

Extraction d'entités dans des collections évolutives

Thierry Despeyroux Eduardo Fraschini
Anne-Marie Vercoustre

Inria, Rocquencourt

EGC 2007, Namur

Au menu...

- 1 Le problème concret
- 2 Méthode utilisée
- 3 Résultats
- 4 Conclusion et perspectives

Le rapport d'activité de l'Inria



- **Annexe technique**
- 168 rapports d'équipes en 2005
- 4611 pages en 2005
- en anglais
- annuel

Le rapport d'activité de l'Inria



- Annexe technique
- 168 rapports d'équipes en 2005
- 4611 pages en 2005
- en anglais
- annuel

Le rapport d'activité de l'Inria



- Annexe technique
- 168 rapports d'équipes en 2005
- 4611 pages en 2005
- en anglais
- annuel

Le rapport d'activité de l'Inria



- Annexe technique
- 168 rapports d'équipes en 2005
- 4611 pages en 2005
- en anglais
- annuel

Le rapport d'activité de l'Inria



- Annexe technique
- 168 rapports d'équipes en 2005
- 4611 pages en 2005
- en anglais
- annuel

Question

- Quels sont les partenaires de l'Inria ?

Réponse... dans le RA

Project Team ACACIA

19

Then we proposed an ontology-based approach for building an information system supporting technology searching implemented by agents. This system facilitates the document searching and automatic track of the watched. Its agents use the ontology to enrich any watcher's query, then formulate a system query and send it to Google to search the Web and finally generate annotations from search result. Thus the watcher can easily access to information associated by exploiting the Chinese semantic search engine.

To do so, we developed and experimented three algorithms using the ontology to search the Web with Google and then generate the RDF annotations from their results of Google automatically. The two first algorithms are formalized of concepts in ontology to search the Web while the third one relies on the balanced situation of downloaded concepts of user's concepts in the original query [10], [11]. This work will be published as *RTA2004* [12].

We are now designing and implementing a sub-library of "annotator" Agents encapsulating this algorithm, working in cooperation with other agents distributed in other tools in the TM system.

7. Contracts and Grants with Industry

7.1. Knowledge Management Platform

Participants: Asian Online Group, L'Oréal Cosmetics, Kawan Databases, Palens Graduate, Thierry Choukret, Nicolas Choukret, Olivier Chagnat

This two-year project started in April 2005. A first prototype was delivered, with the final report of the project [13]. Koaf involved teams specialized in computer science, electronic sciences, management sciences, ergonomics and psychology, namely: Radgo Laboratory (LRSIA-CNRS), Lingua Laboratory (LRSIA-CNRS), Asian Team (INRIA Sophia Antipolis, CSR, Institut Pasteur and 2012 Strategic, Technical Value Association (Sophia Antipolis). The project also involved a set of stake actors, who actively participate to the design of the prototype. The Koaf project led to the construction of a web server facilitating the sharing of competences within a community – the Telecom Value (Sophia Antipolis) which gathers firms, local institutions, and research organizations working in the telecommunications domain. The aim of Koaf is to promote particularly useful and cutting activities the community will [14]. The Asian team coordinated the design and development of the prototype and of its underlying ontologies. The prototype is based on Google, hence its name Koaf-Ceasar. The ROSET Koaf project was to assess that it will continue through follow-up projects.

noted by Semantic Systems.

7.2.

8.4. International Actions

8.4.1. WIC

Palens Choukret is a member of the Semantic Web Best Practices and Deployment Working group and experts on the activities of this group in the INRIA DuSRI. This working group discusses issues and delivers reports on best practices for different aspects of the semantic web: schema engineering, design patterns, vocabulary management, natural language applications and datasets, etc. Palens Choukret participated in the technical advisory sessions of WIC (March 2007) and regularly participates to bi-weekly teleconferences. We also prepared a contribution to the WIC Workshop on Semantic Web Services (June 2007); it showed our experience in applying Ceasar to cooperate semantic web services and stated our position as a stake in our system shared by the first group in building standards for semantic web services.

8.4.2. Carnegie Mellon University

Palens Choukret continued his collaboration with CMU.

