

Cartes et graphiques

Quelques pistes de représentation de données linguistiques en utilisant R

Timothée Premat | Univ. Paris 8 & CNRS : UMR 7023 (SFL)

Séminaire des doctorants de SFL (Univ. Paris 8 & CNRS : UMR 7023), 26/01/2021.

Introduction

Problèmes de visualisation de données en linguistique

- La linguistique est souvent une affaire de variation
 - *Faire émerger l'ordre du magma chaotique des données*
- Pour ça, rien de mieux que de visualiser (correctement) les données !
- Présenter quelques problèmes de visualisation de données que j'ai pu rencontrer
- Présenter les solutions que j'ai pu trouver, les discuter avec vous
- Et éventuellement vous donner des pistes si vous avez rencontré des problèmes analogues

Plan

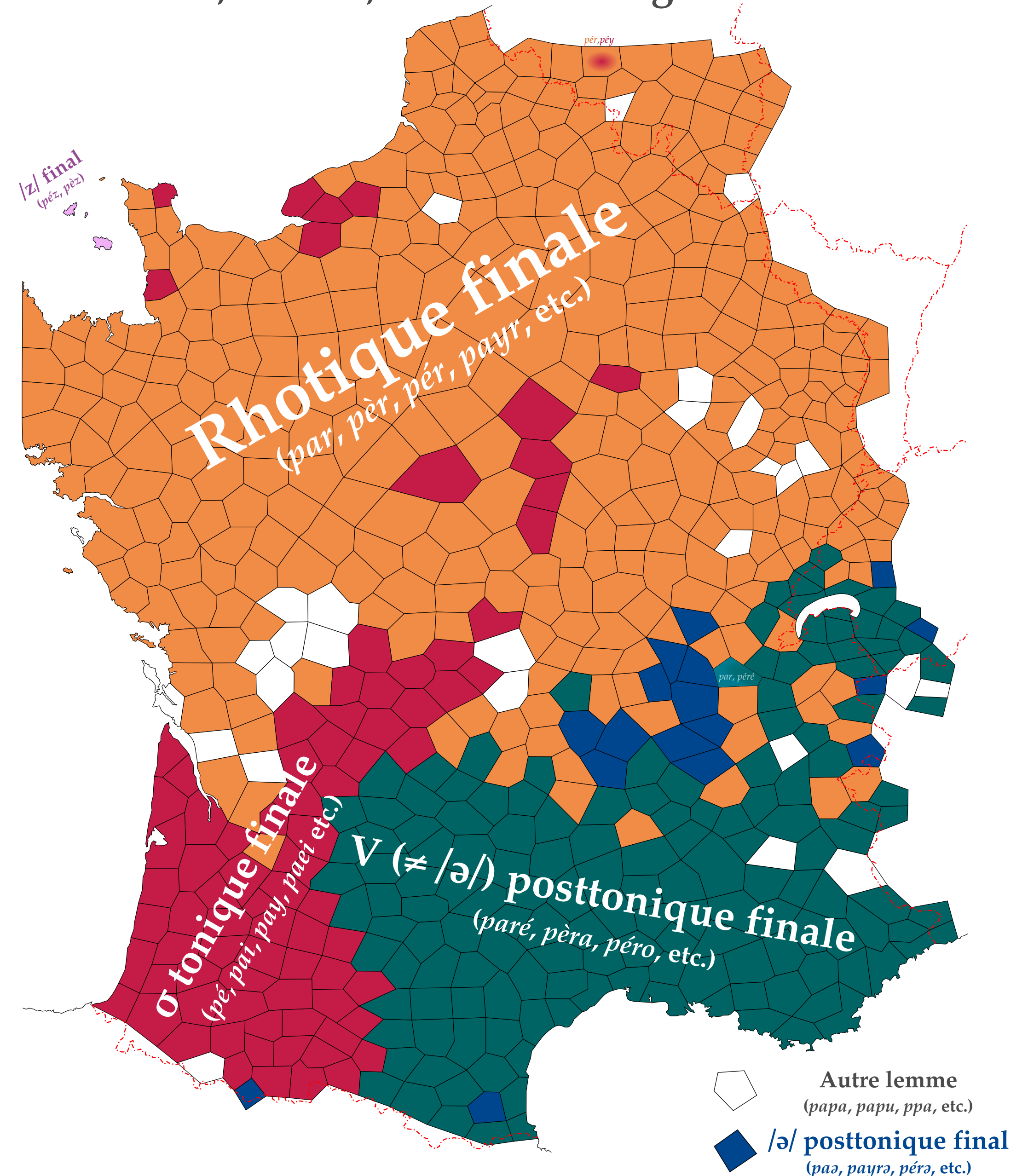
Objectif : introduction à R pour les données variationnelles

- Contexte : mes recherches
- Outil : présentation de R
- Applications
 - Variation dialectale : cartographie linguistique
 - Variation dialectale : représentation non cartographique
 - Distribution de données non dialectale
- Conclusion

Avant R

Faire des cartes sans R...

- À la main
- Avec des programmes de graphismes
 - (Inskape ci-contre)
- ➔ Très mauvaise idée !
 - Très long
 - Très compliqué à modifier en cas d'erreur
 - Aucun traitement *automatisé*
- Les applications SIG (*système d'information géographique*) ne sont pas nécessairement adaptées au besoin du linguistique.



Contexte

Contexte

Ma thèse

- Traitement du /ə/ final en ancien et moyen français (IX^e-XVI^e s.)
 - Dir. Michela Russo & Sophie Wauquier
- Le schwa français est un objet fondamental pour la théorie phonologique, y compris à l'échelle *cross-linguistic* (Anderson 1982)
- Français contemporain 'standard' : schwa résulte d'une insertion phonotactique (insertion de schwa là où il faut une voyelle/un noyau de syllabe et qu'il n'y a pas de voyelle disponible)
 - *belle fille* : /bɛl fiʎ/ [bɛl fiʎ]
- Ancien français : schwa final est voyelle lexicale, supprimée uniquement en hiatus (= élision)
 - *belle fille* : /bɛlə fiʎə/ [bɛlə fiʎə]
 - *belle amie* : /bɛlə amiə/ [bɛl amiə]
- Moyen français : schwa final commence à être supprimé progressivement (→ fr.cont.)

Contexte

Dimensions

- Phonologie représentationnelle : qu'est-ce que schwa en a.fr. et m.fr. ?
- Phonologie prosodique : quel est le domaine de l'élosion ?
- Métrique : comment ça marche dans le discours versifié ?
- Musication : comment ça marche en musique ?
 - métrique et musication \in linguistique
 - donnent accès à des information prosodique en l'absence d'enregistrements.
- Dialectologie :
 - a.fr. : langue non standardisée = multiplicité de normes
 - m.fr. : langue en voie de standardisation = unification progressive des normes
 - Donc : grosse variation diatopique.
- Diachronie

Variations diachronique et diatopique...

= données confuses

Comment y voir clair ?

→ Représentation spatiale :

- cartes
- graphes

R

R

Langage de programmation et logiciel libre

- langage informatique : permet d'interagir avec un ordinateur pour lui faire faire des opérations
- dédié aux statistiques et à la science des données
- permet d'effectuer des opérations sur un jeu de données :
 - faire des calculs (moyenne, coefficient de corrélation (r , ρ , τ), test de probabilités χ etc.)
 - produire des représentations spatiales
 - 'graphiques'
 - produire des transformations

R

Langage de programmation et logiciel libre

- langage ‘simple’ à apprendre
 - syntaxe explicite, faite pour être écrite et lue par un humain
 - larges documentation et communauté en ligne
 - Beaucoup de tutoriels/MOOC en ligne
- permet de faire des choses *vraiment pénibles* à faire dans des programmes grand public comme Excel.
- il existe des alternatives propriétaires et payantes *a priori* plus simple à manier, et des alternatives libres plus compliquées (p.ex. on peut faire tout ça directement en Python)

R

Workflow sur R

1. Préparer les données (saisie manuelle, automatique, extraction depuis une base de données, etc.)
 2. Rédaction du *script* R (ensemble d'instructions R) dans un éditeur de texte
 - Je conseille *Atom* (gratuit et libre)
 3. Exécution de ces instructions dans un logiciel spécialisé
 - Je conseille *RStudio* (gratuit), qui propose une interface graphique pour un certain nombre d'actions
- ➔ Ce qui produit le résultat des instructions (2) sur le jeu de données (1) : graphiques, résultats de calculs, etc.

Dataset

Exemple de jeu de données cartographiques

name;item	name	Lat	Long	description	CODE_DEPT	Nom_dept	Pays	name
1;[õ,õ]	493	48.508533	-2.916611	Plouvara	22	Côtes-dArmor	France	493
3;[õ,õ]	494	48.280427	-2.84085	Uzel	22	Côtes-dArmor	France	494
4;[õ,õ]	485	48.03455	-2.765412	Crêdin	56	Morbihan	France	485
5;[ã,ã]	486	47.837868	-2.639578	Plumelec	56	Morbihan	France	486
6;[ã,ã]	399	49.448627	-2.553337	Saint-Peter-Port	0	Guernesey (Guernsey)	Normandie (Couronne	
7;∅	482	48.448	-2.486	Noyal	22	Côtes-dArmor	France	482
8;[ã,ã]	475	47.591915	-2.456789	Noyal-Muzillac	56	Morbihan	France	475
10;[ã,ã]	484	47.988586	-2.383094	Loyat	56	Morbihan	France	484
11;[ã,ã]	476	47.356155	-2.369867	Guérande	44	Loire-Atlantique	France	476
12;[ã,ã]	398	49.430737	-2.360516	Sark	0	Serq (Sark)	Normandie (Couronne britannique)	398
13;∅	479	46.72451	-2.348529	L'Ile d'Yeu	85	Vendée	France	479
14;[ã,ã]	481	48.655113	-2.331348	Plévenon	22	Côtes-dArmor	France	481
16;[ã,ã]	478	47.024269	-2.303349	Noirmoutier-en-l'Île	85	Vendée	France	478
17;[ã,ã]	483	48.183334	-2.233333	Loscouët-sur-Meu	22	Côtes-dArmor	France	483
19;[õ,õ]	396	49.722038	-2.201772	Sainte-Anne	0	Aurigny (Alderney)	Normandie (Couronne britanni	
20;[ã,ã]	465	47.636829	-2.126369	Saint-Jean-la-Poterie	56	Morbihan	France	465
21;[ã,ã]	466	47.396751	-2.09179	Besné	44	Loire-Atlantique	France	466
22;[õ,õ]	397	49.235168	-2.090838	Trinity	0	Jersey	Normandie (Couronne britannique)	397
23;[ã,ã]	463	47.876583	-2.084128	Comblessac	35	Ille-et-Vilaine	France	463
24;N/A	471	48.383331	-2.066667	Trévron	21	Côte-dOr	France	471
25;N/A	394	49.711658	-1.932204	Auderville	50	Manche	France	394
26;N/A	467	47.128164	-1.910011	Chéméré	44	Loire-Atlantique	France	467
27;∅	470	48.604076	-1.893617	La Gouesnière	35	Ille-et-Vilaine	France	470
28;∅	459	46.671261	-1.883768	Givrand	85	Vendée	France	459
30;∅	458	46.883331	-1.833333	La Garnache	85	Vendée	France	458
31;[a,a]	462	48.12746	-1.817488	L'Hermitage	35	Ille-et-Vilaine	France	462
32;[ã,ã]	453	47.824821	-1.806856	Messac	35	Ille-et-Vilaine	France	453
33;[ã,ã]	395	49.402164	-1.780064	Les Moitiers-d'Allonne	50	Manche	France	395
35;N/A	461	48.309021	-1.668098	Montreuil-sur-Ille	35	Ille-et-Vilaine	France	461
36;[ã,ã]	540	46.466759	-1.617908	Talmont-Saint-Hilaire	85	Vendée	France	540
38;[õ,õ]	387	49.200645	-1.562132	Créances	50	Manche	France	387
40;[õ,õ]	690	43.483093	-1.558613	Biarritz	64	Pyrénées-Atlantiques	France	690
41;[õ,õ]	460	48.516666	-1.55	Vieuxviel	35	Ille-et-Vilaine	France	460
42;[õ,õ]	446	47.340271	-1.526231	Sucé-sur-Erdre	44	Loire-Atlantique	France	446

Script

Exemple de Script R

```
#Timothée Premat | 25 octobre 2020, based on original script by:  
#Mathieu Avanzi | 22 avril 2020  
  
#packages that need to be installed prior to run the code  
library(maps)  
library(plyr)  
library(dplyr)  
library(rgdal)  
library(rgeos)  
library(ggplot2)  
library(ggsn)  
library(scales)  
library(ggmap)  
  
#load data  
dataALF = read.table("commencent_ALF_simple.txt", header=T, sep=";", quote="", dec=".")  
dataALF2 = read.table("ALF_point.txt", header=T, sep="\t", quote="", dec=".")  
dataALF = merge(dataALF,dataALF2,by="name")  
dataALF[!complete.cases(dataALF),]  
  
head(dataALF)  
|  
#get map  
world <- get_googlemap("Paris", zoom = 6, style = s, size=c(640,640), scale=2)  
  
#plot it  
g= ggmap(world)+  
  geom_point(data=dataALF, aes(x=Long,y=Lat,shape=item), size=3) +  
  
  scale_x_continuous(limits = c(-3.5,8), expand = c(0, 0)) +  
  scale_y_continuous(limits = c(45,51.15), expand = c(0, 0))+  
  
  scale_shape_manual(values=c(15,16,17,18,14,0,1,2,5,43,42), "",  
                      limits=c("[ě,ě]", "[ã,ã]", "[õ,õ]", "[ü]", "[õ]",
```

