



# Module Web Intelligence (WI)



# Chapitre I : Introduction

**Pr. Okba KAZAR**

**Professeur des universités  
Directeur du Laboratoire d'INFormatique Intelligente  
LINFI**

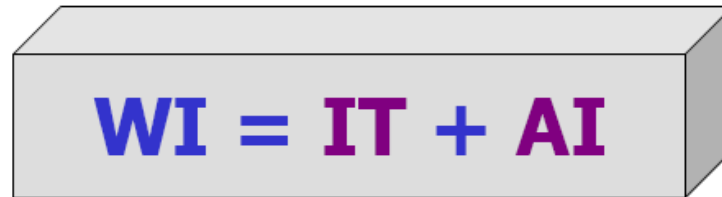
**Département d'informatique  
Université de biskra**

# Contenu

- Une brève histoire du Web Intelligence
- Motivations pour WI
- Définition et perspectives de WI
- Programme de recherche
- Principaux outils de Web Intelligence
- conclusion

# C'est quoi "Web Intelligence"

- Web Intelligence (WI) exploite l'impact fondamental et pratique que la **technologie de l'information** de pointe (TI) et l'innovante **intelligence artificielle** (IA) auront sur le Web


$$WI = IT + AI$$

- Integration de IT avec IA
- Applications de IA sur le Web

# Une brève histoire du Web Intelligence

- 1999: initiatives de recherche Collaborative
  - Ning Zhong, Data Mining et Systemes à base de connaissances
  - Jiming Liu, agents Intelligent et système multi-agents
  - Yiyu Yao, La recherche d'information et les systèmes d'information intelligents
- Les efforts de recherche combinés avec objectif commun: créer une nouvelle sous-discipline couvrant les théories et les techniques liées à l'information Web .

# Une brève histoire du Web Intelligence

- 2000: Publication d'un document de deux pages sur la position sur WI (Zhong, Liu, Yao, Ohsuga, COMPSAC 2000)

**Web Intelligence (WI)**

Ning Zhong Machihashi Inst. Tech.	Jining Liu Hong Kong Baptist U.	Y.Y. Yao U. Regina	Satoshi Ohsuga Waseda U.
--------------------------------------	------------------------------------	-----------------------	-----------------------------

**1 Introduction**

The 21st century is the age of Internet and World Wide Web. The Web revolutionizes the way we gather, process, and use information. At the same time, it also revolutionizes the meanings and processes of business, commerce, marketing, finance, publishing, education, research, development, as well as other aspects of our daily life. The revolution is just beginning. Although individual Web-based information systems are constantly being deployed, advanced means and techniques for developing and for benefiting from Web intelligence still remain to be systematically studied.

This position paper defines a new research field, namely Web Intelligence (WI for short) by giving a complete picture of WI related topics for systematic study on advanced Web technology and developing Web-based intelligent information systems. Regularly updating WI expertise AI and advanced information technology on the Web and Internet. It is the key and the most urgent research field of IT for business intelligence.

**2 WI Related Topics**

What are issues and research topics on WI? In order to study advanced Web technology systematically, and develop advanced web-based intelligent information systems, we give an overview of WI related topics as shown in Figure 1 and list several major subtopics in each topic below.

- **Web Human-Media Engineering:** the art of Web page design, multimedia information representation, multimedia information processing, visualization of Web information, and Web-based human-computer interface.
- **Web Information Management:** data quality management, information transformation, Internet and Web-based data management, multidimensional Web databases and OLAP (on-line analytical processing), multimedia information management, new data models for the Web, object-oriented Web information management, personalized information management, semi-structured data management,

use and management of metadata, Web knowledge management, Web page automatic generation and updating, as well as Web security, integrity, privacy and trust.

- **Web Information Retrieval:** approximate retrieval, conceptual information extraction, image retrieval, multi-linguistic information retrieval, multimedia retrieval, new retrieval models, ontology-based information retrieval, as well as automatic Web content cataloging and indexing.
- **Web Agents:** dynamics of information sources, e-mail filtering, e-mail semi-automatic reply, global information collecting, information filtering, navigation guides, recommender systems, recommendation agents, reputation mechanisms, resource intermediary and coordination mechanisms, as well as Web-based cooperative problem solving.
- **Web Mining and Farming:** data mining, and knowledge discovery, hypothesis analysis and transformation, learning user profiles, multimedia data mining, regularities in Web surfing and Internet navigation, text mining, Web-based ontology engineering, Web-based reverse engineering, Web farming, Web-log mining, and Web watchmaking.
- **Web Information System Environment and Foundations:** competitive dynamics of Web sites, emerging Web technology, network community formation and support, new Web information description and query languages, theories of small world Web, Web information system development tools, and Web protocols.
- **Web-Based Applications:** business intelligence, computational sciences and methods, conversational systems, customer relationship management (CRM), direct marketing, electronic commerce and electronic business, electronic library, information markets, price dynamics and pricing algorithms, measuring and analyzing Web merchandising, Web-based decision support systems, Web-based



Figure 1. A schematic diagram of WI related topics

distributed information systems, Web-based electronic data interchange (EDI), Web marketing and Web publishing

**3 WI Related Case Studies:**

WI presents an excellent opportunity as well as a challenge to the research and development of new generation of information processing technology, as well as for exploiting business intelligence. Specifically, e-commerce activity data involves the real time in investigating a significant revolution [5]. The ability to track users' browsing behavior down to individual mouse clicks has brought the vendor and vendor customers closer than ever before. It is now possible to send a personalized product message for individual customers at a massive scale. This is called *Targeted Marketing*.

Hackathon proposed Web Farming that is the systematic mining of information resources on the Web for business intelligence [3].

Ahmed et al. systematically investigated the data on the Web and the issues of semi-structured data [1].

Zhong, Yao et al. proposed a way of mining search data and quality rules that can be used for Web-log mining [11]. They proposed ways for targeted marketing by mining identification rules and content rules (content [5], [6]). They are also working on text mining on the Web including automatic construction of ontology, content mining systems, and Web-based "business systems [14], [6].

**4 WI Conferences**

We initiated a new high-quality, high impact biennial conference series, namely the Asia-Pacific Conference on

Web Intelligence (WI). The first meeting in this new series, WI'2000, will be held in Machihashi City, Japan, October 23, 26, 2000 (<http://www.machihashi-it.ac.jp/wi2000/>).

WI 2001 is an international forum for researchers and practitioners to present the state-of-the-art in the development of Web intelligence, to examine performance characteristics of various approaches in Web-based intelligent information technology, and to cross-fertilize ideas on the development of Web-based intelligent information systems among different domains. By idea-sharing and discussions on the underlying formalizations and the modeling methods, goals of Web intelligence, WI 2001 is expected to stimulate the future development of new models, new methodologies, and new tools for building a variety of embodiments of Web-based intelligent information systems.

