

Traitemen^t automatique du français en diachronie, retour d'expérience sur le projet Presto

Sascha Diwersy, U. Montpellier, Praxiling

Achille Falaise, ENS de Lyon, ICAR

Gilles Souvay, U. Lorraine, ATILF

CLPS 2016
9-11 mars 2016, Lyon



Deutsche
Forschungsgemeinschaft

Le lexique PRESTO

- ▶ Lexique adapté à la langue du XVI^e-XVII^e
 - orthographe non stabilisée : variation et flexion
 - segmentation
- ▶ Archaïsation du lexique moderne et projection sur un corpus textuel
 - archaïsation par règles
 - base de connaissances du lemmatiseur LGeRM
 - vérification du taux de couverture sur Frantext
 - forme, fréquence, tranche chronologique
 - expérience projet Européen IMPACT

Le lemmatiseur LGeRM

- ▶ LGeRM : Lemmes Graphies et Règles Morphologiques
 - lemmatiseur et environnement de lemmatisation
 - gestion de différents états de langues
 - <http://www.atilf.fr/LGeRM>
- ▶ Liste de formes connues : lexique morphologique
 - gestion de « flexion et variantes » dans Frantext
- ▶ Règles d'analyse des formes inconnues
 - 6 500 règles (4/5 flexion verbale)
 - si (en finale) alors ES → EFS finsi *nes* → *nefs*, NEF
 - Y → I *fayre* → *faire*, FAIRE
 - si (entre voyelles) alors C → SS finsi *mesfacent* → *mesfassent*, MÉFAIRE

Construction du lexique

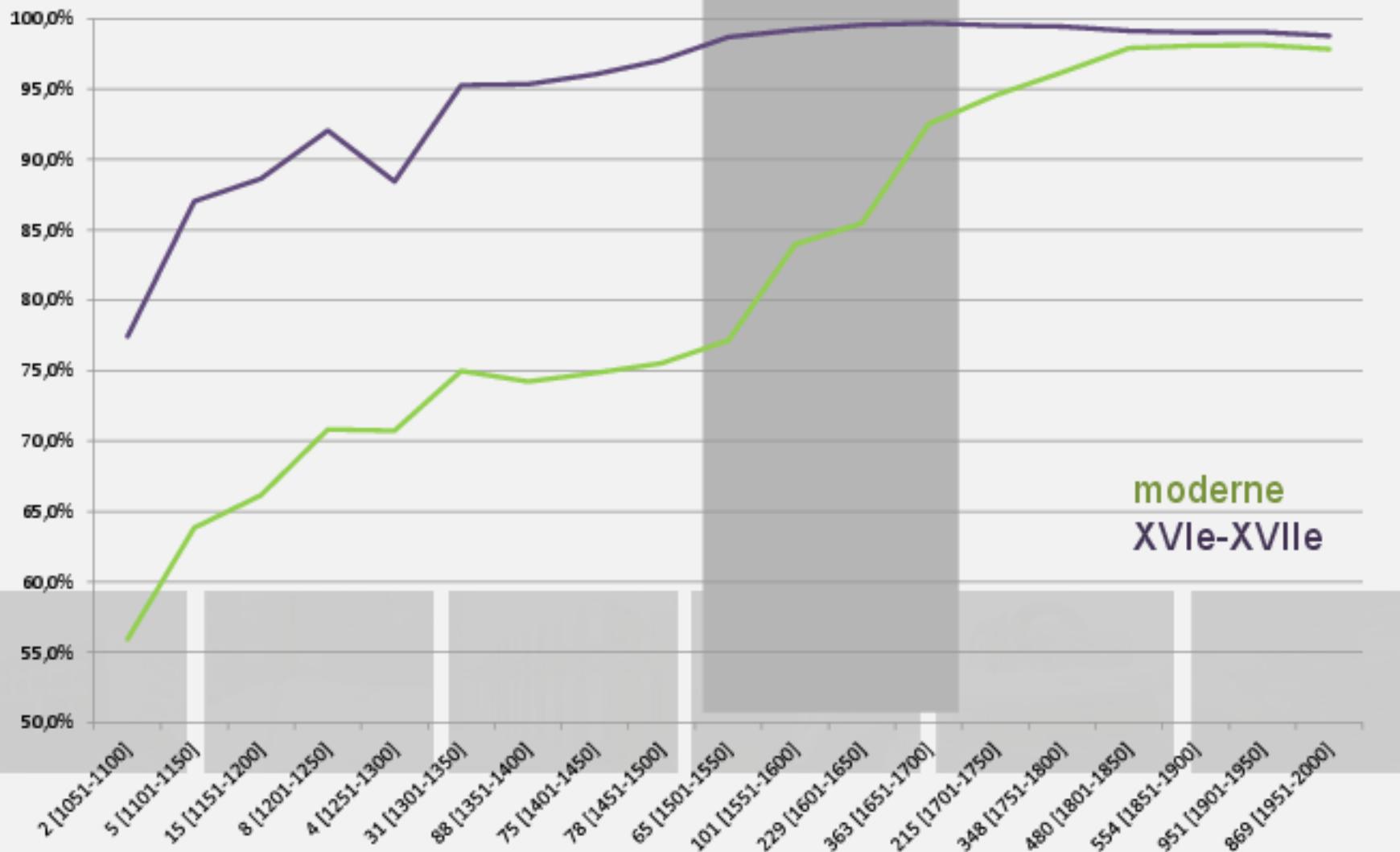
► Lexique moderne LEFFF

- étiquettes spécifiques au projet
- lemmes complémentaires

► Construction

- appliquer les règles d'archaïsation (3 boucles)
- regarder les formes absentes
- réitérer le processus
- compléter avec les formes absentes

Taux de couverture du lexique



Perspectives

- ▶ On n'aura jamais toutes les variantes
 - intégration du lemmatiseur dans la chaîne de traitement ?
- ▶ Produire un lexique adapté à un état de langue, à un type de typographie
 - origine de la forme
 - fréquence des mots

Le corpus Presto

- Besoins
 - Représentation de toutes les périodes de l'histoire du français : 9^{ème} s. au 21^{ème} s.
 - Dans un premier temps : 16^{ème} s. au 21^{ème} s.
 - Présence de différents types de textes et de différents genres discursifs : narratif, poésie, théâtre, traité
 - Enrichissement linguistique : étiquetage morpho-syntaxique et lemmatisation
- Collaboration avec diverses bases textuelles existantes
 - *Frantext*
 - *BVH (Bibliothèques Virtuelles Humanistes)*
 - *Université de Cologne*
 - *ARTFL (American and French Research on the Treasury of the French Language)*
 - *CÉPM (Corpus Électroniques de la Première Modernité)*
 - ...

Le corpus Presto

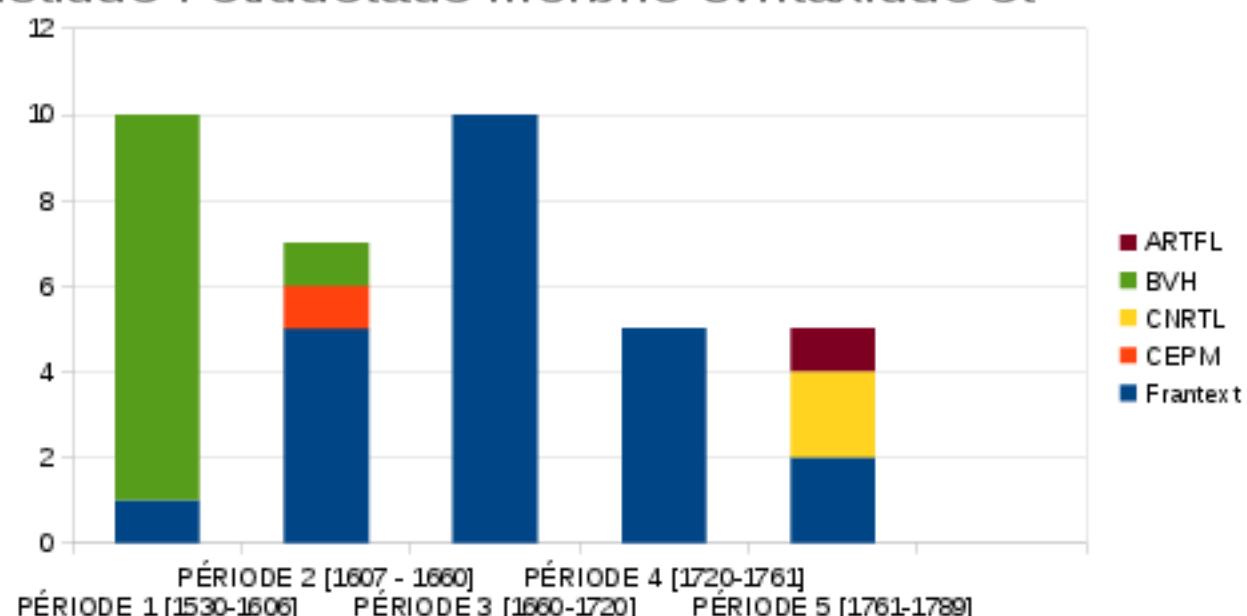
- Besoins
 - Représentation de toutes les périodes de l'histoire du français : 9^{ème} s. au 21^{ème} s.
 - Dans un premier temps : 16^{ème} s. au 21^{ème} s.
 - Présence de différents types de textes et de différents genres discursifs
 - Enrichissement linguistique : étiquetage morpho-syntaxique et lemmatisation
- Collaboration avec diverses bases textuelles existantes
 - Frantext
 - BVH (*Bibliothèques Virtuelles Humanistes*)
 - Université de Cologne
 - ARTFL (*American and French Research on Texts in Linguistics*)
 - CÉPM (*Corpus Électroniques de la Première Guerre mondiale*)
 - ...