RStudio

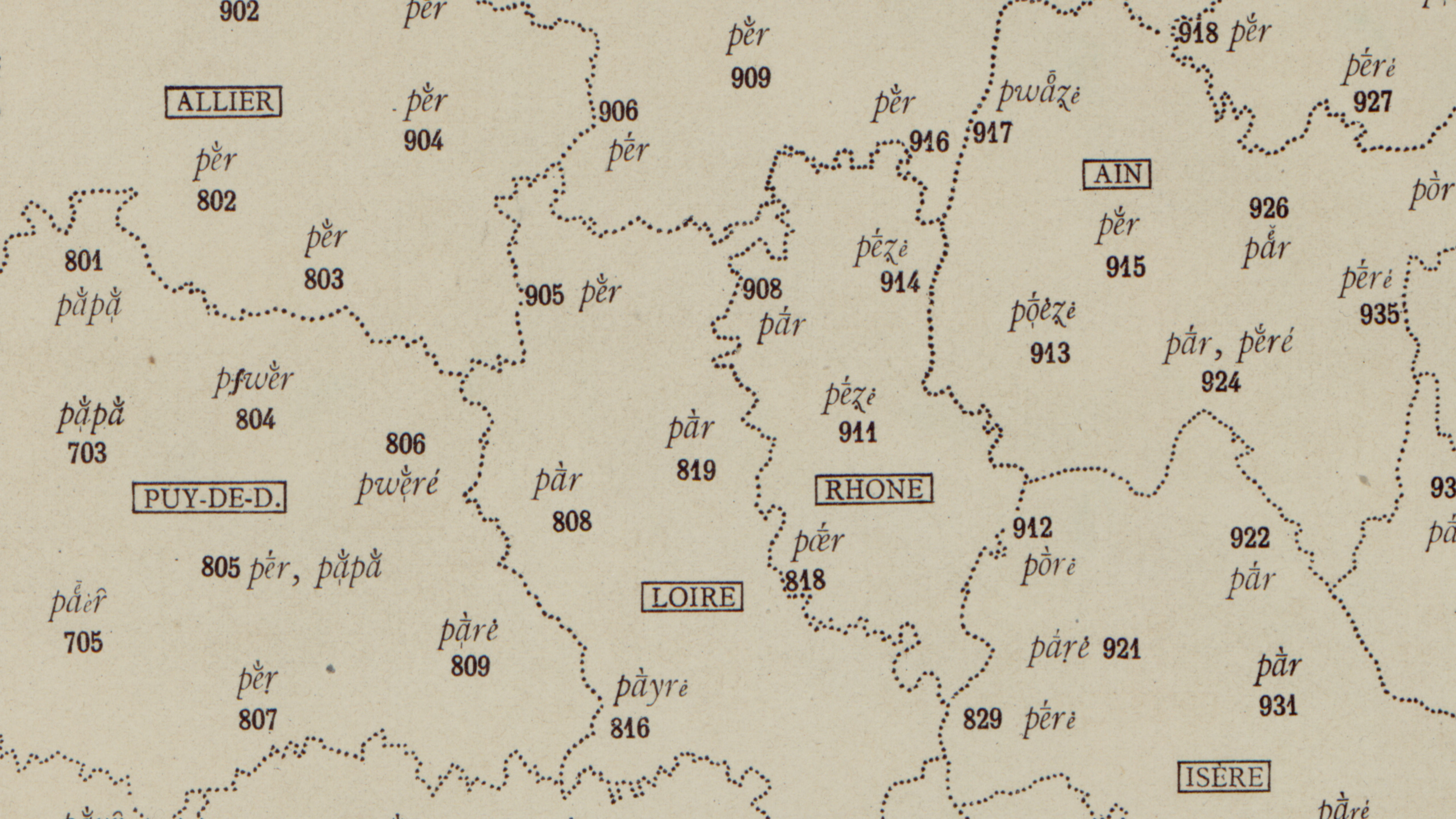
Exécution du script dans RStudio

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for loading packages, reading data, and plotting a map. The code includes comments in French and uses `ggplot2` for visualization.
- Environment:** Lists loaded data frames such as `HOM_R_2`, `HOM_R_3`, `HOM_R_3_Prose`, `HOM_R_3_Vers`, `HOMME`, `opar`, `p`, `paIA_et_U`, `Scatter_accents`, `Scatter_voix`, `VIDE`, and `Voix_accents`.
- Console:** Shows the execution of the script, including warnings about missing values and the saving of the plot as `Aacomment_ALF.png`.
- Plots:** A map titled "Commentent" showing the presence of a posttonic vowel in the region of Les pommiers. The map includes a legend with symbols for different vowel categories: `[ē,ē]` (square), `[ā,ā]` (circle), `[ō,ō]` (triangle), `[ū]` (diamond), `[ê]` (square with X), `[e,e]` (square), `[a,a]` (circle), `[o,o]` (triangle), `[u]` (diamond), `∅` (plus), and `N/A` (asterisk).

Applications

Cartographie avec R

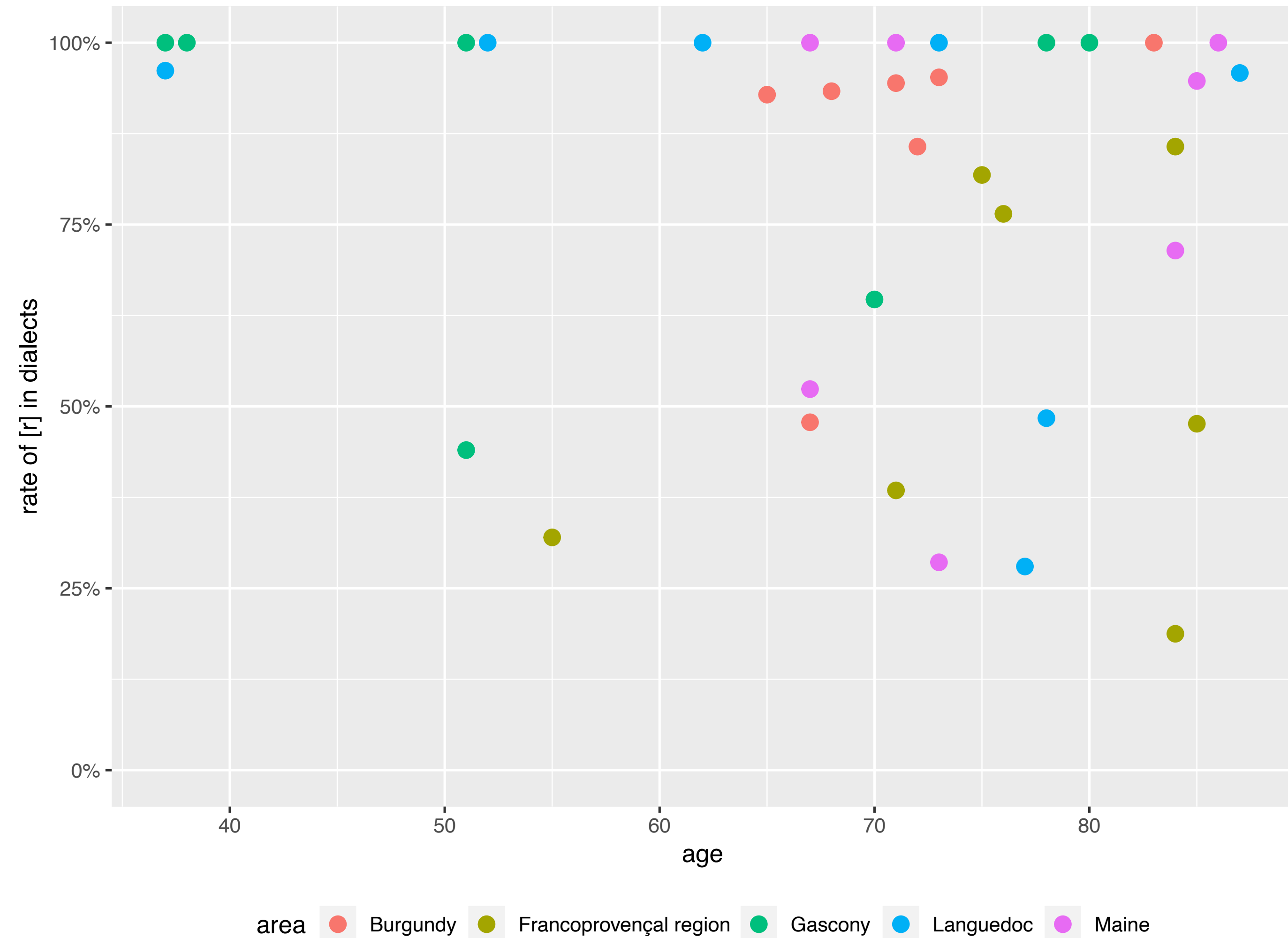


**Comment produire des cartes
interprétatives de manière automatisée ?**

Variation dialectale

Des cartes dans R ?

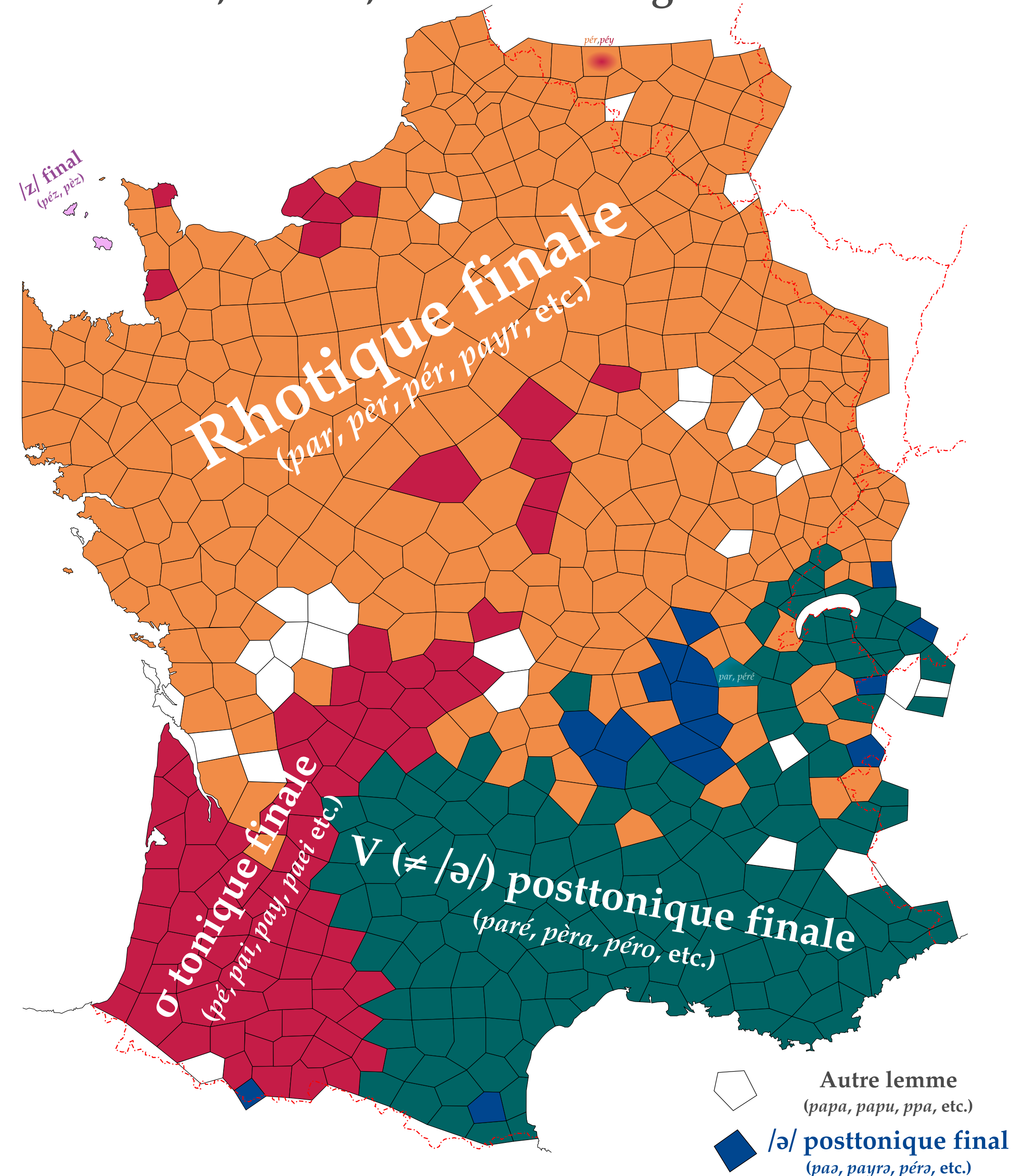
- R est un programme qui peut produire des graphiques
- Nuage de points (*dotplot*)
 - Chaque point est défini par :
 - des coordonnées x, y
 - une esthétique
 - Exemple :
 - $x =$ âge des locuteurs
 - $y =$ taux de [r] sur l'ensemble des /R/
 - couleur : région du locuteur



Variation dialectale

Des cartes dans R ?

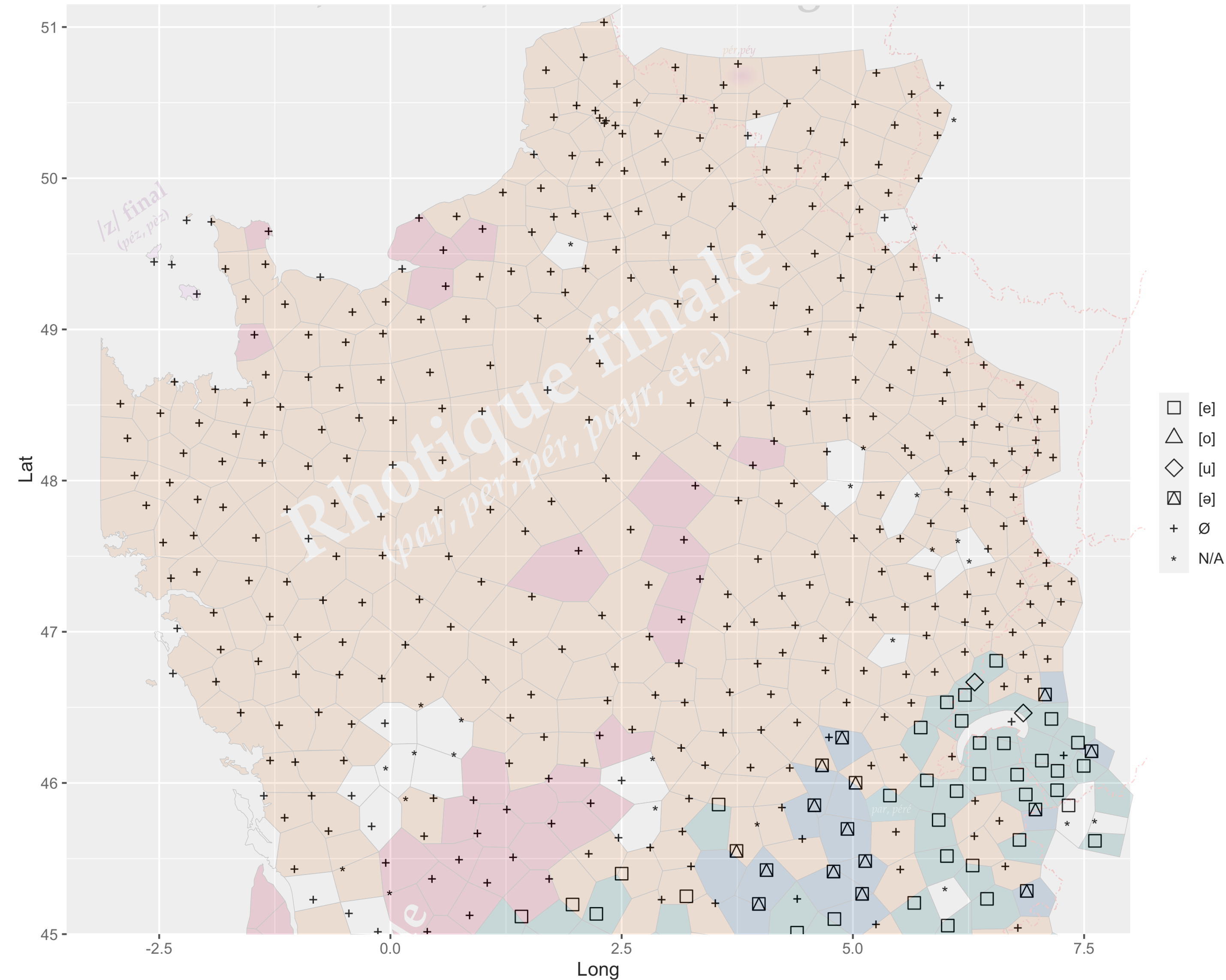
- Une carte est un graphique
- Dans le monde réel, chaque point/région est défini·e par :
 - des coordonnées (GPS, etc.)
 - ici, un polygone tracé autour du village dont vient la forme
 - une esthétique
 - ici, par type de réalisation de la V finale atone
 - + un label



Variation dialectale

Des cartes dans R ?

