



Annotation du corpus PRESTO: création de ressources pour l'analyse du français classique

Sascha Diwersy, Universität zu Köln
Achille Falaise, ENS de Lyon, ICAR

Journées d'étude *Presto*
Cologne, 12-13 février 2015

Plan de la présentation

1. Le corpus
2. Chaîne de traitement
 1. Stratégie de traitement
 2. Création d'un corpus d'apprentissage
 3. Création d'un modèle de langage
 4. Analyse du corpus
3. Exploitation du corpus
 1. Principes d'échantillonnage
 2. Modèles de données et méthodes
 3. Outils (démonstration PrimeStat)
 4. Perspectives – vers une linguistique diachronique outillée

Ressources

Le corpus Presto

- Constitution d'un corpus dans le cadre de PRESTO :
 - Prérequis pour la réalisation du projet
 - Apport majeur du projet
- **Besoins** :
 - Représentation de toutes les périodes de l'histoire du français : 9^{ème} s. au 20^{ème} s.
 - Présence de différents types de textes et de différents genres discursifs
 - Enrichissement linguistique : étiquetage morpho-syntaxique et lemmatisation

Le corpus Presto

- **Objectifs :**

- Disposer d'annotations linguistiques fiables mais également de gros volumes de textes
- Pouvoir rendre disponibles autant que possible le corpus constitué et les outils élaborés tout en respectant les contraintes juridiques

- **Contrainte :**

- Disposer de versions numérisées de textes de bonne qualité

→ Collaboration avec diverses bases textuelles existantes :

- Frantext
- BVH
- Cologne
- ARTFL
- CEPML
- ...

Le corpus Presto

- **Objectifs :**

- Disposer d'annotations linguistiques fiables mais également de gros volumes de textes
- Pouvoir rendre disponibles autant que possible le corpus constitué et les outils élaborés tout en respectant les contraintes juridiques

- **Contrainte :**

- Disposer de versions numérisées de textes de bonne qualité

→ Collaboration avec diverses bases textuelles existantes :

- Frantext
- BVH
- Cologne
- ARTFL
- CEPM
- ...

392 textes :

La Concorde du genre humain (1509)... La carte et le territoire (2010)

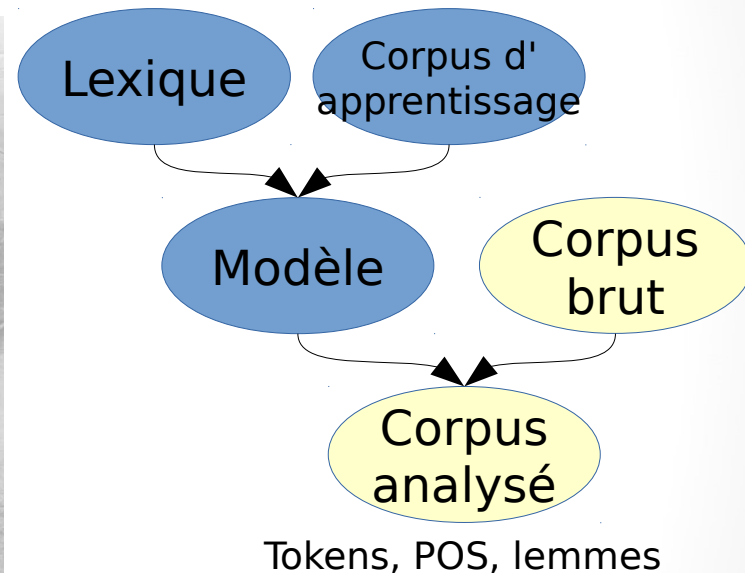
Gargantua, L'Astrée, l'Encyclopédie (tome 7)...

Chaîne 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

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

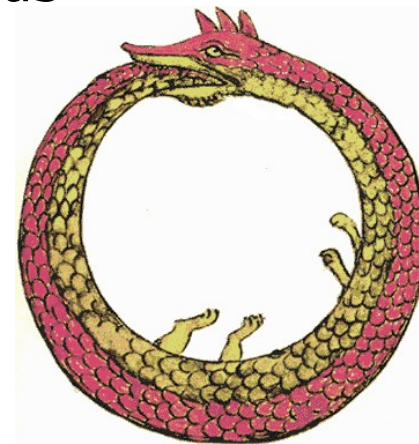
*Pour créer ce corpus annoté,
il faut donc...*

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 !*



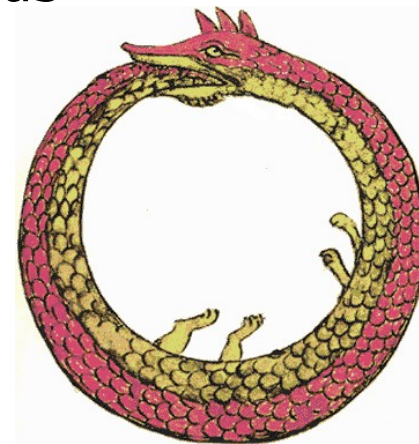
2. Création d'un corpus de référence

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

Choix des textes

- Sélection de 5 textes
 - 5 périodes
 - 5 genres

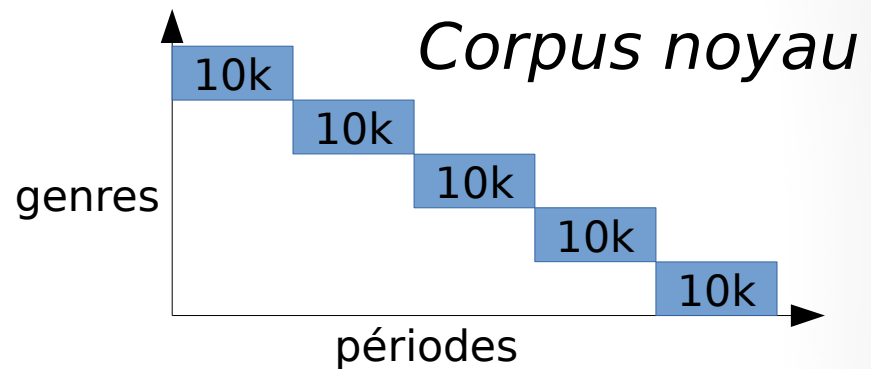
Saulsaye (1547)

