

Semantic Web Technologies for Knowledge Management in Large Distributed Organisations

Professor Fabio Ciravegna

Web Intelligence Technology Lab,
Department of Computer Science,
University of Sheffield

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



image © Rolls-Royce

Sponsored by



www.x-media-project.org



www.3worldt.org

Copyright Notice

2

- These slides were presented during the Semantic Web Days in Stavanger, Norway, April 2008
(<http://www.posccaesar.com/en-GB/PortalObject/2803/POSCCaesar.aspx>)
- Condition of use:
 - the use is limited to personal or educational purposes
 - the copyright footnote must always be visible when slides are presented
 - full recognition is given to me for the paternity of the slides and information contained
 - the context in which the slides are used/presented must be appropriate and not damaging of the image of the University of Sheffield or mine.
- Fabio Ciravegna, University of Sheffield,
fabio@dcs.shef.ac.uk
<http://www.dcs.shef.ac.uk/~fabio/>



Outline of Tutorial

3

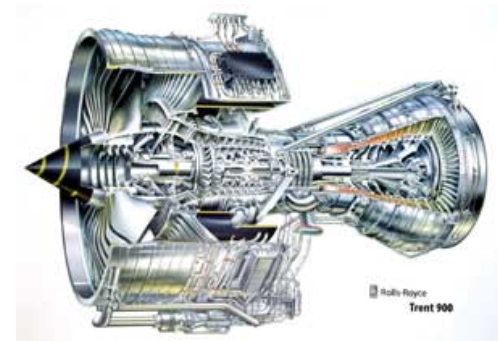
- 13.00-13.15 Introduction and scene setting
 - Issues in knowledge management in large organisations
- 13.15-13.40 Ontologies
- 13.40-14.40 Semantic web technologies for knowledge acquisition
 - 14.40-15.10 Coffee break
- 15.00-16.10 Semantic web technologies for knowledge sharing, reuse and retrieval
- 16.10-16.40 ontology engineering
- 16.40-17.00 Conclusion and future work
- 17.00- Discussion



- Large Scale Knowledge Management across Media
- 2nd largest project funded on Knowledge and Content Technologies by the EU:
 - €13.6M from European Commission
- 2006-2010
- 15 partners
 - Users: Fiat, Rolls Royce
 - a2mac1 new partner at zero cost
- In industrial board:
 - Kodak, Philips, Telenor, CESI, DS&S



- Designing Integrated Products And Services in Manufacturing
- Funded: Rolls Royce plc (50%), DTI (UK) 50%
- January 2005-December 2008
- Value: £2.5M (~€3.3M)
 - £1.2M provided by industrial partners
 - Coordinator: Rolls Royce plc (jet engines)
 - £225,000 (+ £20,000) for Sheffield
- Partners:
 - Universities of Aberdeen, Cambridge, Sheffield, Southampton
 - Rolls Royce, DS&S, Epistemics



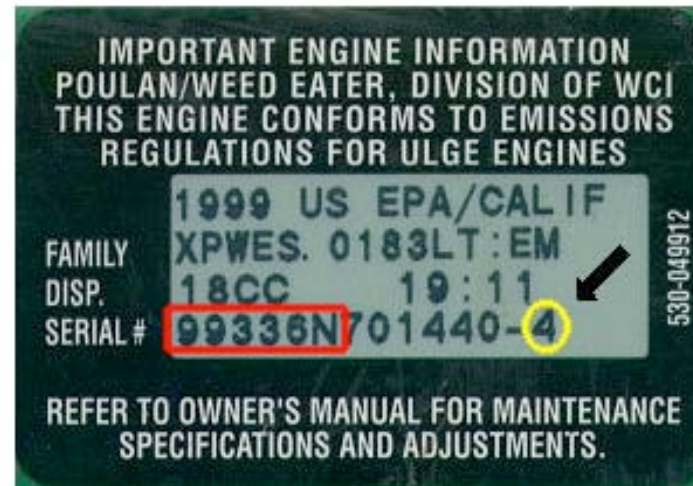
- Gathering knowledge relevant to a task or problem
 - it may be distributed across different storage systems and different media
- Analysing the knowledge they have gathered and make sense of it
- Sharing knowledge with their colleagues
- Keeping track of the process
 - by being aware of what one is 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
└──┬──┘
└──┬──┘
└──┬──┘
Day of the Year
Year of Production

4 = Product Type

- a jet engine can produce ~1 Gbyte 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

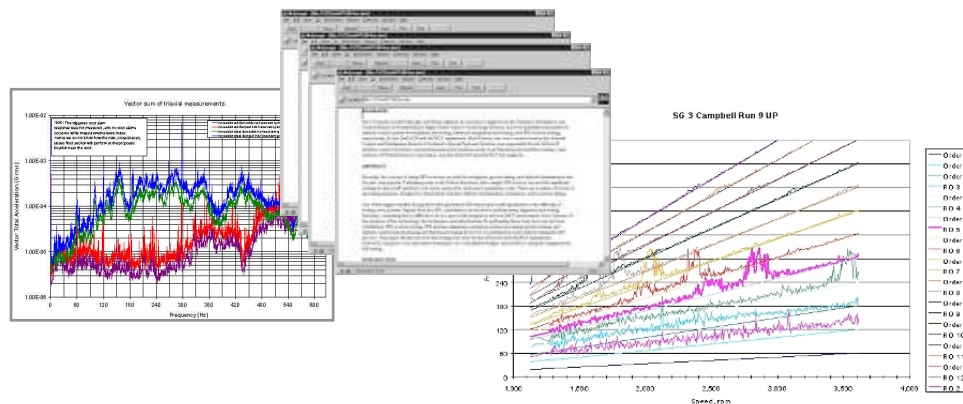
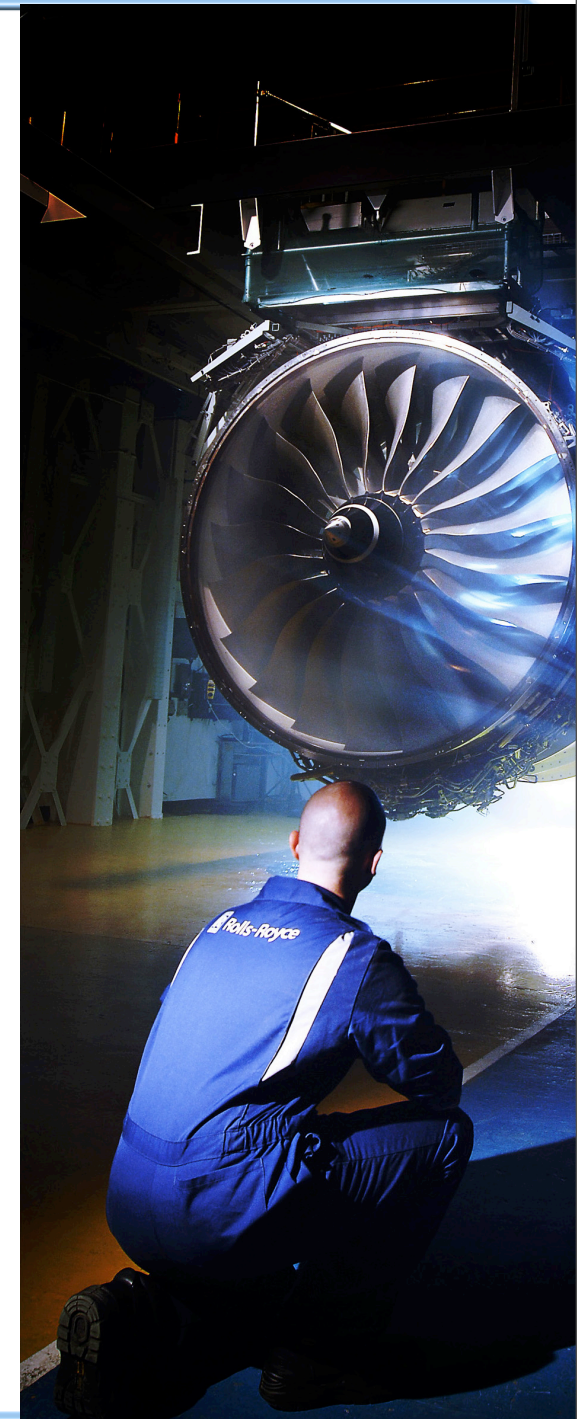


image © www.rolls-royce.com

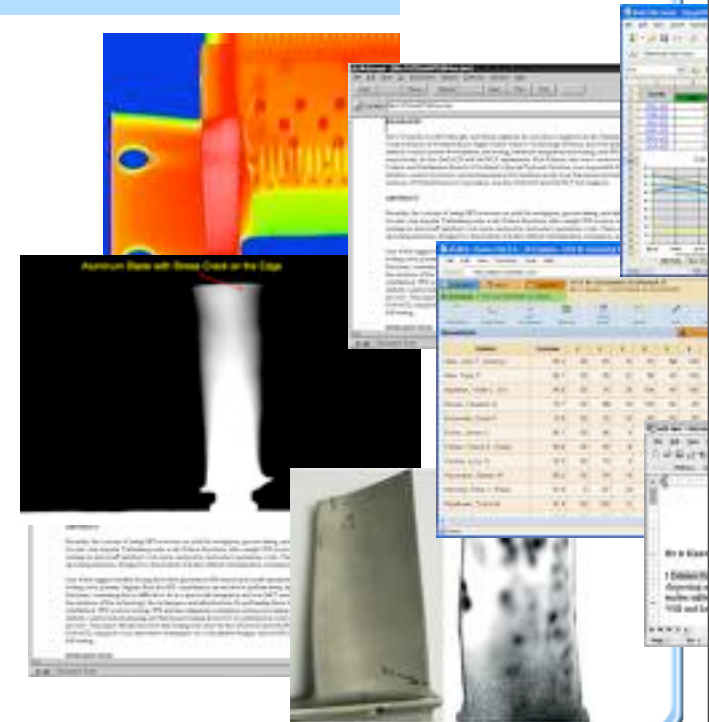


When engine is serviced (e.g. overhaul)

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



image © Rolls-Royce plc



- 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

Document Type

AROC proforma

AROC results

Development

EHM data

Emails

ONWING emails

Images

Lab findings

Monitoring Requirements

Presentations

Procedures

RCP

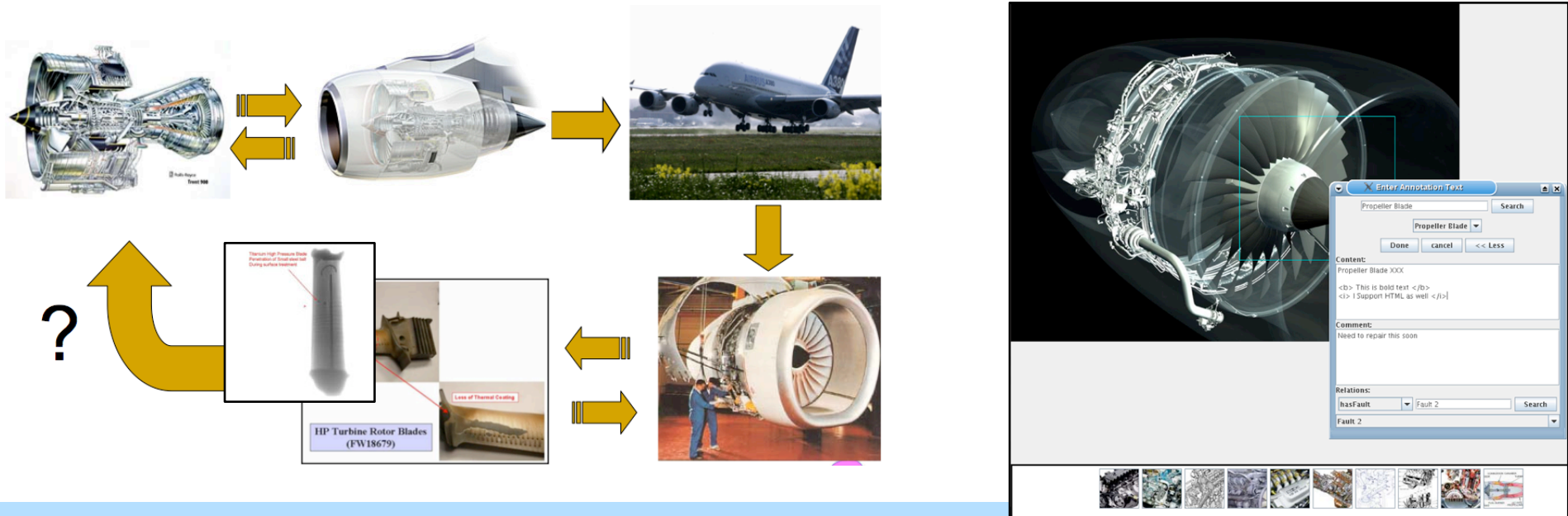
Risk Assessment

Solution Reports

Technical Reports

TS&O Reports

- 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

- Organisation archives are following a Web-like trend
 - Massive shift towards multi- and cross-media
 - Text, Images, Data
 - Videos
 - Large scale
 - Dramatic reduction of memory and storage cost
 - Increase in speed and capacity
 - Number and size of distributed archives of information
 - Large organisations' Intranets as mini webs
 - Thousands of computers
 - Hundreds of repositories
 - Hundreds of millions of documents

An image is worth 1,000 words

– But it is very difficult to index

in 2007: 4 billion cameras

New technical information doubles every 2 years:

• Expected to double every 72 hours by 2010

Web: the largest source of data

13

- User Generated Content
 - With 185+ million registered users (April 2007) MySpace would be the 6th largest country (between Brazil and Pakistan)
 - The average MySpace page is visited 30 times a day, to access (see on teh right)

Images

1+ billion

Millions uploaded / day

150,000 requests / second

Songs

25 million

250,000 concurrent streams

Videos

60 TB

60,000 uploaded / day

15,000 concurrent streams

Servers

6,000 web

650 ad

250 database

Content (ctd)

14

Source: Michael Brodie, Verizon Inc. Sept 2007

Facebook has ...

1.8 billion photos, 31 million active users

10^5 new users / day and 1,800 applications

YouTube Videos

1.7 billion served / month

1 million streams / day

= 75 billion e-mails

The world's 4+ billion devices

- cameras, phones, PCs, CCTVs,

They will increase 50% by 2010.

Source: Michael Brodie, Verizon Inc. Sept 2007

Originality of Information

15

- The data worldwide is
 - 25% original; 75% replicated
 - 25% from the workplace; 75% not
 - 95% unstructured and growing



This presentation is made with 95% of recycled material



Ask yourself:
Is your company
following
a similar path?

Why manage knowledge?

17

- To enable easy timely and effective reuse
 - We need: to enable sharing
 - Requirements: easy and effective sharing
- To enable sharing
 - we need to: acquire knowledge
 - Desiderata:
 - Easy acquisition (do not get in the way of the user's work!)
 - Comprehensive acquisition (do not miss important facts!)
- To enable acquisition:
 - We need modelling the domain and process in an appropriate way

Please note: most books and tutorial work the other way around.

They start with modelling (e.g. ontology building) then move to acquisition, then to sharing (if they do!). This often generates confusion: modelling seems the most important issue!!



- 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



■ Features

- Single Conceptual Schema for official agreed view
 - supporting communication between different parts of organisation
- Large homogeneous knowledge or document repositories for collection and organisation of knowledge
- Enterprise Knowledge Portal providing unique standardized access to proprietary knowledge

■ Effect:

- Many portals are deserted by users
 - replacements: non-official tools such as shared directories, personalized and local databases, etc.
- Reason: difficulty in adopting models, schemas and procedures that are unsuitable to specific communities of users.



- Modern KM is based on dynamic communities
 - that acquire and share knowledge according to dedicated schemas
 - existing across traditional organisational boundaries
 - ill fit pre-determined standard schemas
 - require rapidly tailoring knowledge for their specific ad-hoc uses
 - often outside the company (outsourcing)
- Requirement:
 - independence of ontological views
 - communities must share their knowledge with rest of organisation



What we know and what we do not

21

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

- Lack of efficient publishing [of] 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

A.T. Kearney, Network Publishing Study, 2001
<http://www.computeruser.com/news/01/04/17/news1.html>

- 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.

- 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 inexistent
- 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



- 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...



- 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

- Corporate archives are difficult to cope with:
 - 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, csu are present in 50% of jet engine event report
 - Synonyms are not captured by keyword matching
 - Fuel Metering Unit, FMU, Metering Unit, fmu701mk5, S/N3332223



- 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)



- Can we trust a system answer?

$$Precision = \frac{COR}{\min(ACT, \max No)}$$

- COR= correct answer by system
- ACT= no results returned by system

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

- Are we finding all relevant documents?
 - EXP= no of documents relevant to the query in archive



Results for keyword matching

30

- 56% of documents in first 20 hits are relevant
 - Precision=56%
- 57% of relevant documents are in 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



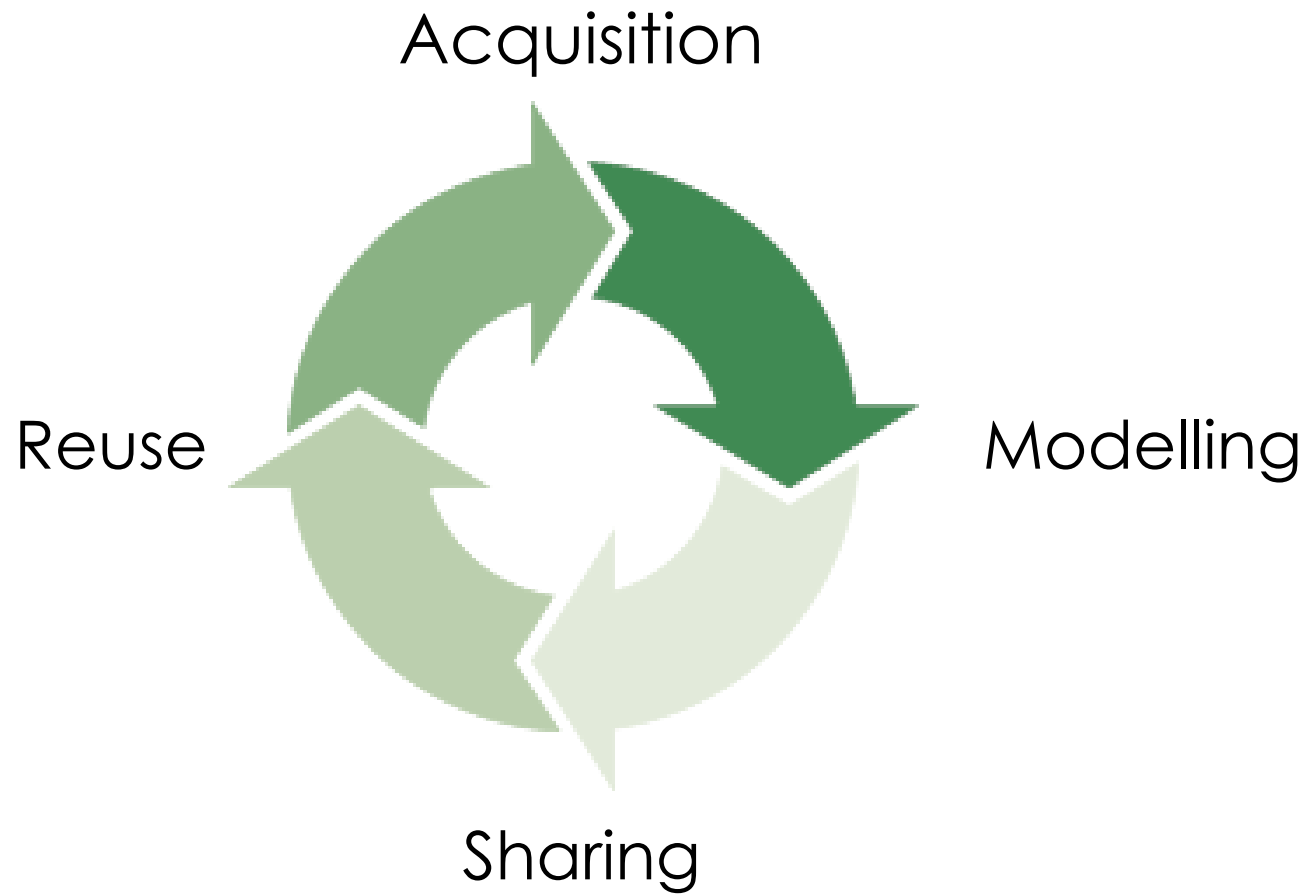
- Most resources available over Intranets.
 - The Web is medium for access to multiple archives in a seamless way
 - It enables remote access independently from geographic distribution
 - it provides a common protocol for communication
 - it provides an interaction modality familiar to users



■ However:

- Scale: dozen to hundred million documents
- Security: risk of information leaks and hackers' attacks limits access
- User disorientation: multiple points of access





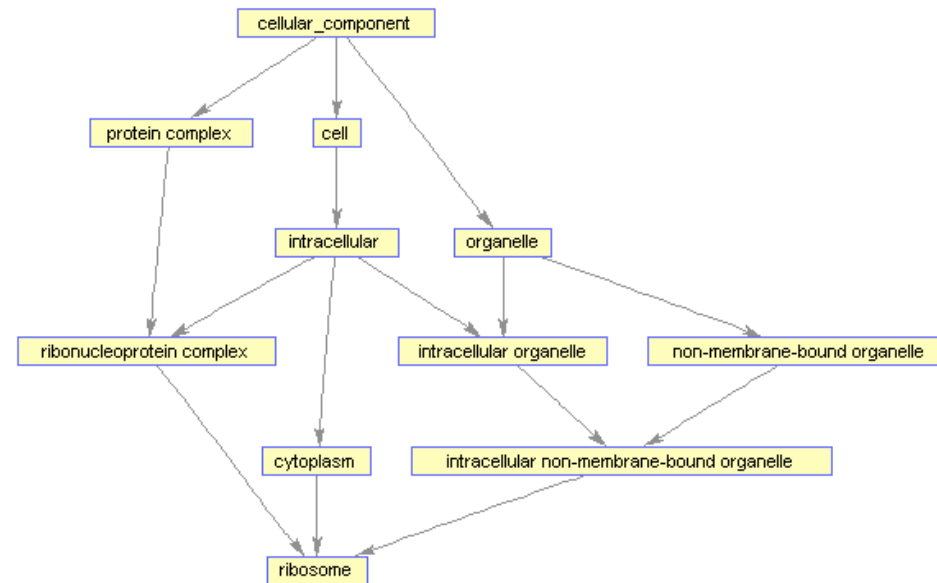
The Knowledge Life-cycle

Steps in Knowledge Life-cycle we will see

34

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





Ontologies

- 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

1. It describes a domain

2. A formal specification

3. Agreed by a community

4. 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

- 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?