8.4.3. University Goshaw Energy

The Asian team visited Maxime Le (University of Goshaw Energy of St. Louis, Stregal), in addition this visit included an exchange between ACACIA and the Computer Science Department of the University of

7.3.

• Les parties “contrats” et “international” sont les plus pertinentes

Réponse... dans le RA, oui mais...

18

Activity Report INRIA 2005

6. Contracts and Grants with Industry

6.1. Industrial Contracts

The Algorithm Project and [Wanted-Magic Inc. \(WMI\)](#) have developed a collaboration based on reciprocal interests. It is obviously interesting for the company to integrate functionalities at the forefront of the current research in computer algebra. Reciprocally, this integration makes our program and our research visible to a very wide audience.

Technical exchanges have thus taken place between the project and the company over the years. After more than 3 years with the project, J. Curien has been for around seven years Product Development Director at [WMI](#), before going back to the academic world. Similarly, B. Moury, who worked for two years at the project developing the software tool packages is now working at [WMI](#).

Thanks to all this activity, the company [WMI](#) considers Inria as a special partner and grants it a free license for all of its research units. Moreover, a cooperation agreement has been signed between [WMI](#) and Inria in 2001. In particular, one of the objectives is to replace all the realizers dealing with asymptotic series

agprojects

7. Contracts and Grants with Industry

7.1. Microsoft Research (2003-2006)

Participants: Jean-Benoît, Patrick, Valérie.

The objective is to contribute to the development of the AMBA model management framework and foster the dissemination of our results in Open Source. We were under a non-restrictive license. In particular, we are adapting the AMBA framework to the principles and tools of the Microsoft Software Project approach (Visual Studio 2005, Team Explorer). Activities leading to tasks in ACT, AMBA, AMBY should be made available to the Microsoft environment with the help of technical space projects.

7.2. IBM Eclipse (2004-2005)

Participants: Jean-Benoît.

The objective was to port the ACT platform to the Eclipse Open Source environment. This was the only French project granted by [IBM](#) (Eclipse Release and Client in 2004). A first version of the prototype has been presented at the OOTEA conference in October 2004 in Vancouver.

7.3. RNTI Modathèque (2004-2005)

Participants: Jean-Benoît, Patrick, Valérie.

In this project, we work with [Thales](#) (project leader), [France Telecom](#) ([RNTI](#), [L3E](#)) and software engineering tool vendors in France. The objective is to define the conceptual and practical basis for model engineering. In particular, using components of the Model Driven Environment (MDE) of the OMO. In this project, we use our ACT platform for MDE components.

7.4. Coroll Motor (2003-2006)

Participants: Jean-Benoît, Patrick, Valérie.

- Style très hétérogène
- Souvent peu rédigé
- Sigles plus ou moins développés
- Noms de réseaux, de labos, de localisation...

Réponse... dans le RA, oui mais...

18

Activity Report INRIA 2005

6. Contracts and Grants with Industry

6.1. Industrial Contracts

The Algorithm Project and [Wanted-Magic Inc. \(WMI\)](#) have developed a collaboration based on mutual interests. It is obviously interesting for the company to integrate functionalities at the forefront of the current research in computer algebra. Conversely, this integration makes our program and our research visible to a very wide audience.

Financial exchanges have taken place between the project and the company over the years. After more than 3 years with the project, J. Curien has been for around seven months Development Director at [WMI](#), before going back to the academic world. Similarly, B. Morley, who worked for two years at the project, developing the software part package is now working at [WMI](#).

Thanks to all this activity, the company [WMI](#) considers Inria as a special partner and grants it a free license for all of its research units. Moreover, a cooperation agreement has been signed between [WMI](#) and Inria in 2001. In particular, one of the objectives is to replace all the realizers dealing with asymptotic series

agprojects

7. Contracts and Grants with Industry

7.1. Microsoft Research (2003-2006)

Participants: Jean-Benoît, Patrick Valduriez.