392 textes (1509-2010), dont 53 libres :

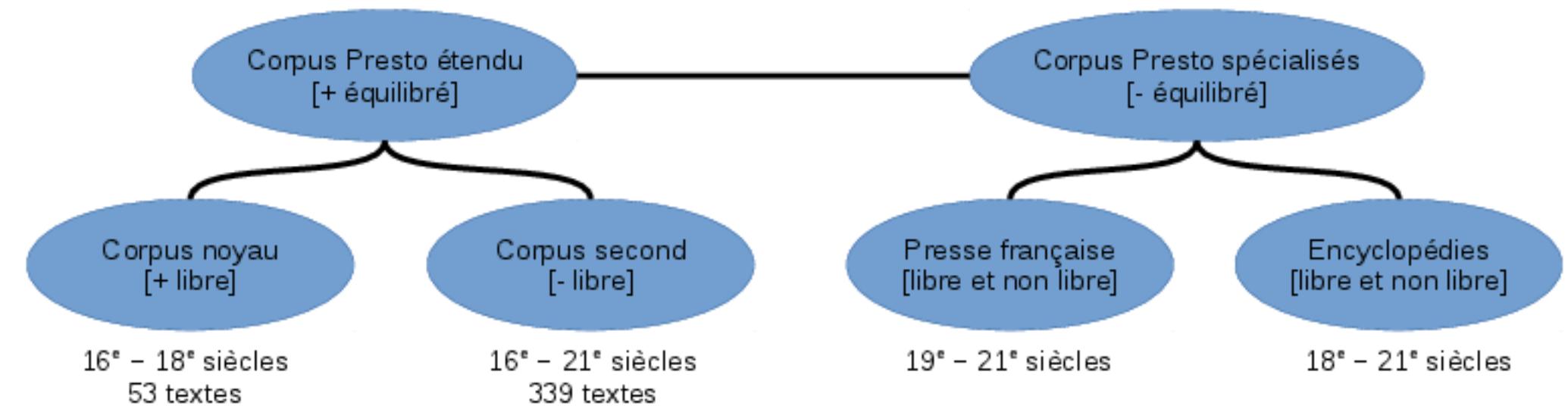
L'Astrée, Gargantua, l'Encyclopédie (vol. 7)...

Le corpus Presto

- Besoins
 - Représentation de toutes les périodes de l'histoire du français : 9^{ème} s. au 21^{ème} s.
 - Dans un premier temps : 16^{ème} s. au 21^{ème} s.
 - Présence de différents types de textes et de différents genres discursifs
 - Enrichissement linguistique : étiquetage morpho-syntaxique et lemmatisation
- Collaboration avec diverses bases
 - Frantext
 - BVH (*Bibliothèques Virtuelles Humaines*)
 - Université de Cologne
 - ARTFL (*American and French Research Project in Linguistics*)
 - CÉPM (*Corpus Électroniques de la Philosophie*)
 - ...



Le corpus Presto



Extrait

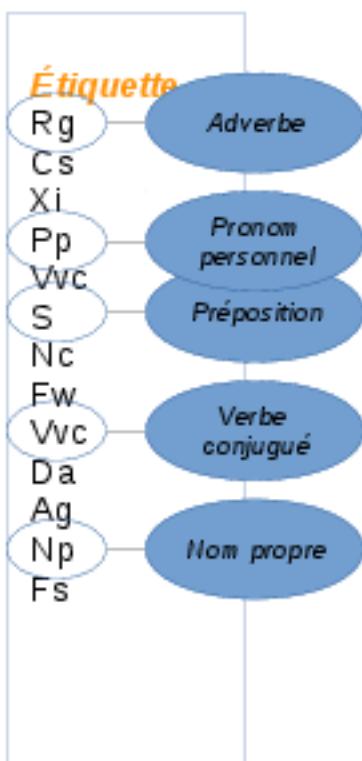
(*Sottie pour le cry de la bazoche, Anonyme, 1549*)

C'est assez dict pour ceste foys.
Quand sçavoir en vous s'assocye,
Monsieur Rien, l'on vous remercye
Du bien qu'avons aprins de vous.
Bazochiens, entendez tous :
Je veulx en triumphant arroy
Eslire et faire ung nouveau roy,
Comme il est coustume de faire ;
Pourtant chacun pense a l'affaire,
Autant les grandz que les petitz,
Et faire les preparatifz ;
Car, ainsi comme liberalle,
Je tendz a monstre generalle
Qui, l'esté qui vient, sera faicte.
En honneur du triumphe et feste,
Ne faillez monstrer vos bons cueurs
Qui font de la vertu approche,
Tant que l'on dye par honneurs :
Vive l'excellente Bazoche !

Annotation

Forme

Tant
que
l'
on
dye
par
honneurs
:
Vive
l'
excellente
Bazoche
!



Désambiguation par
TreeTagger.
Précision 94 %.

Lemme

TANT
QUE
L
ON
DIRE
PAR
HONNEUR
:
VIVRE
LE
EXCELLENT
BAZOCHE
!

Lemme moderne.
Désambiguation et
évaluation à venir.

Sème

Analysé sémantique: en
cours pour les noms propres
(noms de personnes vs noms
de rivières vs noms de villes,
etc.)

