



HAL
open science

Manuel d'annotation linguistique pour le français moderne (XVI^e -XVIII^e siècles)

Simon Gabay, Jean-Baptiste Camps, Thibault Clérice

► To cite this version:

Simon Gabay, Jean-Baptiste Camps, Thibault Clérice. Manuel d'annotation linguistique pour le français moderne (XVI^e -XVIII^e siècles) : Version B. 2022. hal-02571190v2

HAL Id: hal-02571190

<https://hal.science/hal-02571190v2>

Preprint submitted on 18 Apr 2022

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.



Distributed under a Creative Commons Attribution 4.0 International License

Manuel d'annotation linguistique pour le français moderne (XVI^e-XVIII^e siècles)

Simon Gabay¹, Jean-Baptiste Camps² et Thibault Clérice²

¹Université de Genève

²Centre Jean Mabillon (EA 3624), École nationale des chartes, Université
PSL

Version B (« Béate Béatrice »)
18 avril 2022

Sommaire

1	Segmentation en mots (<i>tokenisation</i>)	3
2	Lemmatisation	4
3	Étiquetage morpho-syntaxique (parties du discours)	7
4	Morphologie	11
5	Entités nommées	12
	Annexes	21
A	Exemple d'annotation	21
B	Abréviations pour le tokeniseur (extrait)	22
C	Jeu d'étiquette CATTEX	23
	Références	25

*Lassé de Mars j'aspire aux douceurs de la paix.
J'habite l'Austrasie aux bois les plus épais.
Là je consacre un temple à ce puissant Mercure,
Qui m'ouvre les clartez d'une science obscure :
Qui m'apprend loin du bruit les secrets curieux
Des enfers, de la mer, de la terre, et des cieux*

DESMARETS, *Clovis*, 1657

*La conversation finirait mal, ne l'entamons
point, tirons nos chausses.*

DANCOURT, *Les Vacances*, 1697

À plus!

Cri de guerre de la maison de Rohan

Introduction

Ce manuel a pour objectif d’accompagner la constitution de corpus d’entraînement pour la lemmatisation, l’étiquetage morpho-syntaxique et morphologique, et la reconnaissance d’entités nommées. Il s’intéresse principalement aux problèmes soulevés par les textes français modernes, de la Renaissance aux Lumières. Il propose des règles simples et s’appuie sur des standards nationaux (par ex. CATTEX ou *LGeRM*) et internationaux (par ex. le format BIO) pour permettre la production de données du type de celles présentées en table 1.

Token	Lemme	POS	Morphologie	Entités
Michel	Michel	NOMPro	NOMB.=s GENRE=m	PER-B
annotate	annoter	VERcjk	MODE=ind TEMPS=pst PERS.=3 NOMB.=s	0
les	le	DETdef	NOMB.=p GENRE=m	0
textes	texte	NOMcom	NOMB.=p GENRE=m	0

TABLE 1 – Exemple d’annotation

Son objectif est d’encoder des états de langues anciens sur la longue durée. Il doit donc tenir compte des particularités de ces états de langue, et assurer une interopérabilité minimale entre différentes époques tout en conservant une annotation assez riche pour permettre des requêtes fines.

1 Segmentation en mots (*tokenisation*)

La segmentation en mots obéit à quelques règles simples. Sauf erreur manifeste, qui doit être corrigée :

- Les tokens sont toutes les chaînes de caractères séparées par des espaces et des retours chariot en l’absence de tiret de fin de ligne ;
- La non agglutination (*ce pendant* vs *cependant*) n’est pas corrigée, et les locutions ne sont pas regroupées (*bien que* = 2 tokens) ;
- Les agglutinations ne sont pas non plus corrigées (*tresobeissant*, *parce* = 1 token) ;
- Les signes de ponctuation (point, virgule...) et équivalents (tiret, cadratin...) constituent, sauf exception (apostrophe, cas spéciaux mentionnés *infra*), un token.

Ces choix radicaux ont pour objectif pratique de simplifier la tokenisation d’une part, et pour objectif linguistique de simplifier l’analyse diachronique en évitant un quelconque parti pris sur les phénomènes de figement et d’agglutination, qui sont difficiles à dater précisément. Revers de la médaille : l’analyse grammaticale n’est en revanche pas toujours aisée, notamment pour les locutions – dont nous avons multiplié les exemples dans ce document afin d’aider l’utilisateur.

- Les locutions adverbiales (*à gauche*) ne forment pas un token mais deux (à PRE + gauche NOMcom) ;
- Tant qu’une espace graphique est maintenu (loc. *ce pendant*) nous conservons plusieurs tokens (*ce* PROdem + *pendre* VERppa).

Concernant l’apostrophe, elle n’est pas considérée comme indépendante et se trouve rattachée au token élidé (*c’est* → *c’* + *est*, *quelqu’un* → *quelqu’* + *un*, *l’on* → *l’* + *on*). Quelques rarissimes cas d’élisions internes au mot peuvent exister pour des raisons stylistiques, surtout en littérature contemporaine, comme *V’la c’te voiture* : ils doivent être maintenus en un seul token quand cela est possible.

Concernant le trait d’union, la situation est plus complexe. Il est considéré comme un token à part entière (*beau-frère* → *beau* + *-* + *frère*) sauf dans les cas où il est suivi :

- d’un pronom personnel (*-je*, *-tu*, *-il*, *-elle*, *-on*, *-nous*, *-vous*, *-ils*, *-elles*) ;
- d’un pronom (régime) tonique (*-moi*, *-toi*, *-lui*, *-eux*) ;
- d’un pronom (régime) atone (*-me*, *-m’*, *-te*, *-t’*, *le*, *la*, *-les*, *-leur*, *-leurs*) ;
- d’un pronom démonstratif (*-ce*)
- d’un pronom adverbial (*-en*, *-y*) ;
- d’un pronom indéfini (*-un*, *-uns*, *-une*, *-unes*) ;
- de certains adverbes (*-ci*, *-là*, *-aussi*) ;
- d’un adjectif indéfini (*-même(s)*).

Attention au cas particulier du *t* euphonique (*-t-*) uni avec le verbe qui précède : *a-t-il* est ainsi segmenté *a-t* / *-il*.

Nous considérons en effet que, dans les cas précédemment mentionnés, la présence du trait d’union est un phénomène syntaxique (par ex. la postposition) et non lexical (*grand-mère*). Il en va logiquement de même pour les noms propres, qui restent en plusieurs tokens (*Jean-Baptiste*, *Aulu-Gelle*, *Saint-Étienne-du-Mont...*).

Concernant les abréviations, il est impératif de distinguer les points qui terminent les phrases de ceux terminant une abréviation. Ainsi, dans la phrase *C’est D. Juan.*, le premier point indique que *D* est la forme abrégée de *Don*, tandis que le second marque la fermeture de la phrase. Une liste (évidemment non-exhaustive) des abréviations possibles a été définie en annexe (cf. annexe B), à laquelle il convient d’ajouter les lettres uniques suivies d’un point, qui peuvent indiquer le prénom (*D.* pour *Damien*), les points cardinaux (*N.* pour *nord*) ou des mots précis (*M.* pour *monsieur*, *S.* pour *saint*, *P.* pour *père*, *v.* pour *voir*). Il est à noter que dans certains cas des interférences sont inévitables, car *lit.* peut autant être la formée abrégée de *littérature* que le substantif *lit* en fin de phrase (*Il va au lit.*). Le point qui marque l’abréviation n’est pas considéré comme un token et reste avec le mot abrégé (*M. le prince* → *M. le prince* et non *M . le prince*).

2 Lemmatisation

Le choix du lemme n’est pas libre : il doit autant que possible se trouver dans le référentiel qui dérive de *LGeRM MODE* (i.e. *Moderne étendu*)¹. Ce choix permet de garantir une interopérabilité minimum en amont avec *LGeRM AF* (i.e. *Ancien Français*)²

1. Sascha Diwersy, Achille Falaise, Marie-Hélène Lay et Gilles Souvay, « Ressources et méthodes pour l’analyse diachronique », *Langages*, N° 206-2 (août 2017), p. 21-44, URL : <https://www.cairn.info/revue-langages-2017-2-page-21.htm> (visité le 03/12/2018).