38

- To share common understanding of 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



- 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 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

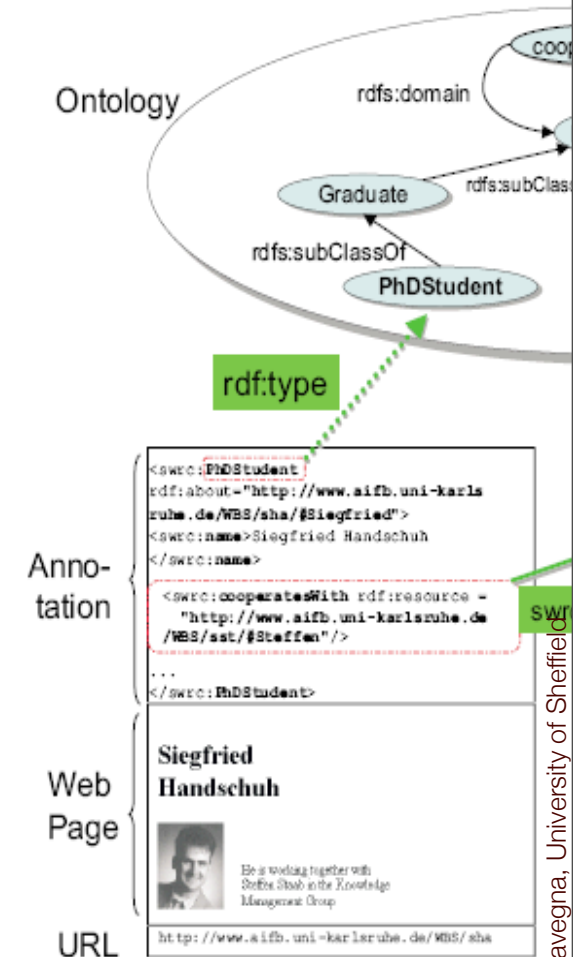


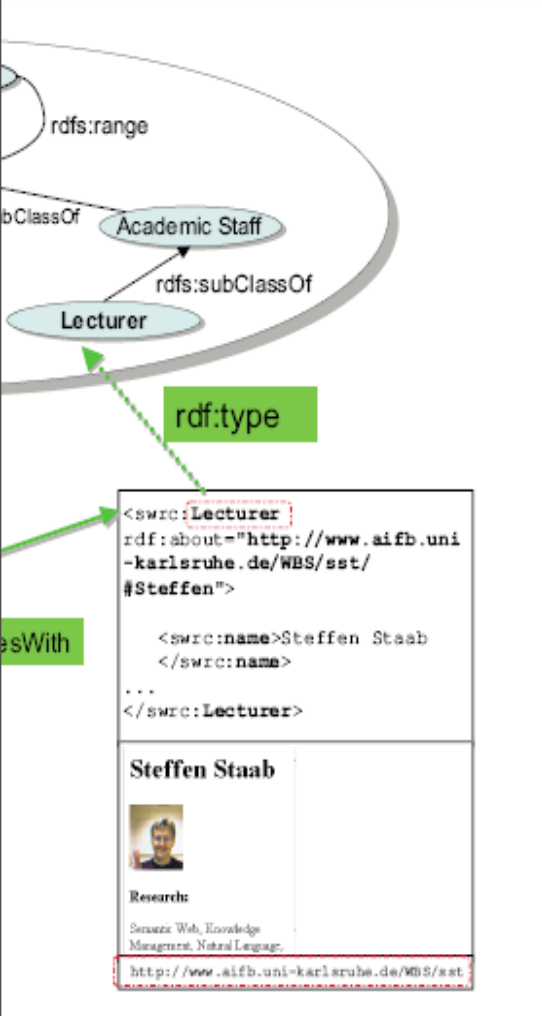
Ontologies and Knowledge Management

42

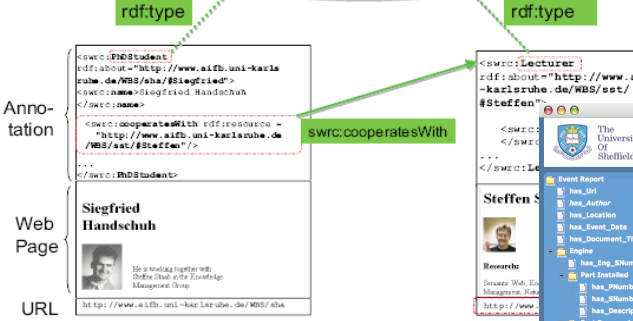
Motivations for use:

- 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





- 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

[illegible]

AK-234-I

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

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

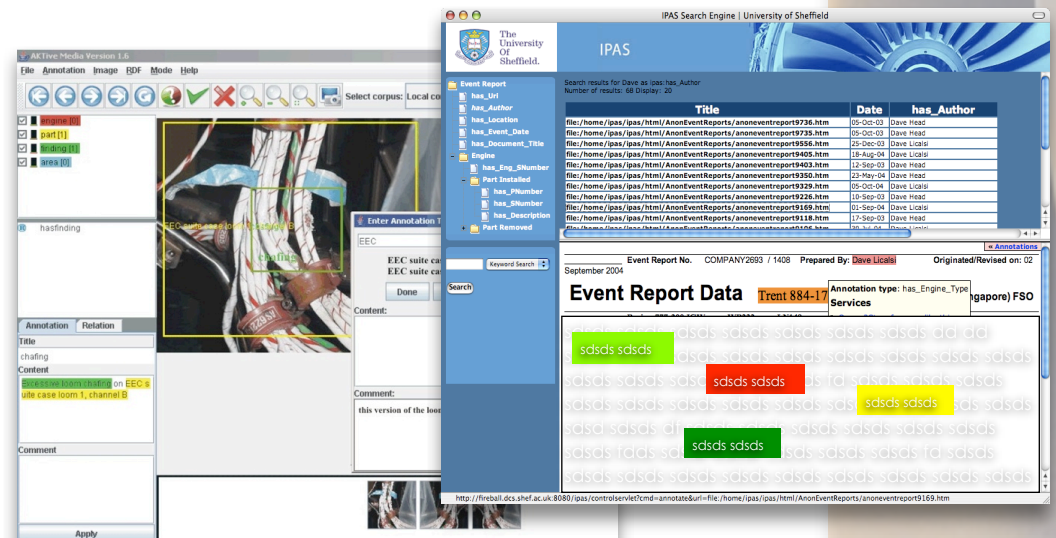
Knowledge Acquisition

45

- Collecting and aggregating multimedia knowledge to make it available for
 - sharing and reuse
 - From document management to knowledge management
 - for integration

In ontological terms knowledge acquisition consists in capturing instances!

- Approaches
 - at source: helping people capturing knowledge when produced
- On legacy documents, pictures, data:
 - Annotation services



- 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, 12-inch PowerBook G4 and 15-inch PowerBook G4.

Apple is voluntarily recalling the affected batteries and has initiated a worldwide exchange program to replace eligible customers with a new replacement battery. This program is being conducted in accordance with the U.S. Consumer Product Safety Act (CPSC) and other international safety regulations.

Identifying your battery

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



- 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 Fiat 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 options.

Every inch of space has been used to provide a range of great 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.

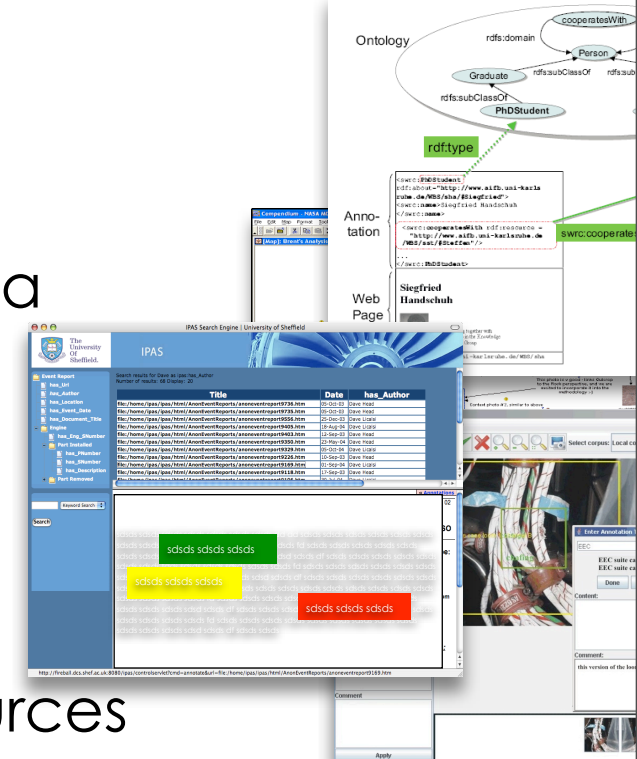
The Fiat - neatly taking care of everything.



Requirements for SW: Robustness

48

- 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

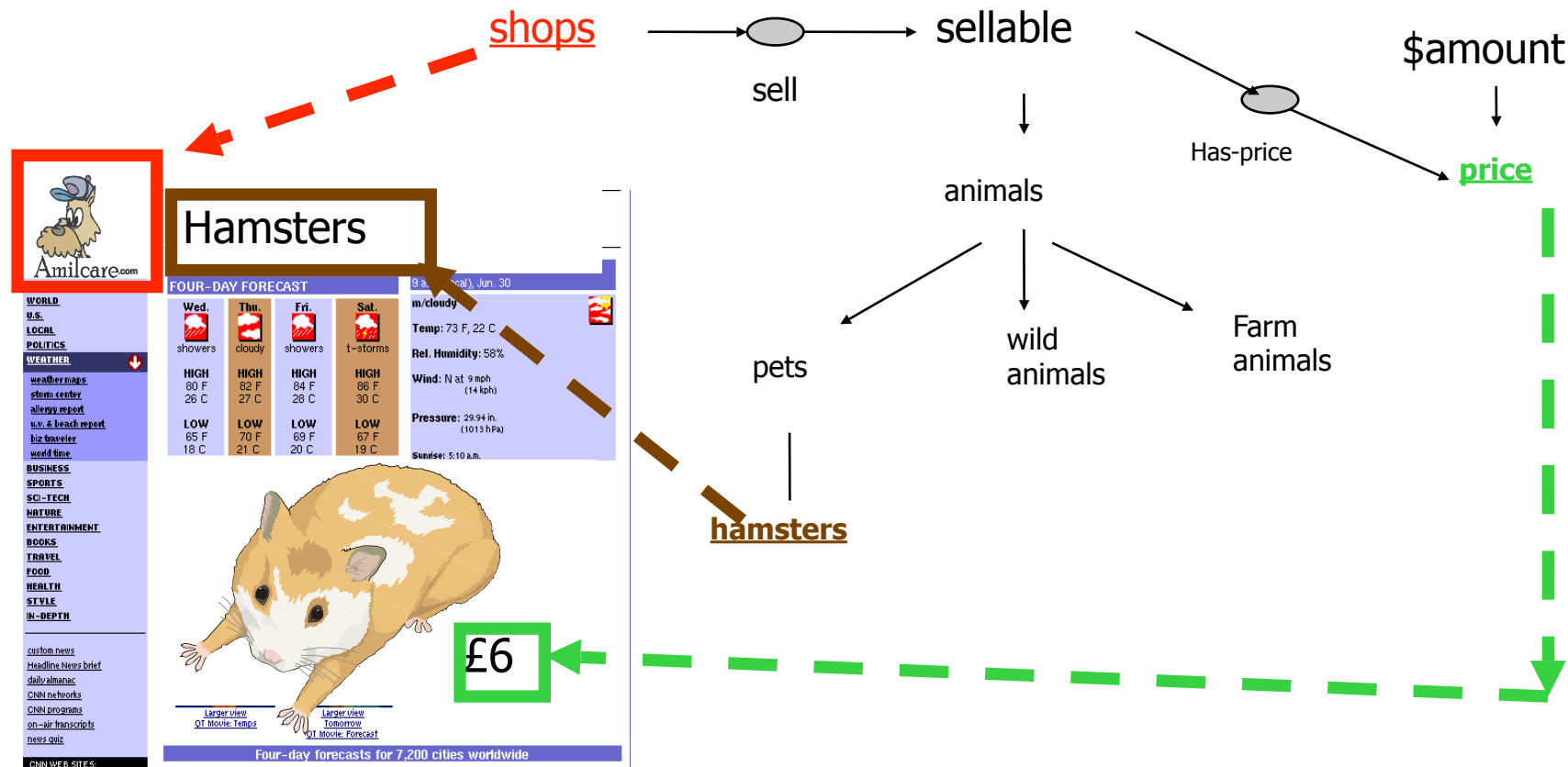


- Acquiring knowledge over large scale resources
 - several thousands of documents
 - hundreds of archives



- Community's requirements
 - independence of ontological views
 - ability to reuse external knowledge
- Organisation's requirements
 - ability to reuse proprietary knowledge
- Knowledge lifecycle:
 - definition of community-specific views of the world;
 - capture and acquisition of knowledge according to them;
 - integration of captured knowledge with the rest of the organisation's knowledge;
 - sharing of knowledge across communities





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



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

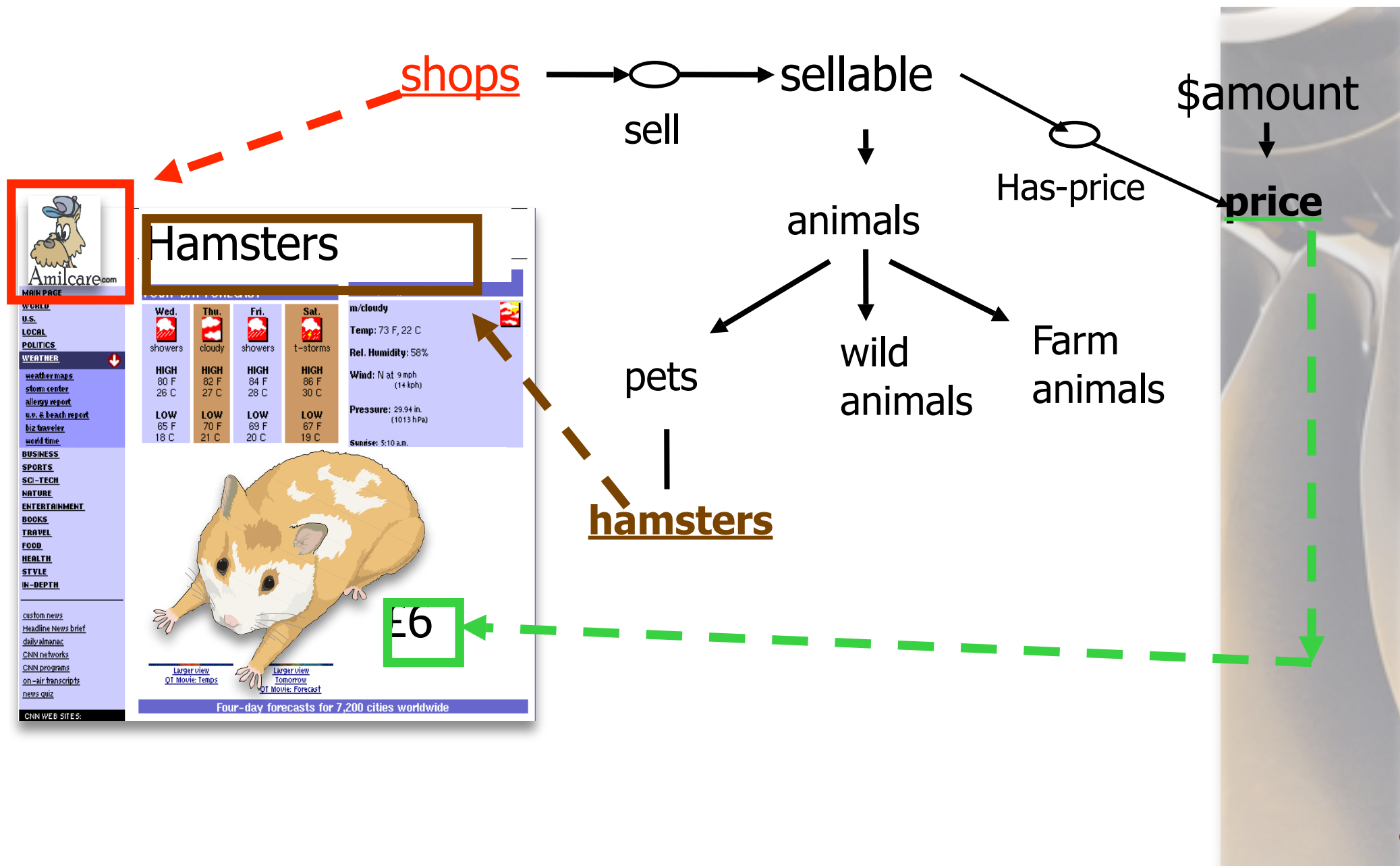


- Marking up contained information
 - Portions of documents associated to objects in ontology
 - Allows:
 - Ontology-driven processing
 - Services based on ontology will be able to use information
 - Ontomat/CREAM (Staab et al 2001)
 - Melita (Ciravegna et al. 2002)
 - SemTag and Seeker (Dill et al.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



Ontology-based Annotation

55



- 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 in simple, and therefore the main body of the resolution is usually clear. [redacted] resolutions 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 make up only 15% of the Legislature. This means six out of [redacted] 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.

Statutes2

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

Why we used these references

Objective	[Click here and type objective]		
Experience	1990-1994	Arbor Shoe	Southridge, SC
National Sales Manager			
<ul style="list-style-type: none"> Increased sales from \$50 million to \$100 million. Doubled sales per representative from \$5 million to \$10 million. Suggested new [redacted] mings by 23%. 			

Objective	[Click here and type objective]		
Experience	1990-1994	[redacted]	Southridge, SC
National Sales Manager			
<ul style="list-style-type: none"> Increased sales from \$50 million to \$100 million. Doubled sales per representative from \$5 million to \$10 million. Suggested new products that increased earnings by 23%. 			

Why we DID NOT use other references

Objective	[Click here and type objective]		
Experience	1990-1994	Arbor Shoe	Southridge, SC
National Sales Manager			
<ul style="list-style-type: none"> Increased sales from \$50 million to \$100 million. Doubled sales per representative from \$5 million to \$10 million. Suggested new products that increased earnings by 23%. 			



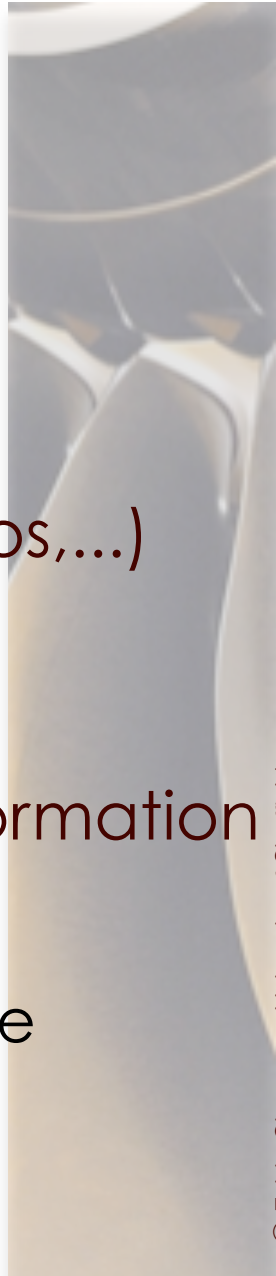
(note: this is not a real legal document!)