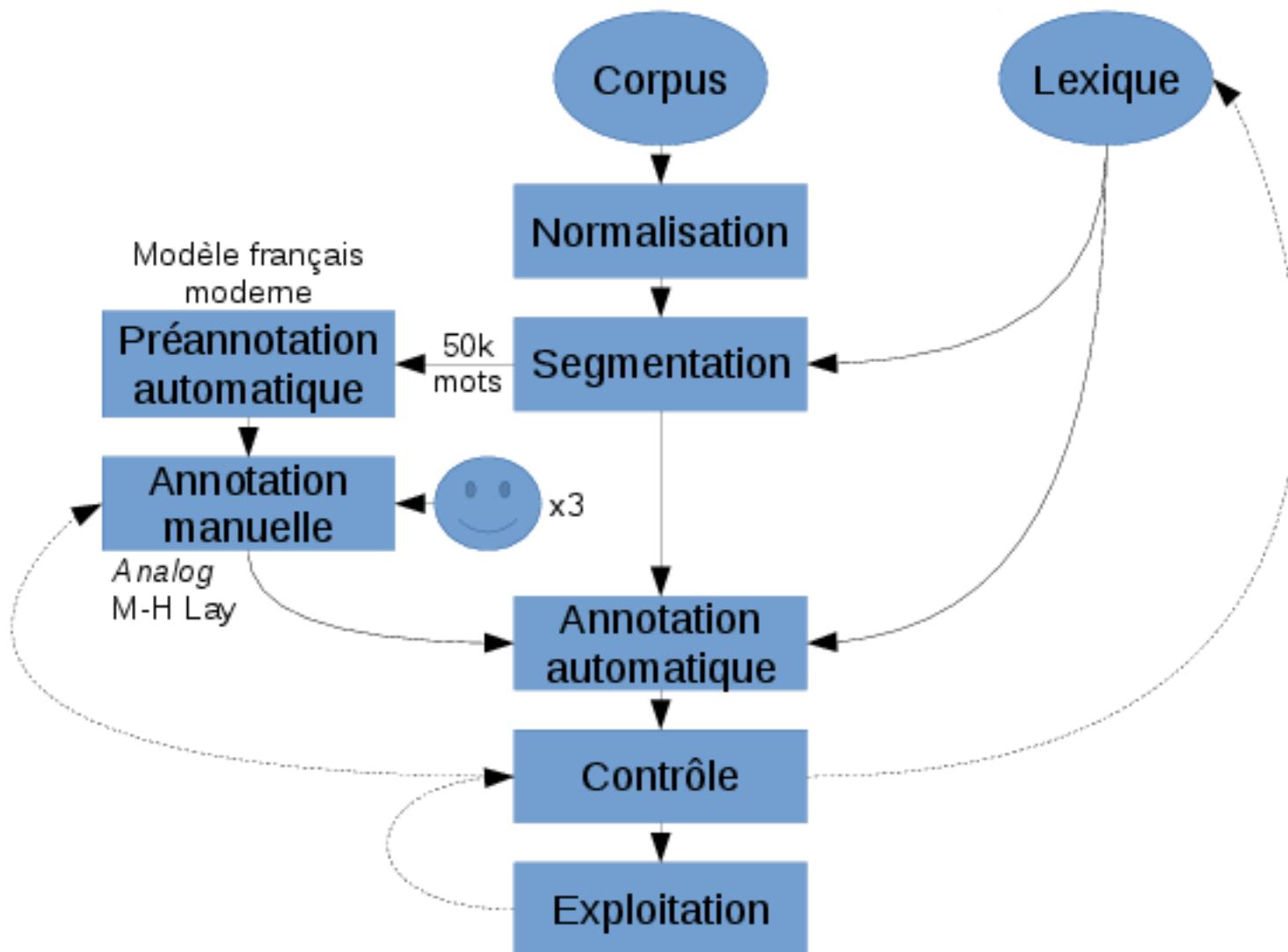
Dépendance

Analysé syntaxique de
dépendances:
à venir

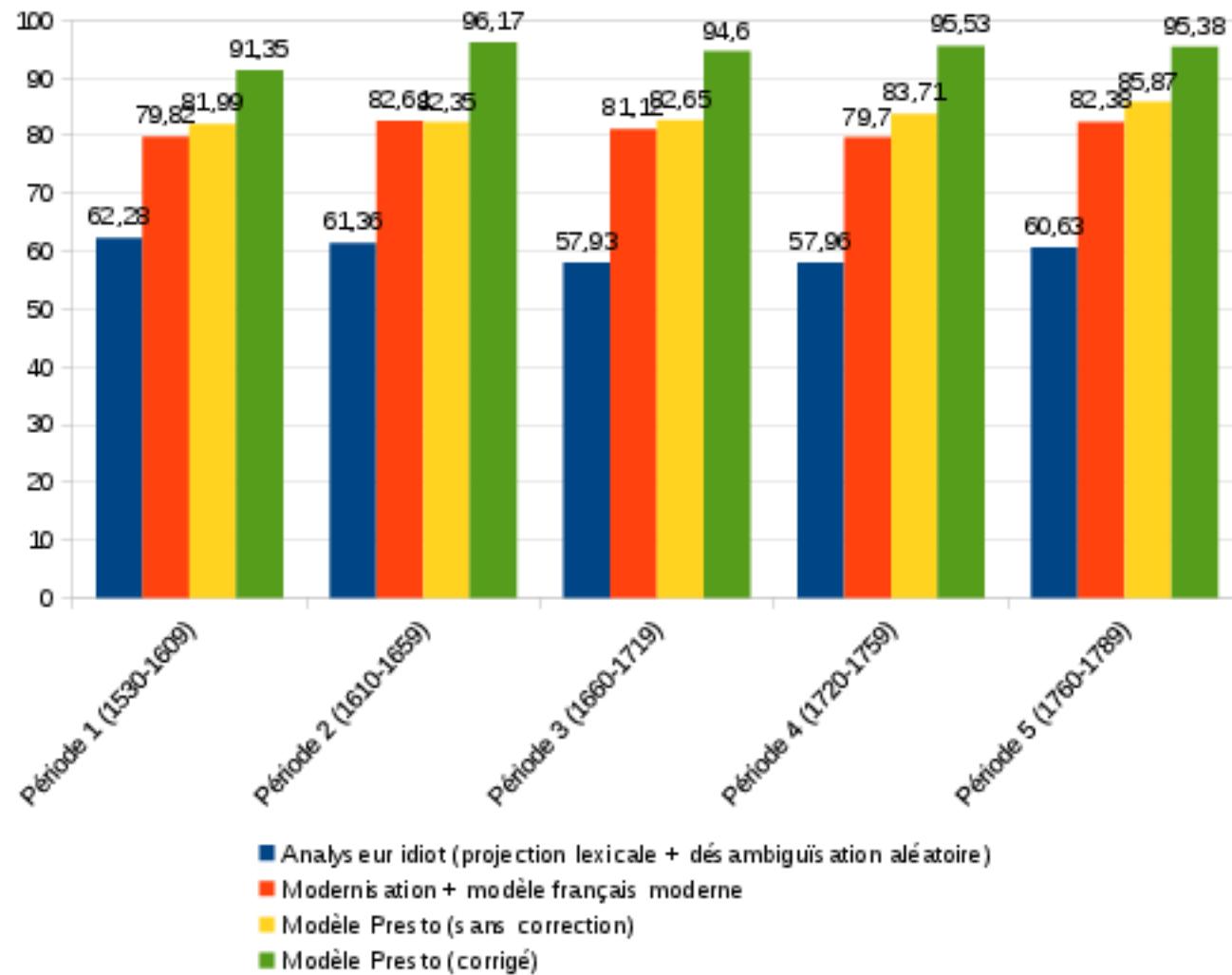
Jeu d'étiquettes

Étiquette (partie du discours 1)	Étiquette (partie du discours 2)	Flexion
Nom (N)	commun, propre	
Verbe (V)	être/avoir, autre	conjugué, infinitif
Adjectif (A)	général, possessif	
Pronom (P)	personnel, démonstratif, indéfini, possessif, interrogatif, relatif	
Déterminant (D)	article défini, démonstratif, article indéfini, article partitif, indéfini, relatif, interrogatif/exclamatif	
Participe-Adjectif-Gérondif (G)	part_présent/adjectif_verbal/gérondif, part_passé/adjectif_verbal	
Adverbe (R)	général, particule, interro-exclamatif	
Adposition (S)		
Conjonction (C)	coordination, subordination	
Numéral (M)	cardinal, ordinal	
Interjection (I)		
Résidu (X)	abréviation, mot étranger, symbole, préfixe, consonne intercalée	
Ponctuation (F)	forte, faible, autre	

Traitement



Précision



Rappel : Principes d'échantillonage



Rappel : Principes d'échantillonage

- Sélection équilibrée de textes selon différents critères (période, genre, auteur)
 - → [+ équilibré] : Corpus de référence PRESTO
 - → [- équilibre] : Corpus spécialisés
- Principes supplémentaires (Corpus de référence PRESTO):
 - statut juridique ([+ libre] / [- libre] de droits)
- Principes supplémentaires (Corpus spécialisés)



Modèles de données et méthodes



Modèles de données

Corpus de référence PRESTO

- Unités structurelles: texte ou recueil de textes, phrase
- Métadonnées:
 - texte : titre, auteur, date, genre
- Propriétés lexicales (mod. PRESTO): mot-forme, lemme, catégorie (2 degrés de granularité), position textuelle (mot → phrase [WS])



Modèles de données

Corpus diachronique de la presse française

- Echantillons de textes journalistiques parus à la fin du 19e et du 20e siècle; différenciation supplémentaire entre presse nationale et régionale
 - 19e s.: **Le Figaro**, le JDD (19e s.), **Le Reveil de Lyon**, etc.
 - 20e s.: **Le Figaro**, **Le Monde**, **L'Est Républicain**, **Sud Ouest**
- Unités structurelles: numéro de journal, phrase
- Métadonnées:
 - numéro : journal, année
- Propriétés lexicales (mod. *TreeTagger* fr.contemp.): mot-forme, lemme, catégorie (2 degrés de granularités), position textuelle (WS)



