



Complexités d'un corpus en diachronie du français. Le cas de Presto.

Achille Falaise
Denis Vigier

Atelier CCC
Lyon, 4 juin 2015

Plan de la présentation

1. Le corpus
2. Chaîne de traitement
 1. Stratégie de traitement
 2. Création d'un lexique
 3. Création d'un corpus d'apprentissage
 4. Création d'un modèle de langage
 5. Analyse du corpus
3. Quelques problèmes épineux
 1. Jeu d'étiquettes : participes, gérondifs, adjectifs
 2. Les amalgames

Le corpus Presto

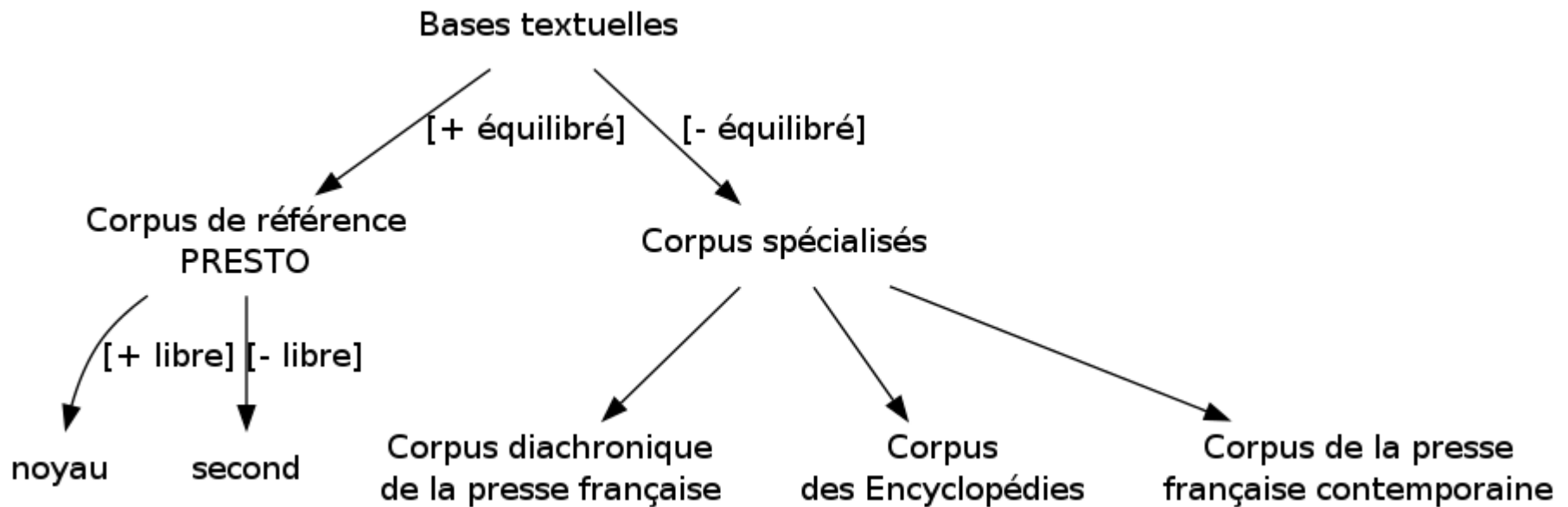
Le corpus Presto

- Constitution d'un corpus dans le cadre de PRESTO :
 - Prérequis pour la réalisation du projet
 - Un des apports majeurs du projet
- Objectifs linguistiques :
 - Proposer une couverture satisfaisante de toutes les périodes de l'histoire du français : 9^{ème} s. au 20^{ème} s.
 - Partir des bases existantes et construire un corpus « contrôlé »
 - Critère temporel : éviter les « trous »
 - Critère générique : variété contrôlée
 - « Qualité » des éditions + orthographe non modernisée
 - => *a minima* : chaque tranche décennale : 3 txt litt (GN, Th, P) + 3 traités
 - Enrichissement linguistique : étiquetage morpho-syntaxique et lemmatisation

Le corpus Presto

- **Objectifs** en termes d' « ouverture » :
 - Rendre disponible (licences CC)
 - Une partie du corpus constitué
 - La totalité des outils élaborés
- Collaborations avec les Bases textuelles : Frantext, BVH, ARTFL, CEPM, Cologne
 - Textes déjà sous licence CC : BVH, CEPM
 - Conventions passées
 - FRANTEXT : mise à disposition de la communauté sous licence CC d'une trentaine de textes jusqu'ici « fermés » ; mise à disposition de lexiques pour les outils PRESTO sous licence CC
 - ARTFL : mise à disposition d'un volume de l'Encyclopédie Diderot & D'Alembert

Le corpus Presto



Chaîne de traitement

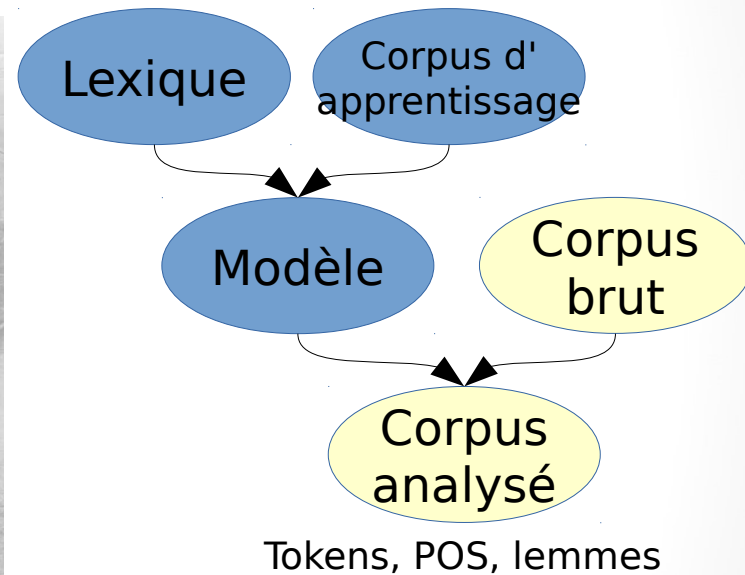
Chaîne de traitement

Stratégie de traitement

1. Stratégie de traitement

Une approche classique

- Approche « classique » : chaîne de traitement



- Contraintes :
 - « Massif »
 - Pas de *feedback*
 - Flux de *tokens*
=> pas d'ambiguïtés de segmentation

1. Stratégie de traitement

Des problèmes originaux

- Français classique
 - Variantes dialectales
 - Orthographe (y compris segmentation) variable, néologismes, etc.
 - Langue peu dotée
- Méthode
 - construction de ressources à partir du français moderne

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 triumpant 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 generale
Qui, l'esté qui vient, sera faicte.
En honneur du triumphe et feste,
Ne faillez monstrier vos bons cueurs
Qui font de la vertu approche,
Tant que l'on dye par honneurs :
Vive l'excellente Bazoche !