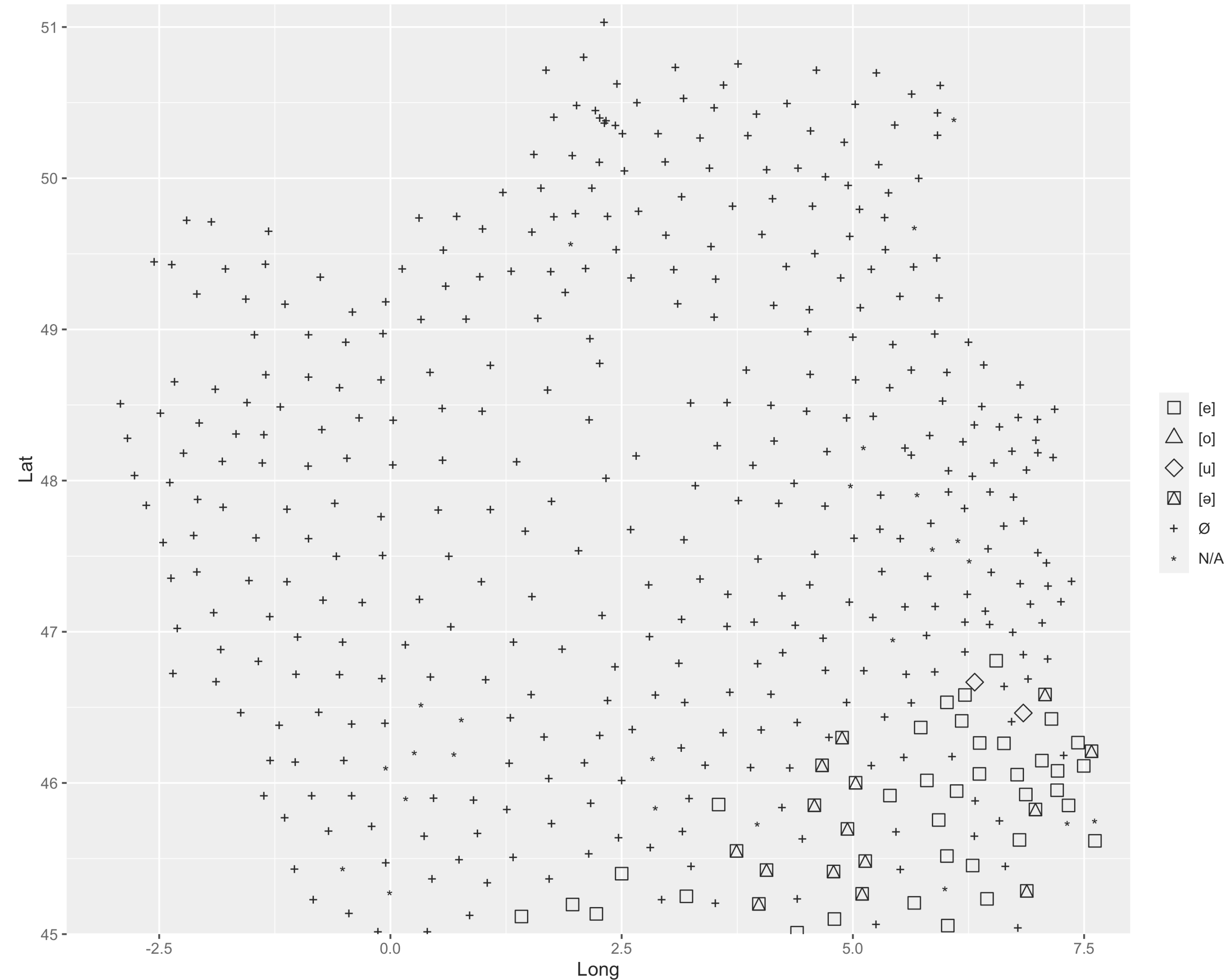
- On peut donc superposer un nuage de point à une carte
- Si :
 - x = longitude
 - y = latitude
 - forme = typologie *ad hoc*



Variation dialectale

Des cartes dans R ?

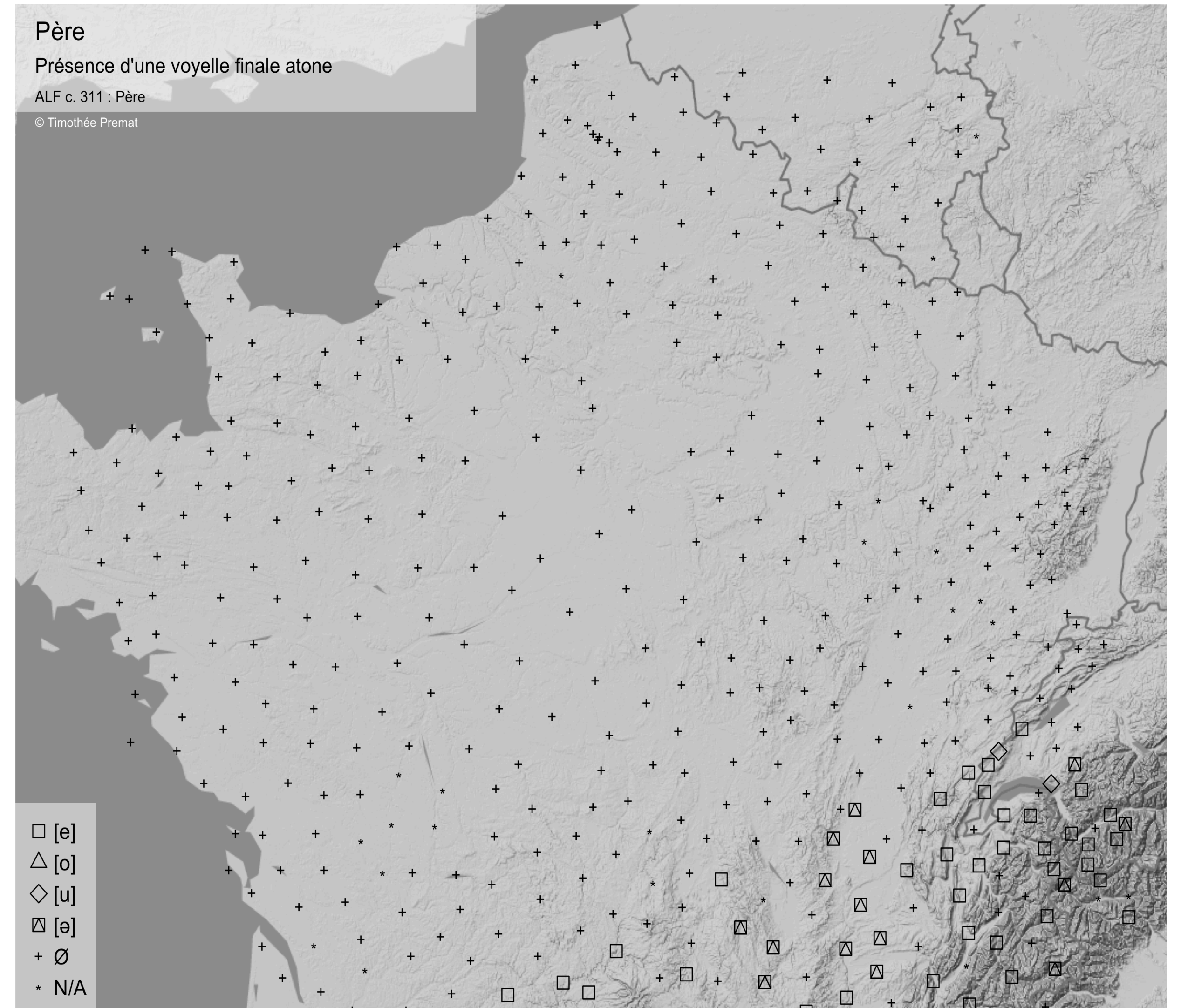
- On peut donc superposer un nuage de point à une carte
- Si :
 - x = longitude
 - y = latitude
 - forme = typologie *ad hoc*



Variation dialectale

Des cartes dans R ?

- Il suffit ensuite d'ajouter un fond de carte
- aligné automatique par des *packages* R sur les coordonnées Long, Lat.
- Ici, fond de carte Google Maps personnalisé. Énormément de fonds de cartes différents possibles, paramétrables, etc.



Variation dialectale

Workflow (procédure de travail)

- Saisie des données dans un tableur
 - d'une part : chaque point : long., lat. (1)
 - chaque point : catégorie de la typologie *ad hoc* (2)
- Écriture du script R
 - définition du fond de carte et des esthétiques
- Exécution du script R
- Pour faire une nouvelle carte, il suffit de préparer les données dans un tableur et de régler les esthétiques, et... c'est tout !

```
name Lat Long description CODE_DEPT Nom_dept Pays
493 48.508533 -2.916611 Plouvara 22 Côtes-dArmor Fran
494 48.280427 -2.84085 Uzel 22 Côtes-dArmor Fran
485 48.03455 -2.765412 Crêdin 56 Morbihan France
486 47.837868 -2.639578 Plumelec 56 Morbihan Fran
399 49.448627 -2.553337 Saint-Peter-Port 0 Guernese
482 48.448 -2.486 Noyal 22 Côtes-dArmor France
475 47.591915 -2.456789 Noyal-Muzillac 56 Morbihan
484 47.988586 -2.383094 Loyat 56 Morbihan France
476 47.356155 -2.369867 Guérande 44 Loire-Atlantiq
398 49.430737 -2.360516 Sark 0 Serq (Sark) Normandi
479 46.72451 -2.348529 L'Ile d'Yeu 85 Vendée Fran
(1) 48.655113 -2.331348 Plévenon 22 Côtes-dArmor
478 47.024269 -2.303349 Noirmoutier-en-l'île 85 Ve
```

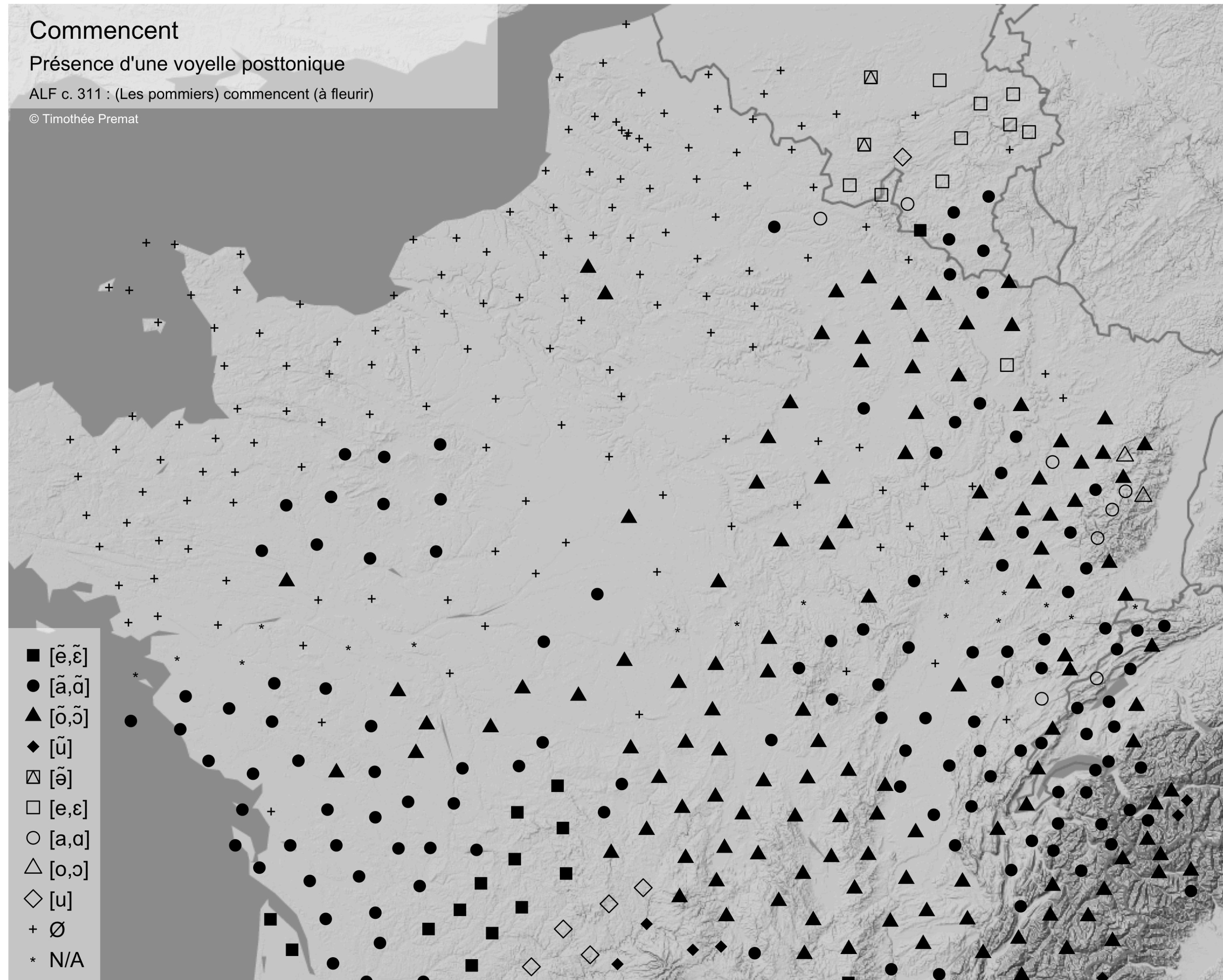
```
name;item
1; [õ, õ]
3; [õ, õ]
4; [õ, õ]
5; [ã, ã]
6; [ã, ã]
7; Ø
8; [ã, ã]
10; [ã, ã]
11; [ã, ã]
12; [ã, ã]
13; Ø
14; [ã, ã]
16; [ã, ã]
```

(2)