The objective is to contribute to the development of the AMBA model management framework and foster the dissemination of our results in Open Source. We were under a non-restrictive license. In particular, we are adapting the AMBA framework to the principles and tools of the Microsoft Software Project approach (Visual Studio 2005, Team Explorer). Activities leading to tasks in ACT, ABEL, AMBY should be made available to the Microsoft environment with the help of technical space projects.

7.2. IBM/Eclipse (2004-2005)

Participants: Jean-Benoît.

The objective was to port the ACT platform to the Eclipse Open Source environment. This was the only French project granted by [IBM/Eclipse Research Cloud](#) in 2004. A first version of the prototype has been presented at the OOTEA conference in October 2004 in Vancouver.

7.3. RNTI, Modathèque (2004-2005)

Participants: Jean-Benoît, Patrick Valduriez.

In this project, we work with [Thales XE](#) (project under), [France Telecom R&D](#), [L3H](#) and software engineering tool vendors in France. The objective is to define the conceptual and practical basis for model engineering, in particular, using components of the Model Driven Environment (MDE) of the OMO. In this project, we use our ACT platform for MDE components.

7.4. Coroll Motor (2003-2006)

Participants: Jean-Benoît, Patrick Valduriez.

- Style très hétérogène
- Souvent peu rédigé
- Sigles plus ou moins développés
- Noms de réseaux, de labos, de localisation...

Réponse... dans le RA, oui mais...

18

Activity Report INRIA 2005

6. Contracts and Grants with Industry

6.1. Industrial Contracts

The Algorithm Project and [Wanted-Magic Inc. \(WMI\)](#) have developed a collaboration based on mutual interest. It is obviously interesting for the company to integrate functionalities at the forefront of the current research in computer algebra. Conversely, this integration makes our program and our research visible to a very wide audience.

Technical exchanges have taken place between the project and the company over the years. After more than 3 years with the project, J. Currie has been for around 6 years Product Development Director at [WMI](#), before going back to the academic world. Similarly, B. Morley, who worked for two years at the project developing the software test package is now working at [WMI](#).

Thanks to all this activity, the company [WMI](#) considers Inria as a special partner and grants it a free license for all of its research units. Moreover, a cooperation agreement has been signed between [WMI](#) and Alcatel in 2001. In particular, one of the objectives is to replace all the realizers dealing with asymptotic series

agprojects

7. Contracts and Grants with Industry

7.1. Microsoft Research (2003-2006)

Participants: Jean-Benoît, Patrick Valdurio.

The objective is to contribute to the development of the AMBA model management framework and foster the dissemination of our results in Open Source. We were under a non-restrictive license. In particular, we are adapting the AMBA framework to the principles and tools of the Microsoft Software Project approach (Visual Studio 2005, Team Explorer). Activities leading to tasks in ACT, AMBY should be made available to the Microsoft environment with the help of technical space projects.

7.2. IBM/Eclipse (2004-2005)

Participants: Jean-Benoît.

The objective was to port the ACT platform to the Eclipse Open Source environment. This was the only French project granted by [IBM/Eclipse Research Cloud](#) in 2004. A first version of the prototype has been presented at the OOTEA conference in October 2004 in Vancouver.

7.3. RNTL Modathèque (2004-2005)

Participants: Jean-Benoît, Patrick Valdurio.

In this project, we work with [Thales XE](#) (project leader), [France Telecom R&D](#), [L3H](#) and software engineering tool vendors in France. The objective is to define the conceptual and practical basis for model engineering, in particular, using components of the Model Driven Environment (MDE) of the OMO. In this project, we use our ACT platform for MDE components.

7.4. Coroll Motor (2003-2006)

Participants: Jean-Benoît, Patrick Valdurio.

- Style très hétérogène
- Souvent peu rédigé
- Sigles plus ou moins développés
- Noms de réseaux, de labos, de localisation...

Réponse... dans le RA, oui mais...

18

Activity Report INRIA 2005

6. Contracts and Grants with Industry

6.1. Industrial Contracts

The Algorithm Project and [Wanted-Magic Inc. \(WMI\)](#) have developed a collaboration based on mutual interest. It is obviously interesting for the company to integrate functionalities at the forefront of the current research in computer algebra. Reciprocally, this integration makes our program and our research visible to a very wide audience.

