

Contribution à la création d'un moteur de recherche sémiotique : application aux manuscrits latins médiévaux

Y. Leydier¹

F. LeBourgeois¹

H. Emptoz¹

¹Laboratoire d'InfoRmatique en Images et Systèmes d'information
INSA de Lyon

20 av. Albert Einstein, 69621 Villeurbanne cedex

{yann.leydier, frank.lebourgeois, huber.emptoz}@liris.cnrs.fr

Résumé

Cet article présente une méthode de recherche de mots par similarité de formes (word-spotting) dédiées aux manuscrits latins médiévaux. Nous proposons une nouvelle méthode de comparaison des formes qui tire avantage de la robustesse du gradient et tolère les variations spatiales. Nous testons notre algorithme sur plusieurs manuscrits latins médiévaux.

Mots clefs

Manuscrits médiévaux, recherche de mots par similarité de formes, *word-spotting*.

1 Introduction

Le travail présenté dans cet article s'inscrit dans le projet « Formes et couleurs des manuscrits médiévaux : élaboration d'un outil de recherche », piloté par l'IRHT¹. L'objectif du projet était de définir un outil de recherche facilitant l'accès aux gigantesques bases de données contenant les manuscrits médiévaux de l'IRHT. L'une des pistes explorées dans ce contexte concerne l'accès au contenu textuel des documents.

Certains mots jouent un rôle important dans les manuscrits médiévaux, comme par exemple « incipit » et « explicit » qui bornent les chapitres ou les livres dans un même volume. D'autres mots peuvent avoir un intérêt particulier pour les chercheurs, entre autres certains noms propres, ou, dans le cas de traduction, des mots dont le sens est inconnu et dont on cherche les différents contextes d'utilisation afin de déduire leur signification.

Il n'existe actuellement aucun système automatique capable de lire le texte des manuscrits médiévaux. Les principaux problèmes rencontrés ont pour origine la difficulté de segmenter le texte en mots, l'irrégularité de l'écriture, la complexité et la diversité des styles typographiques et le vocabulaire ouvert pour lequel nous ne disposons pas de dictionnaire.

Les systèmes de reconnaissance de texte utilisés industriellement aujourd'hui, les OCR (*Optical Character Recognition*), ne sont pas prévus pour traiter les images de do-

cuments manuscrits (voir figure 1). En fait, les OCR ne donnent des résultats corrects que sur les documents imprimés contemporains utilisant des polices de caractères usuelles (Times, Arial, etc.).



Figure 1 – Segmentation et reconnaissance d'un manuscrit médiéval par Fine Reader.

Une solution alternative consiste à limiter le processus à la reconnaissance d'un petit nombre de mots définis par les utilisateurs. Dans ce cas, il est possible de localiser avec précision toutes les occurrences d'un mot dans une image car nous en possédons toujours un modèle précis (défini par l'utilisateur).

La *recherche de mots par similarité de formes* (le terme *word-spotting* est plus souvent employé) est une technique permettant de localiser des mots choisis par un utilisateur dans un texte, écrit ou parlé, sans aucune contrainte [1, 2, 3, 4]. Cette approche générique peut être appliquée à tout type de document écrit, quel que soit son langage et qu'il utilise un alphabet, un syllabaire ou des idéogrammes... Il n'est pas nécessaire de créer une base d'apprentissage adaptée à chaque document ou à chaque scripteur.

Cette technique est utilisée lorsque la reconnaissance de mots est mise en échec, comme par exemple sur les documents très détériorés ou les manuscrits.

Dans notre cas, il s'agira de rechercher toutes les occurrences de l'image d'un mot. Le prototype, ou *mot-clé* est sélectionné par l'utilisateur en l'entourant sur une image du document à l'aide d'une interface graphique. Le mot-clé est comparé à des parties des images des pages du document en utilisant une mesure de similarité ou une distance. En fin de traitement, le système propose une liste d'images de mots triée par ressemblance avec le mot-clé. Cette dernière contient inévitablement des fausses détections que l'utilisa-

¹Institut de Recherche et d'Histoire des Textes.

teur élague manuellement. Cette technique est, en fait, plus proche du domaine de la recherche d'image par le contenu (*Content Based Image Retrieval*, CBIR) que de la reconnaissance de mots.

