

Ontologies and Ontology alignment

Patrick Lambrix
Linköpings universitet

HIT-MSRA 2008 Summer School, Harbin, China, July 2008.

Outline

- Part I: Semantic Web and Ontologies
- Part II: Ontology alignment

Part I Semantic Web and Ontologies

Part I: Semantic Web and ontologies

- Semantic Web
- Ontologies
 - Definition
 - Use
 - Components
 - Knowledge representation

GET THAT PROTEIN!

Locating relevant information

Vision: Web services

- Databases and tools (service providers) announce their service capabilities
- Users request services which may be based on task descriptions
- Service matchers find relevant services (composition) based on user needs and user preferences, negotiate service delivery, and deliver results to user

Retrieving relevant information

Vision:

Based on the meaning of the query:

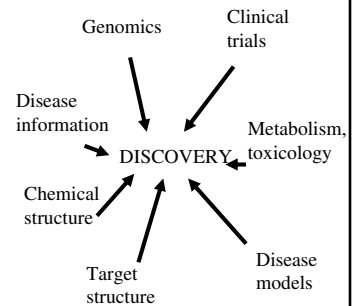
- only relevant information is retrieved
- all relevant information is retrieved



Integrating information

Vision:

Integrate data sources that are heterogeneous in content, data quality, data models, access methods, terminology



Today: syntactic Web

- A library of documents (web pages) interconnected by links
- A common portal to applications accessible through web pages, and presenting their results as web pages

A place where computers do the presentation (easy) and people do the linking and interpreting (hard).

Semantic Web

W3C: Facilities to put machine-understandable data on the Web are becoming a high priority for many communities. The Web can reach its full potential only if it becomes a place where data can be shared and processed by automated tools as well as by people. For the Web to scale, tomorrow's programs must be able to share and process data even when these programs have been designed totally independently. The Semantic Web is a vision: the idea of having data on the web defined and linked in a way that it can be used by machines not just for display purposes, but for automation, integration and reuse of data across various applications.

What is the problem?

Example based on example on slides by P. Patel-Schneider

What information can we see...

Date: 13-15 June, 2005
 Location: Linköping
 Sponsors: IEEE, CERC, LiU
 14th IEEE International Workshops on
 Enabling Technologies: Infrastructures for
 Collaborating Enterprises (WETICE-2005)
 Welcome to WETICE-2005

...

What information can a machine see...

2005 6 13-15
 Linköping
 IEEE, CERC, LiU
 14th IEEE International Workshops on Enabling Technologies: Infrastructures for Collaborating Enterprises (WETICE-2005)
 Welcome to WETICE-2005

Use XML markup with “meaningful” tags

```

<date> 13-15 June 2005 </date>

<location> Linköping </location>

<sponsors>IEEE, CERC, LiU </sponsors>

<name> 14th IEEE International Workshops on Enabling
Technologies: Infrastructures for Collaborating
Enterprises (WETICE-2005) </name>

<welcome> Welcome to WETICE-2005 </welcome>
  
```

Machine sees ...

```

<date>2005 6 13-15</date>

<location>Linköping </location>

<sponsors>IEEE, CERC, LiU </sponsors>

<name>14th IEEE International Workshops on Enabling
Technologies: Infrastructures for Collaborating
Enterprises (WETICE-2005) </name>

<welcome>Welcome to WETICE-2005 </welcome>
  
```

But what about ...

```

<date> 13-15 June 2005 </date>

<place> Linköping </place>

<sponsors>IEEE, CERC, LiU </sponsors>

<conf> 14th IEEE International Workshops on Enabling
Technologies: Infrastructures for Collaborating
Enterprises (WETICE-2005) </conf>

<introduction> Welcome to WETICE-2005 </introduction>
  
```

Machine sees ...

```

<date>2005 6 13-15</date>

<place>Linköping </place>

<sponsors>IEEE, CERC, LiU </sponsors>

<conf>14th IEEE International Workshops on Enabling
Technologies: Infrastructures for Collaborating
Enterprises (WETICE-2005) </conf>

<introduction>Welcome to WETICE-2005 </introduction>
  
```

Adding “Semantics” – first approach

External agreement on meaning of annotations

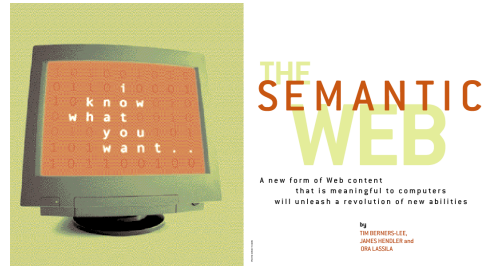
- Agree on the meaning of a set of annotation tags
- Problems with this approach:
 - Inflexible
 - Limited number of things can be expressed