Variation dialectale

Autres exemples

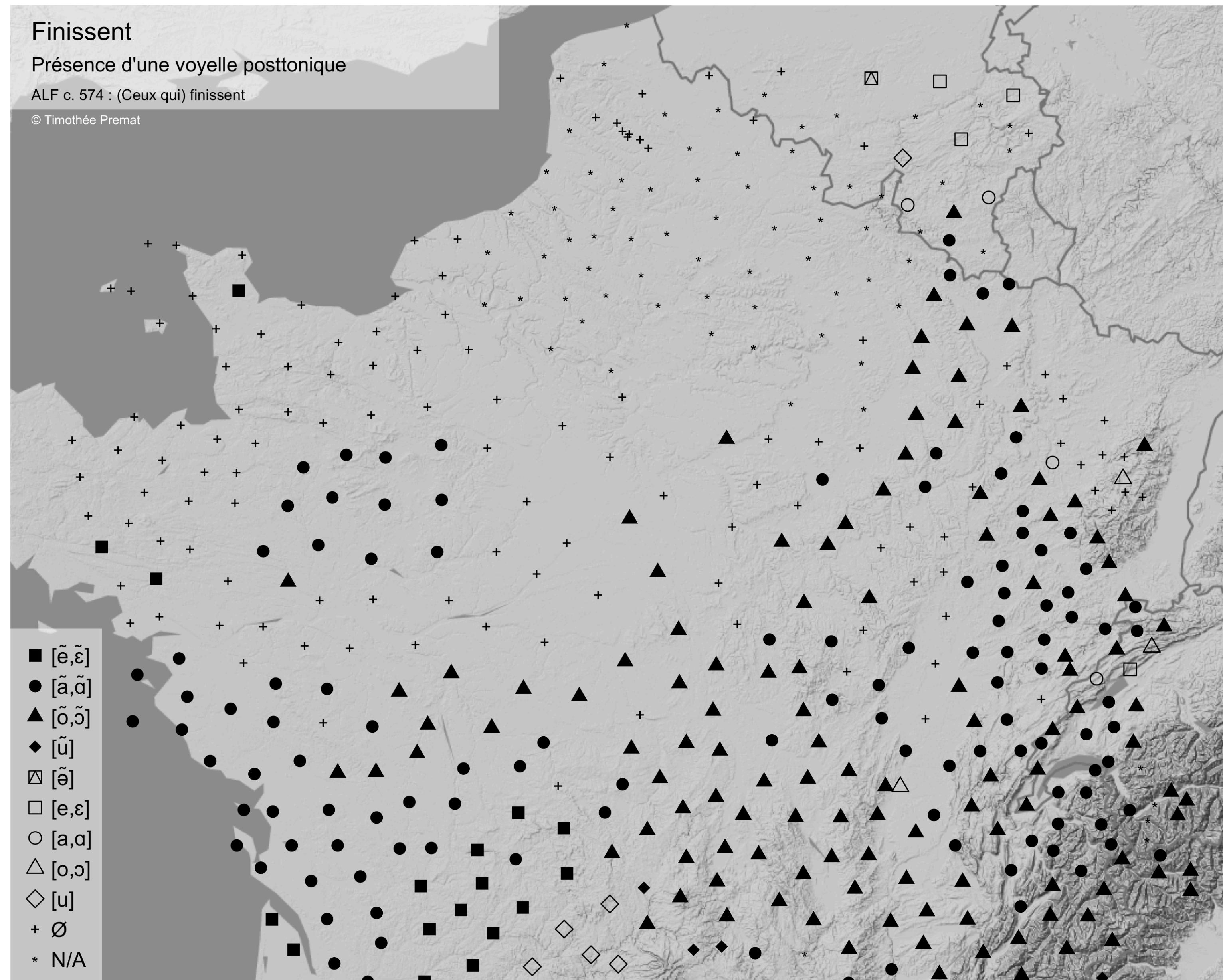
- Présence/absence et type de voyelle pour *-ent* dans *ils commencent*
- Données de l'ALF



Variation dialectale

Autres exemples

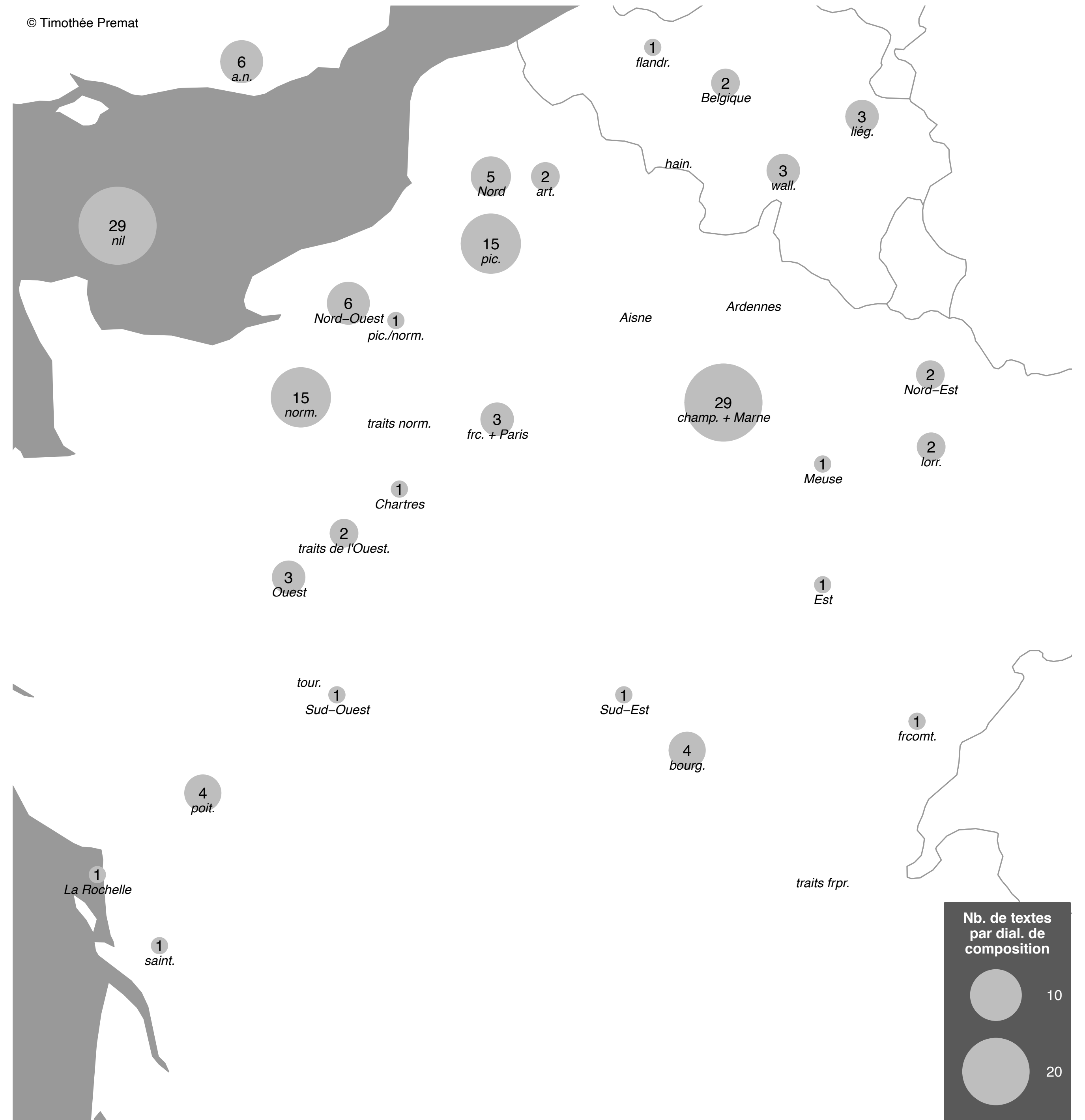
- Présence/absence et type de voyelle pour *-ent* dans *ils finissent*
- Données de l'ALF



Composition d'un corpus

NCA

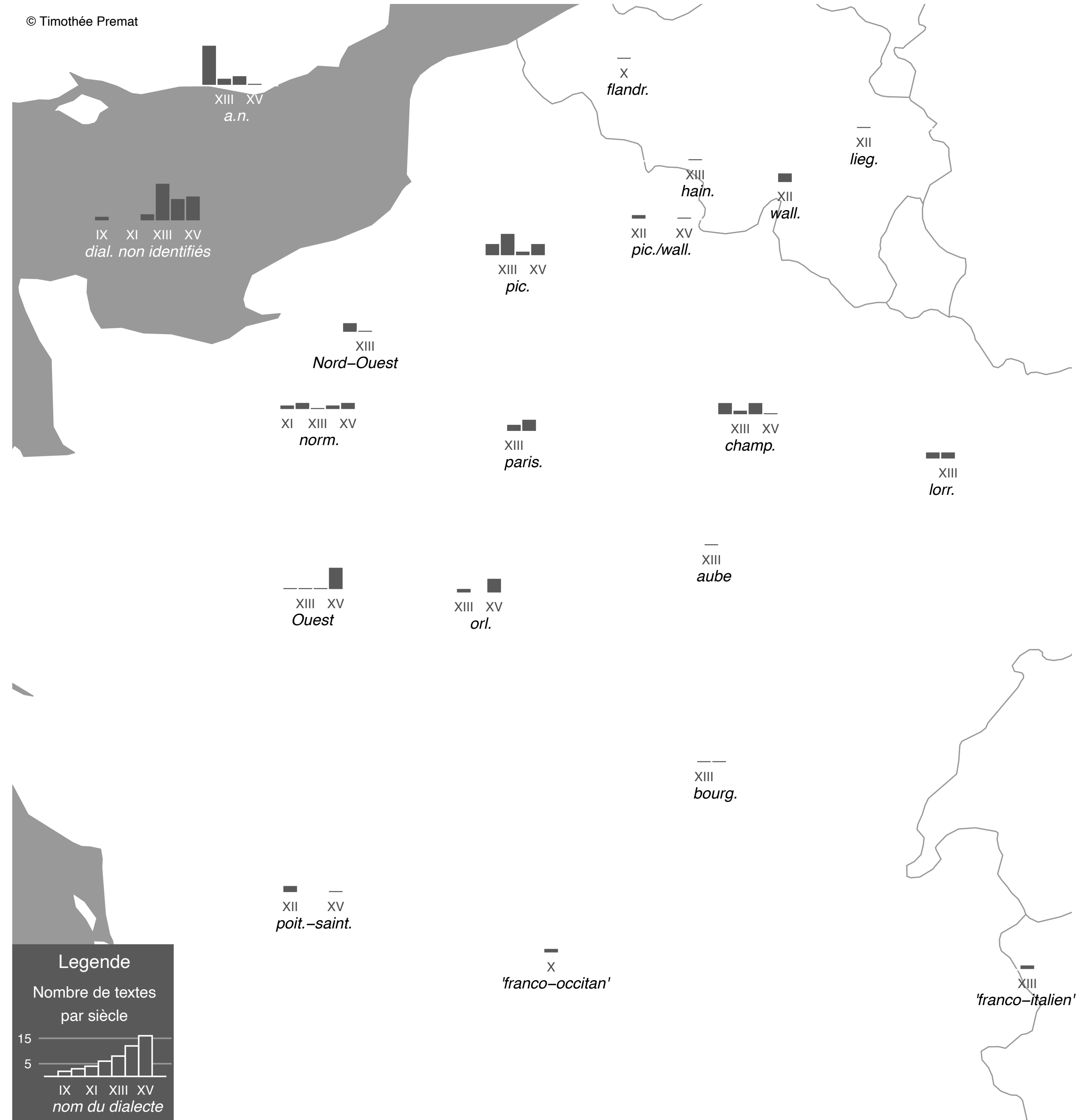
- *Nouveau Corpus d'Amsterdam*
 - Stein, Kunstmann & Gleßgen (2006)
- Objectif : représenter la distribution des textes du corpus
 - par dialecte
- Solution : placer un point pour chaque dialecte, dont la taille représente le nombre de textes



Composition d'un corpus

BFM

- *Base de Français Médiéval*
 - Guillot-Barbance, Heiden & Lavrentiev (2017)
- Objectif : représenter la distribution des textes du corpus
 - par dialecte **et**
 - par siècle
- Solution : placer un histogramme (graphique en barres) représentant la date pour chaque dialecte