Les documents qui nous ont été fournis sont écrits dans différents styles typographiques [5]. Ces derniers correspondent à plusieurs grandes classes (gothique, carolingien...) mais de nombreux manuscrits présentent des variations et des mélanges de ces styles si bien qu'il est impossible de créer des modèles qui pourraient représenter la totalité de notre corpus.

Bien que l'écriture dans les manuscrits médiévaux puisse sembler assez stable au profane, la production d'un même scripteur s'avère parfois très irrégulière. La présence de réglure sur les pages ne suffit d'ailleurs pas toujours à assurer un alignement correct des mots et de fortes courbures de lignes dues aux conditions de prise de vue sont parfois observées. L'espace entre les lettres et les mots est aussi très irrégulier sur une même ligne, souvent guidé par le souci de justifier le texte et d'éviter les césures.

De plus, un livre peut avoir été rédigé par plusieurs copistes. Cela ajoute aux irrégularités de l'écriture, le changement de scripteur étant souvent visible à l'œil nu.

2 État de l'art

La plupart des travaux effectués sur la recherche de mots dans les images de documents repose sur une segmentation du document en mots. Quelques rares auteurs ne segmentent les documents qu'en lignes, opération bien mieux maîtrisée que la segmentation en mots dans le cas des images de manuscrits [3, 6].

Une méthode de segmentation en mots basée sur la théorie de l'espace multi-échelle a été proposée [7] et a donné de bons résultats sur un corpus de manuscrits de George Washington. Elle n'a cependant été testée que sur ce corpus qui possède des caractéristiques très différentes du nôtre dont les documents présentent une structure physique complexe et sont souvent endommagés. Si nous tentons de segmenter nos documents, les résultats ne seront pas optimaux et la qualité des traitements suivants sera compromise.

Il est possible de comparer les pixels directement en utilisant des méthodes sophistiquées comme la distance SLH (*Scott et Longuet-Higgins*) [1] qui est robuste vis à vis des transformées affines, ou bien encore la distance de Hausdorff [4]. La plupart de ces méthodes nécessitent cependant une binarisation de l'image qui cause inévitablement une perte d'information importante.

Des mots segmentés peuvent être représentés par des vecteurs de caractéristiques. Ces caractéristiques sont très proches de celles utilisées pour la reconnaissance de caractères : profils, projections, lissage gaussien [8], etc.

Lesdits vecteurs sont souvent comparés grâce à l'algorithme DTW (*Dynamic Time Warping*) [3, 8, 9] qui est robuste face aux variations spatiales. L'association de ces

vecteurs de caractéristiques et de cet algorithme de comparaison est très efficace mais suppose une segmentation parfaite des mots, ce qui ne peut pas être réalisé sur nos documents.

D'autres caractéristiques comme les concavités [2] ou les coins [10] nécessitent des structures de données plus complexes pour les décrire (en général, des arbres ou des graphes). Ces caractéristiques sont très stables sur des documents propres mais elles sont peu adaptées aux documents de notre corpus.

Dans les études précédentes, certains résultats sont présentés en termes de précision [8] ou à l'aide de courbes ROC (*receiver operating characteristics*) qui représentent le taux de succès en fonction des fausses détections.

Certains auteurs préfèrent donner des valeurs numériques (par exemple le nombre de bonnes réponses) plutôt que des taux mais le manque d'information sur ces valeurs les rend peu facile à interpréter et comparer. Certaines études présentent la liste des images des n meilleures réponses pour certaines requêtes en guise d'exemples [11].

Récemment, des auteurs ont commencé à utiliser des courbes Précision-Rang [12] et Précision-Rappel [13], déjà utilisées dans le domaine de la recherche d'images par le contenu.

Les caractéristiques utilisées dans les méthodes présentées sont très classiques et nécessitent pour la plupart une binarisation des images, voire une segmentation en mots. Bien que la possibilité de segmenter certains manuscrits du Moyen Âge a été démontrée, tous ne sont pas segmentables. Afin de garantir que notre méthode pourra être appliquée à l'ensemble des documents de notre corpus, nous avons décidé de développer une méthode qui n'applique aucune segmentation, que ce soit en lignes, en mots ou même en graphèmes. De plus, nous n'effectuerons aucune binarisation afin de conserver un maximum d'information.

3 Comparaison élastique cohésive

Nous avons testé différentes caractéristiques directement applicables sur les images en niveaux de gris et sans segmentation telles que le gradient, les valeurs propres de la matrice hessienne, la courbure des isophotes, la courbure des lignes de flux, etc.