Modèles de données

Corpus de la presse française contemporaine

- Echantillons de textes journalistiques parus à la fin du 19e et du 20e siècle; corpus de test / de démo
 - Le Monde (+ Le Figaro, Libération, Sud Ouest, Ouest France etc)
- Unités structurelles: année de journal, article, paragraphe, phrase
- Métadonnées:
 - article : titre, auteur, date, section thématique
 - paragraphe : position textuelle (PT)
 - phrase : position textuelle (ST, SP)
- Propriétés lexicales (mod. *Connexor*): mot-forme, lemme, catégorie (2 degrés de granularités), traits morphologiques, relation de dépendance syntaxique, fonction syntaxique, positions textuelles (WT, WP, WS)



Modèles de données

Corpus des Encyclopédies

- Articles de l'*Encyclopédie de Diderot et d'Alembert* (18e s.) et de l'*Encyclopaedia Universalis* (20e s.)
- Unités structurelles: ouvrage, article, paragraphe, phrase
- Métadonnées:
 - article : lemme (entrée), auteur, domaine thématique
 - paragraphe : position textuelle (PT)
 - phrase : position textuelle (ST, SP)
- Propriétés lexicales (mod. PRESTO): mot-forme, lemme, catégorie (2 degrés de granularités), positions textuelles (WT, WP, WS)



Modèles de données

Bilan – Propriétés lexicales

Base textuelle	Forme	Catégorie	Lemme	Traits	Relation syntaxique	WT	WP	WS
Corpus de référence PRESTO	X	X	X	-	-	-	-	X
Corpus diachronique de la presse française	X	X	X	X	-	-	-	(x)
Presse française contemporaine	X	X	X	X	X	X	X	X
Encyclopédies	X	X	X	-	-	X	X	X
Emolex-FR	X	X	X	X	X	X	X	X
CoVaNa-FR	X	X	X	X	X	-	-	-



Visées méthodologiques

Lexical Priming (cf. Hoey 2005) – Modèle descriptif

Niveau cooccurrentiel	Niveau de schématicité distributionnelle
cooccurrences lexicales	
associations sémantiques / pragmatiques	Profil combinatoire
colligations	
collocations textuelles (chaînes lexicales)	
associations sémantiques textuelles (argumentatives)	Profil d'intégration textuelle
colligations textuelles	



Visées méthodologiques

Lexical Priming – Mise en oeuvre du modèle descriptif

Base textuelle	LCOLLOC	LSEM	LCOLLIG	TCOLLOC	TSEM	TCOLLIG
Corpus de référence PRESTO	x	(x)	(x)	-	-	(x)
Corpus diachronique de la presse française	x	(x)	(x)	-	-	(x)
Presse française contemporaine	x	x	x	-	-	x
Encyclopédies	x	(x)	(x)	-	-	x
Emolex-FR	x	x	x	-	-	x
CoVaNa-FR	x	x	x	-	-	-



Outils (*PrimeStat.BTLC*)



3.1 Boîte à outils *PrimeStat*

- Concordances KWIC
- Index (fréquences, spécificités)
- Analyse coocurrentielle (cooccurrences lexicales, lexico-syntaxiques, colligations syntaxiques, colligations textuelles)
- Analyses multivariées (CAH, K-Means, AFC)



3.2 Plateforme BTLC

Adresse du site: <http://persan.rom.uni-koeln.de/btlsc/>

BTLC1.0

BTLC

Se connecter à la base

Crédits

Se connecter à la base textuelle

Nom d'utilisateur :

Mot de passe :

Daten absenden

Fermer la fenêtre

(c) 2013, L'équipe BTLC



BTLC.PrimeStat – Interface graphique

- Définition d'objets de requête
 - Corpus de travail (sous-corpus / partition)
 - Expression de requête (→ pivot)
- Applications
 - Concordances KWIC
 - Calculs de fréquence (spécificités)
 - Cooccurrences lexico-syntactiques
 - etc.



3.2 BTLC.PrimeStat – Interface graphique

Définition du corpus de travail

The screenshot shows the 'Choix du corpus' (Choice of corpus) tab selected in the top navigation bar. Below it, there is a section titled 'Définir le corpus de travail' with a link to '[Aide]'. Two buttons are visible: 'Créer ou Ajouter un sous-corpus' and 'Créer un échantillon à partitions'. A red box highlights the main input field where 'Sous-corpus: test' is entered, along with a small trash icon. Another red box highlights the 'Selectionner un ou plusieurs fichiers de corpus' (Select one or more corpus files) button and its associated text input field, which contains the placeholder 'Veuillez choisir un fichier de corpus...'.



3.2 BTLC.PrimeStat – Interface graphique

Définition du corpus de travail

Sous-corpus : test X

Sélectionner un ou plusieurs fichiers de corpus [\[Aide\]](#) Presse ivoirienne X Presse sénégalaise X

Critères de sélection [\[Aide\]](#)

- ▶ Code corpus
- ▶ Sous-échantillon
- ▶ Code sous-échantillon
- ▶ Titre du texte
- ▶ Auteur
- ▶ Date
- ▶ Section

Liste des descripteurs disponibles



3.2 BTLC.PrimeStat – Interface graphique

Définition de l'expression de requête

Définir l'expression de requête [Aide]

Mot [+ Mot]
Lemme
[+ Crit.]

← Boîte de paramétrage



3.2 BTLC.PrimeStat – Interface graphique

Définition de l'expression de requête

Définir l'expression de requête [\[Aide\]](#)

Mot	Mot
Catégorie 	Catégorie 
 <input type="text" value="DET"/>	<input type="text" value="N"/>
[+ Crit.]	[+ Crit.]
Multiplier : <input type="text" value="1"/> à <input type="text" value="1"/> fois.	



3.2 BTLC.PrimeStat – Interface graphique

Applications: Concordances KWIC

Concordances KWIC Fréquences Cooccurrences

Définir les paramètres de tri de la concordance [Aide]

Trier par:	Mot-forme
Empan:	3G...1G
Afficher la concordance KWIC	



3.2 BTLC.PrimeStat – Interface graphique

Applications: Calculs de fréquence

The screenshot shows a software interface with a top navigation bar containing three tabs: "Concordances KWIC" (selected), "Fréquences" (highlighted in blue), and "Cooccurrences". Below the tabs is a section titled "Définir les paramètres de calcul" with a link "[Aide]". A dropdown menu under "Attributs:" contains the value "Lemme + Catégorie". At the bottom is a large button labeled "Calculer la table de contingence".



3.2 BTLC.PrimeStat – Interface graphique

Applications: Cooccurrences lexico-syntactiques

Définir les paramètres de calcul [\[Aide\]](#)