Sottie pour le cry de la bazoche, 1549

Chaîne de traitement

Création d'un lexique

2. Création d'un corpus de référence

Lexique (Gilles Souvay, ATILF)

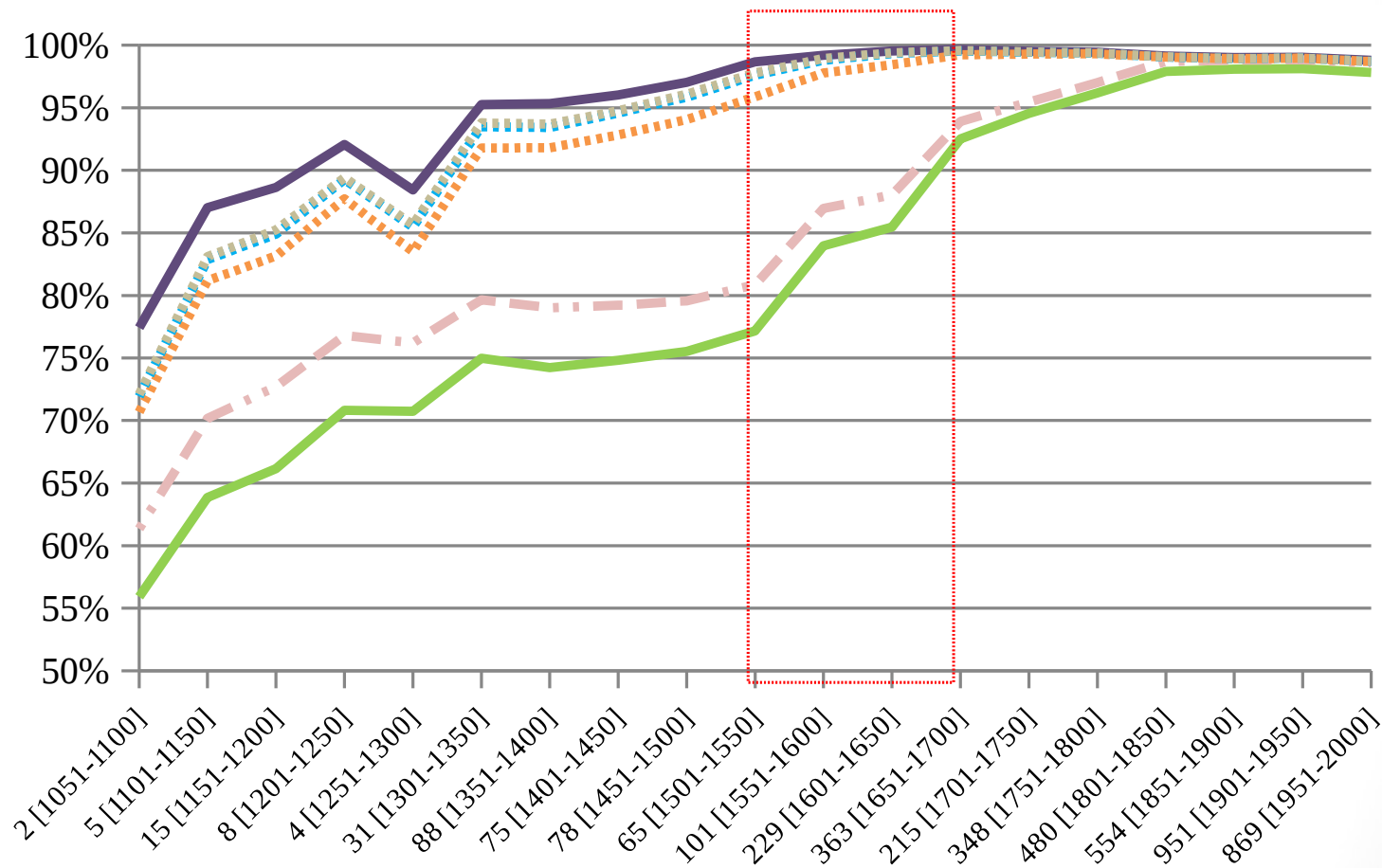
Forme ancienne => lemme moderne + POS

```
aimas;AIMER;VMIS2S0  
aimasmes;AIMER;VMIS1P0  
aimasse;AIMER;VMSI1S0  
aimassent;AIMER;VMSI3P0  
aimasses;AIMER;VMSI2S0  
aimassies;AIMER;VMSI2P0  
aimassiez;AIMER;VMSI2P0  
aimassions;AIMER;VMSI1P0  
aimassiés;AIMER;VMSI2P0  
aimast;AIMER;VMSI3S0  
aimastes;AIMER;VMIS2P0
```

2 735 843 entrées

2. Création d'un corpus de référence

Lexique (Gilles Souvay, ATILF)



Taux de couverture lexicale du lexique moderne (lexique de départ) et du lexique archaïsé (lexique final), avec étapes intermédiaires, mesurée sur le corpus Frantext (XI^{ème} - XX^{ème} siècle).

Chaîne de traitement

Création d'un corpus de référence

2. Création d'un corpus de référence

Choix des textes

- Sélection de 5 textes

- 5 périodes
- 5 genres

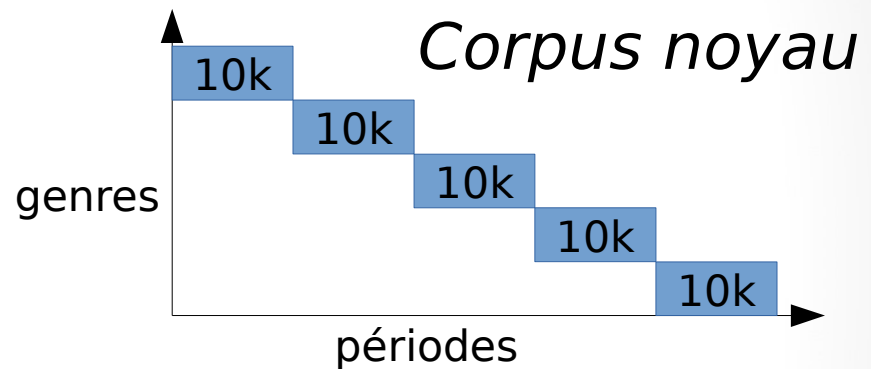
Saulsaye (1547)

Lisandre et Caliste (1631)

Les Lettres de messire Roger de Rabutin, comte de Bussy (1681)

Essay sur l'histoire generale et sur les moeurs et sur l'esprit des nations (1756)

Le Paysan perverti ou les Dangers de la ville (1776)



Total : 62k tokens

2. Création d'un corpus de référence

Objectifs et défis

- Objectifs
 - 50k mots annotés manuellement (token, POS, lemme)
 - Par au moins 2 annotateurs
- Défis
 - C'est très chronophage !
- Méthode
 - Pré-annotation automatique du corpus

2. Création d'un corpus de référence

Objectifs et défis

- Objectifs
 - 50k mots annotés manuellement (token, POS, lemme)
 - Par au moins 2 annotateurs
- Défis
 - C'est très chronophage !
- Méthode
 - Pré-annotation automatique du corpus

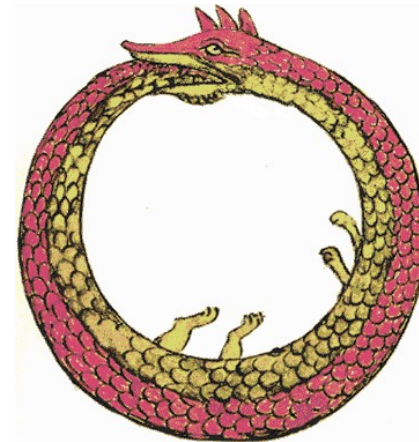
*Pour créer ce corpus annoté,
il faut donc...*

2. Création d'un corpus de référence

Objectifs et défis

- Objectifs
 - 50k mots annotés manuellement (token, POS, lemme)
 - Par au moins 2 annotateurs
- Défis
 - C'est très chronophage !
- Méthode
 - Pré-annotation automatique du corpus

*Pour créer ce corpus annoté,
il faut donc... un corpus
annoté !*



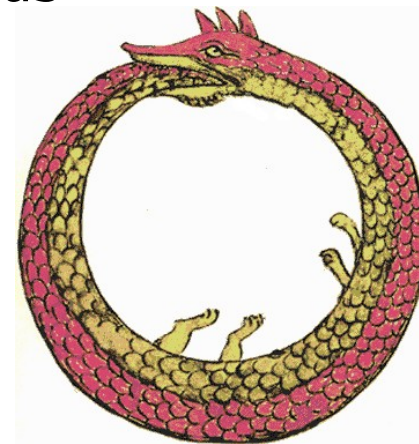
2. Création d'un corpus de référence

Objectifs et défis

- Objectifs
 - 50k mots annotés manuellement (token, POS, lemme)
 - Par au moins 2 annotateurs
- Défis
 - C'est très chronophage !
- Méthode
 - Pré-annotation automatique du corpus

*Pour créer ce corpus annoté,
il faut donc... un corpus
annoté !*

Mais il peut être incomplet.