2. G. Souvay et Jean-Marie Pierrel, « LGeRM Lemmatisation des mots en Moyen Français », *Traitement Automatique des Langues*, 50-2 (2009), p. 149-172, URL : <https://halshs.archives-ouvertes.fr/halshs-00396452>.

et en aval avec le *TLFi*³. La lemmatisation obéit à quelques grandes règles simples :

- Le lemme est, autant que possible, la forme contemporaine du mot et non une ancienne forme (*avecque* → *avec*), le masculin singulier pour les adjectifs et les substantifs, l’infinitif pour les verbes ;
- Dans le cas où il existe une forme masculine et féminine d’un même mot, le lemme est dans la très grande majorité des cas la forme masculine pour les substantifs (*comtesse* → *comte*) comme pour les adjectifs (*grande* → *grand*). Ce n’est en revanche pas le cas pour les noms propres (*Jeanne* → *Jeanne* et non *Jean*). Il existe quelques exceptions qui possèdent deux lemmes différents (*dame* vs ancien français *don*) malgré une racine commune (< *dominus*, *a*) pour conserver un interopérabilité dans le temps (*don* ayant disparu en français) ;
- Si le lemme est absent de LGeRM, l’existence d’une entrée dans un dictionnaire historique (préférentiellement celui de Furetière, mais aussi ceux de Richelet ou même de l’Académie) ou dans un dictionnaire scientifique (comme l’*Altfranzösisches Wörterbuch* de Tobler et Lommatzsch ou le *Dictionnaire du Moyen Français*) est le principal critère pour l’ajout d’un nouveau lemme dans le référentiel ;
- Dans le cas de doublon dans LGeRM, on se réfère aux entrées du *DMF* et du *TLFi* (en privilégiant le premier au second) ;
- Certains tokens posent problème :
 - Ils ne sont pas analysable hors contexte, comme *parce* (→ dans *parce que*), *afin* (dans *afin que*), ou *ledit*. Nous créons alors un lemme (*parce*, *afin*, *ledit*) ;
 - Ce sont des enclises (*du*, *des*, *au*, *dudit*, *auquel*) : les deux lemmes originels sont alors conservés et séparés par un tiret bas (*de_le* ou *à_le*) ;
 - Ce sont des mots-valises (*tresobeissant...*) : nous utilisons la même méthode que pour les enclises (*tres_obeissant...*) ;
- Comme pour les autres tokens, nous considérons que le lemme des noms propres est leur forme contemporaine. Nous privilégions des formes communes malgré des variations diachroniques (*Jehanne* → *Jeanne*) ou graphiques à la marge (*Denys* → *Denis*, *Remus* → *Rémus...*). Si la forme est clairement dans la langue étrangère, le lemme est alors dans cette langue (*Jan* → *Jan...* et non *Jean*, *Vespasianus* → *Vespasianus...* et non *Vespasien*, *Demosthenes* → *Demosthenes* et non *Démosthène*) ;
- Nous considérons comme nom propre un token qui commence par une majuscule : ainsi *la mer Noire* → *le mer Noir*, *madame de Sévigné* → *madame de Sévigné* (mais *madame De Sévigné* → *madame De Sévigné*). L’annotation des entités nommées au format BIO permet de capturer des ensembles plus larges, contenant titres, prénoms, etc. (voir sect. 5) ;
- Le lemme retenu pour les pronoms personnels (sujet, réfléchi, objet direct, indirect ou disjoint) est la forme du pronom personnel sujet (*moi* → *j_e*), le cas échéant singulier (*eux* → *il*) masculin (*elles* → *il*). L’annotation morphologique (voir sect. 4) fournit le complément d’information nécessaire pour des requêtes fines ;
- Nous ne connaissons pas l’adjektivation des participes, car le passage d’une catégorie à l’autre, surtout en diachronie, est beaucoup trop difficile à identifier précisément.

3. J.M. Pierrel, Jacques Dendien et Pascale Bernard, « Le TLFi ou Trésor de la Langue Française informatisé », dans *Proceedings of the 11th EURALEX International Congress*, dir. Geoffrey Williams et Sandra Vessier, Lorient, France, 2004, p. 165-170.

Si un infinitif existe, nous considérons donc qu'il s'agit d'un participe et l'infinitif sert de lemme (*le retour éclatant* → *éclater*, *une âme affligée* → *affliger*). Notons que la présence d'une marque de flexion (*des traits charmants*) n'est pas considérée comme une raison suffisante pour en faire une forme autonome (*charmant* → *charmer* et non *charmant*) pour des raisons de cohérence ;

- Le lemme d'un chiffre en toutes lettres est ce chiffre en toutes lettres (*quatre* → *quatre* et non 4). Le lemme d'un nombre en chiffres arabes ou romains écrit (par ex. *XII*, *12* → 12) est la version arabe du chiffre ;
- Pour des langues étrangères, il est recommandé de ne pas proposer de lemme et d'utiliser une forme générique formée de @ et du nom de la langue en anglais (*animalium* → @latin, *prima* → @italian). Si un lemme doit être déterminé, en accord avec les principes éditoriaux que nous exposons ici, on privilégiera l'infinitif pour le verbe, le singulier masculin (nominatif) pour les substantifs ou adjectifs. On aura de préférence recours à un référentiel existant (par exemple, pour le latin, un dérivé de *Forcellini*⁴), ou à défaut d'un dictionnaire faisant autorité (comme le *DMF*⁵ pour le moyen français) ;
- On respecte les ligatures pour les lemmes, peu importe si elles sont dans l'occurrence analysée ou pas (*noeud* → *nœud* et non *noeud*, *voeu* → *vœu* et non *voeu*, etc.) ;
- Les majuscules sont accentuées (*Etiopie* → *Éthiopie*) ;
- Les lemmes des verbes pronominaux ne contiennent pas le pronom, car cela serait redondant avec le pronom qui est tokenisé à part (*s'enfuir* → *se* + *enfuir*), et que cela permet de ranger sous une même étiquette différents emplois d'un même verbe (emploi transitif vs pronominal d'*abaisser*) ;
- Dans le cas où il existe deux formes très proches dont l'une s'est imposée, notamment pour les états de langue les plus anciens (par ex. *finablement* vs *finale-ment*), on conserve bien deux lemmes distincts (en l'occurrence *finablement* et *finale-ment*) ;
- Pour les mots abrégés, nous ne développons le token que s'il s'agit d'un substantif ou un adjectif (*M.* → *monsieur*), mais pas s'il s'agit d'un nom propre (*le Père R.* avec *R.* pour *Rapin* → *le père R.* et non *le père Rapin*, mais *le P. Rapin* → *le père Rapin*).

Quelques exemples pour le débutant

Des exemples simples pour commencer :

- *je vais à Genève* → *je aller à Genève* ;
- *un sort assez propice* → *un sort assez propice* ;
- *il m'a sauvé toutefois des ravages du temps* → *il je avoir sauver toutefois de_le ravage de_le temps* ;
- *c'est que Vespasian me regardait pour lui* → *ce être que Vespasien je regarder pour il* ;

4. Thibault Clérico, *Référentiel du Latin pour Pyrrba*, d'après le dictionnaire et travaux du LASLA de D. Longrée et al, mai 2020, doi : 10.5281/zenodo.3822040.

5. ATILF-CNRS et Université de Lorraine, *Dictionnaire du Moyen Français (1330-1500)*, 2015, URL : <http://www.atilf.fr/dmf>.

- *vous sentez-vous impropre au **matrimonium*** → matrimonium;
- *Lettre 12. 16 septembre 1676. À M. R.* → lettre 12 . 16 septembre 1676 . à monsieur R.

Il est essentiel de faire attention aux homographes :

- Dans *soit...soit* ou *tu le veux?* **soit** le lemme est soit. En revanche, dans *ainsi soit-il* il s'agit évidemment du verbe être;
- Dans *il ne leur manquera rien* il s'agit du pronom il mais dans *leur vif éclat* il s'agit du déterminant possessif leur;
- Dans *la fin* il s'agit du déterminant le mais dans *je la vois* il s'agit du pronom il.

Attention aussi aux lemmes composés :

- *dudit* → de_ledit;
- *au(x)* → à_le;
- *duquel* → de_lequel;
- *auxquels* → à_lequel;
- *du* → de_le;
- *des* → de_le (ou un, comme dans *des beaux garçons sont arrivés* : attention au contexte!).

Les entités nommées sont complexes à gérer :

- *don Carlos* → don Carlos;
- *la mer Rouge* → le mer Rouge;
- *Monsieur de La Rochefoucauld* → monsieur de La Rochefoucauld;
- *François, duc d'Enghien* → François , duc de Enghien;
- *M. Du Plessis* → monsieur Du Plessis (et non de_le pour *Du*, qui a d'ailleurs une majuscule en français contemporain);
- *Denys d'Halicarnasse* → Denis de Halicarnasse.