Financial exchanges have taken place between the project and the company over the years. After more than 3 years with the project, J. Curien has been for around 3 years Product Development Director at [WMI](#), before going back to the academic world. Similarly, B. Morley, who worked for two years at the project, developing the software test package now working at [WMI](#).

Thanks to all this activity, the company [WMI](#) considers Inria as a special partner and grants it a free license for all of its research units. Moreover, a cooperation agreement has been signed between [WMI](#) and Alcatel in 2001. In particular, one of the objectives is to replace all the realizers dealing with asymptotic series

agprojects

7. Contracts and Grants with Industry

7.1. Microsoft Research (2003-2006)

Participants: Jean-Benoît, Patrick-Vaiteara.

The objective is to contribute to the development of the AMBA model management framework and foster the dissemination of our results in Open Source. We were under a non-restrictive license. In particular, we are adapting the AMBA framework to the principles and tools of the Microsoft Software Project approach (Visual Studio 2005, Team Explorer), Automata being by tools as ACT, ABEL, ANAFY should be made available to the Microsoft environment with the help of technical space projects.

7.2. IBM/Eclipse (2004-2005)

Participants: Jean-Benoît.

The objective was to port the ACT platform to the Eclipse Open Source environment. This was the only French project granted by [IBM](#) (Eclipse Release and Client in 2004). A first version of the ported type has been presented at the OOTEA conference in October 2004 in Vancouver.

7.3. RNTL Modathèque (2004-2005)

Participants: Jean-Benoît, Patrick-Vaiteara.

In this project, we work with [Thierry PO](#) (project leader), [Francois Tesson](#), [Ralf L](#) and software engineering tool, modules in France. The objective is to define the conceptual and practical basis for model engineering, in particular, using components of the Model Driven Development (MDE) of the OMO. In this project, we use our ACT platform for MDE components.

7.4. Coroll Motor (2003-2006)

Participants: Jean-Benoît, Patrick-Vaiteara.

- Style très hétérogène
- Souvent peu rédigé
- Sigles plus ou moins développés
- Noms de réseaux, de labos, de localisation...

Extraction d'entités nommées

- Entités nommées : noms de personnes, d'organisations, de lieux, dates, valeurs monétaires...
- Extraction : tester l'existence ou trouver toutes les occurrences

Extraction d'entités nommées

- Entités nommées : noms de personnes, d'organisations, de lieux, dates, valeurs monétaires...
- Extraction : tester l'existence ou trouver toutes les occurrences

Existant

- Outils commerciaux : Rosette Entity Extractor, Inxight SmartDiscovery, Convera-RetrievalWare, Xerox...
- Large communauté, ex : univ. de Sheffield (ANNIE, GATE)
- Utilisent généralement des ressources linguistiques importantes (dictionnaires), et de grandes collections pour faire l'apprentissage (méthodes statistiques)
- Nécessitent une adaptation manuelle pour un corpus particulier (par exemple écriture de règles)
- Xeros : 250 règles manuelles pour extraire des entités biologiques

Existant

- Outils commerciaux : Rosette Entity Extractor, Inxight SmartDiscovery, Convera-RetrievalWare, Xerox...
- Large communauté, ex : univ. de Sheffield (ANNIE, GATE)
- Utilisent généralement des ressources linguistiques importantes (dictionnaires), et de grandes collections pour faire l'apprentissage (méthodes statistiques)
- Nécessitent une adaptation manuelle pour un corpus particulier (par exemple écriture de règles)
- Xeros : 250 règles manuelles pour extraire des entités biologiques

Existant

- Outils commerciaux : Rosette Entity Extractor, Inxight SmartDiscovery, Convera-RetrievalWare, Xerox...
- Large communauté, ex : univ. de Sheffield (ANNIE, GATE)
- Utilisent généralement des ressources linguistiques importantes (dictionnaires), et de grandes collections pour faire l'apprentissage (méthodes statistiques)
- Nécessitent une adaptation manuelle pour un corpus particulier (par exemple écriture de règles)
- Xeros : 250 règles manuelles pour extraire des entités biologiques