Adding “Semantics” – second approach

Use on-line ontologies to specify meaning of annotations

- Ontologies provide a vocabulary of terms
- New terms can be formed by combining existing ones
- Meaning (semantics) of such terms is formally specified

Scientific American, May 2001:



- First step towards the vision:
adding semantic annotation to web resources

Semantic annotations based on ontologies

- Locating information
 - Web service descriptions use ontologies
 - Users use ontologies when formulating requests
 - Service matchers find services based on meaning
- Retrieving relevant information
 - Reduce non-relevant information (precision)
 - Find more relevant information (recall)
- Integrating information
 - Relating similar entities in different databases

Part I: Semantic Web and ontologies

- Semantic Web
- Ontologies
 - Definition
 - Use
 - Components
 - Knowledge representation

Ontologies

“Ontologies define the basic terms and relations comprising the vocabulary of a topic area, as well as the rules for combining terms and relations to define extensions to the vocabulary.”

(Neches, Fikes, Finin, Gruber, Senator, Swartout, 1991)

Definitions

- Ontology as specification of a conceptualization
- Ontology as philosophical discipline
- Ontology as informal conceptual system
- Ontology as formal semantic account
- Ontology as representation of conceptual system via a logical theory
- Ontology as the vocabulary used by a logical theory
- Ontology as a meta-level specification of a logical theory
(Guarino, Garetta)

Definitions

- An ontology is an explicit specification of a conceptualization (Gruber)
- An ontology is a hierarchically structured set of terms for describing a domain that can be used as a skeletal foundation for a knowledge base. (Swartout, Patil, Knight, Russ)
- An ontology provides the means for describing explicitly the conceptualization behind the knowledge represented in a knowledge base. (Bernaras, Lasergoiti, Corra)
- An ontology is a formal, explicit specification of a shared conceptualization (Studer, Benjamins, Fensel)

Example

```

GENE ONTOLOGY (GO)
immune response
  f- acute-phase response
  f- anaphylaxis
  f- antigen presentation
  f- antigen processing
  f- cellular defense response
  f- cytokine metabolism
    i- cytokine biosynthesis synonym cytokine production
      p- regulation of cytokine biosynthesis
      ...
  ...
  i- B-cell activation
    i- B-cell differentiation
    i- B-cell proliferation
  i- cellular defense response
  ...
  i- T-cell activation
    i- activation of natural killer cell activity
  ...
  
```

Example Ontologies

- Knowledge representation ontology: frame ontology
- Top level ontologies: TLO, Cyc
- Linguistic ontologies: GUM, WordNet
- Engineering ontologies: EngMath, PhysSys
- Domain ontologies: CHEMICALS, Gene Ontology, Open Biomedical Ontologies

Ontologies used ...

- for communication between people and organizations
- for enabling knowledge reuse and sharing
- as basis for interoperability between systems
- as repository of information
- as query model for information sources

Key technology for the Semantic Web

Ontologies in biomedical research

- many biomedical ontologies
e.g. GO, OBO, SNOMED-CT
- practical use of biomedical ontologies
e.g. databases annotated with GO

```

GENE ONTOLOGY (GO)
immune response
  i- acute-phase response
  i- anaphylaxis
  i- antigen presentation
  i- antigen processing
  i- cellular defense response
  i- cytokine metabolism
    i- cytokine biosynthesis synonym cytokine production
      ...
  ...
  i- B-cell activation
    i- B-cell differentiation
    i- B-cell proliferation
  i- cellular defense response
  ...
  i- T-cell activation
    i- activation of natural killer cell activity
  ...
  
```

Components

- concepts
 - represent a set or class of entities in a domain
immune response
 - organized in taxonomies (hierarchies based on e.g. *is-a* or *is-part-of*)
immune response is-a defense response
- instances
 - often not represented in an ontology (instantiated ontology)

Components

- relations

R: C1 x C2 x ... x Cn

Protein hasName ProteinName

*Chromosome hasSubcellularLocation
Nucleus*

Components

- axioms

'facts that are always true'

*The origin of a protein is always of the type
'gene coding origin type'*

Each protein has at least one source.

A helix can never be a sheet and vice versa.

Different kinds of ontologies

- Controlled vocabularies

Concepts

- Taxonomies

Concepts, is-a

- Thesauri

Concepts, predefined relations

- Data models (e.g. EER, UML)

Concepts, relations, axioms

- Logics

Concepts, relations, axioms

Taxonomy - GeneOntology

id: GO:0003674 name: molecular_function
def: "Elemental activities, such as catalysis or binding, describing the actions of a gene product at the molecular level. A given gene product may exhibit one or more molecular functions."

