



The  
University  
Of  
Sheffield.

**ISWC 2007**

# Semantic Web Technologies for Knowledge Management in Large Distributed Organisations

Professor Fabio Ciravegna

Intelligent Web Technology Lab, Natural Language Processing Group

Department of Computer Science, University of Sheffield

<http://www.dcs.shef.ac.uk/~fabio/>

Sponsored by



[www.3worldt.org](http://www.3worldt.org)

Sponsored by



[www.x-media-project.org](http://www.x-media-project.org)



# Copyright Notice

- These slides were presented during the International Semantic Web Conference (ISWC 2007) in Busan, Korea, November 2007 (<http://www.iswc07.org/>)
- They can be reused for personal and educational reasons only
- Updated versions and full conditions of use can be found at:
  - <http://www.dcs.shef.ac.uk/~fabio/ISWCTutorial/>



# Outline of Tutorial

- 9.00-9.15 Introduction and scene setting
  - Issues in knowledge management in large organisations
- 9.15-9.45 Ontologies and ontology engineering
- 9.45-11.30 Semantic web technologies for knowledge acquisition
  - 10.30-11.00 Coffee break
- 11.30-12.15 Semantic web technologies for knowledge sharing, reuse and retrieval
- 12.15-12.30 Conclusion and future work
- 12.30- Discussion

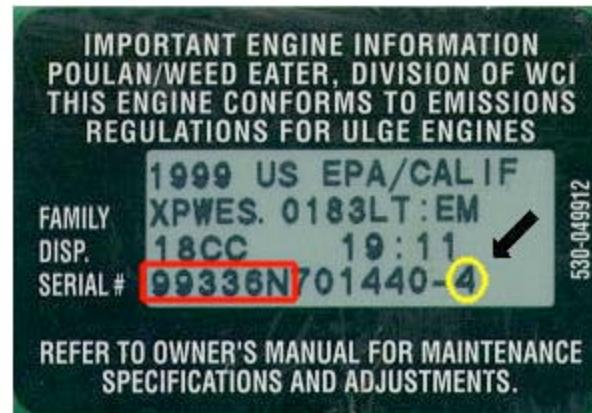
- Gathering knowledge relevant to a task or problem
  - it may be distributed across different storage systems and different media
- Analysing the knowledge they have gathered and make sense of it
- Sharing knowledge with their colleagues
- Keeping track of the process
  - by being aware of what one is doing, what one needs to do next, and what others are doing
- What to search for, what analysis is needed and who to share with
  - depend on the task in hand and the current stage of the process

jet engines are moving towards complete serialisation

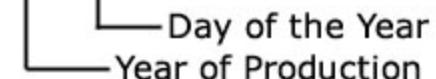
- every piece has a serial number (excepts nuts and bolts)
- the history of each part is recorded
  - e.g. part transferred between engines



© Rolls-Royce plc



99336N = Date Code



4 = Product Type

## Jet engine example

- a jet engine can produce ~1Gbyte of vibration data per hour of flight;
  - if irregularities are found, part of the data can be stored
  - reports can be written (event reports)
  - pictures can be taken

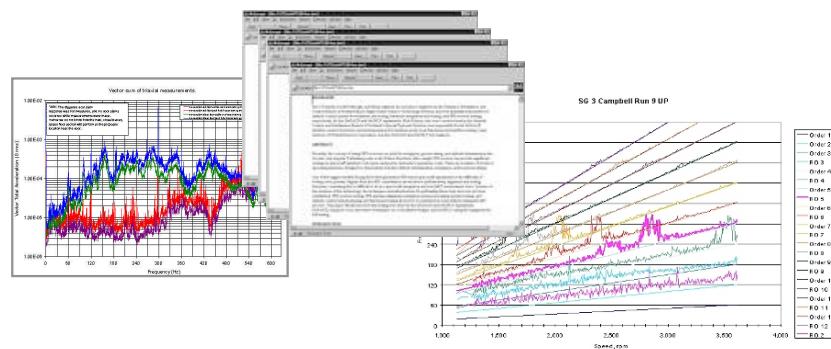


image © [www.rolls-royce.com](http://www.rolls-royce.com)



## Jet engine example (3)

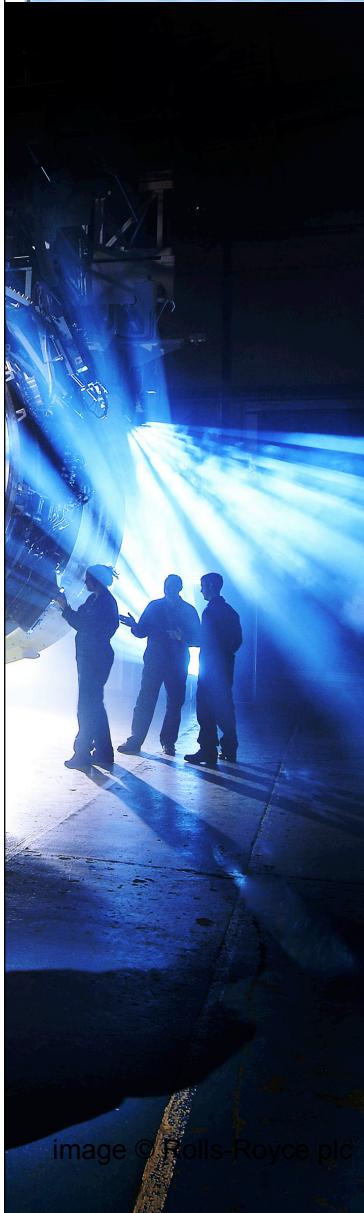
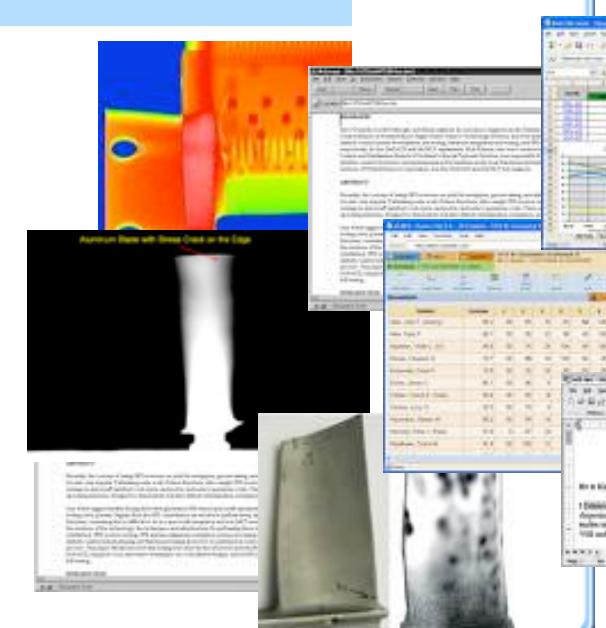


image © Rolls-Royce plc

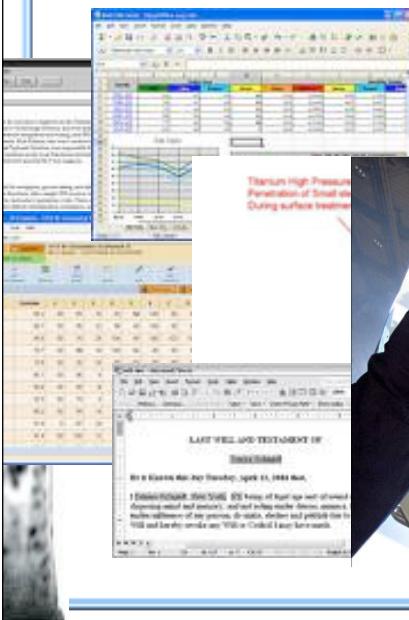
When engine is serviced (e.g. overhaul)

- financial information is produced.
- if issues are found,
  - pictures are taken
  - reports are written
  - engine is tested



## Jet engine example (4)

- If problem is recurring (or suspected so)
  - a problem resolution group is established
    - existing evidence is retrieved
    - further evidence is collected
    - a learned lesson is generated
    - same problem is investigated across models



images © [www.rolls-royce.com](http://www.rolls-royce.com)

### Document Type

AROC proforma

AROC results

Development

EHM data

Emails

ONWING emails

Images

Lab findings

Monitoring Requirements

Presentations

Procedures

RCP

Risk Assessment

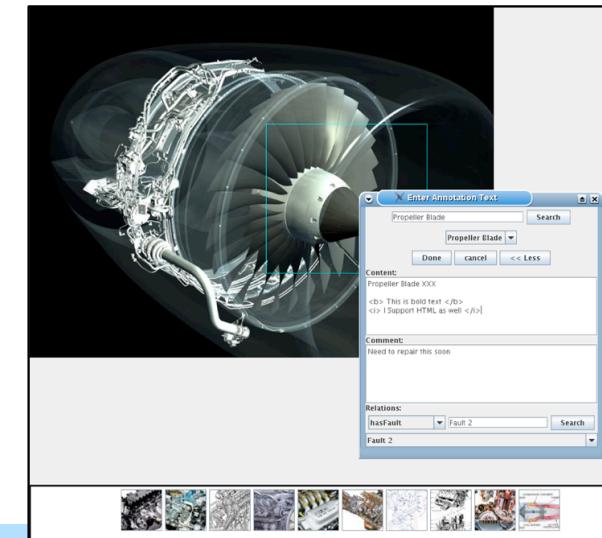
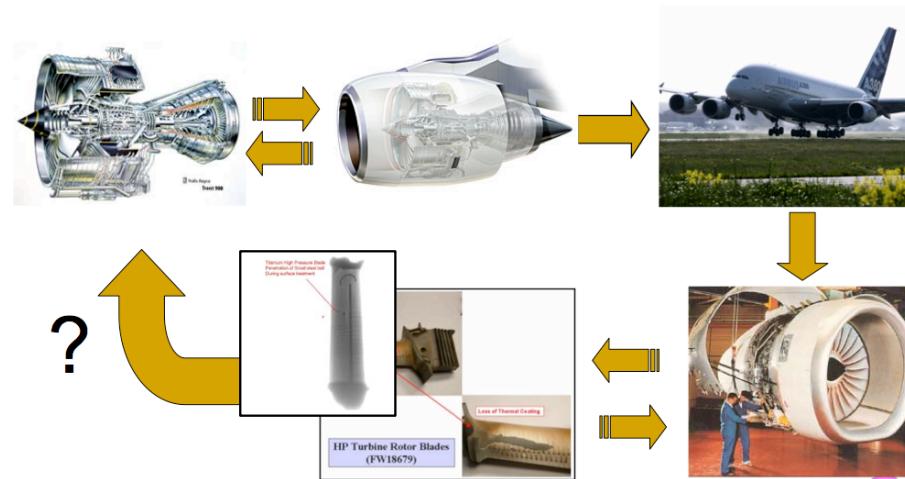
Solution Reports

Technical Reports

TS&O Reports

## Jet engine example (4)

- Lifecycle “folder” will easily sum up to several Terabytes
- Folder will contain highly interrelated information stored in different media



- Goal for Knowledge Management:
- Making information available independently from
  - Data format (structured/unstructured)
  - The archive
- Making it available for automatic processing
- Making it easily accessible and manageable despite its size



# What do we know and what we do not

- As we know, there are known knowns
  - that are things we know we know.
- We also know there are known unknowns;
  - that is to say we know there are some things we do not know.
- But there are also unknown unknowns
  - the ones we don't know we don't know

Donald Rumsfeld



ISWC 2007



# Impact of Limited Sharing

- “The lack of efficient publishing capabilities for digital content costs organisations \$750 billion annually due to wasted time spent by knowledge workers
  - seeking and capturing information necessary for them to do their jobs”
- According to S. Feldman:
  - 15%-35% knowledge worker time spent searching information
  - 50% of searches are successful (=50% fail)
  - 21% knowledge workers find information they need 85-100%

S. Feldman The high cost of not finding information. KMWorld Volume 13, Issue 3, March 2004.



# Failing factors: Technical Issues

- Information scattered in multiple repositories
  - No one really knows which information is available and/or where
  - There isn't a single access point to information
  - Even a company-wide keyword searching facility is often nonexistent
- 80-85% of a company's knowledge is unstructured
  - i.e. expressed in some forms of natural language or images/videos
- Information overload
  - Growing archives
  - Cost of storing very low
    - Video and 2D/3D image storing a reality



# Failing Factors: Human Issues

- Everyone is a database designer
  - Everyone can create a database in some hours
    - Typically ill-designed
  - Some companies are Excel-based
    - Difficult searches
    - Archives do not scale to large size
      - time bomb (!)
- Everyone is a searcher! ...but with no training
  - Where to look, how to search, how to judge quality, when to stop...



- Internal Knowledge
  - Need: capturing and sharing
  - e.g. How to design a product
- Focused external knowledge
  - Need: capturing, understanding, digesting, trusting and sharing
  - e.g. report of faults written by car garages
- External information
  - Need: capturing, understanding, contextualising, digesting, trusting and sharing
  - e.g. Information in Web pages
  - e.g. pictures provided by citizens in an emergency scenario

- Lack of Contextualisation for People, Processes and Technology
  - Current technologies tend to provide functionality in isolation from the processes and teams in which an individual knowledge worker plays a role
- Lack of Support for Cross Media and Cross Resource Sharing
  - Typically knowledge from a wide range of resources in different formats has to be brought together to solve a problem

- Knowledge Generation Requires Initial Investment
  - Systems for producing rich metadata typically require a lot of user effort, for example by annotating documents to provide training data
- The False Assumption that Knowledge is Certain
  - Metadata is usually handled as if it were certain, ignoring the possibility that it may be incorrect, inaccurate or out of date
- Knowledge Gathering Lacks Expressivity and Contextualisation
  - Simple keyword based search engines do not support the needs of knowledge workers either in the sophistication of search formulation or in longer term and exploratory aspects of knowledge gathering



# Issues with Keyword Matching

- Corporate archives are difficult to cope with because
  - Ranking cannot use document interlinking as search engines do
    - Risk of random order:
      - The % of Excite users who examined only one page of results per query in 2001: 50.5%
      - By 2001, more than 70 percent of Excite users looked at two pages or fewer
    - Documents can be very short and keyword matching has been proven not to work effectively on short documents
    - Vocabulary is reduced
      - Relevant terms tend to be very frequent
        - installed, engine, aircraft, removed, hazard, category, nrep, pse, blade, replaced, hkg, esn, csn are present in 50% of jet engine event report
      - Synonyms are not captured by keyword matching
        - Fuel Metering Unit, FMU, Metering Unit, fm701mk5, S/N3332223



## An Experiment on Jet Engine Event Reports

- 21 topics of search, e.g.
  - "How many events were caused during maintenance in 2003?"
  - "What events were caused during maintenance in 2003 due to control units?"
  - 'Find all the events associated with damage to acoustic liners following bird strike"
- Queries:
  - "what events caused during maintenance in 2003 were due to control units?"
- Translated into a set of queries given by all the possible combinations of:
  - "maintenance + 2003 + control + unit" (24 queries)



# Measures of Evaluation

- Can we trust a system answer?

$$Precision = \frac{COR}{min(ACT, maxNo)}$$

- COR= correct answer by system
- ACT= no results returned by system
- Are we finding all relevant documents?

$$Recall = \frac{COR}{EXP}$$

- EXP= no of documents relevant to the query in archive



## Results for keyword matching

- 56% of documents in the first 20 hits are relevant
  - Precision=56%
- 57% of relevant documents are in the first 2 pages
  - Recall=57%
- Keyword matching implies
  - Reading a large amount of irrelevant documents
  - Risking missing documents
  - It is impossible to count the events



# Use of Web Technologies for KM

- Most resources are available over Intranets.
  - The Web is the medium for access to multiple archives in a seamless way
    - It enables remote access independently from geographic distribution
    - it provides a common protocol for communication
    - it provides an interaction modality familiar to users
  - However:
    - Scale: dozen to hundred million documents
    - Security: risk of information leaks and hackers' attacks limits access
    - User disorientation: multiple points of access

# Knowledge Acquisition

Knowledge Reuse



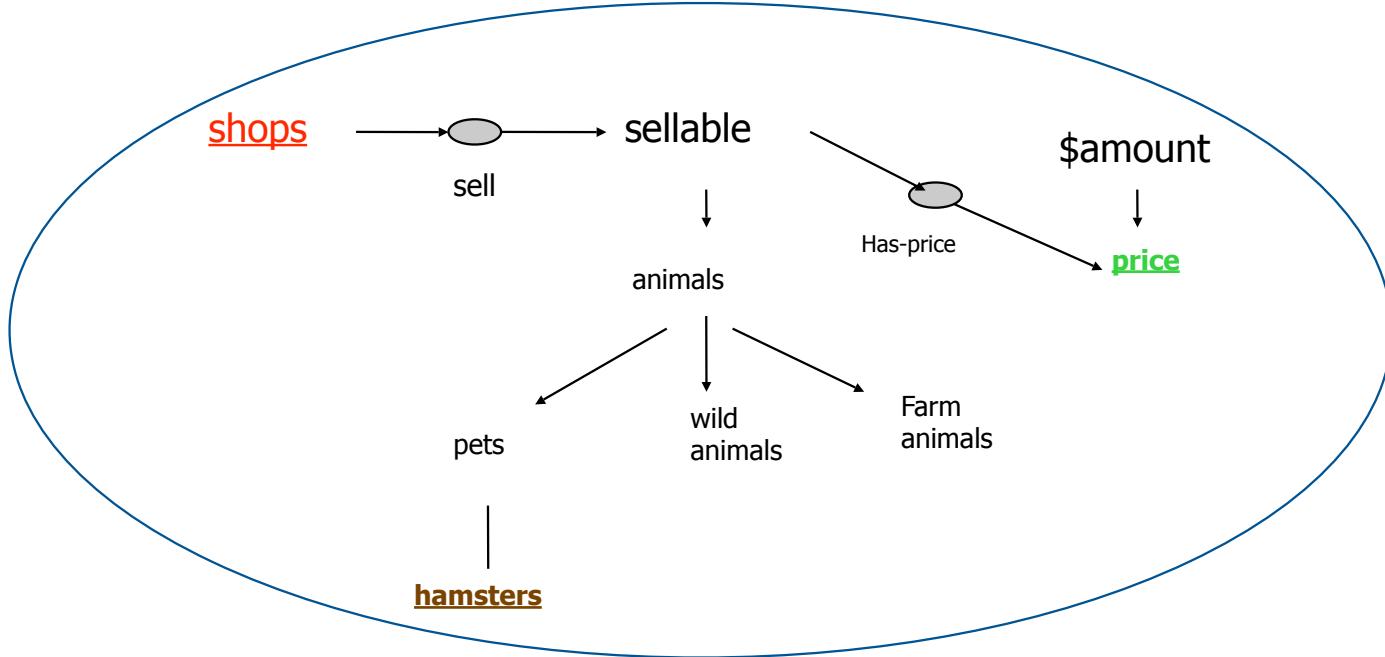
Knowledge Modelling

Knowledge Sharing

The Knowledge Life-cycle



- Knowledge Modelling
  - Ontology Engineering
    - Including also forms of: Acquisition + Reuse
- Knowledge Acquisition
  - Acquisition of information and knowledge from documents
  - Extracting and integrating information from existing archives
- Knowledge Sharing and Reuse
  - Enabling knowledge searching and process support



## Ontologies and Knowledge Engineering

- Semantic web technologies: ontologies
  - use and role of ontologies: motivations and issues
  - cost of ontologies
  - scale of ontologies



# What is an ontology

- An
    - explicit
    - shared
    - formal specification
    - of the terms in the domain
    - and relations among them
- 
- The diagram consists of four numbered points on the right side of the slide:
1. It describes a domain
  2. A formal specification
  3. Agreed by a community
  4. No implicit information
- Arrows point from each of these four points to specific items in the list of characteristics on the left. Point 1 points to 'formal specification'. Point 2 points to 'explicit' and 'shared'. Point 3 points to 'of the terms in the domain' and 'and relations among them'. Point 4 points to 'No implicit information'.