2. Création d'un corpus de référence

Normalisation des textes

C'est assez dict pour ceste foys.<lb/>
Quand sçavoir en vous s'assocye,<lb/>
Monsieur *Rien, l'on vous remercye<lb/>
Du bien qu'avons aprins de vous.<lb/>
Bazochiens, entendez tous :<lb/>
Je veulx en triumpant arroy<lb/>
Eslire et faire ung nouveau roy,<lb/>
Comme il est coustume de faire ;<lb/>
Pourtant chacun pense a l'affaire,<pb n="267"/>
Autant les grandz que les petitz,<lb/>
Et faire les preparatifz ;<lb/>
Car, ainsi comme liberalle,<lb/>
Je tendz a monstre generale<lb/>
Qui, l'esté qui vient, sera faicte.<lb/>
En honneur du triumphe et feste,<lb/>
Ne faillez monstrez vos bons cueurs<lb/>
Qui font de la vertu approche,<lb/>
Tant que l'on dye par honneurs :<lb/>
Vive l'excellente *Bazoche !</p>

2. Création d'un corpus de référence

Normalisation des textes

C'est assez dict pour ceste foys. Quand sçavoir en vous
s'assocye, Monsieur *Rien, l'on vous remercy 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 generale Qui, l'esté qui vient, sera
faicte. En honneur du triumphe et feste, Ne faillez monstrez vos
bons cueurs Qui font de la vertu approche, Tant que l'on dye par
honneurs : Vive l'excellente *Bazoche !

2. Création d'un corpus de référence

Normalisation des textes

```
<lb/>Mais par telle legierete ne convient esti
<lb rend="hyphen"/>mer les oeuvres des humains. Car
  <lb/>vous mesmes dictes, que
<choice><orig>lhabit</orig><reg>l&#x2019;habit</reg></choice> ne faict
  <lb/>point le moine: & amp; tel est vestu
<choice><orig>dhabit</orig><reg>d&#x2019;habit</reg></choice>
  <lb/>monachal, qui au dedans
<choice><orig>nest</orig><reg>n&#x2019;est</reg></choice> rien
moins
  <lb/>que moyne: & amp; tel est vestu de cappe
  <fw place="bot-center" type="sig">A iij</fw>
```

2. Création d'un corpus de référence

Normalisation des textes

Mais par telle legierete ne convient estimer les oeuvres des humains. Car vous mesmes dictes, que l'habit ne faict point le moine: & tel est vestu d'habit monachal, qui au dedans n'est rien moins que moyne: & tel est vestu de cappe (...)

2. Création d'un corpus de référence

Projection lexicale

quelques	QUELQUE:AQ0CP0 QUELQUE:DI0CP0
remarques	REMARQUE:NCFP000 REMARQUER:VMIP2S0 REMARQUER:VMP00PM REMARQUER:VMSP2S0
sur	SUR:AQ0CS0 SUR:SPS00 SÛR:AQ0MS0 SÛR:RG
les	LE:DA0CP0 LES:PP3CPA00 LÈS:SPS00 LÉ:NCMP000
groupements	GROUPEMENT:NCMP000

↑

Texte
normalisé

↑

Lexique projeté

2. Création d'un corpus de référence

Projection lexicale

quelques	QUELQUE:AQ0CP0 QUELQUE:DI0CP0
remarques	REMARQUE:NCFP000 REMARQUER:VMIP2S0 REMARQUER:VMP00PM REMARQUER:VMSP2S0
sur	SUR:AQ0CS0 SUR:SPS00 SÛR:AQ0MS0 SÛR:RG
les	LE:DA0CP0 LES:PP3CPA00 LÈS:SPS00 LÉ:NCMP000
groupements	GROUPEMENT:NCMP000

Que d'ambiguïtés !

- simplification du jeu d'étiquettes
- désambiguïstation à l'aide d'un modèle moderne

2. Création d'un corpus de référence

Simplification des étiquettes

quelques	QUELQUE: AQ0CP0 QUELQUE: DI0CP0
remarques	REMARQUE: NCFP000 REMARQUER: VMIP2S0 REMARQUER: VMP00PM REMARQUER: VMSP2S0
sur	SUR: AQ0CS0 SUR: SPS00 SÛR: AQ0MS0 SÛR: RG
les	LE: DA0CP0 LES: PP3CPA00 LÈS: SPS00 LÉ: NCMP000
groupements	GROUPEMENT: NCMP000

Presto_0

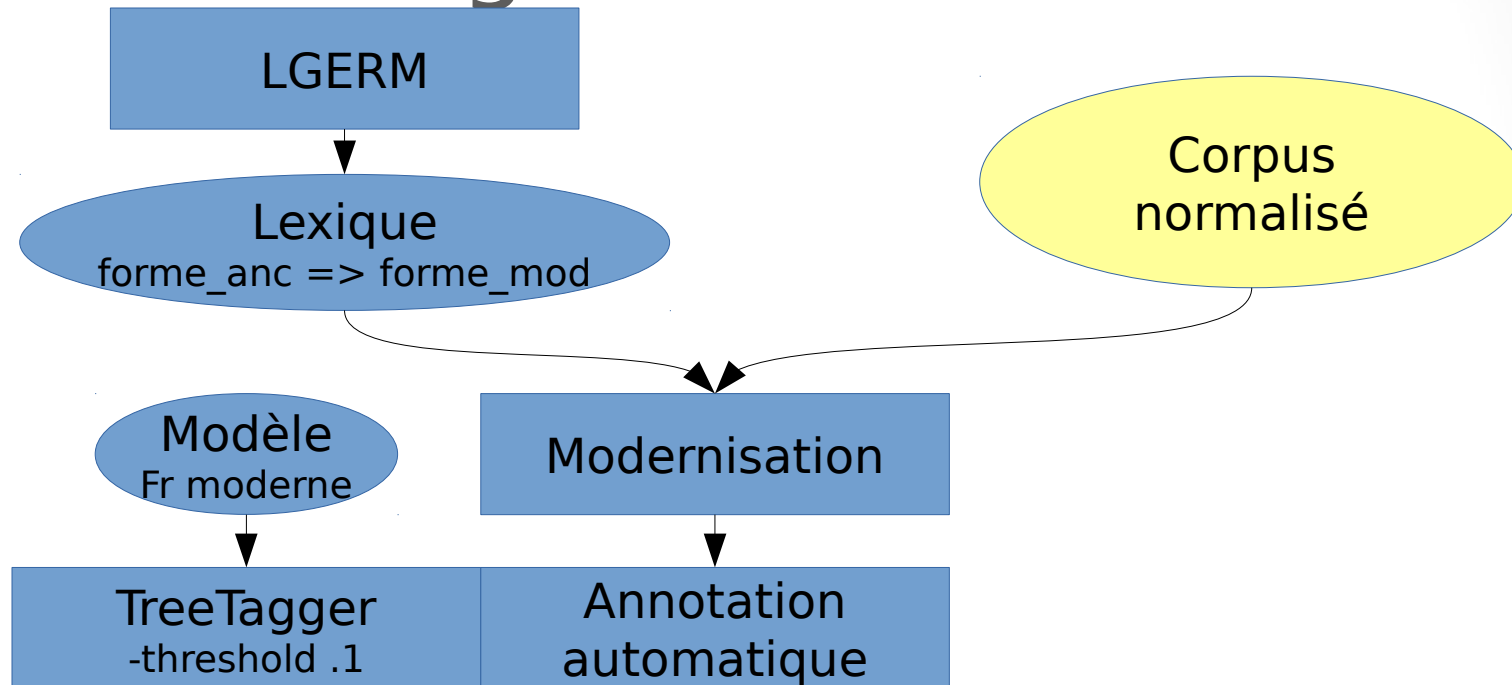


quelques	QUELQUE: Ag QUELQUE: Di
remarques	REMARQUE: Nc REMARQUER: Vvc
sur	SUR: Ag SUR: S SÛR: Ag SÛR: R
les	LE: Da LES: Pp LÈS: S LÉ: Nc
groupements	GROUPEMENT: Nc