Applications

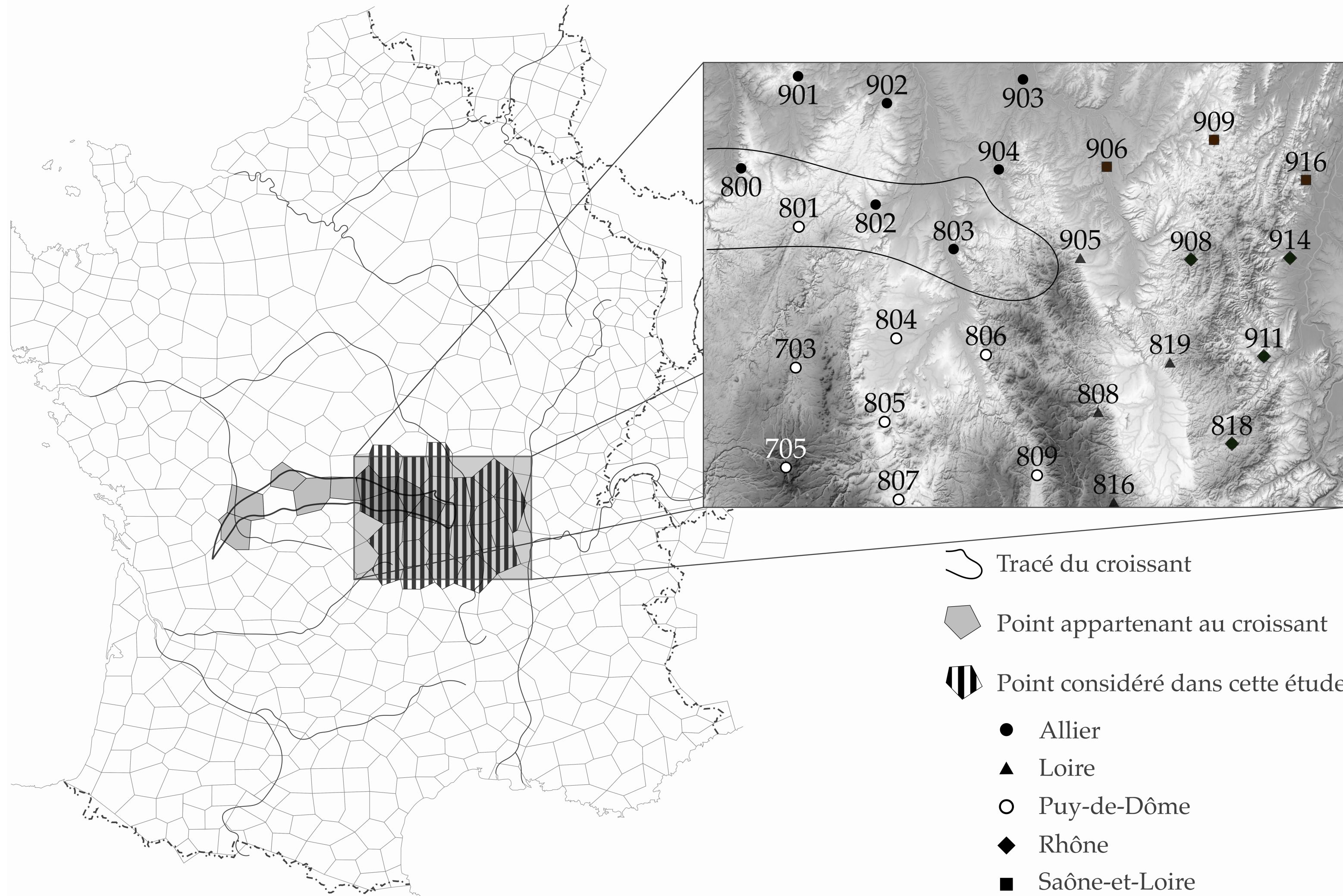
Applications non cartographiques de R

Variation dialectale en phonologie

Variation dialectale

Représentation non géographique

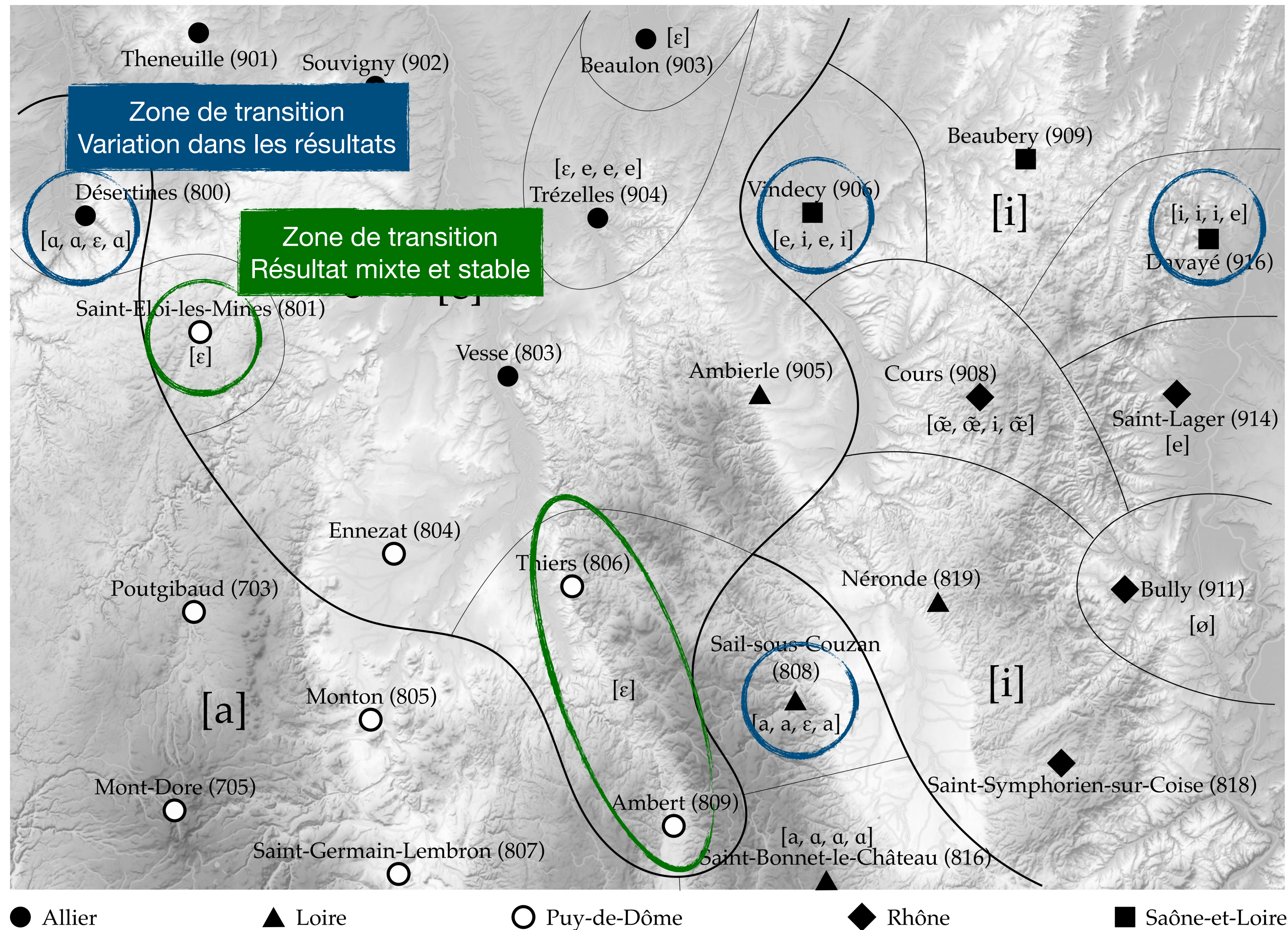
- Russo & Premat :
 - article sur la variation des voyelles finales là où se rencontrent l'oil, l'occitan et le francoprovençal.
 - Données ← ALF
 - 4 mots par points par type de voyelle
 - Simple à cartographier si les points sont cohérents
 - Plus difficiles si 1 points a différents traitements de la V, ou s'il a des traitements intermédiaires



Variation dialectale

Représentation non géographique

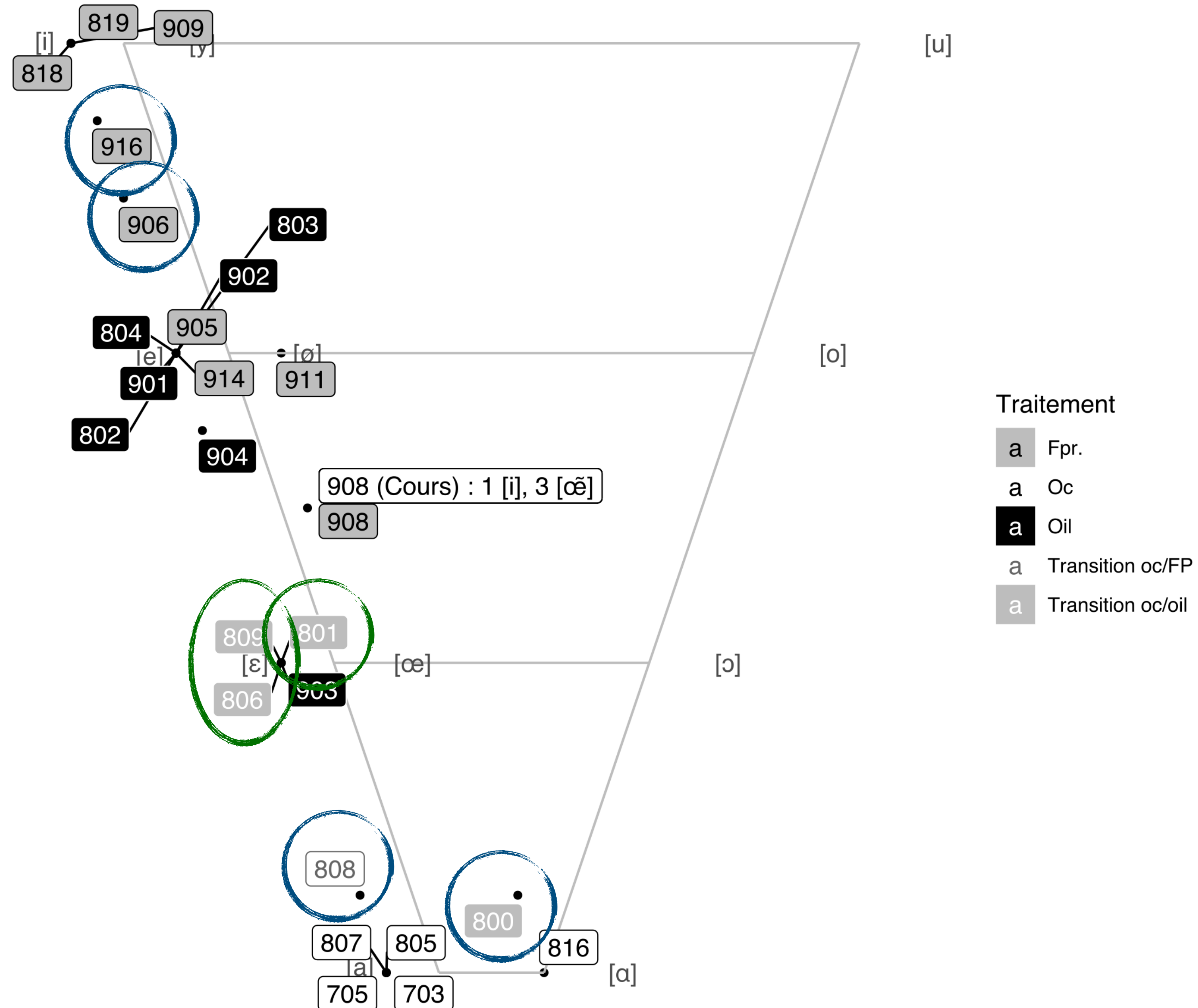
- Exemple : Traitement de -ARE suivi d'un élément palatalisant
- Mots :
 - *nettoyer* *NITIDIĀRE
 - *pêcher* *PISCĀRE
 - *pisser* *PISSIĀRE
 - *purger* PURGĀRE
- La carte montre l'espace géographique, pas la proximité dialectale
- Objectif : montrer à quel diasystème se rattache chaque point



Variation dialectale

Représentation non géographique

- Comment *projeter dans l'espace* des voyelles ?
 - Sur le trapèze vocalique
 - Représente les caractéristiques acoustiques des voyelles (F1, F2)
- Les **résultats mixtes stables** sont naturellement *entre* les diasystèmes
- Comment représenter **les points qui ont une variation** ? En faisant la moyenne de leurs résultats.
- Peut permettre de mesurer à quel point un point est proche de chaque diasystème...



Applications

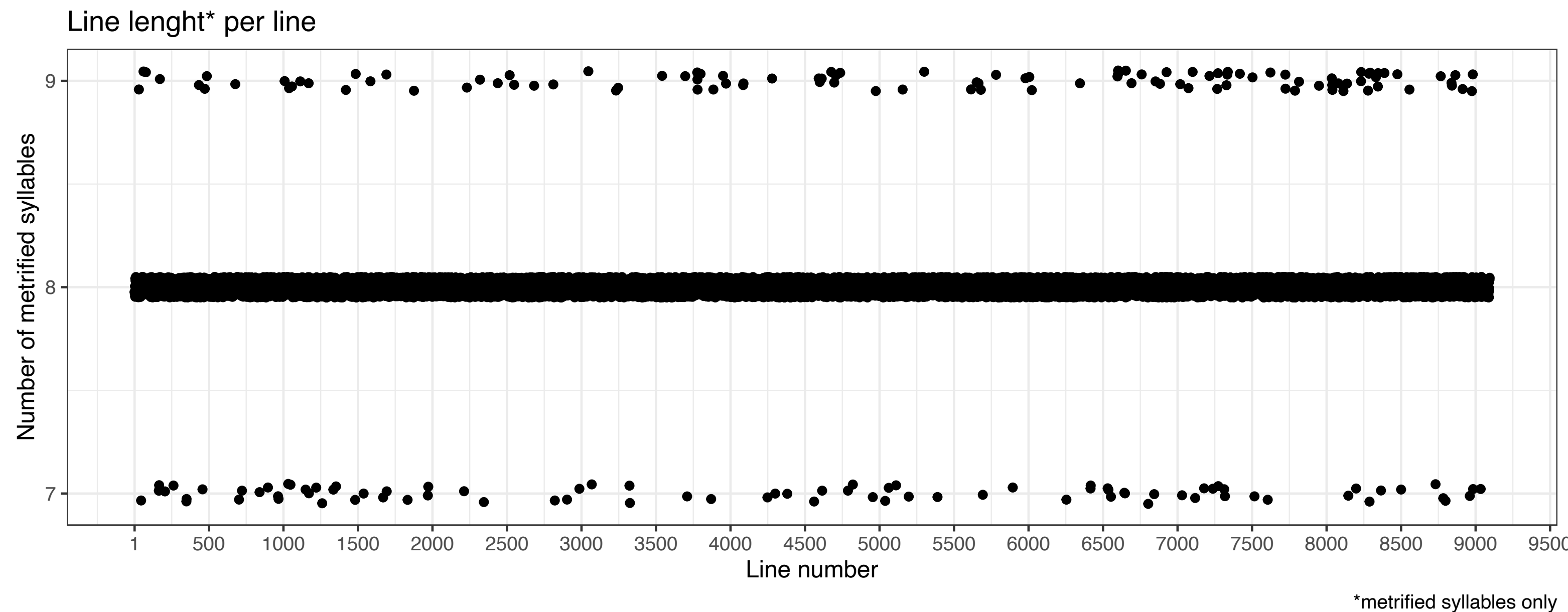
Applications non cartographiques de R

Autres types de variations

Métrique

Analyse des vers d'un texte

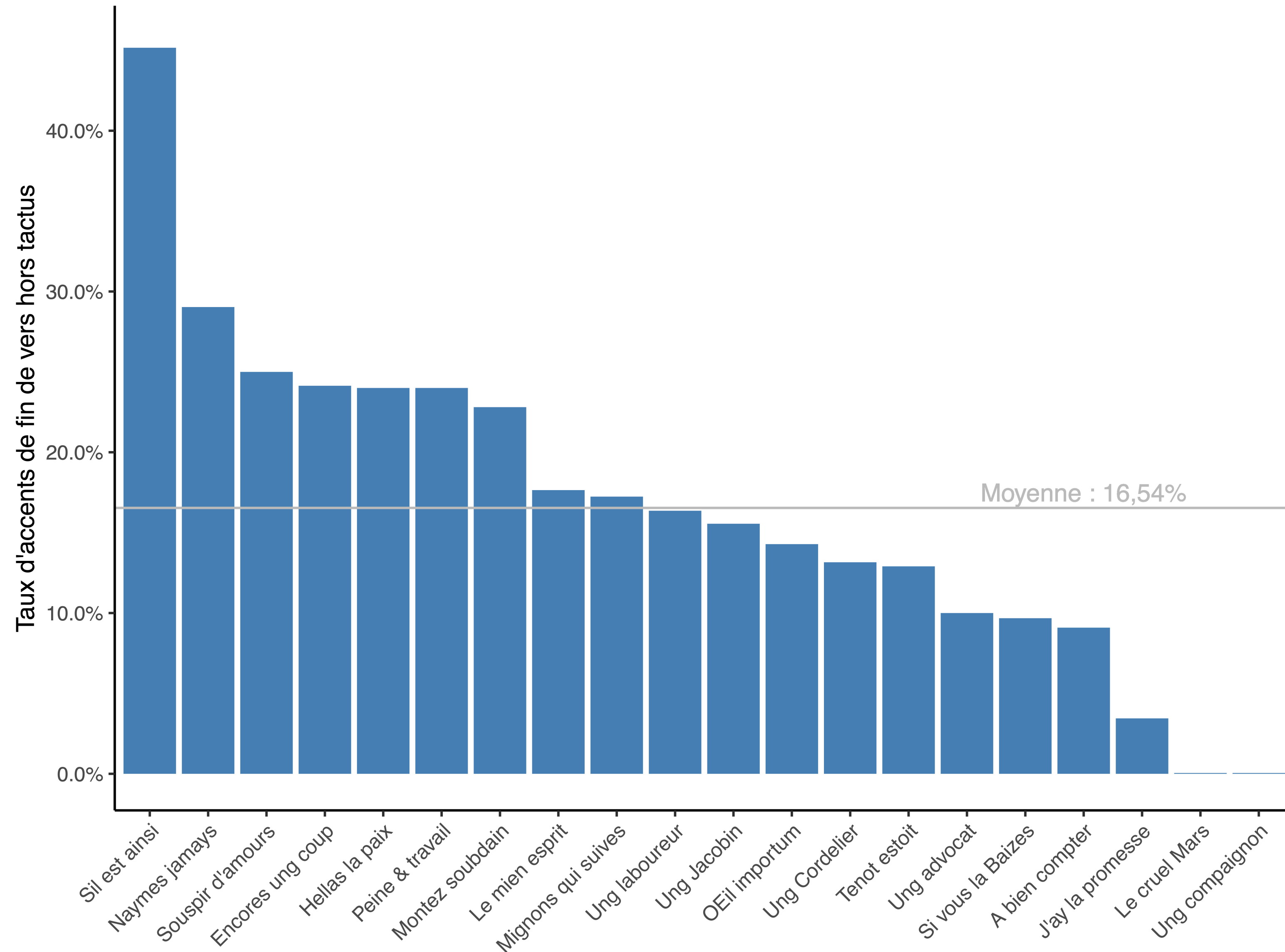
- PAM : Programme d'Analyse Métrique (Poggio & Premat 2019)
- Passe en revue un texte et signale quels sont les vers qui sont 'faux' (qui n'ont pas le bon nombre de syllabes)
- Objectif : avoir une vue d'ensemble de où sont situés les vers 'faux' dans le texte



