



**HAL**  
open science

## Une grammaire fondée sur un corpus numérique

Sophie Prévost

► **To cite this version:**

Sophie Prévost. Une grammaire fondée sur un corpus numérique. Christiane Marchello-Nizia; Bernard Combettes; Sophie Prévost; Tobias Scheer. Grande Grammaire Historique du Français, De Gruyter Mouton, pp.37-53, 2020, 978-3-11-034553-7. halshs-03095100

**HAL Id: halshs-03095100**

**<https://shs.hal.science/halshs-03095100>**

Submitted on 9 Nov 2023

**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.

## Chapitre 3

### Une grammaire fondée sur un corpus numérique

#### 3.1 Une histoire de la langue française

Une langue, même appréhendée dans une perspective synchronique, n'est pas un objet simple : elle revêt des formes diverses en fonction de différents paramètres, tels que le registre, la situation de communication, mais aussi la forme des textes, leur domaine, et, le cas échéant, le dialecte. Parler de LA langue française est donc une commodité de langage : sous cette appellation se dissimulent des variétés, des usages, tous un peu différents, mais néanmoins suffisamment similaires pour que l'on puisse les considérer comme autant de mises en œuvre d'une même langue, le français.

Une grammaire peut vouloir dégager les points communs entre les différents usages d'une langue, mettre au jour ce qui les réunit en une sorte de *koinè*, et laisser à la marge ce qui les distingue les uns des autres ; elle peut même, dans une perspective plus prescriptive que descriptive, et comme cela fut le cas pendant des siècles, ne s'intéresser qu'à certains usages, représentatifs de la norme, et tout simplement ignorer les autres ou ne les mentionner que comme étant déviants. Elle peut aussi, à l'inverse, prendre en compte la diversité des usages.

L'émergence et le développement de la linguistique, en particulier de la sociolinguistique, ont largement modifié la conception et l'approche de la langue. Beaucoup en sont désormais convaincus : faire la grammaire d'une langue, c'est tout autant dire expliciter ce qui devrait être que constater ce qui est réellement, dans toute sa variété, et tenter de l'expliquer. Il reste cependant difficile de déterminer – c'est-à-dire identifier de manière contrastive – les classes grammaticales à recenser et à représenter : sur ce point, notre connaissance, même si elle progresse sans cesse, reste encore incomplète, aussi bien pour les états de langue modernes que pour les états plus anciens (et ce pour des raisons différentes : exploration inachevée d'un ensemble fini de documents, dont certains sont encore à découvrir, et disparition de certains autres pour la langue ancienne, prolifération continue de documents pour les états de langue modernes).

Une grammaire qui a pour objectif de restituer l'évolution de la langue française doit, plus encore qu'une grammaire synchronique, nécessairement prendre en compte la variété de la langue. Il est en effet désormais admis que la plupart des changements résultent d'un fait de variation préalable (► 2.4) : une forme ou une construction nouvelle, un paradigme nouveau, correspondent le plus souvent à une variante qui s'est imposée. On ne saurait donc considérer les seules normes (qui n'ont en outre pas été édictées avant le 16<sup>e</sup> s.), il convient au contraire d'envisager aussi les usages déviants. Cela signifie que, à chaque étape, se trouve répétée la complexe identification de l'ensemble des usages, afin de redéfinir l'objet d'étude.

Une grammaire qui adopte une perspective diachronique couvrant plusieurs siècles est par ailleurs confrontée à la complexité du rapport à la langue. Celui-ci est en effet fort différent selon que l'on considère les états de langue anciens, pour lesquels il n'existe plus de locuteurs, ou au contraire le français moderne, pour lequel nous bénéficions de notre compétence de locuteur. L'absence de locuteur n'a certes pas le même impact selon la période considérée. En effet, alors qu'un locuteur moderne ne peut comprendre un texte du 13<sup>e</sup> s. s'il n'a pas *appris* l'ancien français, il peut lire et comprendre un texte du 18<sup>e</sup> s. même s'il ne

possède pas à proprement parler la compétence de la langue de cette époque, et n'en maîtrise donc pas les nuances (d'où de probables contresens, souvent légers). Se pose ici la question de déterminer la frontière temporelle à partir de laquelle un locuteur moderne ne peut plus comprendre un état de langue passé, et a besoin de l'apprendre comme une langue étrangère. La perte de compréhension est progressive : alors qu'on lit parfaitement Zola, très bien Musset, un peu moins bien La Bruyère et Marguerite de Navarre, la compréhension des *Cent nouvelles nouvelles* (15<sup>e</sup> s.) devient difficile, et celle de la *Chanson de Roland* (début 12<sup>e</sup> s.) impossible pour les non-initiés. Ainsi notre sentiment de familiarité avec la langue décline au fur et à mesure que l'on recule dans le temps, mais, quelle que soit l'époque considérée, l'absence de compétence pour l'état de langue qu'on se propose de décrire oblige à se fonder sur des données attestées : le recours aux textes, au sens large de ce terme, s'impose tout autant pour décrire la langue du 18<sup>e</sup> s. que celle du 13<sup>e</sup> s.

Il est possible de rédiger une grammaire du français moderne en se fondant principalement sur sa compétence de locuteur enrichie de l'expérience que l'on a des différents usages et de leur confrontation, et en se passant donc de l'appui sur des données attestées (même si l'on perd ainsi de précieux éléments de diversité). Une telle démarche est exclue pour les états de langue anciens puisque seuls les textes nous fournissent les données langagières à décrire ; en effet, les *Manières de langage*, qui apparaissent vers le 14<sup>e</sup> s., et les grammaires, qui se développent à partir du 16<sup>e</sup> s., contiennent des éléments d'information mais il est difficile d'évaluer dans quelle mesure ces ouvrages reflètent vraiment l'usage, la probabilité étant grande que, comme encore souvent en français moderne, leur valeur soit plus prescriptive que descriptive. Le corpus utilisé peut ne pas apparaître de manière explicite, ou ne pas présenter un caractère raisonné (en particulier lorsque la description et / ou l'analyse des faits langagiers s'est appuyée sur des études antérieures et que les sources textuelles de celles-ci ne sont pas mentionnées, ou bien quand les textes fonctionnent comme un simple réservoir à exemples) : les données attestées, et donc les textes, ont néanmoins toujours joué un rôle essentiel dans la découverte et l'appréhension des états de langue passés, ainsi que dans la compréhension de leur *grammaire*.

En revanche, ce qui est relativement récent, c'est l'importance accordée à la quantification des données, et à leur caractérisation selon des critères fixes, la fréquence plus ou moins élevée d'une construction ne pouvant en effet s'apprécier de manière intuitive, même d'ailleurs pour sa propre langue. C'est une démarche décisive pour pleinement rendre compte de l'émergence et de la propagation, ou au contraire du recul, d'une construction, ainsi que des faits de concurrence entre constructions, c'est-à-dire, d'une manière générale, de la variation et des changements qui en résultent.