Existant

- Outils commerciaux : Rosette Entity Extractor, Inxight SmartDiscovery, Convera-RetrievalWare, Xerox...
- Large communauté, ex : univ. de Sheffield (ANNIE, GATE)
- Utilisent généralement des ressources linguistiques importantes (dictionnaires), et de grandes collections pour faire l'apprentissage (méthodes statistiques)
- Nécessitent une adaptation manuelle pour un corpus particulier (par exemple écriture de règles)
- Xeros : 250 règles manuelles pour extraire des entités biologiques

Existant

- Outils commerciaux : Rosette Entity Extractor, Inxight SmartDiscovery, Convera-RetrievalWare, Xerox...
- Large communauté, ex : univ. de Sheffield (ANNIE, GATE)
- Utilisent généralement des ressources linguistiques importantes (dictionnaires), et de grandes collections pour faire l'apprentissage (méthodes statistiques)
- Nécessitent une adaptation manuelle pour un corpus particulier (par exemple écriture de règles)
- Xeros : 250 règles manuelles pour extraire des entités biologiques

Idées de base

- Utiliser une méthode similaire aux méthodes de wrapping
- On remplace les balises HTML par les syntagmes du langage (nom, verbe, etc.)
- Apprentissage à partir de peu d'exemples
- Ne pas utiliser notre connaissance de la langue
- Réinjecter les résultats d'une année l'année suivante

Idées de base

- Utiliser une méthode similaire aux méthodes de wrapping
- On remplace les balises HTML par les syntagmes du langage (nom, verbe, etc.)
- Apprentissage à partir de peu d'exemples
- Ne pas utiliser notre connaissance de la langue
- Réinjecter les résultats d'une année l'année suivante

Idées de base

- Utiliser une méthode similaire aux méthodes de wrapping
- On remplace les balises HTML par les syntagmes du langage (nom, verbe, etc.)
- Apprentissage à partir de peu d'exemples
- Ne pas utiliser notre connaissance de la langue
- Réinjecter les résultats d'une année l'année suivante

Idées de base

- Utiliser une méthode similaire aux méthodes de wrapping
- On remplace les balises HTML par les syntagmes du langage (nom, verbe, etc.)
- Apprentissage à partir de peu d'exemples
- Ne pas utiliser notre connaissance de la langue
- Réinjecter les résultats d'une année l'année suivante

Idées de base

- Utiliser une méthode similaire aux méthodes de wrapping
- On remplace les balises HTML par les syntagmes du langage (nom, verbe, etc.)
- Apprentissage à partir de peu d'exemples
- Ne pas utiliser notre connaissance de la langue
- Réinjecter les résultats d'une année l'année suivante

Schémas de phrases

- ...by Texas Instruments because...
- IN NNP NNPS IN
- `<in>by</in>`
`<org><nnp>Texas</nnp>`
`<nnps>Instruments</nnps></org>`
`<in>because</in>`
- IN ~ NNPS* % ~ IN

Schémas de listes

16

Activity Report INRIA 2005

6.1.3. JCDecaux

- Number: Inria 401
- Title: Typed- π -calc.
- Related research activity: see section 3.2
- Partner: [ZFC](#)
- Funding: ACI/Gates
- Starting: 01/11/2004, ending: 31/03/2005

The goal of this contract was to implement a pilot (named Typed- π), in order to help the JCDecaux company to evaluate the potential of subject computing services based on ACE3 technologies in the area of smart literature and software.

7. Other Grants and Activities

7.1. European actions

7.1.1. Coordinated Action: Embedded WISENet

- Title: Supporting Embedded Systems for Embedded and Control Enabling Wireless Sensor Networks
- Partners: [Technische Universität Berlin](#) (Germany), [University of Cambridge](#) (UK), [Aalborg University](#) (Denmark), [Birkbeck Institute of Computer Science](#) (London), [University of Twente](#) (Netherlands), [Bielefeld University](#) (Germany), [Delft University of Technology](#) (The Netherlands), [University of Padova](#) (Italy), [Fraunhofer Institute of Technology](#) (Zurich), [Ghent University](#) (Belgium), [Université de Bourgogne](#) (France), [Institut National de Recherche en Informatique et en Automatique](#) (INRIA), [University of Stuttgart](#) (Germany)
- Starting: September 2004, ending: August 2006

Embedded WISENet aims to increase the awareness and to find out a vision as well as a research roadmap towards sensor networks of cooperating embedded systems, within the academic community and, most importantly, within the manufacturers of proper technologies as well as potential users community.