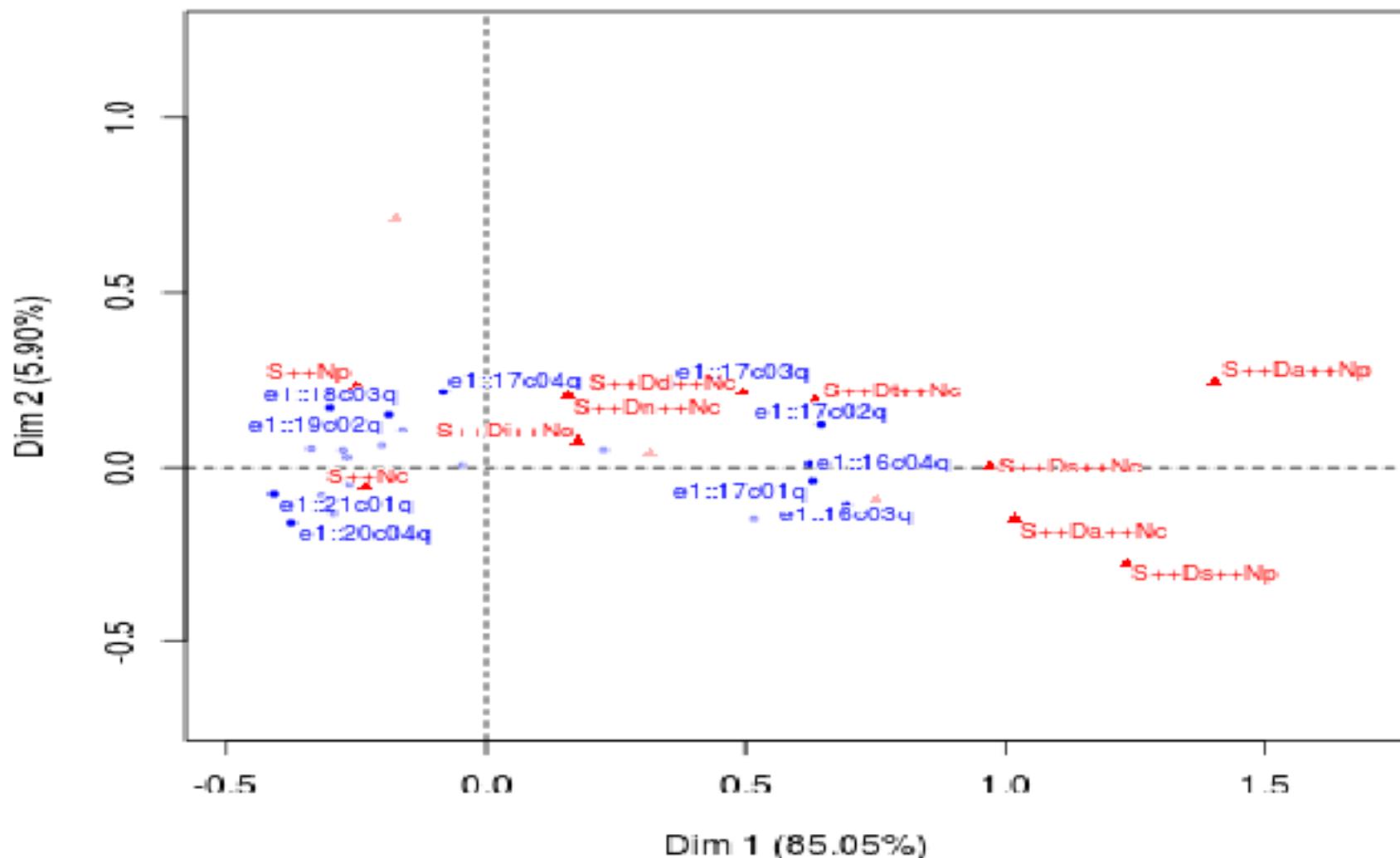
Attributs des collocatifs :	Lemme + Catégorie
Type de relation coocurrentielle :	dépendance
Calculer la table de contingence	



3.2 BTLC.PrimeStat – Applications

AFC – *en* (+ D) + N

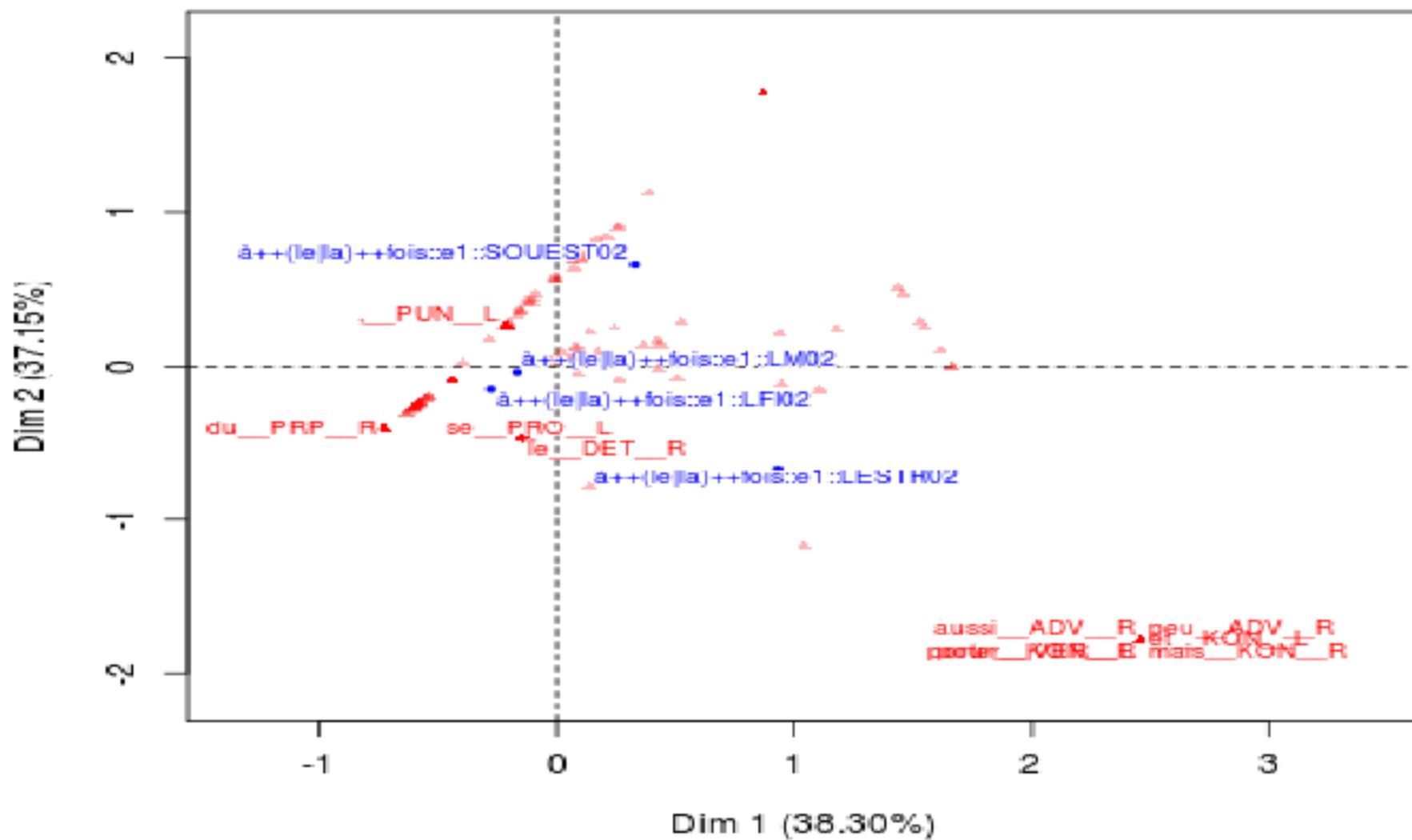
CA factor map



3.2 BTLC.PrimeStat – Applications

AFC – à la fois (Presse 20e s.)