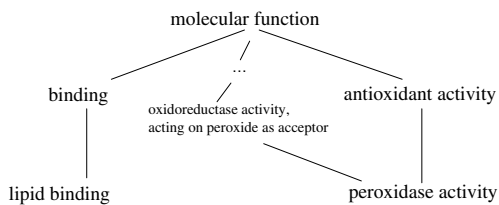
id: GO:0015643 name: binding
def: "The selective, often stoichiometric, interaction of a molecule with one or more specific sites on another molecule."
is-a: GO:0003674 ! molecular_function

id: GO:0008289 name: lipid binding
is_a: GO:0015643 ! binding

id: GO:0016209 name: antioxidant activity
def: "Inhibition of the reactions brought about by dioxygen (O2) or peroxides. Usually the antioxidant is effective because it can itself be more easily oxidized than the substance protected."
is_a: GO:0003674 ! molecular_function

id: GO:0004601 name: peroxidase activity
def: "Catalysis of the reaction: donor + H2O2 - oxidized donor + 2 H2O."
is_a: GO:0016209 ! antioxidant activity
is_a: GO:0016684 ! oxidoreductase activity, acting on peroxide as acceptor

Taxonomy - GeneOntology



Thesaurus

- graph

- fixed set of relations

(synonym, narrower term, broader term, similar)

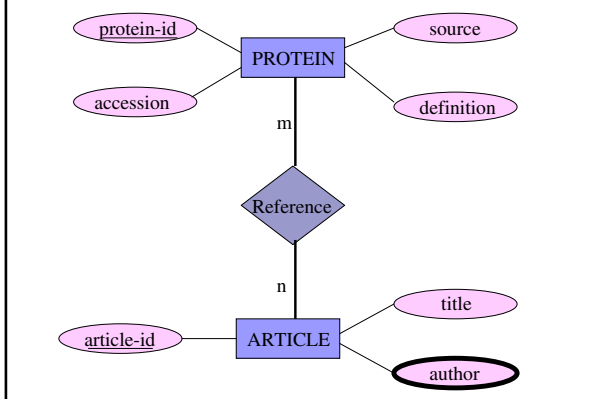
Thesaurus - WordNet

thesaurus, synonym finder
 => workbook
 => reference book, reference, reference work, book of facts
 => book
 => publication
 => print media
 => medium
 => means
 => instrumentality, instrumentation
 => artifact, artefact
 => object, inanimate object, physical object
 => entity
 => work, piece of work
 => product, production
 => creation
 => artifact, artefact
 => object, inanimate object, physical object
 => entity

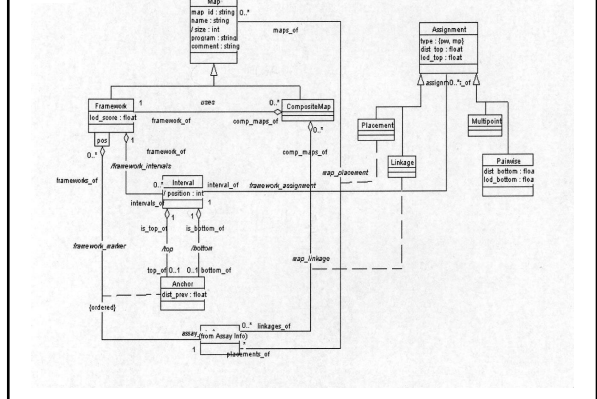
OO Data models

- EER
entity types, attributes, relationships, cardinality constraints, taxonomy
- UML
classes, attributes, associations, cardinality constraints, taxonomy, operations
- Taxonomy/inheritance – semantics?
- Intuitive, lots of tools, widely used.

Entity-relationship



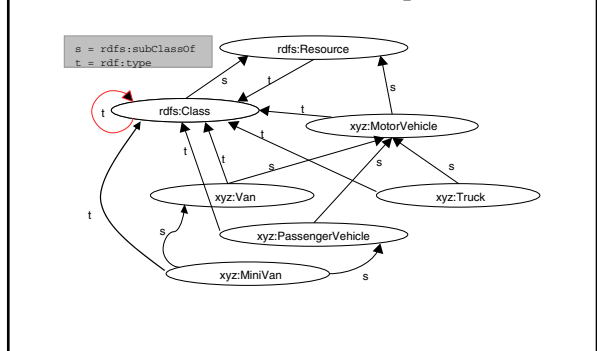
UML



RDF + RDF Schema

- Basic construct: sentence: *Subject Predicate Object*
 - Encoded in XML
 - Can be seen as ground atomic formula
 - Represented as graph
- RDF Schema
- Editors, query tools exist

RDF Schema - example



Logics

- Formal languages
- Syntax, semantics, inference mechanisms

Logics

Reasoning services used in

- **Ontology design**
Check concept satisfiability, ontology satisfiability and (unexpected) implied relationships
- **Ontology aligning and merging**
Assert inter-ontology relationships.
Reasoner computes integrated concept hierarchy/consistency.
- **Ontology deployment**
Determine if a set of facts are consistent w. r. t. ontology.
Determine if individuals are instances of ontology concepts.
Query inclusion.
Classification-based querying.