7.1.2. NoE ResNet

- Title: Resilience and Survivability for IST
- Head: LAAS
- Starting: beginning of 2004

The NoE ResIST (Resilience and Survivability for IST) will focus on the following four objectives in addressing the resiliency of dependability and security via resilience:

- Integration of teams of researchers so that the fundamental topics concerning scalability resilient ubiquitous systems are addressed by a critical mass of co-operative, multi-disciplinary research
- Identification, in an international context, of the key research directions indicated on the supporting ubiquitous systems by the requirement for trust and confidence in AoS
- Production of significant research results that pave the way for scalability resilient ubiquitous systems
- Promotion and propagation of a resilience culture in university curricula and in engineering best practices

- Certaines structures régulières (les listes) revêtent une importance particulière : ORG, NNPS*

Extraction de schémas

- Étant donné un nom d'organisme, on construit tous les schémas contenant jusqu'à 5 syntagmes à droite ou à gauche
- Seuls les "bons" schémas seront ensuite conservés en fonction de leur performance lors d'une phase d'apprentissage

Extraction de schémas

- Étant donné un nom d'organisme, on construit tous les schémas contenant jusqu'à 5 syntagmes à droite ou à gauche
- Seuls les “bons” schémas seront ensuite conservés en fonction de leur performance lors d'une phase d'apprentissage

Méthode

- Un ensemble de départ : L
- Un ensemble d'apprentissage : A
- Un ensemble de test : B
- Ces ensembles de documents sont annotés à la main
- Essais sur différentes combinaisons

Jeu de test 1	L	A	B	total
Documents	4	8	8	20
Entités différentes	74	238	144	456
Occurrences d'entités	238	418	271	827

Méthode

- Un ensemble de départ : L
- Un ensemble d'apprentissage : A
- Un ensemble de test : B
- Ces ensembles de documents sont annotés à la main
- Essais sur différentes combinaisons

Jeu de test 1	L	A	B	total
Documents	4	8	8	20
Entités différentes	74	238	144	456
Occurrences d'entités	238	418	271	827

Méthode

- Un ensemble de départ : L
- Un ensemble d'apprentissage : A
- Un ensemble de test : B
- Ces ensembles de documents sont annotés à la main
- Essais sur différentes combinaisons

Jeu de test 1	L	A	B	total
Documents	4	8	8	20
Entités différentes	74	238	144	456
Occurrences d'entités	238	418	271	827

Méthode

- Un ensemble de départ : L
- Un ensemble d'apprentissage : A
- Un ensemble de test : B
- Ces ensembles de documents sont annotés à la main
- Essais sur différentes combinaisons

Jeu de test 1	L	A	B	total
Documents	4	8	8	20
Entités différentes	74	238	144	456
Occurrences d'entités	238	418	271	827

Méthode

- Un ensemble de départ : L
- Un ensemble d'apprentissage : A
- Un ensemble de test : B
- Ces ensembles de documents sont annotés à la main
- Essais sur différentes combinaisons

Jeu de test 1	L	A	B	total
Documents	4	8	8	20
Entités différentes	74	238	144	456
Occurrences d'entités	238	418	271	827

Algorithme

- Soit les organismes qui sont annotés dans L
- Construire la liste des schémas de A+L qui les contiennent
- Appliquer ces schémas sur A ; supprimer ceux qui génèrent peu de résultats corrects ou trop d'incorrects ; classer par performance (nb de corrects/nb d'incorrects)
- Apprentissage sur A ; les schémas sont appliqués un à un dans l'ordre ; arrêt quand :
précision < seuil (la précision est initialement à 1)
rappel > seuil
- Test sur B (contrôle de la précision et du rappel)
- On applique sur tout les documents (éventuellement calcul de la précision)
- Utiliser le résultat "nettoyé" comme base de départ l'année suivante

