

وزارة التعليم العالي والبحث العلمي

BADJI MOKHTAR UNIVERSITY-ANNABA  
-UNIVERSITE BADJI MOKHTAR-ANNABA-



جامعة باجي مختار  
- عنابة -

Faculté : Sciences de l'Ingénieur

Année 2008

Département : Informatique

## MÉMOIRE

Présentation en vue de l'obtention du diplôme de magister

### Recherche et filtrage d'information basés sur le Text Mining sous la technologie GRID

Option : Texte Parole et Image

Présenté par : BOUALI Ali

DEVANT LE JURY

Pr BADACHE Nadjib, Professeur à l'USTHB Alger, Président

Pr LASKRI Mohamed Tayeb, Professeur à l'Université de Annaba, Rapporteur

Dr FARAH Nadir, Maitre de Conférence à l'Université de Annaba, Examineur

Dr BOUHADADA Tahar, Maitre de Conférence à l'Université de Annaba, Examineur

# Résumé

La quantité d'information textuelle augmente de façon exponentielle aussi bien comme archives que documents de travail dans les organisations académiques, dans les administrations et dans les entreprises.

Il est très difficile aux engins de recherche de retrouver l'information adéquate dans cet important volume de données textuelle, reflétant les intérêts des utilisateurs qui changent avec le temps. Nous avons donc besoin de techniques d'apprentissage permettant de reconnaître les intérêts des utilisateurs en donnant uniquement une simple requête, et de filtrer l'information disponible suivant ces intérêts, et donc de mettre en place un système de recherche et de filtrage d'informations.

Dans ce travail, on propose une architecture du système de recherche et filtrage d'information basé sur le Text Mining sous la technologie GRID, capable de constituer et consulter une base de données d'une part, et de filtrer l'information disponible suivant les besoins des utilisateurs d'une autre part.

Le mécanisme de filtrage est à la base du Text Mining, pour le quel on a appliqué les techniques du GRID pour améliorer les performances de calcul, il réalise la modélisation des intérêts des utilisateurs et, le filtrage d'information; en effet, ces informations sont d'abord proposés par la rétroaction de l'utilisateur, et des mots sont ensuite sélectionnés pour créer un profil. En utilisant ce profil, l'information entrante est filtrée et ainsi, plus d'informations pertinentes sont présentée à l'utilisateur associé.

L'architecture du système de recherche et filtrage d'information basé sur le Text Mining sous la technologie GRID (RFITM-GRID) que nous avons proposé se compose de plusieurs modules :

**Un module d'indexation**, qui, en balayant un ou plusieurs fichiers, il construit un index avec les mots trouvés dans les documents.

**Un module moteur de recherche**, qui permet de réaliser une recherche d'information en consultant les fichiers de la base de donnée créer par le module d'indexation.

**Un module de filtrage**, qui, à base de la rétroaction de pertinence et du Text Mining, en exploitant les hautes performances du GRID, crée des profils, qu'il exploite par la suite pour filtrer l'information entrante.

## **Mots clés**

Recherche d'information, Grille de calcul (GRID), Indexation, Rétroaction de pertinence, filtrage d'information, Profil utilisateur, découverte des connaissances dans les bases de données (DCBD), Text Mining.

# Abstarct

The amount of textual information increases exponentially as well as archive papers in academic organizations, government and enterprise.

It is very difficult for search engines to find adequate information in this important volume of text data, reflecting the interests of users change with time. We need learning techniques to recognize users' interests by giving only a simple query and filter information available following these interests, and thus to establish a system of search and filtering information .

In this work, it proposes a system architecture search and filtering information based on Text Mining in the GRID, able to constitute and consult a database on the one hand and filters information available following user needs on another hand.

The filtering mechanism is the basis of Text Mining, for which we applied the GRID technology to improve performance calculates, it models the user's interests and achieves the filtering information, because this information are first proposed by feedback from the user, and key terms are then selected to create a profile. By using this profile, the incoming information is filtered and thus more relevant information is displayed to the user partner.

The system architecture search and filtering information based on Text Mining in the GRID technology (RFITM-GRID) that we have proposed consists of several modules:

**A module indexing**, this, by browsing one or several files, built an index with words found in the documents.

**A search engine module**, which allows a search for information by consulting the files of the database created by the indexing module.