CA factor map





Perspectives

-

Vers une linguistique diachronique
outillée



Vers une linguistique diachronique outillée

Méthodes d'exploitation de corpus dans une perspective diachronique

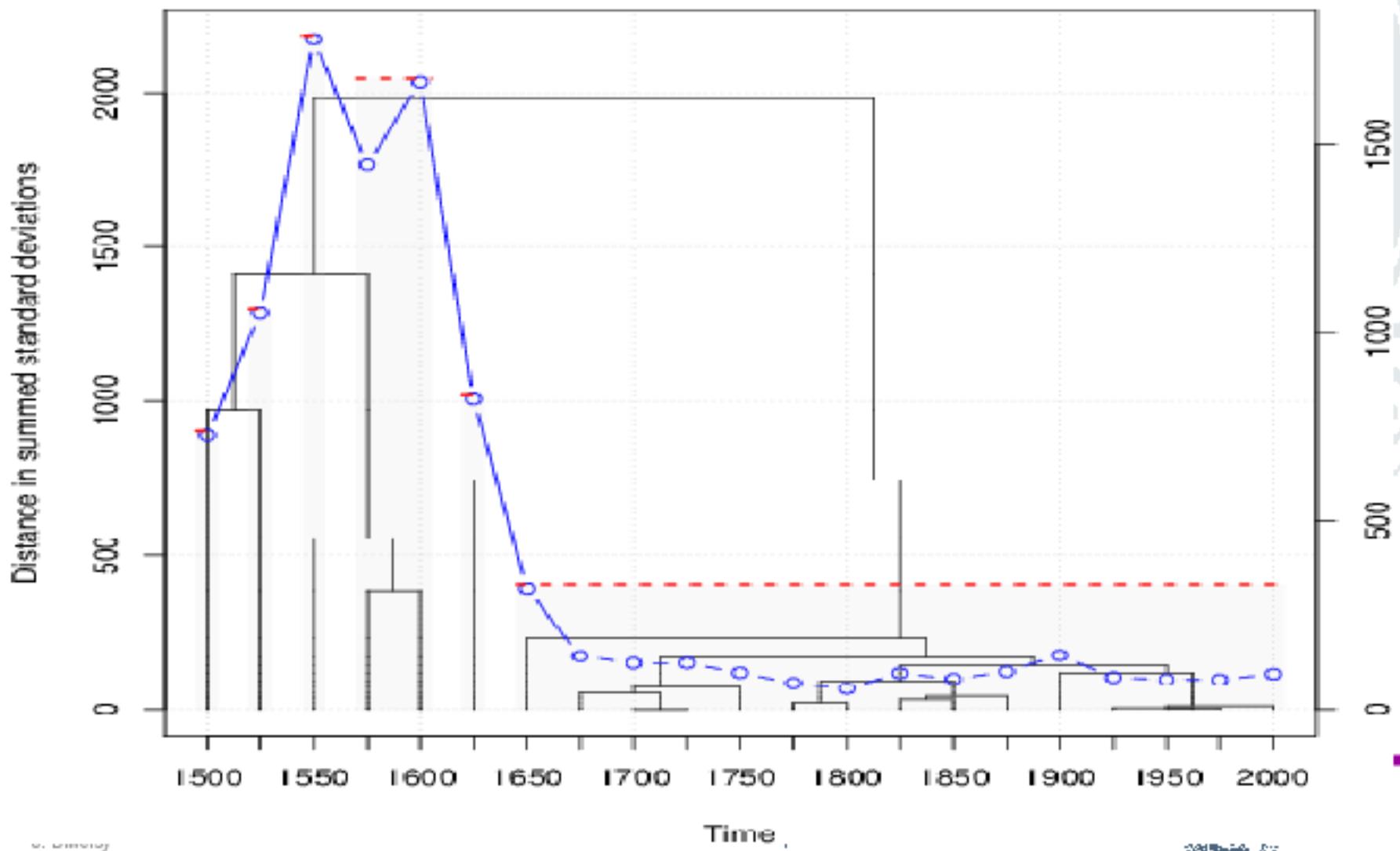
- VNC (*Variability-based Neighbour Clustering*: Gries & Hilpert 2008, Hilpert & Gries 2009, Hilpert 2013)
- Calculs de productivité (Hilpert 2012)
- HCFA (*Hierarchical Configural Frequency Analysis*: Gries 2004, Hilpert 2013)



Vers une linguistique diachronique outillée

VNC appliqué sur un motif : *en + D + N*

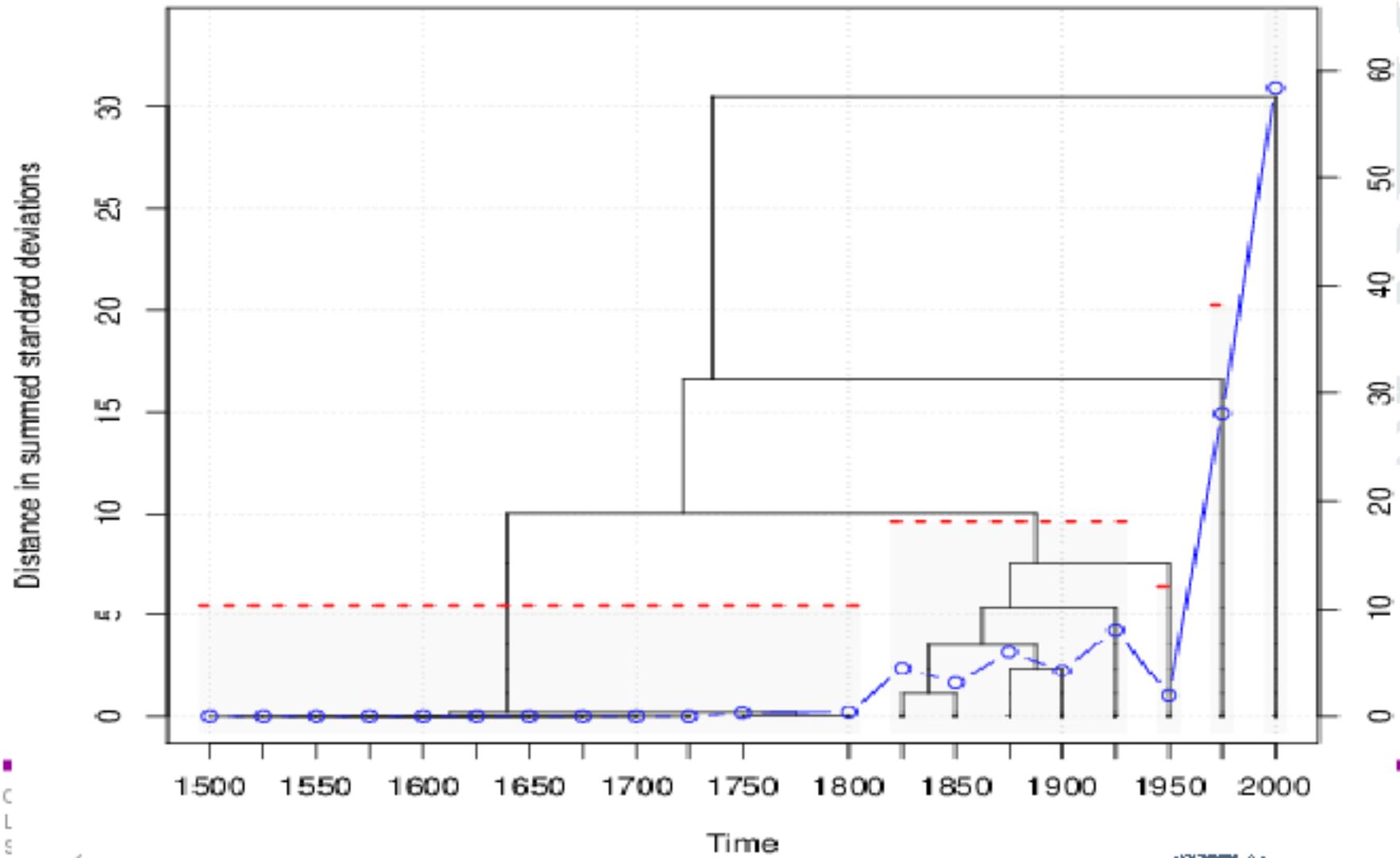
VNC dendrogram



Vers une linguistique diachronique outillée

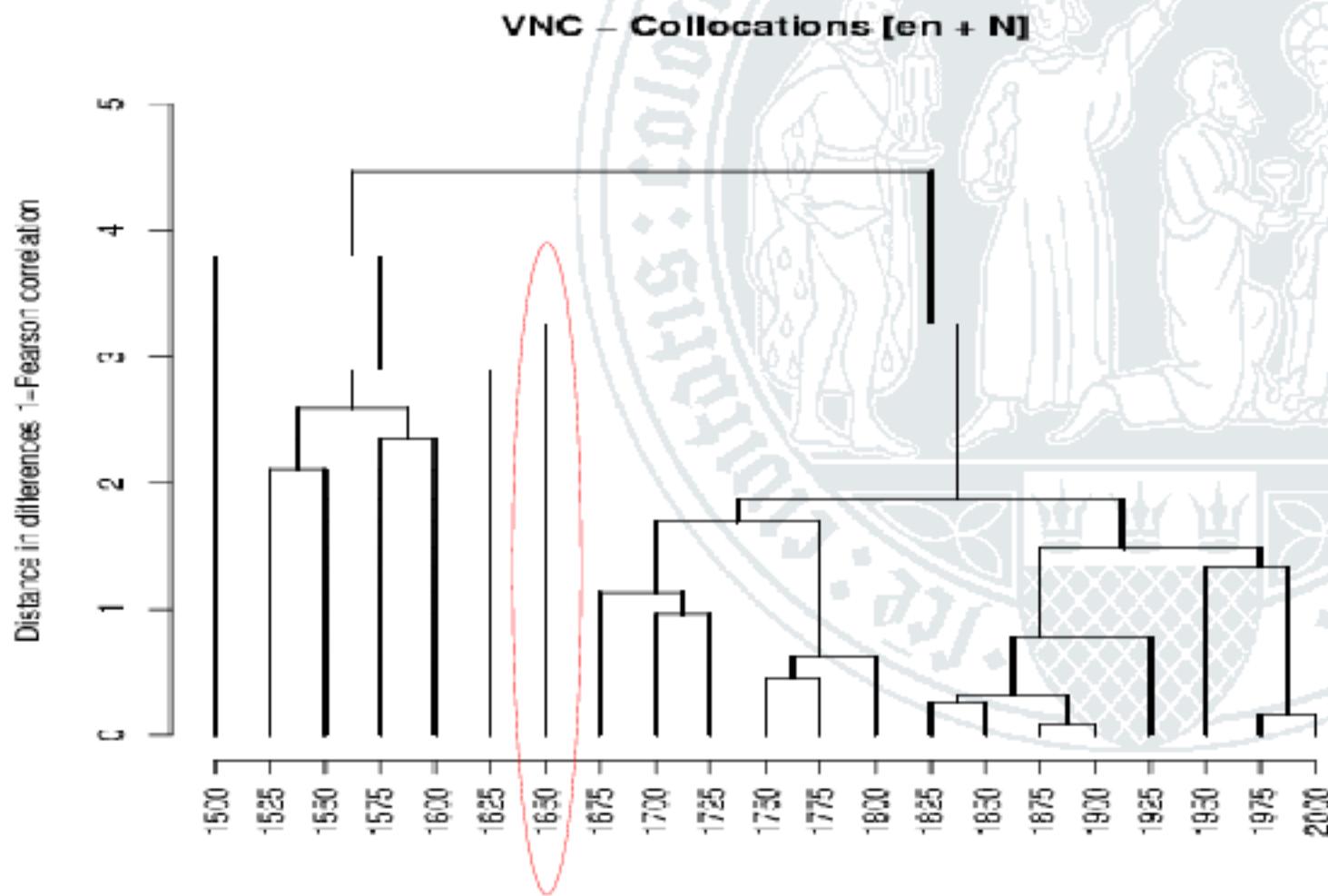
VNC appliqué sur une locution – *dans le cadre de*