Musication

Représenter la variation entres les voix

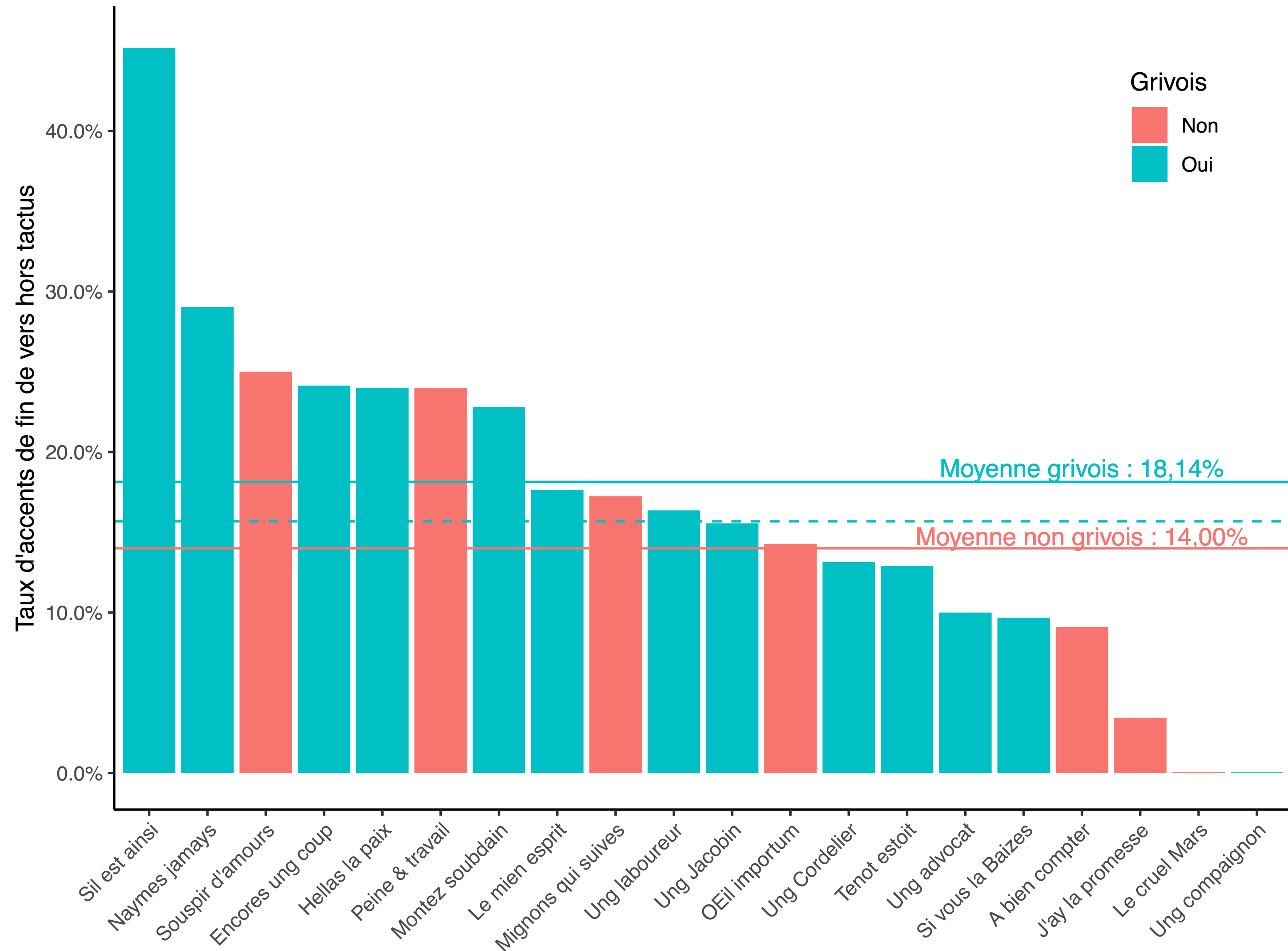
- Henry Fresneau : compositeur du XVIe s.
- Polyphonie à 4 voix
- Objectif : mesurer le taux d'alignement entre accents de fin de vers et temps forts musicaux
- En distinguant voix par voix
- En en fonction de différents facteurs



Musication

Représenter la variation entres les voix

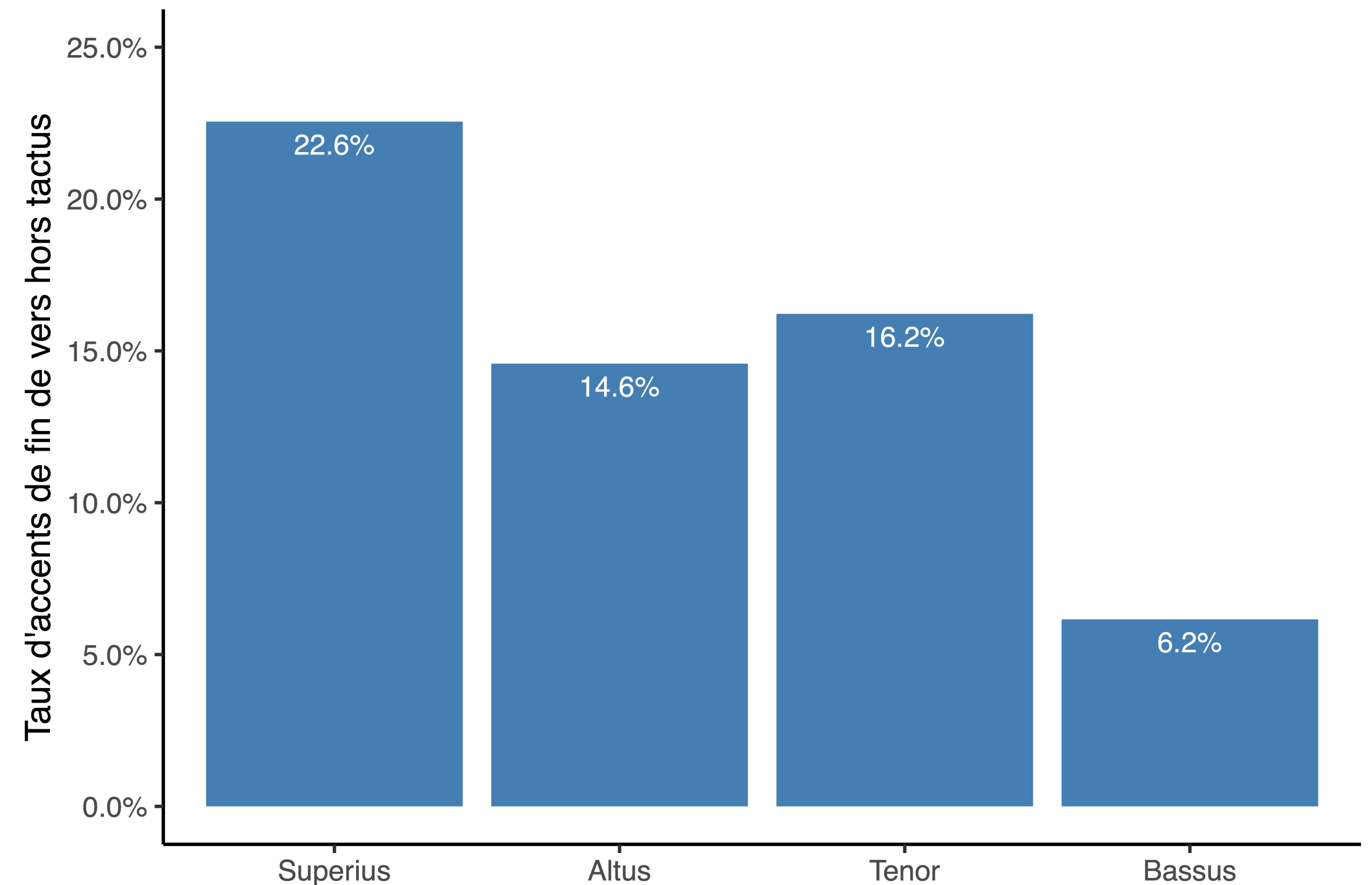
- Facteur :
 - Caractère grivois ou non du texte
 - Ligne pointillée : chansons grivoises sauf *Sil est ainsi*
 - Conclusion : le registre ne change pas significativement le taux d'alignements



Musication

Représenter la variation entres les voix

- Facteur :
 - Voix
- Conclusion :
 - La voix de basse a bien moins d'accents de vers non alignés avec des temps forts
 - Mais est-ce juste une moyenne, ou est-ce vrai dans chaque chanson ?



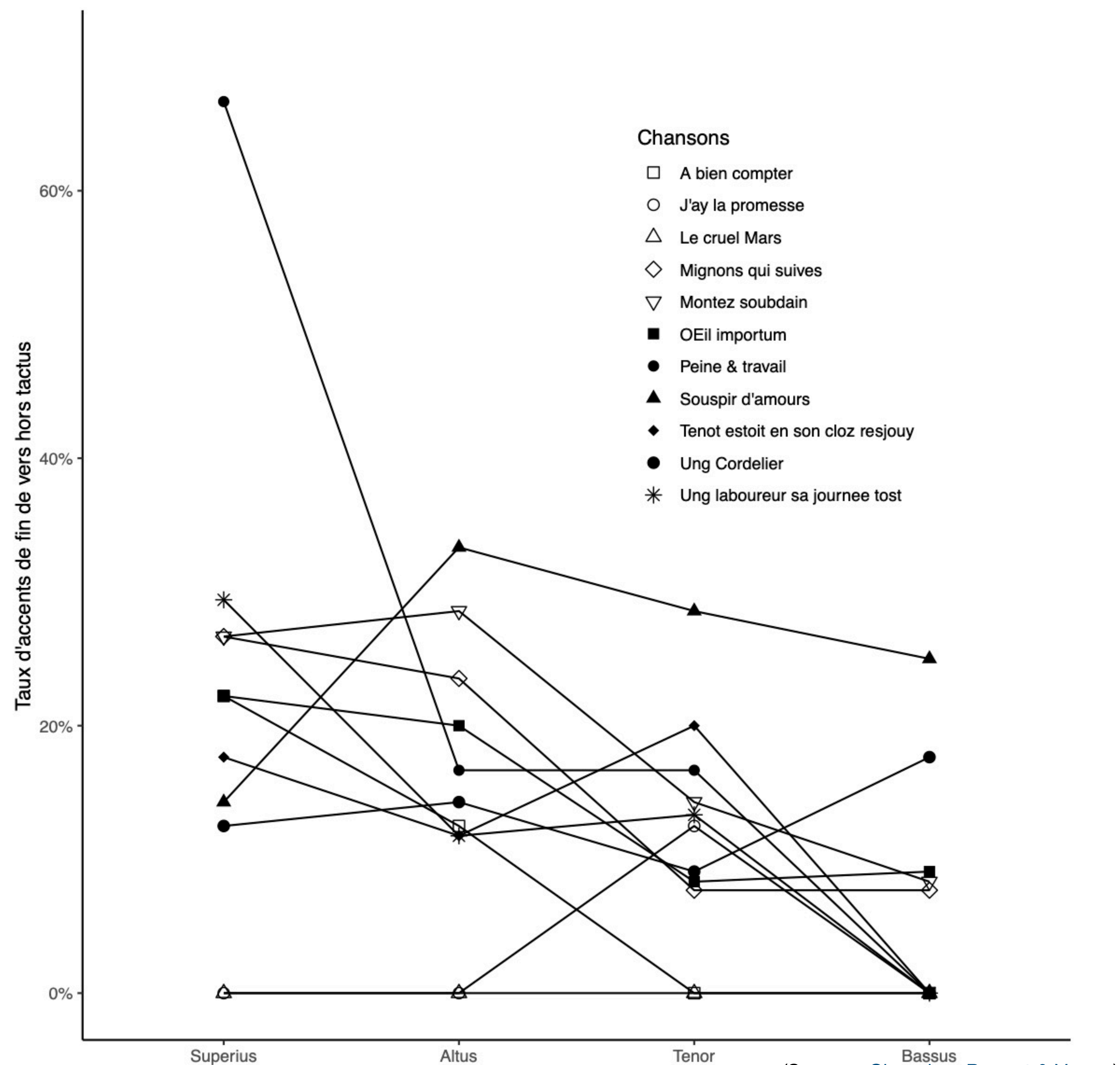
Musication

Représenter la variation entres les voix

Sur 11 chansons

- Oui :
 - 6 chansons
 $\% \text{ bassus} < \% \text{ tenor}$
 - 4 chansons
 $\% \text{ bassus} < \% \text{ autres voix}$
- Non :
 - 1 chansons
 $\% \text{ bassus} > \% \text{ tenor}$
 - 4 chansons
 $\% \text{ bassus} \approx \% \text{ tenor}$

Donc **non**, la basse n'a pas nécessairement moins d'accents de vers non alignés avec des temps forts que les autres voies !