- 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)



¹⁰Be/⁷Be ratio calculated in the GISS general circulation model during January and March. Circled areas indicate maximum

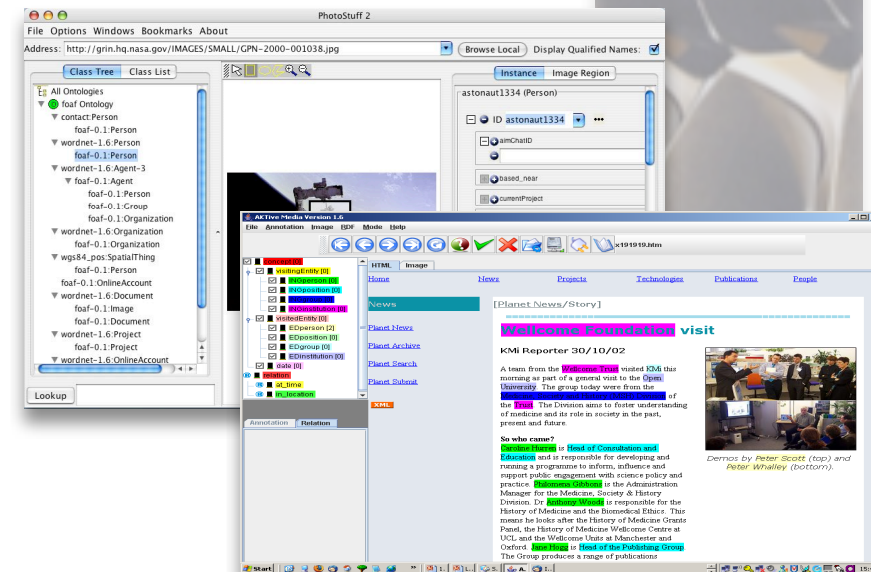
- Input to the KA technologies
 - Ontologies (MMO, domain ontology),
 - Background knowledge (gazetteers, etc.)
 - Normalised document representation
 - Medium to extract from (text, images, data, videos,...)
- Output
 - Evidence represented in terms of conceptual information
 - Evidence used by other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



Ontology-based Annotation

61

- 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

AKTive Media Version 1.6

File Annotation Image RDF Mode Help

Navigation icons: back, forward, search, etc. x191919.htm

HTML Image

Home News Projects Technologies Publications People

News

Planet News

Planet Archive

Planet Search

Planet Submit

XML

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

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?

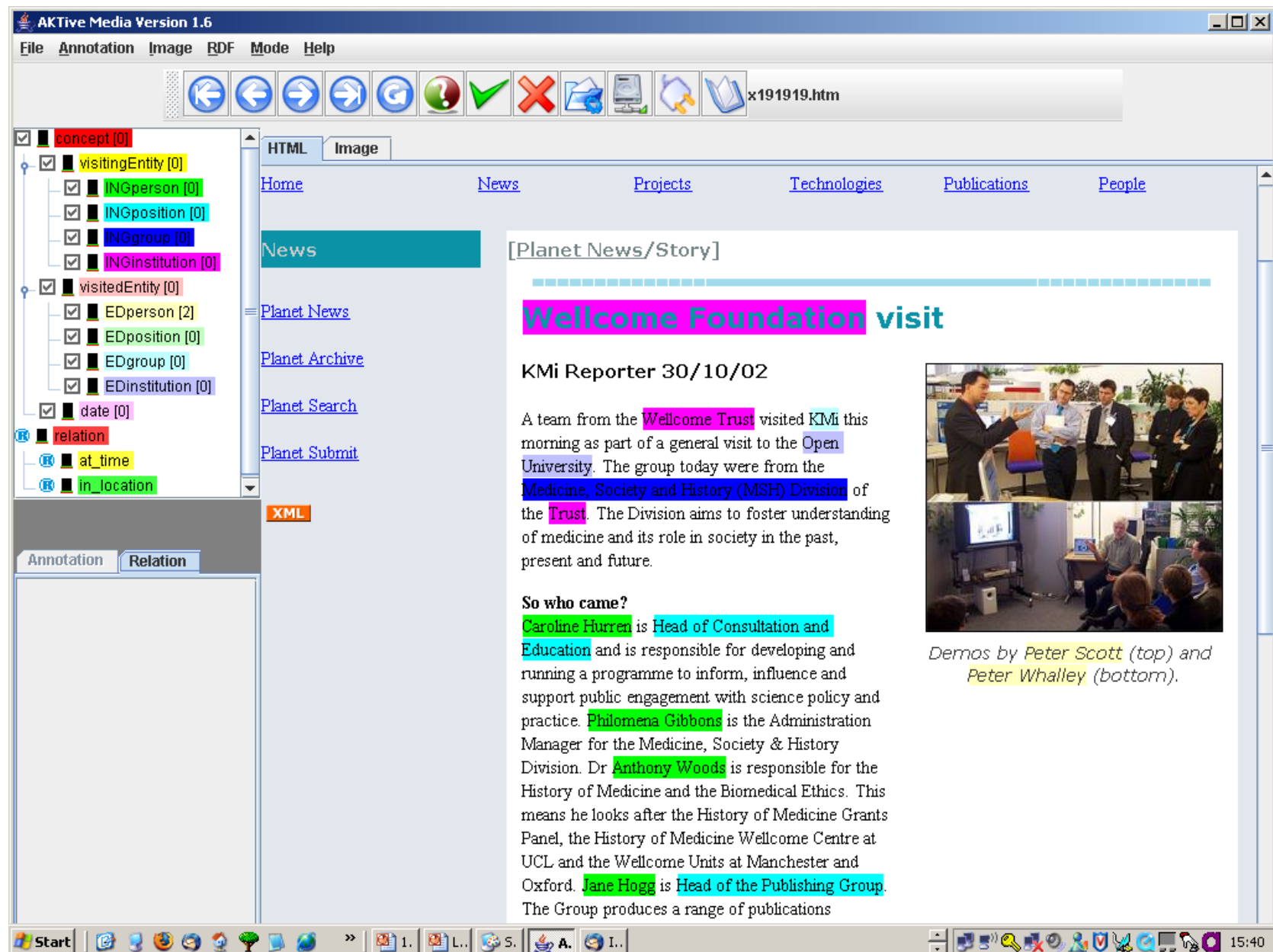
Caroline Hurren is Head of Consultation and Education and is responsible for developing and running a programme to inform, influence and support public engagement with science policy and practice. Philomena Gibbons is the Administration Manager for the Medicine, Society & History Division. Dr Anthony 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 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).

© Fabio Oravegnin, University of Sheffield

15:40

Contextual Annotation of Images and Text



AKTive Media Version 1.6

File Annotation Image RDF Mode Help

Navigation icons: back, forward, search, etc. x191919.htm

Annotations:

- 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**

Home News Projects Technologies Publications People

News

[Planet News/Story]

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 UCL and the Wellcome Units at Manchester and Oxford. **Jane Hogg** is **Head of the Publishing Group**. The Group produces a range of publications

Images:

- Top: A group of people in a meeting room.
- Bottom: A man presenting to a group of people.

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

Annotation | **Relation**

Start | 1. | L.. | S. | A. | I.. | 15:40

Contextual Annotation of Images and Text

AKTive Media Version 1.6

File Annotation Image RDF Mode Help

Navigation icons: back, forward, search, etc.

x191919.htm

HTML Image

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

Annotation Relation

Enter Annotation Text

Martin Dzbór

Search

Martin Dzbór

Martin Dzbór

Simon Buckingham Shum

Objects Technologies Publications People

Story]

visit

this

of

division aims to foster understanding

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

and the Biomedical Ethics. This

er the History of Medicine Grants

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

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

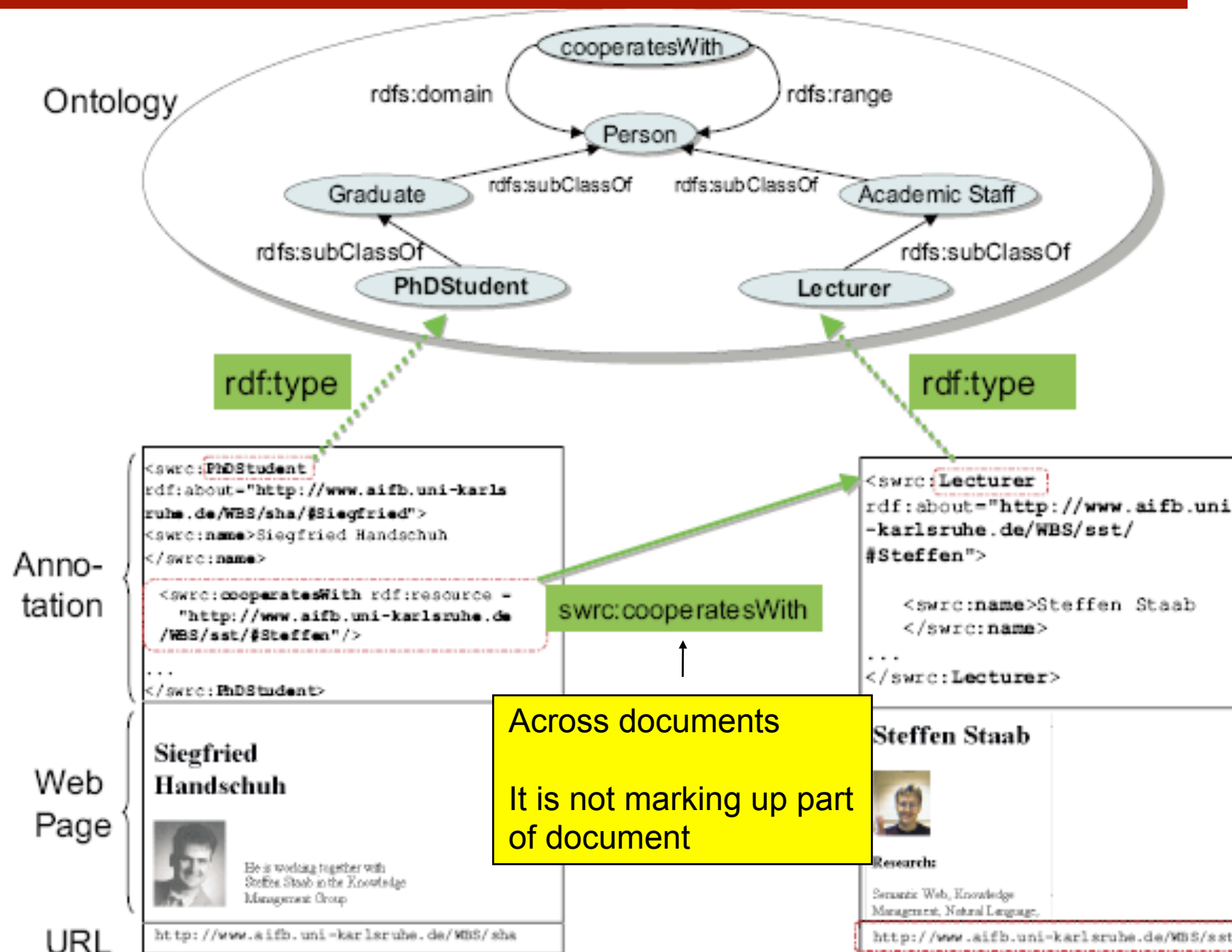
Start

1. L. S. A. I.

15:40

Annotating across documents (CREAM, 2001)

65





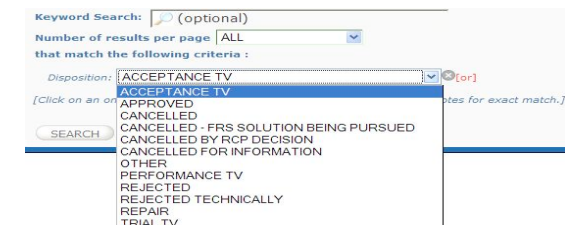
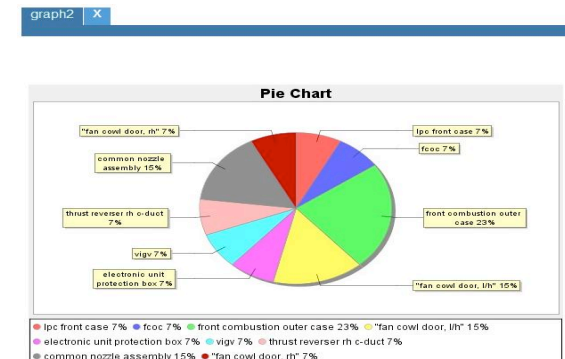
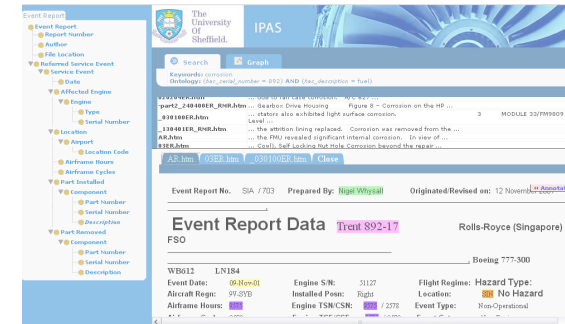
K-Tools

A real world application of annotation
to knowledge management

- Enable multi-media capture and sharing of technical knowledge
 - support distributed networked knowledge acquisition, capturing, retrieval and sharing.
 - enable communities of users to define their own views,
 - while at the same time maintaining the connection with other units' views
 - and with a central schema.
 - They create a Semantic Web of knowledge
 - not a comprehensive integrated knowledge view



- Spin-out company (2008)
 - Solutions for knowledge acquisition and sharing
 - User-centred Knowledge Capture
 - K-Forms
 - Knowledge Capture from Legacy Data
 - K-IE
 - Knowledge Search
 - K-Search
 - Patented Technology
 - Finalist of Rolls-Royce Creativity Award 2007



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 connected to company's ontology or other sources of knowledge
 - e.g. other forms
 - enables sharing and integration with other knowledge sources
 - knowledge is immediately available for search, sharing and reuse
 - based on k-search
 - no additional effort required by user wrt present methodology

Create a New Form Template

Back to options
Add New Section
Save

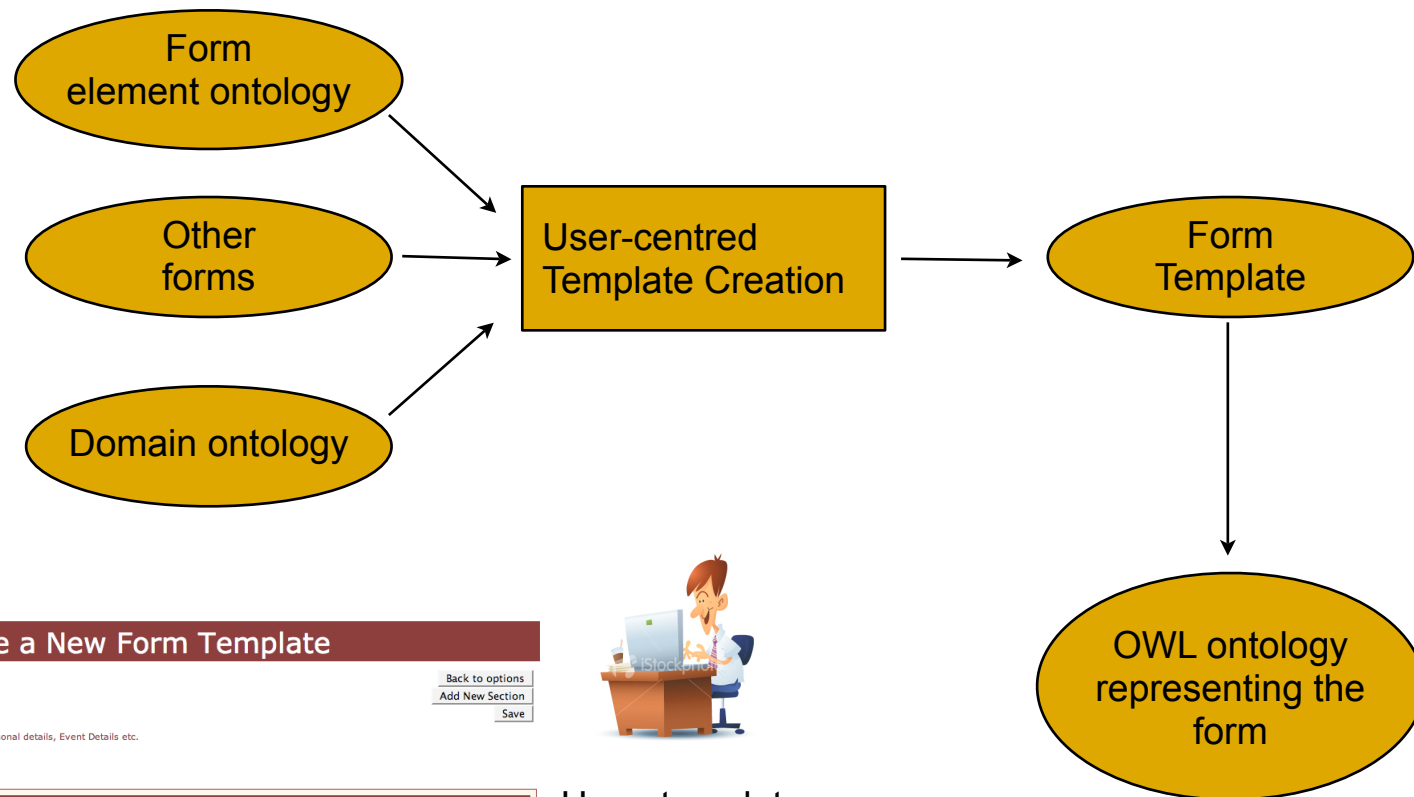
A Form Template can be composed by many sections, e.g. Personal details, Event Details etc.
Click on the button on the right to create a new section

The screenshot shows a web interface for creating a form template. It features a section titled 'Section1' with a light yellow background. Inside this section, there are several fields for configuration: 'Field Type' (set to 'Textfield'), 'Question/Title' (set to 'Engine name'), 'Concept Name' (set to 'engine_name'), 'Validation' (set to 'Mandatory'), 'Help Text(Optional)' (set to 'what is the engine name (e'), and 'Data Type' (set to 'String'). There are 'Add' and 'Delete' buttons at the bottom of the section. Below the main section, there is another field for 'Field Type' set to 'Textarea'.

- 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.
 - Connect form fields to existing ontology or database schema or other forms
- 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



know
RDSX

Create a New Form Template

Back to options
Add New Section
Save

A Form Template can be composed by many sections, e.g., Personal details, Event Details etc.
Click on the button on the right to create a new section

Section1	
Field Type:	Textfield
Question/Title:	Engine name
Concept Name:	engine_name
Validation:	Mandatory
Help Text(Optional):	what is the engine name (e
Data Type:	String
<input type="button" value="Add"/> <input type="button" value="Delete"/>	
Field Type:	Textarea

User: template designer

Forms: Features of template creation

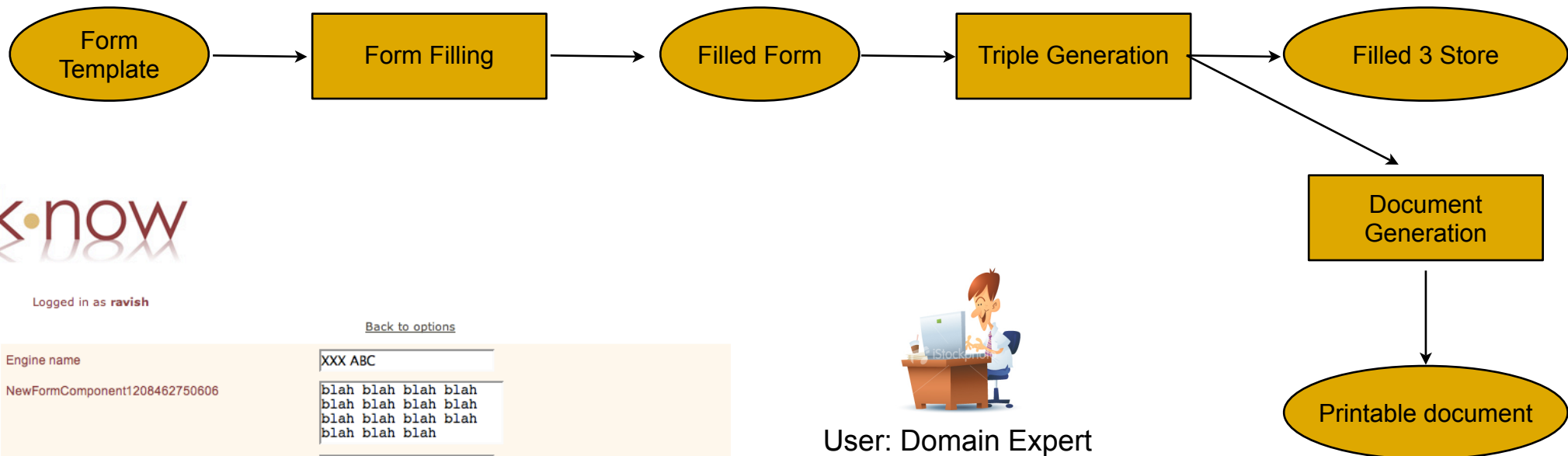
- Easy creation and release of form templates over Intranet and in local
- Easy reuse of other form components
 - e.g. all documents about jet engines will have
 - engine type, model, owner, no of cycles, etc.
- Based on:
 - Use of ontology of potential form elements
 - Templates composed using ontology 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



k.now

Logged in as **ravish**

[Back to options](#)

Engine name	XXX ABC
NewFormComponent1208462750606	blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah
Engine age	1224 cycles
Brithday description	blah blah
picture	<input type="text"/> Browse...
I am fond of it	<input checked="" type="checkbox"/> yes <input type="checkbox"/> not sure <input type="checkbox"/> no
When did you met?	work

Finish



Forms: Features of Knowledge Capturing

74

- 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



- Ontologies used for:
 - modelling forms:
 - forms components are concepts in an ontology
 - rules control combination of elements
 - why Sem Web technology?
 - All system is declarative
 - concepts
 - rules
 - Easy modification of behaviour
 - Sophisticated behaviour
 - e.g. if the value of field X is Y, then present only option Z



- Ontologies used for:
 - Each form is automatically turned into an ontology
 - Importing parts of an existing form
 - implies mapping between two ontologies
 - Advantage: the information is connected
 - it is possible to share across the archives
 - this creates a semantic web of ontologies
 - all the SW technologies used for managing distributed ontologies apply
 - e.g. distributed searching
 - mapping between existing forms can also be done by system administrator



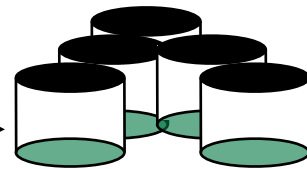
- Filled forms become RDF triples
 - Different types of documents can be generated from the triples
 - to enable user tailoring
 - to adopt a specific terminology
 - different departments use different terminology
 - for confidentiality



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.