Description Logics

- A family of KR formalisms tailored for expressing knowledge about concepts and concept hierarchies
- Based on FOPL, supported by automatic reasoning systems
- Basic building blocks: concepts (concepts), roles (binary relations), individuals (instances)
- Language constructs can be used to define new concepts and roles (axioms).
 - Intersection, union, negation, quantification, ...
- Knowledge base is Tbox + Abox
 - Tbox: concept level - axioms: equality and subsumption (is-a)
 - Abox: instance level - axioms: membership, relations
- Reasoning services
 - Satisfiability of concept, Subsumption/Equivalence/Disjointness between concepts, Classification, Instantiation, Retrieval

Description Logics

Intersection

Signal-transducer-activity \cap binding

Negation

\neg Helix

Quantifiers

\exists hasOrigin.Mitochondrion

\forall hasOrigin.Gene-coding-origin-type

OWL

- OWL-Lite, OWL-DL, OWL-Full: increasing expressivity
- A legal OWL-Lite ontology is a legal OWL-DL ontology is a legal OWL-Full ontology
- OWL-DL: expressive description logic, decidable
- XML-based
- RDF-based (OWL-Full is extension of RDF, OWL-Lite and OWL-DL are extensions of a restriction of RDF)

OWL-Lite

- **Class**, subclassOf, equivalentClass
- intersectionOf (only named classes and restrictions)
- **Property**, subPropertyOf, equivalentProperty
- domain, range (global restrictions)
- inverseOf, TransitiveProperty (*), SymmetricProperty, FunctionalProperty, InverseFunctionalProperty
- allValuesFrom, someValuesFrom (local restrictions)
- minCardinality, maxCardinality (only 0/1)
- **Individual**, sameAs, differentFrom, AllDifferent

(*) restricted

OWL-DL

- **Type separation** (class cannot also be individual or property, property cannot be also class or individual), Separation between DatatypeProperties and ObjectProperties
- **Class –complex classes**, subClassOf, equivalentClass, *disjointWith*
- *intersectionOf*, *unionOf*, *complementOf*
- **Property**, subPropertyOf, equivalentProperty
- domain, range (global restrictions)
- *inverseOf*, *TransitiveProperty* (*), *SymmetricProperty*, *FunctionalProperty*, *InverseFunctionalProperty*
- *allValuesFrom*, *someValuesFrom* (local restrictions), *oneOf*, *hasValue*
- *minCardinality*, *maxCardinality*
- **Individual**, *sameAs*, *differentFrom*, *AllDifferent*

(*) restricted

Defining ontologies is not so easy ...

The Celestial Emporium of Benevolent Knowledge, Borges

"On those remote pages it is written that animals are divided into:

- a. those that belong to the Emperor
- b. embalmed ones
- c. those that are trained
- d. suckling pigs
- e. mermaids
- f. fabulous ones
- g. stray dogs
- h. those that are included in this classification
- i. those that tremble as if they were mad
- j. innumerable ones
- k. those drawn with a very fine camel's hair brush
- l. others
- m. those that have just broken a flower vase
- n. those that resemble flies from a distance"

Slide from talk by C. Goble

Defining ontologies is not so easy ...

Dyirbal classification of objects in the universe

- Bayi: men, kangaroos, possums, bats, most snakes, most fishes, some birds, most insects, the moon, storms, rainbows, boomerangs, some spears, etc.
- Balan: women, anything connected with water or fire, bandicoots, dogs, platypus, echidna, some snakes, some fishes, most birds, fireflies, scorpions, crickets, the stars, shields, some spears, some trees, etc.
- Balam: all edible fruit and the plants that bear them, tubers, ferns, honey, cigarettes, wine, cake.
- Bala: parts of the body, meat, bees, wind, yamsticks, some spears, most trees, grass, mud, stones, noises, language, etc.

