

# Extraction, Exploitation and Evaluation of Document-Based Knowledge

## Différentes approches généralistes du texte

Antoine Doucet  
<http://www.info.unicaen.fr/~doucet>

Habilitation à Diriger des Recherches



## Plan

- Introduction
- Angle Séquentiel
  - Fouille de données textuelles, statistique appliquée
  - Application en RI multilingue
- Angle Discursif
  - Veille Multilingue
- Angle Structurel
  - RI Structurée
  - Clustering
- Évaluation
  - Book Search Track
- Conclusion



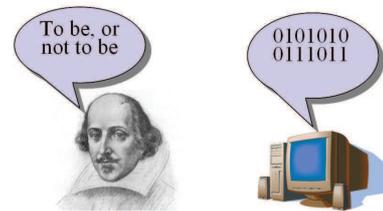
## Enjeux de la description documentaire

- ▶ Le traitement documentaire requiert un grand nombre de comparaisons entre documents
  - ▶ Un moteur de recherche Internet doit par exemple décider quelles pages vont vous convenir parmi quelques milliards, étant donnée votre préférence, définie par la requête soumise
  - ▶ Une réponse en une seconde est considérée comme lente



## Modéliser des documents

- ▶ Il faut donc pouvoir comparer les documents très efficacement
- ▶ Une réponse naturelle est d'opter pour des représentations de document simples



## Pourquoi « générique et multilingue » ?

- ▶ Forte corrélation avec le passage à l'échelle
- ▶ La part de l'anglais sur Internet décroît
- ▶ De très nombreuses langues sont « peu dotées »
- ▶ Ce n'est pas un dogme :
  - ▶ Pose de jalons génériques
    - ▶ dont les performances peuvent souvent être améliorées par des méthodes spécifiques au corpus



## Utilisation de la nature séquentielle de la donnée textuelle

- ▶ Travaux de thèse : *Advanced Document Description : a Sequential Approach*
  - ▶ Extraction, Sélection et Exploitation de séquences d'items (contexte applicatif textuel non restrictif)
- ▶ Pour le texte :
  - ▶ Méthodes entièrement généralistes (notamment multilingues)
  - ▶ Unités lexicales complexes
    - ▶ « cordon bleu » ≠ cordon + bleu
    - ▶ « cordon vert » ? « café puissant » ?
    - ▶ red cell phone



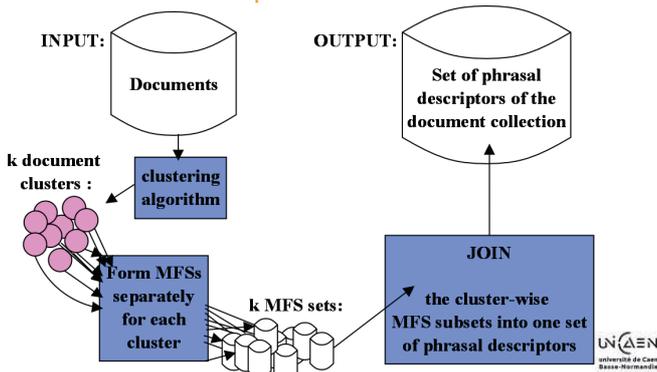
## Extraction de motifs séquentiels discontinus I

- ▶ Séquentiels : la seule contrainte est l'ordre
  - ▶ dans l'état de l'art : contraintes sur la distance, la fréquence, la longueur, des motifs linguistiques...
- ▶ Discontinus : qu'importe la distance les séparant
- ▶ Séquences Fréquentes Maximales (MFS)
  - ▶ Maximales et Fréquentes
  - ▶ Description compacte
    - ▶ Une 10-séquence remplace  $\binom{10}{2} = 45$  paires
    - ▶ Aucune contrainte de longueur
    - ▶ Distance illimitée entre les composants

## Extraction de motifs séquentiels discontinus II

- ▶ On peut noter la similarité entre les deux fragments :
  - ▶ « ancien Président Clinton »
  - ▶ « l'ancien Président des États-Unis Bill Clinton »
- ▶ Méthode
  - ▶ Expansion par combinaison des 2- et 3-séquences fréquentes
  - ▶ Diviser pour régner