Ontology:
 information_bearing
 attending_confer
 generic_agat
 fund_institute
 charitable_organiza
 multimedia_designe
 attending_invent
 organization_unit
 employee
 conferring_on_aver
 social_publication
 learning_centered_ei
 educational_organiz
 event_involving_nc
 book
 operating_system
 higher_educational
 thesis_reference
 event_involving_pri
 city
 article_reference
 industrial_organizat

Name Base

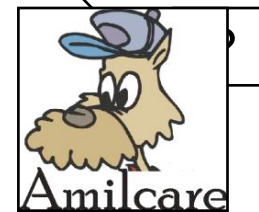


Near Match in Index
 Archive

T. Rex



Disambiguation
 In documents



Automating Annotation

- Solutions like k-forms are very good for annotating new knowledge
 - large repositories of legacy data exist
 - it is important that new management solutions are able to reuse existing data
 - do not require a completely new world to be built for you!!
- Legacy data is generally represented in
 - databases
 - textual documents
 - images
 - ...



■ 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.

Named Entities

Event Recognition

Growing complexity

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 team
as president.

Named Entities

Event Recognition

Growing complexity

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.

Named Entities

Event Recognition

Growing complexity

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.

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

■ 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



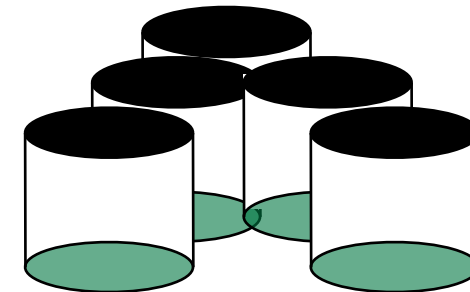
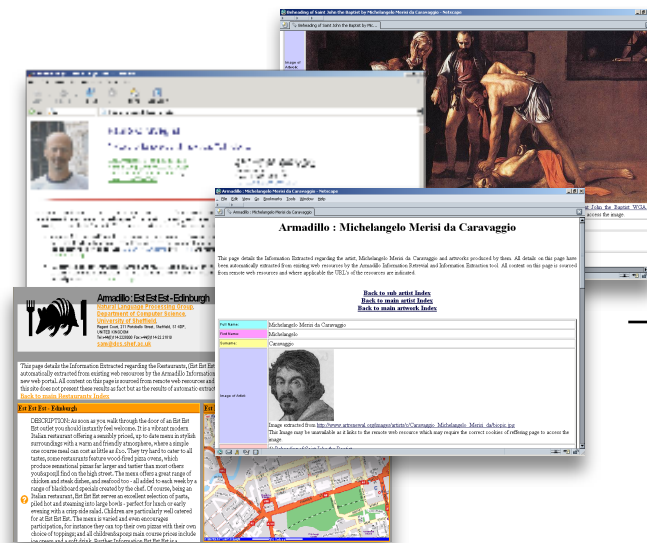
- For large scale (some hundred millions pages) smarter infrastructure is needed
 - Search engine-like indexing infrastructure
 - Faster processing (less processing)
 - Two cases:
 - Recognition of known terms (and their variations)
 - See also information integration
 - Discovery of new names



Large Scale NER: Indexing

85

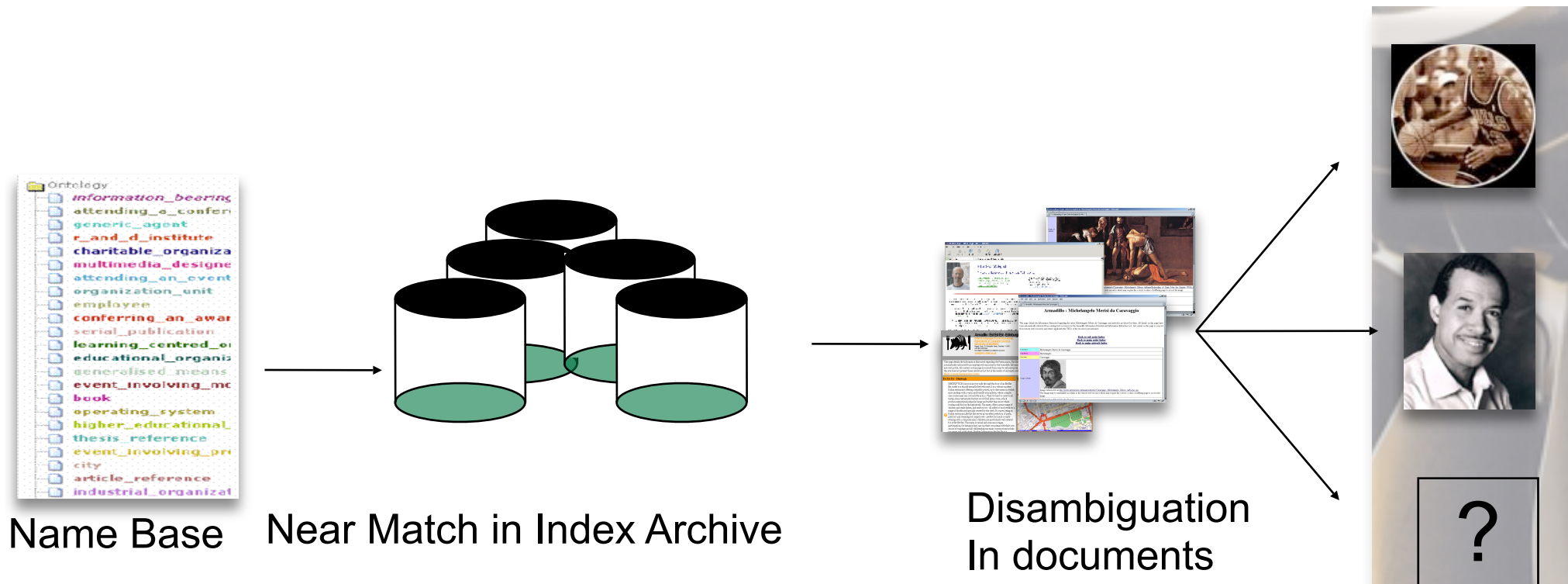
■ Document Indexing as in Search Engines



Distributed Index Archive
(keywords)

Known Name Recognition

86

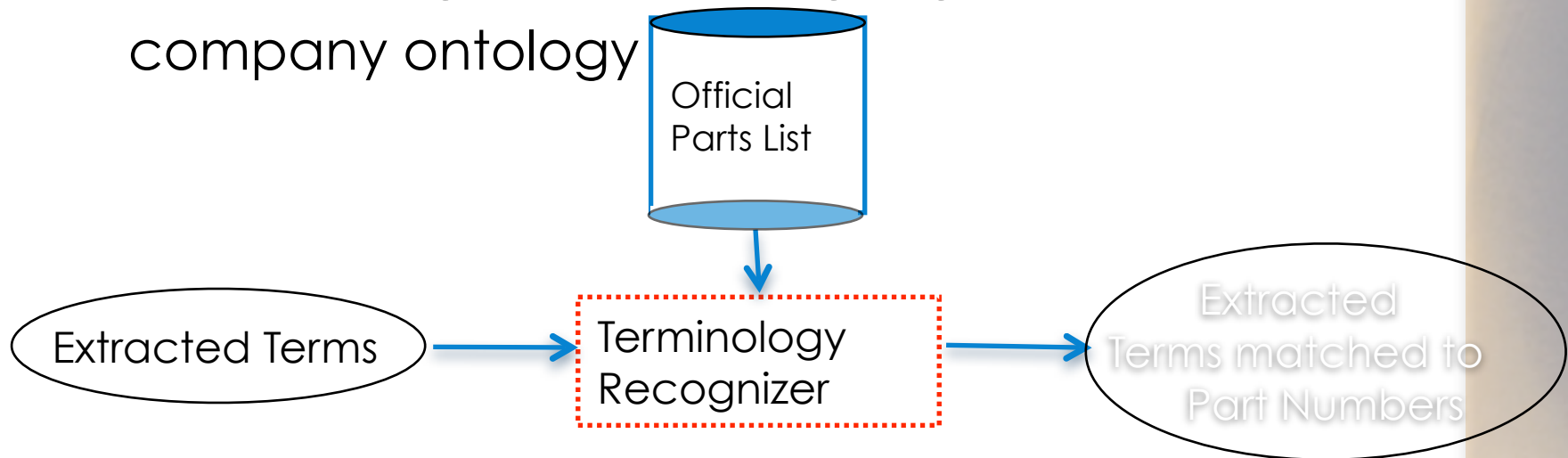


S. Dill, N. Eiron, et al: SemTag and Seeker: Bootstrapping the semantic web via automated semantic annotation. WWW'03

- 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

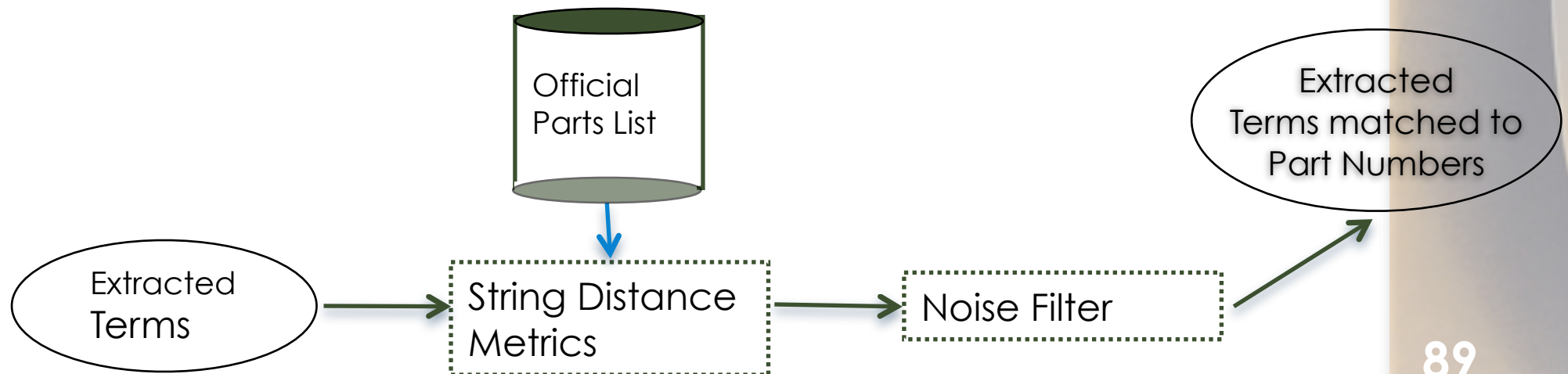


- NER is one example of term recognition
- More useful in technical domains is terminology recognition
 - The task of assigning a URI to a technical description
 - i.e. mapping a natural language description to the official company ontology



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)



- 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



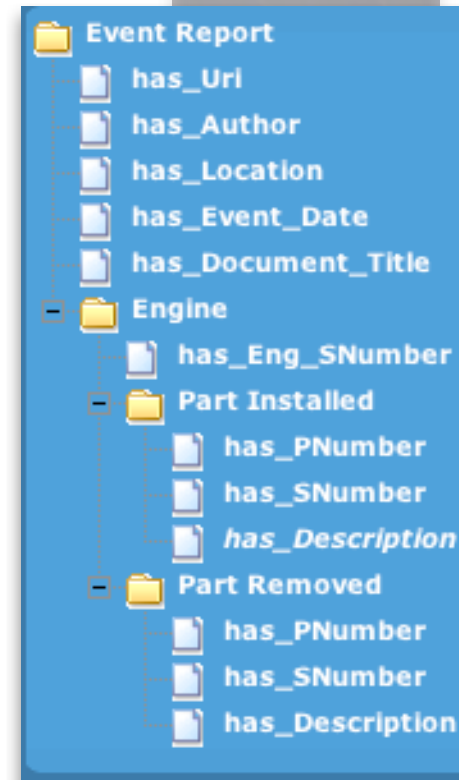
- Tables are an essential part of many documents
 - Most information is represented in tables
- Tables can be represented as forms to fill
 - Semantics is fixed
 - Wrapper writing or wrapper induction (Kushmerick 1997)
- Tables can be created ad hoc in documents (e.g. Word docs)
 - Semantics is unclear
 - 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

92

- 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



Types of tables in Event Reports

93

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

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

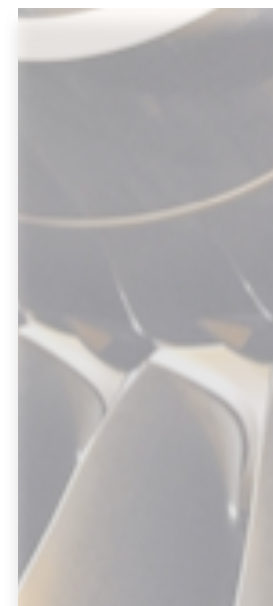
<u>s/n removed</u>	04-0721257 <u>tsn/csn: 268/106</u>
<u>s/n installed</u>	04-1012229 tsn/csn:0/0

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

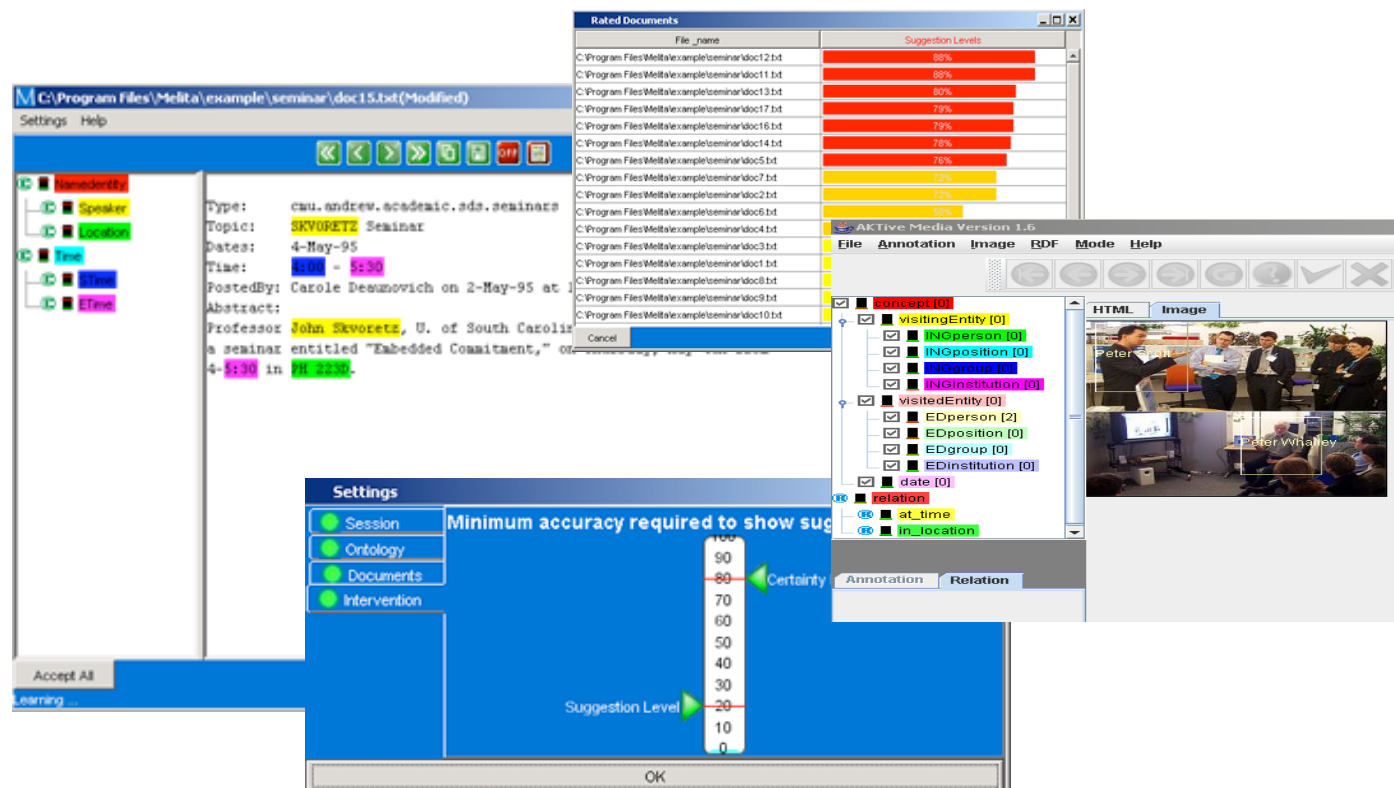
Applying information extraction

94

- AktiveMedia to annotate texts
- TRex system (Jirí et al. 2006) to train and extract
 - <http://tyne.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)



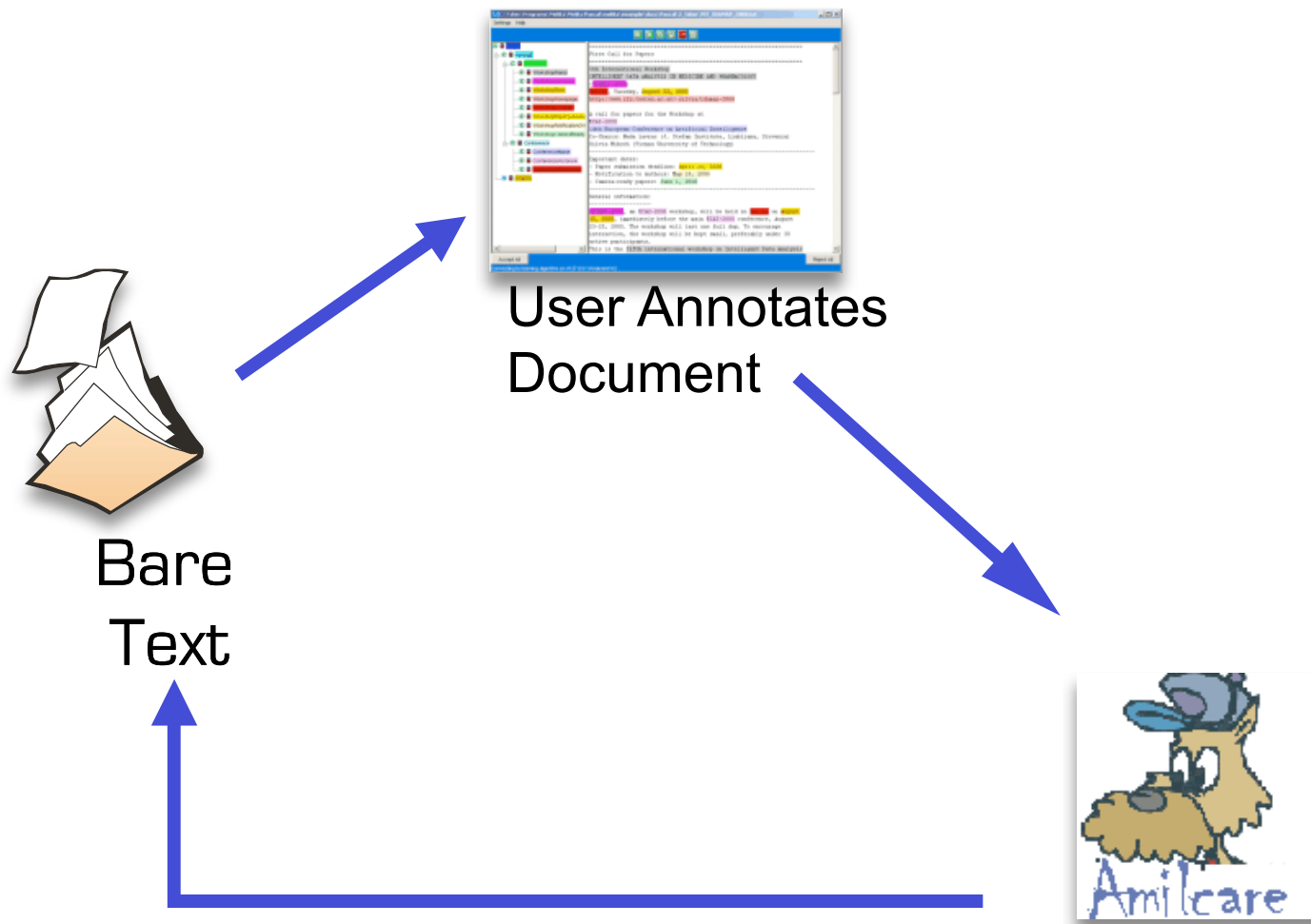
	POS	ACT	CORR	WRONG	MISSED	PREC	REC	F1
airport	120	120	120	0	0	100	100	100
has_airframe_cycles	104	104	104	0	0	100	100	100
has_airframe_hours	104	104	104	0	0	100	100	100
has_author	120	120	120	0	0	100	100	100
has_engine_serial_number	120	120	120	0	0	100	100	100
has_engine_type	120	120	120	0	0	100	100	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
TOTAL	1644	1658	1625	33	19	98	99	98