Natalya F. Noy and Deborah L. McGuinness: Ontology Development 101: A Guide to Creating Your First Ontology  
[http://protege.stanford.edu/publications/ontology\\_development/ontology101-noy-mcguinness.html](http://protege.stanford.edu/publications/ontology_development/ontology101-noy-mcguinness.html)



## Ontology (2)

- An ontology defines a common vocabulary for agents (including people) who need to share information in a domain.
- It includes machine-interpretable definitions of basic concepts in the domain and relations among them.
- It is the main means of knowledge representation and interchange of information for the Semantic Web



# Why build an ontology?

- To share common understanding of the structure of information among people or software agents
  - E.g. for communication among sites in ecommerce
- To make domain assumptions explicit
  - Avoiding hardwiring into code or database schemas
  - Can be changed without changing code
- To enable reuse of domain knowledge
  - Including serendipitous use of knowledge



## Why? (ctd)

- To separate domain knowledge from the operational knowledge
  - Operational knowledge becomes more abstract
  - What works for cars will work also for trucks by just changing underlying ontology
- To analyse domain knowledge



- Elements in ontology
  - Classes or Concepts:
    - concepts in a domain of discourse
  - Slots (or roles or properties)
    - Properties of each concept describing various features and attributes of the concept
  - Facets (or role restrictions),
    - Restrictions on slots
- Knowledge base = instances
  - Instances or Individuals
    - Instances of concepts



# Ontology and Knowledge Base

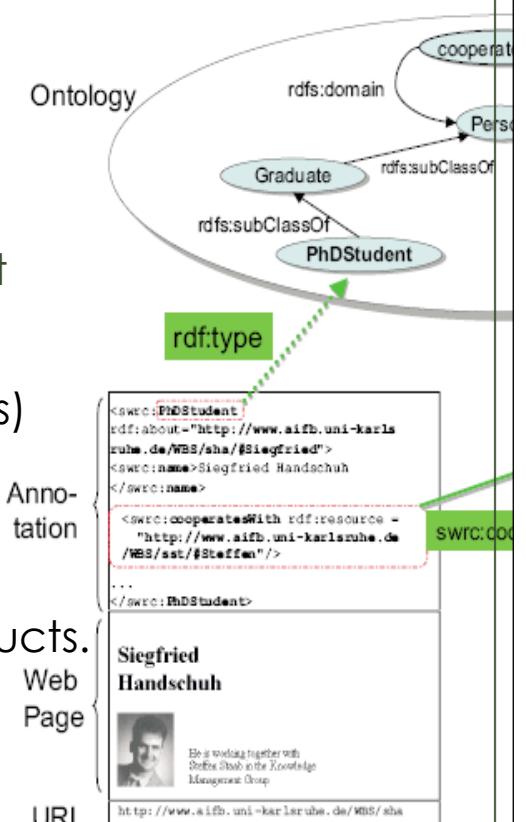
- Ontology defines the domain in abstract terms
  - Types of objects, e.g. person and companies
- Knowledge base adds the specific individuals
  - Joe is-a person,
  - ACME Ltd is-a company
- As an analogy think of
  - a database schema ~ ontology
  - actual content of database ~ instances



# Ontologies and Knowledge Management

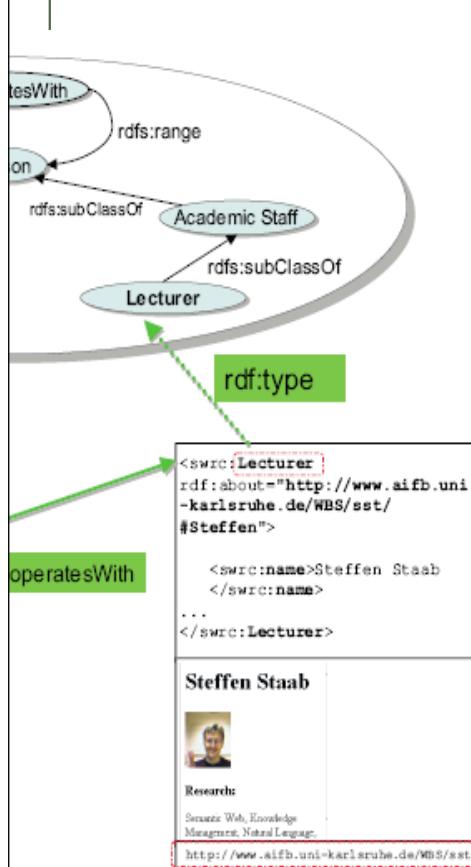
- Motivations for use:

- To represent the company's general view on the domain
  - How does the company work?
  - What is the company's official dictionary?
- As a middle layer to connect information from different information sources
  - The Web of data (as opposed to Web of documents)
- To represent communities' views of domains
  - e.g. marketing dept, customers, design and service departments have different views of the same products.
- Ontology mapping to navigate information sources
  - Mapping enables seamless communication among different worlds





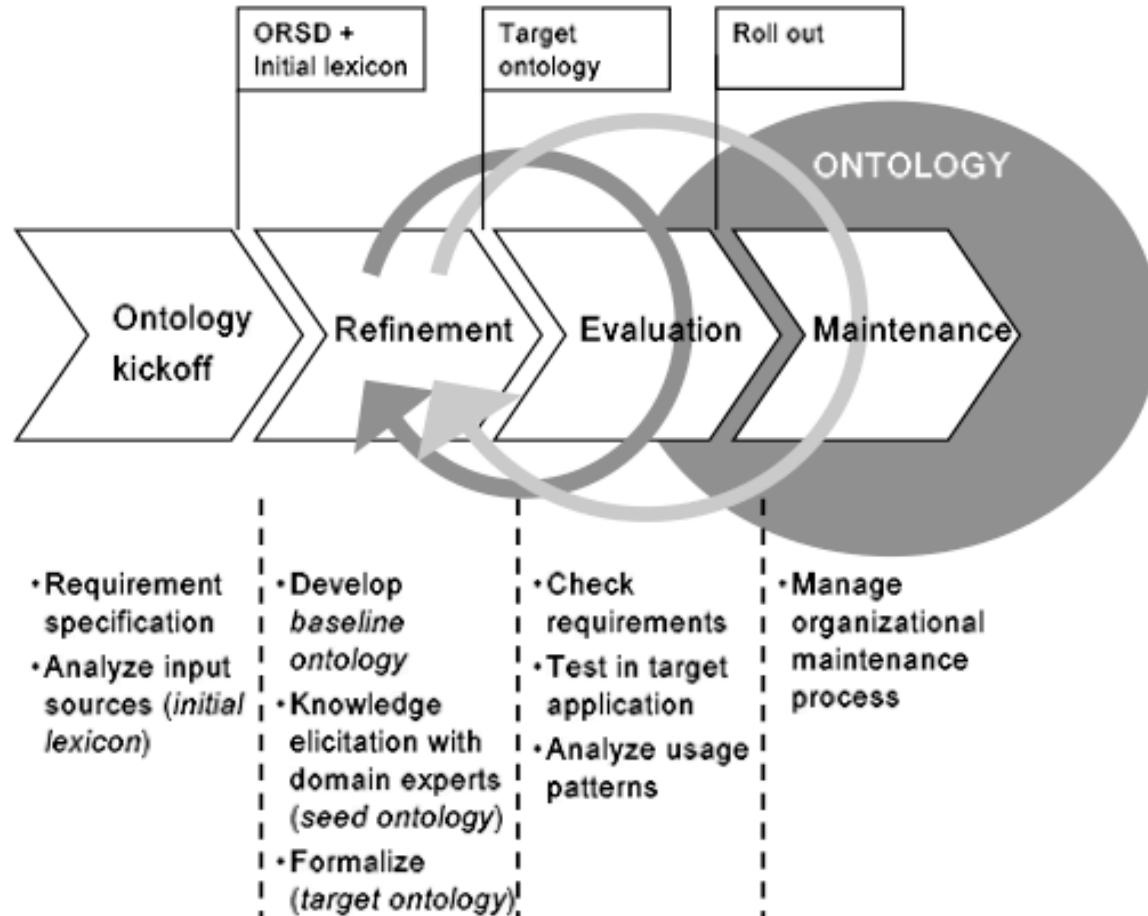
# Ontologies & KM: issues



- Cost of knowledge engineering is very high
  - It requires commitment from company management
    - Chicken-egg problem
- Knowledge engineers are difficult to find
- Lack of engineering methodologies
  - What is the cost of an ontology?
- Cost of mapping information sources to ontology middle layer is high



# Ontology Engineering



Alexander Maedche, Steffen Staab, Nenad Stojanovic, Rudi Studer, York Sure: SEMantic portAL - The SEAL approach  
In D. Fensel, J. Hendler, H. Lieberman, W. Wahlster (eds.), Spinning the Semantic Web, pp. 317-359. MIT Press, Cambridge, MA., 2003.



# Steps in Ontology Engineering

- Requirements Analysis.
  - Domain experts and ontology engineers performs a deep analysis of the project setting w.r.t. a set of pre-defined requirements.
  - Includes:
    - re-use of existing ontological sources
    - extraction of domain information from text corpora, databases etc.
  - Result: ontology requirements specification document
    - Containing competency questions describing the domain and information about its use cases, the expected size, the information sources used, the process participants and the engineering methodology



# Steps in Ontology Engineering (ctd)

- Conceptualisation.

- The application domain is modelled in terms of ontological primitives, e. g. concepts, relations, axioms.

- Implementation.

- The conceptual model is implemented in a (formal) representation language, whose expressivity is appropriate for the richness of the conceptualisation.



- Evaluation.

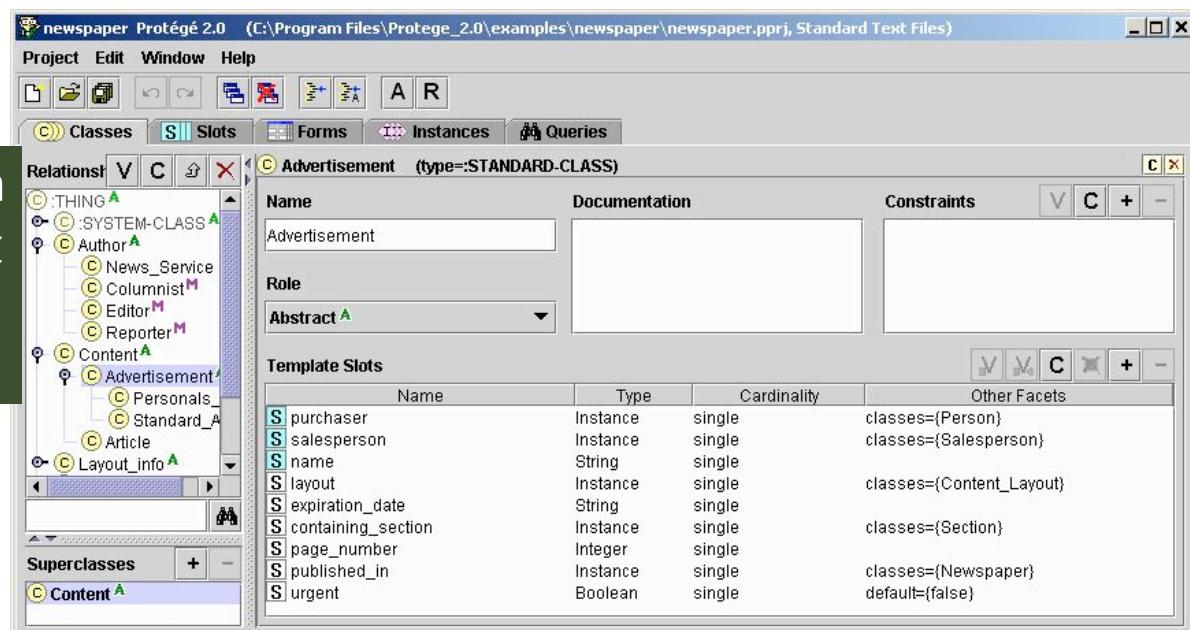
- The ontology is evaluated against the set of competency questions.
- The evaluation may be performed
  - automatically
    - if the competency questions are represented formally,
  - semi-automatically
    - using specific heuristics or human judgement
- Result is a set of modifications/refinements at the requirements, conceptualisation or implementation level



# Editing Ontology Tools: Protégé

- Protégé is a free, open source ontology editor
- Download at <http://protege.stanford.edu/>
- Protégé ontologies can be exported into a variety of formats including RDF(S), OWL, and XML Schema

Please note! Protégé is an editor! It does not support the whole knowledge engineering process





# Issues: Scale of Ontology

- Size of the ontology is a design issue for each application.
  - A jet engine has 30,000 parts;
    - Are they all to be represented as concepts?
      - What is a concept and what is an instance?
    - 20 different engine models.
    - Ontology + (knowledge base if some of them are instances)
      - large to very large if ontology contains representation of hyponymy and meronymy relations.
      - huge if we define precisely all the functional and positional relations among the 30,000 objects

Scale of KB up to billions of triples == size of ontology \* number of engine types  
\* number of events



## Size of Ontology (ctd)

- Large scale has important implications on the definition and management processes:
  - Careful hand-crafting is impossible
    - Large (possibly automatic) reuse of existing resources is necessary
  - Maintenance becomes complex
    - An ontology requires constant maintenance
    - Dedicated ontologists are uncommon in industry
      - Need of enabling users to update ontology

- Ontology in X-Media is defined as a three layer structure:
  - Foundational ontologies (e.g. Dolce)
    - Main requirements: formal precision and generality, no maintenance required by team
  - Infrastructure ontologies
    - Communication between generic tools (like services, email, etc.), including multimedia ontology
      - e.g. an email has always a sender, a recipient and a subject/body
    - Requirements: formal precision, must be updated by expert ontologist
  - Domain ontologies
    - Describe the domain
    - Need not be so precise: they are often sloppy classification schemes (e.g. controlled vocabularies) or folksonomies or thesauri
    - Can be updated by trained experts in the domain



# Cost of Ontologies

- One of the major risks in ontology development is their cost
- Cost is split into:
  - PRODUCT-RELATED COST DRIVERS: account for the impact of the characteristics of the product to be engineered (i.e. the ontology) on the overall costs.
  - PERSONNEL-RELATED COST DRIVERS emphasize the role of team experience, ability and continuity w.r.t. the effort invested in the engineering process:
  - PROJECT- RELATED COST DRIVERS relate to overall characteristics of an ontology engineering process and their impact on the total costs

Elena Paslaru Bontas, Christoph Tempich, York Sure : OntoCom: A Cost Estimation Model for Ontology Engineering  
In: Proceedings of the 5th International Semantic Web Conference (ISWC 2006), November 5-9, 2006, Athens, GA, USA, LNCS. Springer.



# Product Related Cost Drivers

- Domain Analysis Complexity
  - to account for those features of the application setting which influence the complexity of the engineering outcomes,
- Conceptualisation Complexity
  - to account for the impact of a complex conceptual model on the overall costs,
- Implementation Complexity
  - to take into consideration the additional efforts arisen from the usage of a specific implementation language



# Product costs (ctd)

- Instantiation Complexity
  - to capture the effects that the instance data requirements have on the overall process
- Required Reusability
  - to capture the additional effort associated with the development of a reusable ontology
- Evaluation Complexity
  - to account for the additional efforts eventually invested in generating test cases and evaluating test results
- Documentation Needs
  - to state for the additional costs caused by high documentation requirements

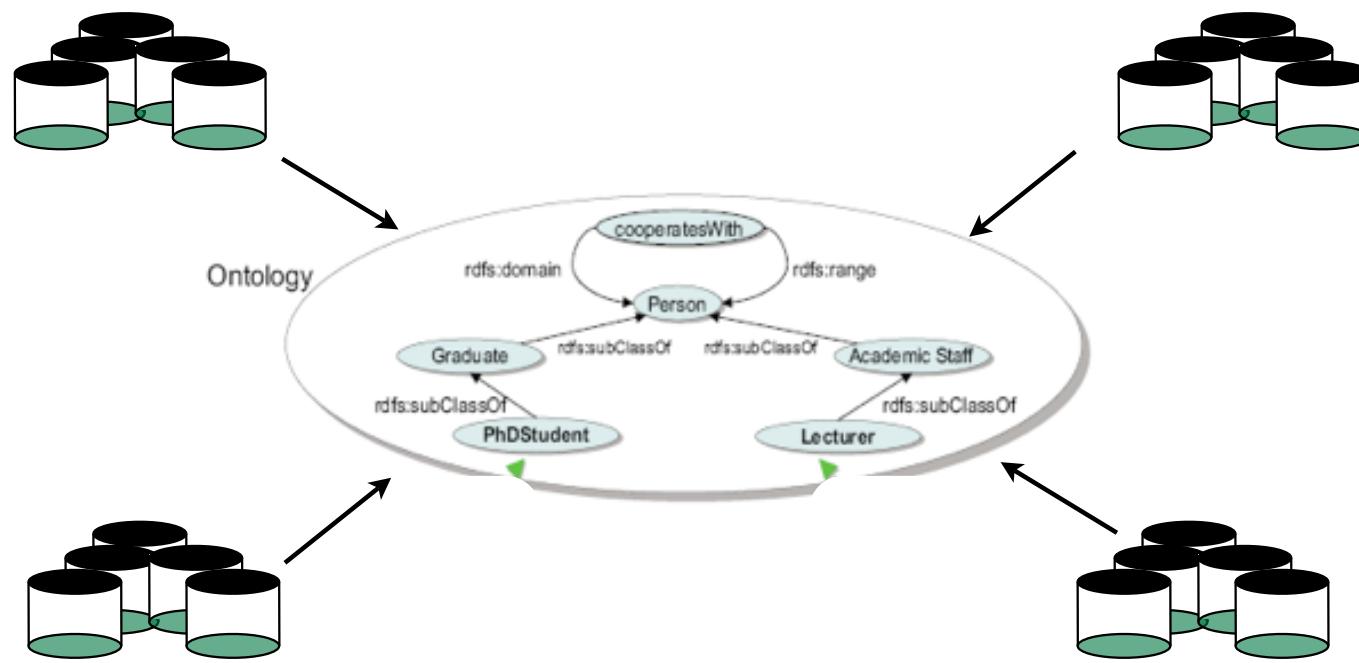


- Ontologist/Domain Expert Capability
  - to account for the perceived ability and efficiency of the single actors involved in the process (ontologist and domain expert) as well as their teamwork capabilities
- Ontologist/Domain Expert Experience
  - to measure the level of experience of the engineering team w.r.t. performing ontology engineering activities,
- Language/Tool Experience
  - to measure the level experience of the project team w.r.t. the representation language and the ontology management tools,
- Personnel Continuity
  - to mirror the frequency of the personnel changes in the team.



# Project-Related Cost Drivers

- Support tools for Ontology Engineering
  - to measure the effects of using ontology management tools in the engineering process
- Multisite Development
  - to mirror the usage of the communication support tools in a location-distributed team.



Towards the Web of Data

ISWC 2007

© Fabio Ciravegna, University of Sheffield



# Web of data and KM

- It is often said that the Semantic Web is rapidly evolving towards a Web of data
  - As opposed to Web of documents
  - Web of documents
    - Feature: Web pages for human reading must be associated with machine readable data
    - Requirement: capturing information from unstructured documents
  - Web of data
    - Feature: a large amount of information is contained in databases
      - e.g. eBay, Amazon.com, etc.
    - Requirement: mapping existing resources (e.g. database schemas) to an ontology and provide metadata together with data automatically
      - e.g. Semantic Amazon.com



# Reusing existing resources

- Reuse of existing sources
  - e.g existing ontologies/KB, Gazetteers or Database schemas or database content
- Existing resources:
  - Uncertainty
    - existing sources were generally invented for human reading
      - Are not 100% certain.
      - Trusting all background knowledge provided may decrease performance
  - Incompleteness
    - Background knowledge may only be available for part of the problem space.
  - Inconsistency
    - Pieces of knowledge from different sources can be conflicting



# (Re-)Using Different Ontologies

- Different communities use different domain representations
  - Design department uses designs and associated part lists
  - Service department uses illustrated part catalogue
    - They largely overlap but:
      - Use different URLs (part numbers)
      - Use different descriptions of the domain
        - In service what cannot be repaired is not described in details
        - In manufacturing what is outsourced is not described in details
  - Ontologies as solution to map different resources (e.g. database schemas) to a common view

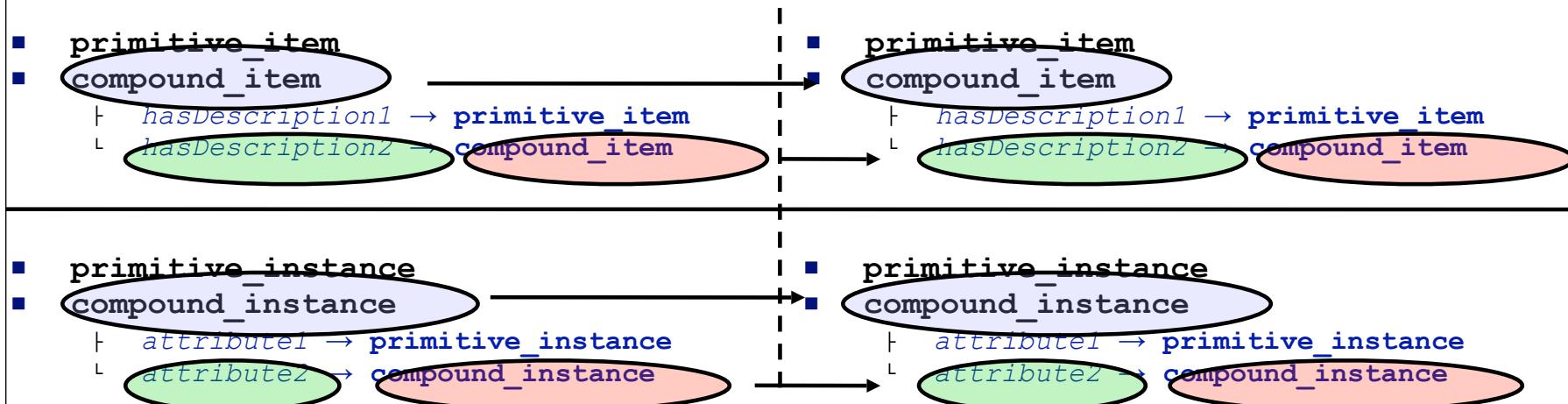


# Issues in mapping schemas

- The main difference between an db schema and an ontology
  - Ontology makes assumptions explicit
  - Database schema leave assumptions implicit
- Issue:
  - Mapping requires:
    - Making assumptions explicit (turn schema into ontology)
    - Mapping the two ontologies
      - Three main approaches
        - Top down (next slides)
        - Bottom up (starting from instances - see slides on integration)
        - Mixed



# Ontology Mapping



- **class to class** mapping (*classMapping*)
- **attribute to attribute** mapping (*attributeMapping*)

A Slide by Adrian Mocan

<http://www.inrialpes.fr/exmo/people/zimmer/SDK-meeting/Presentations/Adrian%20Mocan%20-%20WSMX%20Data%20Mediation.ppt>



# Ontology Mapping: Issues

- Difficulties in mapping concepts and properties
  - Non overlapping
    - Concept/Property X in DB Schema does not exist in Domain Ontology
      - Solution: extension of Domain Ontology to include it
    - Concept/Property X in Domain Ontology does not exist in DB Schema
      - Solution: none; missing information
  - Partially Overlapping
    - Concept/Property exists but with a slightly different definition
    - Next slide



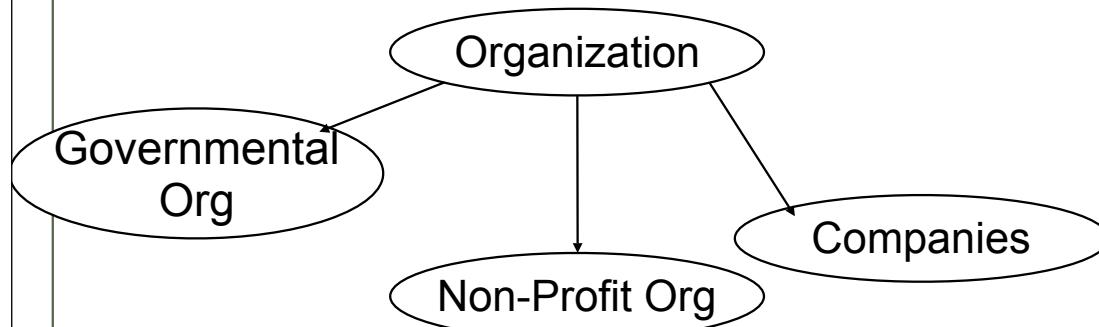
# Partially Overlapping Objects (1)

Specific to generic

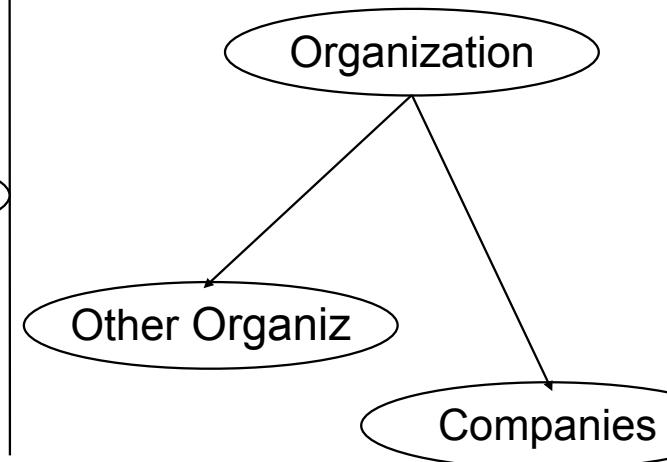
- Easy case

- Collapse all instances and subtypes of Governmental\_Org and Non-profit\_Org into Other\_Organiz

Source (existing ontology)



Target (new ontology)





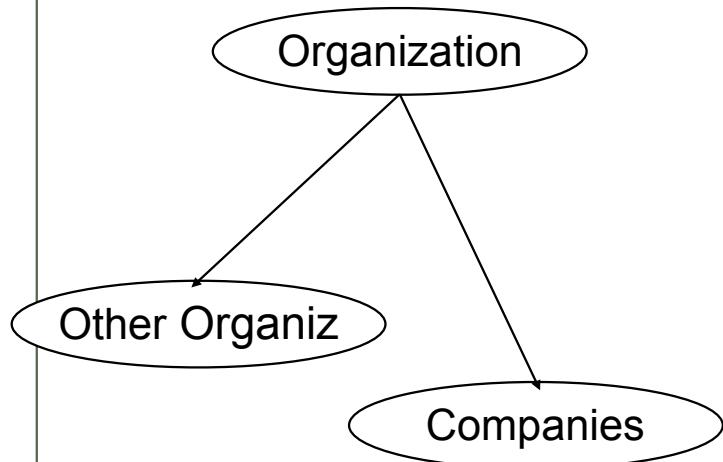
# Partially Overlapping Objects (2)