**References**

- [1] Ahmed, S., Bennett, P. and Sreeni, C. Data on the Web. Morgan Kaufmann, 2000.
- [2] Dong, J.Z., Zhong, N., and Zhang, F. "Probabilistic Rough Inclusion, The CAPM Methodology on Algorithms." J. Int. Res. on A. Science and Management of Information Systems, 1:51-100. Springer Verlag (1999) 821-840.
- [3] Hackathon, ED. Web Farming for the Data Warehouse. Morgan Kaufmann, 2000.
- [4] Liu, J., and Li, C. The Mining of User Markers: methods and applications. Information Systems, 22(4) (1999) 31-37.
- [5] Liu, J. and Zhong, N. web INTELLIGENCE AND SEMI-STRUCTURED SYSTEMS. Morgan Kaufmann, 2000.
- [6] Liu, J., Ohsuga, S., Yao, Y.Y. and Zhang, F.F. "New Web Mining: The Mining of Information on the Web." Springer Verlag (2001).
- [7] Liu, J. and Li, C. Data Analysis in E-commerce Agents. Electronic Commerce Research and Applications, Springer Verlag (2001).
- [8] Bennett, S. et al. "Web Usage Mining: Discovering and Applying Knowledge from Web Data." SIGKDD International Conference on Data Mining and Knowledge Discovery (2000) 122-131.
- [9] Yao, Y.Y. and Zhang, N. "Global Market Value Potentials for Targeted Marketing." Springer, 2001.
- [10] Yao, Y.Y. and Zhang, N. "Customer Clustering using Information Labels." In: Liu, J., Yao, Y.Y. and Gabriel, L. (eds.) "Computational Intelligence and Data Mining: Applications and Methods." Springer Verlag (2000).
- [11] Zhong, N., Yao, Y.Y. and Ahmed, S. "Probabilistic Rough Inclusion Mining." J. Expert and Intelligent Systems: Principles and Applications (Proceedings of the International Conference on Intelligent Systems, Springer Verlag, 1999) 193-198.
- [12] Zhong, N., Bennett, A. and Zhang, S. web-based Business on Rapid Data Mining and Genetic Data Computing. IJMI (11), Springer Verlag (1999).
- [13] Zhong, N. and Zhou, L. web Mining for Knowledge Discovery. Doctoral Dissertation, Waseda U., Japan (1999).
- [14] Zhong, N., Yao, Y.Y. and Bennett, S. "Automatic Construction of Ontology from Data Mining." In: Second International Conference on Data Mining (ICDM) (1999) 309-316, IEEE Press, 2000.
- [15] Zhong, N., Yao, Y.Y. and Ahmed, S. "Web Discovery by Self-organizing Neural Networks." In: Intelligent Systems in Engineering Computing, Springer-Verlag, 1999.
- [16] Zhong, N. et al. "New on-line Mining: a Knowledge Mining and Analysis Framework." Intelligent Systems, Springer-Verlag (2000).
- [17] Zhong, N., Knowledge Discovery and Data Mining, in the Encyclopedia of Information Science and Technology, Elsevier (2001).

# Une brève histoire du Web Intelligence

- 2001: First Asia-Pacific Conference on Web Intelligence
- 2002: Publication of first special issue on WI in IEEE Computer
- 2002: Web Intelligence Consortium
- 2003: First edited book on WI
- 2005: The international WIC Institute

# Motivations

- **La taille de Web**

Difficultés dans le stockage, la gestion, et la récupération (recherche) effective et efficace

- **Complexité du Web**

- Collection hétérogène de documents web structurés, non structurés, semi-structurés, interdépendants et distribués.

- Consiste en textes, images et sons

# Motivations

## Web Intelligence on the Web

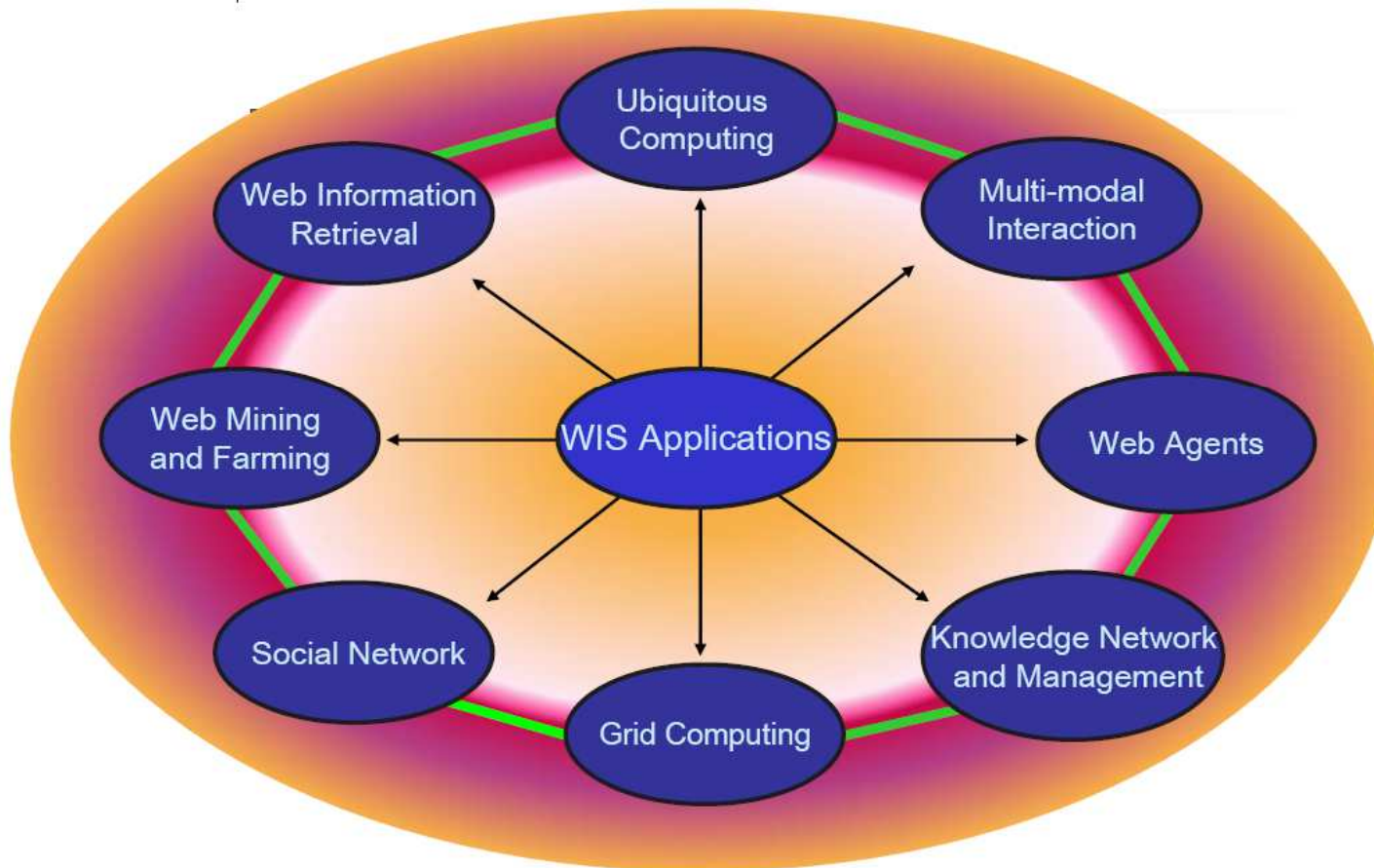
Search Engine	Number of Hits (February 2001)	Number of Hits (February 2003)
Lycos	1,102,279	<b>7,163,922</b>
Google	1,080,000	<b>2,590,000</b>
AltaVista	1,271	<b>1,860,062</b>
Netscape	77	<b>1,900,000</b>
Yahoo	74	<b>2,510,000</b>