Using IE to Support Manual Annotation

Using IE to support annotation: step 1

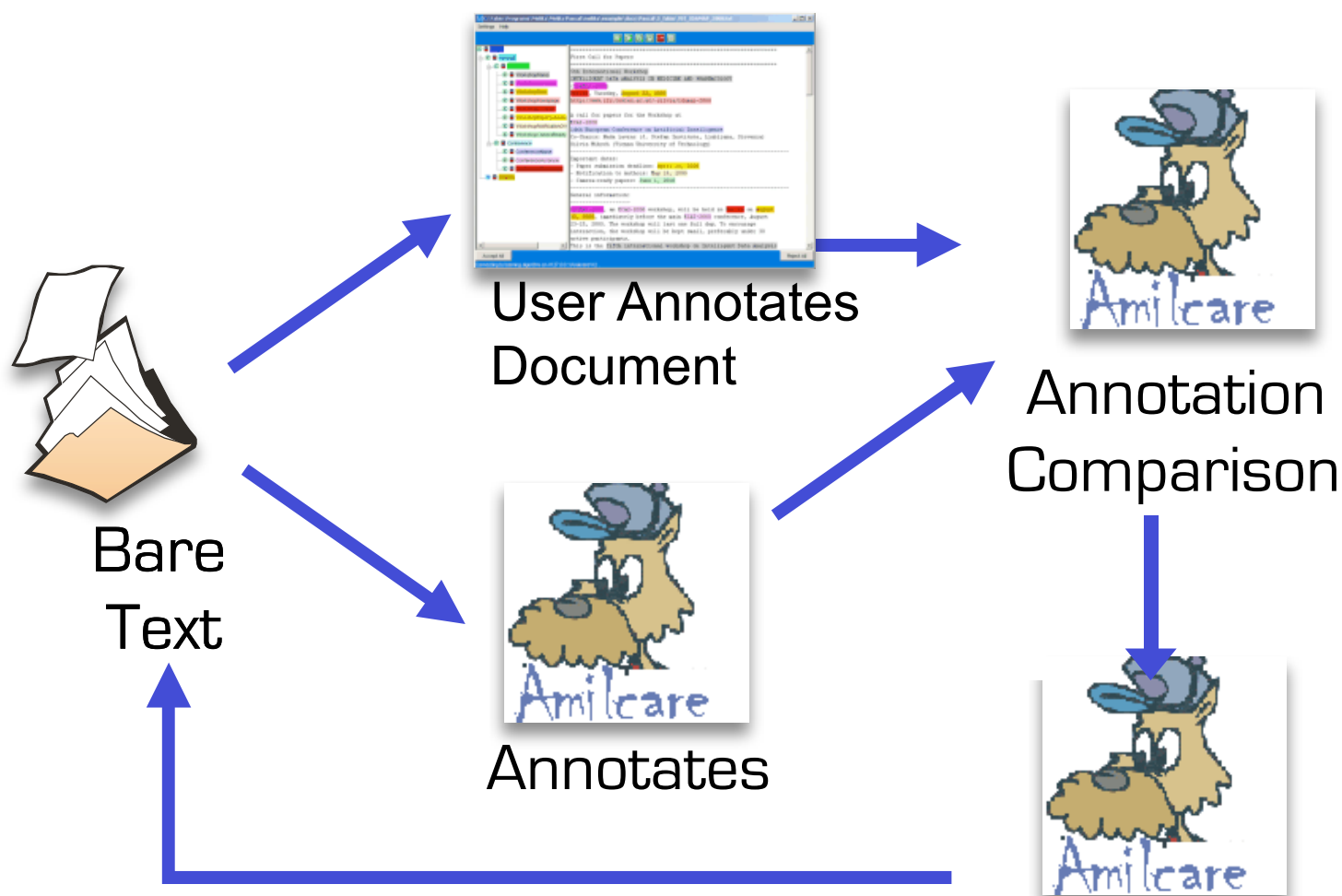
96



Trains on annotated corpus

Using IE to support annotation: step 1

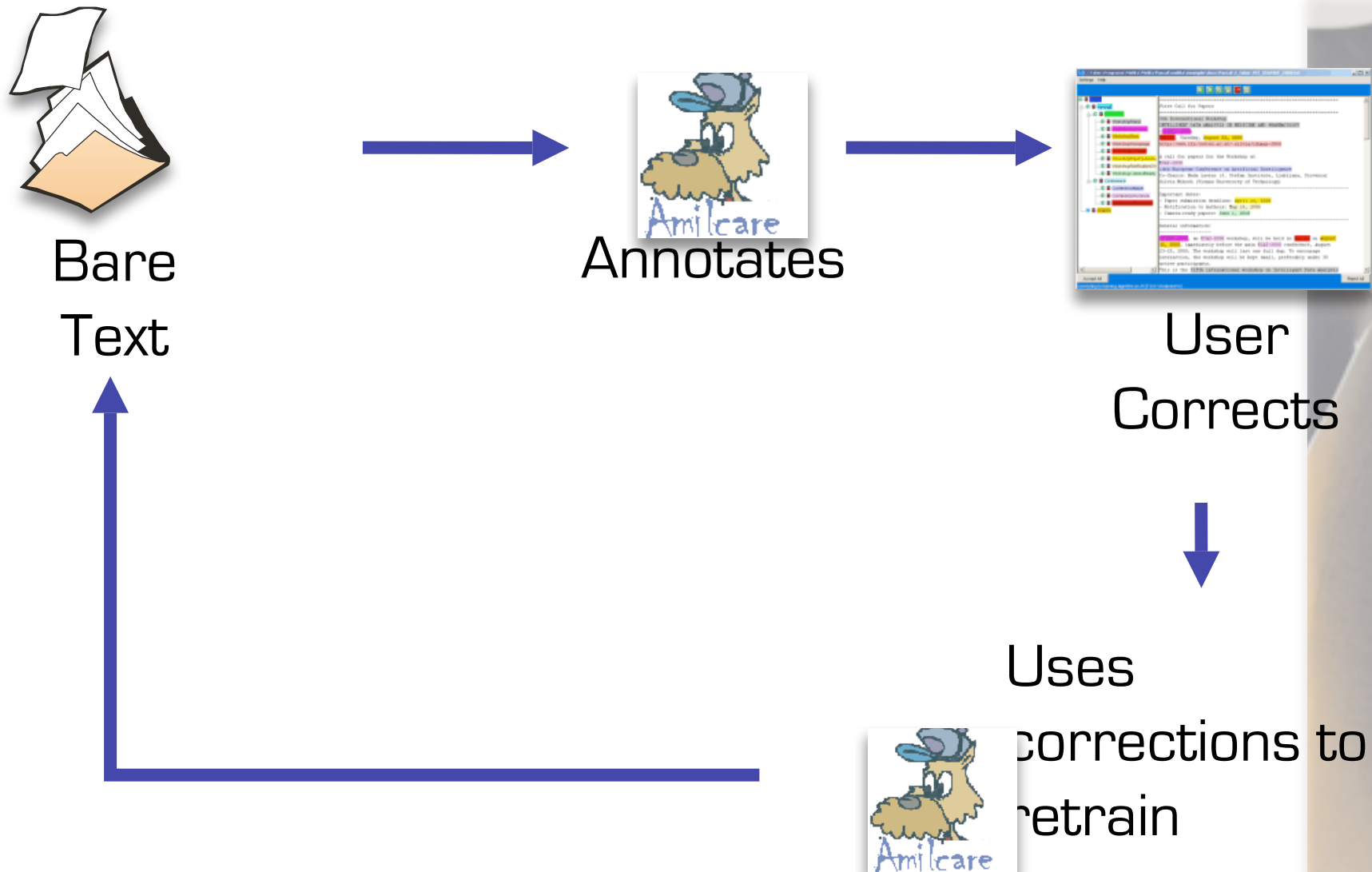
96



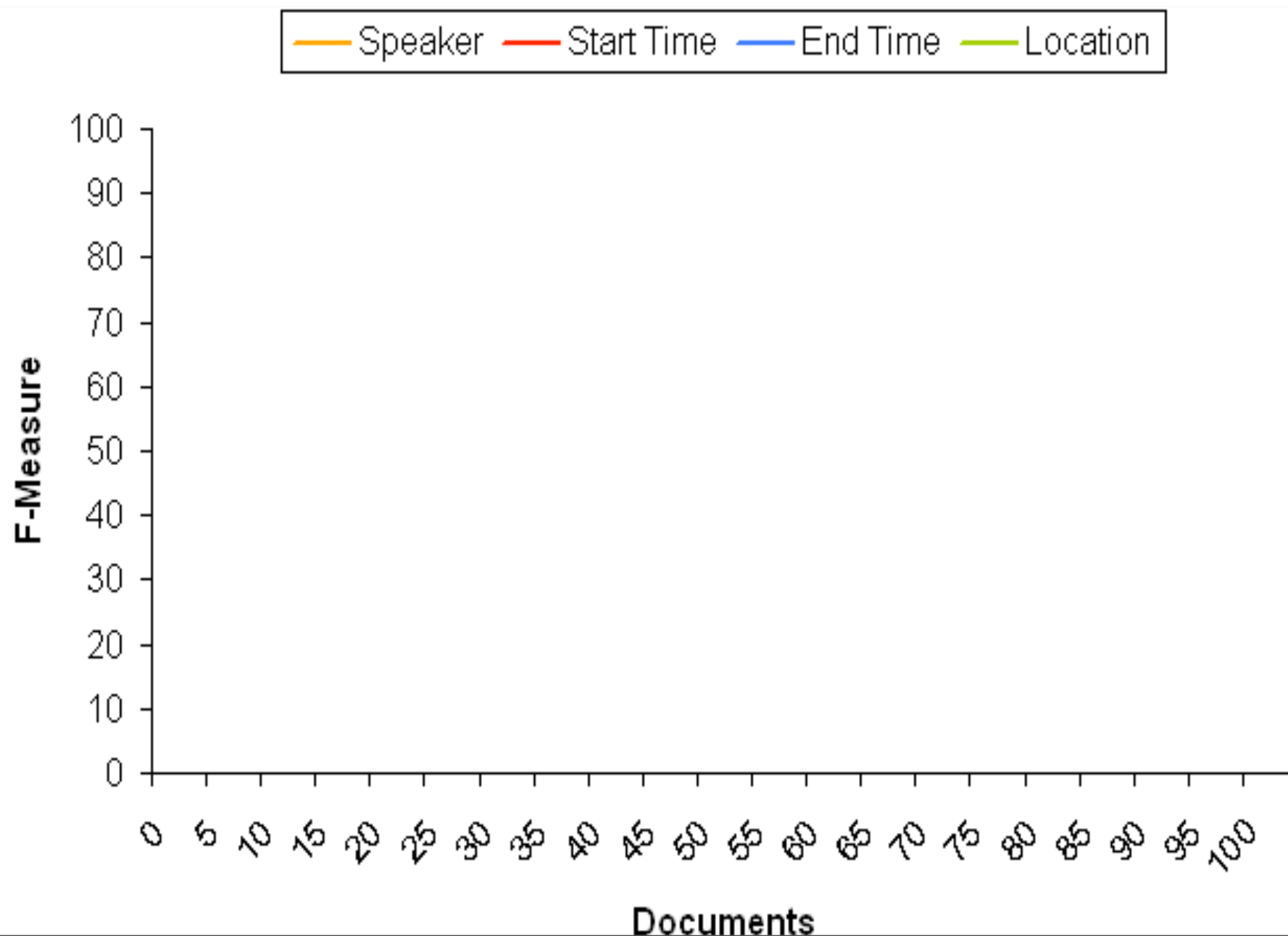
Retrain using errors,
missing tags and mistakes

Using IE to support annotation: step 2

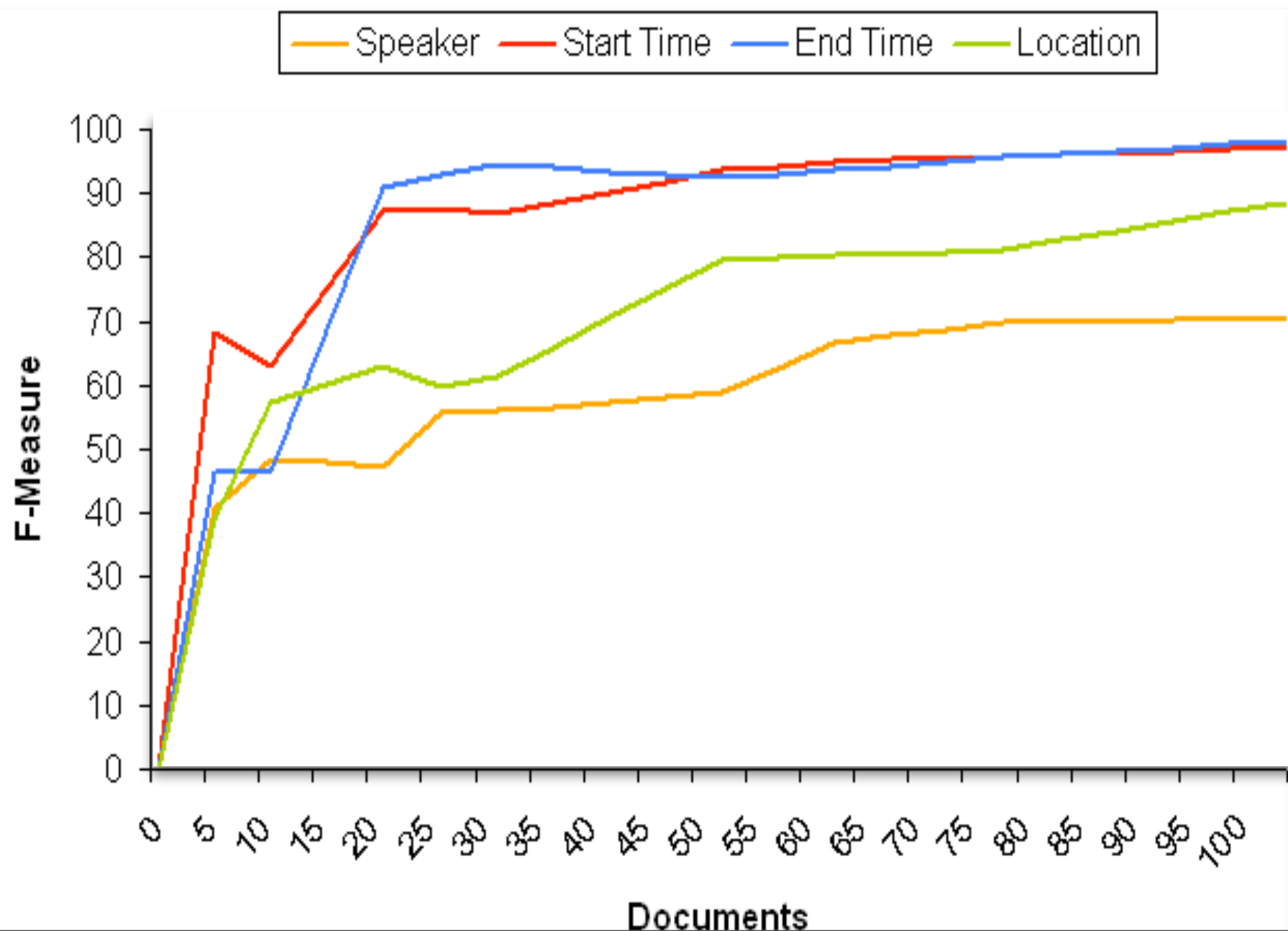
97



Learning curve



Learning curve

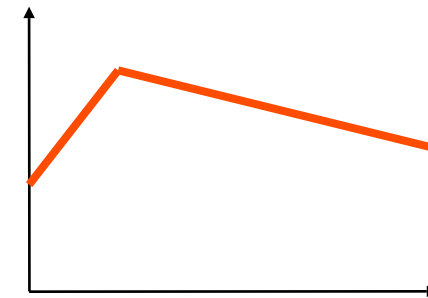


Impact on Annotation

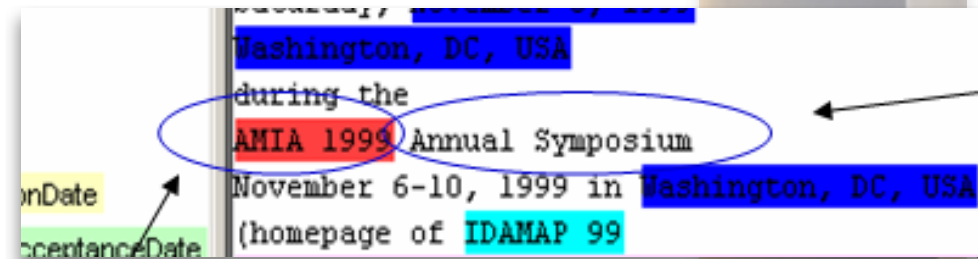
99

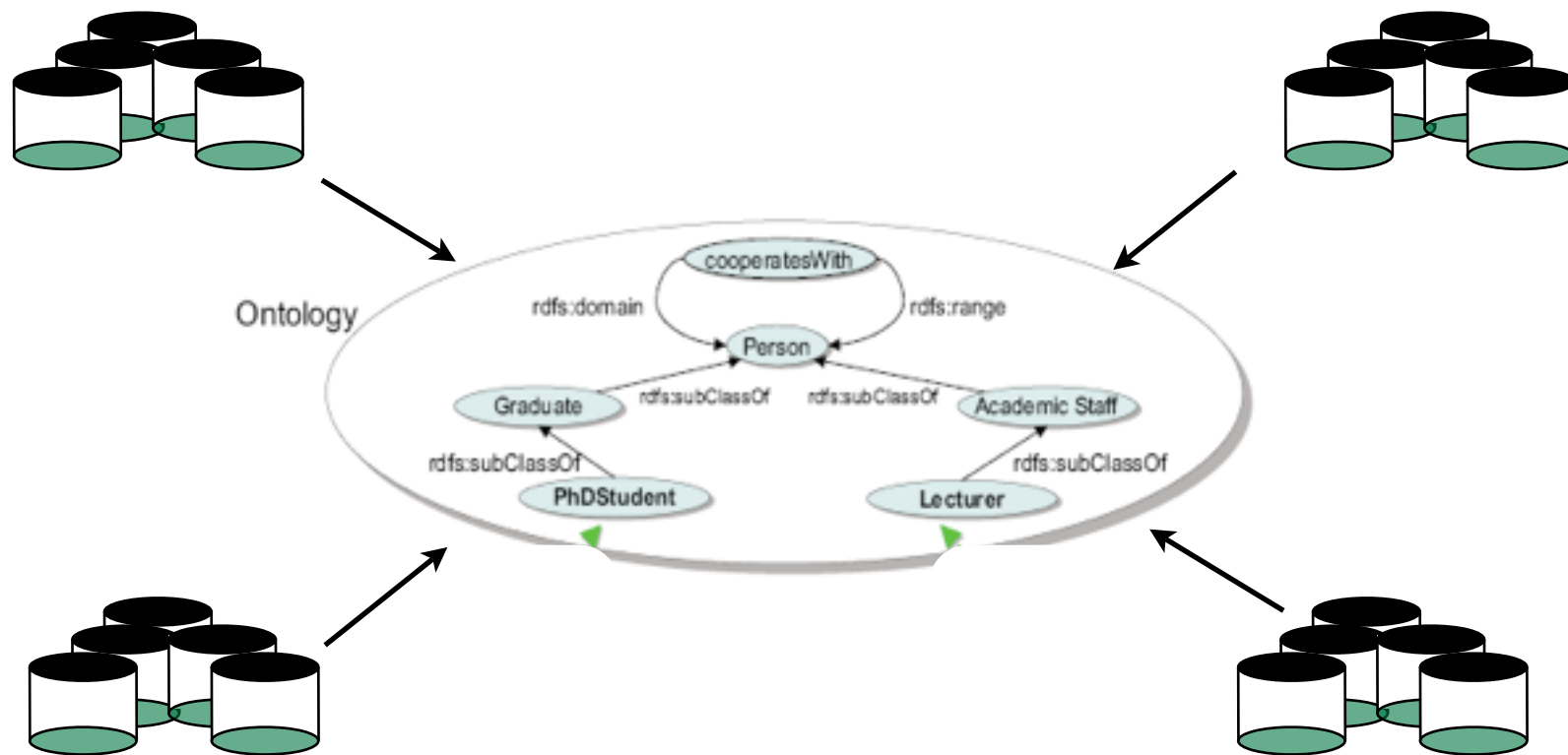
- University of Karlsruhe experiments
 - -80% annotation time
 - +100 interannotator agreement
 - Is this positive?
- Outstanding issue:
 - Impact on annotators of suggestions topping 85% accuracy?
 - Annotation needs to be precise and consistent
 - Otherwise the IE system is confused
 - Can only annotate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

IE accuracy



Amount of annotations



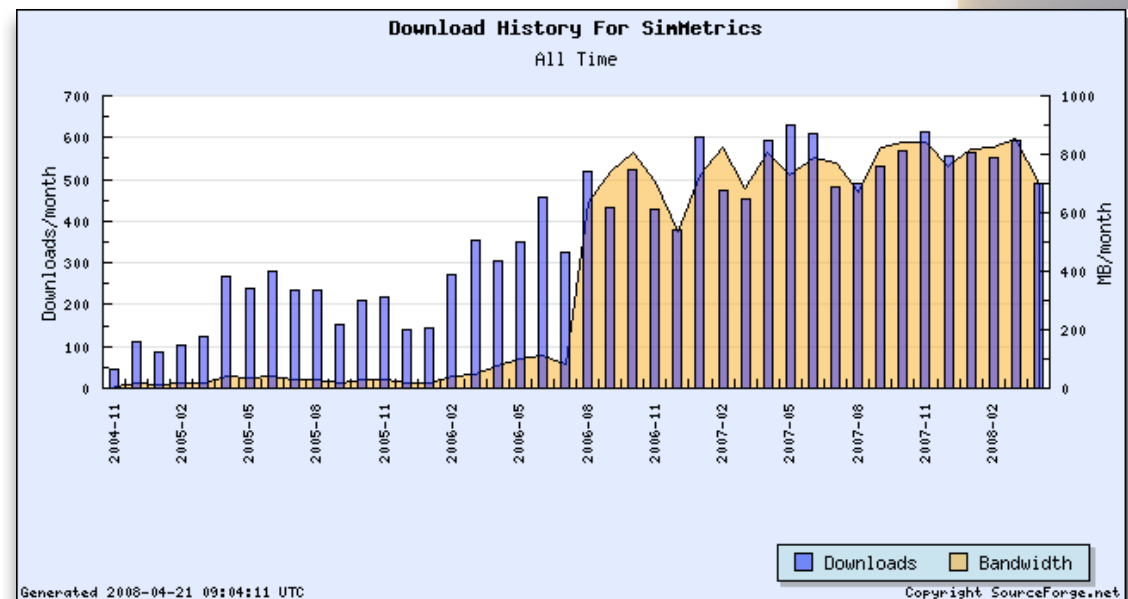


Information Integration

- 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

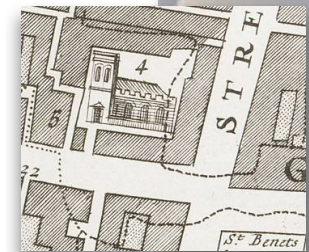
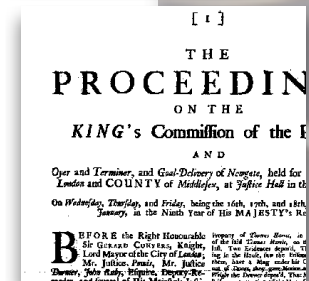