## Generic to Specific

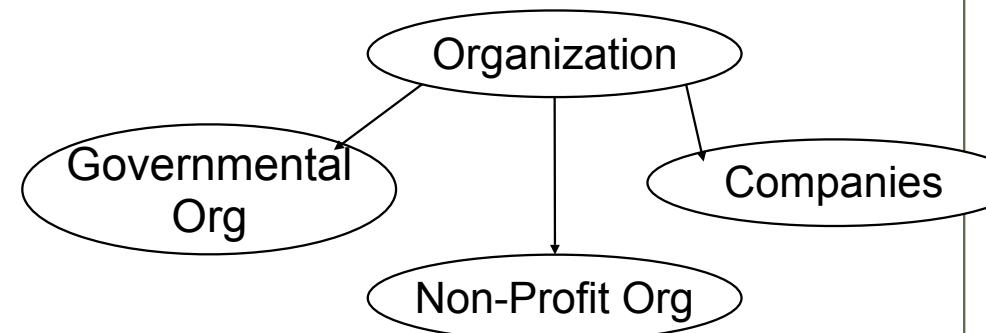
- Difficult case

- How do we divide the Other\_Organiz into Governmental\_Org and Non-profit\_Org into?
- Manual mapping of instances/subtypes or modification of ontology

Source (existing ontology)



Target (new ontology)



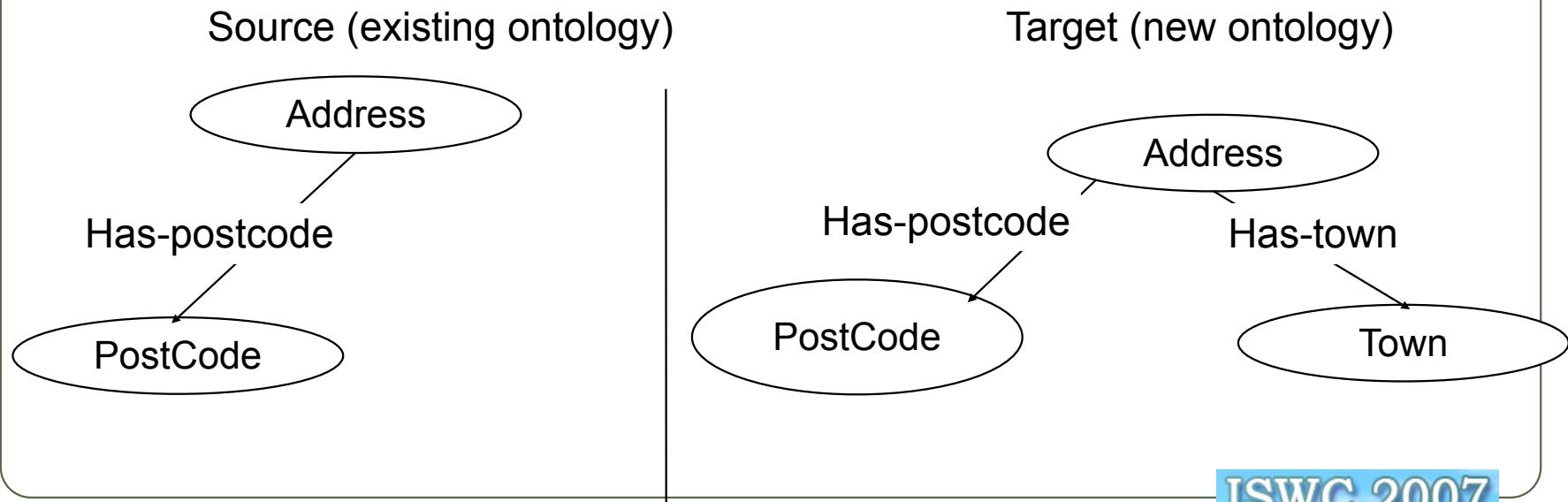


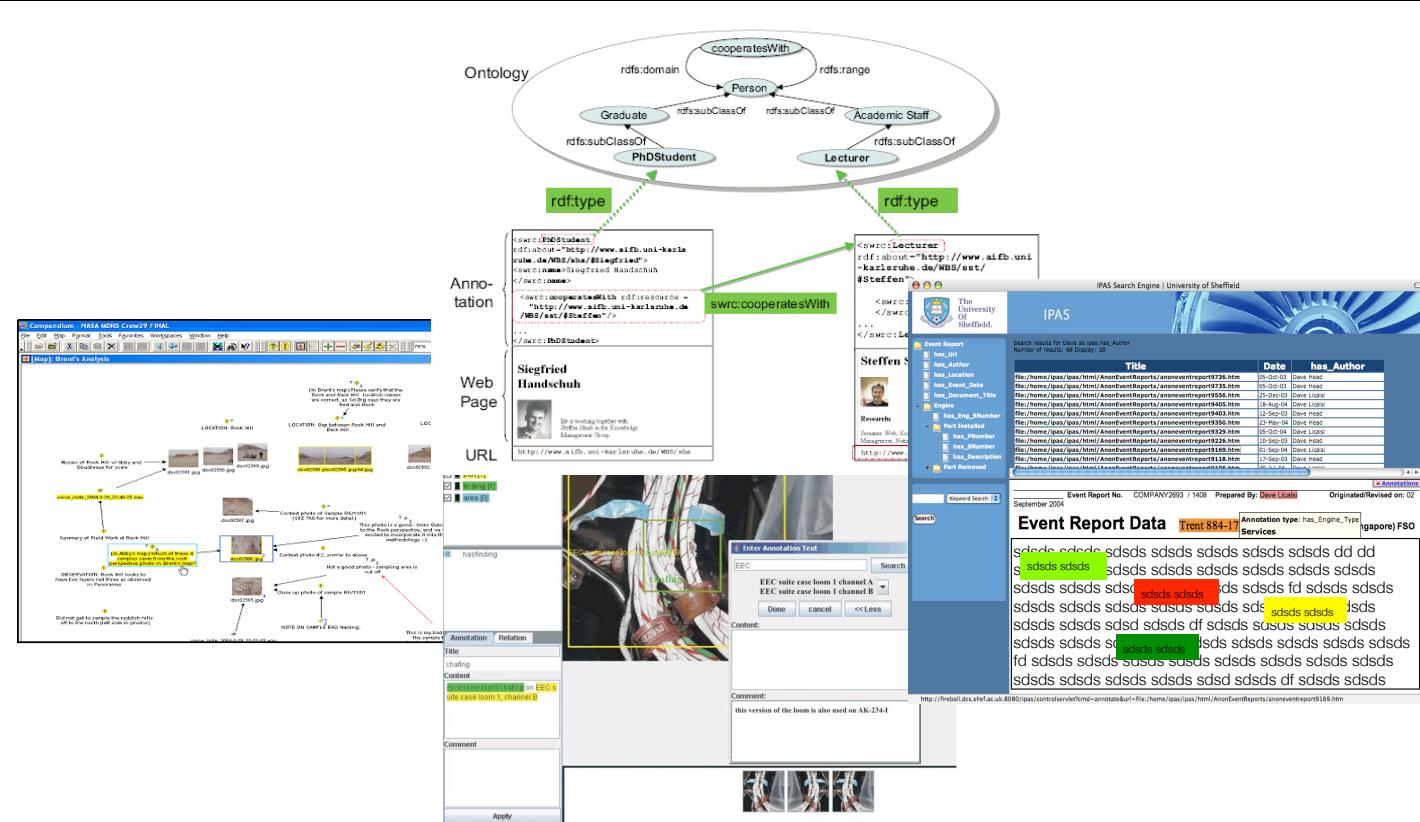
# Partially Overlapping Objects (3)

## A more complex case

- Difficult case

- In the DB the address implicitly includes the town
  - E.g. substring
- In the Domain Ontology the town is explicitly mentioned
- No easy way to map





## Requirements for Knowledge Acquisition

- issues in knowledge acquisition:
  - acquiring: what and what for?



# Knowledge Acquisition

- Collecting and aggregating multimedia knowledge to make it available for
  - sharing and reuse
    - From document management to knowledge management
  - for integration
- Approaches
  - At source: helping people capturing knowledge when produced
  - On legacy documents, pictures, data:
  - Annotation services

The left screenshot shows the Active Media Version 1.8 interface, which includes a video player window displaying a person working on a car engine, overlaid with green and yellow annotation boxes. A floating window titled 'Enter Annotations' contains the text 'EEC suite of EEC suite car' and a search bar. Below the video player, there's a list of annotations: 'has\_annotation on EEC's site base room 1, channel B'. The right screenshot shows the IPAS Search Engine interface, featuring a search results page for 'Dave as (as:has\_Author)'. The results table includes columns for Title, Date, and has\_Author, listing various URLs and dates. Below the table, a specific event report is shown with details like 'Event Report No. COMPANY2693 / 1408', 'Prepared By: Dave Lolis', and 'Originated-Revised on: 02 September 2004'. The report content is heavily redacted with 'sdsds sdsds' placeholder text.



# Requirements for KA: Web Size

- August 2005: Yahoo claimed to cover 20 billion pages
  - 19.2 billion web documents,
  - 1.6 billion images,
  - 50 million audio and video files.

<http://google.weblogsinc.com/2005/08/08/yahoo-upgrades-indexes-claims-20-billion-objects/>

- Almost 375 petabytes (or 787.5 billion photographs) are produced each year
  - almost 2 times all printed material
  - yearly growth rate of 5%
    - Highest growth rate among different data types

Brilakis, I.: Content based integration of construction site images in aec/fm model based systems,  
PhD thesis, University of Illinois at Urbana-Champaign, 2005

Large intranets are following the Web's trends. Going towards hundreds of millions documents (web at end of the 90s)



# Requirements for KA: Cross media

- Evidence is often distributed in different media;
- Knowledge in one medium does not carry the full evidence

## Battery Exchange Program iBook G4 and PowerBook G4

Apple has determined that certain lithium-ion batteries containing cells manufactured by Sony Corporation of Japan pose a safety risk that may result in overheating under rare circumstances.

The affected batteries were sold worldwide from 2003 through August 2006 for use with notebook computers: 12-inch iBook G4 and PowerBook G4 and 15-inch PowerBook G4.

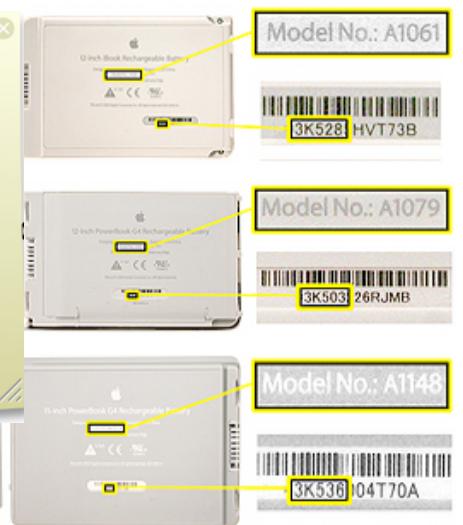
Apple is voluntarily recalling the affected batteries. Apple has initiated a worldwide exchange program for eligible customers with a new replacement battery at no charge. This program is being conducted in cooperation with the U.S. Consumer Product Safety Commission (CPSC) and other international safety agencies.

### Identifying your battery

Please use the chart below to identify the model number and serial numbers that apply to your iBook G4 or PowerBook. If the first 5 digits of your battery's serial number fall within the noted range, replace the battery immediately.

To view the model and serial numbers labeled on the bottom of the battery, you must remove the battery from the computer. The battery serial number is printed in black or dark grey lettering beneath a barcode. See photos below.

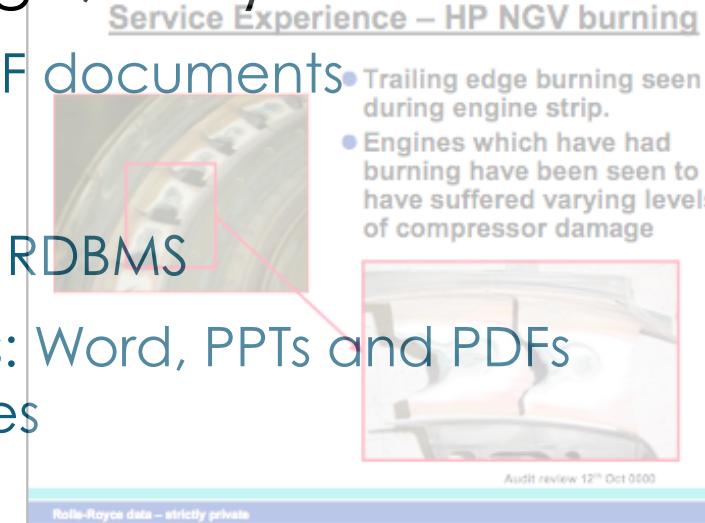
this case is no longer valid because we have introduced Service Note 3445 which requires replacement of component



From Deliverable D8.2

- Typical data objects (text, image, raw)

- Text formats: Word, Excel, PPT and PDF documents
- Images: Jpeg and Gif
- Raw data: Measurements stored in a RDBMS
- Cross-media: Compound documents: Word, PPTs and PDFs containing both text and Jpeg images
  - Portions semantically related to each other within the same physical document
  - Information contained in just one modality is insufficient
  - Cross-media knowledge acquisition techniques needed in order to capture and manage all of the explicit and implicit knowledge



A way of thinking

The inside of every Yaris looks clean and sophisticated. The controls and instruments are ergonomically designed and positioned. And there is a distinct lack of clutter thanks to the innovative storage system.

Every inch of space has been used to provide a range of practical storage compartments that ensure items are neatly and safely stored. The glove box, for example, has two sections, while a useful tray under the passenger seat keeps the contents out of sight.

There are two large storage pockets next to the centre console, two front-door storage bins, cup-holders and a passenger side storage tray. Additionally, the T Sport and T Sport models have pockets in the back of the front seats.

The Yaris – neatly taking care of everything.



Storage different items





# Requirements for SW: Robustness

- Required robustness in:

- Knowledge representation

- e.g. uncertainty and dynamic phenomena modelling

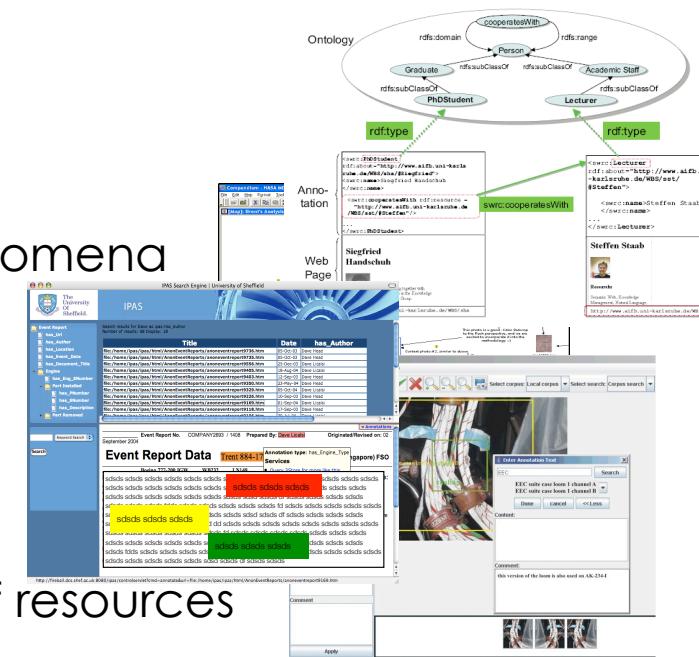
- Against unexpected situations:

- Coping gracefully with downtime of resources

- What if a document disappears/server is down? (reasoning)

- Preventing that a crash of an individual components leads to a whole system down.

- Dealing intelligently with error propagation through the cascade of processors





## Hamsters

FOUR-DAY FORECAST

Wed.	Thu.	Fri.	Sat.
shower	cloudy	shower	t-storms
HIGH 80 F 26 C	HIGH 82 F 27 C	HIGH 84 F 28 C	HIGH 86 F 30 C
LOW 65 F 18 C	LOW 70 F 21 C	LOW 69 F 20 C	LOW 67 F 19 C
Pressure: 29.94 in. (1015 hPa)			
Wind: N at 9 mph (14 kph)			
Sunrise: 5:10 a.m.			



£6

Larger View

OT Movie

Temp2

Larger View

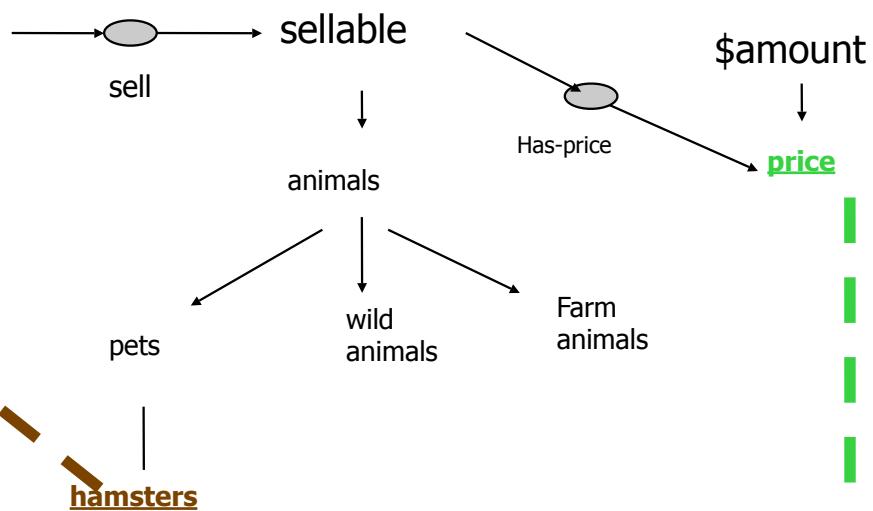
Tomorrow

OT Movie

Forecast

Four-day forecasts for 7,200 cities worldwide

shops



# SW for Knowledge Acquisition

- user centred methodologies and tools for text and image annotation
- automatic methodologies and tools for text annotation



- Aims:

- To acquire knowledge within and across media in a rich, semantically-oriented way
- Outcome of acquisition technologies is a semantic representation of the content (conceptualisation) to be used for knowledge management purposes
- Enrichment of multimedia documents with layers of manually or automatically generated annotation is the main medium of associating conceptualisations to resources



# Making Content Available

- 3 main methods of making the content available:
  - Ontology-based annotations
  - Free text annotations - Braindumps
  - Document enrichment

Vitaveska Lanfranchi, Fabio Ciravegna and Daniela Petrelli: Semantic Web-based Document: Editing and Browsing in AktiveDoc, 2nd European Semantic Web Conference, Crete, June 2005

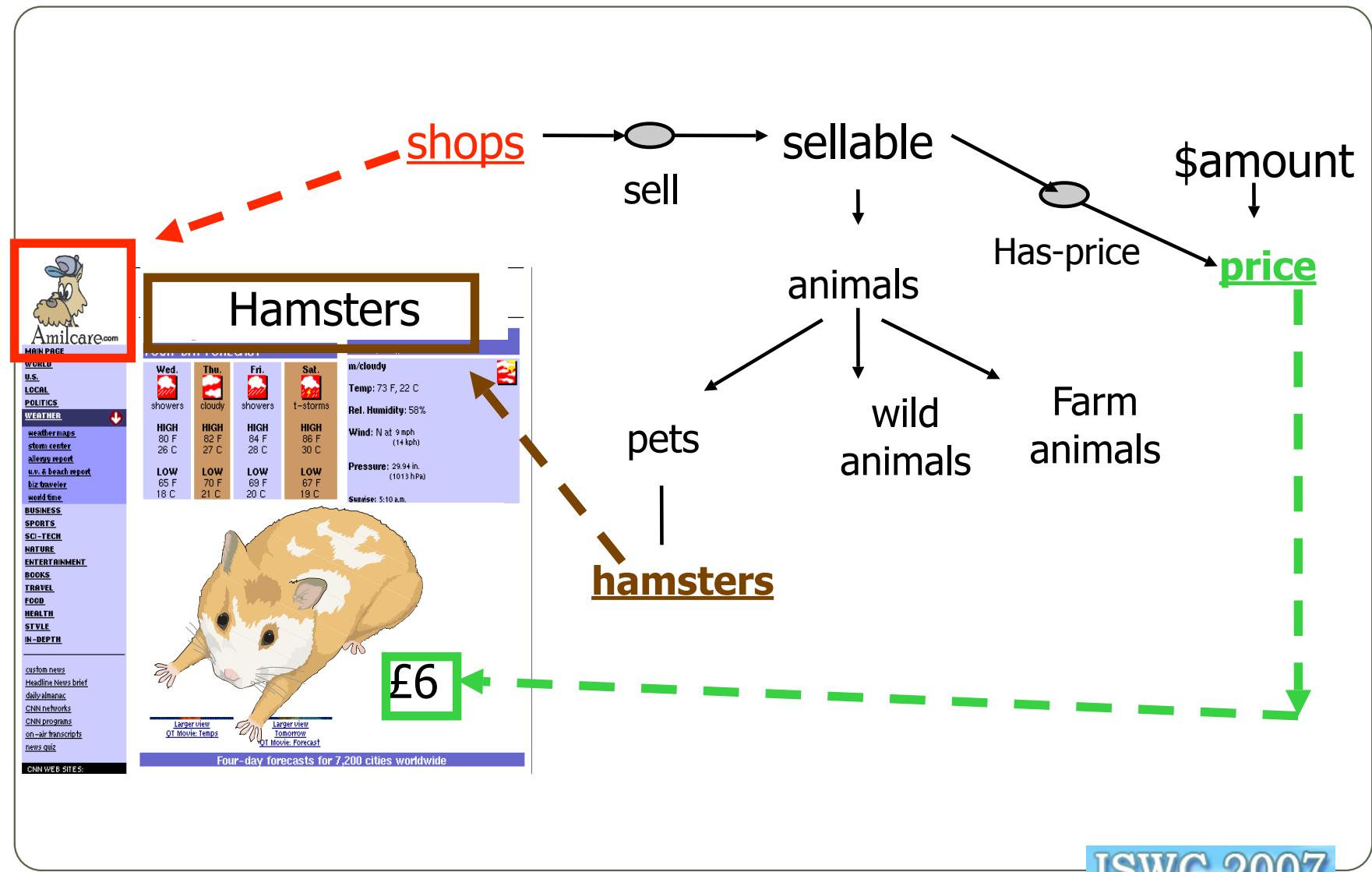


# Ontology-based annotations

- Marking up contained information
  - Portions of documents associated to objects in ontology
    - Allows:
      - Ontology-driven processing
      - Services based on ontology will be able to use information
    - Ontomat/CREAM (Staab et al 2001)
    - Melita (Ciravegna et al. 2002)
    - SemTag and Seeker (Dill et al. 2003)
    - ...and many others...



# Ontology-based Annotation





- Adding knowledge to documents

- Via free – text comments (as in Word)
- This is called braindump
  - The final document is just the final solution
  - E.g. the project for a new Jet Engine
    - During the discussion the working group will consider many alternative solutions
    - Those not selected are not in the final project
    - When next jet engine is designed, the group needs to know
      - What solutions were tried (use of titanium)
      - Why they were not adopted (e.g. too high a cost)
      - If the analysis is still true (titanium cost has decreased)
  - Adding further comments and associate Information
    - Annotea (Barstow et al 2001)
    - Semantik (Gilardoni et al 2004)



# Braindump in a Legal Scenario

## Plain English in the Twenty Types of Legal Documents

By George H. Hathaway

To eliminate legalese from legal documents, you must do more than just criticize legal writing in general. You must zero in on organized way, rather than randomly selecting a sample of legal writing and critiquing it. Therefore, we have grouped all legal Figure 1.