3 Étiquetage morpho-syntaxique (parties du discours)

Nous reprenons le jeu d'étiquettes *CATTEX-max*⁶ (qui inclut la morphologie, sur laquelle nous revenons *infra*) et dont nous rappelons les étiquettes en annexe (voir annexe C). Concernant les principes d'annotation, on se reportera au manuel détaillé conçu par ses créateurs⁷. Nous nous bornons ici à rappeler quelques grandes règles. Le choix de *CATTEX* plutôt que d'un autre système permet de garantir l'interopérabilité avec les corpus médiévaux comme la *Base de français médiéval*⁸ ou ceux développés à l'École des chartes⁹.

6. Sophie Prévost, Céline Guillot, Alexei Lavrentiev et Serge Heiden, *Jeu d'étiquettes morphosyntaxiques CATTEX2009*, rapp. tech., version 2.0. http://bfm.ens-lyon.fr/IMG/pdf/Cattex2009_2.0.pdf, Lyon, École normale supérieure de Lyon, 2013.

7. C. Guillot, S. Prévost et A. Lavrentiev, *Principes d'annotation Cattex09*, rapp. tech., version 2.0. http://bfm.ens-lyon.fr/IMG/pdf/Cattex2009_principes_2.0.pdf, Lyon, École normale supérieure de Lyon, 2013.

8. C. Guillot, S. Heiden et A. Lavrentiev, « Base de français médiéval : une base de référence de sources médiévales ouverte et libre au service de la communauté scientifique », *Diachroniques : Revue de Linguistique française diachronique-7* (déc. 2017), p. 168-184, URL : <https://halshs.archives-ouvertes.fr/halshs-01809581> (visité le 19/12/2019).

9. Jean-Baptiste Camps (éd.), *Geste : un corpus de chansons de geste, 2016-... (Version 02)*, Paris, 2016-2020, URL : <http://doi.org/10.5281/zenodo.2630574>, ainsi que les données provenant des thèses de doctorat, en préparation d'Ariane Pinche (*Édition nativement numérique du recueil hagiographique Li Seint Confessor de*

- L’annotation n’est pas morphologique mais morpho-syntaxique, c’est à dire que l’identification des catégories grammaticales est faite en contexte ;
- Le principe général veut que :
 - un token qui précède un nom sans adjectif est un déterminant ;
 - un token précédé d’un déterminant et suivi d’un substantif est un adjectif ;
 - en l’absence d’un substantif, le token est un pronom ;
- Les adjectifs ou les adverbes peuvent être substantivés – ce que l’on reconnaît à la présence d’un déterminant. Ainsi *beau* est ADJqua mais peut être substantivé *un beau* → NOMcom. De même *paravant* est adverbe, mais *au paravant* → PRE.DETdef NOMcom ;
- Parce que la stabilisation des noms de lieux ou de personnes est trop complexe (est-ce une particule ou une simple préposition ?), nous considérons arbitrairement comme nom propre un token qui commence par une majuscule : ainsi *la mer Noire* → DETdef NOMcom NOMpro, *madame de Sévigné* → NOMcom PRE NOMpro (mais *madame De Sévigné* → NOMcom NOMpro NOMpro). L’annotation des entités au format BIO permet de capturer des ensembles plus larges, contenant titres, prénoms... (voir sect. 5) ;
- La personnification (marquée par la majuscule, voir *supra*) entraîne automatiquement la requalification en nom propre : dans *la Fortune* → NOMpro et pas NOMcom) ;
- Les tokens *Dieu* et *dieu* peuvent être NOMpro si la détermination est non pertinente (*Dieu merci*), potentiellement avec un adjectif (*le bon Dieu*). Dans les autres cas le token est NOMcom (*le dieu Jupiter*) ;
- Il n’y a pas d’adjectivation du participe : si un infinitif existe, nous considérons qu’il s’agit d’un participe (*un retour éclatant* → VERppa, *une âme affligée* → VERppe). Notons que la présence d’une marque de flexion (*des traits charmants*) n’est pas considérée comme une raison suffisante pour en faire une forme autonome (*charmant* → VERppa) par souci de cohérence ;
- Pour les locutions conjonctives : *sauf* ou *sans* suivis de *que* sont toujours PRE, *puis* ou *outre* suivis de *que* sont toujours ADVgen ;
- Les tokens qui sont le résultat d’une soudure (*parce < par ce*, *afin < à fin*) et n’existent que dans une locution sont analysés par analogie : *parce que* → ADVgen car *parce* ≈ *bien* dans *bien que*. Cela permet d’éviter la création de classes d’étiquette d’effectif trop faible du type *afin* → PRE.NOMcom ou *parce* → PRE.PROdem ;
- La graphie joue un rôle important dans l’annotation. Ainsi, *d’avantage* (→ PRE NOMcom) n’est pas la même chose que *davantage* (→ ADVgen). De même, *au par avant* (→ PRE.DETdef PRE ADVgen) n’est pas la même chose que *au paravant* (→ PRE.DETdef NOMcom) ni que *auparavant* (→ ADVgen). Dans des cas aberrants (*n’aguères*) il convient de retokeniser (*n’aguères* → ADVneg VERc jg ADVgen ou *naguère* → ADVgen) ;
- Les tokens *voici* (< *vois ci*) et *voilà* (< *vois là*) sont considérés comme des verbes conjugués (VERc jg). *Idem* pour *vive* (*vive / vivent les vacances*) où *vive(nt)* est VERc jg ;

Wauchier de Denain d’après le manuscrit 412 de la Bibliothèque nationale de France, dir. Corinne Pierreville et Bruno Bureau, Univ. Lyon 3) et Lucence Ing (*Disparitions lexicales en diachronie : traitements automatiques sur le Lancelot en prose*, dir. Frédéric Duval, École nationale des Chartes).

- Pour les chiffres, rappelons que le traditionnel adjectif cardinal est annoté comme déterminant cardinal s'il suivi d'un substantif (**18 ans** → DETcar NOMcom) sauf s'il est précédé d'un déterminant (*ses 18 ans* → DETpos ADJcar NOMcom). Les cardinaux qui suivent une référence éditoriale (*lettre 1, tome 2...*) ou une référence temporelle (*l'an 1671*) sont considérés comme des adjectif cardinaux (ADJcar). En revanche, dans les dates données en style moderne, on considère qu'il y a ellipse des noms et que les cardinaux sont donc des pronoms : *i.e., le dixième jour de mars de l'an 1632* → DETdef ADJord NOMcom PRE NOMcom PRE DETdef NOMcom ADJcar, mais en revanche, *le 10 mars 1632* → DETdef PROcar NOMcom PROcar ;
- Les emprunts ou les passages en langue étrangère sont annotés avec ETR ;
- Les mots abrégés sont étiquetés avec ABR seulement s'il est impossible de connaître leur partie du discours. Dans les cas où il est possible de savoir s'il s'agit d'un nom propre (*le père R.*), d'un substantif (*M. Dupuis*), d'un verbe (*voy. page 12*), on utilise ces étiquettes.

Quelques exemples pour le débutant

Des exemples simples pour commencer :

- *je vais à Metz* → PROper VERcjg PRE NOMpro ;
- *un sort assez propice* → DETndf NOMcom ADVgen ADJqua ;
- *il m'a sauvé toutefois des ravages du temps* → PROper PROper VERcjg VERppe ADVgen PRE.DETdef NOMcom PRE.DETdef NOMcom ;
- *c'est que Vespasian me regardait pour lui* → PROdem VERcjg CONsub NOMpro PROper VERcjg PRE PROper ;
- *Lettre 12. 16 septembre 1676. À M. R.* → NOMcom ADJcar PONfrt PROcar NOMcom PROcar PONfrt PRE NOMcom NOMpro.

Il est essentiel de faire attention aux homographes :

- *soit...soit* → CONcoo mais *tu le veux? Soit* → ADVgen ou bien encore *ainsi soit-il* → VERcjg ;
- *il ne leur manquera rien* → PROper mais *leur vif éclat* → DETpos ;
- *la fin* → DETdef mais *je la vois* → PROper ;
- *il mange* → PROper mais *il semble* → PROimp ;
- *je cours même sort* → DETind mais *faire de même* → ADVgen et *le même sort* → ADJind ;
- *en tel désespoir* → DETind mais *un tel objet* → ADJind ;
- *un jour* → DETndf mais *l'un, l'autre* → PROind ;
- *ce jour* → DETdem mais *c'est* → PROdem ;
- *15 juillet* → PROcar mais *15 jours* → DETcar et *chapitre 15* → ADJcar.