A **filtering module**, which, based on the relevance feedback and Text Mining, exploiting the high performance GRID, create profiles witch will be used to filter the new information.

**Key words**

Information retrieval, GRID, Indexing, Relevance feedback, information filtering, User profile, Knowledge discovering in data base (KDD), Text Mining.

# ملخص

عدد معلومات نصية يزيد أضعافا مضاعفة فضلا عن أوراق في أرشيف الأكاديمية والمنظمات والمؤسسات الحكومية.

ومن الصعب جدا لمحركات البحث للعثور على معلومات كافية في هذا النص المهم من حجم البيانات، والتي تعكس مصالح المستخدمين المتغيرة مع الوقت. ونحن بحاجة إلى تعلم تقنيات الاعتراف مصالح المستخدمين بإعطاء مجرد استفسار والمعلومات المتاحة فلتر بعد هذه المصالح، وبالتالي لإنشاء نظام للبحث وترشيح المعلومات.

في هذا العمل، نقترح هيكل نظام البحث وترشيح المعلومات على أساس نص التعدين في غريد، وتكون قادرة على خلق والتشاور قاعدة بيانات من جهة و ترشيح المعلومات المتاحة وفقا لاحتياجات المستخدمين من جهة أخرى.

فإن آلية الترشيح هو أساس نص التعدين، عليه نطبق تكنولوجيا غريد لتحسين أداء الحساب، يؤدي نمذجة المستخدمين، وترشيح المعلومات، لأن هذه المعلومات هي الأولى التي اقترحها التغذية المرتدة من المستعملين، واختيار الكلمات بعد ذلك لخلق ملف. باستخدام هذا الملف، فإن المعلومات الواردة ترشح وبالتالي المزيد من المعلومات ذات الصلة تعرض إلى المستخدم المرتبطة به.

بنية النظام والبحث وترشيح المعلومات على أساس نص التعدين في غريد (RFITM-GRID) الذي اقترحناها تتكون من عدة وحدات:

**وحدة الفهرسة**، التي، عن طريق المسح واحدة أو أكثر من ملفات، انه يبني فهرس بالكلمات التي تم العثور عليها في وثائق.

وحدة محرك البحث ، الذي يتيح البحث عن المعلومات من خلال التشاور مع ملفات قاعدة البيانات التي أنشأتها وحدة الفهرسة.

وحدة الترشيح ، التي ، على أساس التغذية المرتدة و نص التعدين ، واستغلال الأداء العالي للغريد ، وخلق ملفات ، تعمل بعد ذلك لترشيح المعلومات الواردة .

## الكلمات الدالة

استرجاع المعلومات ، الحوسبة الشبكية ، الفهرسة ، وأهمية التغذية المرتدة ، وترشيح المعلومات ، ملف المستخدم ، واكتشاف المعرفة في قواعد البيانات ، نص التعدين.

# Table des Matières

Introduction Générale.....	1
<b>CHAPITRE 1</b> LES SYSTEMES DE RECHERCHE ET DE FILTRAGE D'INFORMATIONS.....	4
I. INTRODUCTION:.....	5
II. APPROCHE TRADITIONNELLE DES SYSTEMES DE RECHERCHE D'INFORMATION :.....	5
II.1. STOCKAGE :.....	5
II.2. INDEXATION :.....	6
II.3. RECHERCHE :.....	8
III. UNE PRISE EN COMPTE PLUS FINE DES SPECIFICITES DE L'UTILISATEUR :.....	9
III.1. LA RETROACTION DE PERTINENCE : [RF97].....	10
III.2. LE FILTRAGE D'INFORMATION : [NOU ?], [KS00].....	10
IV. RECAPITULATIF : [NOU ?].....	14
V. CONCLUSION :.....	15