Les expérimentations ont montré que les meilleurs résultats sont donnés par l'orientation du gradient lorsque sa magnitude est significative. Les performances baissent si on baisse le seuil sur la magnitude du gradient : plus le seuil est bas, plus le bruit entre en compte dans les calculs.

La distance entre deux gradients est définie ici par l'angle entre eux si les deux magnitudes sont supérieures au seuil. Si l'un des deux gradients a une magnitude trop faible, la distance prend la valeur d'une pénalité. Cette pénalité est égale au double de l'angle maximal possiblement formé par deux gradients. La distance entre deux gradients de faible magnitude est définie comme nulle. Ainsi, dans l'espace d'échelle σ , \mathcal{G} étant l'opérateur gradient ($\|\mathcal{G}\|$ sa

norme et $\theta(\mathcal{G})$ son argument), la distance entre deux pixels a et b est définie par :

$$d_{\sigma}(a,b) = \begin{cases} \min(|\theta(\mathcal{G}^a) - \theta(\mathcal{G}^b)|, \\ 128 - |\theta(\mathcal{G}^a) - \theta(\mathcal{G}^b)|), & \text{si } \|\mathcal{G}^a\| > \epsilon \text{ et } \|\mathcal{G}^b\| > \epsilon \\ 0, & \text{si } \|\mathcal{G}^a\| < \epsilon \text{ et } \|\mathcal{G}^b\| < \epsilon \\ 255, & \text{sinon} \end{cases}$$

De manière générale, pour qu'un algorithme de comparaison soit robuste face aux variations spatiales, chaque pixel du prototype doit être comparé à tous les pixels d'un voisinage sur l'image. Un tel processus est très lourd en calculs mais peut être simplifié. En ne comparant que des zones d'intérêt idoines, le temps de calcul diminue alors drastiquement.

Dans les manuscrits médiévaux, la plupart de l'information du texte est localisée sur des traits verticaux. Illustrons cela avec la lettre « u ». La distance entre les deux traits verticaux de deux occurrences de « u » manuscrits est toujours différente, ce qui met en échec les algorithmes de comparaison naïfs. Il est donc plus pertinent de comparer les pixels autour des traits verticaux que de comparer les deux formes dans leur intégralité.

Nous calculons une signature des images par morphologie mathématique appliquée aux niveaux de gris. Nous effectuons une ouverture avec un élément structurant choisi en fonction du type de document. Un élément structurant idoine pour les manuscrits médiévaux latins est une ligne verticale dont la longueur n dépend de la taille des caractères. n est déterminé visuellement et de façon simple et intuitive par l'utilisateur grâce à une interface graphique. D'autres éléments structurants, obliques par exemple, pourraient être utilisés pour compléter la signature du texte au risque d'y ajouter du bruit.

Nous obtenons ainsi une signature du texte composée de *guides* verticaux. En agrandissant le rectangle englobant de ces guides, nous obtenons des *zones d'intérêt* (voir figure 2). Notons que l'extraction de zones d'intérêt (ou ZI) n'est pas une segmentation en lignes ni en mots. En effet, nous ne cherchons pas à extraire des entités cohérentes d'un point de vue sémantique.

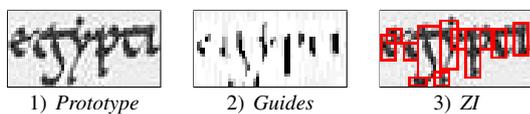


Figure 2 – Guides et zones d'intérêt du mot «egypti».

Le prototype est divisé en morceaux suivant ses zones d'intérêt. La distance et l'orientation entre les centres successifs des dites zones sont stockées (voir figure 3.2). Les liens entre elles sont lâches (voir figure 3) si bien que le prototype peut être déformé afin de mieux s'apparier aux occurrences les plus déviantes (voir figure 4). Le déplacement

possible de chaque zone d'intérêt permet un léger décalage vertical dépendant de la hauteur des caractères. Le déplacement horizontal doit être plus large sans pour autant permettre à des zones d'intérêt de se croiser ; ainsi en général une ZI ne doit pas être décalée de plus la moitié de la largeur moyenne des caractères. Ces paramètres sont réglés par l'opérateur car nous ne pouvons pas segmenter le texte en caractères.

