné sur la décennie

2. Quelles informations apprises ?



- Modèle de dénormalisation, conditionné sur la décennie
- On peut donc créer une version « classique » artificielle pour chaque décennie (1600, 1610, 1620 ...)

2. Quelles informations apprises ?



- Modèle de dénormalisation, conditionné sur la décennie
- On peut donc créer une version « classique » artificielle pour chaque décennie (1600, 1610, 1620 ...)
- L'objectif est de concevoir un corpus parfaitement homogène, où chaque décennie est parfaitement comparable à une autre.

2. Quelles informations apprises ? - exemples

Phrase normalisée (en entrée) :

Un soufflet bien fermé de tous **côtés** fait le **même** effet, **avec une** pareille **préparation**

Phrase denormalisée, conditionnée sur 1620 :

Vn soufflet bien fermé de tous **costez** fait le **mesme** effet, **avec vne** pareille **preparation**

Phrase denormalisée, conditionnée sur 1690 :

Vn soufflet bien fermé de tous **costez** fait le **mesme** effet, **avec une** pareille **preparation**

- Créer des versions « classiques » pour des textes normalisés du corpus D'ALEMBERT
- Aligner le résultat pour chaque décennie avec le texte normalisé et extraire les règles qui permettent de passer d'une version à l'autre, par exemple *une* > *vne*
- Comparer les différentes règles trouvées pour chaque décennie et tenter de dégager une évolution de celles-ci

Quelles informations apprises ? - résultats

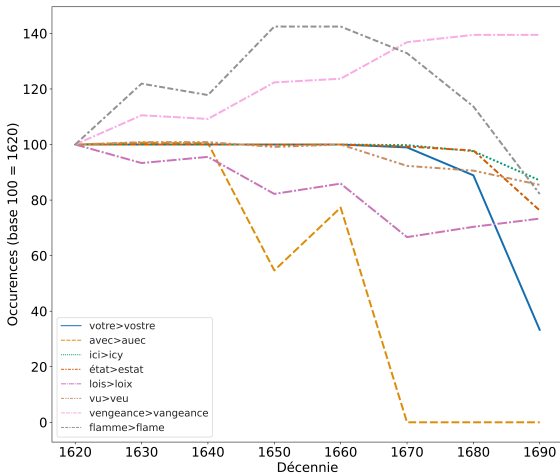


Figure 2 – Evolution de plusieurs règles de dénormalisation appliquées sur les textes normalisés du corpus D'ALEMBERT par périodes de 10 ans entre 1620 (la base 100) et 1700.

Plan de la présentation

Des *scriptae* en français classique

De l'utilité de la scriptométrie

Approche par règles

Approche par réseau de neurones

Conclusion

Conclusion

Conclusions préliminaires

Nous sommes en mesure de tirer des premières conclusions

- ▶ Nous sommes capable de réduire la fenêtre chronologique des analyses à une décennie (contre une vingtaine d'années pour les études de Cl. Vachon)
- ▶ Contrairement aux époques précédentes, la période classique ne semble pas témoigner de phase de reflux
- ▶ Nous pouvons postuler l'existence d'une nouvelle phase graphématique autour des années 1670.

Existence d'un possible angle mort du fait de la méthode même

- ▶ La *scripta* est évaluée à l'aune de sa distance avec le français contemporain
- ▶ Idéal pour étudier la stabilisation du système graphique
- ▶ Méthode valable pour étudier la couleur d'un texte entre deux pôles ?
- ▶ Comment analyser les évolutions de *tempus*, repéré comme moderne si graphié *tems* mais pas comme ancien si graphié *temps* ?

Il reste de multiples pistes à explorer

- ▶ Réintroduire le « pallier textuel », essentiel en philologie, pour caractériser précisément la nature des textes
- ▶ Intégrer les manuscrits dans le corpus, moins sujets que les imprimés à la neutralisation graphématique et à la coexistence de différents systèmes graphiques en diasystème
- ▶ Reprendre l'hypothèse de l'existence d'une *scripta* « habituelle », dont il conviendrait de définir plus précisément la nature de son hétérogénéité
- ▶ Reproduction de la méthode sur des données marquées par des dialectes

Alexandre Bartz, Pedro Ortiz-Suarez... et Myriam pour l'invitation !