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

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

<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’habit</reg></choice> ne faict
 <lb/>point le moine: & amp; tel est vestu
<choice><orig>dhabit</orig><reg>d’habit</reg></choice>
 <lb/>monachal, qui au dedans
<choice><orig>nest</orig><reg>n’est</reg></choice> rien
moins
 <lb/>que moyne: & amp; tel est vestu de cappe hes-
 <fw place="bot-center" type="sig">A iij</fw>

2. Création d'un corpus de référence

Lexique

forme_anc => lemme + 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

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

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:VMPO0PM 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

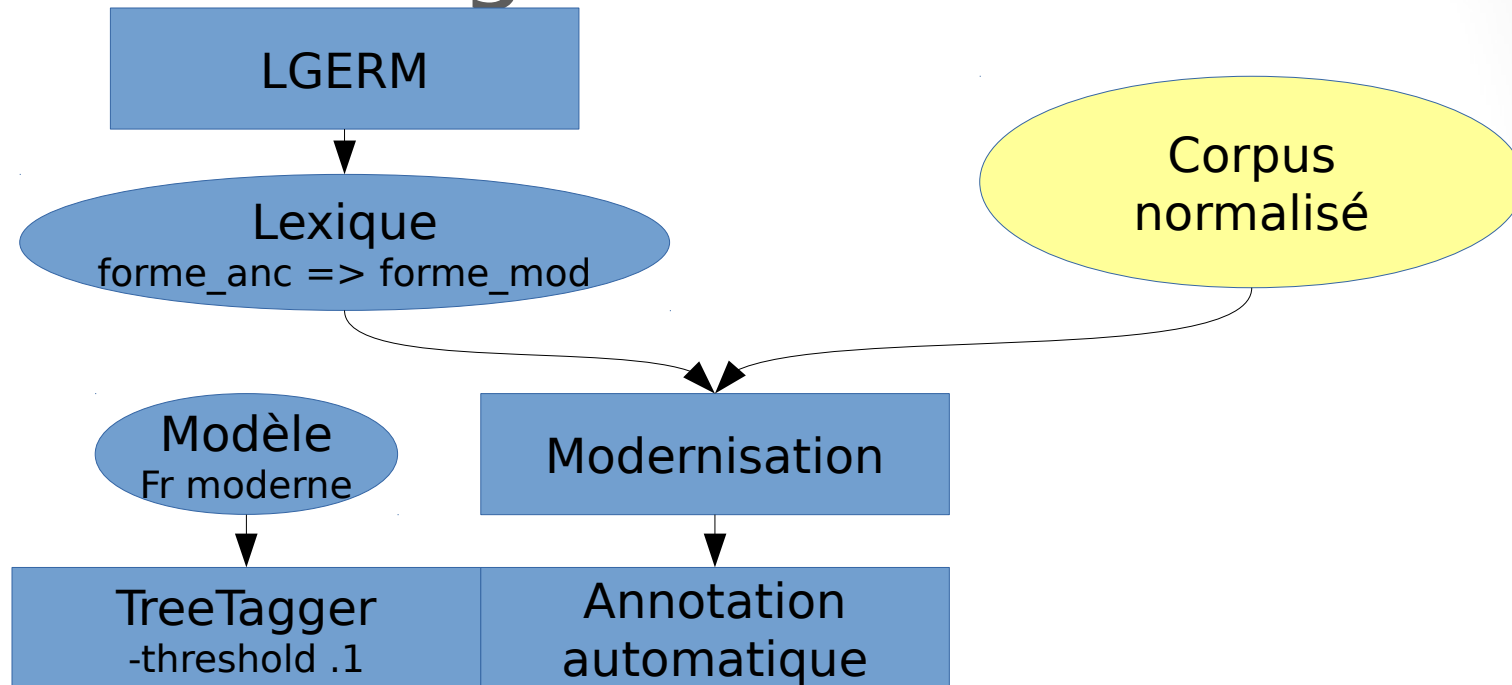


quelques	QUELQUE:Aq QUELQUE:Di
remarques	REMARQUE:Nc REMARQUER:Vvc
sur	SUR:Aq SUR:Sp SÛR:Aq SÛR:R
les	LE:Da LES:Pp LÈS:Sp LÉ:Nc
groupements	GROUPEMENT:Nc