Slide from talk by C. Goble

Ontology tools

- Ontology development tools
- Ontology merge and alignment tools
- Ontology evaluation tools
- Ontology-based annotation tools
- Ontology storage and querying tools
- Ontology learning tools

Part II Ontology Alignment

Part II – Ontology Alignment

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Recommending ontology alignment strategies
- Current issues

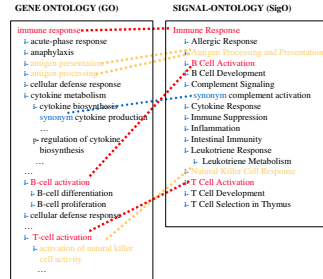
Ontologies in biomedical research

- many biomedical ontologies
e.g. GO, OBO, SNOMED-CT
- practical use of biomedical ontologies
e.g. databases annotated with GO

GENE ONTOLOGY (GO)

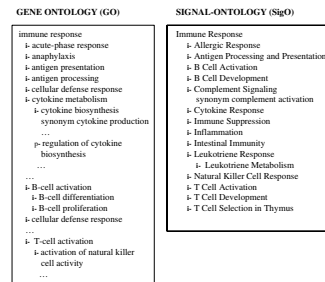
immune response
 i- acute-phase response
 i- anaphylaxis
 i- antigen presentation
 i- antigen processing
 i- cellular defense response
 i- cytokine metabolism
 i- cytokine biosynthesis
 synonym cytokine production
 p- regulation of cytokine biosynthesis
 ...
 i- B-cell activation
 i- B-cell differentiation
 i- B-cell proliferation
 i- cellular defense response
 ...
 i- T-cell activation
 i- activation of natural killer cell activity
 ...

Ontologies with overlapping information

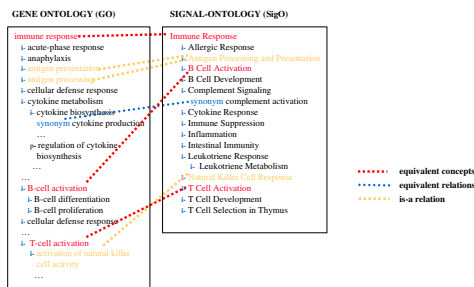


Ontologies with overlapping information

- Use of multiple ontologies
e.g. custom-specific ontology + standard ontology
 - Bottom-up creation of ontologies
experts can focus on their domain of expertise
- important to know the inter-ontology relationships



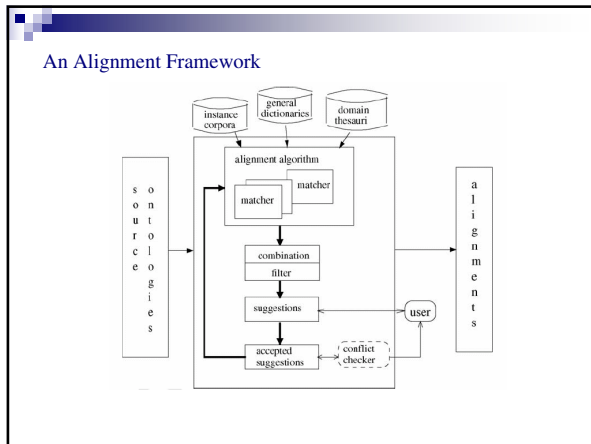
Ontology Alignment



Defining the relations between the terms in different ontologies

Part II – Ontology Alignment

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Recommending ontology alignment strategies
- Current issues



- ### Classification
- According to input
 - KR: OWL, UML, EER, XML, RDF, ...
 - components: concepts, relations, instance, axioms
 - According to process
 - What information is used and how?
 - According to output
 - 1-1, m-n
 - Similarity vs explicit relations (equivalence, is-a)
 - confidence

Matchers

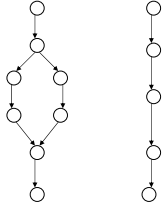
- ### Matcher Strategies
- Strategies based on linguistic matching
 - Structure-based strategies
 - Constraint-based
 - Instance-based
 - Use of auxiliary
-

- ### Example matchers
- Edit distance
 - Number of deletions, insertions, substitutions required to transform one string into another
 - aaaa → baab: edit distance 2
 - N-gram
 - N-gram : N consecutive characters in a string
 - Similarity based on set comparison of n-grams
 - aaaa : {aa, aa, aa}; baab : {ba, aa, ab}

- ### Matcher Strategies
- Strategies based on linguistic matching
 - Structure-based strategies
 - Constraint-based
 - Instance-based st
 - Use of auxiliary
-

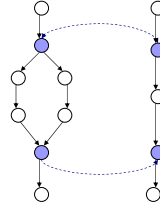
Example matchers

- Propagation of similarity values
- Anchored matching



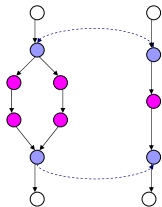
Example matchers

- Propagation of similarity values
- Anchored matching



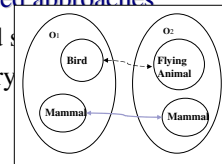
Example matchers

- Propagation of similarity values
- Anchored matching



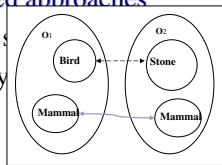
Matcher Strategies

- Strategies based on linguistic matching
- Structure-based strategies
- **Constraint-based approaches**
- Instance-based strategies
- Use of auxiliary