We have published many articles about these 20 types of documents in our Plain English theme issues of November 1983 a present. Now we will discuss: 1) the level of clarity at which each type of document is presently written; 2) improvements, if are still unresolved...and why.

### Resolutions1

Whereas, it is a privilege to congratulate Kyle David Gibbs [REDACTED] the rank of Eagle Scout...

Resolutions are written by the Michigan House and Senate to express a position on an issue. The resolutions are published simple, and therefore the main body of the resolution is usually clear. [REDACTED] solutions have always contained one persi

Trying to eliminate this word from resolutions illustrates the intractability of those who write legalese on purpose. So far the I make up only 15% of the Legislature. This means six out of [REDACTED] legislators are not lawyers. Furthermore, the clerks of the doesn't the Michigan Legislature eliminate Whereas from its resolutions? There are still too many people (lawyers, non-lawye symbol of power and prestige. [REDACTED]

### Statutes2

Michigan statutes are written by the Legal Division of the Legislative Service Bureau. (Of course, the drafters often do not have the Michigan Legislative Service pamphlets and each year in the Public and Local Acts of Michigan. In 1994 we reviewed the Service Bureau for its work.

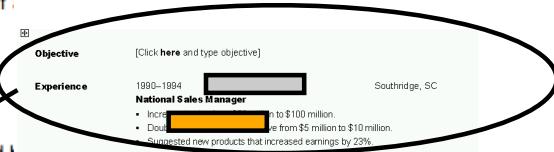
### Executive Orders3

Whereas, Article V, Section 2, of the Constitution of the State of Michigan of 1963 empowers the Governor to make changes units which he considers necessary for efficient administration; and....

These orders are written by the Governor's legal counsel and are published each month in the Michigan Register. They are di format for executive orders has not changed in the last 100 years. It has always contained much legalese. We are told that if orders carried as much weight as legislative statutes, administrative rules, or case opinions. They believe that if the orders a have a better chance of being followed. But this is what critics of legal writing have always chargedNthat lawyers write legale

(note: this is not a real legal document!)

## Why we used these references



## Why we DID NOT use other references





- Adding knowledge to documents (ctd.)
  - Document enrichment: helping connecting the document to the rest of the knowledge
    - Associating Services
      - Magpie (Dzbor et al. 2004)
    - Connected to other documents
      - e.g. Automatic generation of hyperlinks
      - COHSE (Goble et al. 2001)



## NASA GISS: A Stratospheric "Clock" to Measure Upper Atmosphere Circulation - Microsoft Internet Explorer provided by magpie

File Edit View Favorites Tools Help

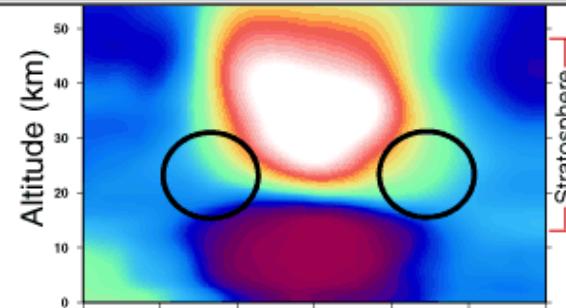
Back Forward Stop Home Search Favorites Media

Address http://www.giss.nasa.gov/research/intro/koch\_01/ Go Links >

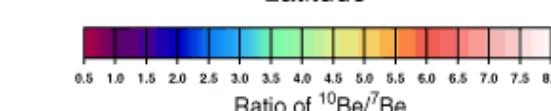
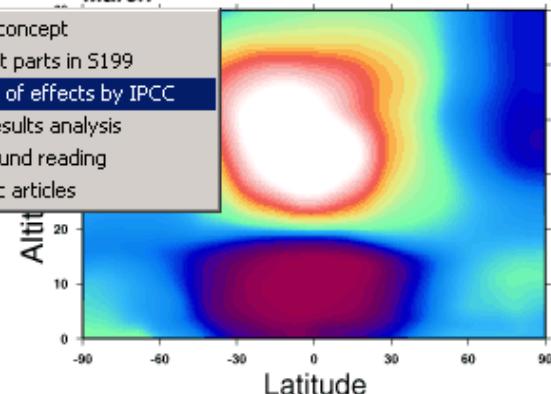
Magpie Climatology Meteorology Physics Chemistry

collision of high-energy particles from space with nitrogen atoms in the atmosphere. Most tracer production occurs between about 30°–70° latitude in both hemispheres of the lower stratosphere, as indicated by the circled regions on the figure. These tracers, which are borne on aerosol particles, are removed from the stratosphere by radioactive decay. While beryllium-7 decays relatively quickly, with a half-life of 53 days,  $^{10}\text{Be}$ 's decay rate is negligible. The only sink for  $^{10}\text{Be}$  occurs after it enters the troposphere, where the radionuclides are efficiently removed by precipitation. Therefore, if we look at the ratio of  $^{10}\text{Be}/^{7}\text{Be}$  as air moves from the midlatitude production region to other parts of the stratosphere, the ratio will generally increase, as  $^{7}\text{Be}$  decays. Thus, the  $^{10}\text{Be}/^{7}\text{Be}$  acts as a "clock" of airmass age.

The figure shows the  $^{10}\text{Be}/^{7}\text{Be}$  ratio calculated in the GISS general circulation model (GCM) during January and March. In the tropical stratosphere, air rises from the troposphere and continues to ascend, but exchange with higher latitudes is inhibited. The  $^{10}\text{Be}/^{7}\text{Be}$  ratio is very high (white region) since slow penetration of air from the mid-latitude production region allows much of the  $^{7}\text{Be}$  to decay. During the early northern hemisphere spring, air from the lower tropical stratosphere moves to higher latitudes relatively quickly. The result is the green blob of relatively high  $^{10}\text{Be}/^{7}\text{Be}$  air at



March



$^{10}\text{Be}/^{7}\text{Be}$  ratio calculated in the GISS general circulation model during January and March. Circled areas indicate maximum

Done

ISWC 2007



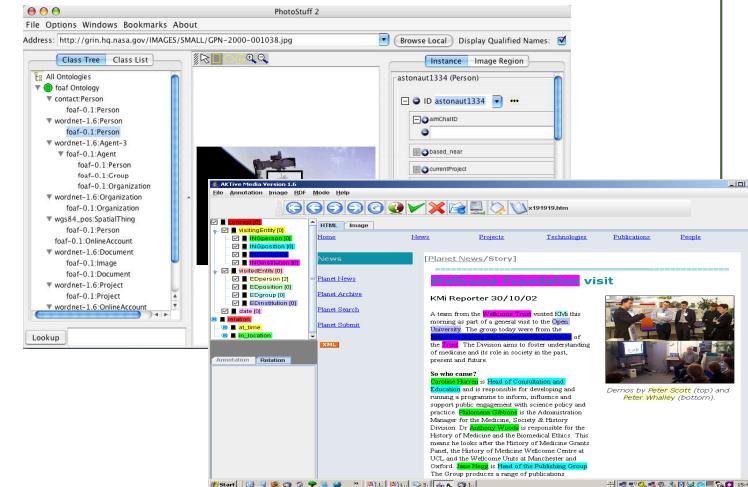
# Input & Output

- Input to the KA technologies
  - Ontologies (MMO, domain ontology),
  - Background knowledge (gazetteers, etc.)
  - Normalised document representation
  - Medium to extract from (text, images, data, videos,...)
- Output
  - Evidence represented in terms of conceptual information
    - Evidence used by other modules as background conceptual knowledge, i.e. pre-existing knowledge
    - Evidence in the form of uncertain output



# Ontology-based Annotation

- The way to annotate pages is to:
  - Select an ontology
  - Define statements to represent meta-data about the document
- Manual Annotation
  - Annotation can be performed by:
    - Domain expert
- User-friendly tools for annotation
  - Cream (Handschuh *et al.* 2002)
  - Melita (Ciravegna *et al.* 2002)
  - Photostuff (Hendler *et al.* 2005)
  - AktiveMedia (Chakravarthy *et al.* 2006)





- Enables semi-automatic annotation across texts and images
- The interface enables
  - HTML editing
  - Annotation of documents in RDF based on an OWL ontology
- Types of annotations
  - Concepts / Relations
- SW: Annotation:
  - Selection of concept/relation and highlighting of text is the way in which annotation is performed

<http://www.dcs.shef.ac.uk/~ajay/html/cresearch.html>

**AKTive Media Version 1.6**

File Annotation Image RDF Mode Help

x191919.htm

**Ontology panel**

concept [0] visitingEntity [0] INGperson [0] INGposition [0] INGgroup [0] INGinstitution [0] visitedEntity [0] EDperson [2] EDposition [0] EDgroup [0] EDinstitution [0] date [0] relation at\_time in\_location

**Document panel**

Text is selected and dropped into a concept in the ontology

HTML Image

News Projects Technologies Publications People

News

Planet News Planet Archive Planet Search Planet Submit

Wellcome Foundation visit

KMi Reporter 30/10/02

A team from the Wellcome Trust visited KMi this morning as part of a general visit to the Open University. The group today were from the Medicine, Society and History (MSH) Division of the Trust. The Division aims to foster understanding of medicine and its role in society in the past, present and future.

So who came?

Caroline Hurren is Head of Consultation and Education and is responsible for developing and running a programme to inform, influence and support public engagement with science policy and practice. Philomena Gibbons is the Administration Manager for the Medicine, Society & History Division. Dr Anthony Woods is responsible for the History of Medicine and the Biomedical Ethics. This means he looks after the History of Medicine Grants Panel, the History of Medicine Wellcome Centre at and the Wellcome Units at Manchester and Oxford. Jane Hogg is Head of the Publishing Group. The Group produces a range of publications

Start

15:40



# Contextual Annotation of Images and Text

AKTive Media Version 1.6

File Annotation Image RDF Mode Help

x191919.htm

concept [0]

- visitingEntity [0]
  - INGperson [0]
  - INGposition [0]
  - INGgroup [0]
  - INGinstitution [0]
- visitedEntity [0]
  - EDperson [2]
  - EDposition [0]
  - EDgroup [0]
  - EDinstitution [0]
- date [0]
- relation
- at\_time
- in\_location

HTML Image

Story

Enter Annotation Text

Martin Dzbor

Search

Martin Dzbor

Martin Dzbor

Simon Buckingham Shum

This vision aims to foster understanding of the role in society in the past,

Head of Consultation and responsible for developing and me to inform, influence and agement with science policy and a Gibbons is the Administration Medicine, Society & History my Woods is responsible for the e and the Biomedical Ethics. This er the History of Medicine Grants of Medicine Wellcome Centre at some Units at Manchester and Oxford. Jane Hogg is Head of the Publishing Group. The Group produces a range of publications

Demos by Peter Scott (top) and Peter Whalley (bottom).

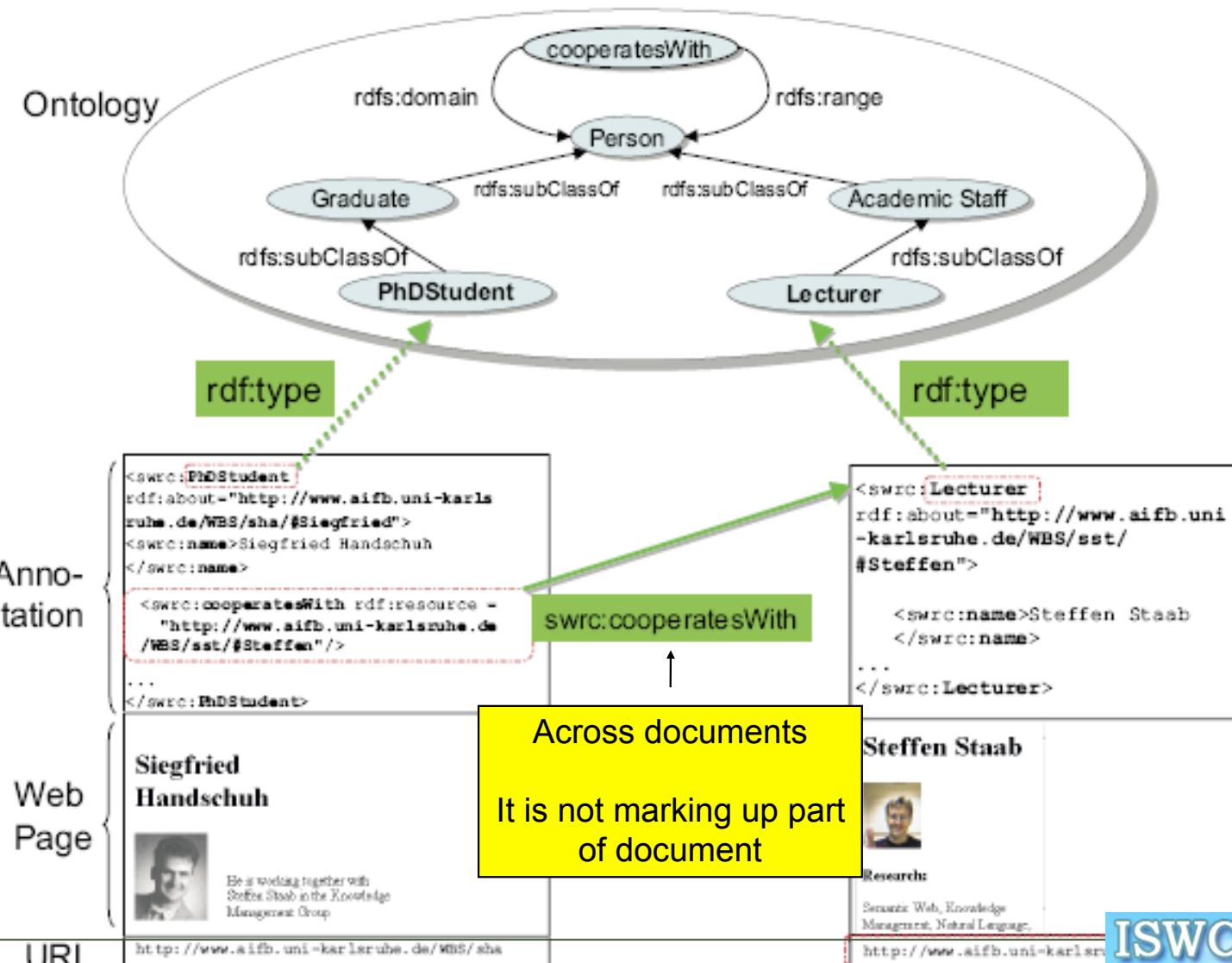
Annotation Relation

Start

ISWC 2007



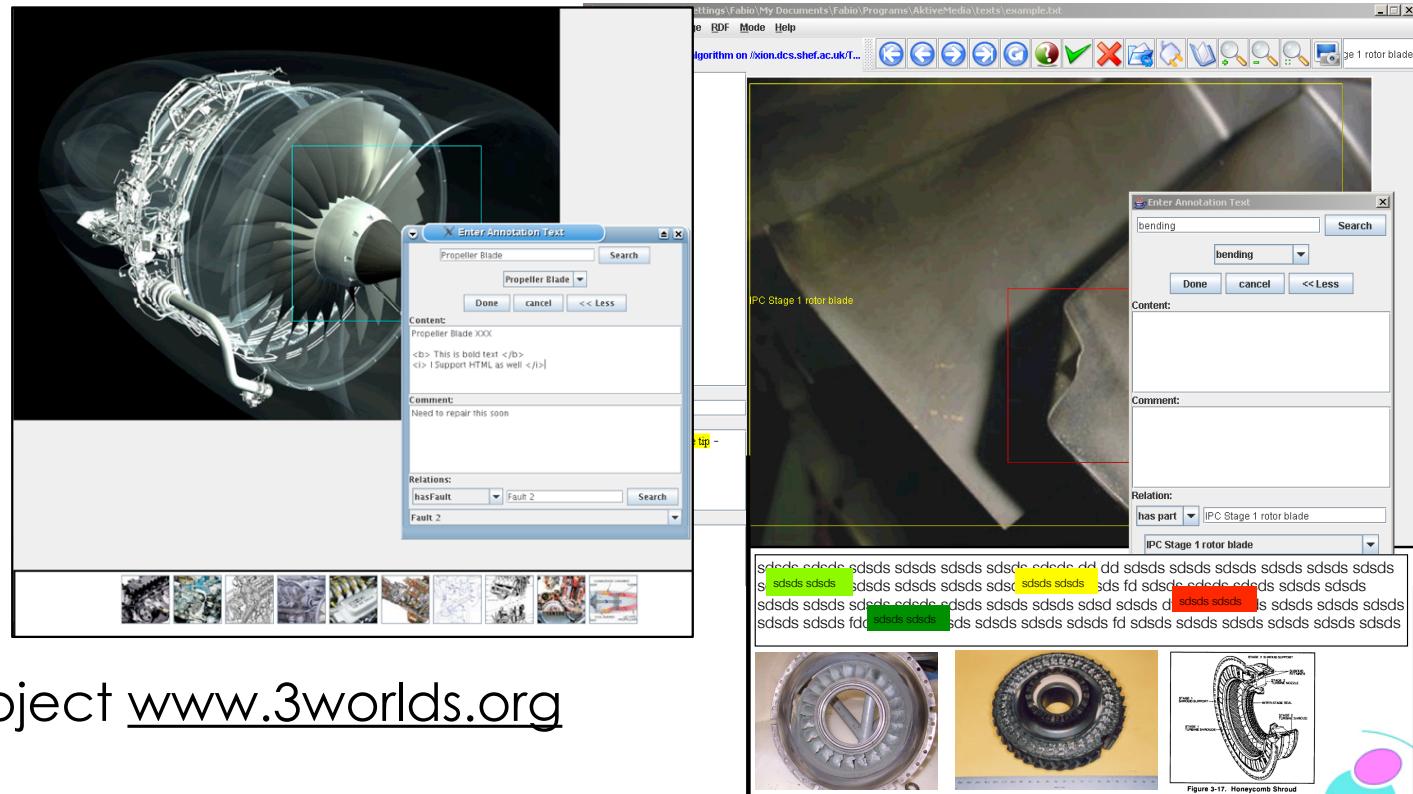
## Annotating across documents (CREAM, 2001)





## Example of Application

- Annotation of compound documents for documenting the overhaul of a jet engine



IPAS project [www.3worlds.org](http://www.3worlds.org)

## Manual Annotation

A real world application  
to knowledge management

**ISWC 2007**

© Fabio Ciravegna, University of Sheffield



# The knowledge capture tension

- Different departments need specific methods and styles to capture knowledge
  - Justified by internal workflows and tradition
- Rest of company needs sharing in ways suitable to reuse
  - The way information is captured is not necessarily the most appropriate





# User-centred Capture in IPAS

- Currently: single departments establish Word or Excel templates to capture knowledge
  - Knowledge is unstructured
    - It requires effort (e.g. Information Extraction) to extract and share knowledge
- K-Forms
  - Easy user-driven creation of Web based forms to capture knowledge
    - Items in forms are manually connected to company's ontology
      - enables sharing and integration with other knowledge sources
    - knowledge is immediately available for search, sharing and reuse
      - based on k-search
    - capture equivalent to that done by IE, but with user in the loop
    - no additional effort required by user wrt present methodology





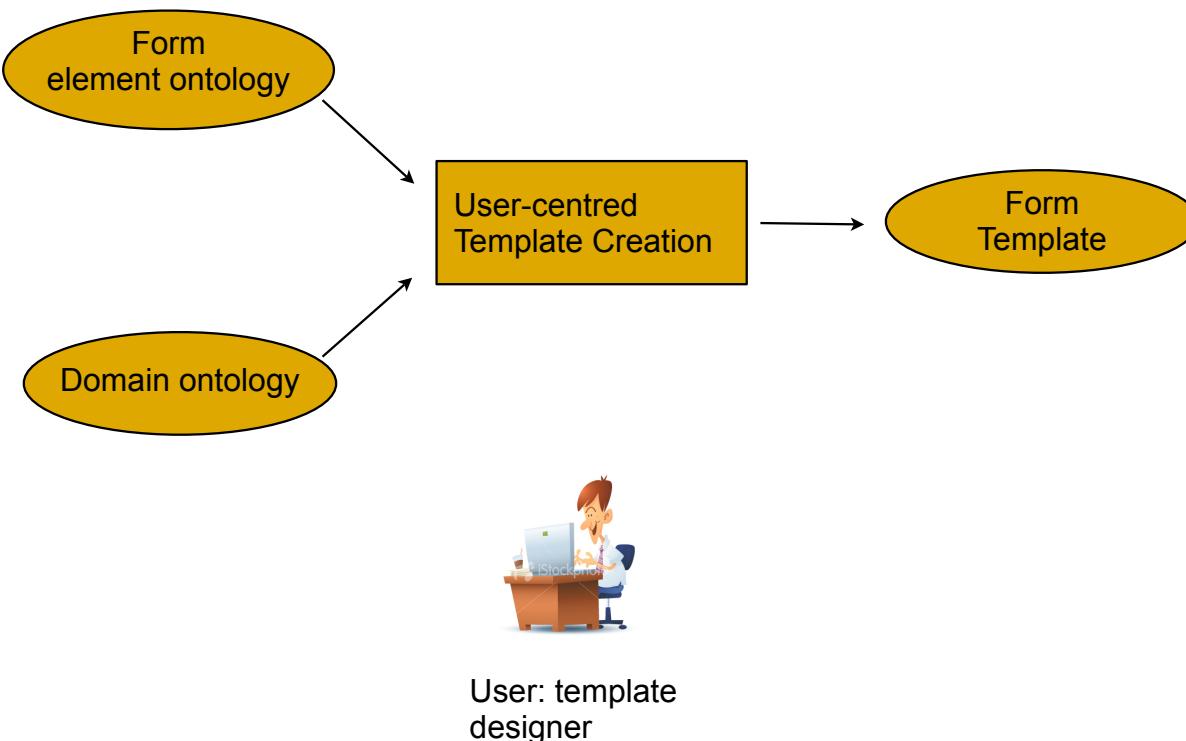
# Forms

- Enables easy definitions of knowledge capturing applications
- Users define Web based forms visually using a Web browser
- Forms, fields, type of information, possible values, etc.
- Connects form fields to existing ontology or database schema
- Applications can be distributed (via intranet), or local to computer (e.g. laptop)
- Upload of data onto central database in second time (if used in local)
- Real world application to Module Condition Reports at Rolls-Royce plc



# Capturing with K-Forms

## 1. Template Design Phase





# Easy creation of Web based forms



[Back to preferences](#)

[Add Form](#)

## Form [0]

[Add Form Field](#)

Name:	Engine_Marks
Field Type:	Textfield
Validation:	Mandatory
Row/Height:	10
Col/Width:	20
Value [Comma seperated for multiple values]:	
Help Text [Optional]:	
Ontology Concept [Optional]:	hasEngine_Marks

[Show/Hide Ontology](#)

### Module Condition Report

#### has Report Detail

##### Report Detail

- has Summary
- has Circulation
- has Issue
- has Additional Keywords
- has Discipline Or Report Series
- has Document Number
- has Authors

[Add Form](#) [Save](#)



## Forms: Features of template creation

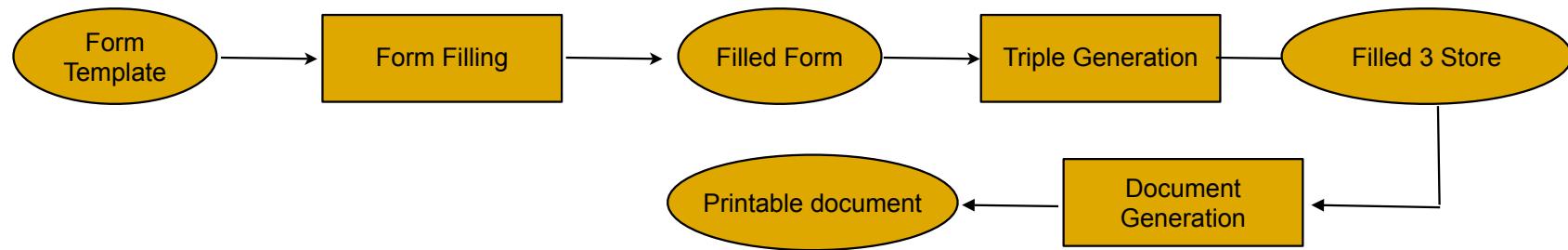
- Easy creation and release of form templates over Intranet and in local
- Automatic generation of final document template for each new template
- Based on:
  - Use of ontology of potential form elements
    - Templates composed using ontology directives compiled using a graphical interface
      - All declarative information
    - Ontology is
      - either created automatically around forms and fields
      - or form fields are mapped to ontology concepts and relations





## Capturing with K-Forms (2)

### 2. Knowledge Capture Phase



User: Domain Expert



## Forms: Features of Knowledge Capturing

- When form is released users receive forms to fill
  - Easy capture in local (no intranet connection)
  - Easy upload to central repository
- Final document automatically generated
  - Can be read and printed and sent by email
- Knowledge immediately searchable with K-Search after uploading to central repository



The  
University  
Of  
Sheffield.