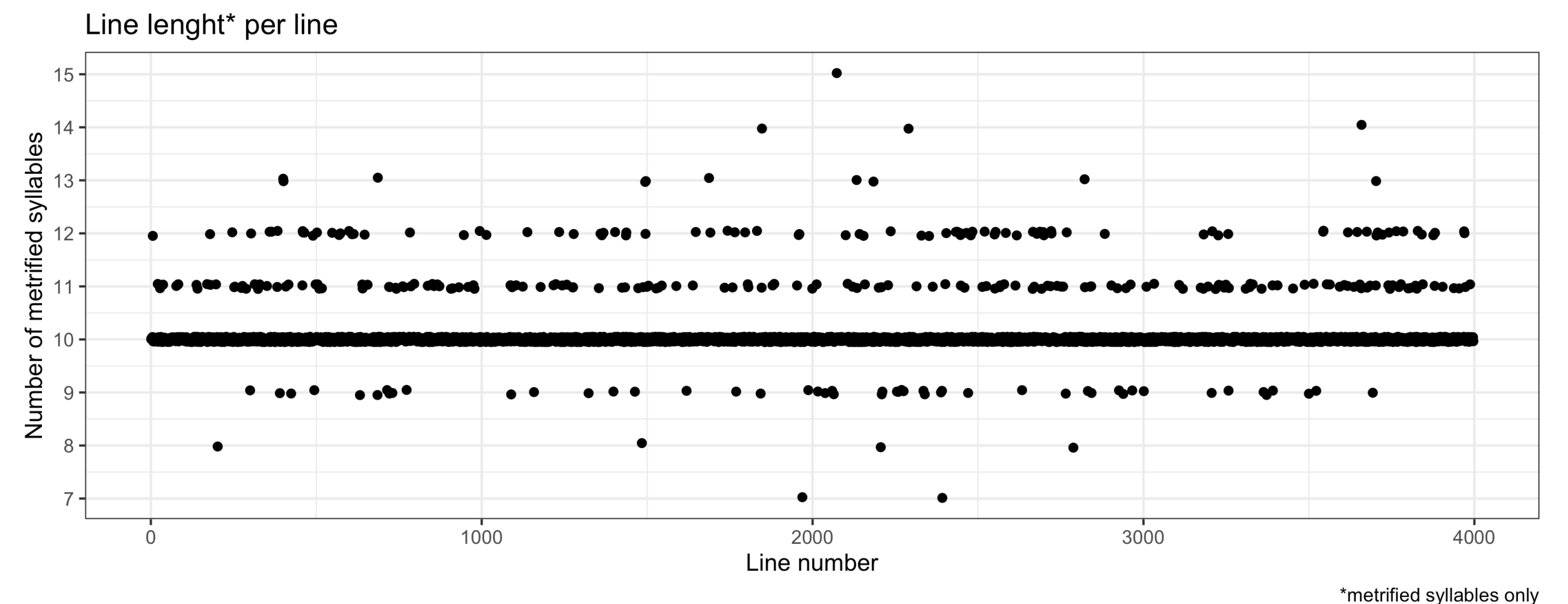
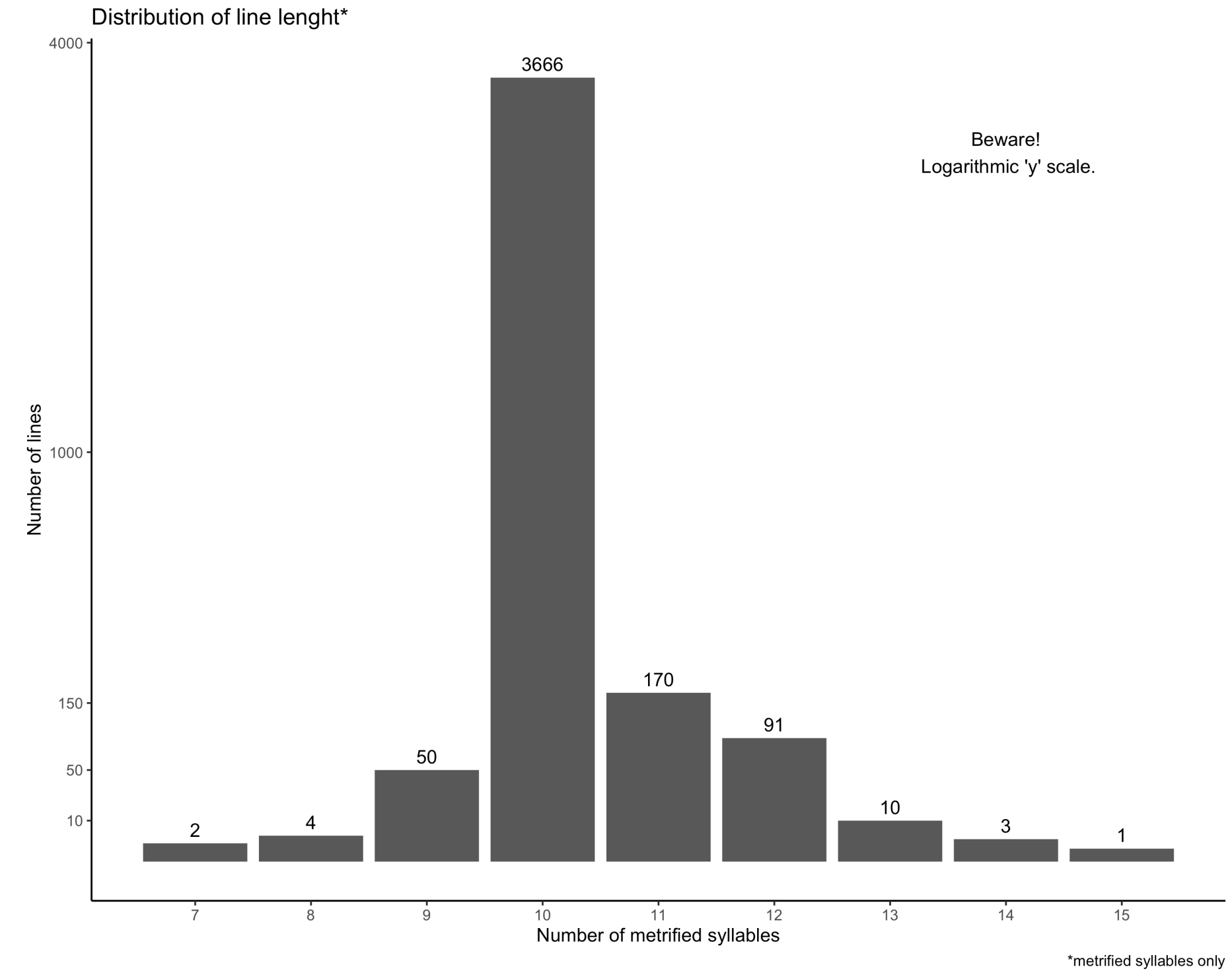
1 *seul* script

Économie du *workflow*

PAM

Un script pour toutes les données

- Le PAM produit un fichier par texte
- Corpus :
 - des dizaines de textes
 - plusieurs graphes à produire pour chaque texte
- Un seul script suffit !



PAM

Un script pour toutes les données

Plusieurs graphiques pour un texte

- Dans *un seul script* :
 - importer toutes les données nécessaires (1)
 - appliquer toutes les transformations nécessaires (2)
 - produire toutes les représentations voulues (3)

(1)

Import df.A

Import df.B

(2)

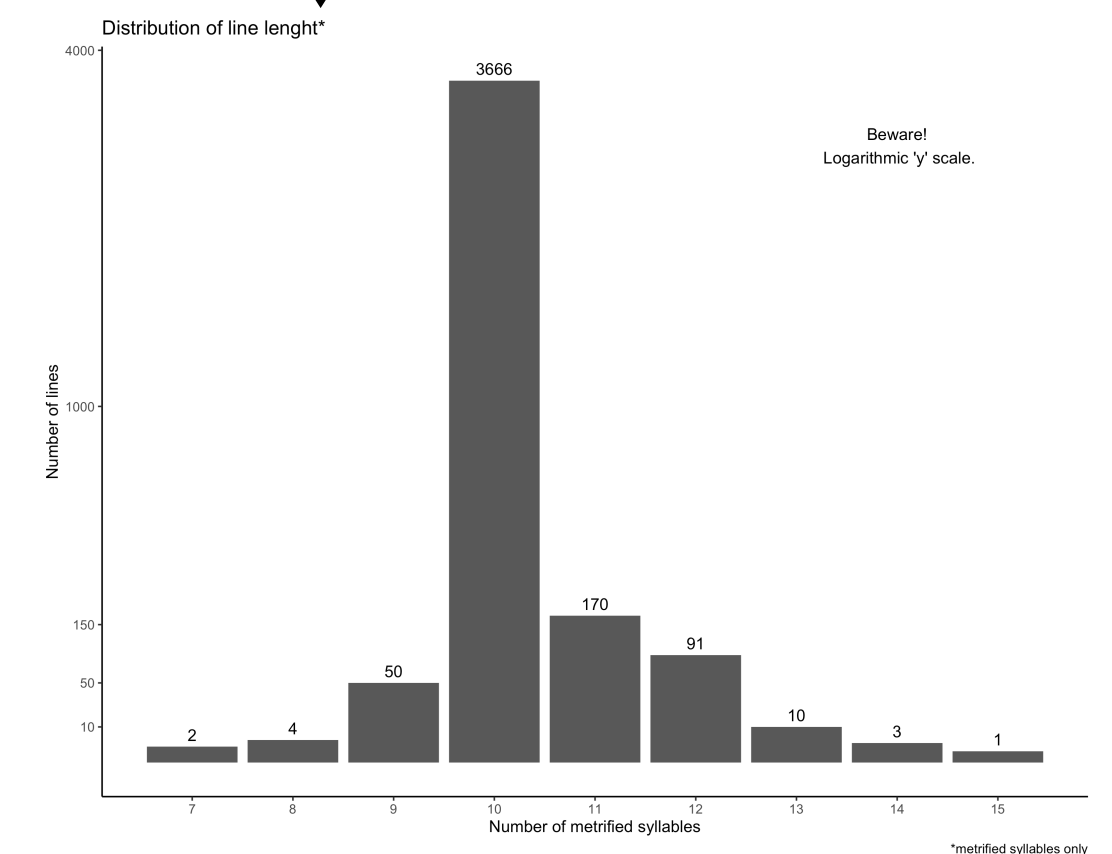
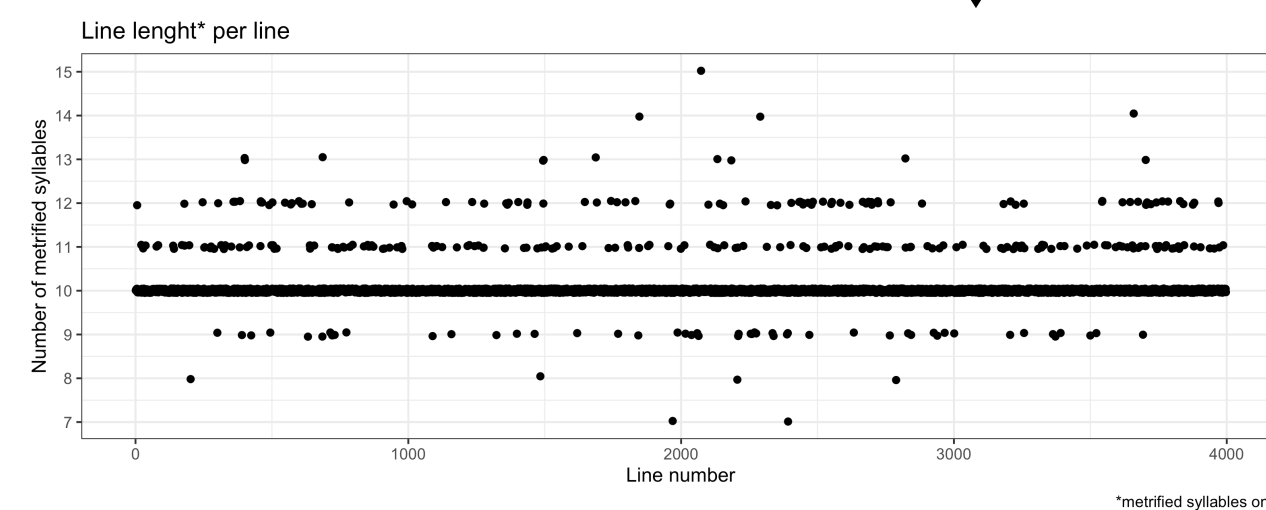
Transform A

Transform B

(3)