<b>CHAPITRE 2 APERÇU SUR LA FOUILLE DE TEXTE (TEXT MINING)</b> .....	16
<b>I. INTRODUCTION:</b> .....	17
I.1. DECOUVERTE DE CONNAISSANCES (KNOWLEDGE DISCOVERY):.....	17
I.2. DATA MINING, APPRENTISSAGE AUTOMATIQUE ET APPRENTISSAGE STATISTIQUE:.....	19
I.3. DEFINITION DE LA FOUILLE DE TEXTE (TEXT MINING):.....	19
I.4. DOMAINES RELATIFS DE RECHERCHES: .....	20
<b>II. CODAGE DU TEXTE:</b> .....	21
II.1. TRAITEMENT DU TEXTE: .....	21
II.1.1. FILTRAGE, LEMMATISATION ET STEMMING:.....	22
II.1.2. SELECTION DES TERMES D'INDEX:.....	23
II.2. LE MODELE D'ESPACE VECTORIEL: .....	23
II.3. TRAITEMENT LINGUISTIQUE: .....	25
<b>III. METHODES DE DATA MINING POUR LE TEXTE:</b> .....	26
III.1. CLASSIFICATION:.....	26
III.2. REGROUPEMENT (CLUSTERISATION):.....	32
III.2. EXTRACTION D'INFORMATION: .....	39
III.4. AUTRES DOMAINES D'APPLICATION:.....	41
<b>IV. APPLICATIONS:</b> .....	41
IV.1. ANALYSE DE BREVETS: .....	41
IV.2. BIO-INFORMATIQUE: .....	42
IV.3. FILTRAGE D'EMAIL, ANTI-SPAM:.....	42
<b>V. CONCLUSION:</b> .....	42

<i>CHAPITRE 3 ETUDE DES GRILLES INFORMATIQUES</i> .....	43
I. INTRODUCTION : .....	44
II. DEFINITIONS: .....	44
III. PRINCIPES COMMUNS DES GRILLES: .....	45
IV. HISTORIQUE : .....	46
IV.1. LES ORIGINES DES GRILLES INFORMATIQUES : .....	46
IV.2. EXEMPLE SUR L'ANALOGIE DES RESEAUX DE DISTRIBUTION ELECTRIQUE ET DES GRILLES INFORMATIQUES: .....	47
IV.3. DEVELOPPEMENT ET EVOLUTION: .....	48
V. LES GRILLES DE CALCUL : .....	49
V.1. PRINCIPES DE BASE DU GRID COMPUTING : .....	49
V.1.2. AVANTAGES : .....	50
V.1.3. INCONVENIENTS : .....	51
V.2. LA TECHNOLOGIE « MASSIVE GRID » : .....	51
V.3. LA TECHNOLOGIE « DESKTOPGRID » OU « INTERNET COMPUTING » : .....	53
VI. LES DIFFERENTS TYPES DE GRILLE INFORMATIQUES : .....	55
VII. QUELQUES INTERGICIELS (MIDDLEWARE) : .....	56
VII.1. INTERGICIELS «MASSIVE GRID» : .....	56
VII.2. INTERGICIELS «DESKTOPGRID» : .....	57
VIII. DOMAINES D'APPLICATION DES GRILLES INFORMATIQUES: .....	57
VIII.1. LE DOMAINE INDUSTRIEL : .....	57
VIII.2. LE DOMAINE SCIENTIFIQUE : .....	58
IX. QUELQUES PROJETS DE GRILLES INFORMATIQUES : .....	59
X. PRESENTATION DU GLOBUS TOOLKIT : .....	60
XI. CONCLUSION : .....	65

CHAPITRE 4 ARCHITECTURE DU SYSTEME RFITM-GRID .....	66
I. INTRODUCTION : .....	67
II. L'ENVIRONNEMENT GRID : .....	67
III. PROCESSUS DE CONCEPTION : .....	69
III.1. DESCRIPTION GENERALE .....	69
III.2. LA THEMATISATION : .....	70
IV. ARCHITECTURE DU SYSTEME : .....	72
VI.1. SYSTEME DE GESTION DES DONNEES (DATA MANAGEMENT SYSTEM, DMS) : .....	73
IV.2. LES ENVIRONNEMENTS CLIENTS ET SERVEUR : .....	73
VII. L'INDEXEUR : .....	76
VIII. LE MOTEUR DE RECHERCHE : .....	78
IX. LE GESTIONNAIRE DES PROFILS ET FILTRAGE : .....	79
X. SOUMISSION DE JOBS : .....	87
XI. CONCLUSION : .....	88
Conclusion générale.....	91
Bibliographie.....	94