

# 1 La problématique

Un manuscrit a besoin d'un

## 2 Les solutions proposées

Plusieurs projets avaient proposé des solutions au manque de cohérence dans la manière par laquelle la communauté scientifique caractérise les composants d'un texte numérisé. Hors de la France, les vocabulaires codicologiques ont été élaborés notamment en anglais pour les manuscrits médiévaux. Un exemple important est la base de données DigiPal (Digital Resource and Database of Palaeography, Manuscripts and Diplomatic), qui n'est plus mis à jour mais qui a été développé au sein du département des humanités numériques à King's College London (citation). L'un de ses auteurs, Peter Stokes, travaille actuellement en France et continue dans la même veine en collaborant avec des chercheurs français sur le projet *SegmOnto* (citation). Avant la création des *guidelines* de *SegmOnto*, le développement des vocabulaires codicologiques en France était toujours basé sur le modèle de Denis Muzerelle et son *Vocabulaire codicologique*, que nous expliquons prochainement.

### 2.1 Le Vocabulaire international de la codicologie

En 1985, Denis Muzerelle a conçu un vocabulaire codicologique qui avait pour but de fournir des médiévistes avec des termes uniformes pouvant décrire les aspects d'un manuscrit. Depuis l'apparition de son vocabulaire en français, des autres chercheurs sont venus pour adapter les termes de Muzerelle en d'autres langues. Marilena Maniaci a publié une version du *Vocabulaire codicologique* pour l'italien en 1997 (citation). Pilar Ostos, Luisa Pardo, et Elena Rodríguez en ont créé un pour l'espagnol l'année suivante (citation). Parfois appelé le *Vocabulaire international de la codicologie*, l'édition multilingue du *Vocabulaire codicologique* que Muzerelle a commencé en 1985 était maintenue jusqu'à l'édition d'une version 1.1. en 2002-2003 (citation).

### 2.2 La Codicologia

Aujourd'hui, la paléographie et l'étude des manuscrits peuvent profiter de l'application web *Codicologia* qui réunit le *Vocabulaire codicologique* ainsi que deux autres bases de données similaires: le projet multilingue *Lexicon* et le *Glossaire codicologique arabe*. Ses trois bases de données spécialisent dans divers écritures. Comme précisé avant, le *Vocabulaire codicologique* fournit une liste de termes en français, italien, espagnol, et anglais. Piloté par Philippe Bobichon, le projet *Lexicon* présente un vocabulaire en français pour décrire les manuscrits écrits en latin, roman, grec, hébreu, et arabe (citation). Un vocabulaire spécialisé plus profondément pour l'arabe est élaboré dans le *Glossaire codicologique arabe* d'Anne-Marie Eddé et Marc Geoffroy **noauthor\*glossaire\*2002**. Ce dernier a

été conçu au sein de l’Institut de recherche et d’histoire des textes après les modèles de Muzerelle et le vocabulaire codicologique en arabe d’Adam Gacek (citation).

L’application web *Codicologia* rassemblent ces projets et présente un vocabulaire bien étendu. Par exemple, *Codicologia* fournit 15 termes pour décrire une faute d’écriture dans un manuscrit. Certains de ses termes possèdent eux-même plusieurs définitions que les divers bases de données fournissent. Le terme *caviarder*, par exemple, a une définition courte dans le vocabulaire français de Muzerelle.

Supprimer un mot, un passage..., en le recouvrant largement d’encre, de façon à ce qu’il ne puisse être lu. (citation)

Selon le *Lexicon* de Bibichon, par contre, le *caviarder* se définit d’une manière plus détaillé et vise à expliquer l’étymologie du mot afin de préciser son usage dans le cadre des manuscrits des divers écritures.

Le mot apparaît en 1907 (noircir à l’encre) : il désigne alors un procédé appliqué par la censure russe, sous Nicolas Ier. Dans certains manuscrits grecs, le détail rempli d’encre est surmonté d’un point et d’un trait court destinés à le neutraliser. Ce procédé est très souvent utilisé parmi d’autres, pour la censure des manuscrits hébreux effectué sous l’autorité de l’inquisition, en Italie, à la fin du xvie siècle et au début du xviiie. **bobichon’caviarder’2011**

Étant élaboré à partir d’un corpus très diversifié, le *Lexicon* de Bibichon a moins de termes qu’a le *Vocabulaire codicologique* mais ses termes sont plus généralisés. Le vocabulaire de Muzerelle, par contre, fait plus de distinctions entre les aspects d’un manuscrit et donc a plus de termes distincts par rapport aux deux autres vocabulaires de l’application *Codicologia*.

En réunissant les trois bases de données, sans privilégier aucun, *Codicologia* présente un vocabulaire codicologique vraiment vaste. Cependant, l’application *Codicologia*, comme toutes ses bases de données, vise à répondre au manque de cohérence dans la manière par laquelle la communauté scientifique décrit les manuscrits. Le grandeur de son vocabulaire pose un problème à cet objectif. Ayant plus de deux milles termes en français—certains d’entre eux ont eux-même plusieurs définitions—la solution proposée par *Codicologia* livre un vocabulaire bien harmonisé et documenté mais trop étendu pour être appliqué à l’échelle dans une approche informatique. Sans un corpus d’entraînement gigantesque, qui coûterait une somme énorme, l’apprentissage automatique ne peut pas faire de distinction au niveau des termes conçus par Muzerelle et les autres auteurs des bases de données de *Codicologia*. Aujourd’hui, un modèle ne peut pas s’entraîner sur des milles des étiquettes possibles et arriver à distinguer entre, par exemple, 15 types de faute d’écriture. Un humaine peut le faire, et pour cette raison les bases de données de *Codicologia* sont utiles. Mais leurs vocabulaires ne conviennent pas bien à une approche informatique.

### 3 Les *guidelines* de *SegmOnto*

Le projet *SegmOnto* propose un vocabulaire plus petit et pourtant peut décrire une grande diversité de documents historiques, y compris les manuscrits et les imprimés. Cet objectif est encore plus compliqué à achever qu’un vocabulaire spécialisé aux manuscrits. Décrire les documents d’une diachronie longue, et sans préférence d’une écriture en particulier, exige un équilibre délicat entre la généralité et la particularité. Les *guidelines* du projet *SegmOnto* essaient de limiter le nombre de termes dans son vocabulaire sans priver un terme d’une identité distincte. Les *guidelines* se divisent en deux catégories génériques : les régions ou “zones” d’une page et les lignes de texte ou d’écriture sur la page. Chaque catégorie se compose d’une liste des étiquettes, chacune de lesquels cherche à parvenir à l’équilibre. Une étiquette devrait pouvoir être appliquée à soit un manuscrit, soit un imprimé, de peu importe quelle langue et quelle écriture.

#### 3.1 Les zones

- CustomZone
- DamageZone
- DecorationZone
- DigitizationArtefact
- DropCapitalZone
- MainZone
- MarginTextZone
- MusicZone
- NumberingZone
- QuireMarksZone
- RunningTitleZone
- SealZone
- StampZone
- TableZone
- TitlePageZone

### 3.2 Les lignes

- CustomList
- DefaultLine
- DropCapitalLine
- HeadingLine
- InterlinearLine
- MusicLine