## Extraction de motifs séquentiels discontinus III



## Filtrage des séquences extraites I

- ▶ Un ensemble de séquences descriptives est ainsi relié à chaque document
- ▶ Beaucoup sont très peu discriminantes : « *the be the a* », « *in the of* », « *barrel dollar* »...
- ▶ Probabilité d'occurrence d'une séquence discontinue
  - ▶ Expectative de fréquence documentaire
    - ▶ Test statistique permettant l'évaluation directe de l'intérêt des séquences (mesure de cohésion lexicale)
  - ▶ Réalisation d'un classement automatique
    - ▶ Indépendant du contexte applicatif (lexicographie, RI)
    - ▶ Combinant les séquences de différentes tailles

## Filtrage des séquences extraites II

- ▶ Calcul de la probabilité d'occurrence d'une séquence discontinue  $p(A_1 \rightarrow \dots \rightarrow A_n, l)$ 
  - ▶ Approche naïve :  $O(ln^{l-n})$ 
    - ▶ Disjonction par position initiale de succès
    - ▶ Soit  $E_i$  l'ensemble des documents contenant la  $n$ -séquence après exactement  $(n+i)$  item :
 
$$E_i = \{\bar{A}_1^{k_1} A_1 \bar{A}_2^{k_2} A_2 \dots \bar{A}_n^{k_n} A_n W^{l-n-1} \mid \sum_{j=1}^n k_j = i\}$$
  - ▶ Or les  $(E_k)$  sont disjoints, donc :
 
$$p(A_1 \rightarrow \dots \rightarrow A_n, l) = \prod_{i=1}^n p_i \sum_{l_n=0}^{l-n} \dots \sum_{l_1=0}^{l-n-(l_n+\dots+l_2)} q_1^{l_1} q_2^{l_2} \dots q_n^{l_n}$$

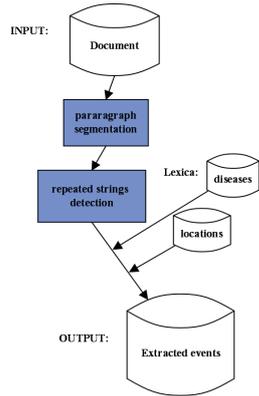
## Filtrage des séquences extraites III

- ▶ Calcul de la probabilité d'occurrence d'une séquence discontinue  $p(A_1 \rightarrow \dots \rightarrow A_n, l)$ 
  - ▶ Approche naïve :  $O(ln^{l-n})$
  - ▶ Modèle Markovien :  $O(ln^3)$
  - ▶ Modèle Markovien + Algèbre linéaire :  $O(ln)$ 
    - ▶ Propriétés spécifiques à la matrice de transition, notamment Jordanisation à coût algorithmique réduit



## L'approche MultiPULS

- ▶ Approche simplifiée au grain texte
- ▶ avec ressources minimales pour le filtrage des événements



## Application d'un modèle rhétorique

- ▶ Les positions sont exploitées
- ▶ Les répétitions sont exploitées
  - ▶ Au grain caractère
  - ▶ Par extraction de chaînes répétées maximales
- ▶ Les paires significatives sont extraites
  - ▶ Maladie - Lieu
  - ▶ Maladie - Nombre de cas

## Exemple : texte brut

WHO checks smallpox reports in Uganda

The World Health Organisation said today it was investigating reports of suspected cases of the previously eradicated disease smallpox in eastern Uganda.

Smallpox is an acute contagious disease and was one of the world's most feared sicknesses until it was officially declared eradicated worldwide in 1979. "WHO takes any report of smallpox seriously" Gregory Hartl, a spokesman for the Geneva-based United Nations health agency, told Reuters via email. "WHO is aware of the reports coming out of Uganda and is taking all the necessary measures to investigate and verify." [...]

## Exemple : répétitions

WHO checks **smallpox reports** in Uganda

The World Health Organisation said today it was **investigating reports** of suspected cases **of the previously eradicated disease smallpox** in eastern Uganda.

**Smallpox** is an acute contagious **disease** and was one of the **world's** most feared sicknesses until it was officially declared **eradicated** worldwide in 1979. "**WHO** takes any **report of smallpox** seriously" Gregory Hartl, a spokesman for the Geneva-based United Nations **health** agency, told Reuters via email. "**WHO** is aware **of the reports** coming out of **Uganda** and is taking all the necessary measures to **investigate** and verify." [...]