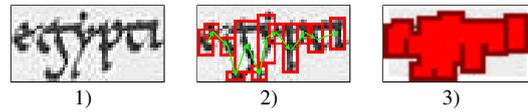


Figure 3 – 1) Un prototype. 2) Les liens entre les ZI. 3) En rouge, les ZI, en rouge foncé leur aire de déplacement possible.



Figure 4 – Un déplacement possible des ZI du prototype «egypti». Le prototype est déformé afin d'être comparé à des mots dont la forme diffère légèrement.

Afin de ne pas comparer le prototype naïvement à la totalité de l'image, nous avons décidé d'orienter la comparaison grâce aux guides de l'image. Ces derniers servent de points de départ, ensuite les liens entre les zones d'intérêt du prototype sont utilisés pour déterminer l'ensemble des zones à comparer

Le processus est illustré sur la figure 5. Les rectangles transparents sont les zones d'intérêt de prototype, les rectangles gris pointillés sont les guides de l'image. Seul la première ZI du prototype doit être calée sur un guide de l'image.

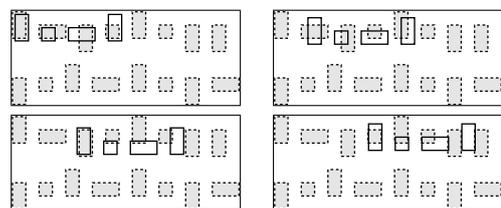


Figure 5 – Exemple de comparaison cohésive.

pour chaque guide G de l'image
 pour chaque déplacement d dans le voisinage
 comparer la 1ère ZI Z du prototype
 [avec l'image en coord(G) + d]
 $score(G) = \min(resultats)$
 $X = coord(Z)$

pour chaque autre zone d'intérêt Z du
 [prototype
 pour chaque déplacement d dans le
 [voisinage
 comparer Z avec l' image en
 [coord(G) + coord(Z) - X + d
 score(G) = score(G) + min(resultats)
 X = coord(Z)

4 Résultats

Nous avons créé une interface permettant de lancer une recherche sur un groupe d'images. Le mot clé est sélectionné par l'utilisateur en l'encadrant sur une image avec la souris.

4.1 Test MS14

La première série de tests présentée a été effectuée sur 24 pages du manuscrit MS14 de la bibliothèque d'Amiens (voir figure 6.1), soit environ 11000 mots. Nous avons choisi les mots-clés selon deux critères : leur pertinence sémantique, pour simuler une utilisation réelle du système, et leur fréquence d'apparence, pour des raisons statistiques. Étant donné que le texte est en latin, nous avons coupé les déclinaisons des mots-clés afin d'en garder les racines. Des recherches ont été effectuées pour les mots-clés suivants : « aaron », « quod », « terra », « ego », « manu », « moyse- », « dño » et « pharao- ». Nous avons compté le nombre d'occurrences de chaque mot-clé et noté les occurrences césurées, avec une majuscule ou une lettrine séparément. En effet, notre méthode n'est pas prévue pour gérer ce genre de cas et nous avons voulu en tenir compte dans nos statistiques. Ainsi notre algorithme est sensé ne pouvoir localiser que 488 occurrences au mieux sur les 530 des huit mots-clés, soit 92,1%.

Au final, notre algorithme a localisé 437 occurrences de mots-clés (82,5%) soit 9,6 points de moins que la maximum théorique (voir table 1). La courbe P-R est donnée figure 6.3.

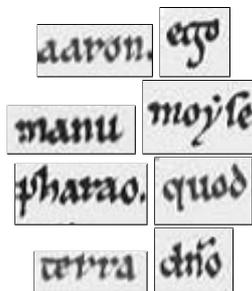
	Total	Localisé	%age
Complet	488	433	88,7%
Lettrine	1	1	Ø
Majuscule	34	1	2,9%
Césure	7	2	Ø
Total	530	437	82,5%

Tableau 1 – Résultats sur le MS14.

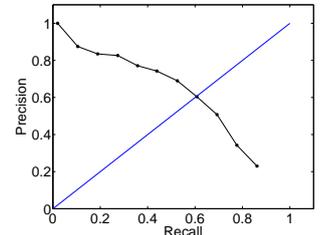
Les mots-clés donnant les plus mauvais résultats en termes de rappel sont « aaron » et « manu ». Le premier s'explique par la présence de deux « a ». Cette lettre est probablement celle qui a la plus grande variabilité dans le texte de ce document. De plus elle ne possède aucune hampe ou jambe pour la caractériser. Le second mot-clé contient un « m », un « n » et un « u ». Ces trois lettres peuvent facilement être confondues avec d'autres à cause de la faible résolution des images (par exemple, « m » et « ni », « nu » et « mi », etc.).