Presto_2

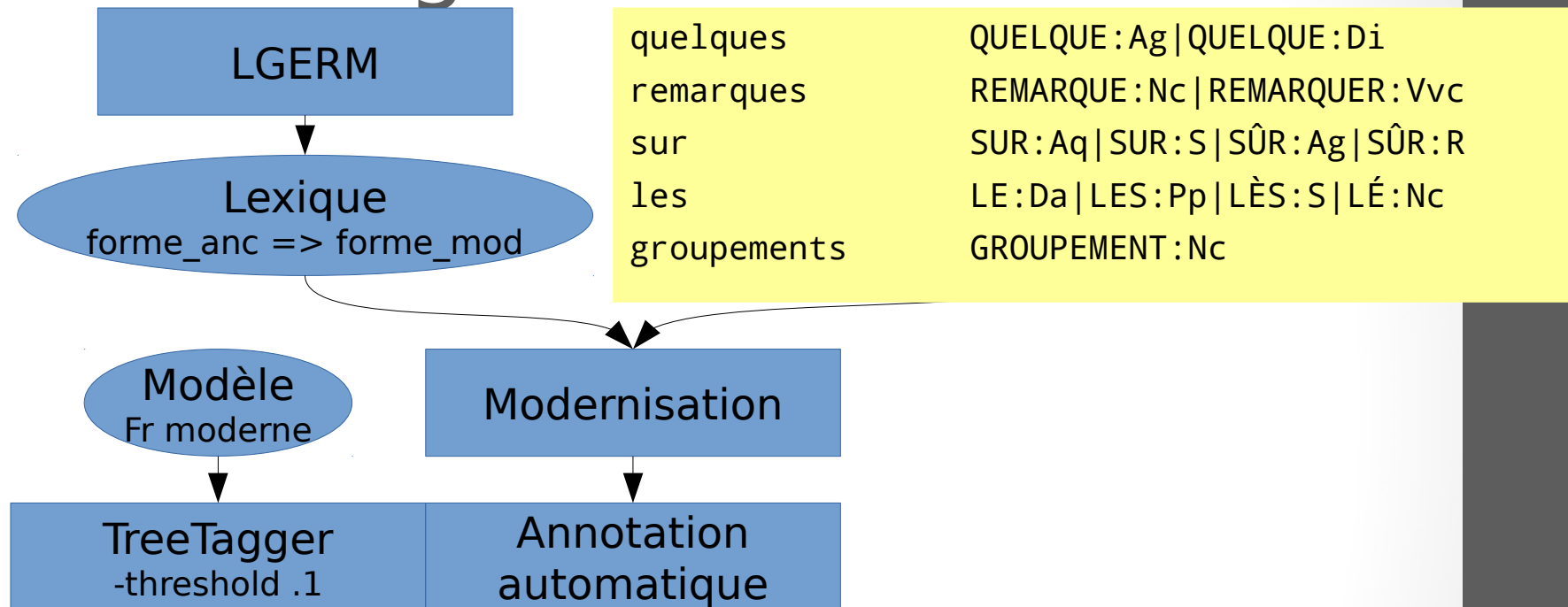
2. Création d'un corpus de référence

Désambiguïsation



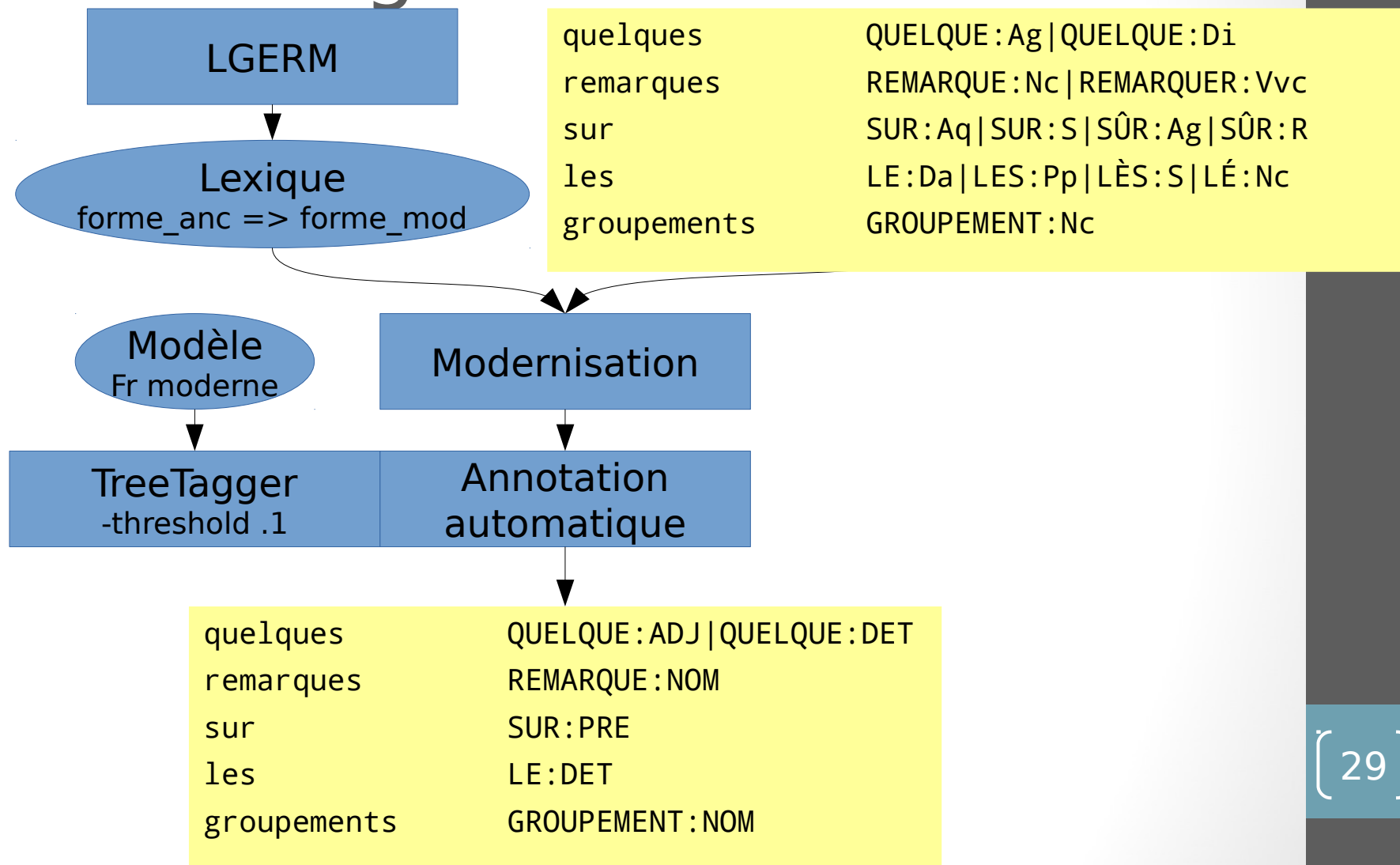
2. Création d'un corpus de référence

Désambiguïsation



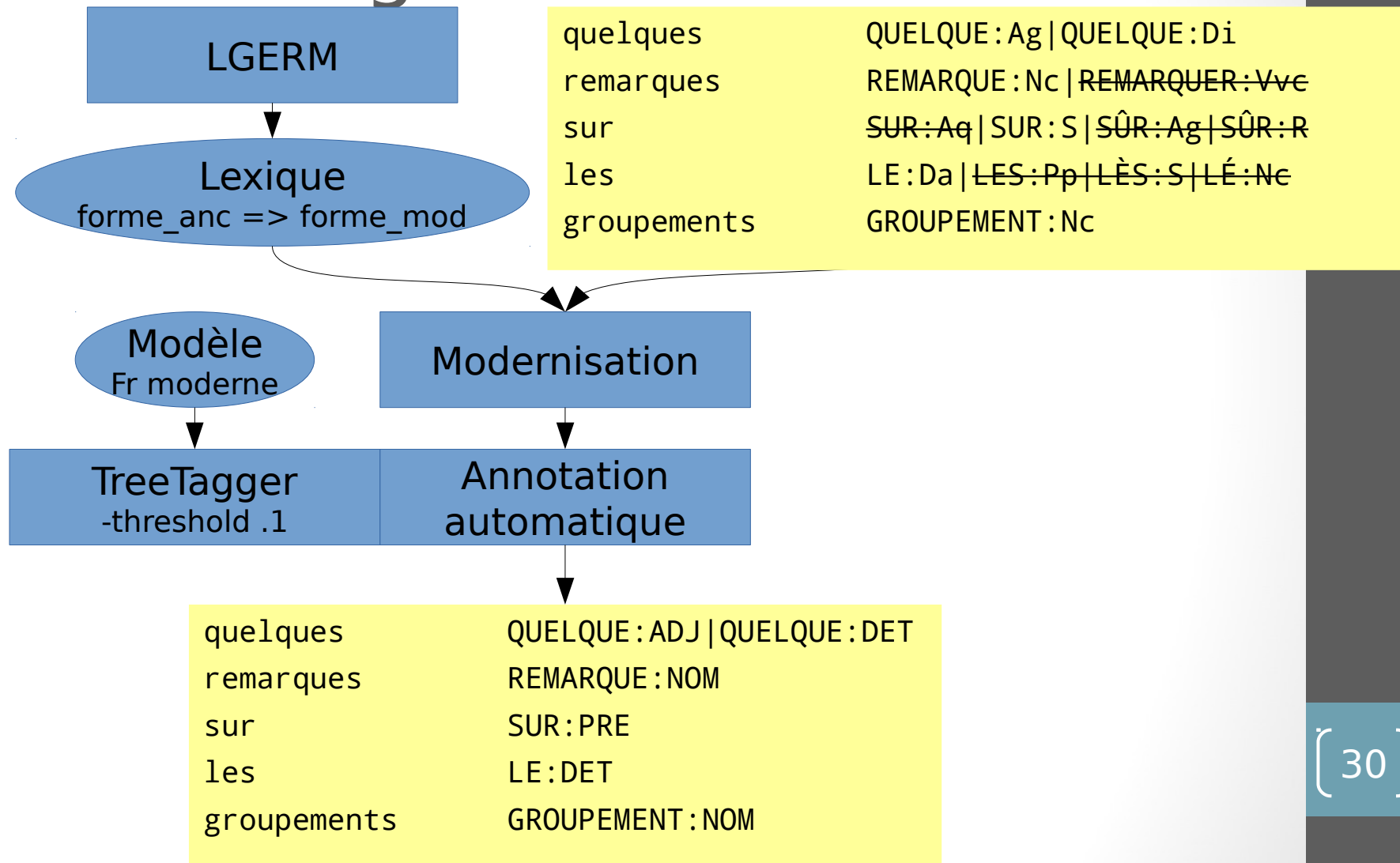
2. Création d'un corpus de référence

Désambiguïsation