## Exemple : répétitions + filtrage lexical maladie/lieu

WHO checks **smallpox reports** in Uganda

The World Health Organisation said today it was **investigating reports** of suspected cases **of the previously eradicated disease smallpox** in eastern Uganda.

**Smallpox** is an acute contagious **disease** and was one of the **world's** most feared sicknesses until it was officially declared **eradicated** worldwide in 1979. "**WHO** takes any **report of smallpox** seriously" Gregory Hartl, a spokesman for the Geneva-based United Nations **health** agency, told Reuters via email. "**WHO** is aware **of the reports** coming out of **Uganda** and is taking all the necessary measures to **investigate** and verify." [...]

## Un exemple en polonais

## Un exemple en russe

Document language : ru

**ВОЗ: свиным гриппом больны 4379 человек в 29-ти странах**

ВОЗ: свиным гриппом больны 4379 человек в 29-ти странах

Нью-Йорк. Количество случаев заболевания свиным гриппом (грипп А/Н1Н1), по данным на 10 мая, увеличилось до 4 тыс. 379. Об этом говорится в сообщении, размещенном на сайте Всемирной организации здравоохранения (ВОЗ).

Заболевание зафиксировано в 29-ти странах. Количество случаев с летальным исходом достигло 49 (45 - в Мексике, два - в США, один - в Канаде, один - в Коста-Рике).

По данным на 9 мая в мире насчитывалось 3 тыс. 440 случаев заболевания свиным гриппом. Большинство заболевших - в Мексике и США.

Нынешняя эпидемия гриппа, начавшаяся в Мексике и США, вызвана мутировавшим вирусом гриппа типа А. У заболевших повышается температура, появляются кашель, насморк, головная и мышечная боль, в некоторых случаях отмечаются рвота и диарея.

Помощник гендиректора ВОЗ по вопросам безопасности в области здравоохранения и окружающей среды Кеджи Фукуда заявил о том, что треть населения Земли может заразиться гриппом А/Н1Н1 в случае пандемии.

Однако ВОЗ приняла решение не повышать с пятого до шестого уровень угрозы пандемии гриппа А/Н1Н1, несмотря на распространение заболевания. Исполняющая обязанности директора программы ВОЗ по контролю за распространением гриппа в мире Сильвия Бриан, отметила, что большинство заразившихся вирусом А/Н1Н1 «привезли» инфекции из Мексики или находились в тесном контакте с теми, кто заразился во время поездки по этой стране. «Мы пока сохраняем пятый уровень угрозы. У нас нет доказательств».

UNCAEN université de Caen Basse-Normandie

## Un exemple en arabe

Document language : ar

**تحصين 72 ألف طالب وطالبة ضد الحصبة في بيشة**

تضمنت لجنة التحصين ضد الحصبة في بيشة اجتماع اللجنة التحصين ضد الحصبة في بيشة

بيشة : عبدالله المعاري

انطلقت أمس في محافظة بيشة الحملة الوطنية للتحصين ضد الحصبة في مراحليها الأولى مستهدفة أكثر من 72 ألف طالب وطالبة في جميع المراحل الدراسية الابتدائية والمتوسطة والثانوية الحكومية والخاصة ، وقد كانت مهمة بيشة استهدافها لتفادي المرحلة الأولى بواسطة فريق من المراكز الصحية لكل إقطاع حيث يغطي كل مركز صحي المدارس الواقعة في منطقة عمارة ، وتحت إشراف اللجنة العليا للحملة الوطنية للتحصين ضد الحصبة و الحصة الإجمالية والكاف بالمحافظة اجتماع في وقت سابق بقرعة بيشة بحضور بنو الشؤون الصحية بالمحافظة الدكتور عبدالله صالح الأعرجي

أوضح بذلك الناطق الإعلامي بصحة بيشة عبدالله سعيد العمادي مبيناً أنه قد تم التنسيق مع تعليم بيشة لتجهيز طاقم خاصة للتطعيم في كل مدرسة ، على أن يشارك في تنفيذ المرحلة الأولى التي ستبدأ على مدار خمسة أسابيع ، إدارة التربية والتعليم وفرع جامعة الملك خالد والنوون الاجتماعية ومركز صحي قوى الأمن إضافة إلى مشاركة القطاع الصحي الخاص في المحافظة

مختصات:عائشة

قمر هذا الموضوع 123456

هذا الخبر من موقع جريدة (بيشة اليومية)