The  
University  
Of  
Sheffield.



[View Data](#) [New Report](#) [Publish](#)

Thank you for using AKTive IPAS

[Logo](#)

This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
- <rdf:RDF>
- <rdf:Description rdf:about="http://xion.dcs.shef.ac.uk:8080/aktiveipas/ontologies/generic/module#1.1194256780635">
  <j:Repair_0_4>1</j:Repair_0_4>
  <j:text_MD_Mechanism_0_1_0>Cracked</j:text_MD_Mechanism_0_1_0>
  <j:Benchmark_2_0/>
  <j:TV_And_Accepted_0_6/>
  <j:text_Mechanism_1_1_1>Other</j:text_Mechanism_1_1_1>
  <j:Measurement_data_0_4_2>NONE NONE NONE </j:Measurement_data_0_4_2>
  <j:Measurement_data_2_4_1>Cracks beyond SP50 limit</j:Measurement_data_2_4_1>
  <j:Inspection_Locations_0_0>Assembly</j:Inspection_Locations_0_0>
  <j:IPC_Reference>725221-01-060</j:IPC_Reference>
  <j:Mechanism_0_1_2>on</j:Mechanism_0_1_2>
  <j:Invoice_Value_Paid_Repair>NONE</j:Invoice_Value_Paid_Repair>
  <j:Operator>SIA</j:Operator>
  <j:SAP_Identification_Number>TR0010708-725221-01-060</j:SAP_Identification_Number>
  <j:Condition_Details_Pre_Clean_0_3_1/>
  <j:text_MD_Mechanism_2_1_1>Other</j:text_MD_Mechanism_2_1_1>
  <j:text_MD_Mechanism_1_1_0>Corrosion</j:text_MD_Mechanism_1_1_0>
  <j:Measurement_data_1_4_2/>
  <j:Mechanism_0_2>Loose shank nuts</j:Mechanism_0_2>
  <j:Mechanism_3_2/>
  <j:Condition_Pre_Clean_0_2_2>Light frettage to NGV location rails</j:Condition_Pre_Clean_0_2_2>
  <j:Accept_As_Is_1_5/>
```

87



# Annotations: Where From?

- SW relies on document annotation
  - Current state of art requires manual annotation
- Manual Annotation for the global Web (Internet)
  - Very few people will annotate web pages by hand
  - What if they did?
    - Isn't the web based on hype?
      - Do people really need to publish their girlfriend photos?



- Expensive/time consuming/difficult
  - Chicken-egg problem
  - If it adds time to page editing, users will not do it unless there is really something for them
    - Usefulness and hype
- Inefficient and never ending
  - Every new document needs to be annotated
- Difficult
  - if two people annotate the same documents have 15-30/100 probabilities to annotate them differently
    - Risk is that the same information is annotated differently
      - Disagreement between annotators means data sparsity
      - Information becomes difficult to retrieve



## An Example

- 10 annotators
- Emails about workshop announcements
  - Name, acronym, date of workshop
  - Name, acronym, URL of associated conference (if any)
  - Submission dates.
- 15% inter-annotator disagreement
  - Especially on name of conference/workshop

C:\Fabio\Projects\Pascal\Challenge Data\Example of Corpus\train>IDAMAP\_1999.txt

File Settings Commands Help

concept

- Workshop
  - WorkshopName
  - WorkshopAcronym
  - WorkshopDate
  - WorkshopHomepage
  - WorkshopLocation
  - WorkshopPaperSubmissionDate
  - WorkshopNotificationOfAcceptanceDate
  - WorkshopCameraReadyCopyDate
- Conference
  - ConferenceName
  - ConferenceAcronym
  - ConferenceHomepage
- relation

Is this the name or the acronym?

ontology

```
*** Workshop: ***
*** Intelligent Data Analysis in Medicine and Pharmacology ***
*** (IDAMAP 99) ***
*****
Saturday, November 6, 1999
Washington, DC, USA
during the
AMIA 1999 Annual Symposium
November 6-10, 1999 in Washington, DC, USA
(homepage of IDAMAP 99
http://www.ifs.tuwien.ac.at/~silvia/idemap99/
(homepage of AMIA 1999
http://www.amia.org/meetings/f99/call/cover.htm
)

-----
-----
```

Important dates

- \* Submission deadline: July 26, 1999
- \* Notification to authors: September 6, 1999
- \* Camera-ready paper: October 11, 1999
- \* Conference: November 6-10, 1999
- \* Workshop: Saturday, November 6, 1999

GENERAL INFORMATION:

-----

IDAMAP-99 is a Workshop at the AMIA 1999 Annual Symposium - November 6-10, 1999 - Washington, DC prior to the start of the main AMIA conference.

Gathering in an informal setting, workshop participants will have the opportunity to meet and discuss selected technical topics in an atmosphere, which fosters the active exchange of ideas among researchers and practitioners. To encourage interaction and a broad exchange of ideas, the workshop will be kept small, preferably under 30 active participants, although registered AMIA 99 Fall Symposium members are welcome to attend. The workshop is intended to be a genuinely interactive event and not a mini-conference, thus ample time will be allotted for general discussion. The workshop will last a half-day.

This is the fourth workshop on Intelligent Data Analysis in Medicine and Pharmacology (IDAMAP). The former IDAMAP Workshops were held in Budapest in 1996, in Nagoya in 1997, and in Brighton in 1998.

Why not including Annual/Fall symposium?

Missing workshop location!

Connecting to learning algorithm on /I27.0.0.1/AmilcareVr2



## Problems in the example

- The previous example contains
  - Three doubtful cases (conference name/acronym)
  - One mistake
- It was annotated by two people and a third one checked their annotations
  - It was in a corpus of hundreds of documents!



- Tedium & Tiring
  - Error prone
- Legacy with the past
  - Ontologies are living objects, new version produced
    - Which version of the ontology is used for annotation?
- Dispersed information
  - Annotation largely unfeasible for large diverse repositories
    - E.g. a Web site
      - Department of CS of the University of Southampton: 1,600 pages
      - How many relevant ontologies are there for that department?



# Problems with Manual Annotation (3)

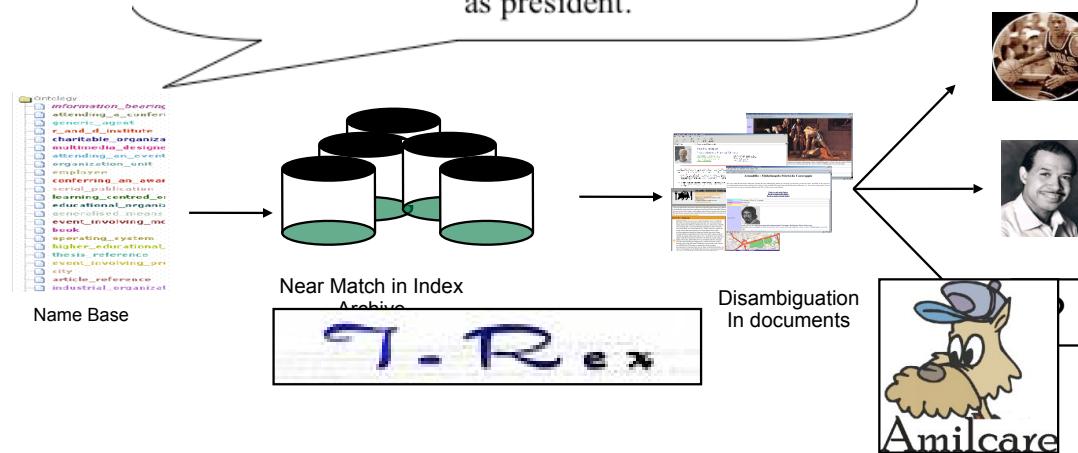
- How many annotation schemas?
  - The Semantic Web is expected to be composed of
    - [Many] small ontological components [Hendler 2001] will be created, mainly related to different domain and applications
    - University of Sheffield web site:
      - What ontology for annotation?
        - Universities/Education, Research life, Scientific Papers,
        - Sport, computer network organization....
        - You name what...



## Annotation for use...

- If annotation is to be chosen by author/owner
  - Selection of Annotation Schema may reflect world model of the creator, not of the user
    - E.g. education is the main goal of the university, so the central Uni will probably choose an ontology on Education
    - Most of my time is actually devoted to research
    - Most of my colleagues look for scientific information on our web site
    - To us, Uni's annotation would be largely unuseful
    - Question:
      - Who (and how!) is going to introduce the annotation for us?
      - Where is the annotation to be inserted?

WASHINGTON, D.C. (October 5, 1999) -  
 nQuest Inc. today announced that Paul Jacobs, former Vice-President of E-Commerce at SRA International, has joined the company's executive management team as president.



## Automating Annotation



# Annotation Engines

- Manual document annotation is still largely expected to be the main SW vehicle creation
  - Especially for trusted environment (e.g. within a company) this is expected to provide high quality material
- But manual annotation is largely unfeasible over large scale
  - Unavailability of users
  - Risks related to incompetence or spamming
  - Large amount of legacy data
- Automatic annotation
  - To help manual annotation OR
  - To replace human annotators (e.g. on legacy data)



# Tasks for KA: Extraction

- Text:

- Entity Extraction
- Table Fields Extraction
- Relation Extraction
- Event Extraction

- Data:

- Similarity of Data Instances
- Functions and relation
- Finding patterns and (ir-)regularities in data

- Images:

- Semantically driven Image analysis using ontologies, for retrieval and annotation
- Image classification/ clustering with respect to the dominant visual trends



# Information Extraction from Text

- Automatically extracting pre-specified information from textual documents
  - salient facts about pre-specified types of events, entities or relationships.
- Populating a structured information source from a semi-structured, unstructured, or free text, information source.

WASHINGTON, D.C. (October 5, 1999) - nQuest Inc. today announced that Paul Jacobs, former Vice-President of E-Commerce at SRA International, has joined the company's executive management as president.

**Company:** nQuest Inc.

**Date:** today

**InPerson:** Paul Jacobs

**InRole:** president

**Company:** SRA International

**OutPerson:** Paul Jacobs

**OutRole:** Vice-President of E-Commerce,

Named Entities

Event Recognition

Growing complexity



# Named Entity Recognition

- Tasks:

- Recognition and classification of named entities
  - E.g. people's names, companies, locations, etc.
- Unique identification of named entities (URI assignment)
  - Including disambiguation
    - Michael Jordan as basketball player Vs lawyer
    - London UK Vs London USA
- Integration with other sources
  - E.g. positioning on a map



- Two steps:
  - Training phase
    - Input: annotated set of representative documents
    - Output: trained system
  - At runtime
    - One-by-one document analysis
- Expected accuracy:
  - 80-95% (free texts)
  - Web documents tend to require additional processing to get equivalent results (but doable to some extent)
- Medium Scale: up to hundreds of thousands of documents

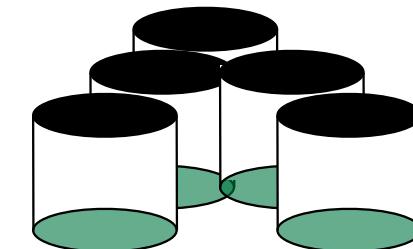
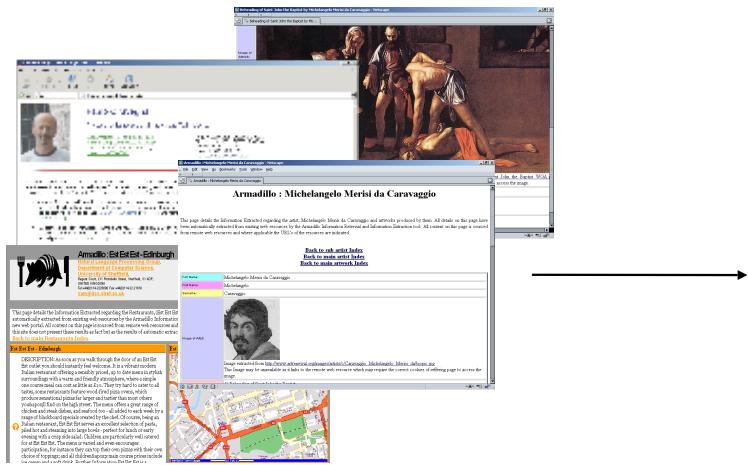


- For large scale (some hundred millions pages) smarter infrastructure is needed
  - Search engine-like indexing infrastructure
  - Faster processing (less processing)
  - Two cases:
    - Recognition of known terms (and their variations)
      - See also information integration
    - Discovery of new names



# Large Scale NER: Indexing

- Document Indexing as in Search Engines



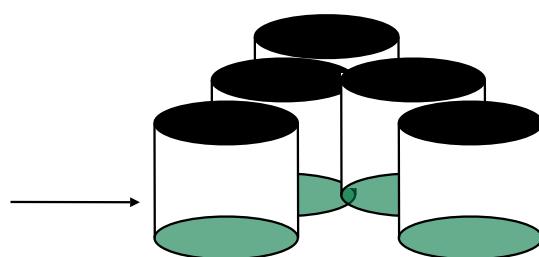
Distributed Index Archive  
(keywords)



# Known Name Recognition