VNC dendrogram



Vers une linguistique diachronique outillée

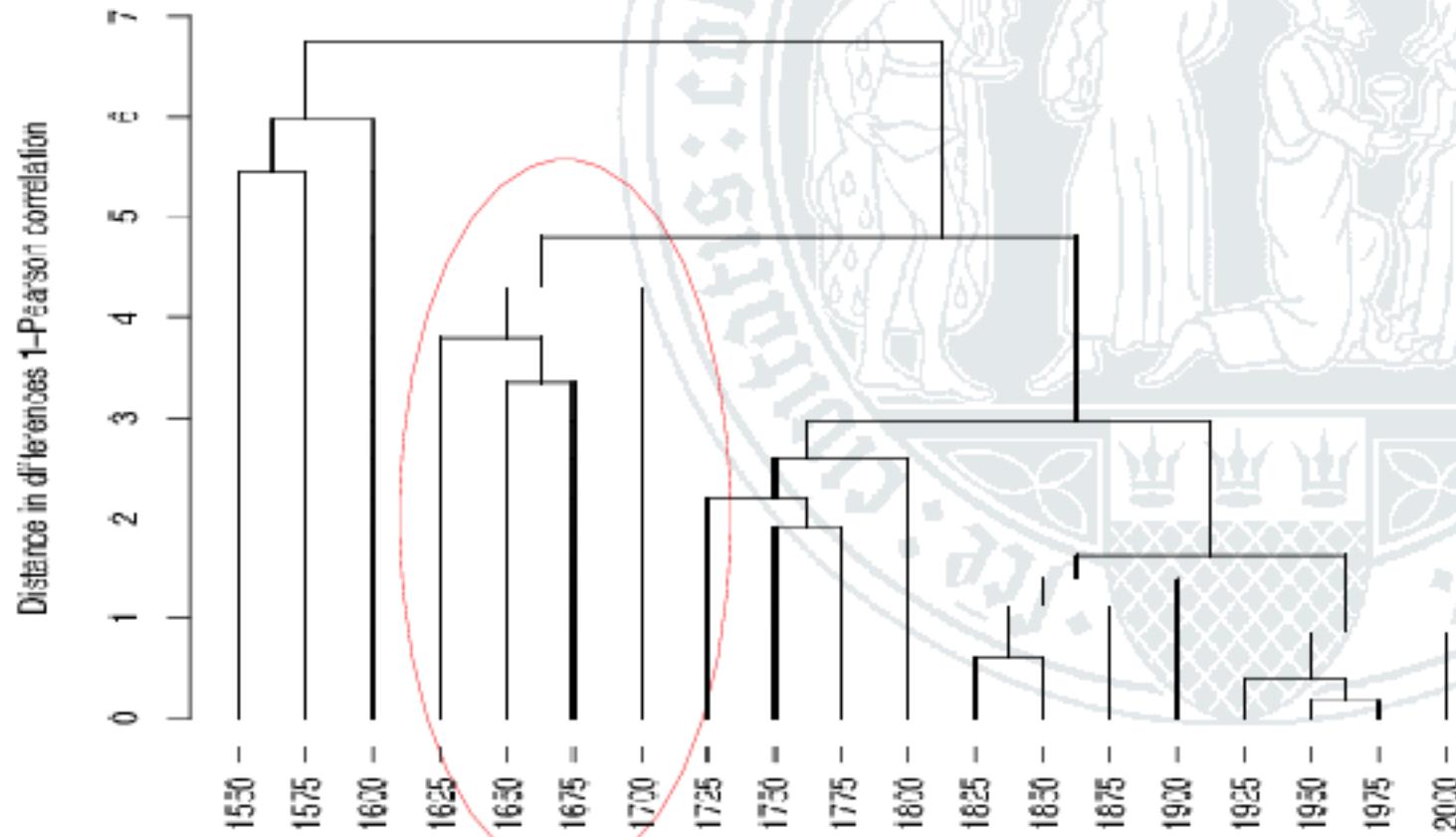
VNC appliqué sur un inventaire cooccurrentiel – *en + N*



Vers une linguistique diachronique outillée

VNC appliqué sur un inventaire cooccurrentiel – *dans + N*

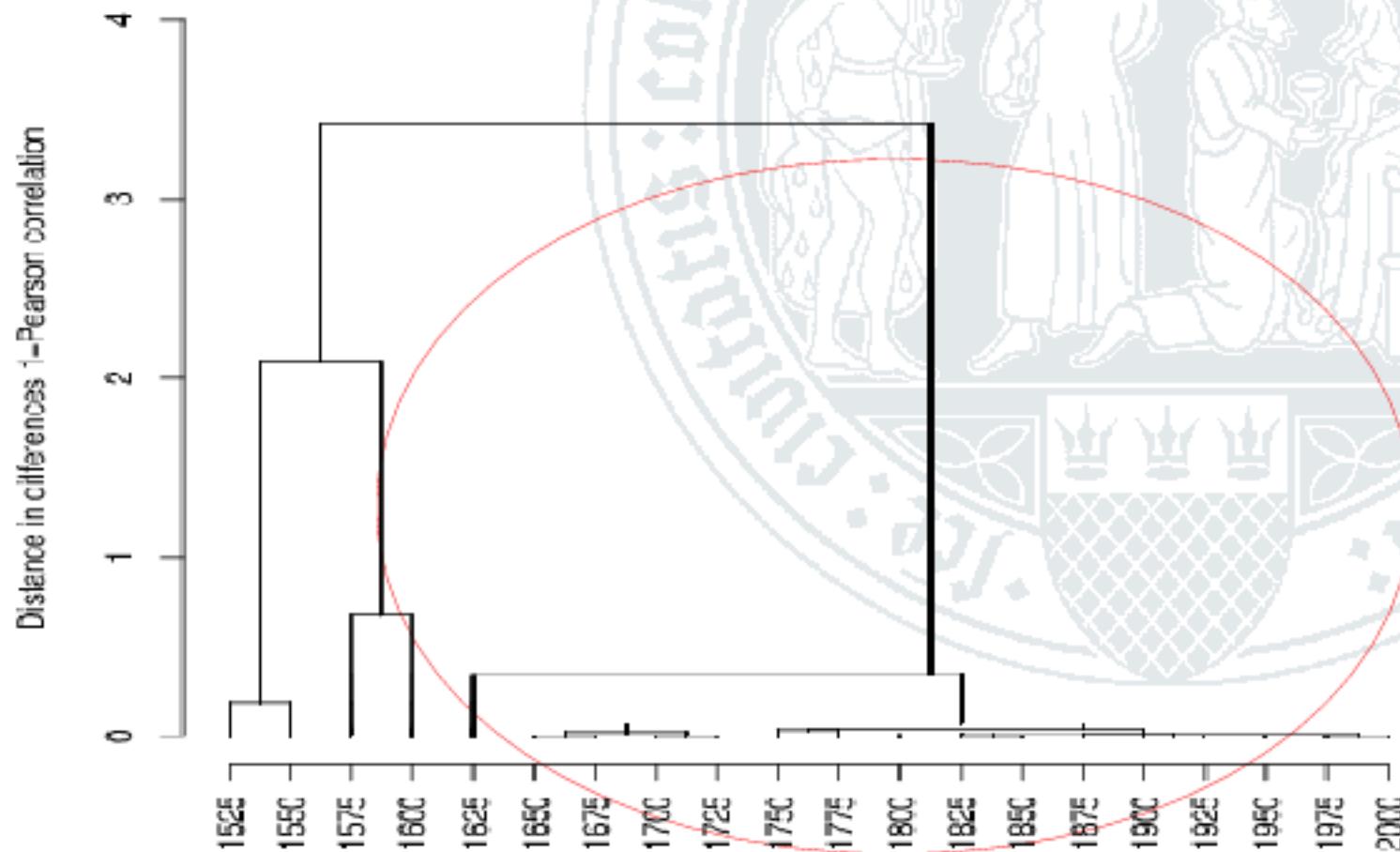
VNC – Collocations [*dans + N*]