Algorithme

- Soit les organismes qui sont annotés dans L
- Construire la liste des schémas de A+L qui les contiennent
- Appliquer ces schémas sur A ; supprimer ceux qui génèrent peu de résultats corrects ou trop d'incorrects ; classer par performance (nb de corrects/nb d'incorrects)
- Apprentissage sur A ; les schémas sont appliqués un à un dans l'ordre ; arrêt quand :
précision < seuil (la précision est initialement à 1)
rappel > seuil
- Test sur B (contrôle de la précision et du rappel)
- On applique sur tout les documents (éventuellement calcul de la précision)
- Utiliser le résultat "nettoyé" comme base de départ l'année suivante

Algorithme

- Soit les organismes qui sont annotés dans L
- Construire la liste des schémas de A+L qui les contiennent
- Appliquer ces schémas sur A ; supprimer ceux qui génèrent peu de résultats corrects ou trop d'incorrects ; classer par performance (nb de corrects/nb d'incorrects)
- Apprentissage sur A ; les schémas sont appliqués un à un dans l'ordre ; arrêt quand :
précision < seuil (la précision est initialement à 1)
rappel > seuil
- Test sur B (contrôle de la précision et du rappel)
- On applique sur tout les documents (éventuellement calcul de la précision)
- Utiliser le résultat "nettoyé" comme base de départ l'année suivante

Algorithme

- Soit les organismes qui sont annotés dans L
- Construire la liste des schémas de A+L qui les contiennent
- Appliquer ces schémas sur A ; supprimer ceux qui génèrent peu de résultats corrects ou trop d'incorrects ; classer par performance (nb de corrects/nb d'incorrects)
- Apprentissage sur A ; les schémas sont appliqués un à un dans l'ordre ; arrêt quand :
précision < seuil (la précision est initialement à 1)
rappel > seuil
- Test sur B (contrôle de la précision et du rappel)
- On applique sur tout les documents (éventuellement calcul de la précision)
- Utiliser le résultat "nettoyé" comme base de départ l'année suivante

Algorithme

- Soit les organismes qui sont annotés dans L
- Construire la liste des schémas de A+L qui les contiennent
- Appliquer ces schémas sur A ; supprimer ceux qui génèrent peu de résultats corrects ou trop d'incorrects ; classer par performance (nb de corrects/nb d'incorrects)
- Apprentissage sur A ; les schémas sont appliqués un à un dans l'ordre ; arrêt quand :
précision < seuil (la précision est initialement à 1)
rappel > seuil
- Test sur B (contrôle de la précision et du rappel)
- On applique sur tout les documents (éventuellement calcul de la précision)
- Utiliser le résultat "nettoyé" comme base de départ l'année suivante

Algorithme

- Soit les organismes qui sont annotés dans L
- Construire la liste des schémas de A+L qui les contiennent
- Appliquer ces schémas sur A ; supprimer ceux qui génèrent peu de résultats corrects ou trop d'incorrects ; classer par performance (nb de corrects/nb d'incorrects)
- Apprentissage sur A ; les schémas sont appliqués un à un dans l'ordre ; arrêt quand :
précision < seuil (la précision est initialement à 1)
rappel > seuil
- Test sur B (contrôle de la précision et du rappel)
- On applique sur tout les documents (éventuellement calcul de la précision)
- Utiliser le résultat "nettoyé" comme base de départ l'année suivante

Algorithme

- Soit les organismes qui sont annotés dans L
- Construire la liste des schémas de A+L qui les contiennent
- Appliquer ces schémas sur A ; supprimer ceux qui génèrent peu de résultats corrects ou trop d'incorrects ; classer par performance (nb de corrects/nb d'incorrects)
- Apprentissage sur A ; les schémas sont appliqués un à un dans l'ordre ; arrêt quand :
précision < seuil (la précision est initialement à 1)
rappel > seuil
- Test sur B (contrôle de la précision et du rappel)
- On applique sur tout les documents (éventuellement calcul de la précision)
- Utiliser le résultat "nettoyé" comme base de départ l'année suivante