```
ontology.owl
└── information_beering
    ├── attending_a_confer
    ├── generic_agent
    ├── r_and_d_institute
    ├── charitable_organiza
    ├── multimedia_designe
    ├── attending_an_event
    ├── organization_unit
    ├── employee
    ├── conferring_an_awar
    ├── serial_publicatio
    ├── learning_centred_oi
    ├── educational_organiz
    ├── generalised_mean
    ├── event_involving_mc
    ├── book
    ├── operating_system
    ├── higher_educational_
    ├── thesis_referenc
    ├── event_involving_pr
    ├── city
    ├── article_referenc
    └── industrial_organiz
```

Name Base



Near Match in Index Archive



Disambiguation  
In documents



S. Dill, N. Eiron, et al: SemTag and Seeker: Bootstrapping the semantic web via automated semantic annotation. WWW'03



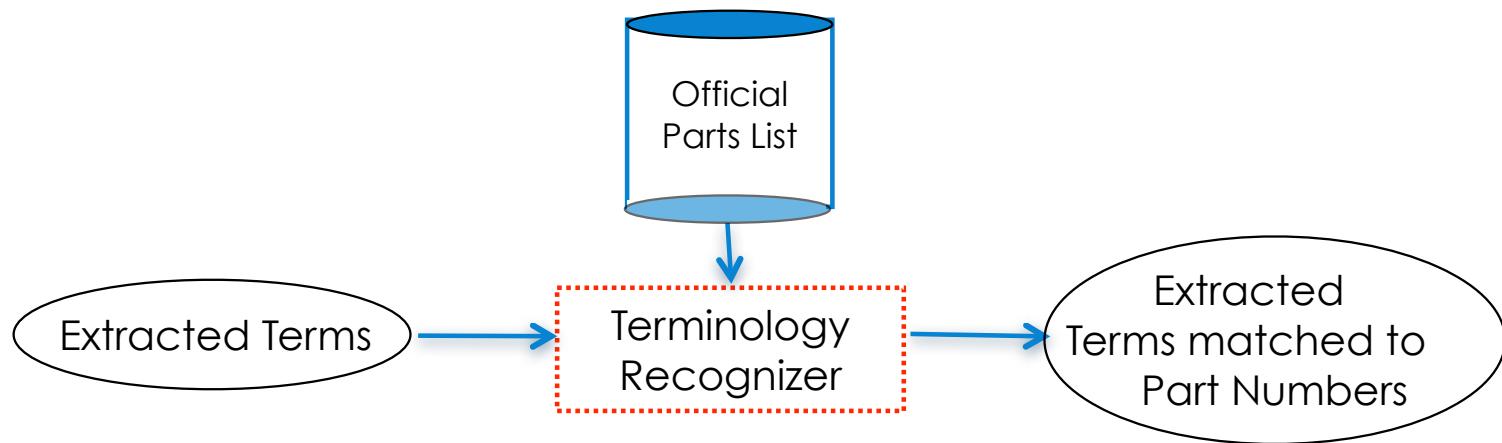
# Discovery of New Names

- Modified Indexing of documents to recognize potential names
  - Traditional NER
    - On the window of words (not the whole doc!!!)
      - Fast and effective
  - Web specific strategies
    - To identify names without context



# Terminology Recognition

- NER is one example of term recognition
- More useful in technical domains is terminology recognition
  - The task of assigning a URI to a technical description
    - i.e. mapping a natural language description to the official company ontology



ISWC 2007



# Terminology Recognition

- Possible approaches

- Linguistic approaches

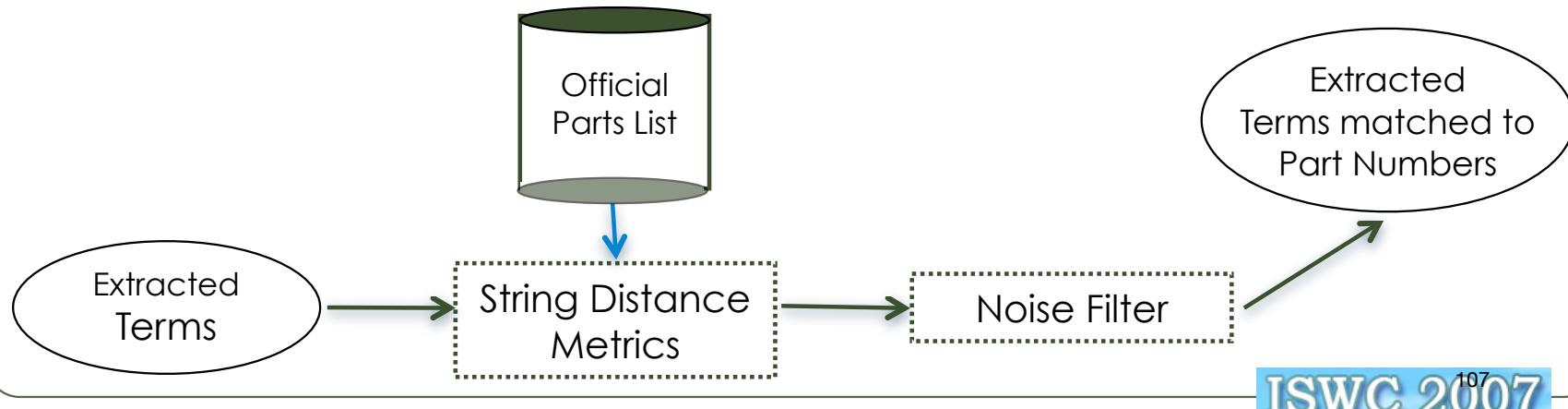
- Based on linguistic analysis of terms (Gaizauskas *et al* 2003)

- Statistical approaches

- Based on frequency analysis and detection

- Other approaches

- Distance metrics based (Butters 2007)



ISWC 2007<sup>107</sup>



## More complex IE: event modelling

- Not just NER but also relation among elements in a document
  - More complex task
  - Requires some reasoning to bridge the complexity of events to the ontology structure
    - Imprecision in extraction
    - Information non matching the ontology schema
- This is where IE has hit a performance ceiling
  - 60/70 Precision/Recall ratio since 1998

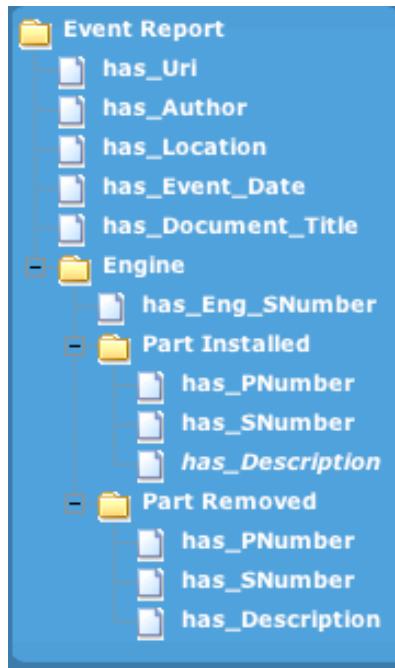


# Table Field Extraction

- Tables are an essential part of many documents
  - Most information is represented in tables
- Tables can be represented as forms to fill
  - Semantics is fixed
  - Wrapper writing or wrapper induction (Kushmerick 1997)
- Tables can be created ad hoc in documents (e.g. Word docs)
  - Semantics is unclear
  - Sometimes documents are created as part of a workflow, therefore they tend to be created using common models
    - e.g. by re-using the previously generated document
    - hence tables evolve, but still semantics can be traced



# An Example of Automatic IE



- Automatic extraction of information from event report
  - 18,000 documents analysed
- Metadata generated according to a simple ontology
- Automatic extraction of metadata and indexing of documents



# An Experiment on Event Report for Jet Engines

1263 Prepared By: Richard Williamson Originated/Revised on: 13 March 2004 Event Report No. COMPANY2698 /

Event Report Data		engine type	company	Aircraft
LN144				
Event Date:	12-Mar-04	Engine S/N:	51179	Flight Regime: Unknown Hazard Type: No Hazard
Aircraft Regn.:	GB-BKA	Installed Power:	Legal	Location: SAV Event Type: Operational
Airframe Hours:	20779	Engine TSN/CSN:	14242 / 4014	Event Category: Basic
Airframe Cycles:	5609	Engine TSF/CSF:	6249 / 1814	
<b>Reactions to Event:</b>				
Primary:	None	ABTO Speed (Knots):	N/A	Operational Effect: No Effect SERAPH Symptom Codes:
Secondary:	None			Delay Time (mins): N/A NREP NREP
Third:	None			Fuel Dumped?: No NREP
Fourth:	None			
<b>EICAS Messages (If Any):</b>				
Parts/Components Removed or Installed (If Any):				
On/Off	Part Number / Serial Number	Part Description	Hours / Cycles	Destiny / Disposition Pull Category / Pull Code
Installed	9-217-62	FUEL FLOW TRANSMITTER	1	SE - Serviceable
Removed	921762 Y403	FUEL FLOW TRANSMITTER	1	R4 - Return to Manufacturer U - Unplanned US - Unserviceable I - Inspection/Investigation
<b>Description of Event:</b>				
a short sentence				

91



# Types of tables in Event Reports

module/accessory details			
<u>item</u>	<u>part number</u>	<u>s/n removed</u>	<u>s/n installed</u>
	p39-401revf	04-0721257 tsn/csn: 268/106	04-1012229 tsn/csn:0/0

Part numbers
04-0721257 tsn/csn: 268/106 off
04-1012229 tsn/csn:0/0 on

<u>s/n removed</u>	04-0721257 tsn/csn: 268/106
<u>s/n installed</u>	04-1012229 tsn/csn:0/0

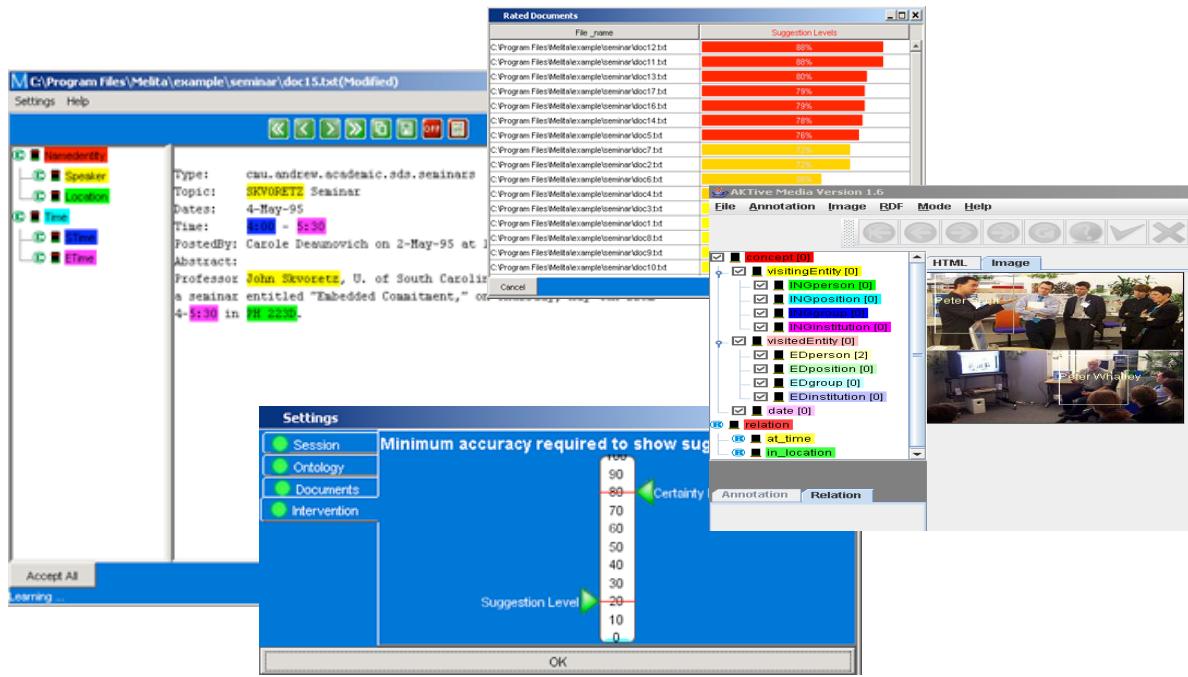
Parts/Components Removed or Installed (If Any):						
On/Off	Part Number / Serial Number	Part Description	Hours / Cycles	Qty	Destiny	Deposit
Installed	FK30840	(TO SB72-C629)	11129 TSN 1954	1		
	RGG12340					
Installed	FK21221		11652 TSN 2119	1		
	EC092					
Installed	FK30840		11129 TSN 1954	1		
	RGG12501					
Installed	FK30840		11129 TSN 1954	1		
	RGG12208					
Installed	FK30840		11129 TSN 1954	1		
	RGG12391					



# Applying information extraction

- AktiveMedia to annotate texts
- TRex system (Jiria et al. 2006) to train and extract
  - <http://tynel.shef.ac.uk/t-rex/>
- IE captures most of the information in tables
  - 99% of the information captured (recall=99)
  - 98% of proposed information is correct (precision=98)

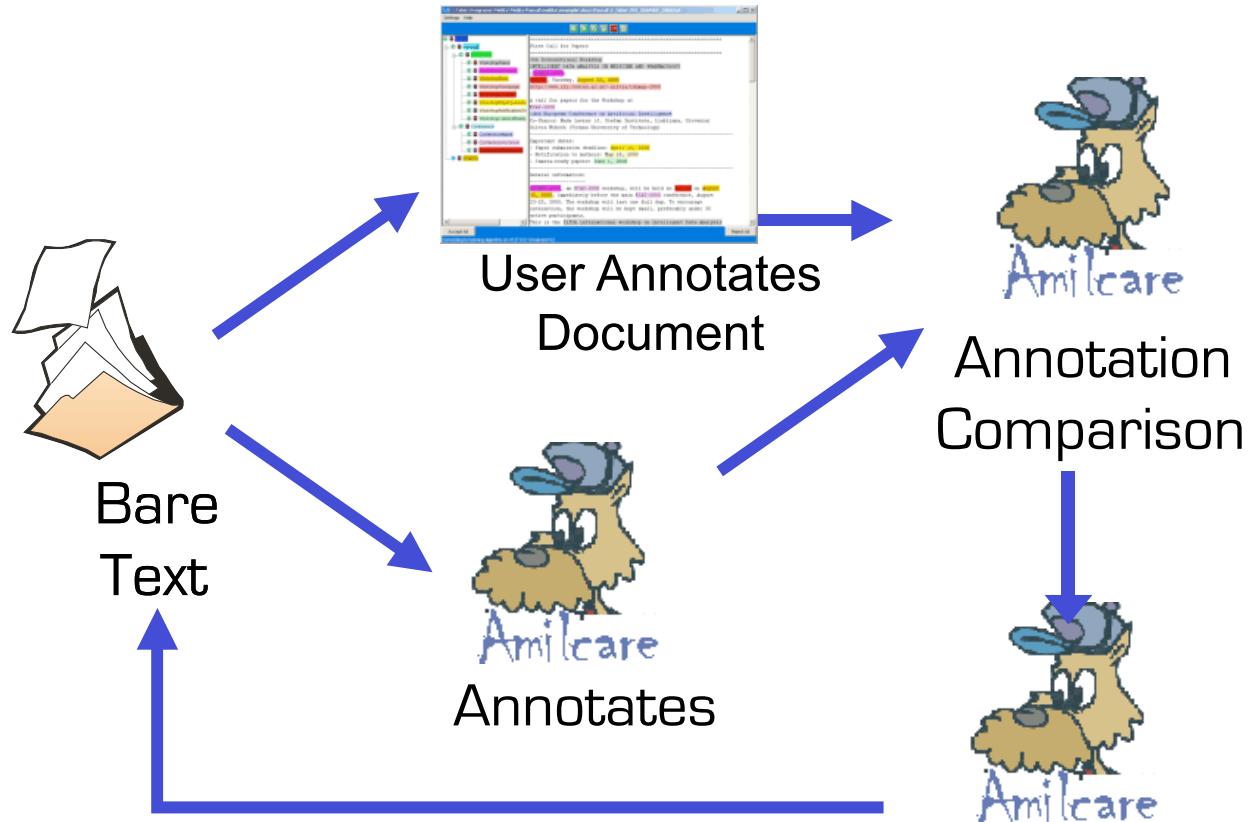
	<b>POS</b>	<b>ACT</b>	<b>CORR</b>	<b>WRONG</b>	<b>MISSSED</b>	<b>PREC</b>	<b>REC</b>	<b>F1</b>
airport	120	120	120	0	0	100	100	100
has_airframe_cycles	104	104	104	0	0	100	100	100
has_airframe_hours	104	104	104	0	0	100	100	100
has_author	120	120	120	0	0	100	100	100
has_engine_serial_number	120	120	120	0	0	100	100	100
has_engine_type	120	120	120	0	0	100	100	100
has_event_date	120	120	120	0	0	100	100	100
has_event_report_no	356	358	356	2	0	99	100	100
has_part_description_installed	120	113	111	2	9	98	93	95
has_part_description_removed	120	133	120	13	0	90	100	95
has_part_number_installed	120	113	111	2	9	98	93	95
has_part_number_removed	120	133	119	14	1	89	99	94
<b>TOTAL</b>	<b>1644</b>	<b>1658</b>	<b>1625</b>	<b>33</b>	<b>19</b>	<b>98</b>	<b>99</b>	<b>98</b>



## Using IE to Support Manual Annotation

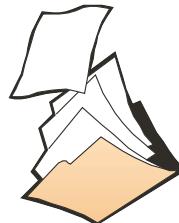


# Using IE to support annotation: step 1





## Using IE to support annotation: step 2



Bare  
Text



Annotates



User  
Corrects

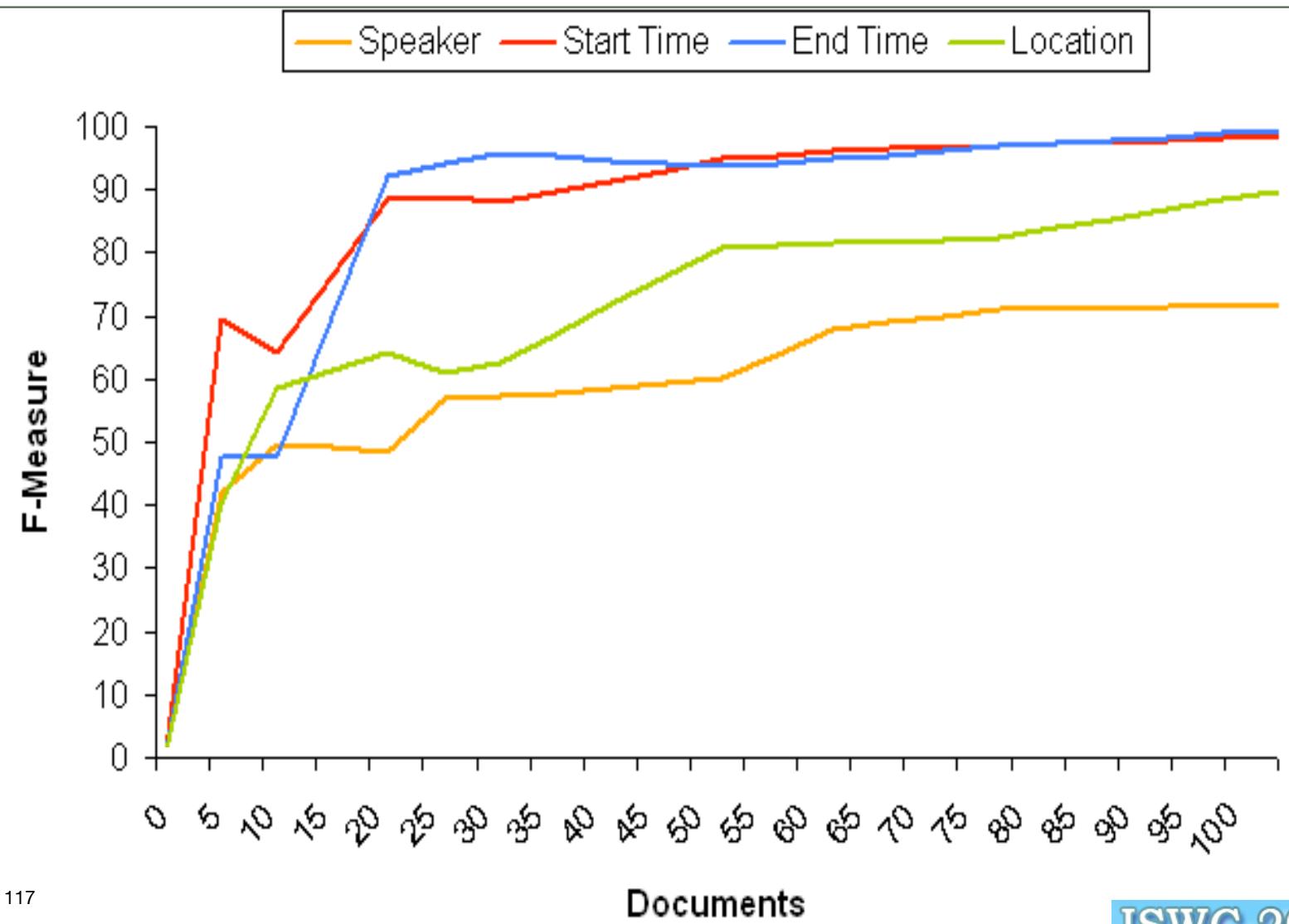


Uses  
corrections to  
retrain





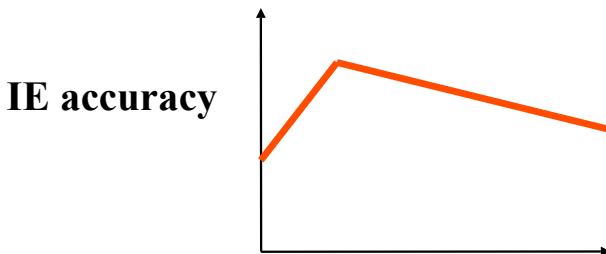
# Learning curve



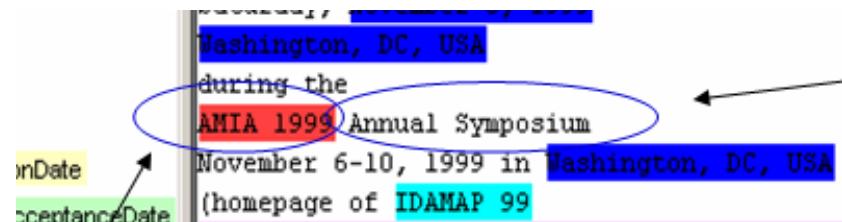


# Impact on Annotation

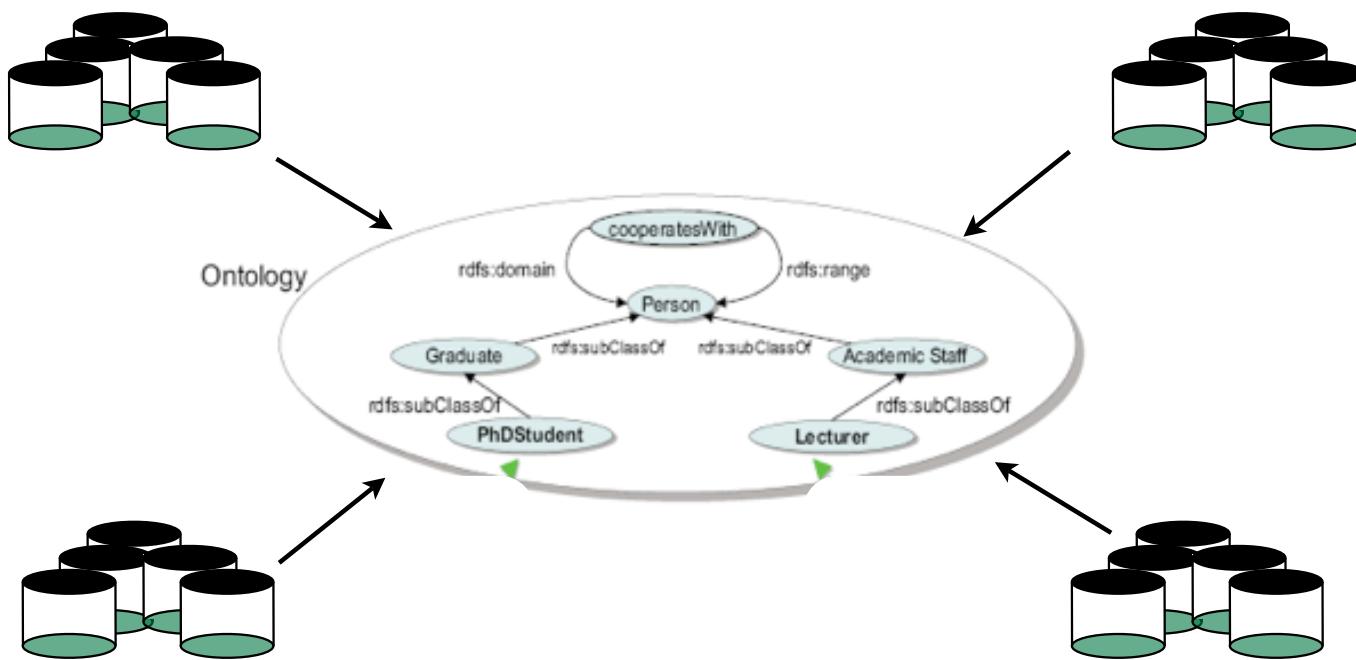
- University of Karlsruhe experiments
  - -80% annotation time
  - +100 interannotator agreement
    - Is this positive?
- Outstanding issue:
  - Impact on annotators of suggestions topping 85% accuracy?
  - Annotation needs to be precise and consistent
    - Otherwise the IE system is confused
  - Can only annotate document content
    - With connections to the rest of the knowledge via information integration



Amount of annotations



ISWC 2007



## Information Integration

ISWC 2007

© Fabio Ciravegna, University of Sheffield

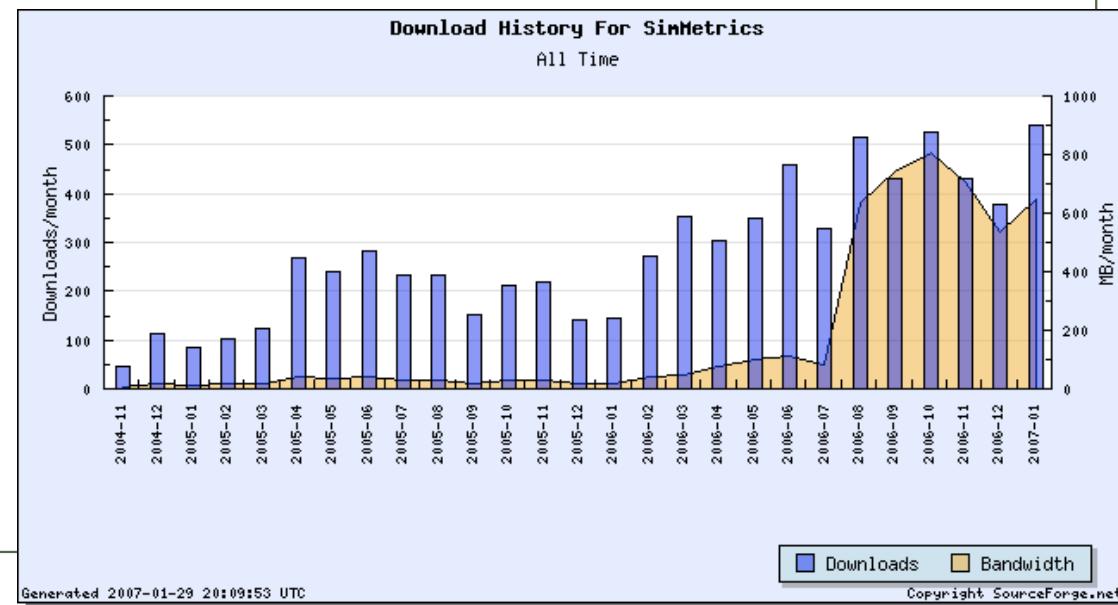


- Facts from different sources need to be integrated
  - To connect information/knowledge across docs
    - Assign unique URI
  - To solve discrepancies and ambiguities
- Steps
  - Unique instance identification (for entities)
  - Record linkage (for events)
- Information Integration strategies
  - Generic
    - Distance metrics  
(Chapman 2004)
    - Using Web bias
  - Statistical matching
  - Application specific
  - Rules



# SimMetrics

- Library of distance metrics released as open source
  - <http://sourceforge.net/projects/simmetrics/>
  - 7,489 downloads in 2 years
  - Most downloaded distance metrics library on the Web
  - Hundreds of applications





The  
University  
Of  
Sheffield.

# Armadillo: Historical Data Mining



Arts & Humanities  
Research Council

**The Marine Society  
Registers**

**The Westminster  
Historical Database**

**Eighteenth Century  
Fire Insurance  
Policies**

**Prerogative Court of  
Canterbury Wills**

**The Proceedings of  
the Old Bailey**

**AHDS Deposits**

**St. Martin's  
Settlement Exams  
Index**  
**WESTCAT**

**Collage image  
database**  
**Guildhall Library**

**Harben's Dictionary  
of London**

**John Strype's  
“Survey...”**

<http://www.motco.com>

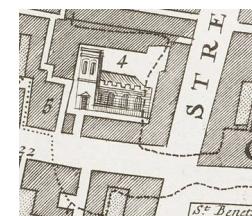
**Metropolitan  
London in the 1690s**  
**IHR**

**Selected Criminal  
Records**  
**PRO**

**House of Lords  
Journals**  
**BOPCRIS**



[+]  
THE  
PROCEEDIN  
ON THE  
KING's Commission of the I  
A N D  
Over and Terminer, and Goal-Delivery of Newgate, held for  
Ledes and COUNTY of Middlesex, at Tyger Hall in the  
On Wednesday, Thursday, and Friday, being the 16th, 17th, and 18th  
January, in the Ninth Year of His MAJESTY's R  
BEFORE the Right Honourable Sirps of Chancery, at  
the Great Contra, Knigge, Mr. Vice-Chancery, Mr. S  
Mr. Justice Peale, Mr. Justice Downe, Mr. Justice Peale,  
Downe, Mr. Justice Peale, Mr. Justice Peale, Mr. Justice Peale,

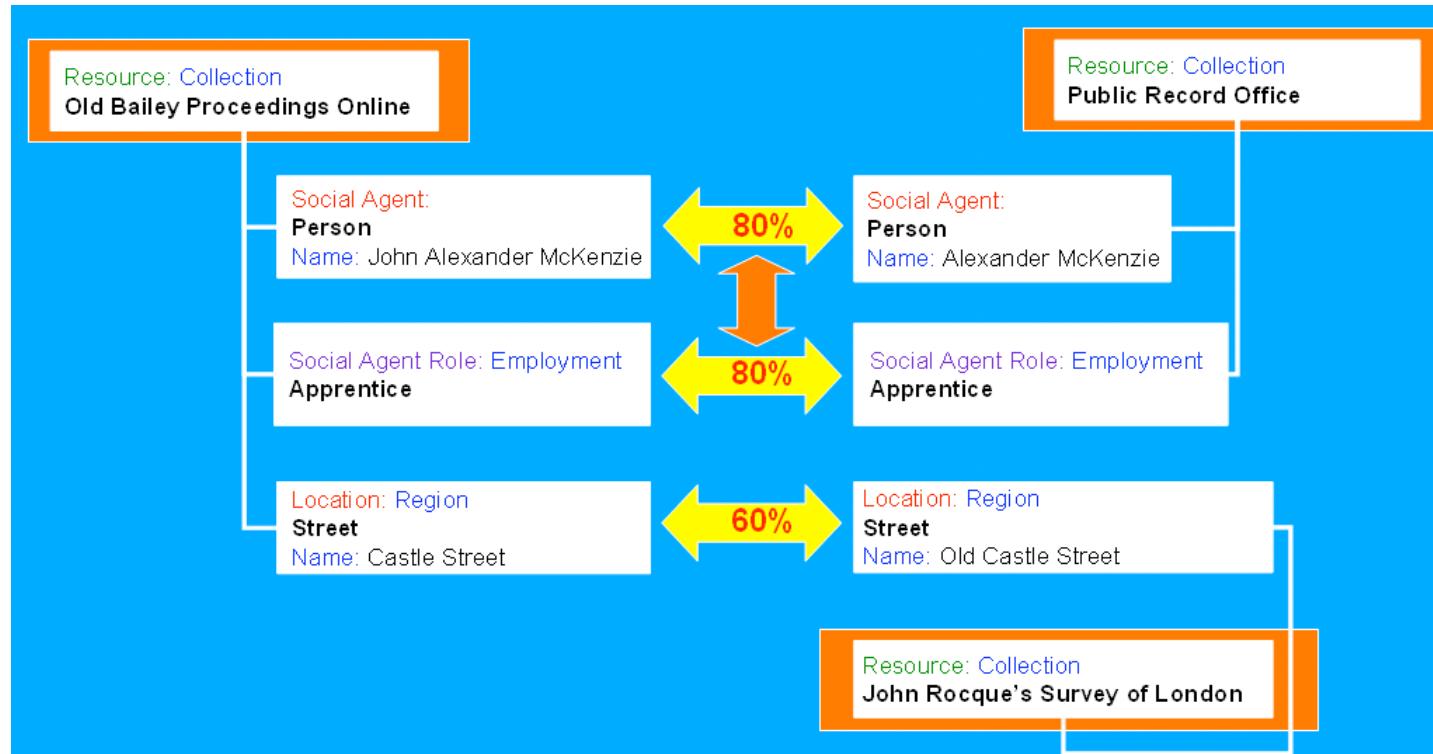


<http://www.hrionline.ac.uk/armadillo/>

**ISWC 2007**



## Armadillo: Historical Data Mining





- Large scale?
  - Ontologies:
    - large ontologies (up to 10k) with simple tasks (SemTag and Seeker, Kim)
    - small/medium scale (up to 100) with more complex tasks
  - KB: large scale
- Portability: most technology difficult to port without experts (Armadillo, KIM)
  - User input well exploited in human-centred acquisition (e.g. Melita, AktiveMedia)
- Cross-Media: exploited in user centred annotation (e.g. AktiveMedia)
- Background Knowledge
  - Used in AktiveMedia, KIM, SemTag and Armadillo to some extent
  - Uncertainty: some use in Armadillo

The screenshot displays the X-ME system interface, which integrates various knowledge management components:

- Top Left:** A map of the Gulf of Mexico and surrounding regions, including states like Mississippi, Alabama, and Florida.
- Top Right:** A diagram illustrating cross-disciplinary issues, cross-jurisdictional issues, shared systems issues, and general comments, connected by arrows.
- Middle Left:** A sidebar with navigation links: New Process, Search, E-mail, Process View, Problem View, Summary View, Annotate, and Add new class.
- Middle Center:** A main workspace divided into several panels:
  - Carbon Formation:** Includes "Image - Engine" and "Engine" sections.
  - Delamination:** Includes "Trailing edge burning", "Surface sintering", and "Report 3207".
  - Loss of coating:** Includes "Cracking on engine surface" and "Erosion -".
  - Other Assertions:** Includes "Burning", "Cracking", and "Loss of coating".
  - Issue no 74:** A detailed view of an issue, showing a summary table with coordinates (-0.3804, 0), (-0.3804, 1), (-0.2743, 0), and (-0.2743, 1), and a text area with "IF development vanes shows burn back".
  - Problem Def:** A search interface for "Delamination" with results "Search 3" and "Search 4".
- Bottom Right:** A screenshot of a desktop environment showing a "Composition NASA MMS Cross2 FINAL" window containing a map titled "Map: Brent's Analysis" and several image files related to Rock Hill analysis.