Attention aussi aux lemmes composés :

- *dudit* → PRE.DETcom ;
- *au(x)* → PRE.DETdef ;
- *duquel* → PRE.DETrel, PRE.PROrel ou PRE.PROint en fonction du contexte ;
- *auxquels* → PRE.DETrel, PRE.PROrel ou PRE.PROint en fonction du contexte ;
- *du* → PRE.DETdef ;

— *des* → PRE.DETdef (ou DETndf, comme dans *il voit des arbres*: attention au contexte!).

Deux types de tokens sont particulièrement problématiques.

— Les noms propres et assimilés (titres...) :

- *don Carlos* → NOMcom NOMpro;
- *la mer Rouge* → DETdef NOMcom ADJqua;
- *Monsieur de La Rochefoucauld* → NOMcom PRE NOMpro NOMpro;
- *François, duc d'Engbien* → NOMpro PONfbl NOMcom PRE NOMpro;
- *M. Du Plessis* → NOMcom NOMpro NOMpro (et non PRE.DETdef pour *Du*, qui a d'ailleurs une majuscule en français contemporain);
- *Denys d'Halycarnasse* → NOMpro PRE NOMpro;
- *Mesnil montant* → NOMpro VERppa;

— Les locutions :

- *afin de* → ADVgen PRE;
- *afin que* → ADVgen CONsub;
- *à fin que* → PRE NOMcom CONsub;
- *puis que* → ADVgen CONsub;
- *puis donc que* → ADVgen ADVgen CONsub;
- *de ce que* → PRE PROdem CONsub;
- *lors que* → ADVgen CONsub;
- *pourvu que* → ADVgen CONsub;
- *vu que* → VERppe CONsub;
- *tandis que* → ADVgen CONsub;
- *selon que* → ADVgen CONsub;
- *bien que* → ADVgen CONsub;
- *pres que* → ADVgen CONsub;
- *parce que* → ADVgen CONsub;
- *par ce que* → PRE PROdem CONsub;
- *pource que* → ADVgen CONsub;
- *pour ce que* → PRE PROdem CONsub;
- *quant à* → ADVgen PRE;
- *par tout* → PRE PROind;
- *à peine* → PRE NOMcom;
- *tout à fait* → ADVgen PRE NOMcom;
- *la plus part* → DETdef ADVgen NOMcom;
- *au paravant* → PRE.DETdef NOMcom;
- *tout à coup* → ADVgen PRE NOMcom;
- *d'avantage* → PRE NOMcom;
- *n'aguères* → à retokeniser;
- *auprès de* → ADV PRE;
- *là-dessus* → ADVgen PONfbl ADVgen;
- *en dessus* → PRE ADVgen;
- *au dessus* → PRE.DETdef NOMcom;
- *au moins* → PRE.DETdef NOMcom;
- *du moins* → PRE.DETdef NOMcom;
- *auprès de* → PRE PRE.

Ajoutons quelques cas particuliers :

- *premier* est ADJord ;
- *dernier* est ADJqua ;
- *16 septembre 1676* → PROcar NOMcom PROcar.

4 Morphologie

Le recours à *CATTEX-max* implique l'étiquetage morphologique précis des tokens en plus de la catégorie grammaticale.

Catégorie	Valeurs possibles
GENRE	<i>m, f, n, x</i>
NOMB	<i>s, p, x</i>
MODE	<i>ind, imp, con, sub</i>
TEMPS	<i>pst, ipf, fut, psp</i>
PERS	<i>1, 2, 3</i>
CAS	<i>n, r, i</i>

TABLE 2 – Valeurs possibles pour la morphologie, auxquelles il faut rajouter MORPH=empty.

Quelques remarques générales :

- Les différents emplois du pronom ayant tous un même lemme (*je, me, moi* → *je*) on utilise le cas pour les distinguer (respectivement CAS=n pour le nominatif, CAS=r pour le régime direct et CAS=i pour le régime indirect) ;
- La question du nombre des possessifs est complexe, car le choix de *CATTEX* de ne retenir que trois personnes (1, 2 et 3, pluriel ou singulier) pose problème : dans ces conditions, le nombre du possessif est-il celui de la personne ou de son référent ? Ainsi *mes* est-ce la première personne du singulier (PERS.=1 | NOMB.=s) ou un déterminant possessif dont le référent est un pluriel (PERS.=1 | NOMB.=p). Nous avons retenu la seconde option ;
- Dans le cas où le contexte immédiat ne permet pas de désambigüiser l'information, on laisse la valeur *x* : *vous êtes odieux* → GENRE=x et *odieux* NOMB.=x, *je cherche un enfant qui joue dehors* → MODE=x. On considère que c'est toujours le cas pour les pronoms personnels sans marque morphologique de genre (*i.e.*, *je, tu, l'* → GENRE=x, mais *la* → GENRE=f) ;
- Il ne faut pas confondre GENRE=x et GENRE=n. Les adjectifs qualificatifs résultant de pronoms impersonnels (*il est clair*) sont décrits comme de genre neutre (NOMB.=s | GENRE=n). Il en va de même pour le pronom impersonnel (*il est clair*), auquel on ajoute qu'il n'a pas de personne (PERS.=0 | NOMB.=s | GENRE=n | CAS=n). C'est aussi le cas du pronom démonstratif neutre *ce* ;
- Le conditionnel est analysé comme une mode et non comme un temps de l'indicatif (MODE=con). Comme l'impératif (MODE=imp), il n'a pas de temps : *Viens ici* → MODE=imp | PERS.=2 | NOMB.=s ;

- Nous ne connaissons pas les temps composés : *J'ai mangé* est composé d'un indicatif présent (VERc jg MODE=ind | TEMPS=pst) et d'un participe passé (VERppe) ;
- Les tokens *voici* et *voilà* sont des indicatifs présent sans personne ni nombre (MODE=ind | TEMPS=pst | PERS.=x | NOMB.=x) ;
- Pour certains tokens (PRE, PON, ADV...) il n'existe pas d'information morphologique : on met la valeur MORPH=empty ;
- Pour les pronoms neutres (*il* et *on*), il sont analysé analysés comme sans genre (GENRE=x) et cas sujet (CAS=n) à la troisième personne du singulier et non personne zéro (PERS.=3 | NOMB.=s) ;
- Pour le pronom relatif, on laisse l'annotation de *que* et *qui*, qui sont invariables vide (MORPH=empty). Pour *quoi*, qui est systématiquement neutre, annoter en conséquence (GENRE=n). Pour les pronoms *où* et *dont*, qui correspondent systématiquement au cas oblique, indiquer que le cas est régime indirect (CAS=i) ;
- Pour les formes en *-ant* : ajouter le genre et le nombre s'il est marqué (*obéissante* ou *obéissants*), sinon (gérondif ou flexion zéro) préférer x (NOMB.=x | GENRE=x).
- Contrairement aux prénoms qui ont un genre, les noms de famille n'en ont pas (*Julien Sorel* → NOMB.=s | GENRE=x) sauf s'ils sont précédés d'un déterminant (*le Sorel* → NOMB.=s | GENRE=m) ou qu'ils renvoient à une personnage précis (*Calvin, Cromwell* → NOMB.=s | GENRE=m) ; Pour les ville, même si un tendance générale porte au féminin, les usages sont fluctuants. On ne met un genre que lorsque le contexte immédiat permet de trancher : *la Rome éternelle* → NOMB.=s | GENRE=f, mais *il revient de Venise* → NOMB.=s | GENRE=x).

Quelques exemples pour le débutant

- *Je suis là, et vous?* → MORPH=empty ;
- *Je suis là* → PERS.=1 | NOMB.=s | GENRE=x | CAS=n ;
- *mes yeux* → PERS.=1 | NOMB.=p | GENRE=m ;
- *Venez ici!* → MODE=imp | PERS.=2 | NOMB.=p ;
- *il y a* → PERS.=3 | NOMB.=s | GENRE=x | CAS=n ;
- *on vous propose* → PERS.=3 | NOMB.=s | GENRE=x | CAS=n ;
- *c'est évident* → NOMB.=s | GENRE=n ;
- *voyant à quoi je me prépare* → NOMB.=s | GENRE=n ;
- *d'où vient cette audace* → NOMB.=x | GENRE=x | CAS=i ;
- *c'est ce que je vois* → MORPH=empty ;
- *des douceurs charmantes* → NOMB.=p | GENRE=f ;
- *en trompant Amarante* → NOMB.=x | GENRE=x ;
- *Voyez un peu Julie* → NOMB.=s | GENRE=f.