Vers une linguistique diachronique outillée

VNC appliqué sur l'association préférentielle à des positions texuelles – *en effet* (en début, milieu, fin de phrase)

VNC – Colligations textuelles (phrase) *en_effet*



Vers une linguistique diachronique outillée

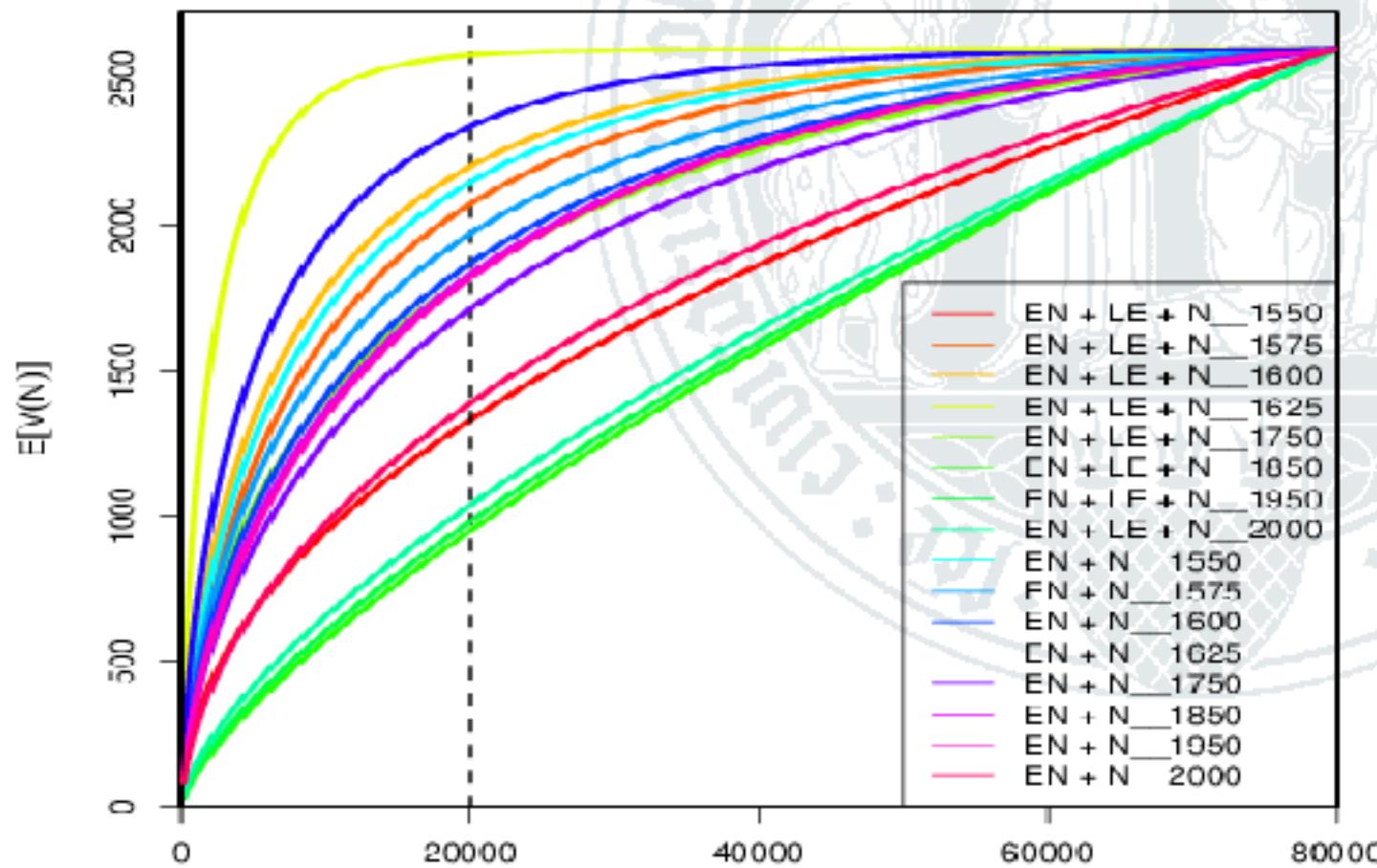
- Etablir, de façon inductive, la saturation d'une position paradigmique au sein d'un motif
- Analyses quantitatives de la productivité
 - morphologique (Baayen & Lieber 1991, Baayen 1992, Evert & Lüdeling 2005)
 - (morpho-)syntaxique (Zeldes 2012, Desaguiler à paraître)
- Méthodes de calcul et outils
 - modélisation LNRE (Baayen 2001, Evert 2004)
 - zipfR (Evert & Baroni 2006, 2007)



Vers une linguistique diachronique outillée

Courbes d'accroissement lexical (Vocabulary Growth Curves) – *en + N* vs. *en + Da + N*

VGCC – *en + N* vs. *en + Da + N*



Vers une linguistique diachronique outillée

HCFA – colligations textuelles de *en effet* (selon période, genre)

11	12	time	genre	Freq	Exp	Cont.chisq	Obs-exp	Dec
EN EFFET	WS_B	2000	encyclopédie	247	41,5909	1014,4713	>	***
EN EFFET	WS_B	1875	traité	518	326,7615	111,9185	>	***
EN EFFET	WS_B	1750	encyclopédie	195	34,2458	754,6009	>	***
EN EFFET	WS_B	1800	traité	377	239,009	79,6685	>	***
EN EFFET	WS_B	2000	presse	178	52,3511	301,5407	>	***
EN EFFET	WS_B	1675	traité	264	141,265	106,6356	>	***
EN EFFET	WS_B	1975	presse	157	46,4102	263,5219	>	***
EN EFFET	WS_B	1900	traité	264	158,7448	69,7891	>	***
EN EFFET	WS_B	1825	traité	330	245,7868	28,8537	>	***
EN EFFET	WS_B	1850	traité	292	204,0495	37,9089	>	***



Vers une linguistique diachronique outillée

HCFA – colligations textuelles de *en effet* (selon période, genre)

I1	I2	time	genre	Freq	Exp	Cont.chisq	Obs-exp	Dec
EN EFFET	WS_B	19/5	traité	195	2/5,7/521	23,6477	<	****
EN EFFET	WS_B	1750	narratif	10	78,5027	50,7765	<	***
EN EFFET	WS_B	10/5	narratif	57	100,151	18,592	<	**
EN EFFET	WS_B	1950	traité	77	125,5689	18,7/86	<	**
EN EFFET	WS_M	1975	traité	30	79,2768	30,6294	<	***
EN EFFET	WS_B	16/5	narratif	3	13,2967	37,5046	<	****
EN EFFET	WS_B	1825	narratif	34	75,3319	22,6773	<	***
EN EFFET	WS_B	1826	presse	5	11,3069	31,9/12	<	***
EN EFFET	WS_E	1975	traité	9	50,4256	34,0319	<	***
EN EFFET	WS_B	1900	narratif	14	48,6542	24,6826	<	***
EN EFFET	WS_M	1/50	narratif	1	22,569	20,6133	<	****



**1 Je
2 vous
3 remercie
4 de
5 votre
6 attention
7 !**

**Pp JE
Pp VOUS
Vvn REMERCIER
S DE
Ds VOTRE
Nc ATTENTION
Fs !**