### 3.2 Les grammaires et les corpus

Dans les ouvrages de grammaire sur la langue ancienne édités depuis un siècle, les textes qui sont utilisés sont toujours cités, ce qui permet d'observer des écarts assez importants quant à leur nombre. Cela est en partie dû, à l'époque moderne, à des choix personnels de la part des auteurs de ces grammaires, mais aussi et surtout aux possibilités nouvelles offertes par la numérisation des textes et à l'automatisation de certaines tâches de recherche et d'extraction.

L'évolution est nette lorsque l'on considère quelques ouvrages de référence, en s'en tenant par exemple à la période médiévale. Ainsi, pour l'ancien français, la *Petite syntaxe de*

*l'ancien français* de L. Foulet s'appuie de manière explicite sur 13 textes (essentiellement du 13<sup>e</sup> s., et quelques-uns de la seconde moitié du 12<sup>e</sup> s.) pour sa première édition en 1919 (19 textes pour la 3<sup>e</sup> édition en 1930), qui ont permis dès cette époque – et c'était alors tout à fait inédit – de calculer la fréquence de certains faits linguistiques. La *Grammaire de l'ancien français* de G. Moignet (<sup>2</sup>1984 [1973]) se fonde quant à elle sur 102 textes, tandis que la *Grammaire nouvelle de l'ancien français* de C. Buridant (2000a) mentionne plus de 350 textes auxquels s'ajoutent des index et des concordances électroniques, ainsi que des fréquences issues de nombreuses études portant sur des thèmes précis. Pour le moyen français, la *Syntaxe du moyen français* (1980) de R. Martin et M. Wilmet, qui couvre une brève synchronie (1455-1465), repose sur 24 textes (dont 3 numérisés et analysés exhaustivement) ; *La langue française aux XIV<sup>e</sup> et XV<sup>e</sup> siècles* de C. Marchello-Nizia (<sup>2</sup>1997a [1979]), qui se place dans une perspective partiellement diachronique (14<sup>e</sup>-15<sup>e</sup> s.), mentionne de son côté 69 textes, qui constituent deux corpus, dont l'un a fait l'objet de dénombrements systématiques. Pour ce qui est des grammaires portant sur les périodes ultérieures, la *Grammaire de la langue française du seizième siècle*, de G. Gougenheim (<sup>2</sup>1974 [1951]) mentionne 33 textes, tandis que la *Grammaire du Français classique*, de N. Fournier, écrite quelques décennies plus tard (1998), en cite plus de 300.

Si le nombre de textes retenus est variable (on aimerait d'ailleurs aussi connaître leur nombre de mots), l'utilisation qui en est faite l'est tout autant. En effet les textes peuvent avoir subi un dépouillement exhaustif, ou avoir été l'objet de sondages réguliers. Ils peuvent au contraire n'être que des réservoirs à exemples (cela semble être les cas dans les grammaires de G. Gougenheim et de N. Fournier, qui dressent la liste des textes « cités »). De plus, les textes utilisés n'ont pas nécessairement été tous soumis au même traitement. Ces points ne sont guère explicités (C. Marchello-Nizia (<sup>2</sup>1997a) distingue néanmoins les « textes dépouillés ou consultés » et ceux « auxquels il est fait occasionnellement référence »). On peut le regretter, car ces choix ne sont pas sans conséquence sur les descriptions proposées.

Il est un autre aspect qui varie assez fortement d'une grammaire à l'autre : il s'agit de la diversité des textes retenus. Ainsi, chez Foulet, les textes sont littéraires, alors que la grammaire de Marchello-Nizia s'appuie sur des genres plus variés, prenant par exemple en compte des textes historiques. Cette caractéristique qualitative est elle aussi susceptible de peser assez lourdement sur les analyses proposées (et cela d'autant plus que les textes ne sont pas de simples réservoirs à exemples mais étaient véritablement l'analyse).

A ce jour, aucune grammaire sur le français moderne ne s'appuie sur un corpus explicite et raisonné, qui serait exploité de manière systématique (ou non) pour l'ensemble des faits décrits et analysés : les données attestées qui sont mentionnées ne viennent qu'exemplifier la description et l'analyse. De ce point de vue, la grammaire publiée il y a vingt ans pour l'anglais contemporain, *The Longman Grammar of Spoken and Written English* de D. Biber et al. (<sup>4</sup>2004 [1999]), constitue un ouvrage remarquable. Trois aspects en particulier méritent d'être signalés : la prise en compte systématique de différents registres, écrits et oraux, l'établissement de données quantifiées, et enfin l'exploitation à géométrie variable des corpus (40 millions de mots), en fonction des points à étudier. La démarche à l'œuvre traduit un traitement raisonné et explicite des textes, qui a été permis par une vaste entreprise de numérisation et de catégorisation des textes.

Il n'était pas possible, dans le cadre du présent projet, d'envisager une entreprise de constitution de corpus analogue à celle réalisée pour la grammaire de Biber et al., et ce d'autant que notre projet ne s'en tient pas à trois décennies, mais couvre une période de 12 siècles. C'est donc, pour une large part, sur les ressources existantes que nous nous

sommes appuyés, à savoir des textes déjà numérisés, et enrichis morpho-syntaxiquement (parfois syntaxiquement) pour certains d'entre eux.

### 3.3 La représentativité du corpus

La difficulté de délimiter l'objet qu'une grammaire se donne à décrire trouve un écho direct dans la constitution du corpus, qui, idéalement, devrait regrouper toutes les variétés, tous les usages recensés dans la langue. Trouver des textes, en nombre suffisant, et les représentant tous, pose un problème particulièrement aigu pour les états de langue anciens. En effet, certains usages ne nous sont parvenus qu'à travers de rares textes, d'autres ne nous sont connus qu'indirectement, par des témoignages. Il est probable aussi que nous ignorons tout simplement certains usages, dont nulle trace n'est arrivée jusque nous (cela tient parfois aussi aux éditeurs, dont certains ont tendance à « corriger » comme déviants des termes ou des constructions non déjà recensés). C'est le cas en particulier des usages oraux, dont nous n'avons, pour la plupart des périodes, aucune trace directe : nous ne pouvons nous fonder que sur d'éventuels témoignages ou sur des représentations écrites. Or il est possible que ces dernières tirent les usages réels vers la caricature ou au contraire vers une normalisation sur le modèle de l'écrit. Ainsi, en adoptant une démarche sur corpus, on explicite une difficulté qui serait sinon éludée : la représentativité du corpus, et du même coup l'aptitude de la *grammaire* à être généralisable au-delà des seuls textes sur lesquels elle se fonde.

La représentativité du corpus est à envisager d'un double point de vue, quantitatif et qualitatif. Sur le plan quantitatif, il faut décider si l'on travaille sur des textes intégraux ou sur des échantillons, ou bien en combinant les deux, selon la taille des textes. On peut ainsi décider d'échantillonner les textes lorsqu'ils excèdent un certain nombre de mots. Il convient dans ce cas de fixer un seuil ainsi que les modalités de l'échantillonnage (un seul échantillon ou au contraire plusieurs échantillons répartis dans le texte). Sur le plan qualitatif, les critères qui s'avèrent décisifs dans un projet de grammaire historique, sont les suivants : la date, le domaine ou le genre textuel, le dialecte, la forme du texte. Tous ne sont pas pareillement importants tout au long de la période envisagée, comme on le verra plus bas en présentant les choix qui ont été faits. D'autres sont importants aussi, tels que le registre et la classe sociale, l'âge de l'auteur, la distinction entre région de l'auteur et région du copiste... Mais l'impossibilité d'accéder pour les périodes les plus reculées à, par exemple, des écrits familiers ou émanant des classes sociales peu lettrées ne permet pas d'en faire des critères de sélection systématiques. Par ailleurs, l'ampleur du projet, qui porte sur 12 siècles, nous a obligés à restreindre l'ensemble des critères qui auraient pu être envisagés, pour des raisons pratiques de faisabilité. En effet, à la difficulté de déterminer (en termes de dates, genres, dialectes...) les textes qu'il faudrait idéalement verser au corpus, s'ajoute celle d'y accéder, et en nombre suffisant.