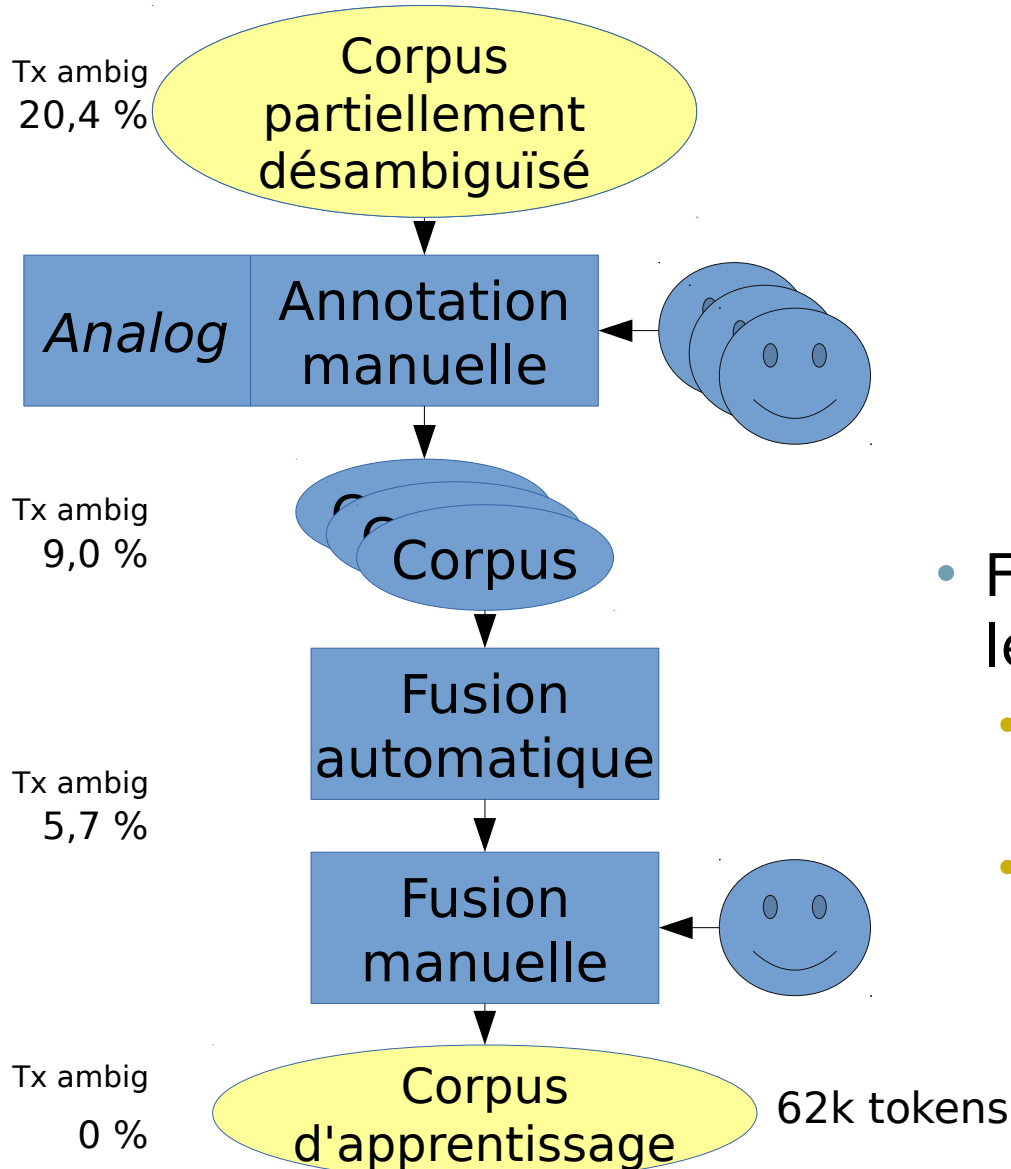
2. Création d'un corpus de référence

Désambiguïsation



2. Création d'un corpus de référence

Annotation manuelle et fusion



- Fusion automatique pour les cas « évidents » :
 - Au moins 2 annotateurs d'accord
 - Diacritiques

2. Création d'un corpus de référence

Analog (MH Lay, Poitiers)