Résultats

Seuil 0,6	Jeu de test 1	Jeu de test 2	Jeu de test 3
Entités dans L	74	93	114
Entités de L trouvées dans A	26	10	23
Occurrences d'entités de L trouvées dans A	116 (R=0,17 ; P=0,92) (MR=0,26 ; MP=0,94)	44 (R=0,17 ; P=0,96) (MR=0,22 ; MP=0,89)	51 (R=0,21 ; P=1) (MR=0,25 ; MP=0,96)
Schémas retenus	325	250	126
Entités dans A à la fin de l'apprentissage	247 (R=0,56 ; P=0,60) (MR=0,67 ; MP=0,63)	183 (R=0,56 ; P=0,60) (MR=0,67 ; MP=0,62)	126 (R=0,61 ; P=0,90) (MR=0,64 ; MP=0,91)
Entités de L trouvées dans B	12	24	17
Entités extraites de B	140	330	240
Comptage simple	R=0,31 ; P=0,46	R=0,45 ; P=0,41	R=0,36 ; P=0,47
Comptage multiple	R=0,35 ; P=0,44	R=0,53 ; P=0,41	R=0,42 ; P=0,49

Résultats

- Expériences avec différents seuils pour l'apprentissage (la précision est prépondérante : pourquoi un seuil sur le rappel ?) ; si on demande plus de précision, on extrait moins d'entités
- Expériences en séparant les parties “contrats” et “collaborations” : résultats moins bons

Résultats

- Expériences avec différents seuils pour l'apprentissage (la précision est prépondérante : pourquoi un seuil sur le rappel ?) ; si on demande plus de précision, on extrait moins d'entités
- Expériences en séparant les parties “contrats” et “collaborations” : résultats moins bons

Conclusion

- Problème difficile
- Ne marche pas très bien car les schémas sont trop génériques
- Aucune utilisation d'une connaissance de la langue

Conclusion

- Problème difficile
- Ne marche pas très bien car les schémas sont trop génériques
- Aucune utilisation d'une connaissance de la langue

Conclusion

- Problème difficile
- Ne marche pas très bien car les schémas sont trop génériques
- Aucune utilisation d'une connaissance de la langue

Perspectives

- Être moins générique
- Utiliser des propriétés lexicales (majuscules...)
- Utiliser des “valeurs” de syntagmes (partners, University of...) et des schémas prédéfinis
- Sélectionner plus localement les schémas (style d’écriture)
- Utiliser plus de connaissance linguistique
- ...

Perspectives

- Être moins générique
- Utiliser des propriétés lexicales (majuscules...)
- Utiliser des “valeurs” de syntagmes (partners, University of...) et des schémas prédéfinis
- Sélectionner plus localement les schémas (style d’écriture)
- Utiliser plus de connaissance linguistique
- ...

Perspectives

- Être moins générique
- Utiliser des propriétés lexicales (majuscules...)
- Utiliser des “valeurs” de syntagmes (partners, University of...) et des schémas prédéfinis
- Sélectionner plus localement les schémas (style d’écriture)
- Utiliser plus de connaissance linguistique
- ...

Perspectives

- Être moins générique
- Utiliser des propriétés lexicales (majuscules...)
- Utiliser des “valeurs” de syntagmes (partners, University of...) et des schémas prédéfinis
- Sélectionner plus localement les schémas (style d’écriture)
- Utiliser plus de connaissance linguistique
- ...

Perspectives

- Être moins générique
- Utiliser des propriétés lexicales (majuscules...)
- Utiliser des “valeurs” de syntagmes (partners, University of...) et des schémas prédéfinis
- Sélectionner plus localement les schémas (style d’écriture)
- Utiliser plus de connaissance linguistique
- ...

Perspectives

- Être moins générique
- Utiliser des propriétés lexicales (majuscules...)
- Utiliser des “valeurs” de syntagmes (partners, University of...) et des schémas prédéfinis
- Sélectionner plus localement les schémas (style d’écriture)
- Utiliser plus de connaissance linguistique
- ...

Fin

- Merci...