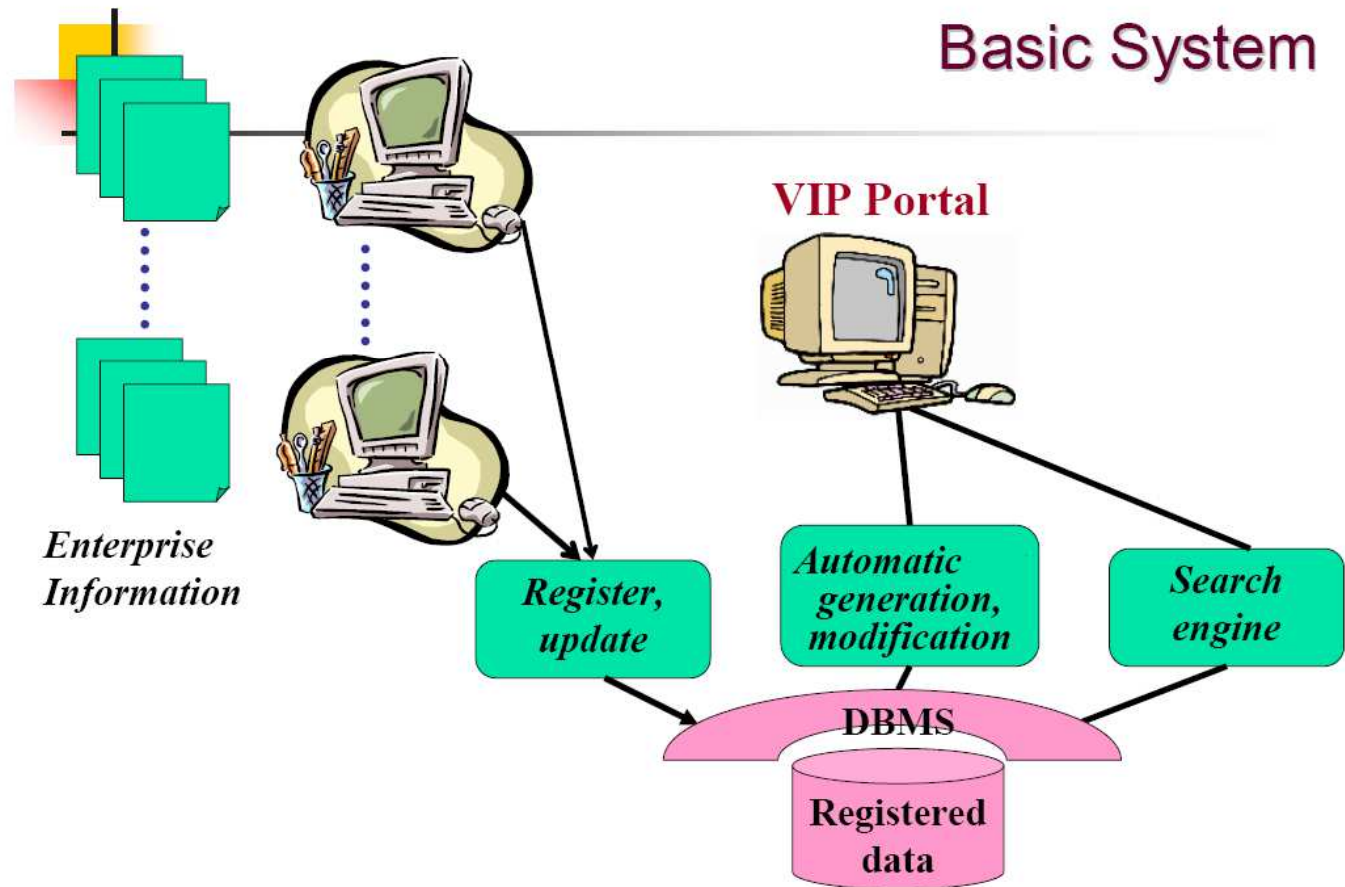
# Motivations

- Production de données sur le Web est à un taux de croissance exponentielle.
- Un intérêt industriel en croissance rapide dans WI.
- Seuls quelques articles **académiques**.
- Nous devons réduire **l'écart** entre les besoins de **l'industrie** et la **recherche universitaire**