Presto_1

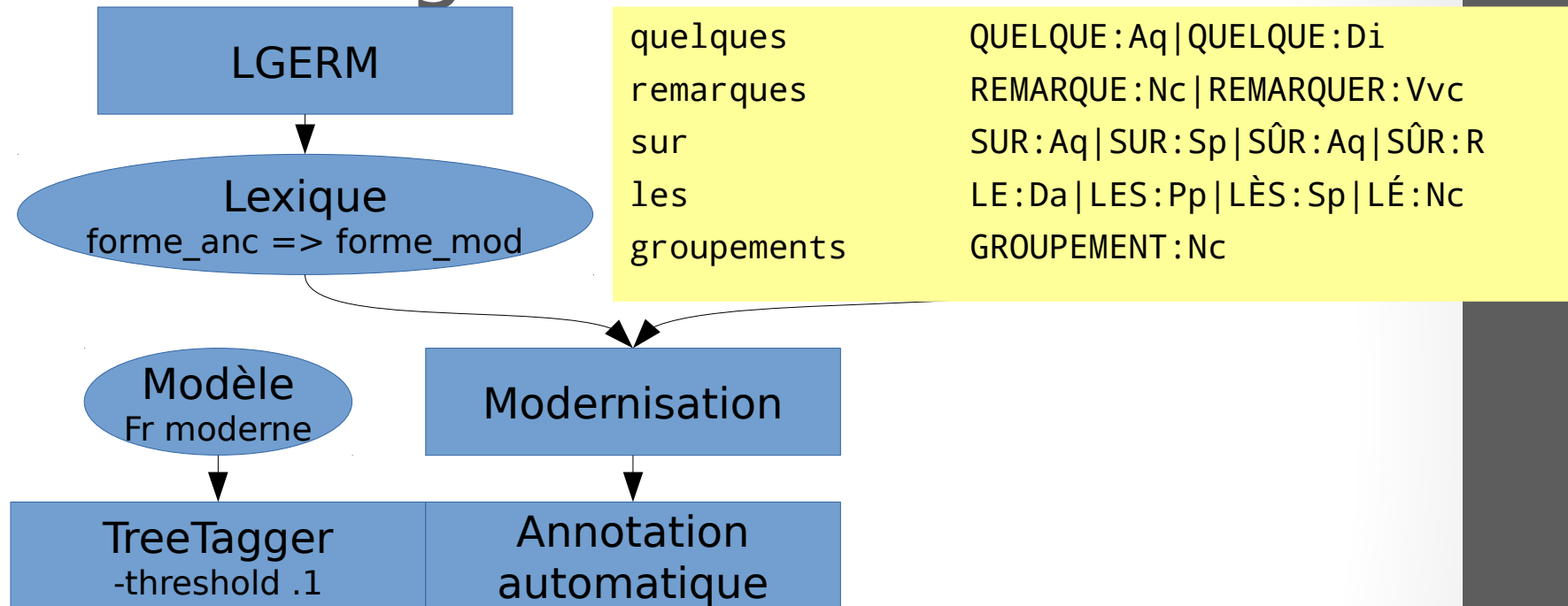
2. Création d'un corpus de référence

Désambiguïsation



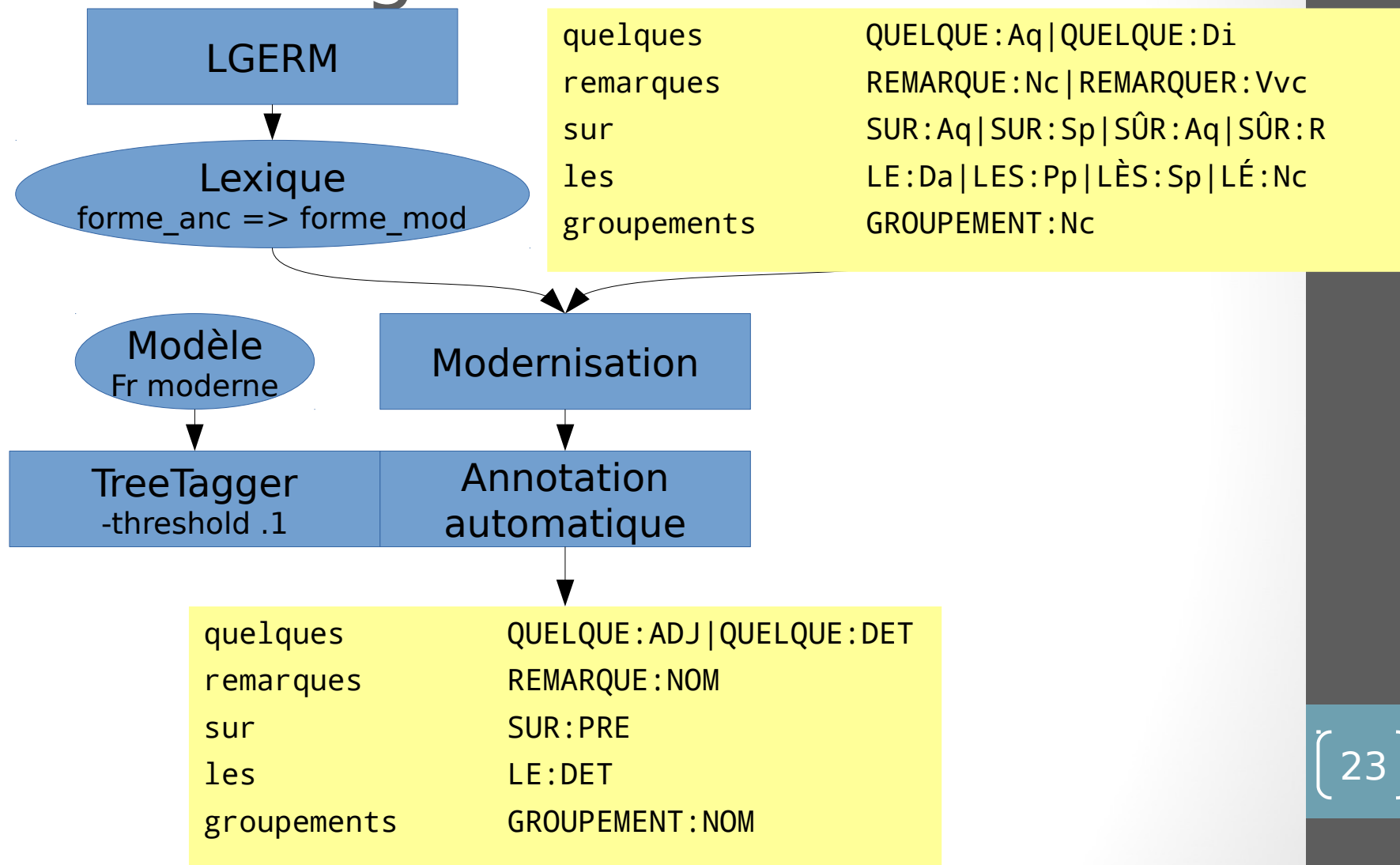
2. Création d'un corpus de référence

Désambiguïsation



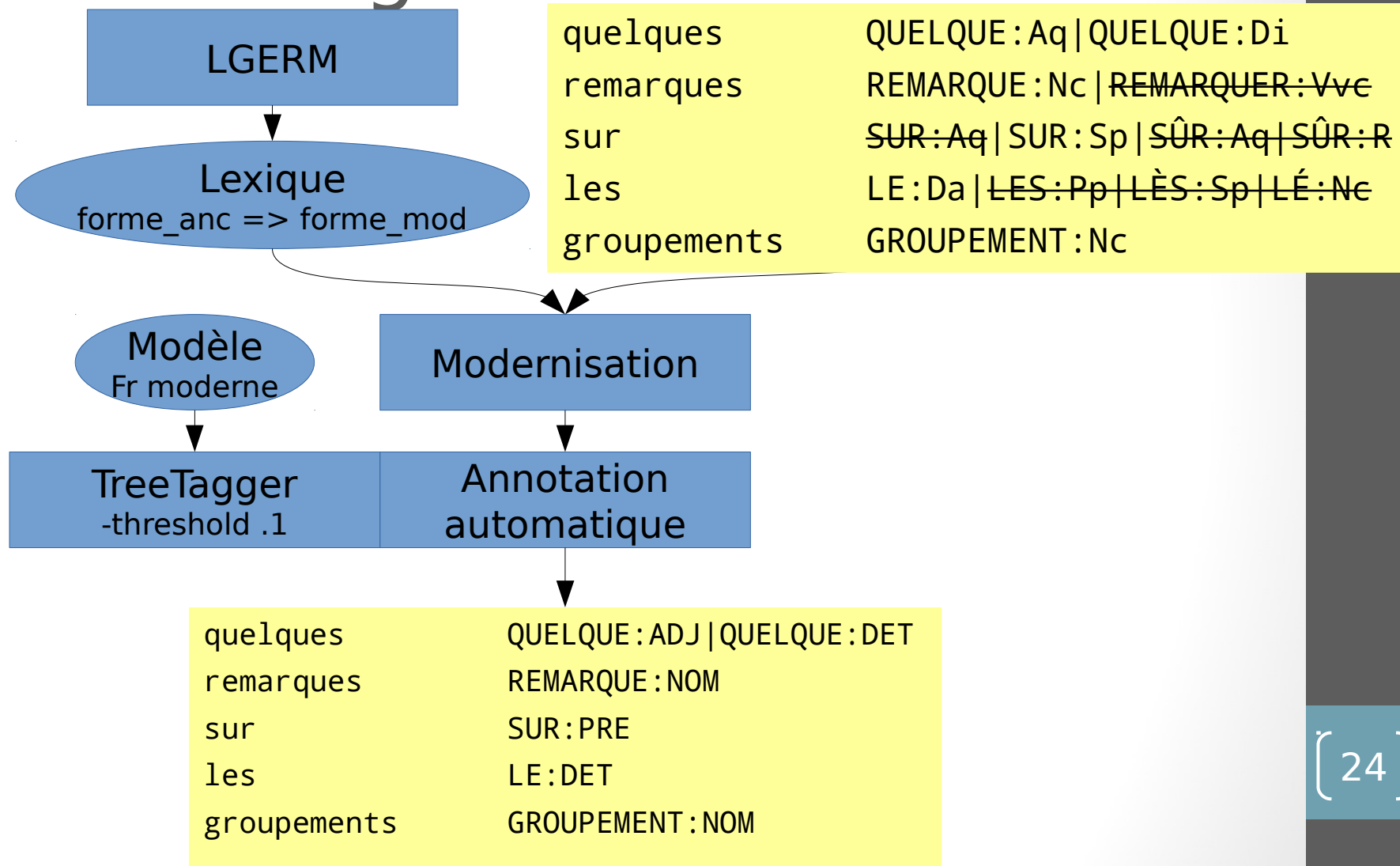
2. Création d'un corpus de référence

Désambiguïsation



2. Création d'un corpus de référence

Désambiguïsation



2. Création d'un corpus de référence

Simplification des étiquettes

quelques	QUELQUE:AQ0CP0 QUELQUE:DI0CP0
remarques	REMARQUE:NCFP000 REMARQUER:VMIP2S0 REMARQUER:VMPO0PM 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