Matcher Strategies

- Strategies based on linguistic matching
- Structure-based strategies
- **Constraint-based approaches**
- Instance-based strategies
- Use of auxiliary

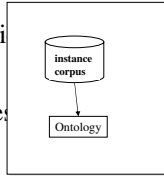


Example matchers

- Similarities between data types
- Similarities based on cardinalities

Matcher Strategies

- Strategies based on linguistics
- Structure-based strategies
- Constraint-based approaches
- Instance-based strategies
- Use of auxiliary information



Example matchers

- Instance-based
- Use life science literature as instances
- Structure-based extensions

Learning matchers – instance-based strategies

- Basic intuition
 - A similarity measure between concepts can be computed based on the probability that documents about one concept are also about the other concept and vice versa.
- Intuition for structure-based extensions
 - Documents about a concept are also about their super-concepts.
 - (No requirement for previous alignment results.)

Learning matchers - steps

- Generate corpora
 - Use concept as query term in PubMed
 - Retrieve most recent PubMed abstracts
- Generate text classifiers
 - One classifier per ontology / One classifier per concept
- Classification
 - Abstracts related to one ontology are classified by the other ontology's classifier(s) and vice versa
- Calculate similarities

Basic Naïve Bayes matcher

- Generate corpora
- Generate classifiers
 - Naïve Bayes classifiers, one per ontology
- Classification
 - Abstracts related to one ontology are classified to the concept in the other ontology with highest posterior probability $P(C|d)$
- Calculate similarities

$$sim(C_1, C_2) = \frac{n_{NBC2}(C_1, C_2) + n_{NBC1}(C_2, C_1)}{n_D(C_1) + n_D(C_2)}$$

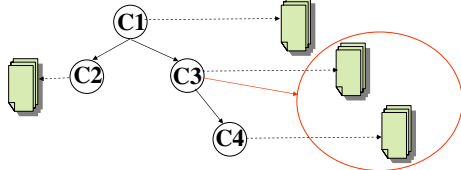
Basic Support Vector Machines matcher

- Generate corpora
- Generate classifiers
 - SVM-based classifiers, one per concept
- Classification
 - Single classification variant: Abstracts related to concepts in one ontology are classified to the concept in the other ontology for which its classifier gives the abstract the highest positive value.
 - Multiple classification variant: Abstracts related to concepts in one ontology are classified all the concepts in the other ontology whose classifiers give the abstract a positive value.
- Calculate similarities

$$\frac{n_{SVMC-C_2}(C_1, C_2) + n_{SVMC-C_1}(C_2, C_1)}{n_D(C_1) + n_D(C_2)}$$

Structural extension 'CI'

- Generate classifiers
 - Take (is-a) structure of the ontologies into account when building the classifiers
 - Extend the set of abstracts associated to a concept by adding the abstracts related to the sub-concepts



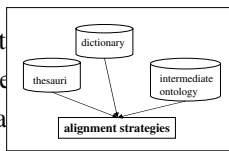
Structural extension 'Sim'

- Calculate similarities
 - Take structure of the ontologies into account when calculating similarities
 - Similarity is computed based on the classifiers applied to the concepts and their sub-concepts

$$sim_{struct}(C_1, C_2) = \frac{\sum_{C_i \subseteq C_1, C_j \subseteq C_2} nNBC_2(C_i, C_j) + \sum_{C_i \subseteq C_1, C_j \subseteq C_2} nNBC_1(C_j, C_i)}{\sum_{C_i \subseteq C_1} nD(C_i) + \sum_{C_j \subseteq C_2} nD(C_j)}$$

Matcher Strategies

- Strategies based linguistic
- Structure-based strategies
- Constraint-based approaches
- Instance-based strategies
- Use of auxiliary information



Example matchers

- Use of WordNet
 - Use WordNet to find synonyms
 - Use WordNet to find ancestors and descendants in the is-a hierarchy
- Use of Unified Medical Language System (UMLS)
 - Includes many ontologies
 - Includes many alignments (not complete)
 - Use UMLS alignments in the computation of the similarity values