1) Une page du manuscrit MS14.



2) Les huit prototypes.



3) La courbe P-R. P = R = 0,6, P(R = 1) = 0,1

Figure 6 – Le test MS14.

Les mots commençant par des majuscules n'ont pas été localisés à cause de la trop grande différence morphologique entre les lettres majuscules et minuscules dans ce type d'écriture.

Les résultats sont globalement très bons compte tenu de la résolution de l'image et de l'irrégularité de l'écriture.

4.2 Test Escalopier 22

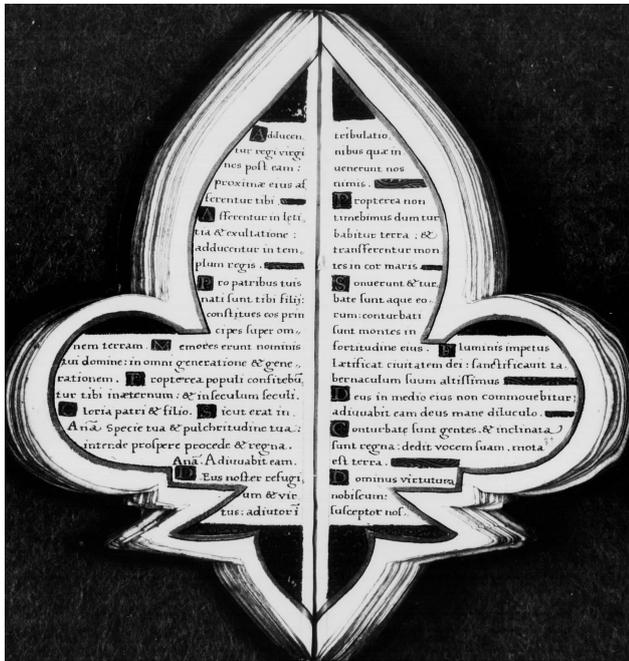
Nous avons ensuite extrait 24 pages (soit environ 2000 mots) dans le manuscrit Escalopier 22 (voir figure 7.1) de la bibliothèque d'Amiens. Nous avons utilisé les mêmes critères que précédemment pour sélectionner les mots-clés et avons choisi : « benedic- », « deus », « domin- », « fili », « gloria », « Maria », « patri », « terra » et « virg- ». À cause des césures, des majuscules et des lettrines, notre système n'est sensé pouvoir localiser que 153 occurrences sur les 195 occurrences des neuf mots-clés, soit 78,5%.

Nous avons pourtant réussi à localiser 171 mots-clés (87,7%) soit 9,2 points de plus que le maximum théorique (voir table 2). La courbe P-R est donnée figure 7.3. Tous les mots-clés ont été localisés à des taux de précision et rappel équivalents.

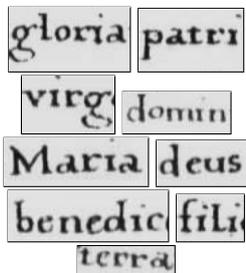
Notons qu'un plus grand nombre de mots en majuscules et de mots césurés ont été localisés sur ce manuscrit que sur le précédent. Cela semble indiquer que le style d'écriture

	Total	Localisé	%age
Complet	153	151	98,7%
Lettrine	27	12	44,4%
Majuscule	3	3	Ø
Césure	12	5	41,7%
Total	195	171	87,7%

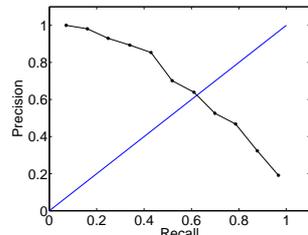
Tableau 2 – Résultats sur Escalopier 22.



1) Une page du manuscrit Escalopier 22.



2) Les neuf prototypes.



3) La courbe P-R. $P = R = 0.62$, $P(R = 1) = 0.19$.

Figure 7 – Le test Escalopier 22.