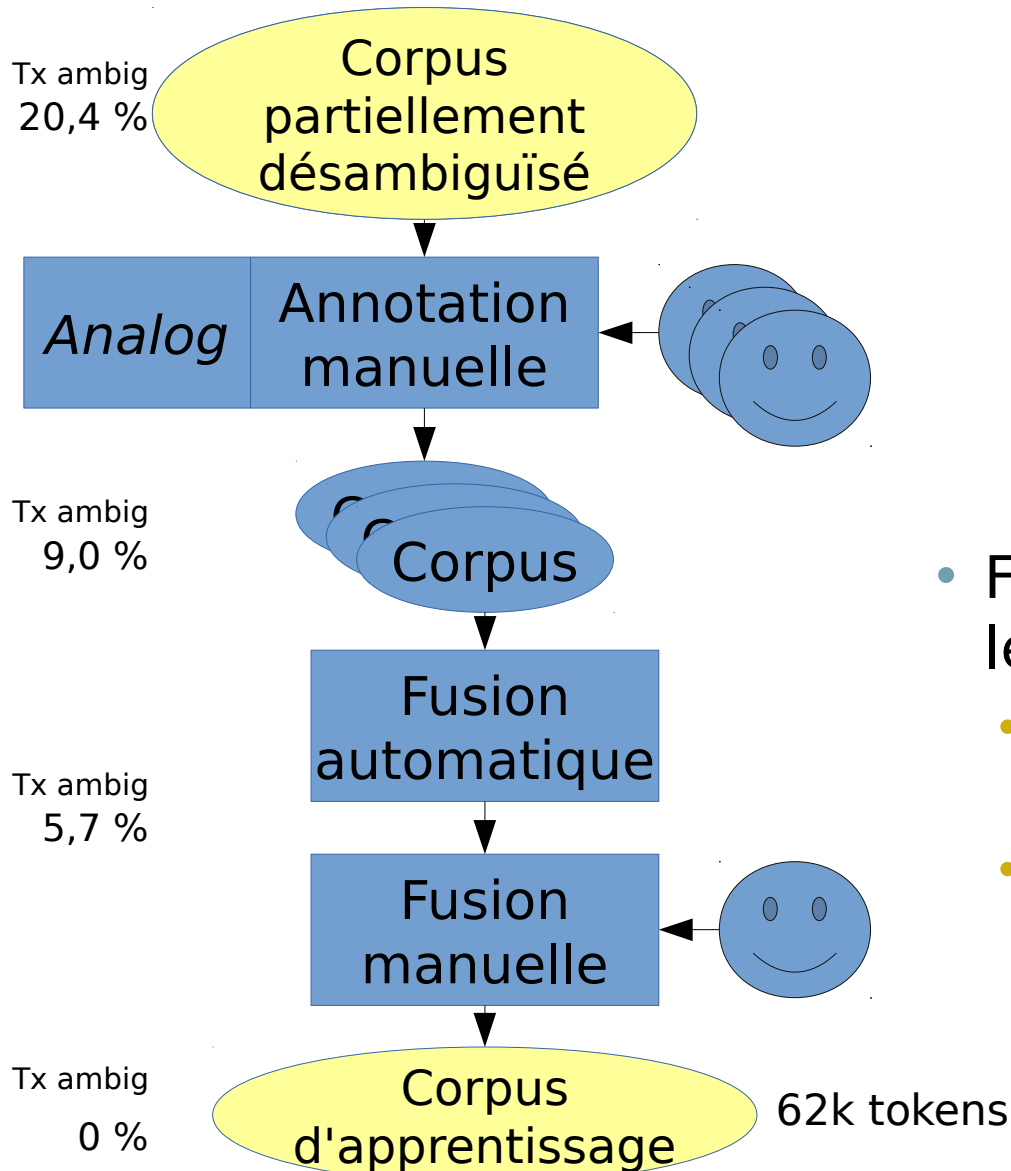


quelques	QUELQUE:Aq QUELQUE:Di
remarques	REMARQUE:Nc REMARQUER:Vvc
sur	SUR:Aq SUR:Sp SÛR:Aq SÛR:R
les	LE:Da LES:Pp LÈS:Sp LÉ:Nc
groupements	GROUPEMENT:Nc