# Systeme en Web Intelligence



# Un exemple



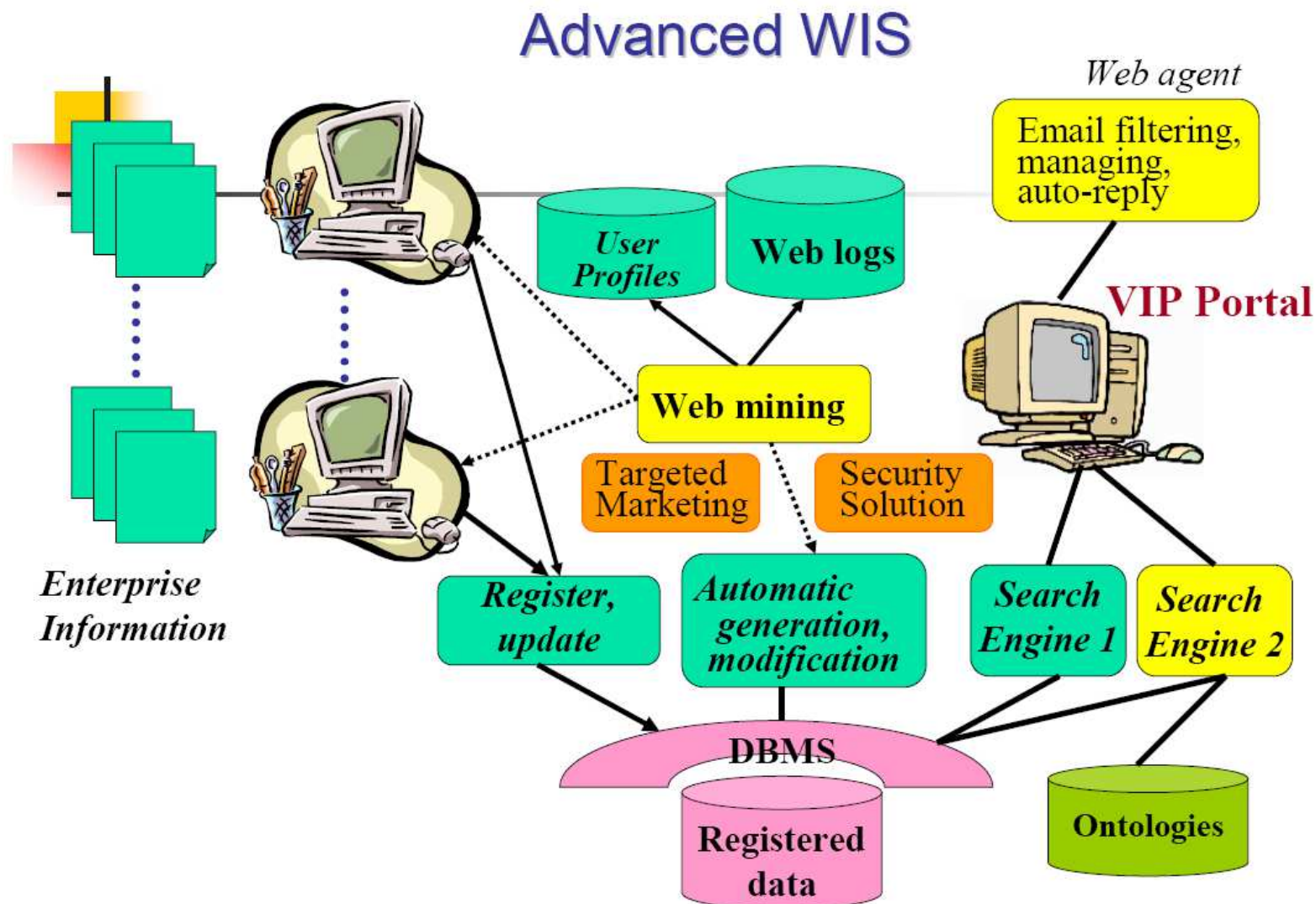
# Questions avancée

- **Comment le client utilise le portail VIP** afin de cibler les produits et gérer les promotions et campagnes marketing?
- Quel est l'association sémantique entre les pages visitées par le client?
- Est-ce que le visiteur est familiariser avec la structure du site Web? Ou est-il un nouvel utilisateur ou quelqu'un au hasard?
- Est-ce que le visiteur est un robot Web ou d'autres utilisateurs?
- ...

# systeme WI avancé

- Faire une recommandation dynamique à un utilisateur Web basé sur le profil de l'utilisateur et les comportements d'utilisation;
- La modification automatique du contenu et de l'organisation d'un site web;
- En combinant les données d'utilisation sur le Web avec des données marketing pour donner des informations sur la façon dont les visiteurs ont utilisé un site Web.

# Systeme avance pour le WI



# Perspectives du WI

- WI peut être classés en quatre catégories
- (based on Russel & Norvig`s scheme)

conception de la philosophie de WIS

capacité,  
la fonctionnalité  
d'un WIS

System that <b>thinks</b> like <b>humans</b>	System that <b>thinks</b> <b>rationally</b>
System that <b>acts</b> like <b>humans</b>	System that <b>acts</b> <b>rationally</b>

# Intérêts industriels dans WI

- **Web Intelligence** [kis-lab.com/wi01/](http://kis-lab.com/wi01/)
- **Web-Intelligence** Home Page
  - [www.web-intelligence.com/](http://www.web-intelligence.com/)
- **Intelligence** on the **Web**
  - [www.fas.org/irp/intelwww.html](http://www.fas.org/irp/intelwww.html)
- WIN: home **WEB INTELLIGENCE NETWORK**,
  - [smarter.net/](http://smarter.net/)
- CatchTheWeb - **Web** Research, **Web Intelligence** Collaboration [www.catchtheweb.com/](http://www.catchtheweb.com/)
- Infonoia: **Web Intelligence** In Your Hands
  - [www.infonoia.com/myagent/en/baseframe.html](http://www.infonoia.com/myagent/en/baseframe.html)

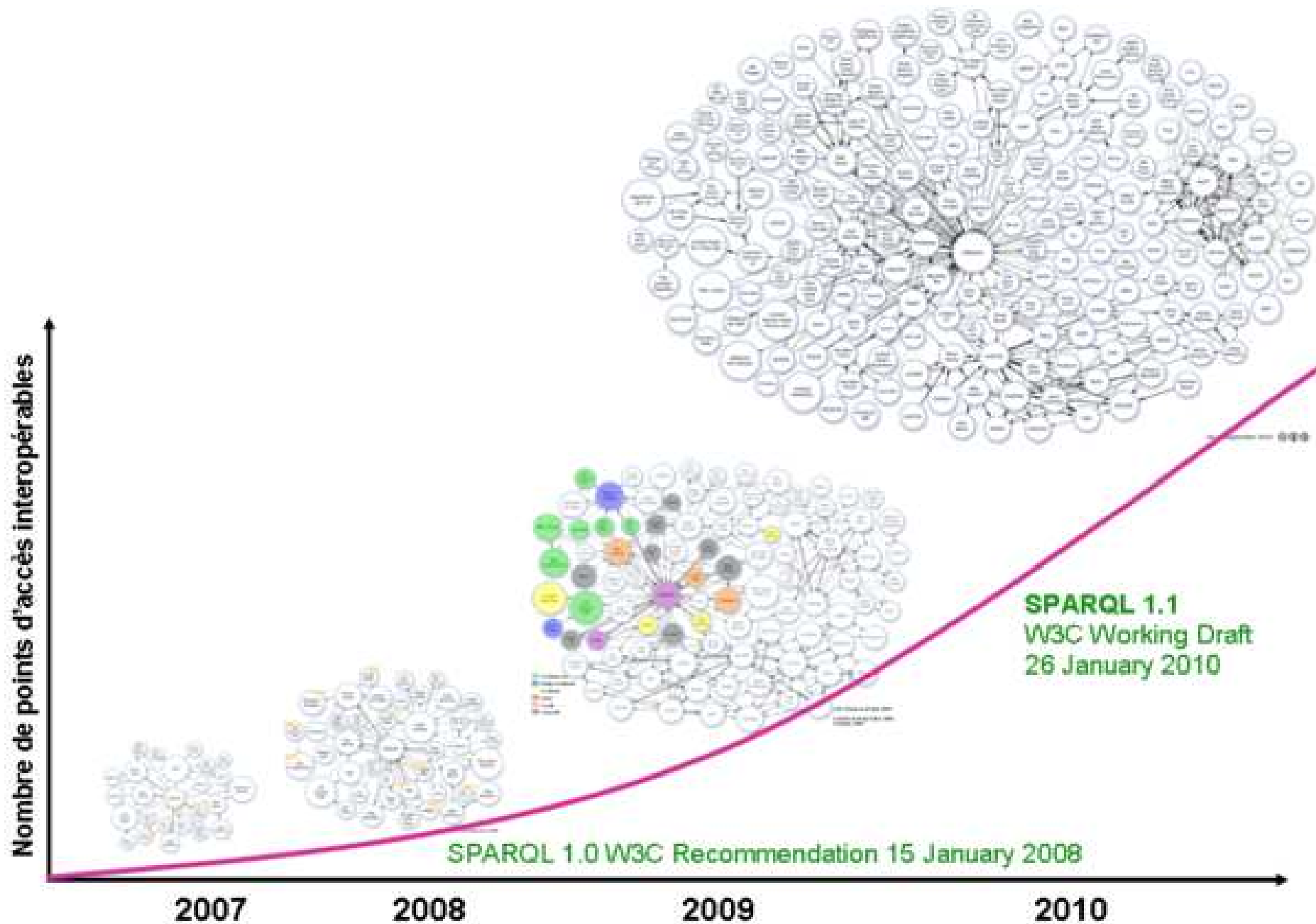


# Programme de recherche sur le WI

- Web sémantique, « webmining » et construction automatique d'ontologies
- réseau social intelligent
  
- Cloud computing...
- Application : e-business, e-learning, e-santé, (mobile, ubiquitous....)
- ...