L'objectif était d'obtenir un corpus aussi représentatif et équilibré que possible de l'objet *langue française*, dans toute la diversité qu'on lui présuppose. Plus un corpus est jugé représentatif, plus il est légitime de généraliser les résultats obtenus au-delà des seuls textes qui le constituent. Mais il convenait par ailleurs de constituer un corpus qui reste maniable, non seulement du point de vue de l'exploration des textes (certains faits sont plus faciles à repérer automatiquement que d'autres), mais aussi du traitement des données extraites. Selon les phénomènes étudiés, on a ainsi affaire à quelques occurrences (faits peu fréquents, hapax), ou à quelques centaines, ou bien encore à des dizaines de milliers :

étudier *quant à, désormais*, ou l'ordre des mots, les déterminants, n'a évidemment pas les mêmes implications de ce point de vue.

Il a donc fallu trouver un compromis acceptable entre le corpus idéal (que l'on sait inaccessible, mais dont il faut se rapprocher autant que possible), le corpus souhaité, et le corpus possible et raisonnable. La constitution d'un corpus à géométrie variable a permis de résoudre en partie les difficultés liées à la variation des modalités d'exploration des corpus et de traitement des résultats.

## 3.4 Nos choix, notre démarche

### 3.4.1 Un corpus à géométrie variable

La *Grande Grammaire Historique du Français* (GGHF) couvre plus de 12 siècles : le rapport du locuteur moderne aux états langagiers successifs n'est pas le même, de même que varient le rôle des textes dans notre accès à la langue, la disponibilité des données textuelles, le rapport des genres entre eux, etc. L'un des défis a consisté, pour élaborer le corpus, à dépasser cette hétérogénéité, et à adopter une démarche aussi homogène que possible à travers les siècles.

La GGHF s'est donné un corpus à géométrie variable, tant du point de vue de sa constitution que de son utilisation. En effet, pour chaque période (voir ci-dessous 3.4.2.2 a. pour la délimitation des périodes), un double corpus a été élaboré : un corpus « noyau » et un corpus « complémentaire ». Le premier répond à des critères de composition stricts quant à la taille des textes et quant à leur diversité.

Pour ce qui est de la taille, nous avons fait le choix de retenir les textes dans leur intégralité lorsqu'ils n'excèdent pas 45 000 « occurrences » (mots et ponctuation, soit un peu plus de 40 000 mots). Pour les textes dépassant ce seuil, nous avons sélectionné trois échantillons d'environ 15 000 occurrences en début, milieu, et fin de texte. Toutefois, pour certains textes, jugés répétitifs du point de vue de leurs structures morphosyntaxiques, la taille de l'échantillon a été réduite à 20 000 mots. C'est le cas, par exemple, du *Registre criminel du Chatelet*, au 14<sup>e</sup> s. Pour chaque période le corpus noyau comprend entre 200 000 et 245 000 mots, hormis pour la période la plus ancienne, avant 1100, pour laquelle la quasi-totalité des textes disponibles a été retenue, l'ensemble ne dépassant pas 10 000 mots. C'est dans le corpus noyau, dont certains textes bénéficient d'un étiquetage morpho-syntaxique, qu'ont prioritairement été effectués les calculs de fréquence. Le corpus noyau échantillonné contient 205 5891 mots (et le corpus noyau non échantillonné 9 millions de mots).

Le corpus complémentaire a été conçu plus particulièrement pour l'étude des faits peu fréquents, susceptibles donc d'être peu représentés dans le corpus noyau, pour confirmer, ou non, une hypothèse développée à partir du corpus noyau, et il a de plus fourni un vaste réservoir d'exemples, permettant ainsi de diversifier les sources citées. A l'image des objectifs qui lui ont été fixés, sa constitution n'a pas été soumise aux mêmes contraintes que celle du corpus noyau : la taille des textes n'a pas été limitée, et les autres critères ont été appliqués avec une rigueur moindre. Le corpus complémentaire contient 4 571 477 mots. Le corpus intégral (corpus noyau non échantillonné et corpus complémentaire) comprend donc 13,5 millions de mots.

Par ailleurs pour certains chapitres et / ou pour certaines périodes, les contributeurs ont parfois fait usage, en plus du corpus noyau, ou même du corpus complémentaire, de corpus

spécifiques. C'est en particulier le cas pour la partie 4, qui porte sur les codes de l'écrit, et pour laquelle il s'est avéré nécessaire de recourir à des textes plus variés ou à des éditions plus sûres, voire aux manuscrits eux-mêmes. C'est aussi le cas, dans une certaine mesure, pour l'analyse des changements à partir du 17<sup>e</sup> s., des écrits relevant de registres plus familiers étant accessibles à partir de cette période.

Le corpus joue un rôle décisif dans la GGHF, et c'est là l'une des innovations de cet ouvrage. Il ne saurait cependant être question de faire table rase des études qui ont précédé : de nombreux phénomènes linguistiques ont déjà été bien décrits, et il ne s'agit pas de tout réécrire. Nous avons donc exploité plusieurs études, ainsi que les données quantifiées qui les accompagnent, le cas échéant. Ces dernières ont parfois été complétées par de nouveaux relevés, opérés dans notre corpus. Les études inédites, ou partiellement inédites, se sont beaucoup plus largement appuyées sur l'exploitation du corpus. Les modalités de ces relevés ont pu varier, de même que le traitement qui a été fait des données collectées. En effet, comme cela a déjà été souligné, la complexité à collecter des constructions, de quelque nature qu'elles soient, varie fortement selon leur caractère plus ou moins abstrait et selon le degré d'enrichissement morphosyntaxique des textes : il est plus aisé d'établir la fréquence des adverbes en *-ment* que celle des sujets nominaux. Les phénomènes qui relèvent de l'énonciation ou de la textualité, en particulier, se prêtent bien plus difficilement à une quantification des faits concernés. Par ailleurs, les modalités de traitement des données ont pu varier, selon leur nombre, leur analyse exhaustive n'étant pas possible au-delà de certains seuils.

### 3.4.2 Les critères de sélection des textes

Différents critères ont été retenus pour la sélection des textes. Certains, les *descripteurs*, ont pour but de caractériser le contenu des textes, sous différents aspects. D'autres, d'ordre en quelque sorte « paratextuels », relèvent davantage du point de vue que le locuteur moderne porte sur ces textes.