Presto_1

2. Création d'un corpus de référence

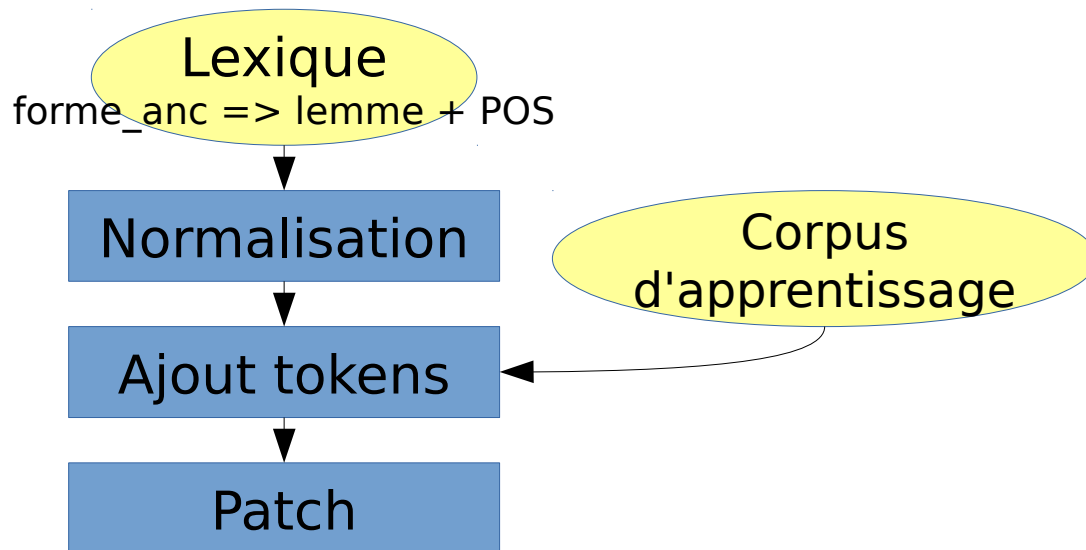
Annotation manuelle et fusion



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

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

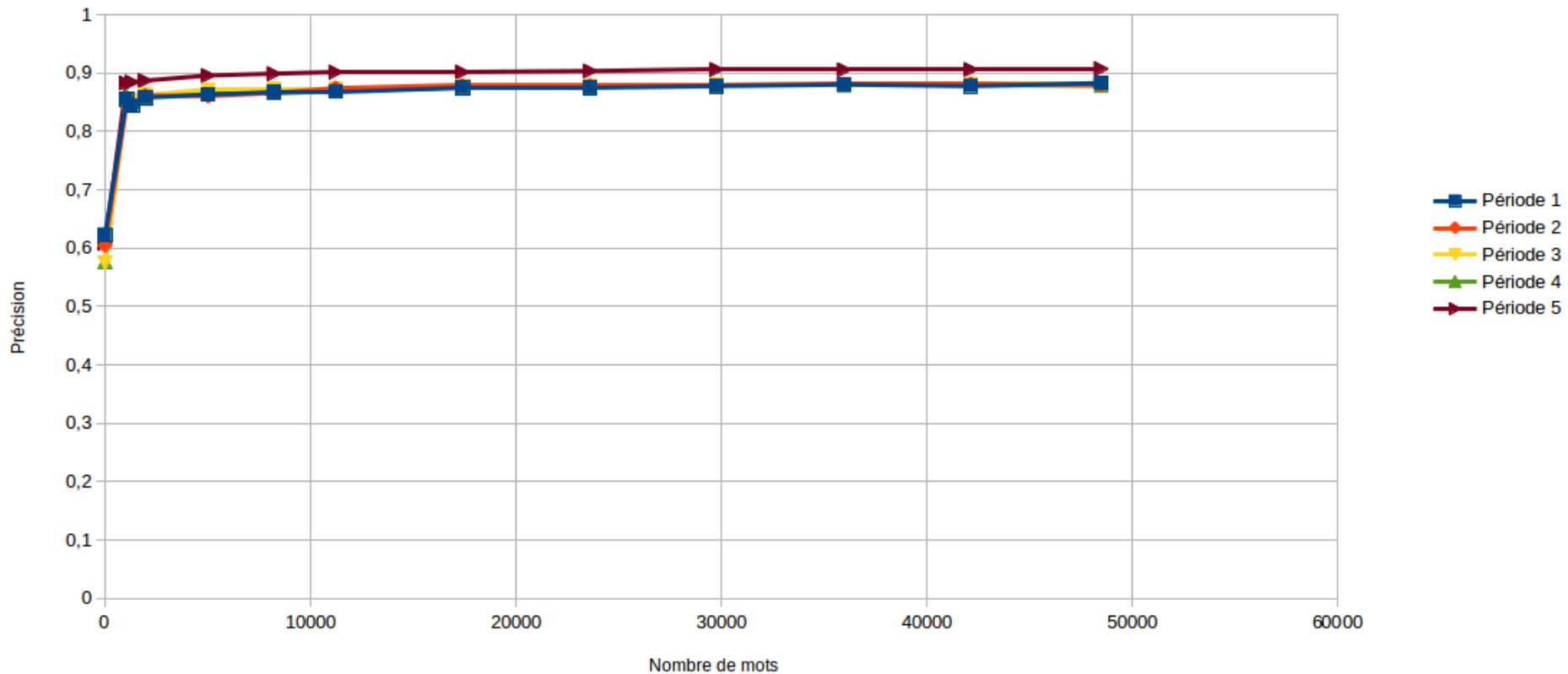
- Division du corpus d'apprentissage en trois
 - Corpus d'entraînement (80% – 49 630 tokens)
 - Corpus de développement (10% – 6 164 tokens)
 - Corpus de référence (10% – 6 110 tokens)
- *Autotuning* pour trouver les meilleurs paramètres pour TreeTagger
 - $cl\ 2 ; dtg\ 0,5 ; sw\ 1 ; ecw\ 0,06 ; atg\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