# Le Web Sémantique (3.0)

- Le Web sémantique est basé sur les **langages** qui font plus de contenu sémantique de la page disponible en **formats lisibles par la machine** pour un traitement automatique basé sur des agents.
- Un langage «sémantique» qui lie les informations sur une page à la sémantique lisibles par la machine s'appelle « ontologie ».



Le **Web** Sémantique ou l'importance des données liées

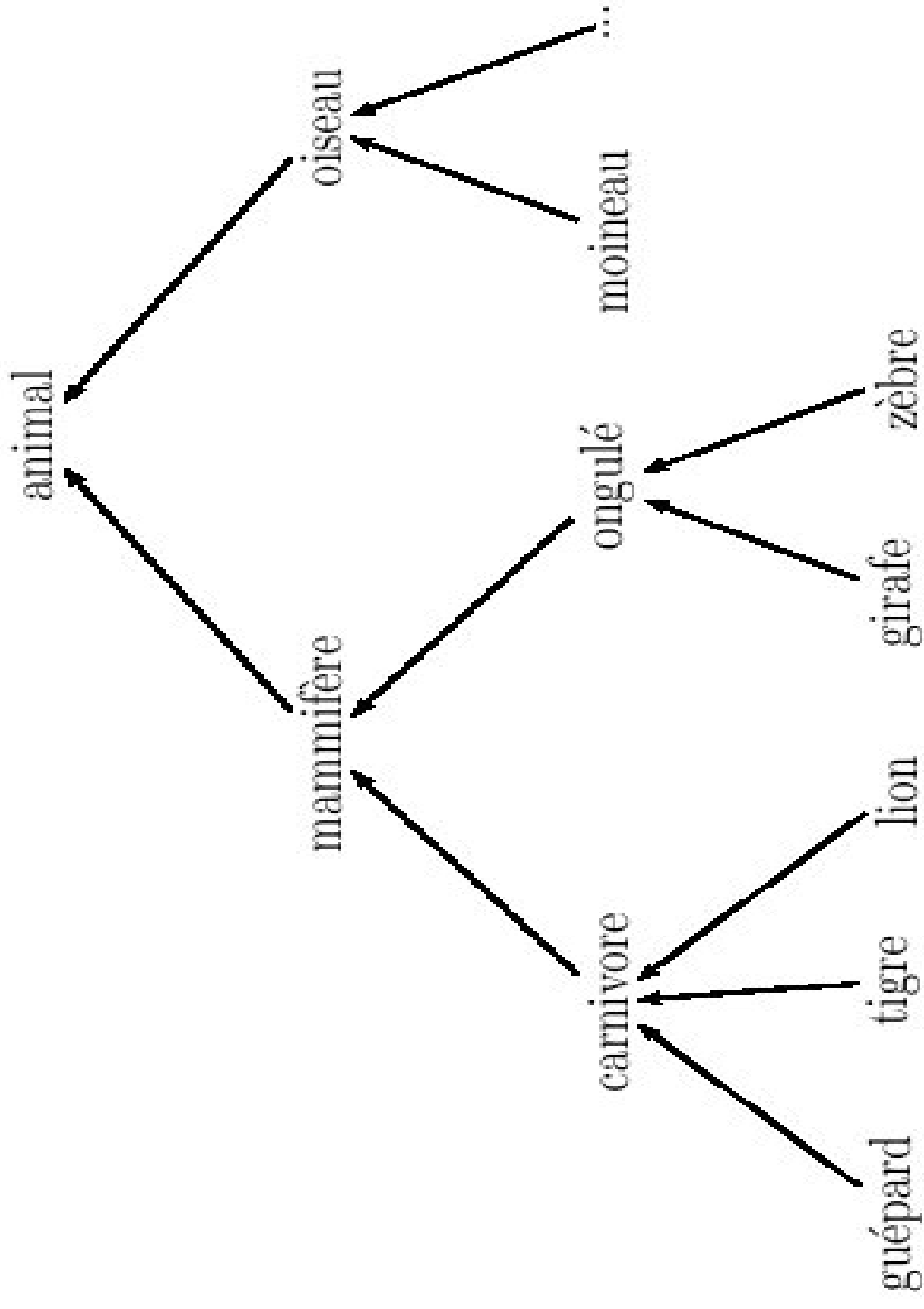
# Compoants du Web Semantique

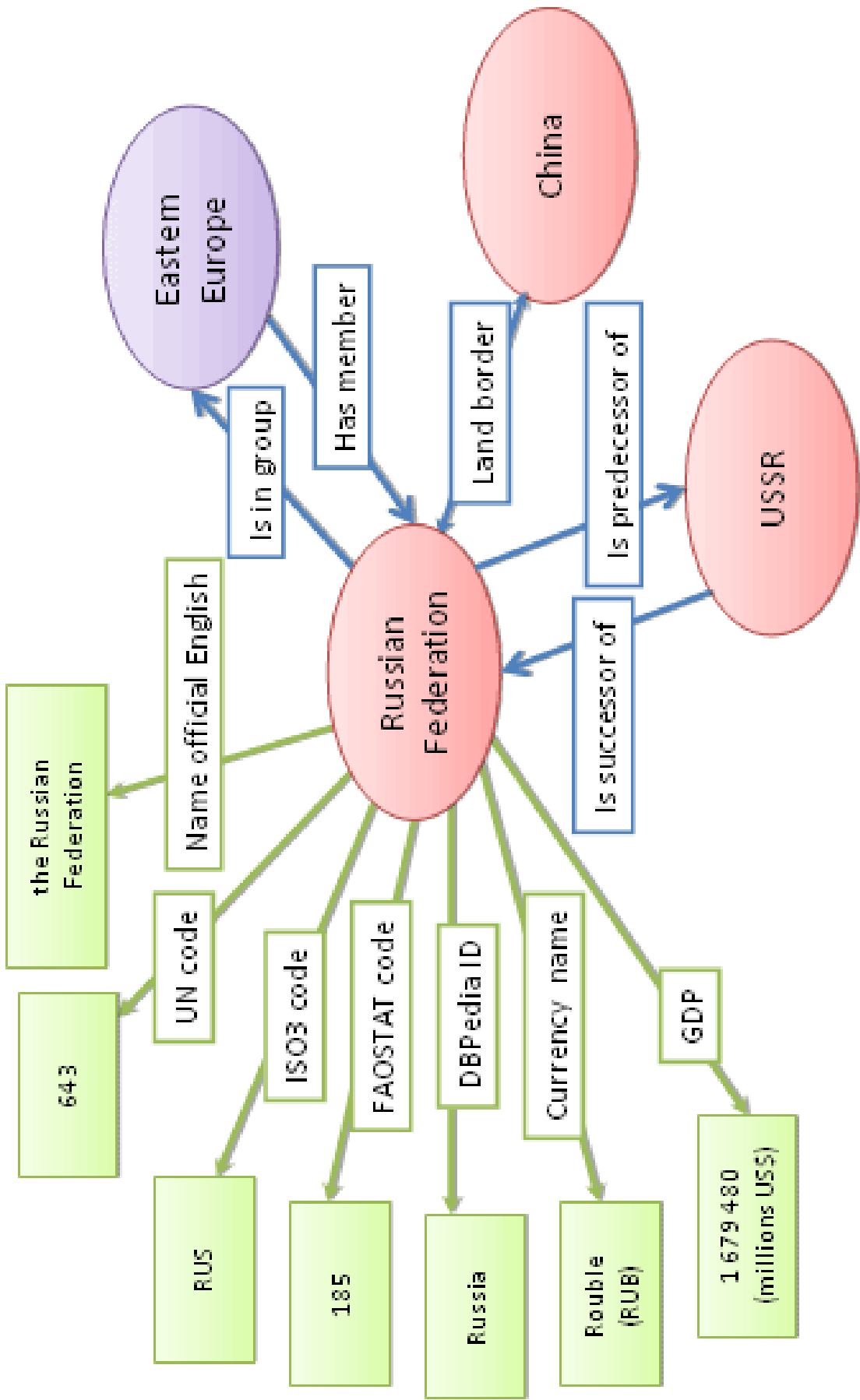
- Un **modèle de données** unifié comme RDF.
- **Langues** avec sémantique définie, construit sur RDF, tels que OWL (DAML + OIL).
- **Ontologies** pour normaliser la terminologie pour le balisage des ressources Web.
- Les **outils** qui aident la génération et le traitement de balisage sémantique.