5 Entités nommées

Une entité nommée est une expression linguistique référentielle, souvent associée aux noms propres et aux descriptions définies. Nous avons décidé de nous inspirer du guide

d'annotation *Quaero*¹⁰, qui a pour particularité de proposer une version étendue de l'annotation traditionnelle¹¹.

L'annotation repose sur une combinaison de huit types (personne/pers, fonction/func, localisation/loc, production/prod, organisation/org, temps/time, montant/amount et événement/event) et de plusieurs sous-types (.ind, .phys, .date...).

Personne		Fonction			
pers.ind	pers.coll	func.ind		func.coll	
Localisation		Production			
loc.adm.town	loc.phys.geo	loc.fac	prod.art	prod.rule	prod.object
loc.adm.reg	loc.phys.hydro	loc.oro			
loc.adm.nat					
loc.adm.sup					
Organisation		Temps	Événement	Montant	
org.adm	org.ent	time.date.abs	event	amount	
		time.date.rel			

TABLE 3 – Types (en gris) et sous-types retenus pour notre annotation dans le guide *Quaero*

- Personne :
 - pers.ind pour les noms d'individus (*La Pérouse, M. de Vaugelas*);
 - pers.coll pour les groupes, comme les maisons aristocratiques (*maison de Priam, maison d'Autriche...*) ou les dynasties (*les Comnènes*).
- Fonction :
 - func.ind pour un individu désigné par sa fonction (*le roi de France*);
 - Attention : une fonction précédée de *monsieur* ou *madame* est classée comme pers.ind (*Monsieur le roi de France* ↯ func.ind);
 - func.coll pour un groupe désigné par sa fonction (*Les évêques du Mans*).
- Localisation :
 - loc.adm :
 - loc.adm.town pour la plus petite unité habitée (*Paris*);
 - loc.adm.nat pour les pays (*La France*);
 - loc.adm.reg pour tout ce qui est entre loc.adm.town et loc.adm.nat : *Le Maine, le pays de Pamphile, la Basse-Lorraine*;
 - Attention, le référentiel reste un référentiel moderne : le *Royaume de Naples* est donc une loc.adm.reg et non une loc.adm.nat, car c'est aujourd'hui une province italienne et plus un royaume indépendant. Ce choix permet de conserver une certaine interopérabilité avec des données contemporaines (et les référentiels afférents);
 - loc.adm.sup pour tout ce qui est supérieur au pays (*L'Europe, L'Asie, L'Arabie, Les Gaules*).
 - loc.phys :

10. Sophie Rosset, Cyril Grouin et Pierre Zweigenbaum, *Entités nommées structurées : guide d'annotation Quaero*, Orsay, 2011, URL : <http://www.quaero.org/media/files/bibliographie/quaero-guide-annotation-2011.pdf>.

11. C. Grouin, S. Rosset, P. Zweigenbaum, Karën Fort, Olivier Galibert et Ludovic Quintard, « Proposal for an Extension of Traditional Named Entities : From Guidelines to Evaluation, an Overview », dans *Proceedings of the 5th Linguistic Annotation Workshop*, Portland, Oregon, USA, 2011, p. 92-100, URL : <https://www.aclweb.org/anthology/W11-0411> (visité le 29/05/2021).

- loc.phys.geo pour les lieux terrestres : bois (*forêt de Bière*), montagne (*le Parnasse*), île (*Majorque*), rivage (*rives du Gange*), cap (*Cap de Comorin*), vallée (*vallée d'Aspe*), massif (*Alpes*) ;
 - Les îles sont toujours encodées comme des lieux physiques terrestres peu importe leur statut politique (*Sardaigne*), notamment pour gérer la complexité du statut des petites îles, parfois compliquées à identifier (cf. récits de voyage) ;
 - loc.phys.hydro pour les lieux aquatiques : rivière (*la Seine*), mer (*mer du sud*, *mer du Japon*), océan (*océan septentrional*), canal (*canal de la Goullette*), golfe (*golfe persique*), embouchure (*embouchure du Rhône*), détroit (*détroit de Gibraltar*), baie (*baie de Noël*)...
 - loc.fac pour les bâtiments : hôtel (*Hostel Dieu*), palais (*Palais de messire Polliando*), port (*port d'Oloné*), édifices religieux (*église Notre-Dame*, *cathédrale Saint-Julien*, *temple d'Apollon*...) ;
 - loc.oro pour les chemins (*chemin de Lorraine*), rue (*rue du Feurre*), voie (*voie de Sion*), place (*place Vendôme*).
 - Production :
 - prod.art pour les œuvres d'art (*Lancelot*, *Iliade*), textes religieux (*Le Nouveau Testament*, *Évangile*) ou philosophiques (*Histoire de Gorgias*) ;
 - prod.rule pour les noms de lois ou de codes juridiques (*coutume de Normandie*) ;
 - prod.object pour les noms d'objets, comme les bateaux (*L'Astrolabe*) ou les épées (*Durandal*).
 - Organisation :
 - org.adm joue un rôle administratif. On l'utilise pour les administrations comme les parlements (*Parlement de Toulouse*) ou les ordres (*ordre de Saint-Michel*) ;
 - Attention, org.adm est aussi utilisé pour les lieux désignés par leur forme politique : royaume (*royaume de France*), empire (*Saint-Empire*, *Empire d'Occident*), duché (*duché de Milan*), diocèse (*diocèse du Mans*) ne sont donc pas des loc.adm.reg ou des loc.adm.nat) ;
 - Une org.ent rend un service comme les bibliothèques (*Bibliothèque du Vatican*) les académies (*Académie royale*, *Académie des sciences*), les lieux d'enseignement (*La Sorbonne*, *collège de Navarre*, *université de Leyde*).
 - Temps :
 - time.date.abs pour les dates absolues (*5 Février 1731*, *carême-prenant*) ;
 - time.date.rel pour les dates relatives (*milieu du dernier siècle*, *an prochain*...).
 - Montant :
 - amount, notamment pour les prix (*vingt livres tournois*, *5 sols*).
 - Événement :
 - event comme les batailles (*bataille d'Ivry*), les traités (*traité de Budé*), les sièges (*siège de Troyes*), les conciles (*concile de Trente*), les congrès (*congrès de Verdun*)...
- Afin de compenser nos choix de lemmatisation et d'étiquetage morpho-syntaxique très minimalistes pour ce qui concerne les noms propres (uniquement les tokens commençant avec une majuscule), nous optons ici pour une approche très maximaliste, dite *fine-grained*, avec des entités « imbriquées ». Cela nous permet de gérer des cas complexes :
- *Le roi Edouard* combine un titre et un nom ;

- *Edouard III* combine un nom et un numéro de génération ;
- *Edouard d'Angleterre* combine un nom de personne et un nom de lieu, le second étant la propriété du premier ;
- *Edouard d'Angleterre* pouvant renvoyer à plusieurs personnes, nous essayons autant que possible de désambiguïser les occurrences via Wikidata.

Quatre éléments transverses permettent de compléter l'encodage en types en en sous-types :

- `name` permet de préciser les tokens renvoyant au nom (*Louis, Philippe...*) ;
- `title` permet d'encoder le titre (*sieur, duc, abbé...*) ;
- `qualifier` permet d'annoter un adjectif (*l'Indes orientale, l'Arabie heureuse, la mer atlantique, l'ancienne Colchide*)... mais aussi le numéro de génération (*Henri IV*) ou une position cardinale ;
- `kind` permet d'encoder l'hyponyme (*l'Empire de Constantinople, la mer du Japon* ;
- `unit` permet d'encoder les unités (lieue, mètres, livres, pouces, brasses...) ;
- `val` permet d'encoder les valeurs (un numéral), en lien avec une unité pour encoder un amount.