Texte Annoté - Pantagruel 1542-UNIC

Choix pour l'affichage Exporter Tri Alphabétique Filtrer Srce Cpt CptG CptG %

CT Mode Validation Validation Auto InVal Concordance Conc.* Re-Annoter ReA-Dico Exporter FF Validées Stat

Mot n°	Forme rencontrée	Variante de ...	Lemme Vali...	CG Validée	Constellation	Mode Valid...	V	NCM	JQua	NPro	NCF	VAux	NC	Autre	Inconnu
117	maintesfoys														INC
118	passee						passer(passer)	passé(passe)			passee(passe)				
119	vostre														INC
120	temps	temps	temps	NCM		VA/DS		temps(temps)							
121	avecques														INC
122	les	le	le	Autre		VA/DS								le(le)	
123	honorables	honorable	honorable	JQua		VA/DS			honora...						
124	Dames						damer(damer)				dame(dame)				
125	et	et	et	Autre		VA/DS								et(et)	
126	Damoyselles														INC
127	,	,	,	Autre		VA/DS								,(,)	
128	leur													leur(leur)/lu...	
129	en	en	en	Autre		VA/DS								en(en)	
130	faisans	faisan	faisan	NC		VA/DS							faisa...		
131	beaux														INC
132	et	et	et	Autre		VA/DS								et(et)	
133	longs							long(long)	long(lo...						
134	narrez	narrer	narrer	V		VA/DS	narrer(narrer)								
135	,	,	,	Autre		VA/DS								,(,)	
136	alors	alors	alors	Autre		VA/DS								alors(alors)	
137	que	que	que	Autre		VA/DS								que(que)	
138	estiez														INC
139	hors	hors	hors	Autre		VA/DS								hors(hors)	
140	de	de	de	Autre		VA/DS								de(de)	
141	propos	propos	propos	NCM		VA/DS		propos(propos)						longs narrez , alors que estiez - hors - de propos : dont estez bien	
142	:	:	:	Autre		VA/DS								:(,)	
143	dont	dont	dont	Autre		VA/DS								dont(dont)	
144	estez														
145	bien							bien(bien)	bien(bi...					bien(bien)	INC
146	dignes	digne	digne	JQua		VA/DS			digne(...						
147	de	de	de	Autre		VA/DS								de(de)	
148	grande								grand(...				gran...		
149	louange						louanger(louanger)				louange(loua...				
150															

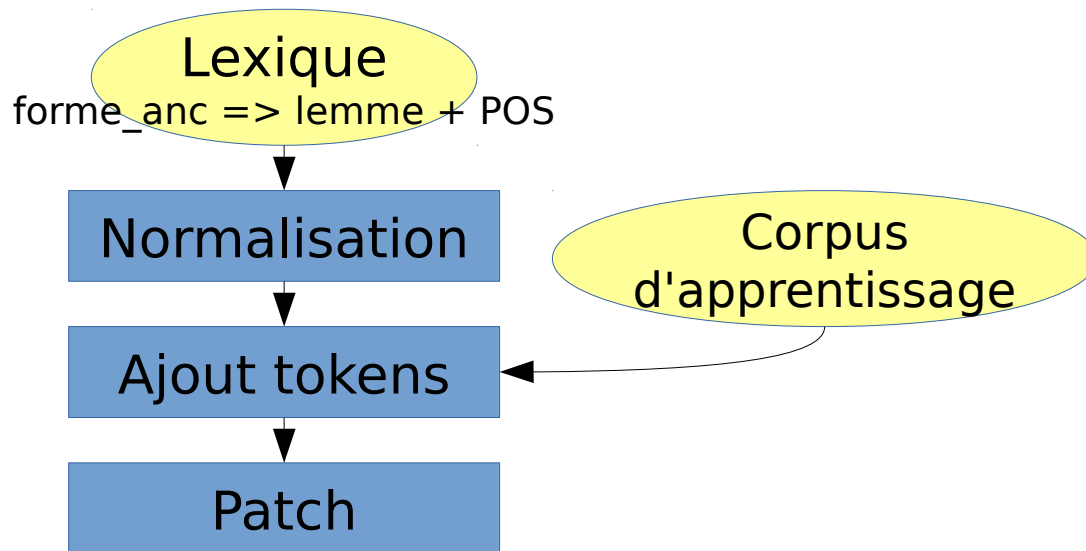
TreeTager=>ANALOG +CG

Chaîne de traitement

Création d'un modèle de langue pour le français classique

3. Création d'un modèle de langue pour le français classique

Préparation du lexique



- Patch :
 - Listes de tokens
 - À ajouter
 - À enlever
 - Règles ad hoc

3. Création d'un modèle de langue pour le français classique

Création du modèle

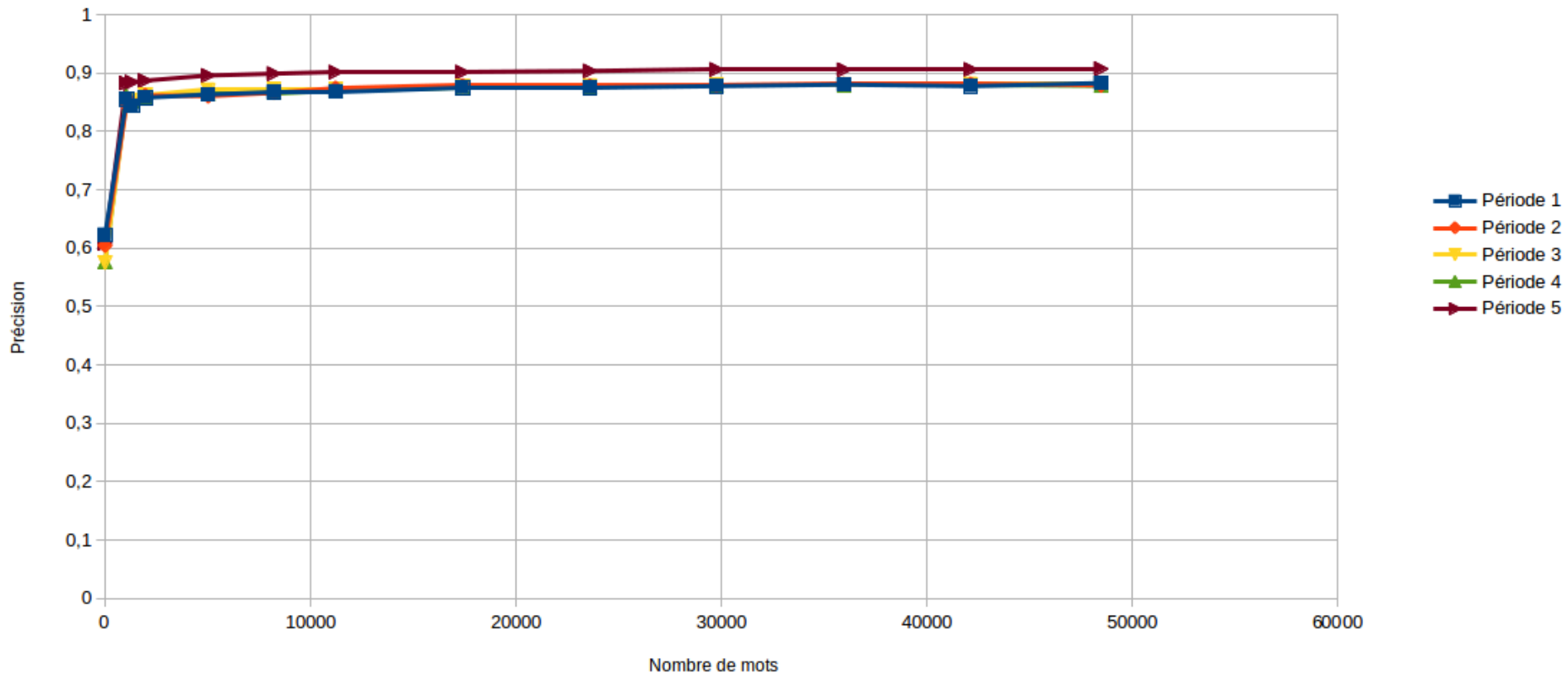
- Division du corpus d'apprentissage en trois
 - Corpus d'entraînement (80% – 49 630 tokens)
 - Corpus de développement (10% – 6 164 tokens)
 - Corpus de référence (10% – 6 110 tokens)
- *Autotuning* pour trouver les meilleurs paramètres pour TreeTagger
 - $cl\ 2 ; dtg\ 0,5 ; sw\ 1 ; ecw\ 0,06 ; atg\ 1,15$
 - Précision +0,05 %
- Précision
 - Corpus d'entraînement : 95,77 %
 - Corpus de développement : 94,28 %
 - Corpus de référence : 94,46 %