- Library of distance metrics released as open source
 - <http://sourceforge.net/projects/simmetrics/>
 - >15,000 downloads since end of 2004
 - Most downloaded distance metrics library on the Web
 - for strings and records
 - Hundreds of applications
 - Developed by Sam Chapman, University of Sheffield



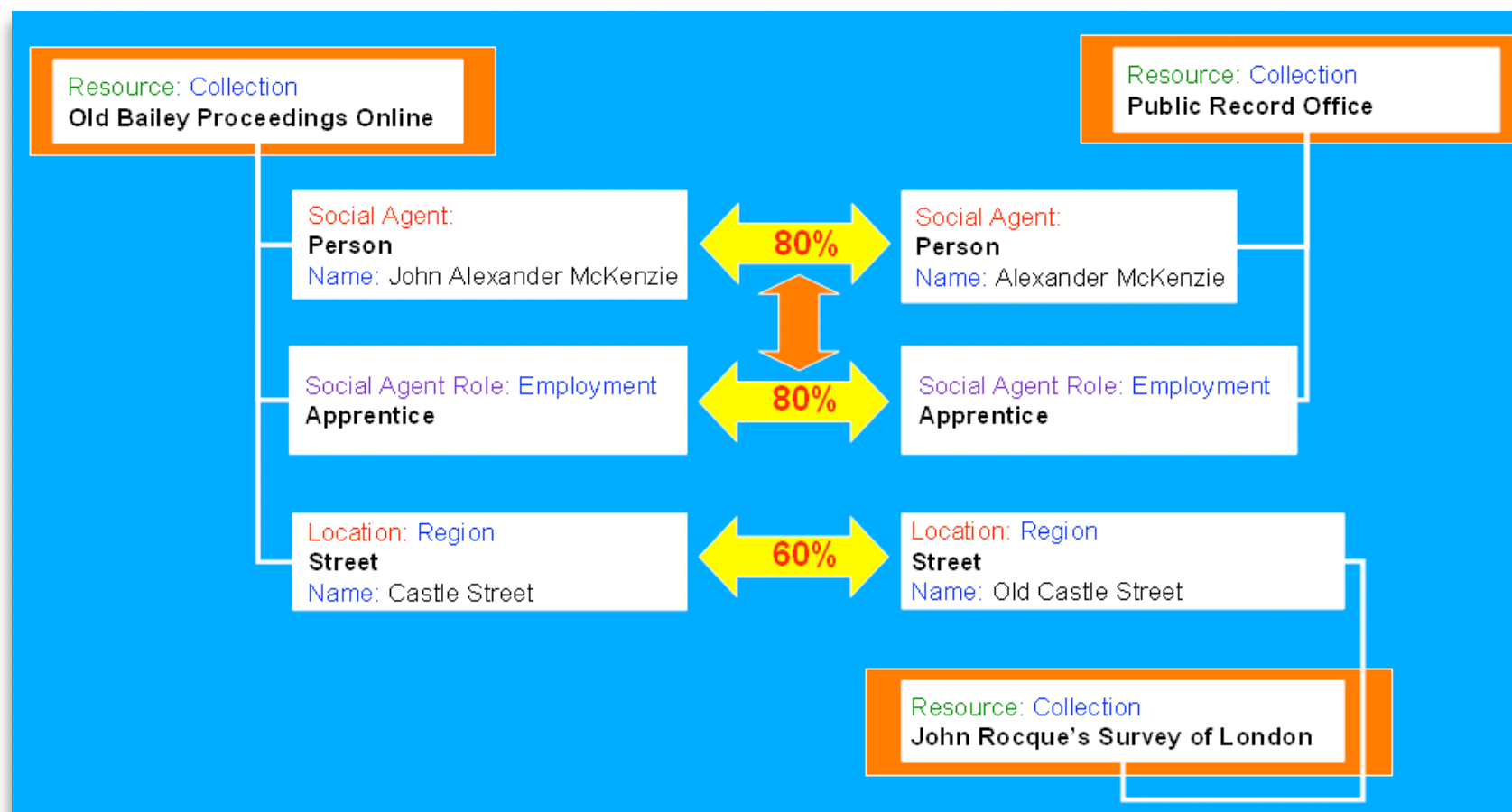
The diagram illustrates various AHDS (Archives and Heritage Deposits) categorized by color-coded boxes:

- Pink Box (Top):**
 - 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** (Red text)
- Grey Box (Middle Left):**
 - St. Martin's Settlement Exams Index
 - WESTCAT** (Green text)
- Grey Box (Middle Right):**
 - Collage image database
 - Guildhall Library** (Green text)
- Blue Box (Bottom Right):**
 - Harben's Dictionary of London
 - John Strype's "Survey..."
 - <http://www.motco.com> (Blue text)
- Grey Box (Bottom Left):**
 - Metropolitan London in the 1690s
 - IHR** (Green text)
- Grey Box (Bottom Middle):**
 - Selected Criminal Records
 - PRO** (Green text)
- Grey Box (Bottom Far Left):**
 - House of Lords Journals
 - BOPCRIS** (Green text)



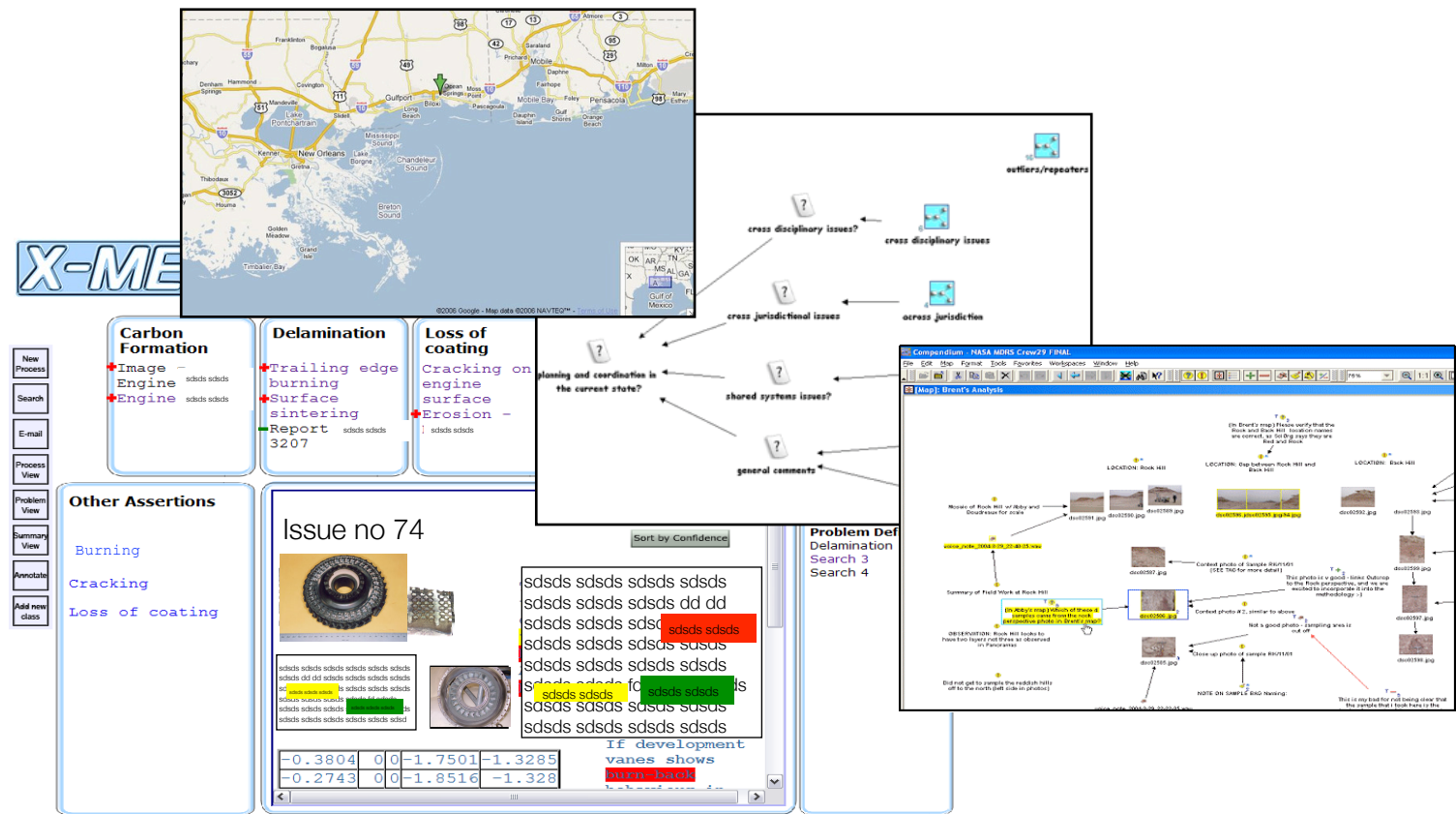
<http://www.hrionline.ac.uk/armadillo/>

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





Knowledge Sharing and Reuse

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

- 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



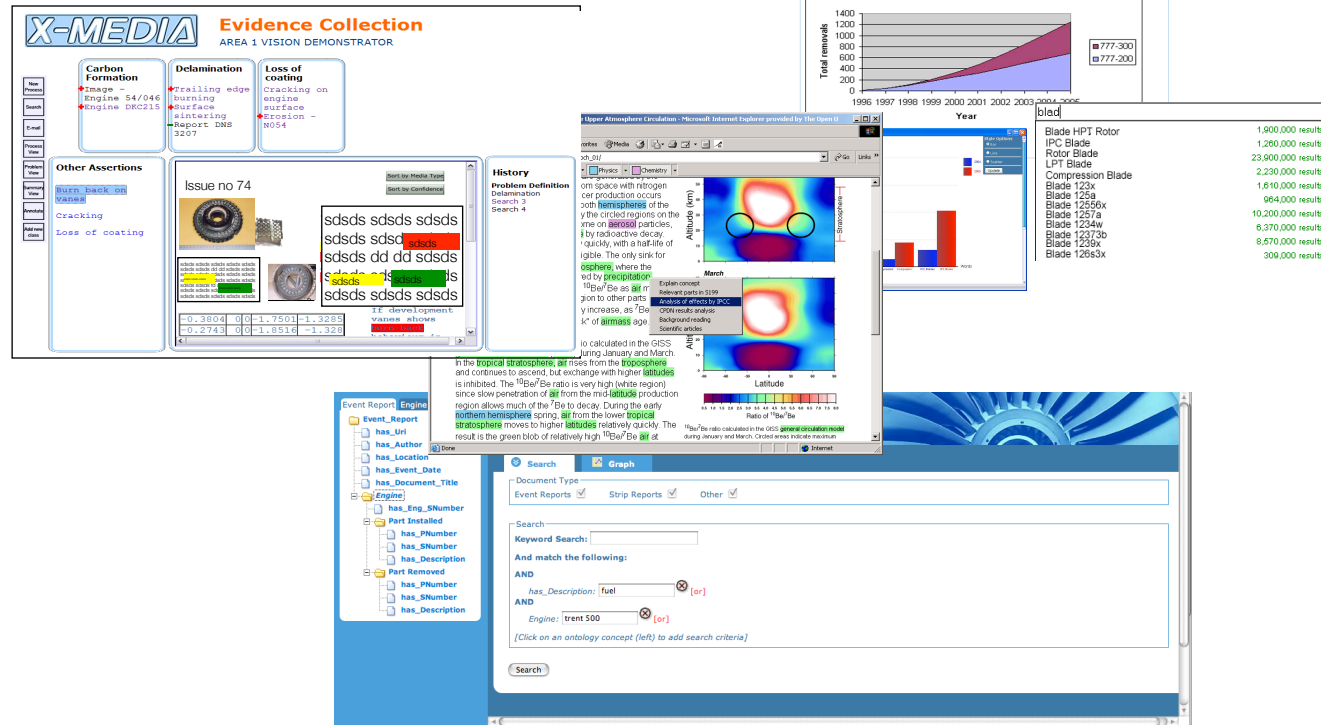
- Large distributed archives require ability to
 - Map distribution of information,
 - Weight every single source
 - Distribute searches carefully;
- Currently: search is performed just in some of the archives, disregarding others that can bring very useful information
- Importance of context and provenance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



- 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

- 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

- 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



- Many types of technologies
 - Search based on structural query languages, such as SPARQL, see, e.g., ARQ, and
 - User-centred search to retrieve ontologies (e.g. Swoogle [Ding et al. 2004] and Watson [d'Aquin et al. 2007])
 - User-centred approaches to retrieve information and knowledge
- We will see the latter



- 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)

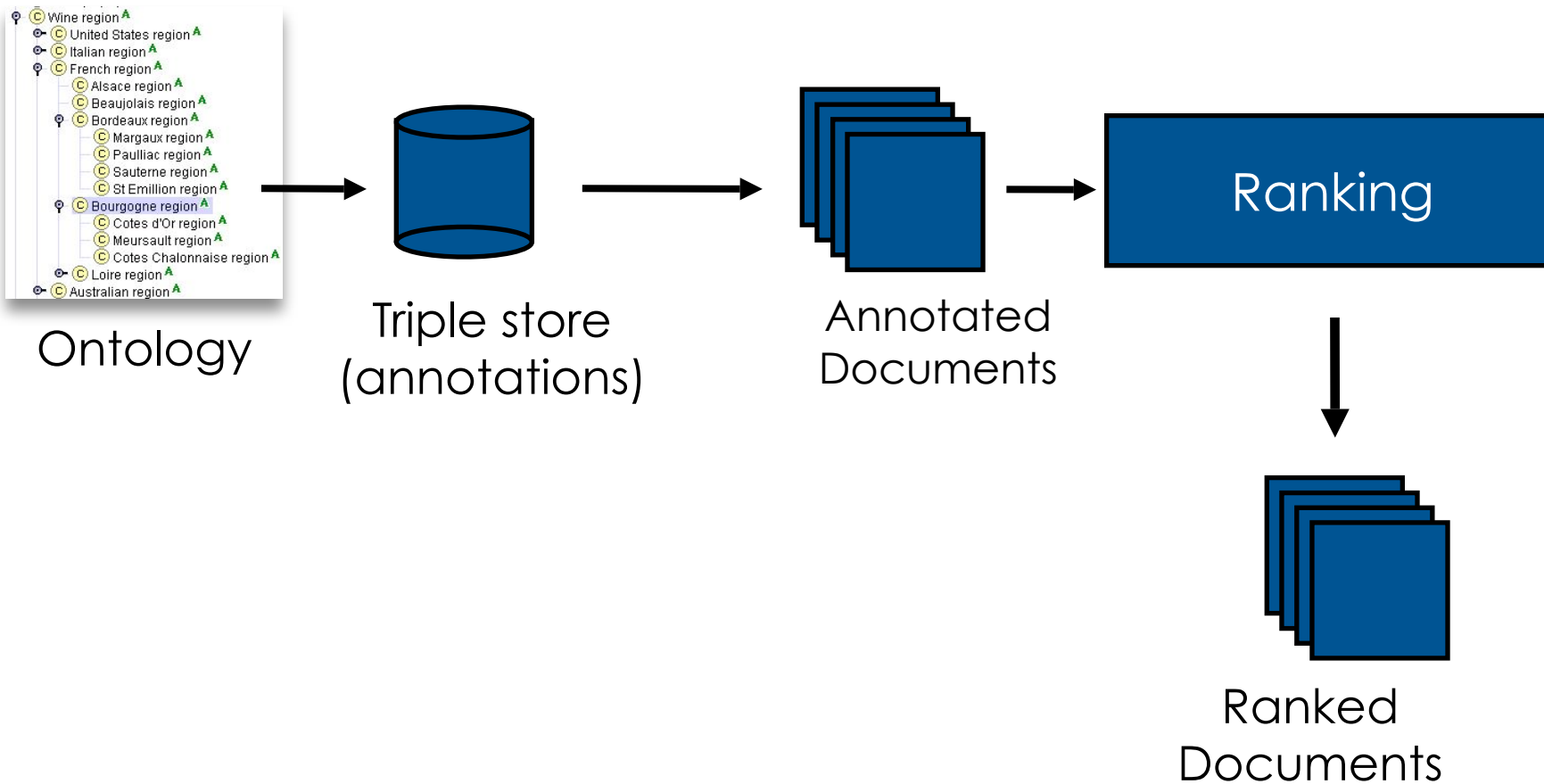


- By merging the definitions in [Uren et al. 2008], [Kaufmann et al. 2007b] and [Baghdev et al. 2008]:
 - Keyword-based approaches considering a natural language query as a bag of words
 - [Kaufmann et al. 2007a] [Lei et al., 2006])
 - Natural language approaches: modelling the linguistics of the query
 - [Lopez et al. 2005],[Bernstein et al. 2005b], [Kaufmann et al. 2006]
 - Graph-based approaches
 - [Bernstein et al. 2005a], SEWASIE, Falcon-S.
 - Form-based approaches (e.g. Corese)
 - Hybrid approaches
 - K-Search [Baghdev et al. 2008])



Classic Ontology-based Querying

116



- Building an ontology is expensive and needs maintenance
- Metadata generation of documents is
 - Expensive and Error prone
- Solutions like k-forms help simplifying knowledge capture and acquisition
- 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



- 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



- 85% of documents in the first 20 hits are relevant
 - Compare with keywords: 56%
- 40% of relevant documents are in the first 2 pages
 - Compare with keywords: 57%
- Ontology matching implies
 - Reading a limited amount of irrelevant documents
 - Risking missing many documents
 - It is possible to count the events



- 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



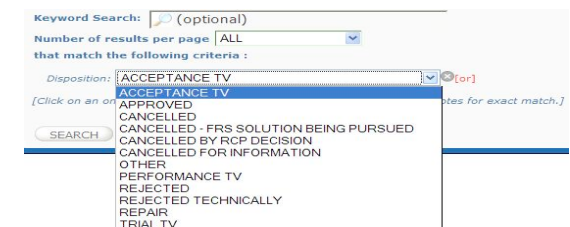
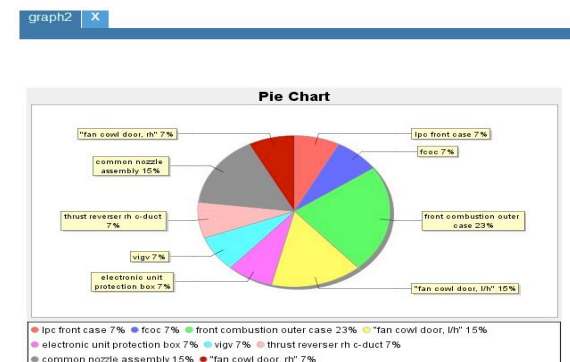
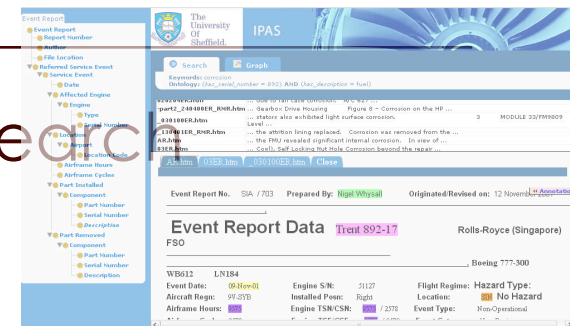
- Mixes keyword and ontology based search
 - Ontology based search
 - Traditional keyword search
 - Keyword in contest of ontology-based annotations
- Potential queries:
 - Return all documents where the word fuel is mentioned
 - Return all documents where the affected part description includes the word 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