UNCAEN université de Caen Basse-Normandie

## Résultats

### ► Résultats proches de l'état de l'art...

Langue	Taille Corpus	Rappel	Précision
Français	1954	92%	84%
Anglais	540	97%	84%
Russe	400	88%	85%
Polonais	439	85%	73%
Chinois	100	92%	85%

### ► ... pour la minorité de langues qui y sont traitées

Langue	Anglais	Français	Russe	Polonais	Chinois
MultiPULS	84%	84%	85%	73%	85%
Biocaster	93%	n/a	n/a	n/a	n/a

## Bilan

- Contexte : Thèse de Gaël Lejeune
  - Encadrante principale : Nadine Lucas
  - Collaboration avec l'université d'Helsinki (PICS Multipuls 2009–2011)
- Une chaîne de traitement alingue complète
- Un faible coût en ressources
  - lexique d'environ 500 termes par langue
  - 100 fois inférieur à l'état de l'art
- Un temps de traitement satisfaisant
  - 1000 documents par minute
  - 10 fois plus rapide que PULS
- Une plateforme d'annotation et de validation en ligne

## Perspective : détection de fraîcheur

- Vers une veille indépendante du domaine
  - Pas de ressources → indépendance du domaine
  - Le choix éventuel d'un domaine se fait *a posteriori*
- Détection d'associations lexicales fortement divergentes
- Thèse d'Oskar Gross
  - Encadrant principal : Hannu Toivonen (U. Helsinki)
  - Projet d'ANR Sucret - *Supporting CREativity from Text*
  - Semestre à Helsinki en cours

## Exploitation de la structure des documents

- Utiliser la structure logique des documents (e.g., sections, sous-sections paragraphes), stylistique (e.g., gras, italique)
- De nouvelles possibilités :
  - Répondre à un besoin d'information par des fragments de documents précis, en tirant profit de leur structure logique
  - Affiner les réponses par granularité structurelle
- La notion de granularité s'ajoute à celle de pertinence
- Deux contributions :
  - RI structurée
  - Partitionnement de documents

## Système EXTIRP

EXacT coverage IR based on static Passage clusters

1. Détection des unités minimales de recherche
2. Calcul de leur pertinence
3. Propagation verticale dans l'arbre XML du document
4. Sélection du grain adéquat

Metric	@1	@5	@10	@100	@1500
ncXG (strict)	8	12	15	34	27
ncXG (generalised)	13	26	29	32	22
inex_eval (strict)	1	2	11	26	22
inex_eval (generalised)	11	19	25	31	24

[rangs INEX 2005 (sur 44)]

## Système EXTIRP : amélioration

- ▶ Notion de « *textitude* » d'un marqueur XML
  - ▶ T/E (*Text node/Element node*) ratio
  - ▶ Différencie les marqueurs XML de structuration de ceux de mise en forme
  - ▶ Duplication en ligne des textes mis en forme

...kernel trick has been applied to several algorithms in  
 <link>machine learning</link> <link>machine learning</link> and  
 <link>statistics</link> <link>statistics</link>, including...

- ▶ Utile notamment pour l'extraction de séquences !

	Original	all Dupl.	Hybrid
iP 0.01 (107 topics)	0.3319	0.3773	<b>0.3815</b>
MAiP (107 topics)	0.0912	0.1024	<b>0.1036</b>

[résultats INEX 2007 – SIGIR 2008]

## Partitionnement de documents XML

- ▶ But : prendre en compte la structure dans le processus de classification
  - ▶ Génération automatique de DTDs
  - ▶ *Cluster hypothesis (a priori)*
- ▶ Méthode de clustering en 2 étapes
  1. Clustering basé sur des descripteurs structurels exclusivement : Détection des « outliers » structurels
  2. Pour les clusters dont la similarité interne est inférieure à un seuil : second clustering basé sur des descripteurs textuels classiques

Features	Text	Tags	Tags + Text	Tags → Text
Entropy	0.633	0.798	0.678	<b>0.630</b>
Purity	0.379	0.228	0.372	<b>0.394</b>
Clustering Time	754s.	11s.	837s.	<b>11+742s.</b>