#### 3.4.2.1 Les critères paratextuels

Il nous a ainsi paru important que le corpus de la GGHF comprenne, pour chaque période, quelques textes de référence, à côté de textes moins connus (et souvent aussi – car les deux sont de fait liés – moins littéraires). La notion de texte de référence peut certes varier, et elle est en partie subjective, mais l'on peut cependant identifier comme tel quelques oeuvres, en particulier pour la période médiévale. Il n'était ainsi pas concevable que le corpus du 12<sup>e</sup> s., par exemple, ne contienne pas la *Chanson de Roland* et un roman de Chrétien de Troyes ; pour le 13<sup>e</sup> s., la *Queste du Graal* et le *Roman de la Rose* se sont d'emblée imposés, bien qu'il ne s'agisse pas des seuls textes de référence pour les périodes concernées.

La sélection s'est révélée plus difficile au fur et à mesure que l'on avance dans le temps et que se multiplie la production écrite. Le choix a nécessairement été partial, mais néanmoins influencé par la prise en compte d'un autre paramètre : la qualité des éditions. Pour les textes les plus anciens, cette exigence nous a conduits à privilégier les éditions les moins interventionnistes, et pour la période suivante (du 16<sup>e</sup> au 19<sup>e</sup> s.), des textes non (ou très peu) modernisés (l'examen des graphies est un bon indice).

Il est enfin un critère pratique qui est intervenu dans nos choix, conjointement à ceux précédemment mentionnés et aux descripteurs qui vont être évoqués ci-après. Il s'agit de

l'existence d'une version numérisée (disponible) des textes, au moins pour ceux qui appartiennent au corpus noyau et qui ont fait l'objet de quantifications. Pour la période médiévale, nous nous sommes très largement appuyés sur les textes de la *Base de Français Médiéval* (BFM), dont certains sont enrichis linguistiquement (étiquette morpho-syntaxique, et syntaxique pour certains) et deux textes proviennent du corpus *Modéliser le Changement : les Voies du Français* (<http://www.voies.uottawa.ca/index.html>). Pour la période suivante, nous avons majoritairement sélectionné les textes dans la base Frantext (<http://www.frantext.fr/>), mais aussi dans la base Epistemon (<http://www.bvh.univ-tours.fr/Epistemon/index.asp>).

Les textes retenus l'ont été aussi, et prioritairement, parce qu'ils contribuaient à construire le corpus diversifié et représentatif que nous souhaitions, au regard des critères qui nous semblaient les plus pertinents, et qui sont présentés ci-dessous.

### 3.4.2.2 Les descripteurs

#### a. La date des textes

La GGHF se distingue d'autres ouvrages diachroniques, en particulier l'*Histoire de la Langue française* de F. Brunot, en ce qu'elle est organisée, en premier lieu, non par siècle ou par grande période mais par grands domaines de la langue (phonétique, morphologie, sémantique, ...) : c'est au sein de chacune des questions abordées qu'intervient la perspective chronologique. Chaque phénomène a sa propre temporalité : il n'est donc pas possible d'établir un découpage chronologique adapté à l'évolution de l'ensemble des phénomènes. Nous avons déterminé, pour le corpus, un cadre chronologique très général, en délimitant des périodes de manière arbitraire, suivant pour cela un simple découpage par siècles, cette division ne correspondant en aucun cas à une quelconque présupposition quant à la périodisation des évolutions individuelles.

Le corpus a donc été organisé par siècles, en sélectionnant des textes qui s'échelonnent du début à la fin de chaque siècle. La période qui précède le 12<sup>e</sup> s. fait exception : en raison du petit nombre de documents qui nous sont parvenus, et de leur brièveté, les quelques témoins dont nous disposons ont été regroupés ensemble. Pour eux, la mise en oeuvre des autres critères n'est donc pas pertinente : la *Séquence de Sainte Eulalie* a été retenue non pas parce que c'est un texte en vers qui relève du domaine religieux, mais simplement parce que c'est, avec les *Serments de Strasbourg*, le seul texte en français du 9<sup>e</sup> s.

Pour chaque siècle envisagé, les textes ont été choisis en fonction de trois critères – forme, domaine et genre, dialecte – en faisant en sorte que l'ensemble composé soit diversifié. On ne peut cependant éviter, pour les périodes reculées, un certain *parasitage* entre les critères, en raison du nombre trop peu élevé de documents, ou simplement de leur absence. Ainsi, jusqu'à la fin du 12<sup>e</sup> s., les textes qui nous sont parvenus sont très majoritairement en vers, et le dialecte anglo-normand ou normand est particulièrement représenté jusqu'au milieu du siècle ; de même les textes qui relèvent du domaine historique ne se rencontrent guère avant le 13<sup>e</sup> s.

#### b. La forme des textes : vers / prose

La distinction entre textes en vers et textes en prose recouvre des réalités différentes selon les périodes considérées. Jusqu'au 12<sup>e</sup> s., la grande majorité des textes s'est écrite en vers (décasyllabes, puis octosyllabes), qu'il s'agisse de récits épiques, de « romans », de récits hagiographiques... Au fil des siècles l'écriture versifiée va reculer, conjointement au développement de la prose à partir du 13<sup>e</sup> s., pour finalement se voir réservée aux textes de poésie et de théâtre,

ainsi qu'aux chansons (types de textes qui peuvent aussi, surtout depuis le 20<sup>e</sup> s., être écrits en prose). La place respective faite à la prose et au vers n'est donc pas la même dans le corpus selon les siècles considérés : les textes en vers sont très largement majoritaires jusqu'au 12<sup>e</sup> s., puis ils cèdent une place croissante aux textes en prose, pour n'être plus associés, à partir du 18<sup>e</sup> s., qu'à certains genres : le théâtre et la poésie (chansons incluses).

#### c. *Les dialectes*

Le critère dialectal occupe une position à part parmi les critères retenus. Tout d'abord, il n'est véritablement pertinent que jusqu'au 15<sup>e</sup> s. environ, et déjà bien moins discriminant à cette époque qu'au 12<sup>e</sup> s. Par ailleurs il n'est pas toujours facile de définir le dialecte d'un texte, et il est fréquent d'opter pour un dialecte « non défini ». Enfin, et cela résulte en grande partie de la remarque précédente, ce n'est pas sur la base de leur dialecte que nous avons prioritairement sélectionné les textes. Il se trouve néanmoins que les textes retenus présentent, pour les périodes où cette distinction est pertinente, une relative diversité. Sont en particulier bien représentés l'Anglo-normand et le Picard, dont on sait qu'ils présentent plusieurs traits linguistiques spécifiques.

#### d. *Les domaines et les genres*

Nous nous sommes appuyés, pour déterminer de grands domaines, sur la classification qui a été proposée pour les textes d'ancien et de moyen français par l'équipe de la *Base de Français Médiéval*. Le *domaine* est défini comme un trait fonctionnel qui correspond à la destination principale du texte et au domaine d'activité auquel il se rattache. Dans cette perspective, les domaines retenus sont les suivants :

- littéraire : divertir
- didactico-scientifique : enseigner, instruire
- religieux : édifier (concerne le rituel et la diffusion du message chrétien)
- historique : consigner / relater les événements du passé
- juridique : réguler la vie sociale

Deux autres domaines ont été ajoutés pour les textes à partir du moyen français : il s'agit des domaines épistolaire et argumentatif, qui ne trouvent pas d'instanciation, dans les textes qui nous sont accessibles, avant le 14<sup>e</sup> s.