Ontology Alignment and Merging Systems

	linguistic	structure	constraints	instances	auxiliary
ArtGen	name	parents, children		domain specific documents	WordNet
ASCO	name, label, description	parents, children, siblings, path from root			WordNet
Chimaera	name	parents, children			
FCA-Merge	name			domain specific documents	
FOAM	name, label	parents, children	equivalence		
GLUE	name	neighborhood		instances	
HCONE	name	parents, children			WordNet
IF-Map				instances	a reference ontology
iMapper		leaf, non-leaf, children, related node	domain, range	instances	WordNet
OntoMapper		parents, children		documents	
(Anchor-) PROMPT	name	direct graphs			
SAMBO	name, synonym	is-a and part-of, descendants and ancestors		domain specific documents	WordNet, UMLS
S-Match	label	path from root	semantic relations codified in labels		WordNet

Combinations

Combination Strategies

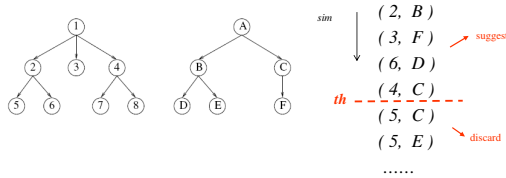
- Usually weighted sum of similarity values of different matchers

Filtering

Filtering techniques

Threshold filtering

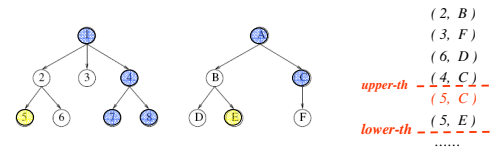
Pairs of concepts with similarity higher or equal than threshold are alignment suggestions



Filtering techniques

Double threshold filtering

- Pairs of concepts with similarity higher than or equal to **upper** threshold are alignment suggestions
- Pairs of concepts with similarity between **lower** and **upper** thresholds are alignment suggestions if they make sense with respect to the structure of the ontologies and the suggestions according to (1)



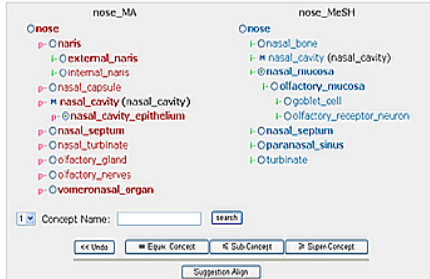
Example alignment system

SAMBO – matchers, combination, filter

Example alignment system

SAMBO – suggestion mode

Example alignment system SAMBO – manual mode



Part II – Ontology Alignment

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Recommending ontology alignment strategies
- Current issues

Evaluation measures

- Precision:

$$\frac{\# \text{ correct suggested alignments}}{\# \text{ suggested alignments}}$$
- Recall:

$$\frac{\# \text{ correct suggested alignments}}{\# \text{ correct alignments}}$$
- F-measure: combination of precision and recall

Ontology Alignment Evaluation Initiative

OAEI

- Since 2004
- Evaluation of systems
- Different tracks
 - comparison: benchmark (open)
 - expressive: anatomy (blind)
 - directories and thesauri: directory, food, environment, library (blind)
 - consensus: conference

OAEI

- Evaluation measures
 - Precision/recall/f-measure
 - recall of non-trivial alignments
 - full / partial golden standard

OAEI 2007

- 17 systems participated
 - benchmark (13)
 - ASMOV: $p = 0.95$, $r = 0.90$
 - anatomy (11)
 - AOAS: $f = 0.86$, $r+ = 0.50$
 - SAMBO: $f = 0.81$, $r+ = 0.58$
 - library (3)
 - Thesaurus merging: FALCON: $p = 0.97$, $r = 0.87$
 - Annotation scenario:
 - FALCON: $pb = 0.65$, $rb = 0.49$, $pa = 0.52$, $ra = 0.36$, $Ja = 0.30$
 - Silas: $pb = 0.66$, $rb = 0.47$, $pa = 0.53$, $ra = 0.35$, $Ja = 0.29$
 - directory (9), food (6), environment (2), conference (6)

OAEI 2007

- Systems can use only one combination of strategies per task
 - systems use similar strategies
 - text: string matching, tf-idf
 - structure: propagation of similarity to ancestors and/or descendants
 - thesaurus (WordNet)
 - domain knowledge important for anatomy task

Evaluation of algorithms

Cases

- GO vs. SigO

$\frac{GO: 70 \text{ terms}}{GO-immune \text{ defense}}$ $\frac{SigO: 15 \text{ terms}}{SigO-immune \text{ defense}}$ $\frac{GO: 60 \text{ terms}}{GO-behavior}$ $\frac{SigO: 10 \text{ terms}}{SigO-behavior}$