## Contributions

- ▶ Généralisme : aucune connaissance préalable requise
  - ▶ fonctionne sans DTD ou *Schema XML*
- ▶ RI structurée – Projet *EXTIRP*
  - ▶ Mené à l'université d'Helsinki 2002–2003 (6 personnes)
  - ▶ Participation annuelle à INEX jusqu'à 2009
- ▶ Classification non supervisée
  - ▶ Système vainqueur *INEX Mining track 2006* (Wikipedia)
- ▶ Exploitation de la *textitude* des nœuds XML
  - ▶ Thèse de Miro Lehtonen (2006)

## Perspectives

- ▶ RI structurée
  - ▶ RI dans des bibliothèques numériques/numérisées
- ▶ Classification non supervisée
  - ▶ Classification de fragments de documents
  - ▶ Résumé multi-documents

## Travaux en méthodologie d'évaluation

- ▶ Objectif
  - ▶ Évaluation de tâches sur des collections d'ouvrages numérisés
  - ▶ Deux contributions :
    - ▶ Évaluation de la performance des SRI
    - ▶ Évaluation de la performance en extraction de structure



## Contexte

- ▶ Contexte : *INEX Book Search Track*
- ▶ Un cadre d'application clé
  - ▶ Accès à des collections d'ouvrage en ligne
    - ▶ Trouver un livre
    - ▶ Trouver de l'information dans des livres
  - ▶ Numérisation massive
    - ▶ *Google Books*
    - ▶ Bibliothèques nationales
  - ▶ Apporter le livre ancien aux *e-books*, liseuses électroniques



## Numérisation

1. Livre (papier)
2. Photos des pages (.jpg)
3. OCR (.xml)
  - ▶ Bruit
  - ▶ Nécessite des standards



Kirtas book scanner

## Difficultés spécifiques

- ▶ Grandes collections de grands documents
- ▶ Structure physique plutôt que logique
- ▶ Absence de références croisées
- ▶ Bruit (livres anciens ← OCR)



## Évaluer la performance des SRI

- ▶ Défi principal : collection des annotations
  - ▶ Développement d'une plateforme
    - ▶ <http://www.booksearch.co.uk>
  - ▶ Compétition d'annotateurs (jeu)
  - ▶ *Crowdsourcing*



## Évaluer la performance de l'extraction de structure

- ▶ Aucune méthodologie adéquate existante
  - ▶ Définition d'un cadre d'évaluation complet
  - ▶ Définition de mesures d'évaluation
  - ▶ Construction de la vérité de terrain
    - ▶ Processus d'annotation collaborative avec plateforme Java distribuée



## Contributions

- ▶ *INEX Book Search track*
  - ▶ Avec Gabriella Kazai (*Microsoft Research Cambridge*)
  - ▶ Atelier de travail annuel pour les participants
- ▶ Mise en place de méthodologies d'évaluation :
  - ▶ RI dans des collections d'ouvrages
    - ▶ Initié en 2007
    - ▶ Tâche principale d'INEX depuis 2011
    - ▶ Intégré à CLEF depuis 2012
  - ▶ Extraction de structure (ouvrages numérisés)
    - ▶ Initié à INEX 2008
    - ▶ Compétition ICDAR depuis 2009



## Perspectives

- ▶ RI dans des collections d'ouvrages
  - ▶ RI sociale
  - ▶ Exploitation de résumés et de recommandations
- ▶ Extraction de structure (ouvrages numérisés)
  - ▶ *Crowdsourcing*
  - ▶ Évaluation indirecte



## Bilan

### PIPELINE DE TRAITEMENT

- ▶ Extraction
- ▶ Exploitation
- ▶ Évaluation

### APPROCHE DU MATÉRIAU TEXTE

- ▶ Séquentielle
- ▶ Discursive
- ▶ Structurale

CONSTANTE : GÉNÉRICITÉ, NOTAMMENT MULTILINGUISME

## Perspectives

- ▶ Détection de fraîcheur
  - ▶ Détection d'associations lexicales fortement divergentes
  - ▶ Généralisation de la veille multilingue
    - ▶ Le choix éventuel d'un domaine se fait *a posteriori*
- ▶ RI contextuelle structurée (personnalisée)
  - ▶ RI structurée sur des collections d'ouvrages
  - ▶ Suggestions adaptées au contexte
- ▶ Recherche temporelle d'images
  - ▶ Datation automatique de photographies
    - ▶ corrélation entre contenu des images, date de prise de vue
    - ▶ environnement textuel et journaux de requêtes