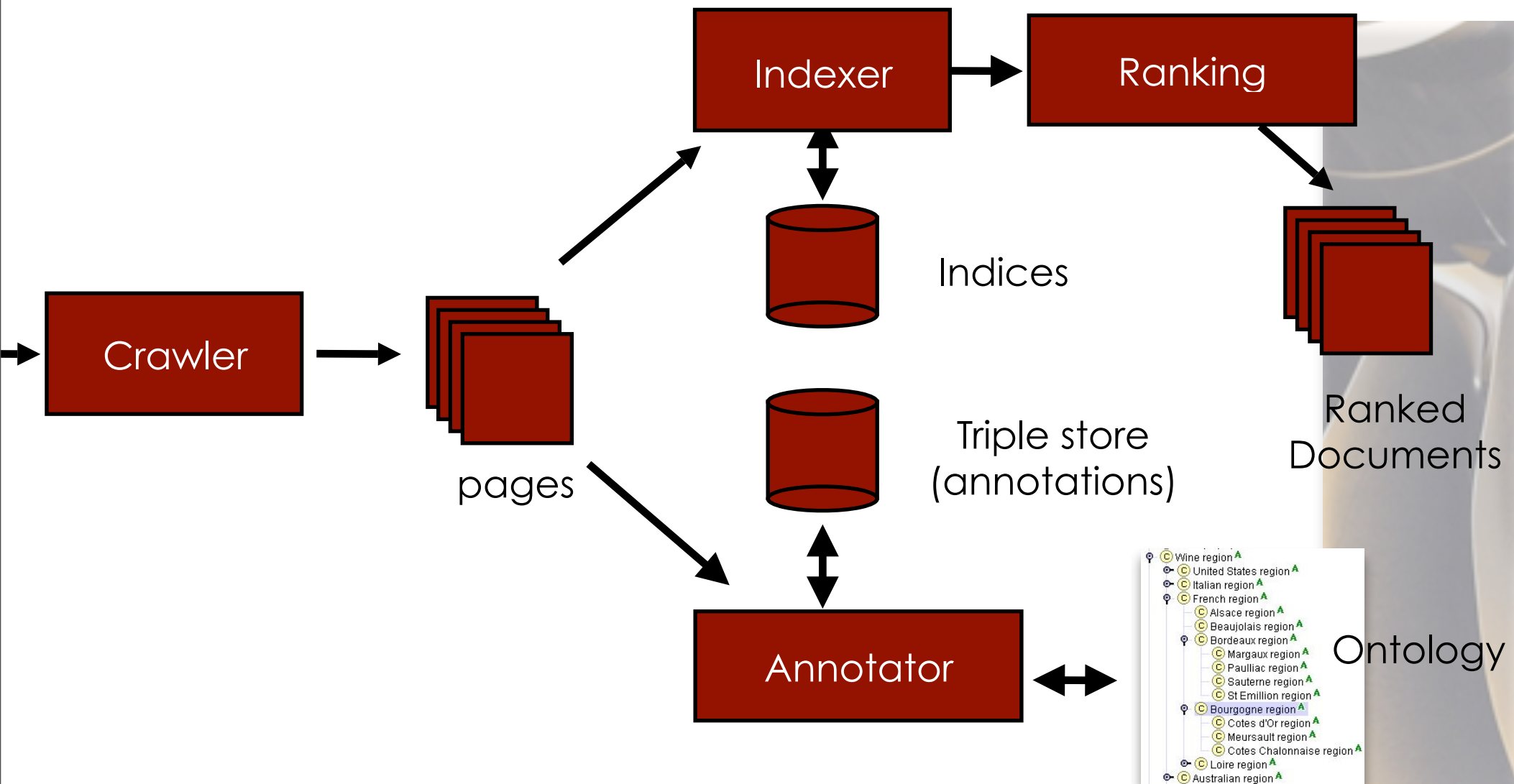
- Users can mix keywords and ontology-based search
 - Accuracy of Ontology-based searching available
 - When metadata covers information
 - Expressiveness of Keyword querying is available
 - For all other cases
 - Keyword-in-context available
 - Keyword matching available for matching concepts names
 - e.g. “fuel” is matched only on snippets of texts annotated as removed parts
- Uses provenance of annotations
 - Portion of document annotated with concepts are stored in 3store



- Enables querying documents using hybrid search
- Enables quantification of unstructured information
- Currently applied to:
 - Event Reports (1998-2004),
 - Technical variances
- Finalist of Rolls-Royce Creativity Award 2007



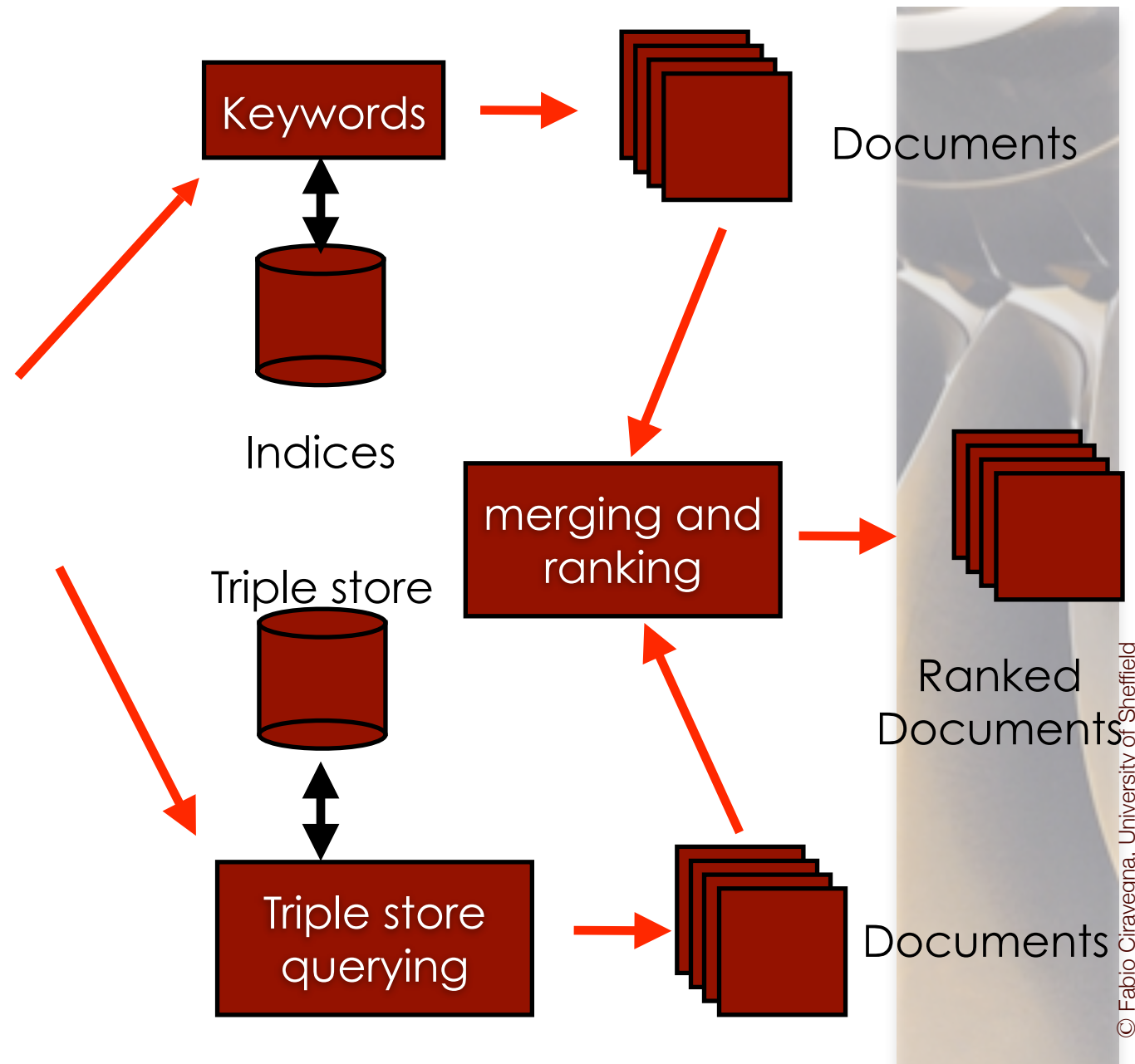
K-Search: indexing



K-Search is the Companion of K-Forms

K-Search: retrieval

The screenshot shows the K-Search web interface. At the top, there are tabs for 'Search', 'Results', and 'Graph'. The 'Search' tab is active. Below the tabs, there is a 'Keyword Search:' field with a magnifying glass icon and the text '(optional)'. To the right of this field is a dropdown menu for 'Number of results per page' with 'ALL' selected. Below these fields, there is a section for 'that match the following criteria :'. It includes two input fields for 'hascolumnB:' and 'hascolumnA:', each with a magnifying glass icon and the text '(optional)'. Between these fields is an 'AND' label. To the right of the 'hascolumnB:' field is an 'OR' label. Below the 'hascolumnB:' field is a 'test' button. To the right of the 'test' button is a 'table' button. At the bottom left, there is a 'know' logo. At the bottom right, there is a 'Available Reports' section with a 'test' button. A 'SEARCH' button is located at the bottom right of the search area.



- 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



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



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



SearchGraph

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 Level ...

_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.htm03ER.htm_030100ER.htmClose

Event Report No. SIA / 703 Prepared By: name of person Originated/Revised on: 12 November 2007

Event Report Data engine name here Rolls-Royce place here

FSO

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: 9573

Engine TSN/CSN: TSN here name

Event Type: Non-Operational

© Fabio Ciravegna, University of Sheffield

Graph visualisation

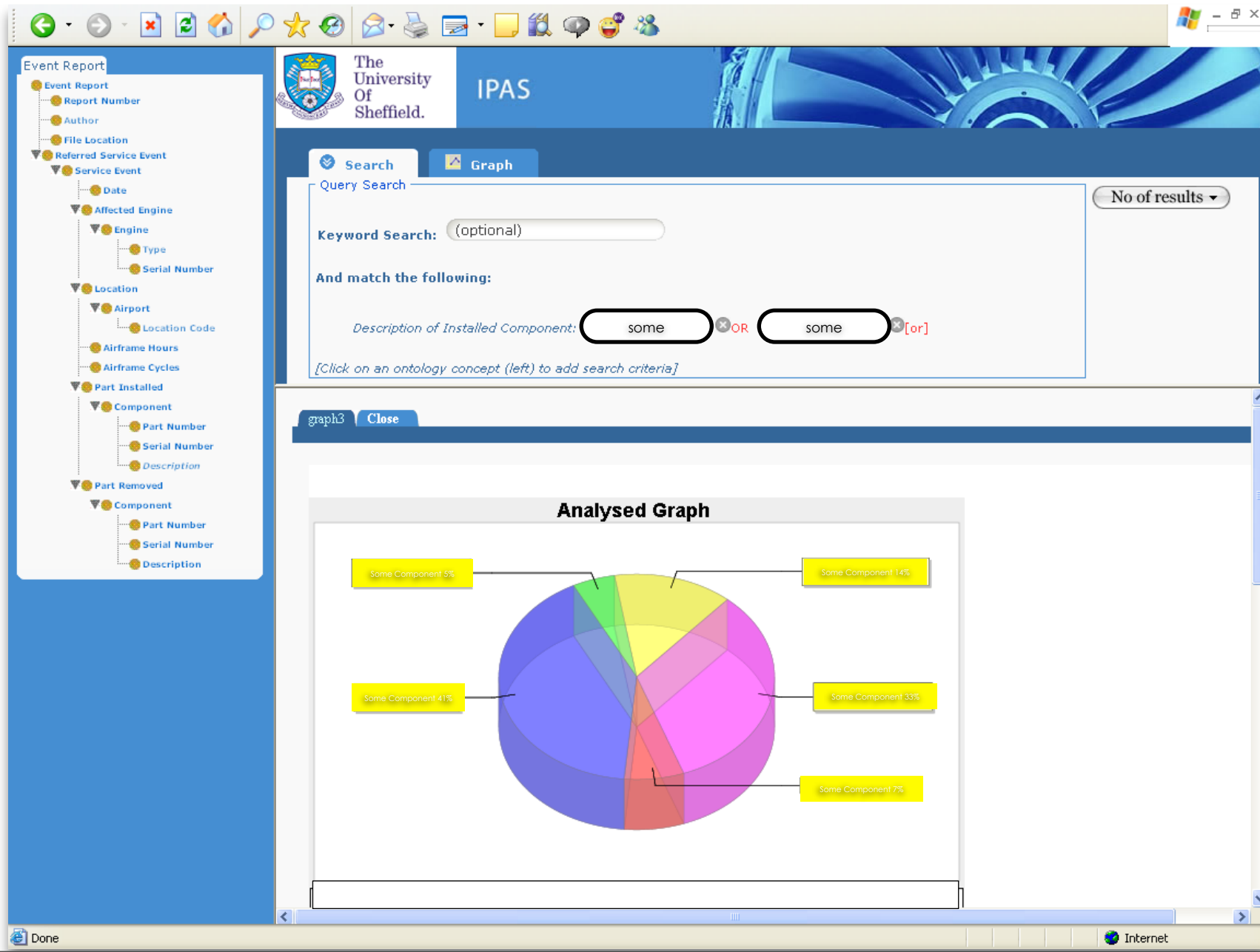
- Query results can be visualised as a graph
- Users can select graph types
 - ▣ Bar chart
 - ▣ Pie chart
- Users can select group and subgroup
 - ▣ For aggregating results
- Graphs are opened in the tab interface alongside other documents



Graph - example

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

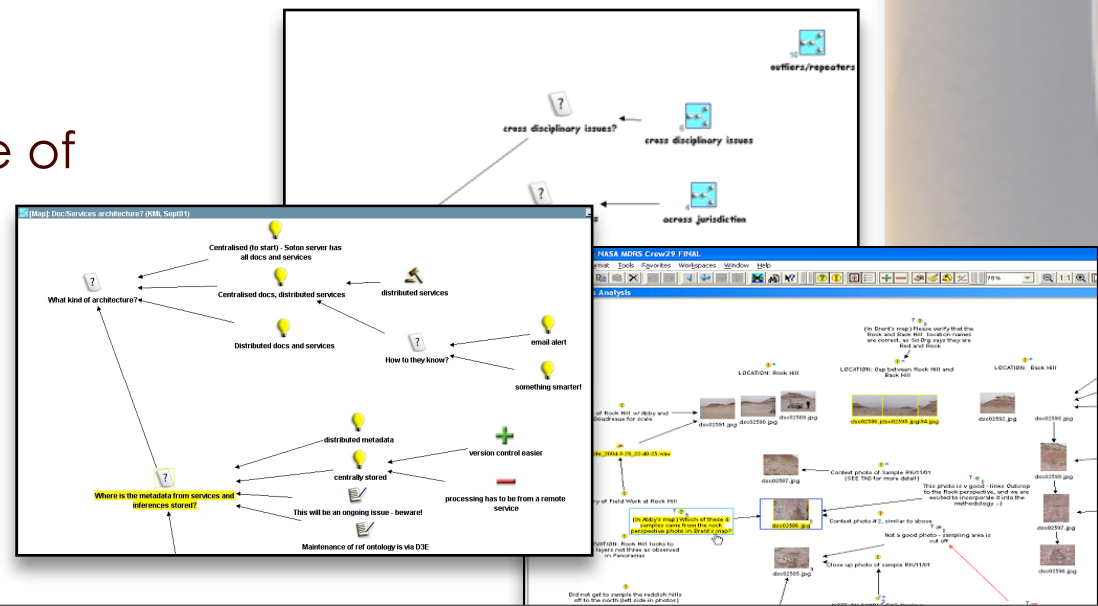
130



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



- 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
 - If the analysis is still true (titanium cost has decreased)
 - Compendium (Buckingham-Shum 2002)
 - D/Red (Wallace *et al.* 2005, 2007)
- What solutions were tried (use of titanium)
 - Why they were not adopted (e.g. too high a cost)



■ Infrastructure:

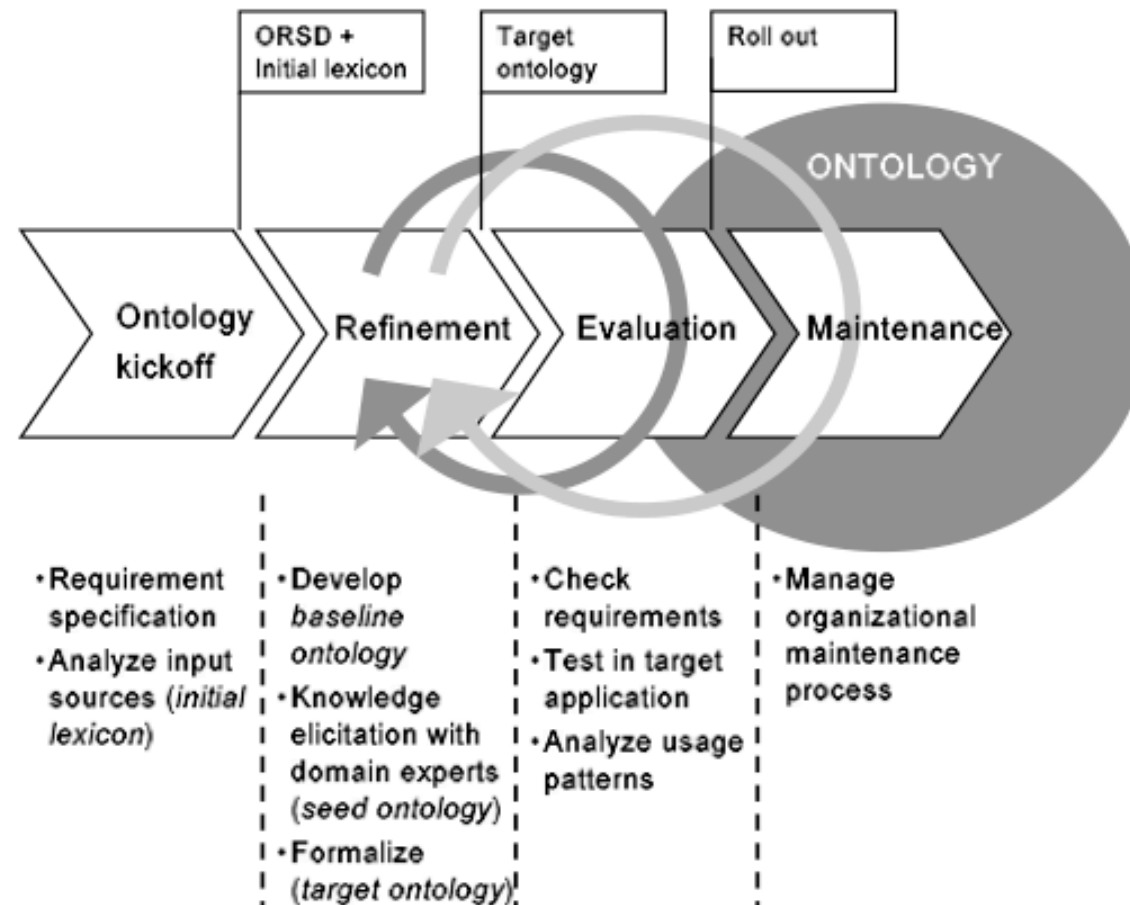
- Different media cannot easily be shared in the same way as knowledge from one format.
 - A folder of text documents may still be easily sent via email, but a folder of image files may not, and may instead require a shared image repository.
 - For 10 GByte of noise sensor data even an upload to a centralized 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.



- 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?



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

- 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

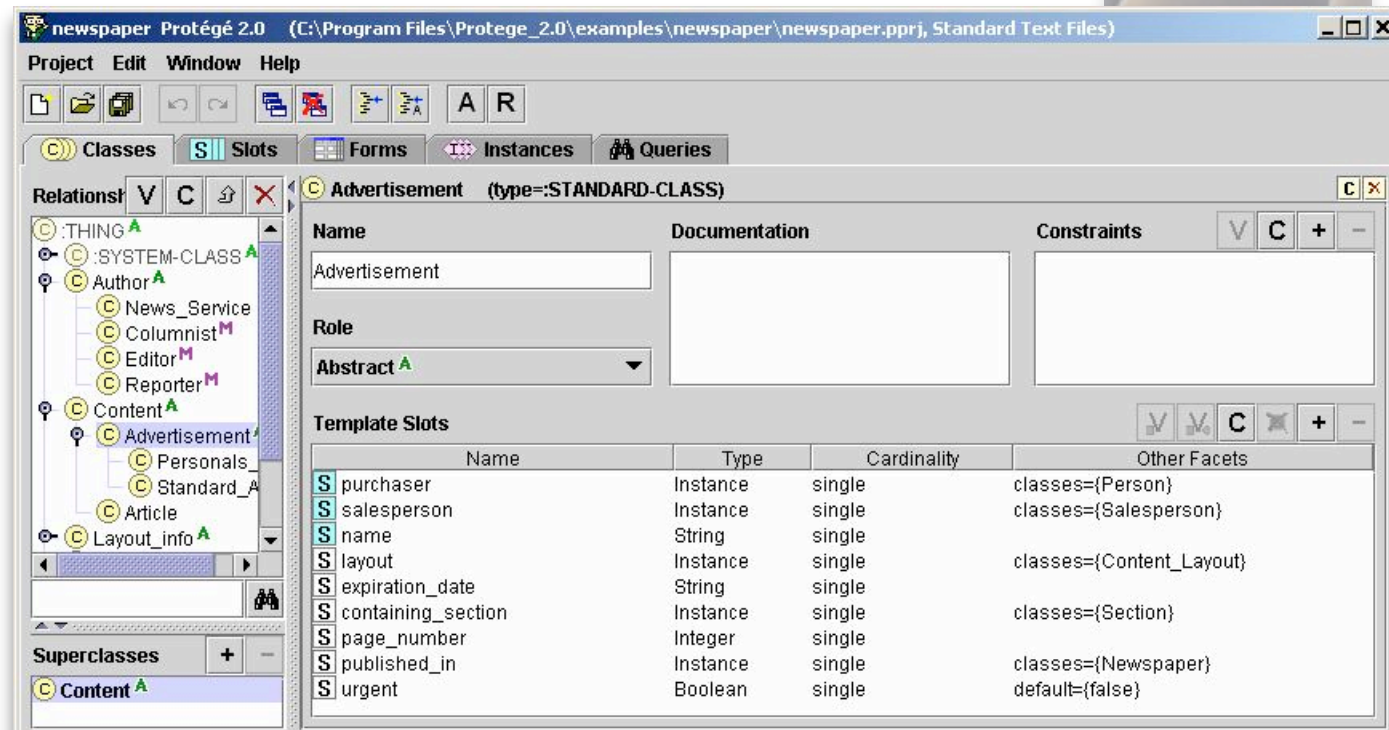


Editing Ontology Tools: Protégé

140

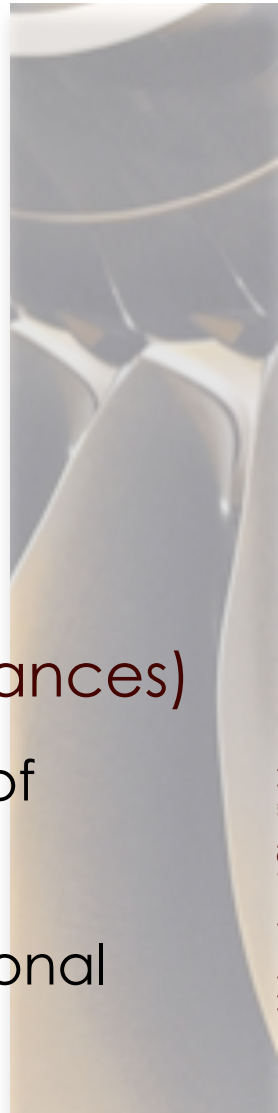
- Protégé is a free, open source ontology editor
- Download at <http://protege.stanford.edu/>
- Protégé ontologies can be exported into a variety of formats including RDF(S), OWL, and XML Schema

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



- 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



- 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

- 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** emphasise 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

- 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



- 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.