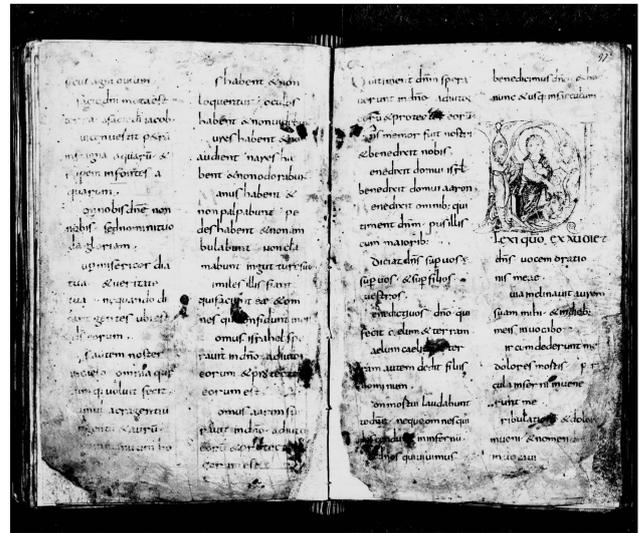
de ce manuscrit est plus discriminant. Cela montre de plus que cette écriture est plus régulière.

Cette fois encore les résultats sont très bons avec une précision relativement élevée pour de très hautes valeurs de rappel.

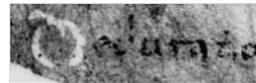
4.3 Test MS18

Nous avons lancé une autre série de tests sur le manuscrit MS18 de la bibliothèque d'Amiens. C'est un document fortement endommagé contenant de nombreuses taches et dont l'encre est souvent estompée.

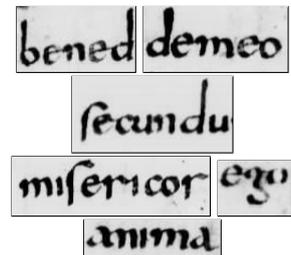
Nous avons extrait 24 pages (soit environ 3100 mots) pour notre test. Malheureusement, ce document ne contient pas beaucoup de mots redondants et nous avons dû limiter notre test à la recherche de six mots (« de meo », « ego », « bened- », « secundu- », « misericor- », « anima »). Étant donné que même ces mots ne sont pas très fréquents, nous ne présenterons pas de courbe P-R car elle ne serait pas significative.



1) Une page du MS18



2) Le microfilmage ne rend pas correctement les couleurs des letrines. Celle-ci n'est visible que parce que le parchemin est sombre à cet endroit.



3) Les six prototypes.

Figure 8 – Le test MS18.

Les résultats sont très bons étant donné l'état du document (voir table 3). Si le rappel est encore une fois très satisfaisant, notons que la précision est assez faible. Certaines occurrences des mots-clés sont si gravement endommagées qu'elles n'apparaissent que très loin dans la liste des réponses (voir figure 9).

	Total	Localisé	%age
Complet	49	48	98%
Lettrine	9	1	11,1%
Majuscule	0	0	Ø
Césure	2	0	Ø
Total	60	49	81,7%

Tableau 3 – Résultats sur le MS18.

Les mots commençant par une letrine n'ont pas été localisés.

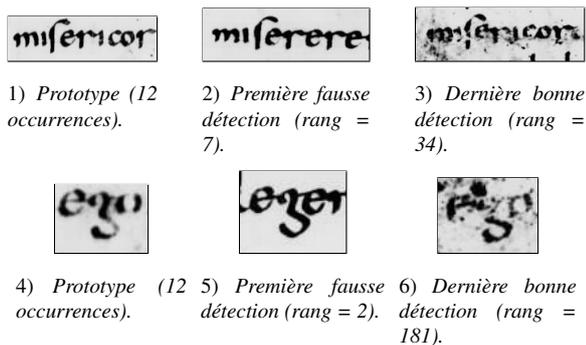


Figure 9 – Prototypes, premières fausses détections et dernières bonnes détections de deux mots-clés sur le MS18.

sés. Dans ce document, les lettrines apparaissent en blanc (voir figure 8.2). Si le parchemin n’est pas taché, les lettrines sont invisibles, le cas échéant, elles sont blanches : aucun guide n’en est jamais extrait. Ainsi, lors de la comparaison du prototype, la première ZI de ce dernier est calée sur le second guide du mot sur l’image. L’élasticité de notre algorithme n’a pas été conçue pour surmonter un tel décalage car si les liens entre les ZI étaient trop lâches, l’algorithme sauterait des morceaux de mots et pourrait, par exemple, apparier « sale » avec « sable ».