En plus de toutes ces informations, nous avons décidé d'ajouter une annotation sémantique :

- Nous utilisons les identifiants wikidata¹² pour la désambiguïstation, car non seulement les données de Wikidata sont liées aux autres principaux référentiels, mais il est aussi possible d'ajouter soi-même de nouvelles entrées à Wikidata (ce qui est particulièrement important pour les documents anciens, riches en personnages inconnus éligibles à rejoindre ce référentiel) ;
- Conformément à notre choix d'encoder d'anciens royaumes comme des régions si ces royaumes ont disparu (*royaume de Bohême* → `loc.adm.reg`) ou que des états actuels sont vassaux d'autres (*Irlande* → `loc.adm.nat`), nous recourons autant que possible aux identifiants contemporains si le nom n'a pas changé : nous utilisons donc l'identifiant wikidata Q39193 (la région de Bohême peu importe l'époque) et non Q42585 (le royaume historique de Bohême) ;
- Dans le cas où le nom a changé dans le temps (*Paris* vs *Lutèce*) nous encodons en avec des identifiants Wikidata différents s'il existe des pages différentes, partant du principe que ce sont deux entités proches, mais différentes ;
- Dans le cas où la désambiguïstation n'est pas immédiatement faisable avec l'entité seule, sans avoir recours à un contexte élargi (par ex. *le pape Jean V*), mais qu'un doute existe (par ex. *le pape Jean*), aucun identifiant n'est mis.

Le résultat ressemble à celui-ci :

12. Denny Vrandečić et Markus Krötzsch, « Wikidata : a free collaborative knowledgebase », *Communications of the ACM*, 57–10 (23 sept. 2014), p. 78–85, DOI : 10.1145/2629489.

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
Edouard	Edouard	Np	B-pers	B-pers.ind	B-comp.name	O	Q129247
III	3	Mc	I-pers	I-pers.ind	B-comp.qualifier	O	Q129247
roi	roi	Nc	I-pers	I-pers.ind	B-comp.title	O	Q129247
d'	de	S	I-pers	I-pers.ind	I-comp.title	O	Q129247
Angleterre	Angleterre	Np	I-pers	I-pers.ind	I-comp.title	B-loc.adm.nat	Q129247

TABLE 4 – Annotation *Fine-Grained*) avec *Entity linking*

Quelques remarques importantes :

- Les articles ne sont pas inclus dans les entités : *le roi de France* ;
- Une même unité de sens séparée en deux par une virgule est annoté comme deux groupes différents : ① *Louis*, ② *roi de France* ;
- Les coordinations ne sont pas prises en compte : ① *duc du Maine* et ② *d'Anjou*. Dans ce cas, *Anjou* est étiqueté *loc.adm.reg* à l'inverse de *duc du Maine* qui étiqueté *func.ind* ;
- L'analyse en composants varie en fonction de l'entité. Elle n'est, par exemple, pas la même dans le cadre d'une personne ou d'une fonction : *duc* est annoté *kind* si c'est un *func*, mais *title* si c'est un *pers* ;
- Nous n'analysons pas les emplois métonymiques. Si *la France part en guerre*, *France* est étiqueté *loc.adm.nat* et non *org.adm*, même s'il est évident que *France* ne désigne pas le pays mais son armée.

Ainsi, un exemple très proche de celui présenté dans le tab. 4 doit être encodé de la manière suivante :

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
le	le	Da	O	O	O	O	–
roi	roi	Nc	B-func	I-pers.ind	B-comp.kind	O	–
d'	de	S	I-func	I-pers.ind	O	O	–
Angleterre	Angleterre	Np	I-func	I-func.ind	B-comp.name	B-loc.adm.nat	–
et	et	Cc	O	O	O	O	–
seigneur	seigneur	Nc	B-func	I-pers.ind	B-comp.kind	O	–
d'	de	S	I-func	I-pers.ind	O	O	–
Irlande	Irlande	Np	I-func	I-func.ind	B-comp.name	B-loc.adm.nat	–
,	,	Fw	O	O	O	O	–
Edouard	Edouard	Np	B-pers	B-pers.ind	B-comp.name	O	Q129247
III	3	Mc	I-pers	I-pers.ind	B-comp.qualifier	O	Q129247

TABLE 5 – Annotation *Fine-Grained* avec *Entity linking*

Quelques exemples pour le débutant

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
roi	roi	Nc	B-func	B-func.ind	B-comp.kind	O	Q312735
d'	de	S	I-func	I-func.ind	O	O	Q312735
Arragon	Aragon	Np	I-func	I-func.ind	B-comp.name	B-loc.adm.reg	Q312735

TABLE 6 – Exemple 1

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
roi	roi	Nc	B-pers	B-pers.ind	B-comp.title	O	Q312735
d'	de	S	I-pers	I-pers.ind	I-comp.title	O	Q312735
Arragon	Aragon	Np	I-pers	I-pers.ind	I-comp.title	B-loc.adm.reg	Q312735
Jacques	Jacques	Np	I-pers	I-pers.ind	B-comp.name	O	Q312735
II	2	Np	I-pers	I-pers.ind	B-comp.qualifier	O	Q312735
dit	dire	Vvc	I-pers	I-pers.ind	O	O	Q312735
Le	Le	Np	I-pers	I-pers.ind	B-comp.name	O	Q312735
Juste	juste	Np	I-pers	I-pers.ind	I-comp.name	O	Q312735

TABLE 7 – Exemple 2

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
ouest	ouest	Nc	B-loc	B-loc.adm.reg	B-comp.qualifier	O	–
de	de	S	I-loc	I-loc.adm.reg	O	O	–
la	le	Da	I-loc	I-loc.adm.reg	O	O	–
nouvelle	nouveau	Ag	I-loc	I-loc.adm.reg	B-comp.qualifier	O	–
Irlande	Irlande	Np	I-loc	I-loc.adm.reg	B-comp.name	O	–

TABLE 8 – Exemple 3

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
Marie	Marie	Np	B-pers	B-pers.ind	O	O	Q182021
De	de	Np	I-pers	I-pers.ind	O	O	Q182021
Medicis	Médicis	Np	I-pers	I-pers.ind	O	O	Q182021

TABLE 9 – Exemple 4

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
reine	reine	Nc	B-pers	B-pers.ind	B-comp.title	O	Q182021
Marie	Marie	Np	I-pers	I-pers.ind	B-comp.name	O	Q182021
De	de	Np	I-pers	I-pers.ind	B-comp.name	O	Q182021
Medicis	Médicis	Np	I-pers	I-pers.ind	B-comp.name	O	Q182021

TABLE 10 – Exemple 5

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
château	château	Nc	B-loc	B-loc.fac	B-comp.kind	O	Q516697
de	de	S	I-loc	I-loc.fac	O	O	Q516697
Compiègne	Compiègne	Np	I-loc	I-loc.fac	B-comp.name	O	Q516697

TABLE 11 – Exemple 6

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
23	23	Mc	B-time	B-time.date.abs	O	O	–
février	février	Nc	I-time	I-time.date.abs	O	O	–

TABLE 12 – Exemple 7

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
l'	le	Da	O	O	O	O	–
an	an	Nc	B-time	B-time.date.rel	O	O	–
passé	passer	Ge	I-time	I-time.date.rel	O	O	–

TABLE 13 – Exemple 8

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
académie	académie	Nc	B-org	B-org.ent	B-comp.kind	O	Q2822388
de	de	S	I-org	I-org.ent	O	O	Q2822388
Dijon	Dijon	Np	I-org	I-org.ent	B-comp.name	B-loc.adm.town	Q2822388

TABLE 14 – Exemple 9

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
rue	rue	Nc	I-loc	I-loc.oro	B-comp.kind	O	Q24932930
des	de_le	S+Da	I-loc	I-loc.oro	O	O	Q24932930
deux	deux	Mc	I-loc	I-loc.oro	B-comp.name	O	Q24932930
écus	écu	Nc	I-loc	I-loc.oro	I-comp.name	O	Q24932930

TABLE 15 – Exemple 10

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
la	le	Da	O	O	O	O	–
coutume	coutume	Nc	B-prod	B-prod.rule	B-comp.kind	O	Q7728503
de	de	S	I-prod	I-prod.rule	O	O	Q7728503
Paris	Paris	Np	I-prod	I-prod.rule	B-comp.name	B-loc.adm.town	Q7728503

TABLE 16 – Exemple 11

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
Dictionnaire	dictionnaire	Nc	B-prod	B-prod.art	B-comp.kind	O	Q1210345
de	de	S	I-prod	I-prod.art	O	O	Q1210345
Trévoux	Trévoux	Np	I-prod	B-prod.art	B-comp.name	B-loc.adm.town	Q1210345