Contrairement à la liste des domaines, celle des *genres* est ouverte : roman, nouvelle, mémoire, chronique, lapidaire, traité, hagiographie, miracle, lyrique ... Cela tient principalement au fait que les genres ne sont pas nécessairement les mêmes d'une période à l'autre : certains apparaissent (« mémoires »), d'autres disparaissent (chanson de geste). De plus, les genres n'ont pas tous le même statut : certains sont en effet emblématiques d'une époque (les mémoires et les chroniques au 15<sup>e</sup> s., les nouvelles au 16<sup>e</sup>...), tandis que d'autres traversent les siècles, mais en connaissant des transformations radicales. Ainsi, la dénomination « roman » recouvre des réalités bien différentes au 13<sup>e</sup> s. et au 20<sup>e</sup> s.

### 3.4.3 La représentation de l'oral

Nous n'avons pas accès à la réalité matérielle, prosodique, de la langue orale avant le début du 20<sup>e</sup> s., mais nous pouvons faire l'hypothèse que, hier comme aujourd'hui, ses différents registres et genres (car, comme la langue écrite, la langue orale n'est pas homogène) présentent des spécificités qui les distinguent de ceux de la langue écrite (► chap. 37).

Nous n'avons pas non plus tenté de constituer, pour la période moderne, un corpus de données orales. L'oral n'en est pas moins présent dans la GGHF, indirectement. Il l'est tout d'abord à travers l'exploitation que nous avons faite des travaux portant sur l'oral contemporain. Par ailleurs, certains phénomènes et leur évolution ont été étudiés en observant de manière spécifique leur actualisation dans le discours direct (et en comparant la langue de ces épisodes en discours direct avec celle du récit dans lequel ils s'insèrent), et pour certains plus largement dans ce que C. Marchello Nizia a proposé d'appeler l'*oral représenté*. Il s'agit d'une partie des discours directs dans un récit (roman, chanson de geste, chronique, etc.), linguistiquement balisés (annonce, incise..), ce qui permet une comparaison entre la langue du récit enchâssant et celle du discours direct enchâssé.

Sans prétendre rendre compte de ce qu'a pu être réellement le français oral dans les siècles passés, nous avons fait l'hypothèse que, comme aujourd'hui, les réalisations écrites et orales de la langue avaient dû différer, et que les secondes avaient pu être pionnières de certains changements. De fait, cette procédure de comparaison permet de voir que de nombreuses innovations sont apparues d'abord dans les épisodes en discours direct enchâssé, c'est-à-dire en « oral représenté ».

Le corpus sur lequel s'est appuyé la GGHF n'est évidemment pas parfait. Il constitue cependant un compromis raisonnable entre l'exigence de représentativité des données et les contraintes liées à la fois à leur accessibilité et à leur traitement.

On le sait : nous n'accéderons jamais à la représentation des états passés du français dans toute leur diversité (si tant est qu'on puisse prétendre y accéder pour le français moderne). Certains aspects nous en resteront probablement inconnus à jamais, c'est une nécessité inhérente au fait d'étudier des états de langue révolus. Notre démarche, tant par la constitution du corpus que par le traitement qui en a été fait, a cependant tenté de réduire au mieux la part des zones d'ombre.

*La constitution de ce corpus n'aurait pas été possible sans l'immense travail réalisé par les responsables de la Base de Français Médiéval (ENS Lyon, anciennement ICAR UMR 5191, désormais IHRIM UMR 5317), plus spécifiquement Céline Guillot-Barbance pour son expertise, le choix des textes de la période médiévale, et Alexei Lavrentiev pour l'échantillonnage puis l'intégration des textes à la plateforme TXM (<http://textometrie.ens-lyon.fr/>), créée par Serge Heiden, qui a permis leur exploitation. Le projet de la Grande Grammaire Historique du Français leur doit beaucoup, et nous les en remercions très sincèrement.*

### 3.5 Liste des textes du corpus de la GGHF

Le tableau 1 ci-dessous présente la liste des textes du corpus et leurs caractéristiques. Les références complètes des textes sont données dans la bibliographie générale.

Les abréviations suivantes sont utilisées :

- *Forme* : P (prose) ; V (vers) ; M (mixte).
- *Domaine* : A (argumentatif) ; D (didactique) ; E (épistolaire) ; H (historique) ; J (juridique) ; L (littéraire) ; R (religieux).
- *Genre* : corresp. (correspondance) ; dramat. (dramatique) ; chroniq. (chroniques) ; hagiog. (hagiographie)
- *Dialecte* : champ. (champanois) ; angl.norm. (anglo-normand) ; ND (non défini) ;

Les textes dont le nombre de mots est suivi d'un astérisque ont été échantillonnés.

Siègle	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
<b>Avant 1100</b>								
<i>Corpus noyau avant 1100 :</i>								9588
Strasbourg	<i>Serments</i>		842	P	J	serments	ND	115
Eulalie	<i>Eulalie</i>		881	V	R	hagiog.	ND	188
Passion	<i>Passion de Clermont</i>		ca 1000	V	R	dramat.	ND	2904
StLegier	<i>Vie saint Léger</i>		ca 1000	V	R	hagiog.	ND	1406
StAlexis	<i>Vie saint Alexis</i>		ca 1050	V	R	hagiog.	normand	4975
<i>Corpus complémentaire avant 1100 :</i>								
Jonas	<i>Sermon sur Jonas</i>		entre 938 et 952	P	R	sermon	Flandre	815
<b>12<sup>e</sup> siècle</b>								
<i>Corpus noyau 12<sup>e</sup> s. :</i>								204980
Roland	<i>Chanson de Roland</i>		ca 1100	V	L	épique	normand	30039
Eneas1	<i>Eneas (1)</i>		ca 1155	V	L	roman	normand	35 152
Beroul Tristan	<i>Tristan</i>	Beroul	entre 1165 et 1200	V	L	roman	franco-picard	27 708
Pont-StMaxence Becket	<i>Vie de saint Thomas Becket</i>	Guernes de Pont Sainte Maxence	1172-1174	V	R	hagiogr.	ouest	39 145*
TroyesYvain	<i>Yvain</i>	Chrétien de Troyes	1177-1181	V	L	roman	champ.	42 331
Lapidaire	<i>Lapidaire en prose</i>		mi.-12 <sup>e</sup>	P	D	lapidaire	angl. norm.	4781
AmiAmil	<i>Ami et Amile</i>		ca 1200	V	L	épique	ND	25 824
<i>Corpus complémentaire 12<sup>e</sup> s. :</i>								160918
Benedeit Brendan	<i>Voyage de St Brendan</i>	Benedeit	déb. 12 <sup>e</sup>	V	R	hagiogr.	angl. norm.	10 955
Thaon Comput	<i>Comput</i>	Philippe de Thaon	1113 ou 1119	V	D	comput	angl. norm.	14 678
Descri Engleterre	<i>Description d'Engleterre</i>		peu ap. 1139	V	H	histoire	angl. norm.	1303
Psaut Cambridge	<i>Psautier de Cambridge</i>		entre 1155 et 1160	P	R	psautier	angl. norm.	4312
Eneas2	<i>Eneas (2)</i>		ca 1155	V	L	roman	angl. norm.	24 965
WaceBrut2	<i>Brut</i>	Wace	achevé en 1155	V	H	chroniq.	angl. norm.	15 675
Adgar Miracles	<i>Collection de Miracles</i>	Adgar	3 <sup>e</sup> tiers du 12 <sup>e</sup>	V	R	miracles	angl. norm.	49 330
SteMaure Chron Normandie	<i>Chronique des ducs de Normandie</i>	Benoît de Sainte Maure	1174	V	H	chroniq.	poitevin	25 285
Béthune Chansons	<i>Chansons</i>	Conon de Béthune	ca 1180-1190	V	L	lyrique	picard	2687
Charte Chièvres	<i>Charte de Chièvres</i>		1194	P	J	charte	traits picards	1282
Bodel Nicolas	<i>Jeu de Saint Nicolas</i>	Jehan Bodel	entre 1191 et 1202	V	L	dramat.	artois	10 446

