

Introduction au Word Spotting

Frédéric Rayar

Partha Pratim Roy

Jean-Yves Ramel



Introduction

- Préservation et accès au patrimoine culturel
 - Numérisation
 - Transcription (OCR)
- Difficultés à retranscrire les documents historiques
- Alternative : le Repérage de mots (*Word Spotting*)
 - On ne cherche plus à reconnaître chaque caractère
 - Listes de régions similaires à la requête

Word Spotting

- Domaine de recherche actif :
 - Documents modernes, anciens
 - Documents manuscrits, imprimés
 - Texte ou graphique simple (logo, ...)
 - Avec ou sans segmentation (lignes, mots)

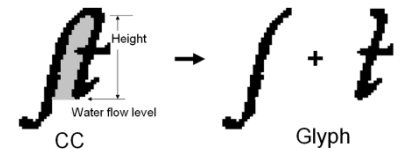
Approche proposée

➤ Objectifs

- Robustesse aux documents historiques
- Indépendance à l'alphabet
- **Rapidité du temps de réponse à une requête**

➤ Principes exploités

- Segmentation en ligne (AGORA)
- Considération de primitives de caractères
- Indexation hors ligne (dictionnaire de formes)
- Recherche reportée au niveau du texte



Résultats expérimentaux

- Centre d'Etudes Supérieures de la Renaissance
- « *Essais – Livre I* » de Montaigne (1580)



peut

faicte de iour. Chacun peut penser,

il y peut auoir du masque, ou ces beaux

Et de la se peut engédrer par fois la de-

semble.

seule, ce me semble, ou il se trouue en fa

ceté: ils semblent auoir d'autant moins

tes, semble se desprendre, se desmeler

Résultats expérimentaux

- « La mendicité spirituelle » de Gerson (1501)

Et l'homme respond que saint desir peut aller par tout
De ie te respons que les sieuy & les places ou ces
Aumosnes se doient pprement trouver / sont les
sieuy spirituelz des cieuy comme iay dit par a-
nant / et la te peult et doit porter ton desir. Car desir est done
a toy en semblance dung cheual Volant come les poetes fai-
gnoient q pegasus estoit. Et en sont les deux elles paour &
esperance. Du est come Dng chariot de feu semblable a cel-
luy qui porta hebe en paradis terrestre. Et sont les quatre
roes les quatre Vertus cardinaulz: prudēce / force / iustice &
atrempance. Les cheuaulz qui se tirent: Les troys Vertus
theologales: foy / esperāce / & charite. Et le charrestier discre-
tion. Icy pourroit estre belle matiere a poursuir ce que ie
passe. Si te peuz porter ou comporte ou tu Vouloiras. Mais
aduisse bien ou tu yras: car la ou tu enuoiras ton desir illec.

autres

iesuchrist et autres furent

moins recoitnēt laumosne. Tu se Voisd aucuns autres qui

pour moy. Et Vne bôte autre requierē.

Résultats expérimentaux

- Journaux en (a) hindi et (b) bengali
- Rongo (*Entretiens de Confucius en japonais*) (1533)

भूगर्भ (a)

भूगर्भ खदानों से कोयला निकालने

रहा है और ओपेन कास्ट से भूगर्भ

के पास भूगर्भ से कोयला निकालने

বন্ধ (b)

নিজস্ব প্রতিনিধি, কলকাতা: এ রাজ্যে ঘন ঘন বন্ধ

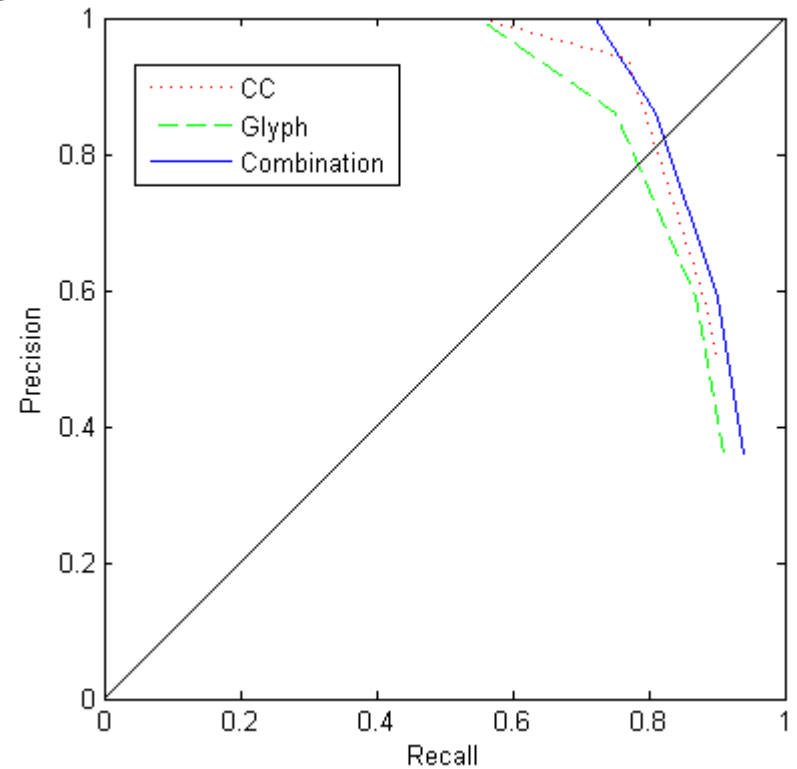
কর্মীর সম্মতি নিয়ে বন্ধ ডাকা এবং চিকিৎসক,

বলেন, বন্ধের ফলে শিল্পের ক্ষতি হচ্ছে বলে যে প্রচার



Résultats expérimentaux

- « Essais – Livre I » de Montaigne (1580)
 - Vérité terrain fournie par le CESR
 - Sous-ensemble de 78 pages
 - 1579 lignes indexées



- 72% de résultats pertinents avec 100% de précision

Résultats expérimentaux

- PC avec un processeur de type i5 à 2.53 GHz et 4G de RAM
- Tests menés sur le sous-ensemble de 78 pages
- La recherche a été effectuée plusieurs fois, et un temps de réponse moyen a été calculé

Image requête	Temps (s)
(a) (a)	1,2
(b) (b)	2,5
(c) (c)	4,5

- Temps de réponse de 2 secondes en moyenne

Projet ReNom