Les Ontologies fournissent l'épine dorsale sémantique pour les applications du Web sémantique.

# Les ontologies offrent

- **Communication**
  - Les modèles normatifs, des réseaux de relations
- **Partage et réutilisation**
  - Spécifications, Fiabilité
- **Contrôle**
  - Classification et découverte, partage, découverte des relations

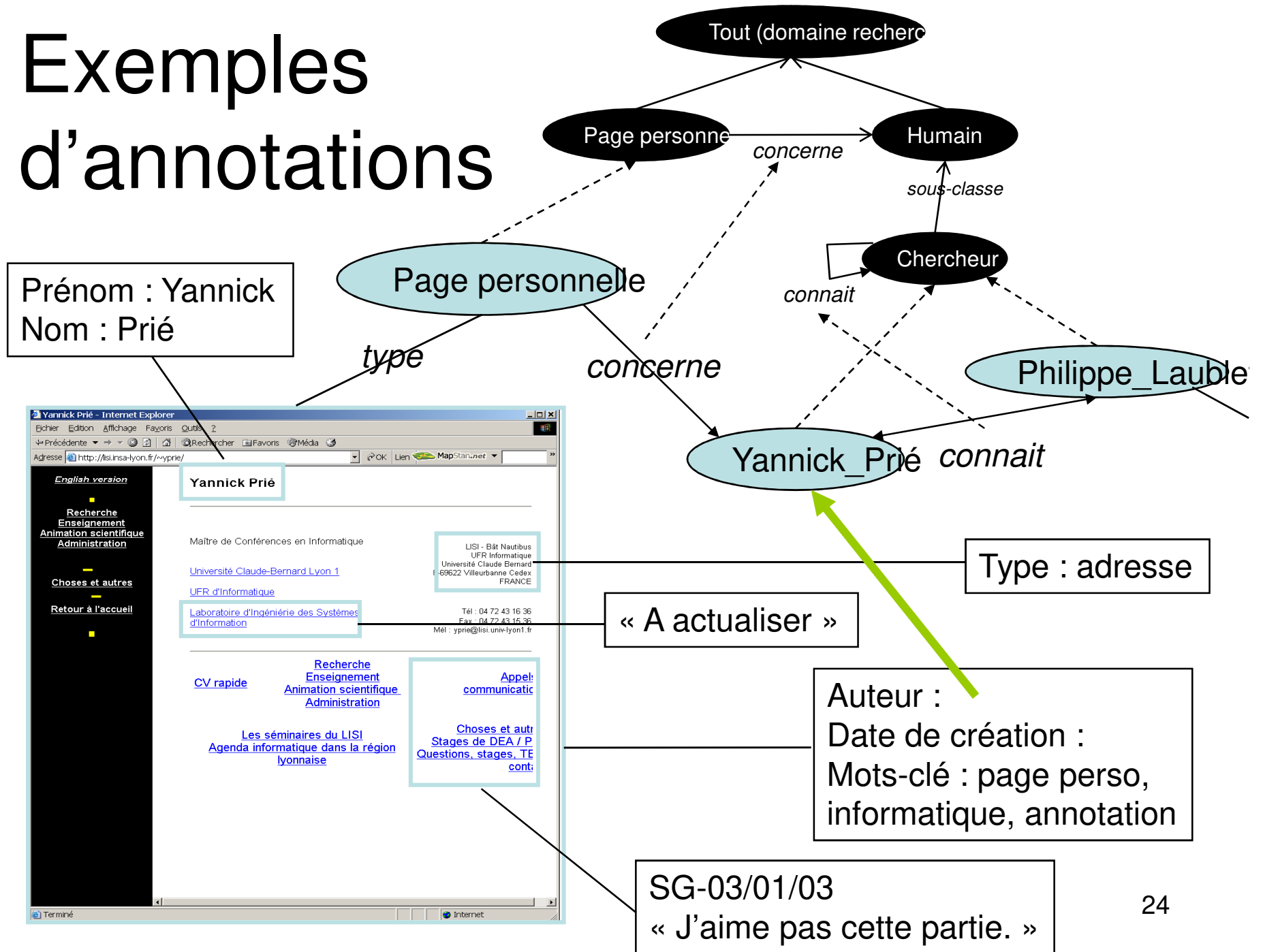




- Data value
- Territory
- Group

→ Relationships with other territories/groups  
→ Associated data

# Exemples d'annotations





# Categories des Ontologies

- Une **ontologie spécifique au domaine** décrit un domaine technique ou commercial bien défini.
- Une ontologie de **tâche** peut être soit spécifique à un domaine ou reconstituées à partir d'un ensemble d'ontologies spécifiques au domaine pour répondre aux exigences d'une tâche.
- Une ontologie **universelle** décrit connaissances à des niveaux supérieurs

# Programme de recherche sur le WI

- Web sémantique, « webmining » et construction automatique d'ontologies
- réseau social intelligent
  
- Cloud computing...
- ...
- ...

# Social Network



Anne Helmond, May 2009

- Adopter et intégrer le web 2.0 sont plus faciles que le définir
- Web 1.0 vs Web 2.0

Facteurs de comparaison	Web 1.0	Web 2.0
Contribution aux contenus	Dépend du Webmestre	Dépend des internautes
Fréquence actualisation	Relativement élevée	Très élevée (participation active d'un grand nombre d'intervenants)
Relation	« <b>Un à plusieurs</b> » : webmestre aux lecteurs	« <b>Plusieurs à plusieurs</b> » : l'internaute est à la fois lecteur et auteur (peut réagir, interagir...)