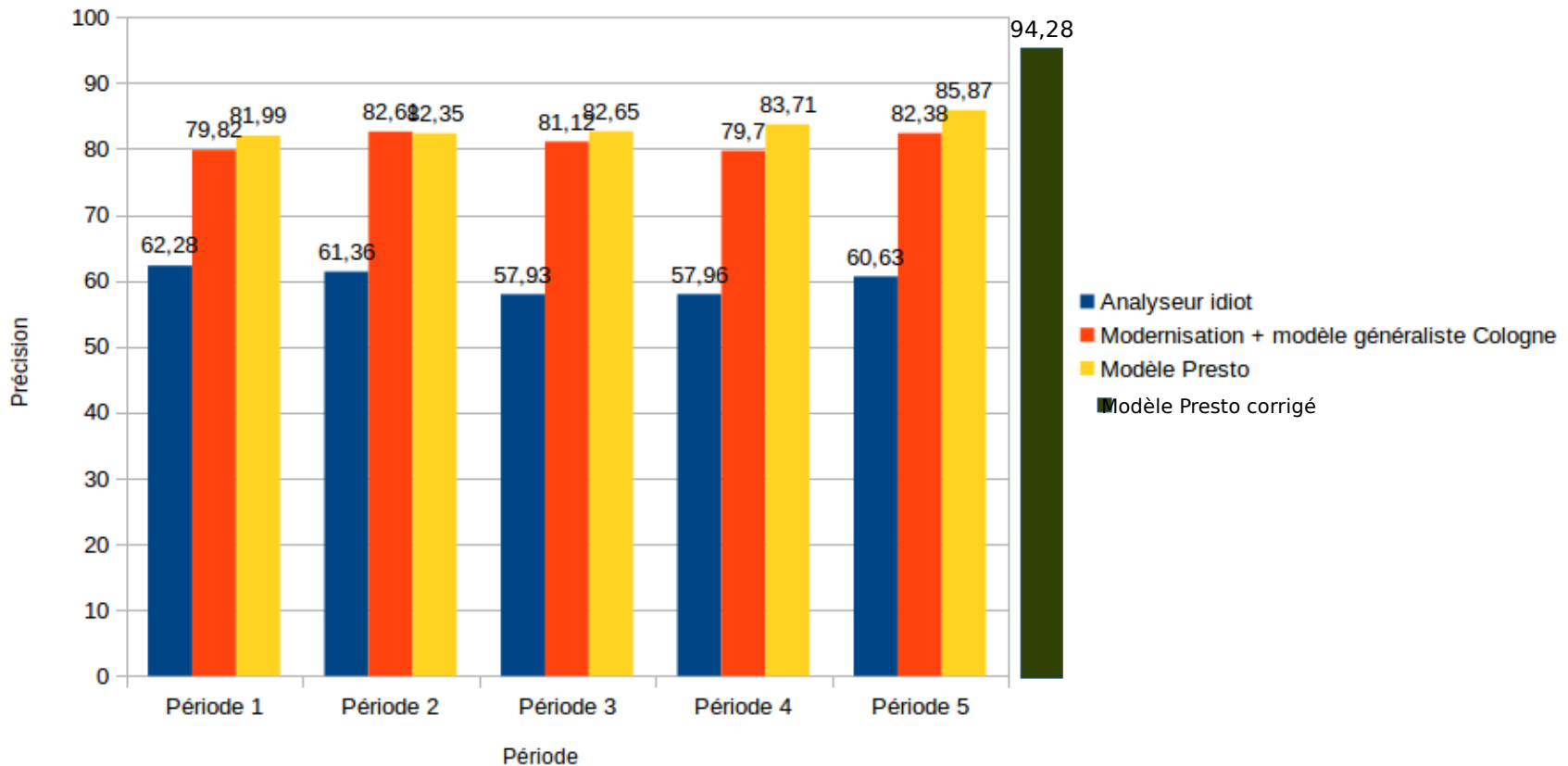
Le corpus d'apprentissage comporte toutes les périodes, on fait varier le nombre de mots.
Le baseline «0 mots» est obtenu, sans modèle, par tirage aléatoire des catégories à partir du lexique d'apprentissage.
Le corpus d'évaluation est différent pour chaque période, et comporte 761 à 1946 mots selon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