Siècle	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
<b>13<sup>e</sup> siècle</b>								
<i>Corpus noyau 13<sup>e</sup> s. :</i>								223 298
ClariConstantinople	<i>Conquête de Constantinople</i>	Robert de Clari	ap. 1205	P	H	chroniq.	picard	34 292
RenartDole	<i>Roman de la rose ou de Guillaume de Dole</i>	Jean Renart	1210 ou 1228	V	L	roman	picard	35 050
Aucassin	<i>Aucassin et Nicolette</i>		dernier ¼ 12 <sup>e</sup> ou 1 <sup>ère</sup> moitié. 13 <sup>e</sup>	M	L	écrits brefs	traits picards	10 078
CoinciMiracles 1/2 3/ 4	<i>Miracles de Notre Dame</i>	Gautier de Coinci	1218-1227	V	R	miracles	picard	37 789*
Graal	<i>Queste del Saint Graal</i>		ca 1225	P	L	roman	ND	41 462*
LorrisRose	<i>Roman de la Rose</i>	Guillaume de Lorris	entre 1225 et 1230	V	D	roman	orléanais	24 325
MeunRose 1/2/3	<i>Roman de la Rose</i>	Jean de Meun	entre 1269 et 1278	V	D	roman	ND	19 563*
Beaumanoir Beauvaisis	<i>Coutume de Beauvaisis</i>	Philippe de Beaumanoir	1283	P	J	traité	traits picards	20 739*
<i>Corpus complémentaire 13<sup>e</sup> s. :</i>								166 046
Renart10/11	<i>Roman de Renart branches X-XI</i>		déb. 13 <sup>e</sup>	V	L	écrits brefs	ND	22 300
TristanProse	<i>Tristan en prose</i>		ap. 1240	P	L	roman	picard	75 186
MenestReims	<i>Récit d'un Ménestrel de Reims</i>		ca 1260	P	H	chroniq.	ND	50 046
CharteParis	<i>Chartes de la région parisienne</i>		1250 (1200-1299)	P	J	charte	Ile de France	18 514
<b>14<sup>e</sup> siècle</b>								
<i>Corpus noyau 14<sup>e</sup> s. :</i>								217 486
Joinville Mémoires	<i>Mémoires ou Vie de saint Louis</i>	Jean de Joinville	entre 1305 et 1309	P	H	mémoires	champ.	40 707*
Machaut Fortune	<i>Remède de Fortune</i>	Guillaume de Machaut	1341	V	L	lyrique	champ.	25 265
Froissart Chroniques	<i>Chroniques</i>	Jean Froissart	entre 1369 et 1400	P	H	chroniq.	picard	40 512*
Registre Chatelet1	<i>Registre criminel du Chatelet</i>		1389	P	J	procès	Ile de France	19 623*
Mesnagier	<i>Mesnagier de Paris</i>		1393	P	D	manuel	ND	19 993*
Griseldis	<i>Estoire de Griseldis en rimes et par personnages</i>		1395	V	L	dramatq.	traits picards	16 249
Manières 1396/1399	<i>Manières de langage</i>		1396, 1399	M	D	manuel	angl. norm.	20 315

Signe	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
QuinzeJoies	<i>Quinze Joyes de Mariage</i>		ca 1400	P	L	nouvelles	ouest	34 822
<i>Corpus complémentaire 14<sup>e</sup> s. :</i>								555 328
Bersuire Décades1/9	<i>Les Décades de Titus Livius, I,1 et I, 9</i>	Pierre Bersuire	1354	P	H	histoire	ND	71 667
Berinus1/2	<i>Berinus</i>		ca 1370	P	L	roman	Ile de France	150 628
Oreste Aristote Commentaire	<i>le Livre de Ethiques d'Aristote, Commentaire</i>	Nicole Oresme	1370	P	A	traité	ND	124 719
Phoebus Chasse	<i>Le Livre de chasse</i>	Gaston Phebus	1387	P	D	traité	picard	77 798
Arras Mélusine	<i>Melusine</i>	Jean d'Arras	1392	P	L	roman	ND	124 929
Deschamps ArtDictier	<i>L'art de dictier</i>	Eustache Deschamp	1392	M	D	traité	champ.	5 587
<b>15<sup>e</sup> siècle</b>								
<i>Corpus noyau 15<sup>e</sup> s. :</i>								201 214
Gerson Sermon	<i>Sermon pour le Fete de la Sainte trinité</i>	Jean Gerson	1402	P	R	sermon	ND	6 915
PizanCité	<i>Le livre de la Cité des dames</i>	Christine de Pizan	entre 1404 et 1405	P	D	exemples	Ile de France	40 734*
Orléans Ballades	<i>Ballades</i>	Charles d'Orléans	1415	V	L	lyrique	orléanais	22 251
Manières1415	<i>Manières de langage</i>		1415	M	D	manuel	angl. norm.	3 156
Pathelin	<i>Farce de Maitre Pathelin</i>		1456-1469	V	L	dramatq.	Ile de France	10 752
Cent Nouvelles	<i>Cent nouvelles nouvelles</i>		1456-1467	P	L	nouvelles	picard	39 449*
LouisXI Lettre223	<i>Lettres de Louis XI</i>	Louis XI	1461-1465	P	E	corresp.	Ile de France	2 362
LouisXI Lettre234	<i>Lettres de Louis XI</i>	Louis XI	1465-1469	P	E	corresp.	Ile de France	3 143
Archier Baignollet	<i>Le franc Archier de Baignollet</i>		1468	V	L	comique	Ile de France	2 500
LouisXI Lettre248	<i>Lettres de Louis XI</i>	Louis XI	1469-1472	P	E	corresp.	Ile de France	3 419
Commynes Mémoires	<i>Mémoires</i>	Philippe de Commynes	ca 1490-1505	P	H	mémoires	ouest	40 435*
JehanParis	<i>Roman de Jehan de Paris</i>		1494	P	L	roman	Ile de France	26 098
<i>Corpus complémentaire 15<sup>e</sup> s. :</i>								364 247
SaleSaintré	<i>Jehan de Saintré</i>	Antoine de la Sale	1456	P	L	roman	Ile de France	92 056
Bueil Jouvencel1/2	<i>Le Jouvencel</i>	Jean de Bueil	1461	P	D	roman	ND	123 452