## Knowledge Sharing and Reuse

- issues in knowledge sharing
- approaches and novel methods to searching, sharing and reuse knowledge

ISWC 2007

© Fabio Ciravegna, University of Sheffield



- In KM mainly means
  - Retrieving information and knowledge
    - At the right time
    - In the right form
      - E.g. independently from where it is stored
      - Or even the form in which it is stored
    - Suitable to the specific users
      - e.g. patients should not receive information using technical terms
    - Suitable to specific interests
      - I am working on social aspects of SW, not interested in engineering aspect of SW
  - In an efficient and effective way
    - Coping with large scale
  - Supporting processes



## Requirements for KS&R: Integration

- Large distributed archives require ability to
  - Map distribution of information,
  - Weight every single source
  - Distribute searches carefully;
- Currently: search is performed just in some of the archives, disregarding others that can bring very useful information
- Importance of context and provenance in searching
  - *Sylvester attacked Tweety and Tweety flew away* (Tweety is a tweety hence tweeties fly)
  - *Tweety came back from hospital with a broken wing. Sylvester attacked her. And Tweety could not fly away* (tweety is a tweety hence tweeties do NOT fly ???)



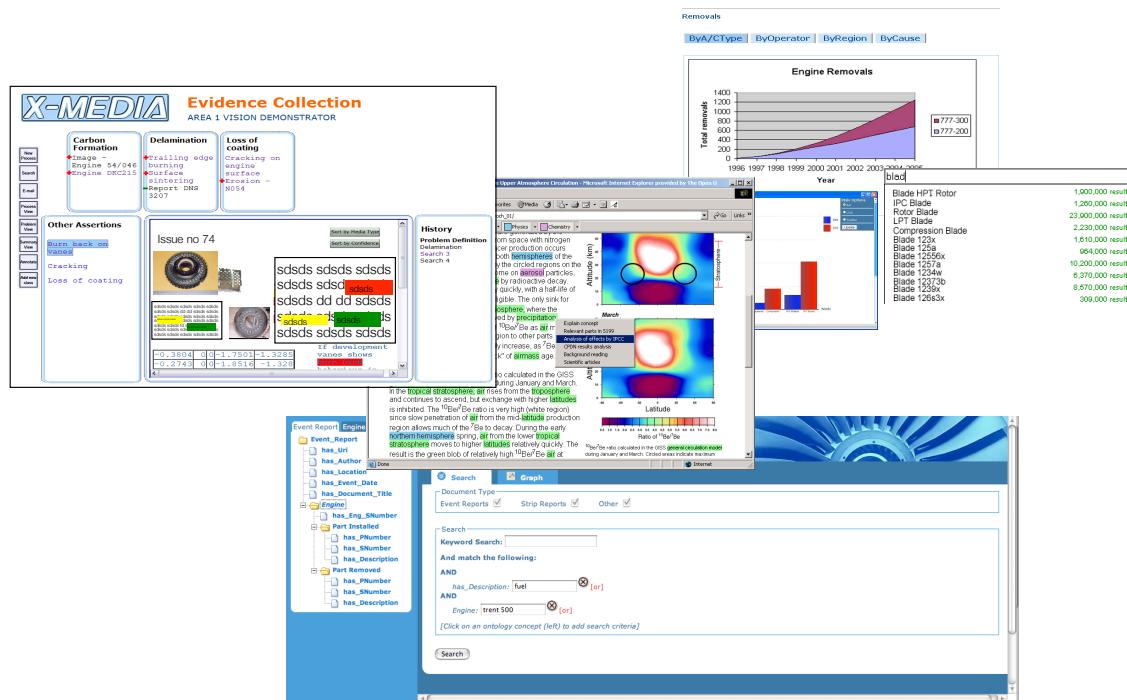
# Requirements for KS&R: Focussing

- Managing knowledge becomes more complex and needs powerful focusing methodologies.
- Focus of searching
  - Changes in time and from user to user,
  - Requires a balanced mixture of exploration and searching;
- Focus on what I know: My knowledge as a basis of what to look for
  - My context rather than everyone's context
    - Tell me what I do not know
      - What is that other people know that I do not
      - Tell me more about what I know



## Requirements of KS&R: Focussing (2)

- What is the user interested in:
  - Most frequent phenomena
    - Redundancy-based approaches to KA can work
  - In less frequent phenomena
    - Redundancy-based approaches do not work
  - Domain specific metrics
    - e.g. disruption caused to customers
    - A mix of the two above



## SW for Knowledge Sharing and Reuse

ISWC 2007



- Ontology based annotation enables
  - Searching using ontologies
    - Searching metadata rather than text
  - Connection of information across documents, media and archives
    - Retrieving information independently from the store/media
  - Reasoning on knowledge
    - Making implicit explicit
  - Workflow support
    - Supporting user actions rather than single searches

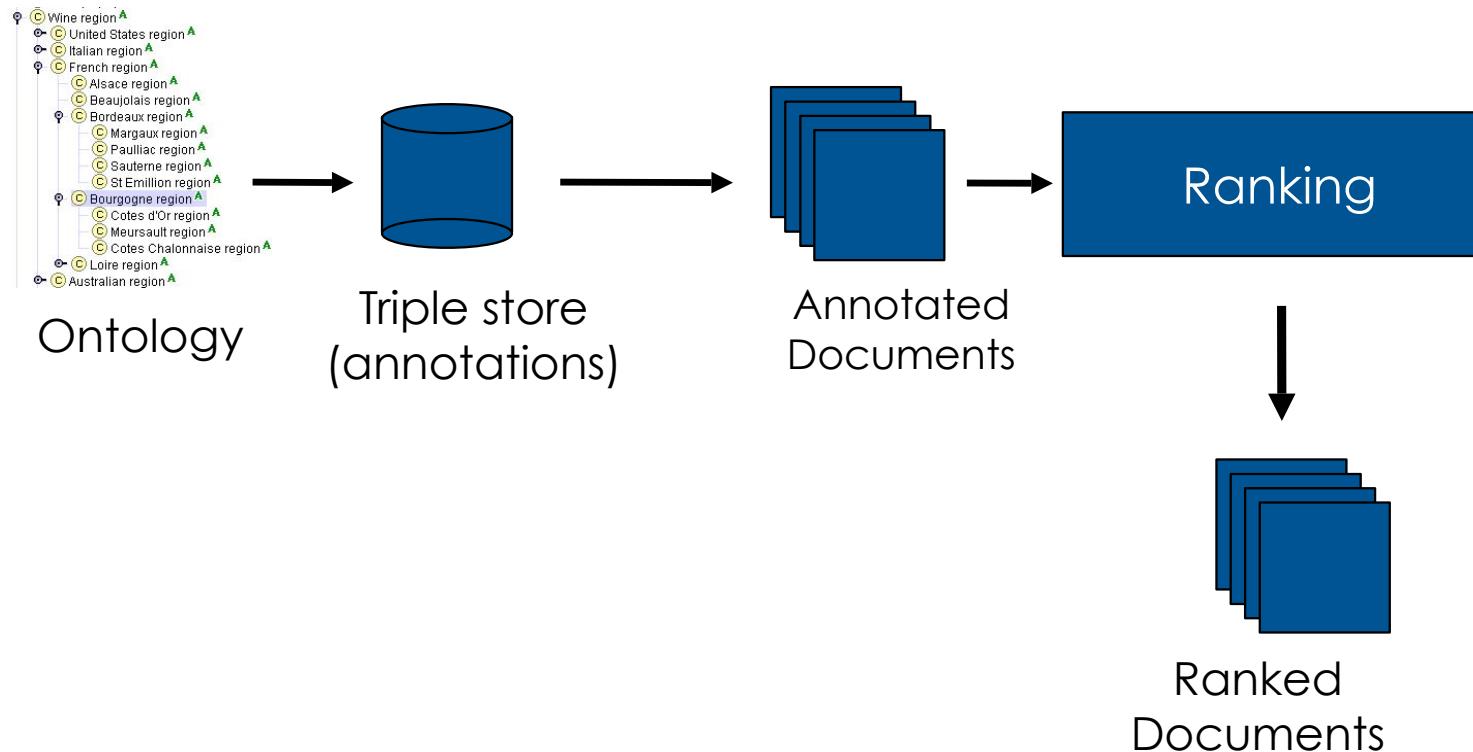


# Ontology-based Search

- Searching metadata rather than texts or images
  - Ontology enables reasoning
    - More flexible than searching using traditional methods
- Searching to...
  - Retrieve documents (images/texts/videos/data)
    - As replacement of traditional document management systems
  - Retrieve information/knowledge
    - Querying the knowledge (e.g. the triple store)



# Searching Documents using ontologies





# Ontology-based Querying: Issues

- Building an ontology is expensive and needs maintenance
- Metadata generation of documents is
  - Expensive
  - Error prone
- Metadata can cover just part of the material of interest to the users
  - The information not annotated using metadata is irretrievable
    - Often the use people will do of information is impossible to foresee
    - Sometimes Information is impossible to retrieve reliably using automatic methods
- If automatic means are used, often some parts of the knowledge is beyond the current technical capabilities





## An Experiment on Jet Engine Event Reports

- 21 topics of search, e.g.
  - "How many events were caused during maintenance in 2003?"
  - "What events were caused during maintenance in 2003 due to control units?"
  - 'Find all the events associated with damage to acoustic liners following bird strike"
- How many topics can we model with Information Extraction?
  - 21 topics/ 14 topics partially or not covered by IE-based annotations
  - given size of corpus there is no way that manual annotations are added





## Results for ontology matching for even reports

- 85% of documents in the first 20 hits a
  - Compare with keywords: 56%
- 40% of relevant documents are in the pages
  - Compare with keywords: 57%
- Ontology matching implies
  - Reading a limited amount of irrelevant dc
  - Risking missing many documents
  - It is possible to count the events





# Issues and Solutions

- Ontology can be extended
  - But increases effort in indexing
    - Equivalent to extending metadata in SDM
  - **But it is impossible to foresee all uses of information**
    - Ontology will always be insufficient somehow
- Information Extraction can be used to reduce burden of annotation
  - But some parts are irretrievable





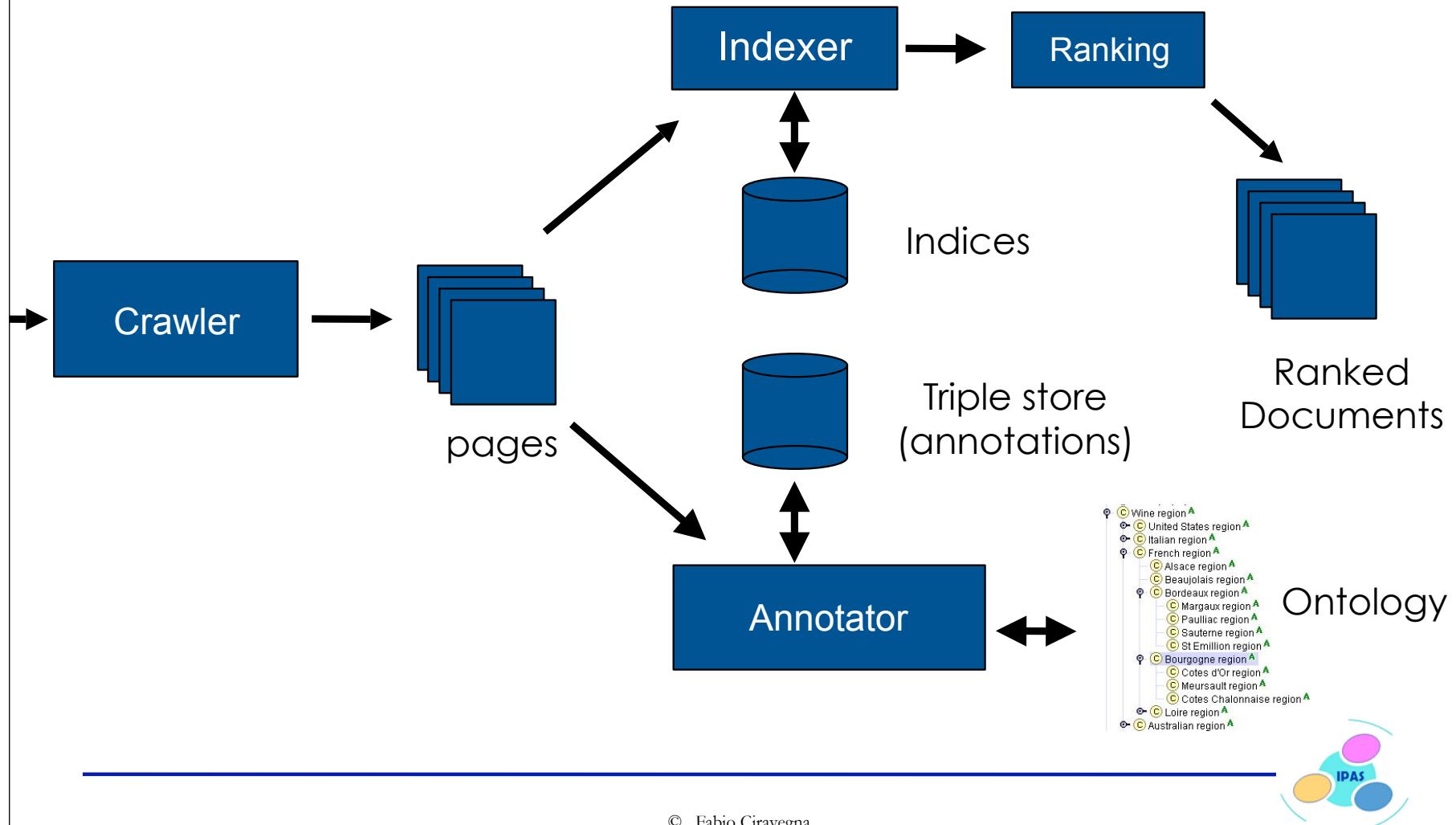
## Hybrid Search (keywords+ontology)

- Mixes keyword and ontology based search
- Ontology based search
- Traditional keyword search
- Keyword in context of ontology-based annotations
- Potential queries:
  - Return all documents where the word fuel is mentioned
  - Return all documents where the affected part description includes the word fuel
  - Return all documents where the affected part description is similar to “fuel duct”
  - Return all documents where the affected part description is equal to “fuel duct” (URI=XXXXX)

affected parts is concept in ontology

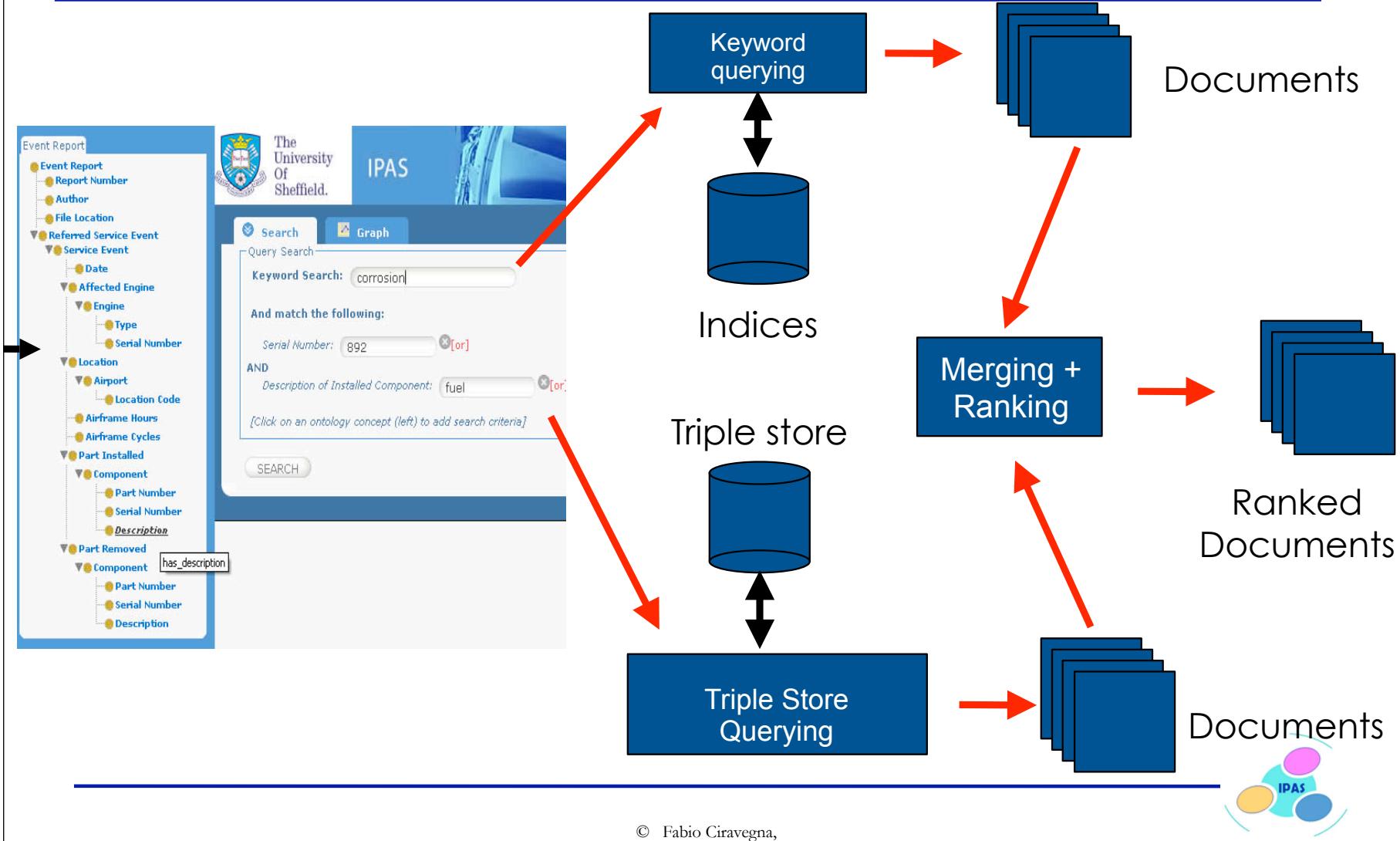


# Hybrid Indexing/Annotation





# Hybrid Search





# Advantages with Hybrid Search

- Accuracy of Ontology-based searching available
  - When metadata covers information
- Expressiveness of Keyword querying is available
  - For all other cases
- Keyword-in-context available
  - Keyword matching available for matching concepts names
    - e.g. match “fuel” in the description of the removed parts
  - Uses provenance of annotations
    - Portion of document annotated with concepts are stored in 3store
    - Keyword matching applied only on the relevant strings
      - e.g. “fuel” is matched only on snippets of texts annotated as removed parts





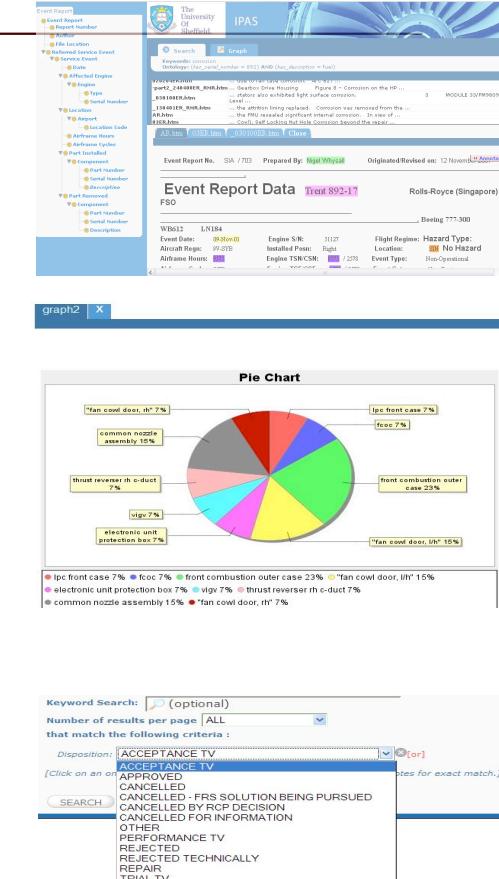
## Results for Hybrid Search

- 83% of documents in the first 20 hits are relevant
  - K:56% O:85%
- 85% of relevant documents are in the first 2 pages
  - K: 57% O:47%
- $F(1)=84\%$ 
  - K:57% O:54%
- Hybrid Search implies
  - Reading a limited amount of irrelevant documents
  - Being able to retrieve easily a very large part of documents



# K-Search

- Enables querying documents based on
  - keywords
  - ontology
  - keywords in context
- Enables quantification of unstructured information
- Currently applied to:
  - Event Reports (1998-2004),
  - Technical variances
- Finalist of Rolls-Royce Creativity Award 2007





## Query results

- Results are displayed as a list
- User can click on a document and open it in the lower frame
- The document will be enriched by annotations with attached services
- Multiple documents can be opened in a tab interface



The  
University  
Of  
Sheffield.

# Query results

Event Report

- Event Report
- Report Number
- Author
- File Location
- Referred Service Event
  - Service Event
    - Date
    - Affected Engine
      - Engine
        - Type
        - Serial Number
    - Location
      - Airport
        - Location Code
      - Airframe Hours
      - Airframe Cycles
    - Part Installed
      - Component
        - Part Number
        - Serial Number
        - Description
    - Part Removed
      - Component
        - Part Number
        - Serial Number
        - Description

The University Of Sheffield. IPAS

Search Graph

Keywords: corrosion  
Ontology: (has\_serial\_number = 892) AND (has\_description = fuel)

020204ER.htm ... due to fan case corrosion. A/C 627 ...  
-part2\_240400ER\_RMR.htm ... Gearbox Drive Housing Figure 8 - Corrosion on the HP ...  
\_030100ER.htm ... stators also exhibited light surface corrosion. 3 MODULE 33/FM9809  
\_130401ER\_RMR.htm ... the attrition lining replaced. Corrosion was removed from the ...  
AR.htm ... the FMU revealed significant internal corrosion. In view of ...  
03ER.htm ... Cowl), Self Locking Nut Hole Corrosion beyond the repair ...

AR.htm 03ER.htm \_030100ER.htm Close

Event Report No. SIA / 703 Prepared By: name of person Originated/Revised on: 12 November 2001 [Annotation](#)

## Event Report Data engine name here Rolls-Royce place here

FSO

Boeing 777-300

WB612 LN184

Event Date: 09-Nov-01

Engine S/N: 51127

Flight Regime: Hazard Type:

Aircraft Regn: 9V-SYB

Installed Posn: Right

Location: SIN No Hazard

Airframe Hours: 9575

Engine TSN/CSN: TSN

Event Type: Non-Operational





# Graph visualisation

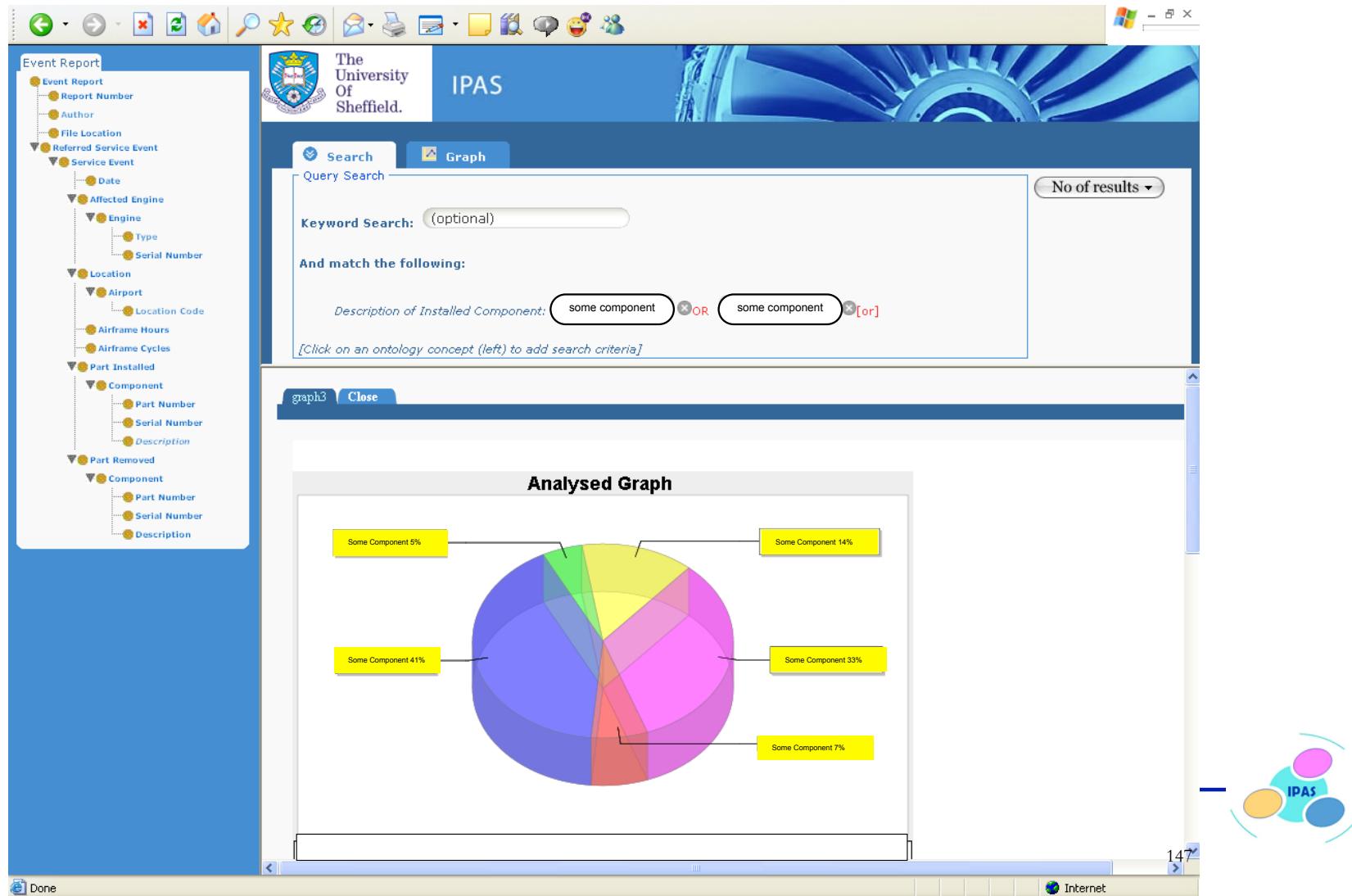
- Query results can be visualised as a graph
- Users can select graph types
  - Bar chart
  - Pie chart
- Users can select group and subgroup
  - For aggregating results
- Graphs are opened in the tab interface alongside other documents



The  
University  
Of  
Sheffield.

# Graph - example

- Percentage of
- occurrences of installations
- of fuel based parts
- Pie chart





# Process Support

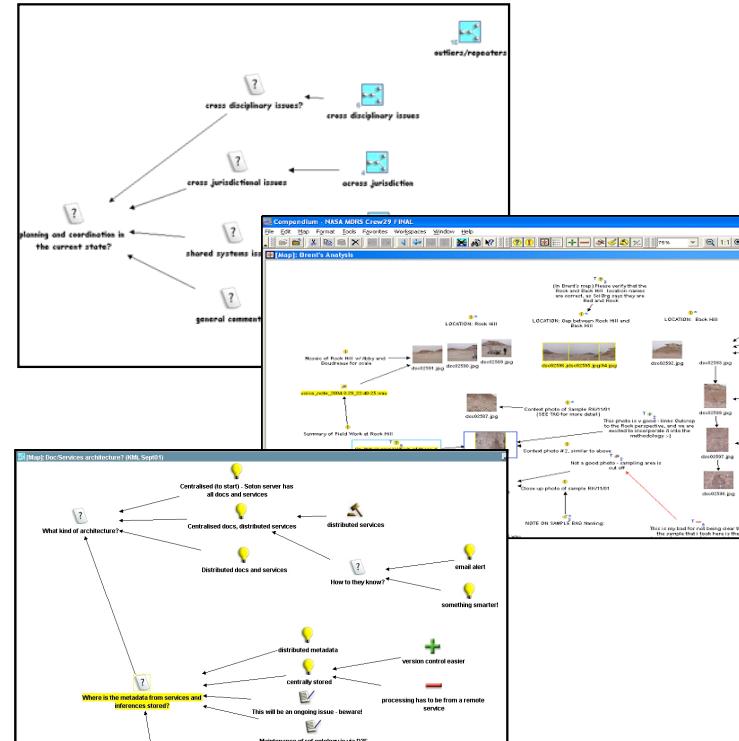
- Goals:

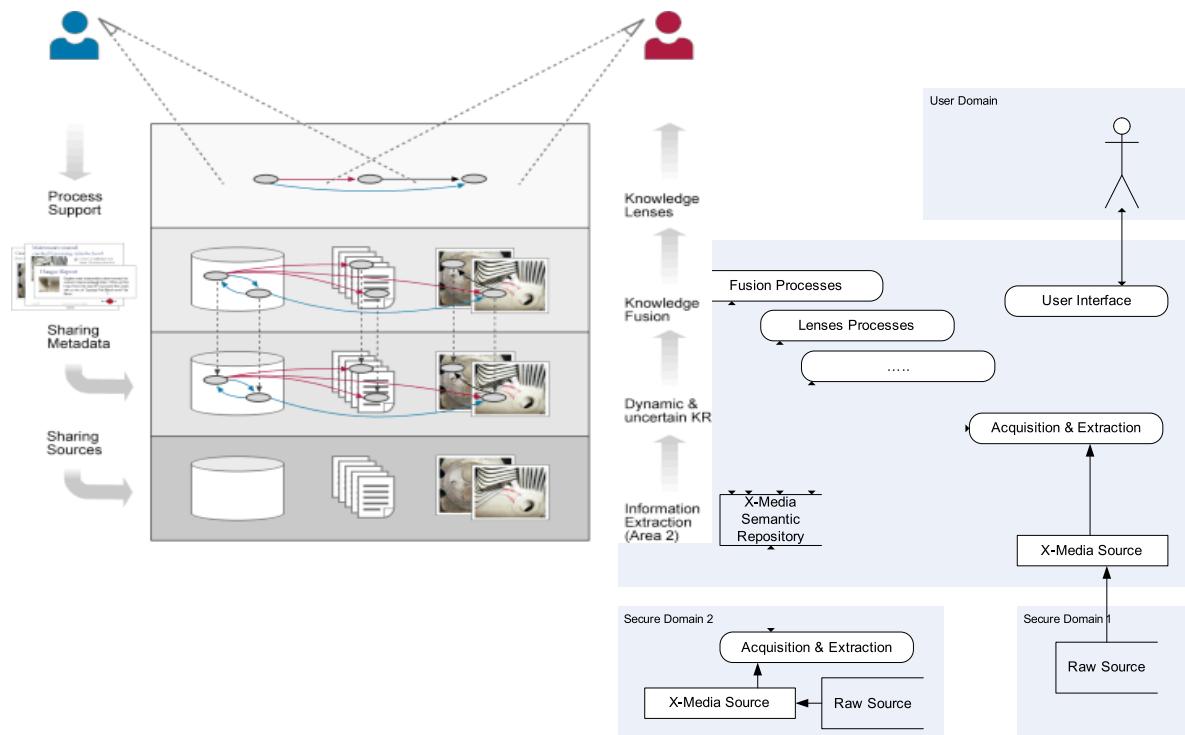
- Supporting users in their tasks
- Enabling capturing knowledge while knowledge is created



# Example: design rationale capturing

- During the discussion, the working group will consider many alternative solutions
  - Those discarded are not in the final document
- When next engine is designed, the group needs to know
  - What solutions were tried (use of titanium)
  - Why they were not adopted (e.g. too high a cost)
  - If the analysis is still true (titanium cost has decreased)
- Compendium (Buckingham-Shum 2002)
- D/Red (Wallace et al. 2005, 2007)





## Other issues: Architectures and privacy

ISWC 2007

© Fabio Ciravegna, University of Sheffield



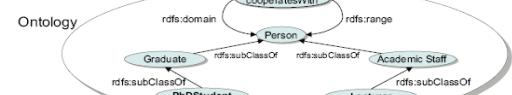
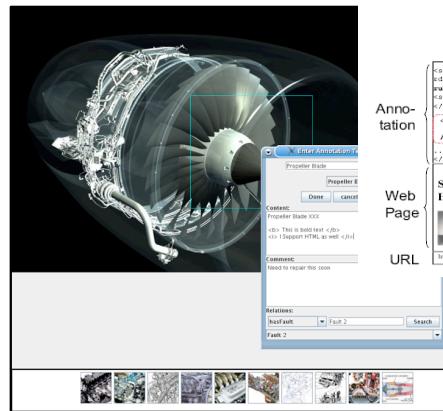
- Infrastructure:

- Different media cannot easily be shared in the same way as knowledge from one format.
  - A folder of text documents may still be easily sent via email, but a folder of image files may not, and may instead require a shared image repository.
  - For 10 GByte of noise sensor data even an upload to a centralized repository may be out of the question and instead remote access to the underlying data base is to be considered.
- A complex infrastructure is needed in order to implement knowledge management across media.



# Privacy and Confidentiality

- Not all information and knowledge is to share
  - Personal knowledge
  - Departmental knowledge
  - Organisational knowledge
  - Public knowledge
- Protecting information is important
- Acquired Knowledge must be shared with care
  - All knowledge must be marked with provenance and confidentiality (Lanfranchi *et al.* 2004)
  - Question: can a piece of knowledge derived also from confidential data (e.g. statistics) be shown to everyone?
    - If not, what do you show to non allowed users? False deductions?



**Annotation**

**rdf:type**