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

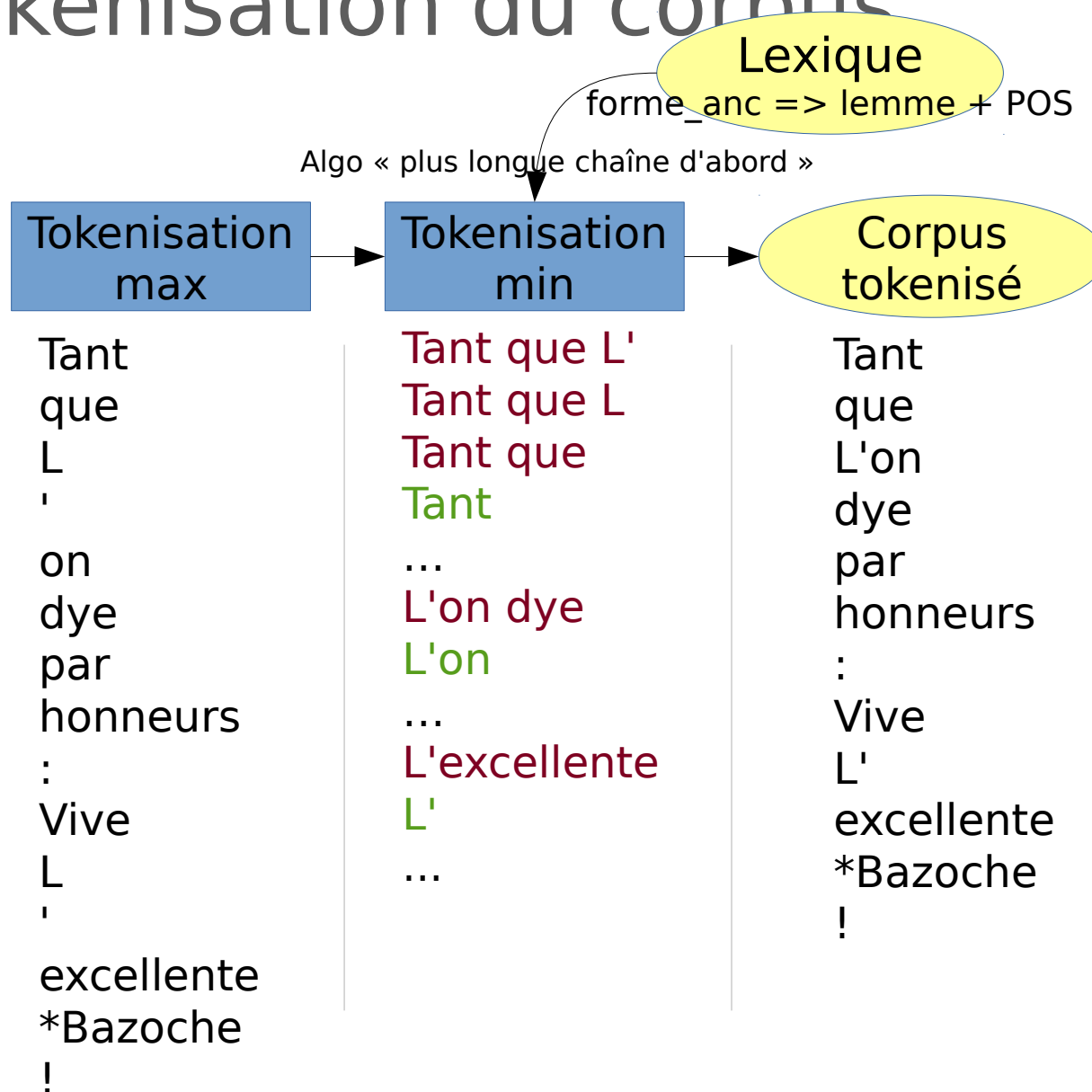
Évaluation du modèle

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



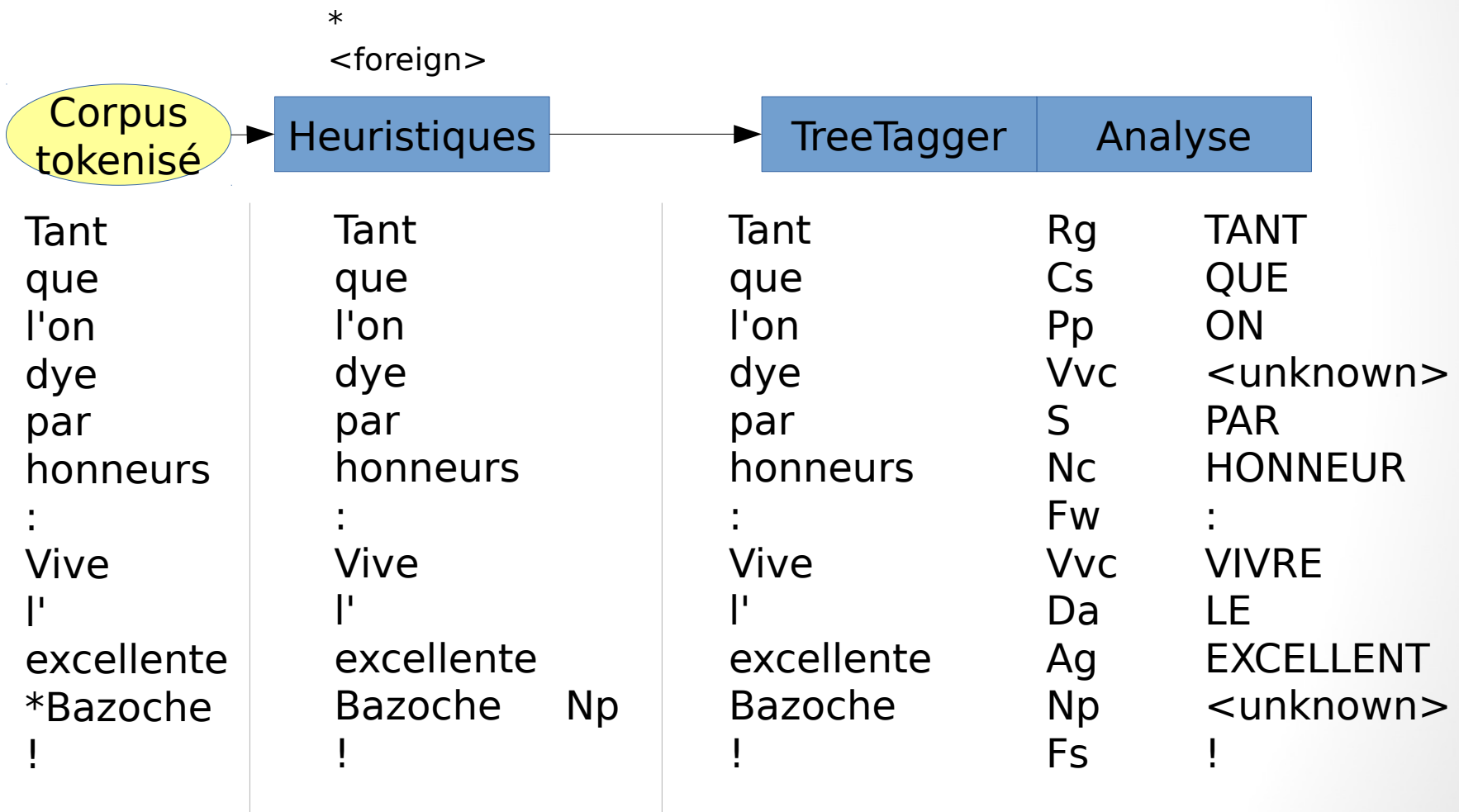
4. Analyse du corpus

Tokenisation du corpus



4. Analyse du corpus

Traitement du corpus



4. Analyse du corpus

Implémentation

- Programmation
 - Chaîne de montage : *Pipes* Unix
 - ~multitâche
 - Souple, facile à déboguer
 - Shell, Perl
- 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



4. Analyse du corpus

Perspectives

- Pêche aux erreurs
 - Vérification des ressources
- Annotation manuelle ciblée
- Règles de segmentation
 - *Trèsgrand, trèsbel*
 - *Par ce que*
- Détection de la langue
- Désambiguïstation sémantique pour les lemmes
 - 1,8 % des lemmes du corpus sont ambigus
 - *fil|fils*
 - *congre|congrès*
- Analyse syntaxique



Annotation du corpus PRESTO: création de ressources pour l'analyse du français classique

Sascha Diwersy, Universität zu Köln
Achille Falaise, ENS de Lyon, ICAR

Journées d'étude *Presto*
Cologne, 12-13 février 2015