- Principe du « push » et du « pull »
- Pas un effet de mode, mais un véritable phénomène :
  - Personnalité de l'année 2006 magazine Time : « YOU » (Vous, les internautes générant le contenu sur Internet)
  - Loin d'être éphémère

# Le Web comme un Graphe

- Nous pouvons voir le Web comme un réseau social orienté qui relie les gens (organisations ou entités sociales).
- Questions de recherche:
  - Quelle est la taille du graphe? (degré sortant et degré entrant)
  - Peut-on naviguer à partir de n'importe quelle page à une autre? (clics)
  - Pouvons-nous exploiter la structure du Web? (recherche et webmining)
  - Comment découvrir et gérer les communautés Web?
  - Est-ce que le graphe du Web révèle sur la dynamique sociale?

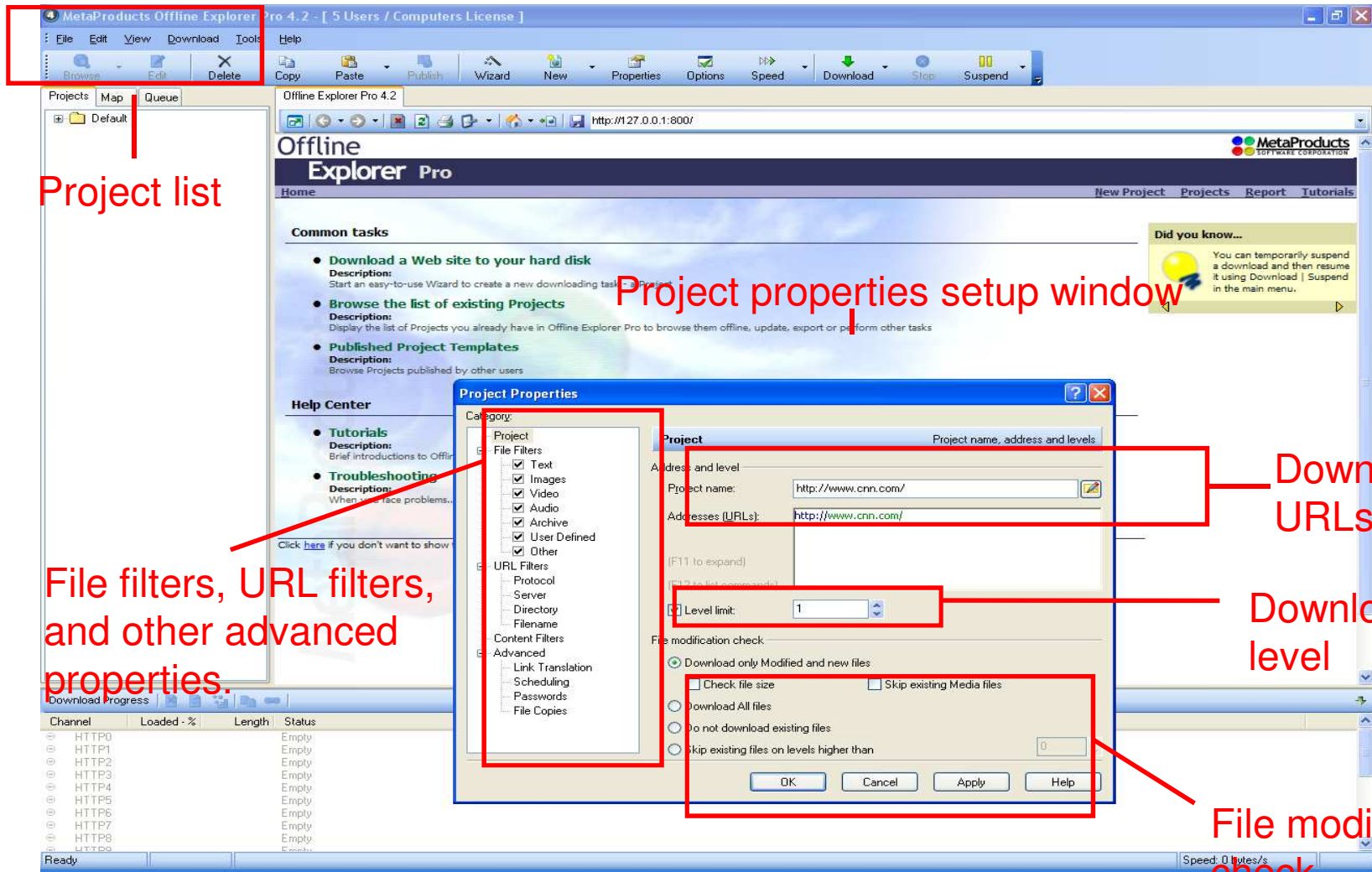
# Outline

- A brief history of Web Intelligence
- Motivations for WI
- Definition and Perspectives of WI
- Agenda de Recherche
- **Les outils Majeurs pour le Web Intelligence**
- Conclusion

# Principaux outils de Web Intelligence

- I. [Collection](#)
  - Offline Explorer
  - SpidersRUs (AI Lab)
  - Google Scholar
- II. [Analyses](#) (Données et Text Mining)
  - Google APIs
  - Google Translation
  - GATE
  - Arizona Noun Phraser (AI Lab)
  - Self-Organizing Map, SOM (AI Lab)
- III. [Visualisation](#)
  - NetDraw
  - JUNG
  - Analyst's Notebook and Starlight

# Collection: Offline Explorer



Project list

Project properties setup window

File filters, URL filters, and other advanced properties.

Download URLs

Download level

File modification check



- Offline Explorer télécharge vos [sites Web](#) et FTP favoris (jusqu'à une centaine simultanément) afin de les consulter hors connexion, de les imprimer ou de les parcourir.
- Une des forces de Offline Explorer est de vous donner accès de manière sélective aux différents serveurs (inclus ou exclus), répertoires et fichiers en utilisant seulement des mots-clé.
- Efficace et rapide, il comprend un support de technologies au standard industriel, comme FTP, différents serveurs proxy, Macromedia Flash, Cookies, un serveur HTTP intégré vous permettant de partager vos fichiers téléchargés sur l'Intranet.

# Analysis: Google APIs

- L'*API Google* (application programming interface) est un ensemble d'outils mis à disposition de tous qui permet d'interroger à distance les serveurs du Moteur de recherches *Google*.  
<http://code.google.com/more/>
- Examples of Google APIs:
  - **Google API for Inlink:** Discovers what pages link to your website.
  - **Google Data APIs:** Provide a simple, standard protocol for reading and writing data on the Web. Several Google services provide a Google Data API, including Google Base, Blogger, Google Calendar, Google Spreadsheets and Picasa Web Albums.
  - **Google AJAX Search API:** Uses JavaScript to embed a simple, dynamic Google search box and display search results in your own Web pages.
  - **Google Analytics:** Allows users gather, view, and analyze data about their Website traffic. Users can see which content gets the most visits, average page views and time on site for visits.
  - **Google Safe Browsing APIs:** Allow client applications to check URLs against Google's constantly-updated blacklists of suspected phishing and malware pages.
  - **YouTube Data API:** Integrates online videos from YouTube into your applications.