Plot α

Plot β

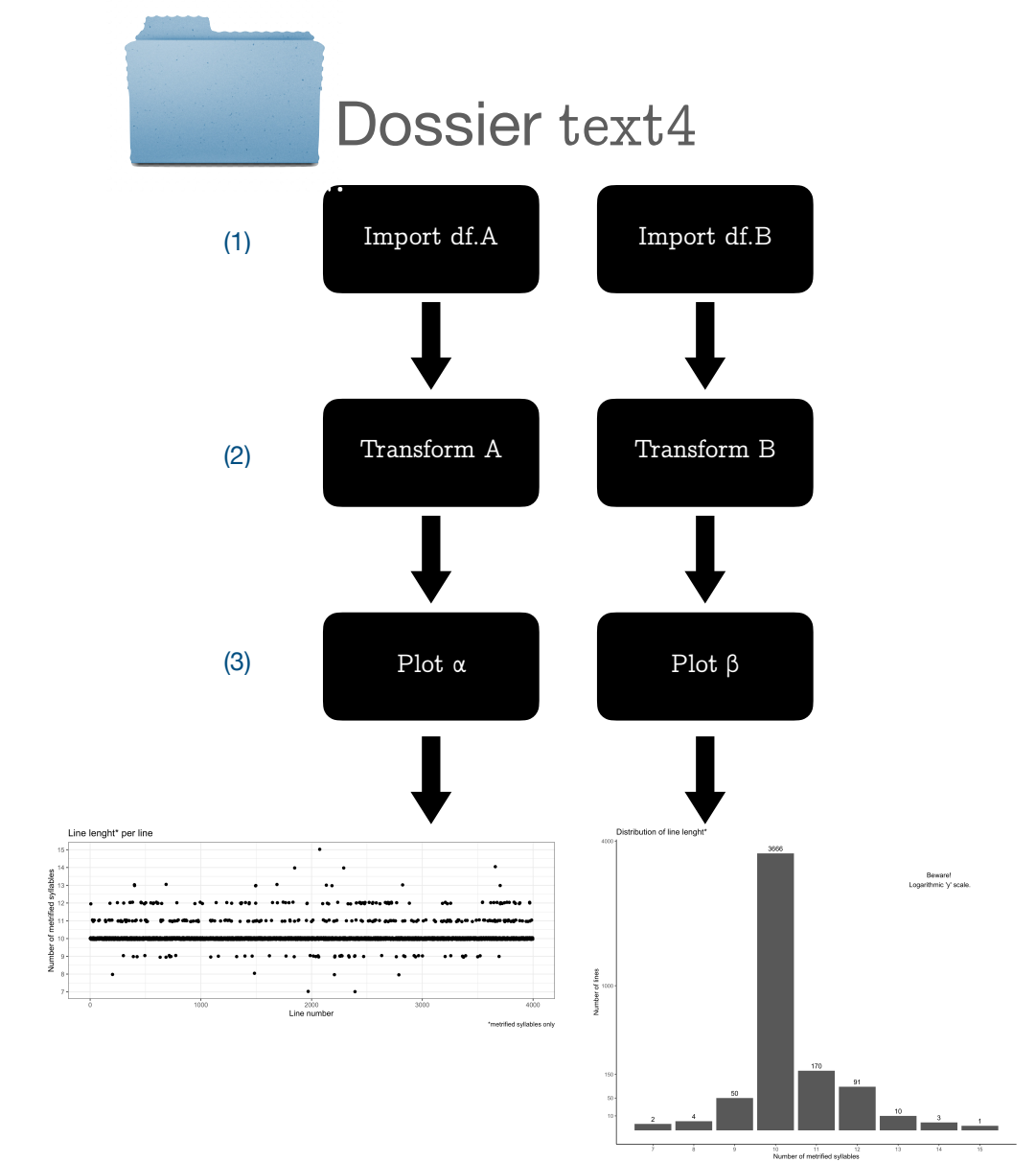
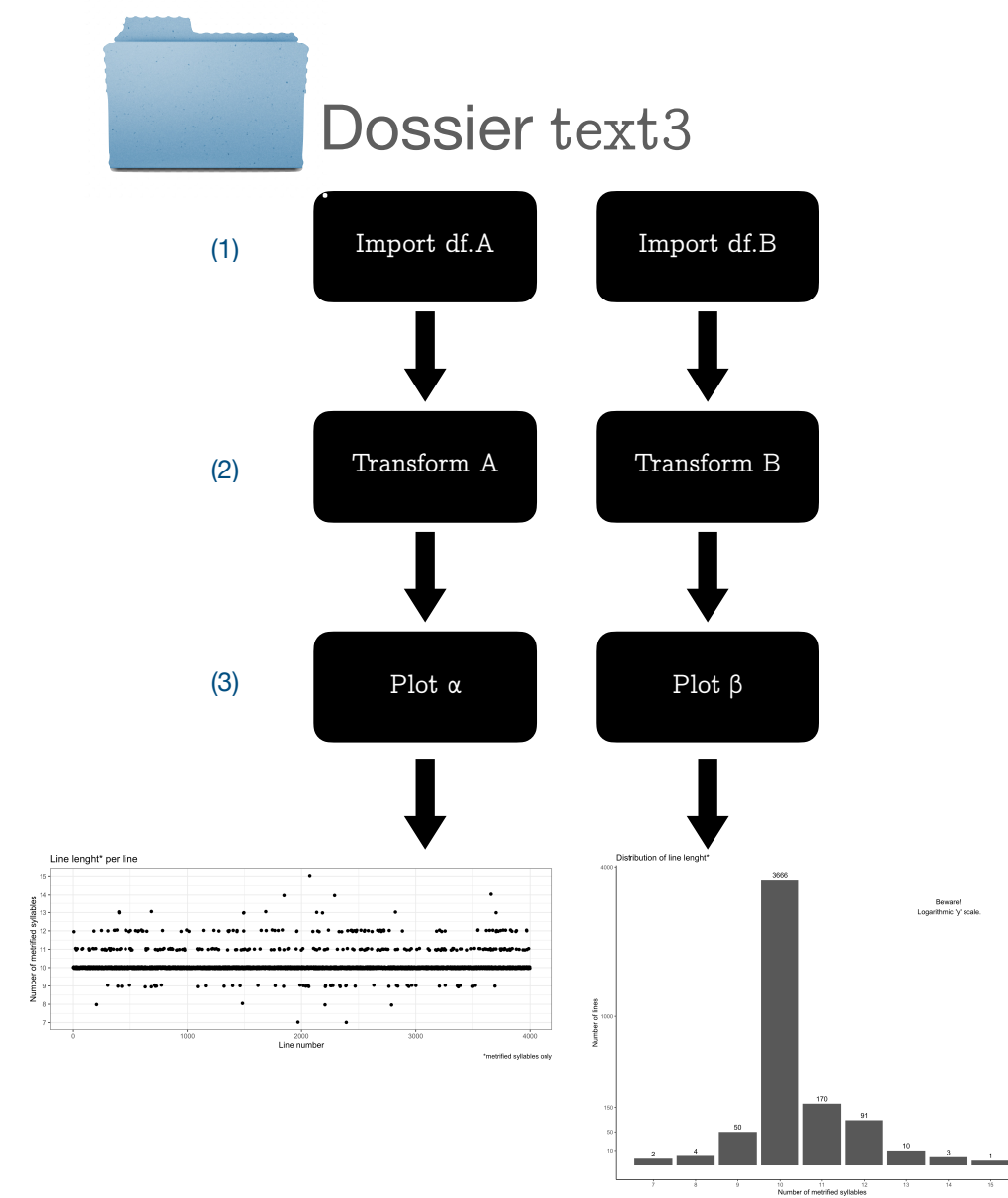
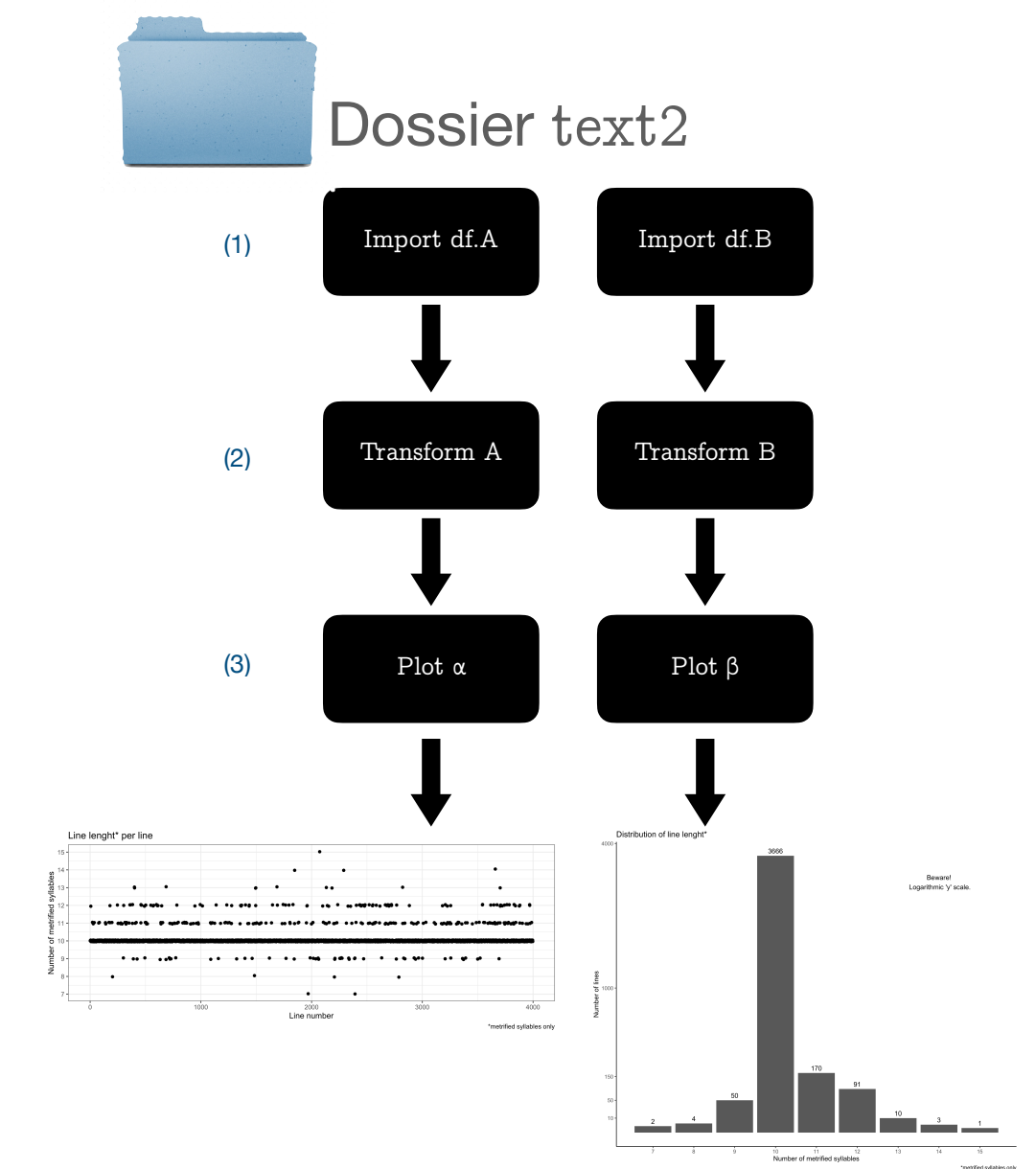
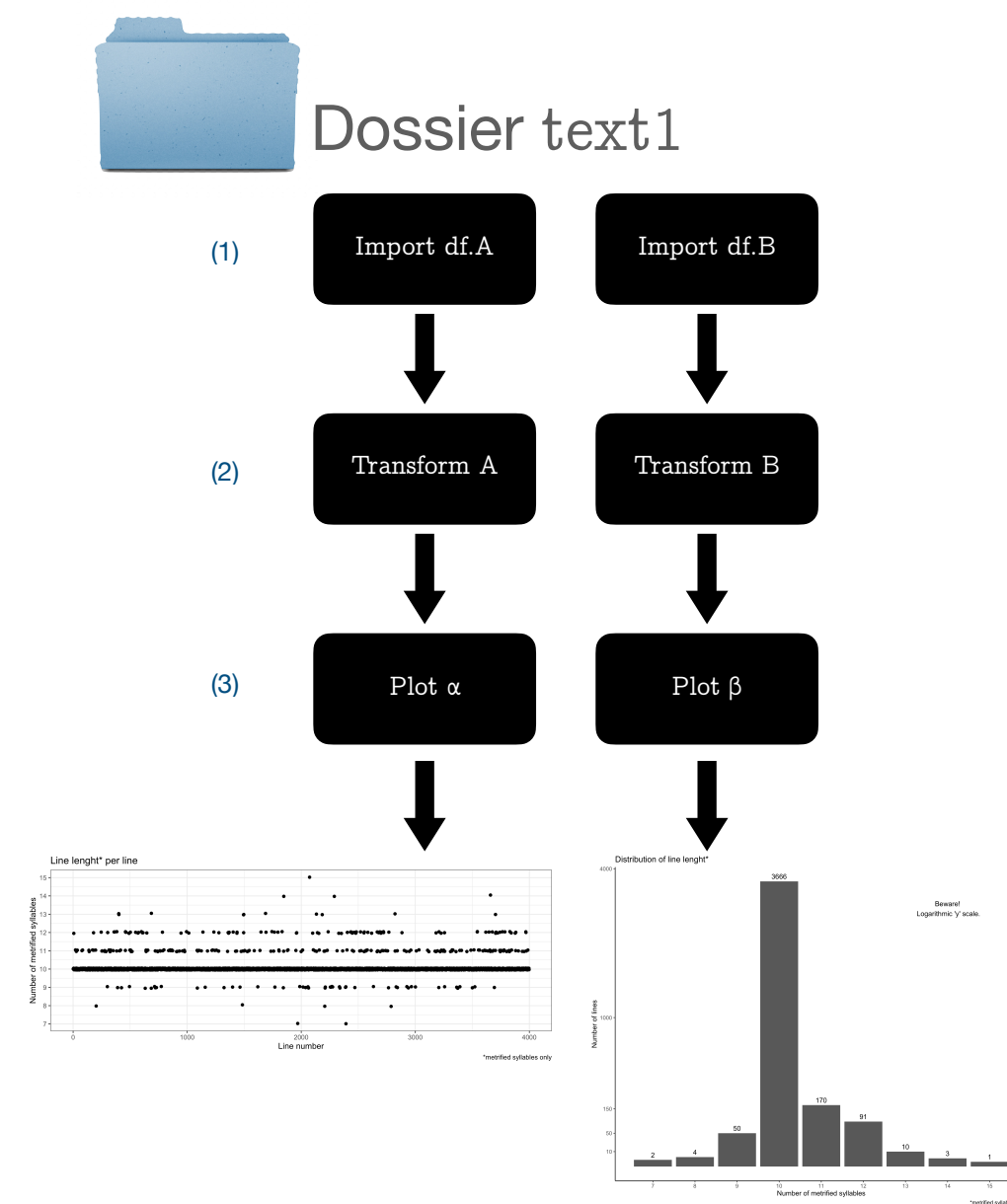


PAM

Un script pour toutes les données

Plusieurs graphiques pour plusieurs textes

- Créer un répertoire par texte
- Placer le script dans chaque répertoire
- Placer les fichiers de données de chaque texte dans son répertoire
 - /! Ils doivent avoir toujours le même nom (p.ex. pas d'ID du texte dans le nom)
- Et lancer le script à chaque fois



PAM

Un script pour toutes les données

Automatisation du processus :

- pas d'erreurs de manipulations
- pas de centaines de manipulations répétitives
- très rapide !

Conclusion

Représenter des données linguistiques avec R

- Représenter la variation de données permet
 - de ne pas se laisser tromper par leur apparence superficielle
 - de ne pas être perdu dans leur détail fin
- R est un bon outil pour faire cela
 - Cartes, graphiques
 - Beaucoup de ressources en ligne, une forte communauté
 - Il permet une automatisation du processus → énorme gain de temps
- Face à des données difficiles à représenter, il faut être inventif, et R permet cela la plupart du temps

Merci de votre attention !

Chanson de Roland, édité par Gérard Moignet, Paris, Bordas, 1972. Publié en ligne par la Base de français médiéval, <http://catalog.bfm-corpus.org/roland>. Dernière révision le 2014-02-24.

ANDERSON, Stephen R. (1982). « The Analysis of French Shwa: Or, how to Get Something for Nothing », in : *Language* 3(58), pp. 534-573.

CHOUVION, Sophie, Timothée PREMAT & Axelle VERNER (éval. en cours). "La musication dans les chansons d'Henry Fresneau, Grammaire d'association, parophonologie et forme polyphonique au XVIe siècle", in : *Textus & Musica*.

GUILLIÉRON, Jules & Edmond EDMONT (1902-1910). *Atlas Linguistique de la France*, Paris : Champion.

GUILLOT-BARBANCE, Céline, Serge HEIDEN & Alexei LAVRENTIEV (2017). "Base de français médiéval : une base de référence de sources médiévales ouverte et libre au service de la communauté scientifique". In : *Diachroniques* 7, pp. 168-184.

POGGIO, Enzo & Timothée PREMAT (2019). "Le PAM, un Programme d'Analyse Métrique pour le français médiéval », in : *Actes des Rencontres Lyonnaises des jeunes chercheurs en linguistique historique*, sous la dir. de Timothée Premat & Ariane Pinche, Lyon : Diachronies contemporaines, pp. 59-70.

PREMAT, Timothée & Philippe BOULA DE MAREÜIL (2019), "The "rolled" /R/ in French and some regional languages of France", présentation à *R-atics* 6, Paris, 07/11/2019.

RUSSO, Michela & Timothée PREMAT (accepté ; à paraître). "Voyelles finales et traits-φ à la rencontre des diasystèmes d'oïl, d'oc et du francoprovençal", in : *Verbum*, "Modélisation diasystémique et typologie", numéro sous la dir. de Jean Léo Léonard.

STEIN, Achim, Pierre KUNSTMANN & Martin-D. GLEBGEN (2006). "Nouveau Corpus d'Amsterdam. Corpus informatique de textes littéraires d'ancien français (ca 1150-1350)". Établi par Anthonij Dees (Amsterdam : 1987), remanié par Achim Stein, Pierre Kunstmann et Martin-D. Gleßgen, Stuttgart : Institut für Linguistik/Romanistik.

Plus de scripts et d'exemples sur mon site :
<https://sites.google.com/view/timothee-premat/programmes>

Diapo : [〈halshs-03116096〉](#)