3. Création d'un modèle de langue pour le français classique

Évaluation du modèle

Précision du modèle TreeTagger générique pour les POS

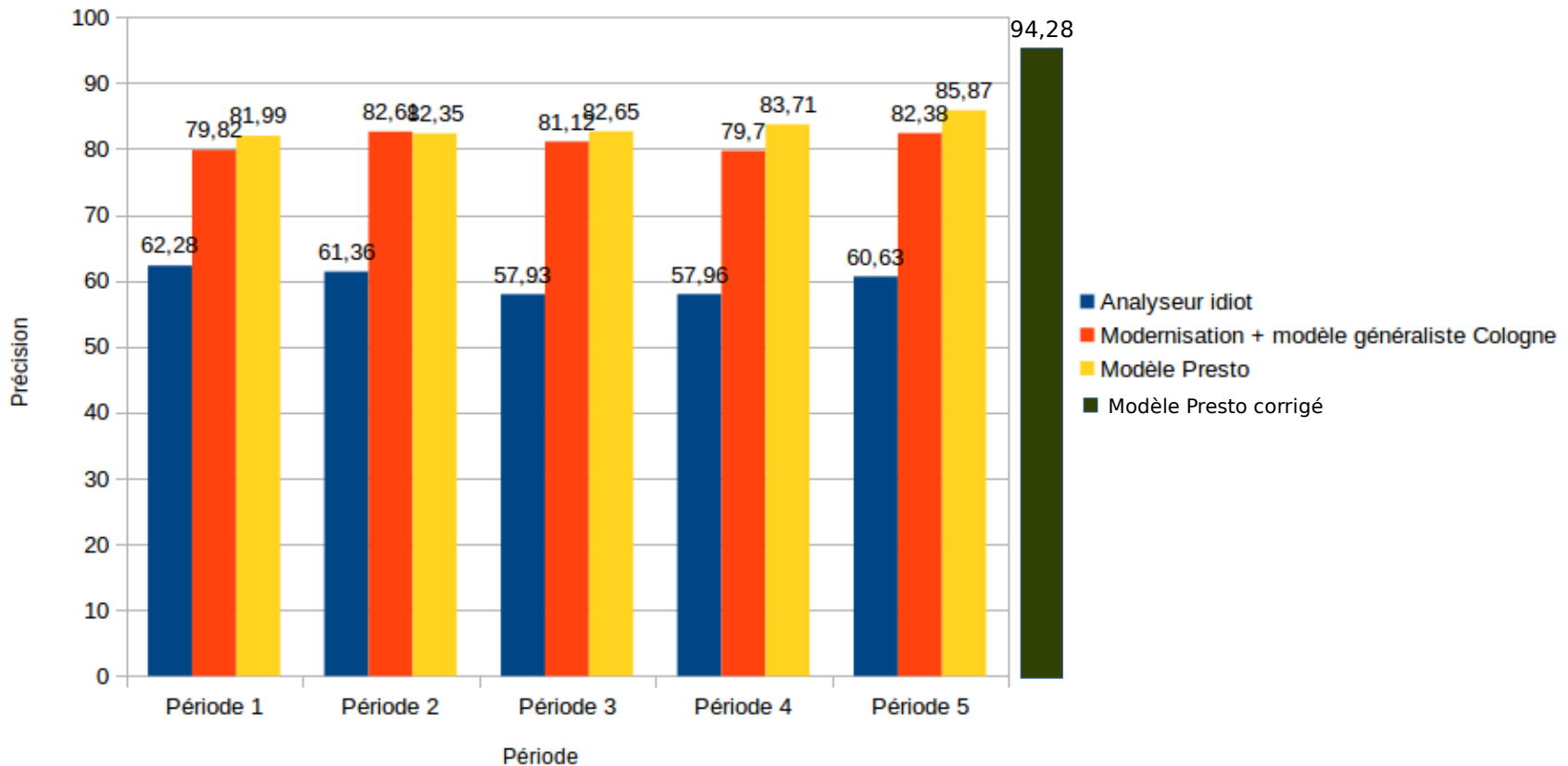
Le corpus d'apprentissage comporte toutes les périodes, on fait varier le nombre de mots.
Le baseline «0 mots» est obtenu, sans modèle, par tirage aléatoire des catégories à partir du lexique d'apprentissage.
Le corpus d'évaluation est différent pour chaque période, et comporte 761 à 1946 mots selon la période.



3. Création d'un modèle de langue pour le français classique

Évaluation du modèle

Précision selon la stratégie d'annotation automatique
(étiquettes uniquement, hors tokenisation)