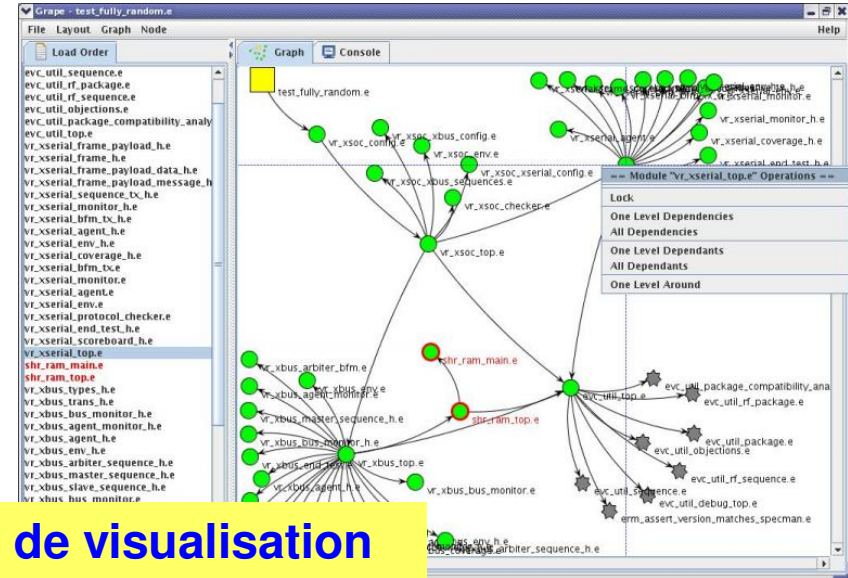
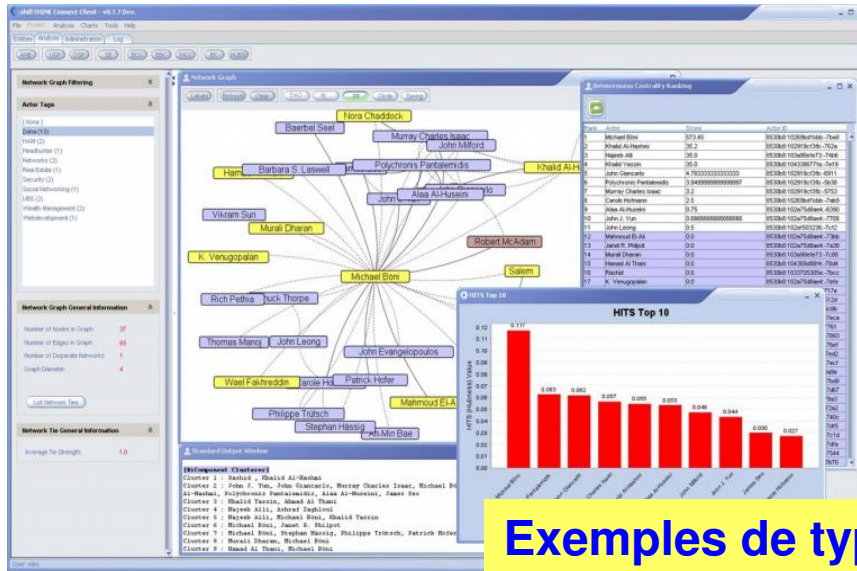
# Visualization: JUNG

- Java Universal Network/Graph Framework (JUNG) est une bibliothèque de logiciels pour la modélisation, l'analyse et la visualisation des données qui peuvent être représentées sous forme de graphique ou de réseau. Il a été développé par l'École de l'information et d'informatique à l'Université de Californie, Irvine

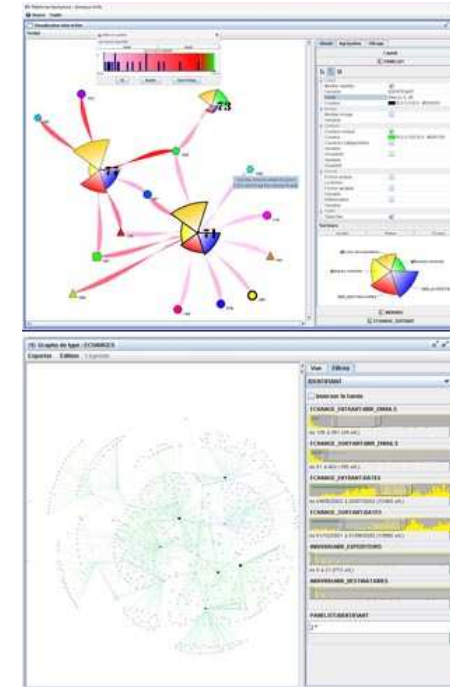
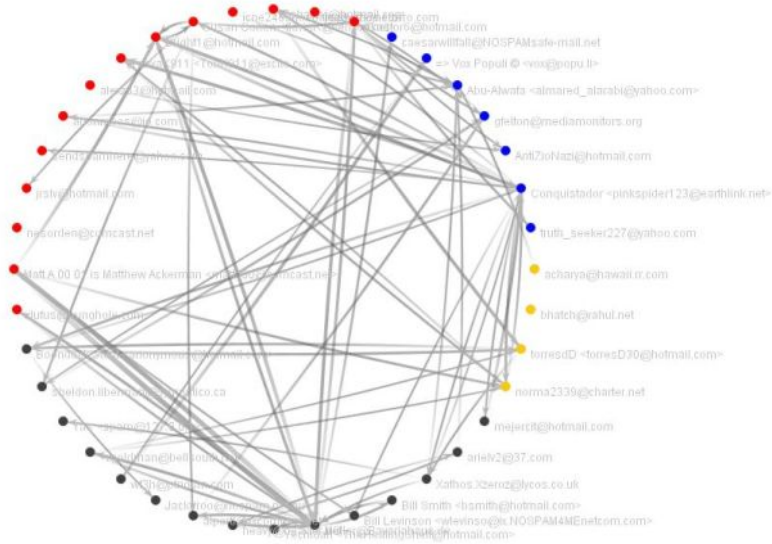
<http://jung.sourceforge.net/index.html>

- La distribution actuelle de JUNG comprend des implémentations d'un certain nombre d'algorithmes de la théorie des graphes, data mining et d'analyse de réseau social
  - Clustering
  - Decomposition
  - Optimization
  - Random Graph Generation
  - Statistical Analysis
  - Calculation of Network Distances and Flows and Importance Measures (Centrality, PageRank, HITS, etc.).

# JUNG



Exemples de types de visualisation



# A retenir

- Le couplage de l'hypertexte d'Internet conduit à une révolution: le Web.
- Le couplage de l'intelligence artificielle et de la technologie d'information avancé sur la plate-forme de Web, mènera à un autre changement de paradigme: le Web Intelligent.

.

**Fin du premier cours**