```

<swc:PhDStudent
  rdf:about="http://www.aifb.uni-kassel.de/nwb/de/PhDStudent"
  rdfs:label="Siegfried Handschuh"
  rdfs:subClassOf="Person"/>

```

**rdf:type**

```

<swc:Lecturer
  rdf:about="http://www.aifb.uni-kassel.de/nwb/test/#Steffen"
  rdfs:label="Steffen Staab"
  rdfs:subClassOf="Person"/>

```

**Web Page**

**URL**

Siegfried Handschuh AKTIVE Media Version 1.6

Annotation Image RDF Mode Help

File

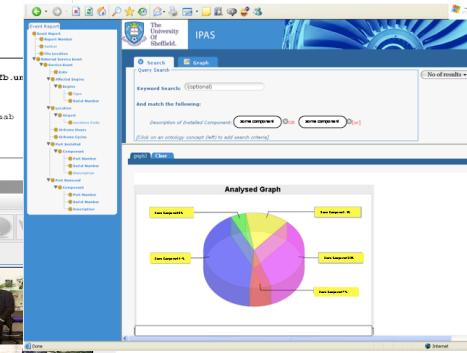
Comment: Need to repair this icon

Relations: hasTask [ Task 2 ] Search

Annotations:

- visitingEntity [ ]
- INposition [ ]
- EDposition [ ]
- EDgroup [ ]
- EDinstitution [ ]
- date [ ]
- relation [ ]
  - st\_time [ ]
  - in\_location [ ]

Annotation Relation



## Conclusions

ISWC 2007



# Conclusions and Future Work

- Knowledge Management is moving towards large scale
  - Initially expected around 2010 now already happening
- The Semantic WEB offers potentially key technologies to the development of future KM
  - More Web than Semantics, but:
    - A little semantics goes a long way (J. Hendler)
- The potential must be exploited addressing real world requirements
  - Rather than in principle AI-oriented requirements (e.g. closed world, small scale, etc.)
- Strong application pull can be obtained
  - Do not sell slogans, sell ideas and applications!



Rolls-Royce

**Kodak**



ty of Sheffield



**ISWC 2007**



- Folksonomies:

- Easier to use and implement than ontologies
  - SW needs to make the best out of them (Specia et al. 2007)
- Work very well as approximations of ontologies in many applications and tasks

- Blogs

- The new frontier of knowledge sharing (e.g. Google)
- Serious risks seen by the companies for information leak and corporate responsibility
  - Whistleblowers and real concerns about putting in writing
- Semantic blog to acquire and share information



## Future Trends: Web 2.0 (ctd)

- Semantic email as a way to trace what is in emails
- Wikis
  - Collaborative working made easier
  - Semantic Wikis a way to acquire and share knowledge in a more effective way
- In general: collaborative thinking of Web 2.0 can potentially impact KM
  - Social aspects:
    - Flink for expert finding
    - Importance of social connection in the current organisation to be enabled, not prevented
  - Semantic Compendium to help capture rationale of design.



# Thank You

- Contact Information

- [www.dcs.shef.ac.uk/~fabio](http://www.dcs.shef.ac.uk/~fabio)
- fabio@dcs.shef.ac.uk

- Intelligent Web Technologies Lab

- <http://nlp.shef.ac.uk/wig/>

- NLP Sheffield

- <http://nlp.shef.ac.uk/>

- University of Sheffield

- [www.shef.ac.uk](http://www.shef.ac.uk)

The image contains three separate screenshots of University of Sheffield websites:

- Top Screenshot:** A personal profile page for Fabio Ciravegna. It features a photo of him, his name, and a "quick links" sidebar with links like "my home page", "publications", and "curriculum vitae".
- Middle Screenshot:** The homepage of the "web intelligence technologies lab". It includes a photo of Fabio, his title as "professor of language and knowledge technologies", and details about the lab's research focus on tomorrow's Web.
- Bottom Screenshot:** The main homepage of the University of Sheffield. It features the university crest, news items, and links to various departments and services.



# A very Incomplete Bibliography

- F. Ciravegna: Challenges in Information Extraction from Text for Knowledge Management, in S. Staab, (ed), "Human Language Technologies for Knowledge Management", IEEE Intelligent Systems and Their Applications (Trends and Controversies), Vol. 16, No. 6, pp 88-90, 2001.
- Fabio Ciravegna. Adaptive information extraction from text by rule induction and generalisation. In Proceedings of 17th International Joint Conference on Artificial Intelligence (IJCAI), 2001. Seattle.
- H. Cunningham, D. Maynard, K. Bontcheva, V. Tablan. GATE: A Framework and Graphical Development Environment for Robust NLP Tools and Applications. 40th Anniversary Meeting of the Association for Computational Linguistics (ACL'02). 2002.
- I. Muslea, S. Minton, and C. Knoblock. 1998. Wrapper induction for semistructured webbased information sources. In Proceedings of the Conference on Automated Learning and Discovery (CONALD), 1998.
- Vitaveska Lanfranchi, Fabio Ciravegna, Daniela Petrelli: Semantic Web-based Document: Editing and Browsing in AktiveDoc, Proceedings of the 2nd European Semantic Web Conference , Heraklion, Greece, May 29-June 1, 2005
- Handschuh, Staab, Ciravegna. S-CREAM - Semi-automatic CREAtion of Metadata (2002) <http://citeseer.nj.nec.com/529793.html>
- F. Ciravegna, A. Dingli, D. Petrelli, Y. Wilks: User-System Cooperation in Document Annotation based on Information Extraction. Knowledge Engineering and Knowledge Management (Ontologies and the Semantic Web), (EKAW02), 2002.
- M. Vargas-Vera, Enrico Motta, J. Domingue, M. Lanzoni, A. Stutt, and F. Ciravegna. MnM: Ontology driven semi-automatic or automatic support for semantic markup. In Proc. of the 13th International Conference on Knowledge Engineering and Knowledge Management, EKAW02. Springer Verlag, 2002



## A very Incomplete Bibliography (ctd)

- Fabio Ciravegna. Designing adaptive information extraction for the Semantic Web in Amilcare. In S. Handschuh and S. Staab, editors, Annotation for the Semantic Web, Frontiers in Artificial Intelligence and Applications. IOS Press, 2003.
- C. Goble, S. Bechhofer, L. Carr, D. De Roure, and W. Hall. Conceptual Open Hypermedia = The Semantic Web? In The Second International Workshop on the Semantic Web, pages 44–50, Hong Kong, May 2001
- Fabio Ciravegna, Sam Chapman, Alexiei Dingli, and Yorick Wilks: Learning to Harvest Information for the Semantic Web, Proceedings of the First European Semantic Web Conference, Crete, May 2004
- A. Kiryakov, B. Popov, et al. Semantic Annotation, Indexing, and Retrieval. 2nd International Semantic Web Conference (ISWC2003), <http://www.ontotext.com/publications/index.html#KiryakovEtAl2003>
- S. Dill, N. Eiron, et al: <http://www.tomkinshome.com/papers/2Web/semtag.pdf> . SemTag and Seeker: Bootstrapping the semantic web via automated semantic annotation. WWW'03.
- Thomas Leonard and Hugh Glaser. Large scale acquisition and maintenance from the web without source access. In Siegfried Handschuh, Rose Dieng-Kuntz, and Steffen Staab, editors, Proceedings Workshop 4, Knowledge Markup and Semantic Annotation, K-CAP 2001, 2001
- Martin Dzbor, John B. Domingue, and Enrico Motta. Magpie - towards a semantic web browser. In Proceedings of the 2nd Intl. Semantic Web Conference, October 2003. Sanibel Island, Florida
- Alexander Maedche, Steffen Staab, Nenad Stojanovic, Rudi Studer, York Sure: SEMantic portAL - The SEAL approach In D. Fensel, J. Hendler, H. Lieberman, W. Wahlster (eds.), Spinning the Semantic Web, pp. 317-359. MIT Press, Cambridge, MA., 2003.



## A very Incomplete Bibliography (ctd)

- Natalya F. Noy and Deborah L. McGuinness: Ontology Development 101: A Guide to Creating Your First Ontology, [http://protege.stanford.edu/publications/ontology\\_development/ontology101-noy-mcguinness.html](http://protege.stanford.edu/publications/ontology_development/ontology101-noy-mcguinness.html)
- Elena Paslaru Bontas, Christoph Tempich, York Sure : OntoCom: A Cost Estimation Model for Ontology Engineering, In: Proceedings of the 5th International Semantic Web Conference (ISWC 2006), November 5-9, 2006, Athens, GA, USA, LNCS. Springer.
- Ajay Chakravarthy, Vita Lanfranchi and Fabio Ciravegna: Cross-media Document Annotation and Enrichment, SAAW2006 - 1st Semantic Authoring and Annotation Workshop, The 5th International Semantic Web Conference (ISWC2006), Athens, GA, USA, Monday, November 6th 2006
- R. Gaizauskas and G. Demetriou and P. Artymiuk and P. Willett: Protein Structures and Information Extraction from Biological Texts: The PASTA System, Journal of Bioinformatics 19(1), 135-143, 2003
- Vitaveska Lanfranchi, Ravish Bhagdev, Sam Chapman, Fabio Ciravegna, Daniela Petrelli: Extracting and Searching Knowledge for the Aerospace Industry, in Proc. of 1st European Semantic Technology Conference, Vienna, May 2007