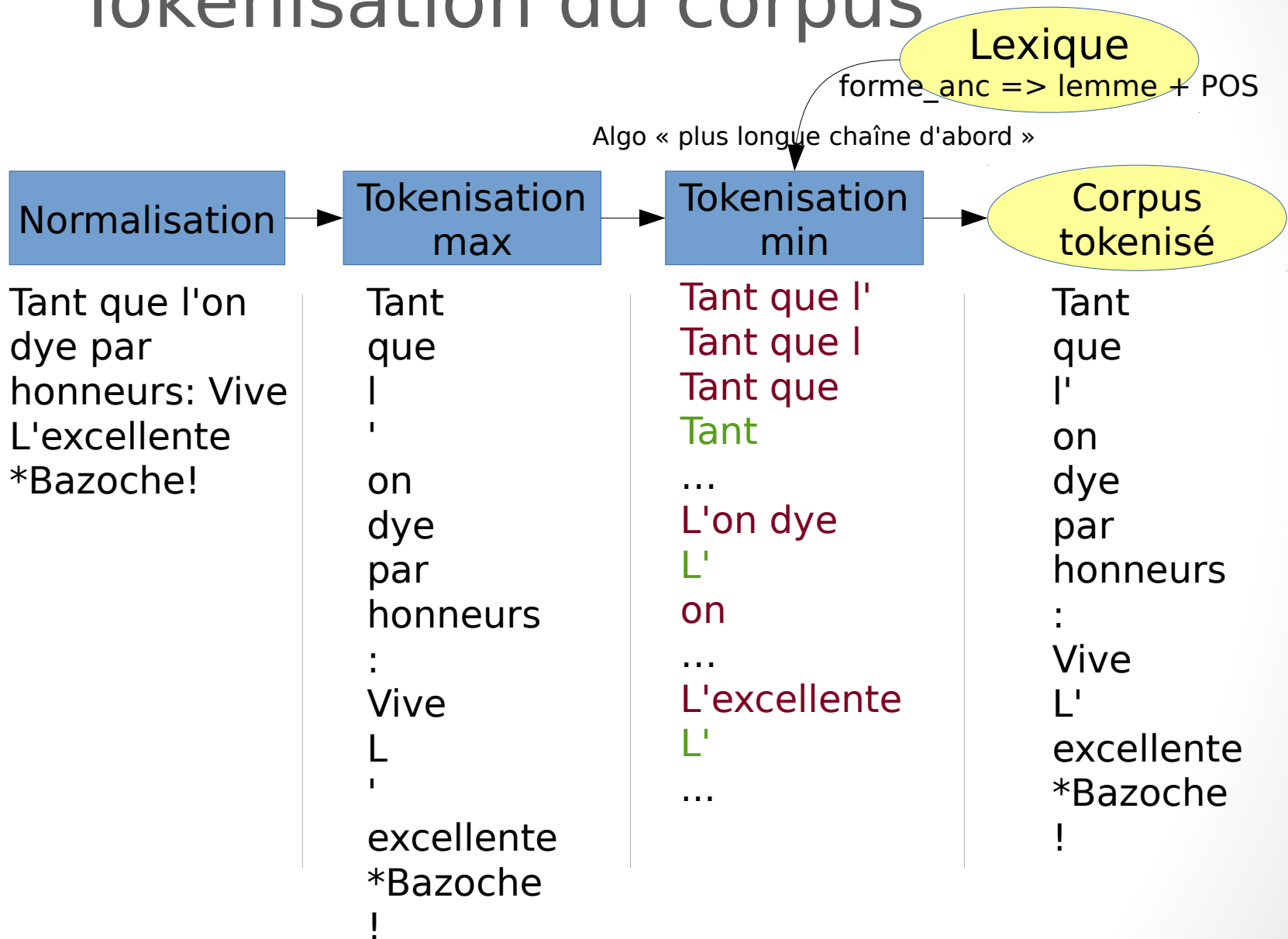


Chaîne de traitement

Analyse du corpus

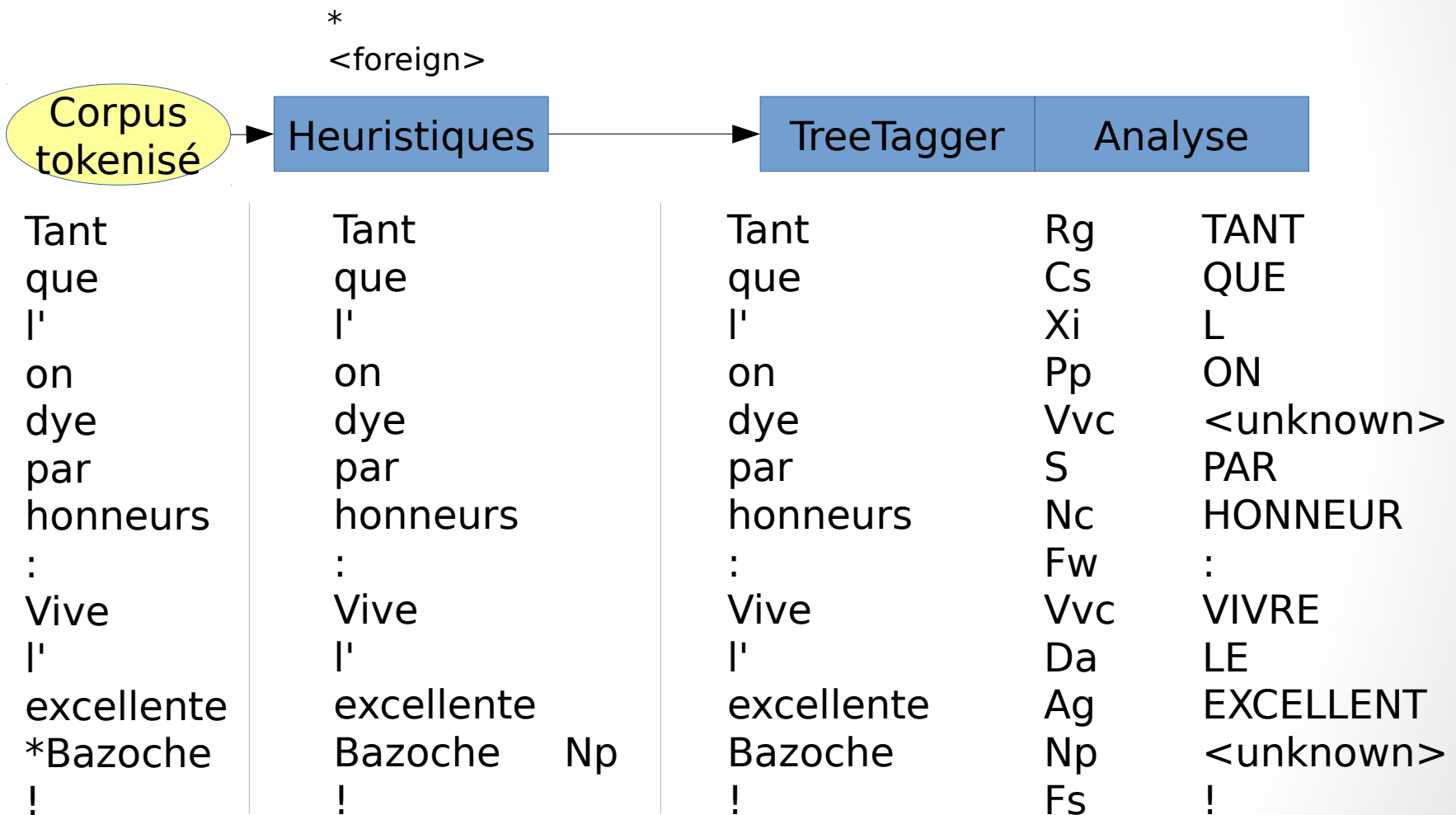
4. Analyse du corpus

Tokenisation du corpus



4. Analyse du corpus

Traitement du corpus



4. Analyse du corpus

Implémentation

- Programmation
 - Format de données tabulaire
 - Commandes Unix (« *Unix for Poets* »), Bash, Perl
 - Chaîne de montage : *Pipes* Unix
 - ~multitâche
 - Souple, facile à déboguer
- **Vitesse** (création du modèle, mäj lexique, tokenisation, analyse)
 - Pour un lexique de 2,7 M tokens
 - Pour un corpus de 28,3 M tokens
 - Processeur Xeon 4*3,2 Ghz
 - => 10 minutes

4. Analyse du corpus

Perspectives

- Pêche aux erreurs
 - Vérification des ressources



Quelques problèmes épineux

Jeu d'étiquettes : La catégorie G (participe, adjectifs, gérondifs)

Ga : Participes présents, adjectifs verbaux, gérondifs

Ge : participes passés, adjectifs

- **En synchronie** : cas pb -> tests (Riegel & al. 2009 : 737-738) -> chances de divergences entre annotateurs. En outre :
 - *Selon les cas (le type de verbe, le contexte), ils [les participes] sont sentis comme plus ou moins «verbaux» ou «adjectivaux» (avec une marge appréciable de liberté d'interprétation) (P. le Goffic 1993, § 134: 201)*
- **Dimension diachronique** :
 - *En français classique, la tripartition des formes en -ant ne va pas de soi (...) le gérondif, invariable, se distingue mal du participe (au masculin singulier) du fait qu'il n'est pas régulièrement précédé de en; le participe qui peut être variable en genre et en nombre, se distingue mal de l'adjectif verbal. (N. Fournier, 2002, § 421: 291-292)*

Quelques problèmes épineux

Les amalgames

- **en+le > enl, el, on, ou, au, hou, hu**
 - spécialement que pour toy estre liee **ou** lien de mariage (PRESTO_1535, La déplourable fin de Flamete, J. de Flores)
 - ainsi que escript cesar en ses commentaires, & jean de gravot **on** mythologies gallicques. (PRESTO_1542_Pantagruel_Rabelais)
- **en+les > ens, ans, eins, ons, es, ès, eus, aus, aux**
 - mais par ce moyen de propagation seminale demoure **es** enfans ce que estoit de perdu **es** parens (PRESTO_1542_Pantagruel_Rabelais)
 - la quele il faudroit chercher en ces vieux grecz, & latins, non point **és** aucteurs francoys. (PRESTO_1549_Deffence_J. du Bellay)

Importance pour les calculs de fréquence, de cooccurrence lexicale ...