TABLE 17 – Exemple 12

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
nord	nord	Nc	B-loc	B-loc.adm.reg	B-comp.qualifier	O	–
de	de	S	I-loc	I-loc.adm.reg	O	O	–
Bar	bar	Nc	I-loc	I-loc.adm.reg	B-comp.name	O	–
-	-	Fo	I-loc	I-loc.adm.reg	I-comp.name	O	–
sur	sur	S	I-loc	I-loc.adm.reg	I-comp.name	O	–
-	-	Fo	I-loc	I-loc.adm.reg	I-comp.name	O	–
Aube	Aube	Np	I-loc	I-loc.adm.reg	I-comp.name	O	–

TABLE 18 – Exemple 13

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
N.	nord	Xa	B-loc	B-loc.adm.reg	B-comp.qualifier	O	–
E.	est	Xa	I-loc	I-loc.adm.reg	I-comp.qualifier	O	–
de	de	S	I-loc	I-loc.adm.reg	O	O	–
la	le	Da	I-loc	I-loc.adm.reg	O	O	–
Rochelle	Rochelle	Np	I-loc	I-loc.adm.reg	B-comp.name	O	–

TABLE 19 – Exemple 14

TOKEN	LEMMA	POS	COARSE	FINE	COMP	NESTED	WIKIDATA
évêques	évêque	Nc	B-func	B-func.coll	B-comp.kind	O	–
de	de	S	I-func	I-func.coll	O	O	–
Grenoble	Grenoble	Np	I-func	I-func.coll	B-loc	B-loc.adm.town	–
et	et	Cc	O	O	O	O	–
d'	de	S	O	O	O	O	–
Angers	Angers	Np	B-loc	B-loc.adm.town	O	O	Q38380

TABLE 20 – Exemple 15

Remerciements

Un grand nombre des cas problématiques commentés dans ce manuel proviennent de discussions avec Ariane Pinche, Maeva Nguyen, Jean-Baptiste Tanguy, Marie Puren, Frédéric Duval, Lucence Ing, Suzanne Duval, et Maxime Cario : leur aide fut donc précieuse lors de la rédaction. Une pensée toute particulière va à Florian Cafiero, frère d’armes lors de la campagne *CornMol*.

Annexes

A Exemple d'annotation

Ami	ami	NOMcom	NOMB.=s GENRE=m
,	,	PONfbl	MORPH=empty
j'	je	PROper	PERS.=1 NOMB.=s GENRE=x CAS=n
ai	avoir	VERcjk	MODE=ind TEMPS=pst PERS.=1 NOMB.=s
beau	beau	ADVgen	MORPH=empty
rêver	rêver	VERinf	MORPH=empty
,	,	PONfbl	MORPH=empty
toute	tout	DETind	NOMB.=s GENRE=f
ma	mon	DETpos	PERS.=1 NOMB.=s GENRE=f
rêverie	rêverie	NOMcom	NOMB.=s GENRE=f
Ne	ne	ADVneg	MORPH=empty
me	je	PROper	PERS.=1 NOMB.=s GENRE=x CAS=r
fait	faire	VERcjk	MODE=ind TEMPS=pst PERS.=3 NOMB.=s
rien	rien	PROind	NOMB.=s GENRE=x
comprendre	comprendre	VERinf	MORPH=empty
en	en	PRE	MORPH=empty
ta	ton	DETpos	PERS.=2 NOMB.=s GENRE=f
galanterie	galanterie	NOMcom	NOMB.=s GENRE=f
.	.	PONfbl	MORPH=empty
Auprès	auprès	ADVgen	MORPH=empty
de	de	PRE	MORPH=empty
ta	ton	DETpos	PERS.=2 NOMB.=s GENRE=f
maîtresse	maîtresse	NOMcom	NOMB.=s GENRE=f
engager	engager	VERinf	MORPH=empty
un	un	DETndf	NOMB.=s GENRE=m
ami	ami	NOMcom	NOMB.=s GENRE=m
,	,	PONfbl	MORPH=empty
C'	ce	PROdem	NOMB.=s GENRE=n
est	être	VERcjk	MODE=ind TEMPS=pst PERS.=3 NOMB.=s
,	,	PONfbl	MORPH=empty
à	à	PRE	MORPH=empty
mon	mon	DETpos	PERS.=1 NOMB.=s GENRE=m
jugement	jugement	NOMcom	NOMB.=s GENRE=m
,	,	PONfbl	MORPH=empty
ne	ne	ADVneg	MORPH=empty
l'	il	PROper	PERS.=3 NOMB.=s GENRE=x CAS=r
aimer	aimer	VERinf	MORPH=empty
qu'	que	CONsub	MORPH=empty
à	à	PRE	MORPH=empty
demi	demi	NOMcom	NOMB.=s GENRE=m
.	.	PONfbl	MORPH=empty

B Abréviations pour le tokeniseur (extrait)

<i>Acad.</i>	<i>académie</i>	<i>Mech.</i>	<i>mécanique</i>
<i>Adj.</i>	<i>adjectif</i>	<i>Med.</i>	<i>médecine</i>
<i>Agricol.</i>	<i>agricole</i>	<i>Med.</i>	<i>médical</i>
<i>Agricul.</i>	<i>agriculture</i>	<i>Mem.</i>	<i>mémoire</i>
<i>Apocal.</i>	<i>Apocalypse</i>	<i>Menuis.</i>	<i>menuiserie</i>
<i>anc.</i>	<i>ancienne</i>	<i>Milit.</i>	<i>militaire</i>
<i>Bot.</i>	<i>botanique</i>	<i>Mod.</i>	<i>moderne</i>
<i>Botan.</i>	<i>botanique</i>	<i>Mor.</i>	<i>moral</i>
<i>Botaniqu.</i>	<i>botanique</i>	<i>Mr.</i>	<i>monsieur</i>
<i>ca.</i>	<i>capitulum</i>	<i>Monsr.</i>	<i>monsieur</i>
<i>cap.</i>	<i>capitulum</i>	<i>nat.</i>	<i>naturel</i>
<i>capi.</i>	<i>capitulum</i>	<i>natur.</i>	<i>naturel</i>
<i>Cf.</i>	<i>confer</i>	<i>N.b.</i>	<i>nota bene</i>
<i>Cha.</i>	<i>chapitre</i>	<i>Orat.</i>	<i>oratoire</i>
<i>Chap.</i>	<i>chapitre</i>	<i>Ornith.</i>	<i>ornithologie</i>
<i>Col.</i>	<i>colonne</i>	<i>Ornythol.</i>	<i>ornithologie</i>
<i>Dic.</i>	<i>dictionnaire</i>	<i>Ornitholog.</i>	<i>ornithologie</i>
<i>Diction.</i>	<i>dictionnaire</i>	<i>Part.</i>	<i>partie</i>
<i>Dictionn.</i>	<i>dictionnaire</i>	<i>Pag.</i>	<i>page</i>
<i>Eccl.</i>	<i>ecclésiastique</i>	<i>Pharm.</i>	<i>pharmacie</i>
<i>Écon.</i>	<i>économie</i>	<i>Phil.</i>	<i>philosophie</i>
<i>Élem.</i>	<i>élément</i>	<i>Philos.</i>	<i>philosophie</i>
<i>Fig.</i>	<i>figure</i>	<i>Pl.</i>	<i>planche</i>
<i>Fr.</i>	<i>français(e)</i>	<i>Pl.</i>	<i>pluriel</i>
<i>Geog.</i>	<i>géographie</i>	<i>Politiq.</i>	<i>politique</i>
<i>Gram.</i>	<i>grammaire</i>	<i>P.S.</i>	<i>post scriptum</i>
<i>Gramm.</i>	<i>grammaire</i>	<i>Phys.</i>	<i>physique</i>
<i>Hist.</i>	<i>histoire</i>	<i>Physiq.</i>	<i>physique</i>
<i>Ibid.</i>	<i>ibidem</i>	<i>Sr.</i>	<i>sieur</i>
<i>Ibid.</i>	<i>ibidem</i>	<i>St.</i>	<i>saint</i>
<i>Inst.</i>	<i>institution</i>	<i>Subst.</i>	<i>substantif</i>
<i>Jard.</i>	<i>jardinage</i>	<i>s.f.</i>	<i>substantif féminin</i>
<i>Jurisprud.</i>	<i>jurisprudence</i>	<i>S.M.</i>	<i>Sa Majesté</i>
<i>Latit.</i>	<i>latitude</i>	<i>s.m.</i>	<i>substantif masculin</i>
<i>Li.</i>	<i>liber</i>	<i>Tab.</i>	<i>tableau</i>
<i>Lib.</i>	<i>liber</i>	<i>Tbât.</i>	<i>théâtre</i>
<i>Libr.</i>	<i>liber</i>	<i>Trév.</i>	<i>Trévoux</i>
<i>Lig.</i>	<i>ligne</i>	<i>Tom.</i>	<i>Tome</i>
<i>Lit.</i>	<i>littérature</i>	<i>Vól.</i>	<i>Volume</i>
<i>Littérat.</i>	<i>littérature</i>	<i>V. n.</i>	<i>Verbe neutre</i>
<i>Liv.</i>	<i>livre</i>	<i>V. a.</i>	<i>Verbe actif</i>
<i>Long.</i>	<i>longitude</i>	<i>V. act.</i>	<i>Verbe actif</i>
<i>Mar.</i>	<i>marin</i>	<i>Zoo.</i>	<i>zoologique</i>
<i>Mat.</i>	<i>mathématiques</i>	<i>Zoolog.</i>	<i>zoologique</i>
<i>Mathém.</i>	<i>mathématiques</i>		