Sigle	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
Villon Testament	<i>Testament</i>	François Villon	1461	V	L	lyrique	Ile de France	13 378
Commynes Lettres	<i>Lettres</i>	Philippe de Commynes	1478-1511	P	E	corresp.	ouest	20 019
Phares Astrologues	<i>Recueil des plus celebres astrologues et quelques hommes doctes</i>	Simon de Phares	1494-1498	P	D	traité	ND	115 342
<b>16<sup>e</sup> siècle</b>								
<i>Corpus noyau 16<sup>e</sup> s. :</i>								245 202
Vigneulles Nouvelles	<i>Cent Nouvelles Nouvelles</i>	Philippe de Vigneulles	1515	P	L	nouvelles	standard	40 321*
Calvin Lettres	<i>Lettres à monsieur et madame de Falais</i>	Jean Calvin	1549 (1543-1554)	P	E	corresp.	standard	33 735
DuBella Défense	<i>Défense et Illustration de la langue française</i>	Joachim du Bellay	1549	P	A	traité	standard	19 691
DuBella Olive	<i>L'Olive</i>	Joachim du Bellay	1550	V	L	lyrique	standard	21 214
Ronsard Misères	<i>Discours des misères de ce temps</i>	Pierre de Ronsard	1563	V	L	lyrique	standard	5 934
LaTailleSaül	<i>Saül le furieux</i>	Jean de la Taille	1572	V	L	dramatiqu.	standard	15 416
L'Estoile Registre 1/2/3/4/5	<i>Registre-journal du regne de Henri III (t.1-5)</i>	Pierre de l'Estoile	1574-75 1576-78 1585-87	P	H	registre	standard	28 842*
Léry Brésil	<i>Histoire d'un voyage fait en la terre du Brésil</i>	Jean de Léry	1578	P	L	récit de voyage	standard	40 278
Montaigne Essais	<i>Essais</i>	Michel de Montaigne	1592	P	A	traité	standard	39 771*
<i>Corpus complémentaire 16<sup>e</sup> s. :</i>								578 220
Consistoire Genève	<i>Minutes du consistoire de Genève</i>		1542	P	J	procès	standard	136 182
DesPériers Récréations	<i>Nouvelles récréations et joyeux devis</i>	Bonaventure des Périers	1561	P	L	nouvelles	standard	75 543
Palissy Recepte	<i>Recepte veritable</i>	Bernard Palissy	1563	P	D	traité	standard	58 418
Etienne Agriculture	<i>L'Agriculture et maison rustique</i>	Charles Estienne	1564	P	D	traité	standard	120 199
Vigénère Décadence	<i>L'Histoire de la decadence de l'Empire grec</i>	Blaise de Vigénère	1577	P	H	traité	standard	187 878
<b>17<sup>e</sup> siècle</b>								
<i>Corpus noyau 17<sup>e</sup> s. :</i>								233 009
Urfé Astrée	<i>L'Astrée, 2<sup>de</sup> partie</i>	Honoré d'Urfé	1610	P	L	roman	standard	39 939*
Béroalde Parvenir	<i>Le moyen de parvenir</i>	François Béroalde de Verville	1616	P	L	roman	standard	39 326*

Sigle	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
SorelBerger	<i>Le berger extravagant</i>	Charles Sorel	1627	P	L	roman	standard	20 308*
CorneilleCid	<i>Le Cid</i>	Pierre Corneille	1637	V	L	dramatiqu.	standard	18 160
Descartes Discours	<i>Discours de la méthode</i>	René Descartes	1637	P	A	traité	standard	23 142
Coëffeteau Histoire	<i>Histoire romaine</i>	Nicolas Coëffeteau	1646	P	H	traité	standard	20 302*
Assoucy Poësies	<i>Poësies et lettres ... contenant diverses pièces héroïques, satiriques et burlesque</i>	Charles Coypeau d'Assoucy	1653	V	L	lyrique	standard	28 855
Molière Précieuses	<i>Les Précieuses ridicules</i>	Molière	1660	V	L	dramatiqu.	standard	7 127
Rabutin Lettres2/3/1/4	<i>Les Lettres de messire Roger de Rabutin, comte de Bussy (t. 1-4)</i>	Roger de Bussy-Rabutin	1672 1681 1686 1692	P	E	corresp.	standard	20 019*
Racine Athalie	<i>Athalie</i>	Jean Racine	1691	V	L	dramatiqu.	standard	15 831
<i>Corpus complémentaire 17<sup>e</sup> s. :</i>								928 801
Serres Agriculture1 Serres Agriculture2	<i>Le Théâtre d'agriculture et mesnage des champs</i>	Olivier de Serres	1603	P	D	traité	standard	518 641
Gerhard Heroard	<i>Journal d'Heroard</i>	Jean Héroard	1601-1610	P	H	journal	standard	195 930
Lafayette Clèves	<i>La Princesse de Clèves</i>	Madame de La Fayette	1678	P	L	roman	standard	65 255
Bossuet Discours	<i>Discours sur l'histoire universelle</i>	Jacques-Bénigne Bossuet	1681	P	H	traité	standard	140 524
Fléchier Oraison	<i>Oraison fuëbre de Marie-Thérèse d'Autriche reine de France</i>	Valentin-Esprit Fléchier	1691	P	R	oraison	standard	8 451
<b>18<sup>e</sup> siècle</b>								
<i>Corpus noyau 18<sup>e</sup> s. :</i>								244 966
Regnard Légataire	<i>Le légataire universel</i>	Jean-François Regnard	1708	V	L	dramatiqu.	standard	17 874
Montesquieu Lois	<i>L'Esprit des lois</i>	Montesquieu	1755	P	A	traité	standard	19 922*
Voltaire Essay	<i>Essay sur l'histoire générale et sur les moeurs et sur l'esprit des nations</i>	Voltaire	1756	P	A	traité	standard	39 938*
Prévoist Mémoires	<i>Le Monde moral ou Mémoires pour servir à l'histoire du coeur humain</i>	L'Abbé Prévost	1760	P	L	roman	standard	26 455*