5 Conclusion

Nous avons décrit une nouvelle méthode de localisation de mots adaptée aux manuscrits latins médiévaux. Cette méthode est particulièrement originale dans le sens où elle ne nécessite aucune segmentation ou binarisation du texte.

Nous avons choisi comme caractéristique l’orientation des gradients dans les zones où ces derniers ont une magnitude significative. Cette mesure est robuste aux faibles variations géométriques et rend compte de la structure locale des formes.

Nous avons proposé un algorithme de comparaison robuste aux irrégularités géométriques du texte manuscrit. Cet algorithme a trois principaux avantages : il est élastique et cohésif afin de gérer de grandes déformations des caractères et il est plus rapide que les méthodes de corrélation simples.

Nous avons testé notre méthode sur des manuscrits latins et avons obtenu de très bons résultats.

Des tests sont en cours sur d’autres types de documents, notamment des documents imprimés très dégradés (binarisés avec tramage), des manuscrits arabes, hébreux, tibétains et japonais.

Nous avons prévu d’améliorer l’extraction des guides afin de la rendre indépendante du type des documents traités.

Une méthode de raffinement des résultats sera mise au point afin de trier la liste des réponses de manière plus fine. Nous réfléchissons par ailleurs à un moyen de rendre notre algorithme de comparaison robuste aux homothéties afin de régler le problème de la normalisation.

Références

- [1] R. Manmatha, C. Han, et E.M. Riseman. Word spotting : A new approach to indexing handwriting. Dans *CVPR*, pages 631–637, San Francisco, Etats-Unis, 1996.
- [2] P. Keaton, H. Greenspan, et R. Goodman. Keyword spotting for cursive document retrieval. Dans *Workshop on Document Image Analysis*, pages 74–81, San Juan, Puerto Rico, 1997.
- [3] A. Kołcz, J. Alspector, M. Augusteijn, R. Carlson, et G. Viorel Popescu. A line-oriented approach to word spotting in handwritten documents. *PAA*, 3 :153–168, 2000.
- [4] Y. Lu et C.L. Tan. Word spotting in chinese document images without layout analysis. Dans *ICPR*, pages 57–60, Quebec, Canada, 2002.
- [5] A. Derolez. *The Paleography of Gothic Manuscript Books*. Cambridge University Press, Cambridge, Grande-Bretagne, 2003.
- [6] J. Edwards, Y.W. Teh, D. Forsyth, R. Bock, et M. Maire. Making latin manuscripts searchable using ghmm’s. Dans *Neural Information Processing Systems*, pages 385–392, Cambridge, Etats-Unis, 2004.
- [7] R. Manmatha et J.L. Rothfeder. A scale space approach for automatically segmenting word from historical handwritten documents. *IEEE TPAMI*, 27(8) :1212–1225, 2005.
- [8] T.M. Rath et R. Manmatha. Features for word spotting in historical manuscripts. Dans *ICDAR*, volume 1, pages 218–222, Edinbourg, Ecosse, 2003.
- [9] T.M. Rath et R. Manmatha. Word image matching using dynamic time warping. Dans *IEEE Computer Vision and Pattern Recognition*, pages 521–527, Madison, Etats-Unis, 2003.
- [10] J.L. Rothfeder, S. Feng, et T.M. Rath. Using corner feature correspondences to rank word images by similarity. Dans *Conference on Computer Vision and Pattern Recognition Workshop*, pages 30–35, Madison, Etats-Unis, 2003.
- [11] R. Manmatha, C. Han, E.M. Riseman, et W.B. Croft. Indexing handwriting using word matching. Dans *ACM First Intl. Conf. on Digital Libraries*, pages 151–159, Bethesda, Etats-Unis, 1996.
- [12] T.M. Rath, R. Manmatha, et V. Lavrenko. A search engine for historical manuscript images. Dans *ACM SIGIR conference on Research and development in information retrieval*, pages 369–376, Sheffield, Royaume-Unis, 2004.
- [13] N.R. Howe, T.M. Rath, et R. Manmatha. Boosted decision trees for word recognition in handwritten document retrieval. Dans *ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 377–383, Salvador, Brésil, 2005.