C Jeu d'étiquette CATTEX

Type	Étiquette	Définition	
Verbes	VERcjpg	Verbe conjugué	
	VERinf	Verbe infinitif	
	VERppe	Verbe p.passé	
	VERppa	Verbe p.présent	
Noms	NOMcom	Nom commun	
	NOMpro	Nom propre	
Adjectifs	ADJqua	Adjectif qualificatif	
	ADJind	Adjectif indéfini	
	ADJpos	Adjectif possessif	
	ADJcar	Adjectif cardinal	
	ADJord	Adjectif ordinal	
Pronoms	PROper	Pronom personnel	
	PROper.PROper	Pronoms personnel composés (<i>jel/jol</i>)	
	PROimp	Pronom impersonnel	
	PROadv	Pronom adverbial	
	PROpos	Pronom possessif	
	PROdem	Pronom démonstratif	
	PROind	Pronom indéfini	
	PROcar	Pronom cardinal	
	PROord	Pronom ordinal	
	PROrel	Pronom relatif	
	PROint	Pronom interrogatif	
	PROcom	Pronom composé (<i>ledict, ladict</i>)	
	Déterminant	DETdef	Déterminant défini
		DETndf	Déterminant non défini
DETdem		Déterminant démonstratif	
DETpos		Déterminant possessif	
DETind		Déterminant indéfini	
DETcar		Déterminant cardinal	
DETrrel		Déterminant relatif	
DETint		Déterminant interrogatif	
DETcom		Déterminant composé	
Adverbes		ADVgen	Adverbe général
	ADVgen.PROper	Adverbe général + pronom personnel (<i>sil, sel</i>)	
	ADVgen.PROadv	Adverbe général + pronom adverbial (<i>sin</i>)	
	ADVneg	Adverbe de négation	
	ADVneg.PROper	Adverbe de négation + pronom personnel (<i>nel</i>)	
	ADVneg.PROadv	Adverbe de négation + pronom adverbial (<i>non = ne + en</i>)	
	ADVint	Adverbe interrogatif	
	ADVing	Adverbe interrogatif négatif	
Prépositions	ADVsub	Adverbe subordonnant	
	PRE	Préposition (<i>sauf, par, de, en, por, sans</i>)	
	PRE.DETdef	Enclise du déterminant défini après préposition	
	PRE.DETcom	Enclise du déterminant composé après préposition	
	PRE.DETrrel	Enclise du déterminant relatif (ou interrogatif en interrogative indirecte) après préposition	

Type	Étiquette	Définition
Conjonctions	PRE.PROper	Enclise du pronom personnel après préposition
	PRE.PROrel	Enclise du pronom relatif (ou interrogatif en interrogative indirecte) après préposition
	CONcoo	Conjonction de coordination
	CONsub	Conjonction de subordination
	CONsub.PROper	Enclise du pronom personnel après conjonction de subordination
Interjections	INT	Interjection
Ponctuations	PON	Ponctuation
	PONfrt	Ponctuation forte (délimite les phrases)
	PONfbl	Ponctuation faible (interne à une phrase)
	PONpga	Guillemet ou parenthèse ouvrants
	PONpdr	Guillemet ou parenthèse fermant
	PONpxx	Guillemet droit (quand on ne sait pas si c'est ouvrant ou fermant)
Redondance	RED	"que" redondant
OUT	OUT	Ce qui ne doit pas être pris en compte dans l'analyse linguistique

Références

- ATILF-CNRS et UNIVERSITÉ DE LORRAINE, *Dictionnaire du Moyen Français (1330-1500)*, 2015, URL : <http://www.atilf.fr/dmf>.
- Jean-Baptiste Camps (éd.), *Geste : un corpus de chansons de geste, 2016-... (Version 02)*, Paris, 2016-2020, URL : <http://doi.org/10.5281/zenodo.2630574>.
- CLÉRICE (Thibault), *Référentiel du Latin pour Pyrrha, d'après le dictionnaire et travaux du LASLA de D. Longrée et al*, mai 2020, DOI : 10.5281/zenodo.3822040.
- DIWERSY (Sascha), FALAISE (Achille), LAY (Marie-Hélène) et SOUVAY (Gilles), « Ressources et méthodes pour l'analyse diachronique », *Langages*, N° 206-2 (août 2017), p. 21-44, URL : <https://www.cairn.info/revue-langages-2017-2-page-21.htm> (visité le 03/12/2018).
- GROUIN (Cyril), ROSSET (Sophie), ZWEIGENBAUM (Pierre), FORT (Karèn), GALIBERT (Olivier) et QUINTARD (Ludovic), « Proposal for an Extension of Traditional Named Entities : From Guidelines to Evaluation, an Overview », dans *Proceedings of the 5th Linguistic Annotation Workshop*, Portland, Oregon, USA, 2011, p. 92-100, URL : <https://www.aclweb.org/anthology/W11-0411> (visité le 29/05/2021).
- GUILLOT (Céline), HEIDEN (Serge) et LAVRENTIEV (Alexei), « Base de français médiéval : une base de référence de sources médiévales ouverte et libre au service de la communauté scientifique », *Diachroniques : Revue de Linguistique française diachronique-7* (déc. 2017), p. 168-184, URL : <https://halshs.archives-ouvertes.fr/halshs-01809581> (visité le 19/12/2019).
- GUILLOT (Céline), PRÉVOST (Sophie) et LAVRENTIEV (Alexei), *Principes d'annotation Cattex09*, rapp. tech., version 2.0. http://bfm.ens-lyon.fr/IMG/pdf/Cattex2009_principes_2.0.pdf, Lyon, École normale supérieure de Lyon, 2013.
- PIERREL (Jean-Marie), DENDIEN (Jacques) et BERNARD (Pascale), « Le TLFi ou Trésor de la Langue Française informatisé », dans *Proceedings of the 11th EURALEX International Congress*, dir. Geoffrey Williams et Sandra Vessier, Lorient, France, 2004, p. 165-170.
- PRÉVOST (Sophie), GUILLOT (Céline), LAVRENTIEV (Alexei) et HEIDEN (Serge), *Jeu d'étiquettes morphosyntaxiques CATTEX2009*, rapp. tech., version 2.0. http://bfm.ens-lyon.fr/IMG/pdf/Cattex2009_2.0.pdf, Lyon, École normale supérieure de Lyon, 2013.
- ROSSET (Sophie), GROUIN (Cyril) et ZWEIGENBAUM (Pierre), *Entités nommées structurées : guide d'annotation Quaero*, Orsay, 2011, URL : <http://www.quaero.org/media/files/bibliographie/quaero-guide-annotation-2011.pdf>.
- SOUVAY (Gilles) et PIERREL (Jean-Marie), « LGeRM Lemmatisation des mots en Moyen Français », *Traitement Automatique des Langues*, 50-2 (2009), p. 149-172, URL : <https://halshs.archives-ouvertes.fr/halshs-00396452>.
- VRANDEČIĆ (Denny) et KRÖTZSCH (Markus), « Wikidata : a free collaborative knowledge-base », *Communications of the ACM*, 57-10 (23 sept. 2014), p. 78-85, DOI : 10.1145/2629489.