- 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.



- Reusing existing resources can help ontology development
 - 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 may decrease performance
 - Incompleteness
 - Background knowledge only available for part of the problem space
 - Inconsistency
 - Pieces of knowledge from different sources can be conflicting



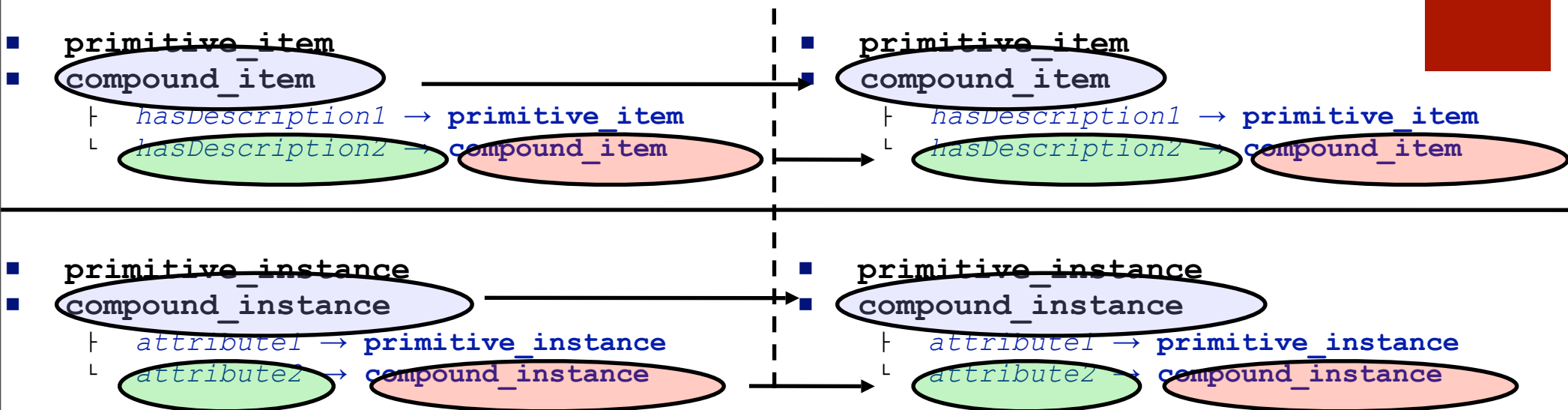
- Different communities use different domain representations
 - Design department uses designs and associated part lists
 - Service department uses illustrated part catalogue
 - The largely overlap but:
 - Use different URIs (part numbers)
 - Use different descriptions of the domain
 - In service what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



- Main difference between db schema and 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>

- 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

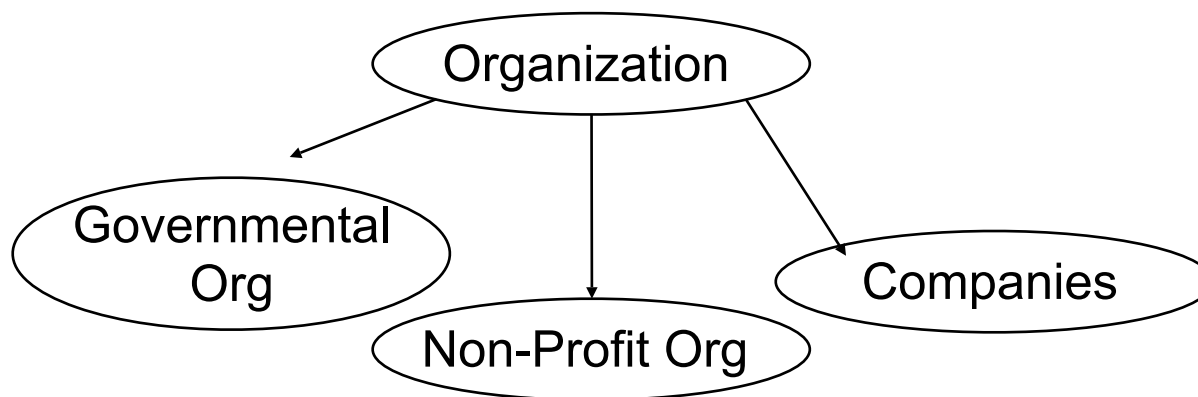


- Easy case

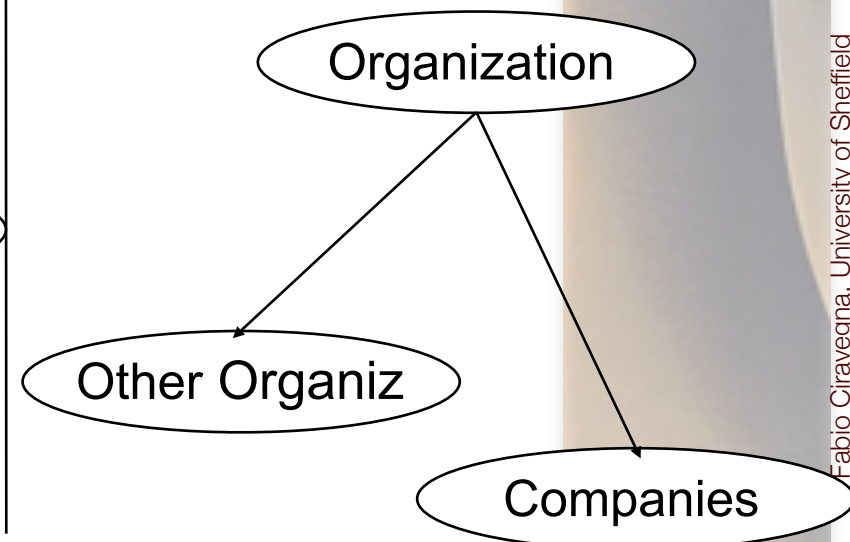
Specific to generic

- Collapse all instances and subtypes of Governmental_Org and Non-profit_Org into Other_Organiz

Source (existing ontology)



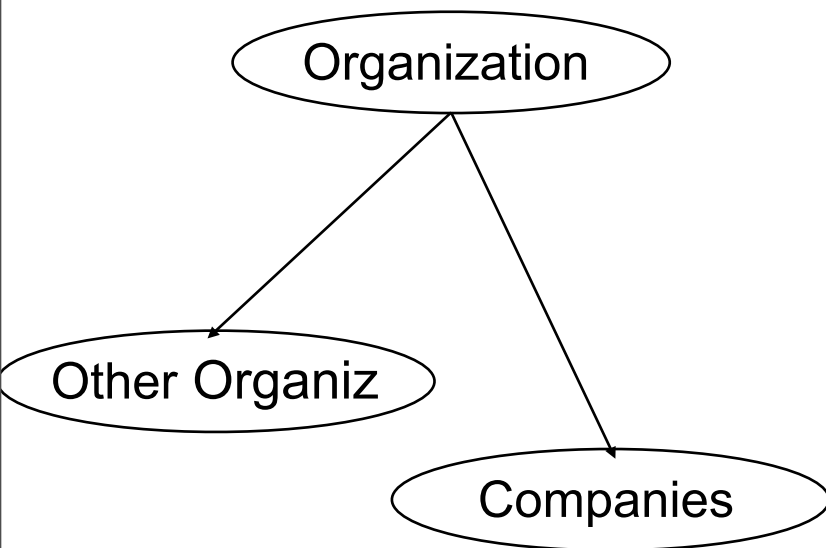
Target (new ontology)



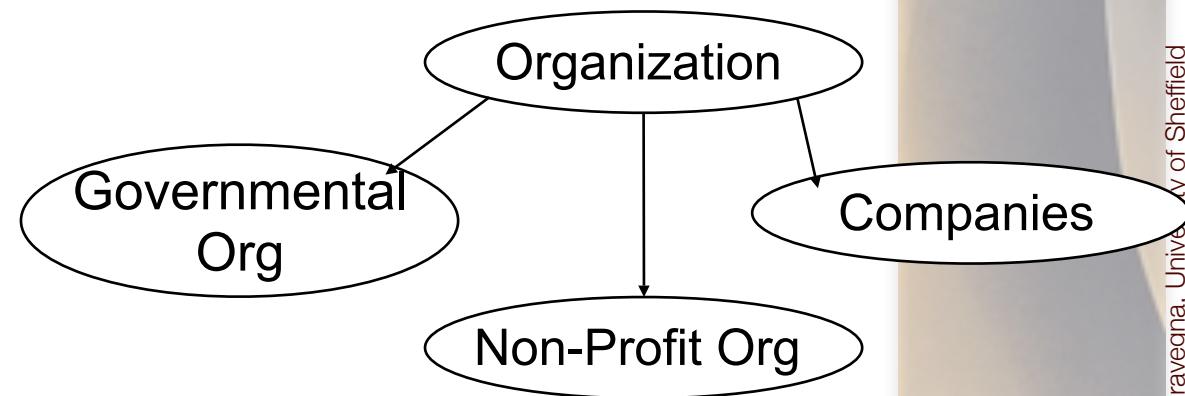
Generic to Specific

- 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)



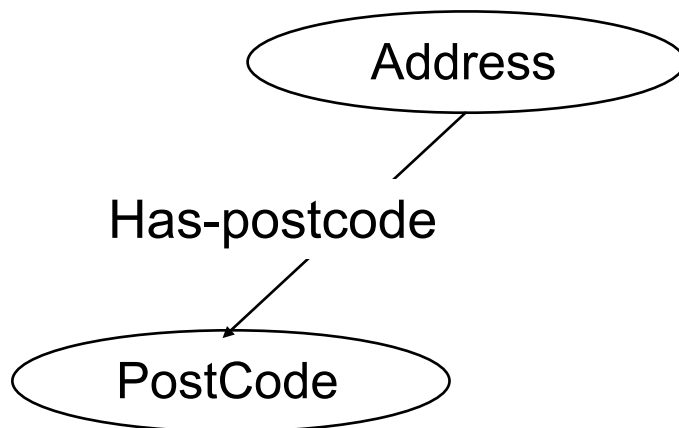
Target (new ontology)



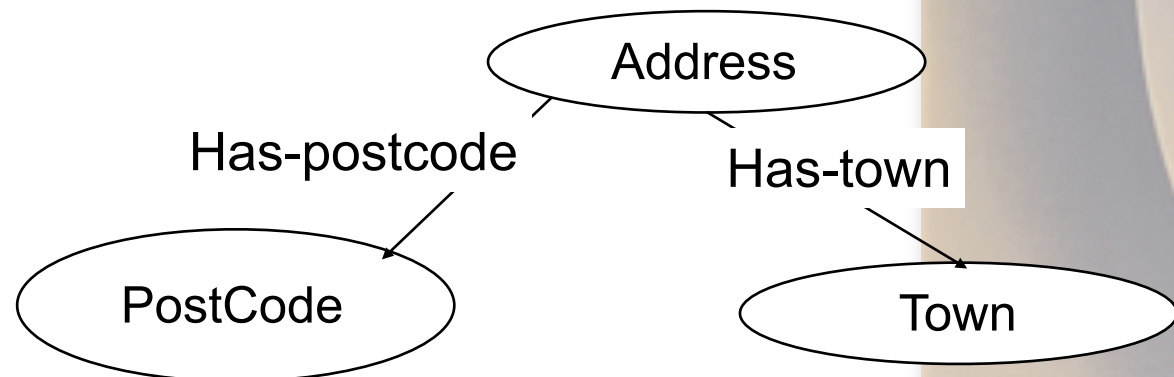
A more complex case

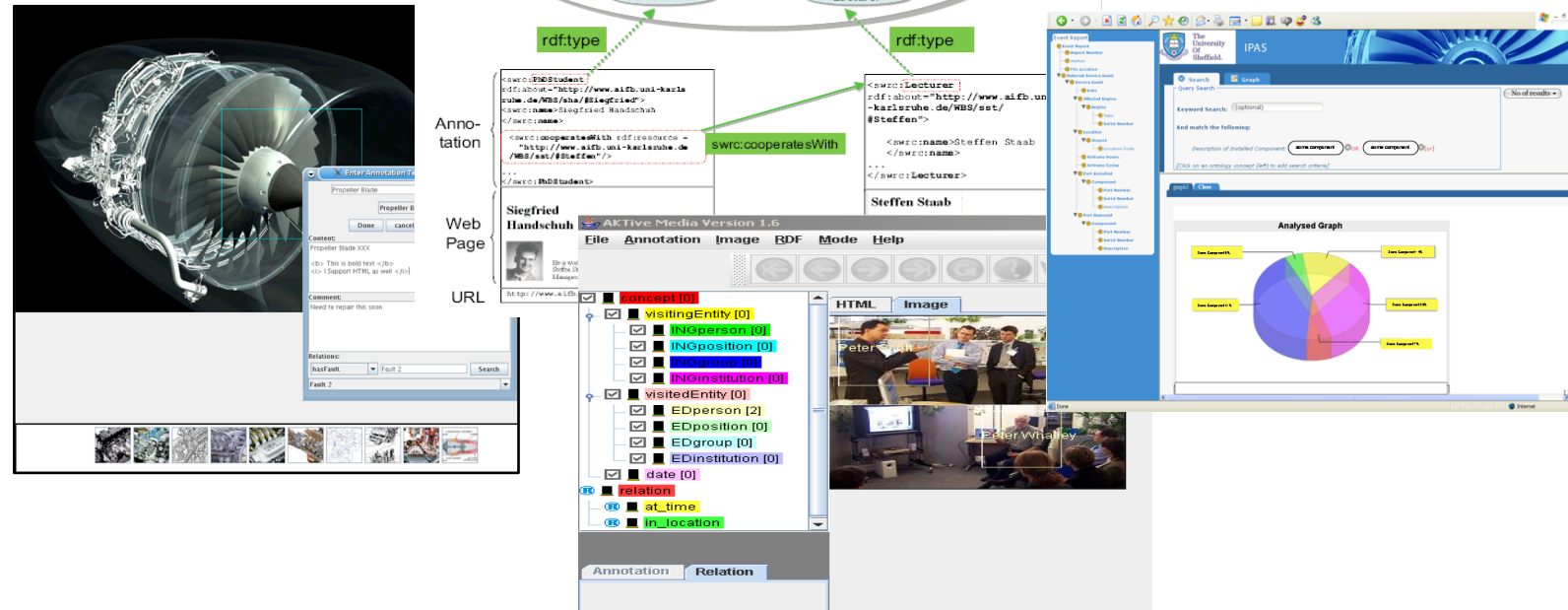
- Difficult case
 - In the DB the address implicitly includes the town
 - E.g. substring
 - In the Domain Ontology the town is explicitly mentioned
 - No easy way to map

Source (existing ontology)



Target (new ontology)





Conclusions

- 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!

Future Trends

159

	less than 2 years	2 to 5 years	5 to 10 years	more than 10
transformational	Web 2.0	Mobile Phone Payments	Collective Intelligence RFID (Case/Pallet) RFID (Item)	DNA Logic Quantum Com
high	Ajax Internal Web Services Location-Aware Technology Social Network Analysis VoIP	Digital Paper/E-Paper Grid Computing Location-Aware Applications	Corporate Semantic Web Event-Driven Architecture Model-Driven Architectures	Augmented R Mesh Network Tera-architect
moderate	Corporate Blogging Mashup Smartphone	Enterprise Instant Messaging Offline Ajax RSS Enterprise Speech Recognition for Mobile Devices Tablet PC Wikis	Biometric Payments Prediction Markets Speech-to-Speech Translation	Telepresence
low	Folksonomies		IPv6	

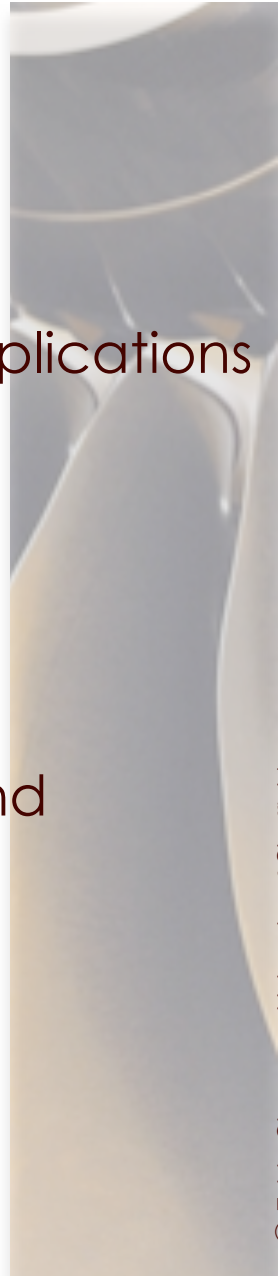
Priority Matrix for Emerging Technologies [Gartner 2006]

■ 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



- 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

162

■ Contact Information

- www.dcs.shef.ac.uk/~fabio
- fabio@dc.s.shef.ac.uk

■ Intelligent Web Technologies Lab

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

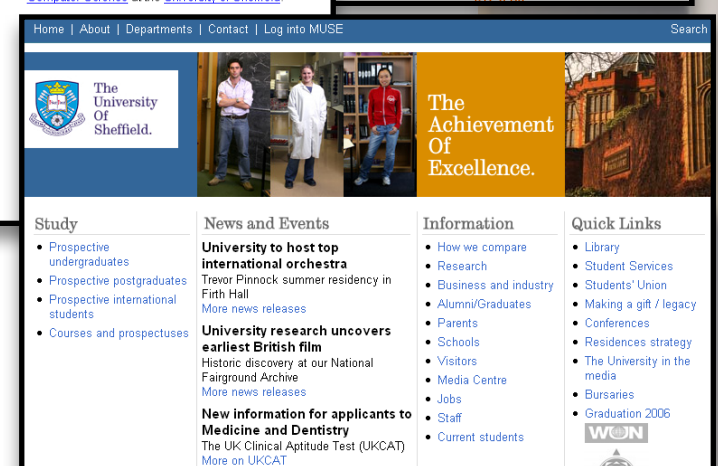
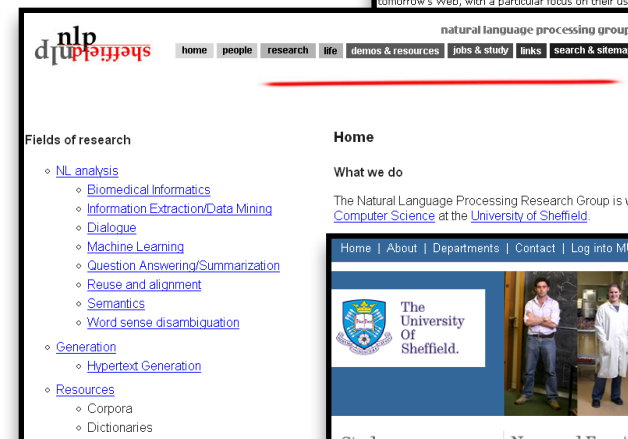
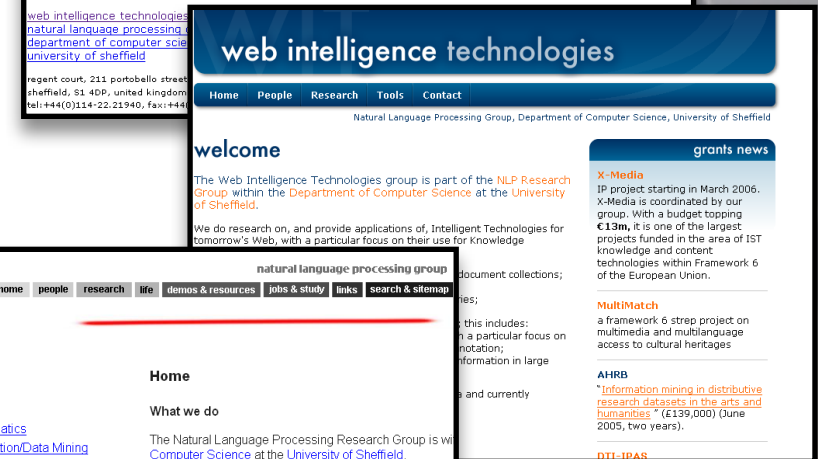
■ NLP Sheffield

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

■ University of Sheffield

- www.shef.ac.uk

■ K-Now Ltd (see next page)



- 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)

164

- 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

- 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
- 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