Sigle	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
RetifBretonne Paysan	<i>Le Paysan perversi ou les Dangers de la ville</i>	Nicolas Rétif de La Bretonne	1776	P	L	roman	standard	19 183*
Mirabeau Lettres	<i>Lettres originales écrites du donjon de Vincennes pendant les années 1777, 1778, 1779, 1780</i>	comte Honoré de Mirabeau	1780	P	E	corresp.	standard	19 546*
Beaumarchais Figaro	<i>La folle journée ou le Mariage de Figaro</i>	Beaumarchais	1785	V	L	dramatiqu.	standard	34 722
Robespierre Discours	<i>Discours</i>	Maximilien de Robespierre	1793	P	A	discours	standard	39 832*
Delille Géorgiques	<i>L'homme des champs ou les Géorgiques françaises</i>	Jacques Delille	1800	V	L	lyrique	standard	27 494
<i>Corpus complémentaire 18<sup>e</sup> s. :</i>								333 748
VaubanDixme	<i>Projet d'une dixme royale qui, supprimant la taille, les aydes, les doüanes d'une province à l'autre, les décimes du Clergé, les affaires extraordinaires...produiroit au Roy un revenu certain et suffisant</i>	Sébastien de Vauban	1707	P	D	traité	standard	39 737
Ramsay Cyrus	<i>Les voyages de Cyrus</i>	André-Michel de Ramsay	1727	P	L	roman	standard	66 578
Condillac Essai	<i>Essai sur l'origine des connaissances humaines</i>	Étienne Bonnot de Condillac	1746	P	D	traité	standard	83 684
Rousseau Discours	<i>Discours sur les sciences et les arts</i>	Jean-Jacques Rousseau	1750	P	A	traité	standard	9 015
Diderot Essais	<i>Essais sur la peinture / Salon de 1759 / Salon de 1761 / Salon de 1763</i>	Denis Diderot	1759-1766	P	A	traité	standard	69 310
Buffon Epoques	<i>Des époques de la nature</i>	Georges-Louis de Buffon	1778	P	D	traité	standard	65 424
<b>19<sup>e</sup> siècle</b>								
<i>Corpus noyau 19<sup>e</sup> s. :</i>								235 688
CodeCivil	<i>Le Code civil des Français</i>		1804	P	D	traité	standard	19 692*
ChâteaubriandGénie	<i>Le Génie du christianisme</i>	François-René de Châteaubriand	1803	P	H	traité	standard	19 596*
Musset Articles	<i>Articles publiés dans la Revue des deux mondes</i>	Alfred de Musset	1832	P	L	presse	standard	43 347

Sigle	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
FlaubertCorrespondance	<i>Correspondance (1830-1839)</i>	Gustave Flaubert	1839	P	E	corresp.	standard	16 888
Hugo Hernani	<i>Hernani</i>	Victor Hugo	1841	V	L	dramatiq.	standard	27 048
DuCampNil	<i>Le Nil, Egypte et Nubie</i>	Maxime Du Camp	1854	P	L	récit de voyage	standard	37 135
Baudelaire Fleurs	<i>Les Fleurs du Mal,</i>	Charles Baudelaire	1861	V	L	lyrique	standard	26 491
Goncourt Journal2/3/4	<i>Journal : Mémoires de la vie littéraire (t. 2,3,4)</i>	Edmond et Jules de Goncourt	1870 / 1890 / 1896	P	L	mémoires	standard	19 286*
Clemenceau Iniquité Clemenceau Réparation	<i>L'iniquité et Vers la réparation</i>	Georges Clémenceau	1899	P	L	mémoires	standard	26 205*
<i>Corpus complémentaire 19<sup>e</sup> s. :</i>								706 196
Procès TribunalSeine	<i>Procès instruit par le tribunal criminel du département de la Seine</i>		1801	P	J	procès	standard	22 873
StaelCorinne	<i>Corinne ou l'Italie</i>	Germaine de Staël	1807	P	L	roman	standard	192 672
Scribe Mariage	<i>Le mariage de raison</i>	Eugène Scribe	1826	P	L	dramatiq.	standard	16 900
Béranger Chansons	<i>Chansons</i>	Pierre-Jean de Béranger	1829	V	L	chanson	standard	22 898
Barbey Memorandum 1/2/3/4	<i>Memorandum 1, 2, 3, 4,</i>	Jules Barbey d'Aurevilly	1838-1858	P	L	mémoires	standard	138 865
DuCamp Hollande	<i>En Hollande, Lettres à un ami</i>	Maxime Du Camp	1859	P	E	corresp.	standard	53 750
Guyot Rapport	<i>Rapport de l'état de l'agriculture en Lorraine : 1789-1889</i>	Charles Guyot	1889	P	D	traité	standard	11 499
Claudé TêteOr	<i>Tête d'or</i>	Paul Claudel	1890	P	L	dramatiq.	standard	40 933
ZolaDébâcle	<i>La débâcle</i>	Emile Zola	1892	P	L	roman	standard	205 806
<b>20<sup>e</sup> siècle</b>								
<i>Corpus noyau 20<sup>e</sup> s. :</i>								239 897
Rolland JChristophe	<i>Jean-Christophe I : L'Aube</i>	Romain Rolland	1904	P	L	roman	standard	30 008*
VidalBlache Tableau	<i>Tableau de la géographie de la France</i>	Paul Vidal de la Blache	1908	P	D	traité	standard	26 679*
Apollinaire Alcools	<i>Alcools</i>	Guillaume Apollinaire	1913	V	L	lyrique	standard	17 121
Feydeau Maxim	<i>La dame de chez Maxim's</i>	Georges Feydeau	1914	P	L	dramatiq.	standard	32 136
Alain BeauxArts	<i>Système des Beaux-Arts</i>	Alain	1920	P	D	traité	standard	26 483*

Sigle	Titre	Auteur	Date	Forme	Domaine	Genre	Dialecte	Nombre de mots
MaletVie	<i>La vie est dégueulasse</i>	Léo Malet	1948	P	L	roman	standard	29451*
Green Journal 1/2/4/5	<i>Journal</i> (t. 1, 2, 4, 5)	Julien Green	1934 1939 1946 1950	P	L	journal	standard	29258*
Sartre Lettres1 Sartre Lettres2	<i>Lettres au castor et à quelques autres</i> (vol.1, 2)	Jean-Paul Sartre	1932 (1926-39) 1951 (1940-63)	P	E	corresp.	standard	19777*
Perec ModeEmploi	<i>La vie mode d'emploi</i>	Georges Perec	1978	P	L	roman	standard	28984*
<i>Corpus complémentaire 20<sup>e</sup> s. :</i>								777158
Poincaré Electricité PoincaréVa- leurScience Poincaré Mécanique Poincaré Leçons	<i>5 traités</i>	Henri Poincaré	1901- 1911	P	D	traité	standard	202654
Fournier Correspon- dance	<i>Correspondance avec Jacques Rivière</i>	Alain-Fournier	1905- 1914	P	E	corresp.	standard	367117
GideCongo	<i>Voyage au Congo</i>	André Gide	1927	P	L	récit de voyage	standard	60870
DabitHôtel	<i>L'Hôtel du Nord</i>	Eugène Dabit	1929	P	L	roman	standard	39758
Césaire Discours	<i>Discours sur le colonialisme</i>	Aimé Césaire	1955	P	A	traité	standard	12873
Zitrone Courses	<i>Léon Zitronne vous emmène aux courses</i>	Léon Zitronne	1962	P	D	traité	standard	7094
Jaccottet Chants	<i>Chants d'en-bas</i>	Philippe Jaccottet	1977	V	L	lyrique	standard	2263
Koltès Solitude	<i>Dans la solitude des champs de coton</i>	Bernard-Marie Koltès	1986	P	L	dramatiq.	standard	10820
IzzoKhéops	<i>Total Khéops</i>	Jean-Claude Izzo	1995	P	L	roman	standard	73709

Tableau 1 : Liste des textes des corpus noyau (échantillonné) et complémentaire de le GGHF

Sophie Prévost