- MA vs. MeSH

$\frac{MA: 15 \text{ terms}}{MA-nose}$ $\frac{MeSH: 18 \text{ terms}}{MeSH-nose}$ $\frac{MA: 77 \text{ terms}}{MA-car}$ $\frac{MeSH: 39 \text{ terms}}{MeSH-car}$
 $\frac{MA: 112 \text{ terms}}{MA-eye}$ $\frac{MeSH: 45 \text{ terms}}{MeSH-eye}$

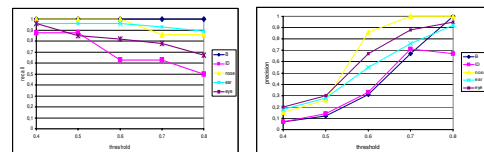
Evaluation of matchers

- Matchers
 - Term, TermWN, Dom, Learn (Learn+structure), Struc
- Parameters
 - Quality of suggestions: precision/recall
 - Threshold filtering : 0.4, 0.5, 0.6, 0.7, 0.8
 - Weights for combination: 1.0/1.2

KitAMO
<http://www.ida.liu.se/labs/iislab/projects/KitAMO>

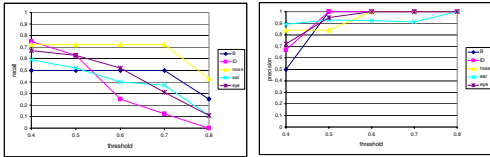
Results

- Terminological matchers



Results

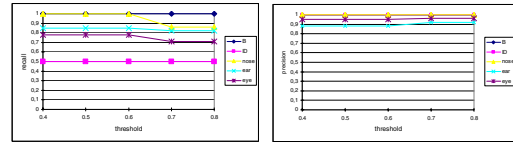
- Basic learning matcher (Naïve Bayes)



Naïve Bayes slightly better recall, but slightly worse precision than SVM-single
SVM-multiple (much) better recall, but worse precision than SVM-single

Results

- Domain matcher (using UMLS)



Results

- Comparison of the matchers

$CS_TermWN \supseteq CS_Dom \supseteq CS_Learn$

- Combinations of the different matchers
 - combinations give often better results
 - no significant difference on the quality of suggestions for different weight assignments in the combinations (but: did not check yet for large variations for the weights)
- Structural matcher did not find (many) new correct alignments (but: good results for systems biology schemas SBML – PSI MI)

Evaluation of filtering

- Matcher

TermWN

- Parameters

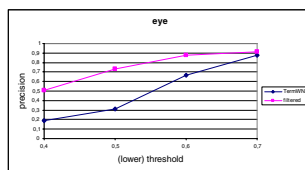
Quality of suggestions: precision/recall

Double threshold filtering using structure:

Upper threshold: 0.8

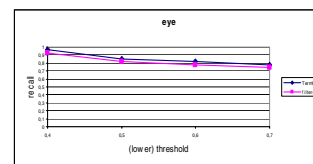
Lower threshold: 0.4, 0.5, 0.6, 0.7, 0.8

Results



- The precision for double threshold filtering with upper threshold 0.8 and lower threshold T is higher than for threshold filtering with threshold T

Results



- The recall for double threshold filtering with upper threshold 0.8 and lower threshold T is about the same as for threshold filtering with threshold T

Part II – Ontology Alignment

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- **Recommending ontology alignment strategies**
- Current issues

Recommending strategies - 1

- Use knowledge about previous use of alignment strategies
 - gather knowledge about input, output, use, performance, cost via questionnaires
 - Not so much knowledge available
 - OAEI

(Mochol, Jentzsch, Euzenat 2006)

Recommending strategies - 2

- Optimize
 - Parameters for ontologies, similarity assessment, matchers, combinations and filters
 - Run general alignment algorithm
 - User validates the alignment result
 - Optimize parameters based on validation

(Ehrig, Staab, Sure 2005)

Recommending strategies - 2

- Tests
 - travel in russia
QOM: $r=0.618$, $p=0.596$, $f=0.607$
Decision tree 150: $r=0.723$, $p=0.591$, $f=0.650$
 - bibster
QOM: $r=0.279$, $p=0.397$, $f=0.328$
Decision tree 150: $r=0.630$, $p=0.375$, $f=0.470$

Decision trees better than Neural Nets and Support Vector Machines.

Recommending strategies - 3

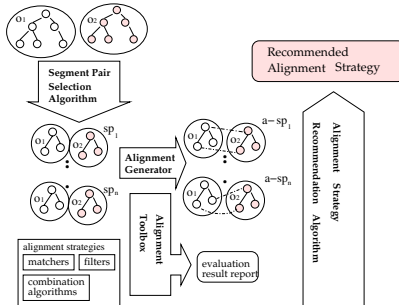
- Based on inherent knowledge
 - Use the actual ontologies to align to find good candidate alignment strategies
 - User/oracle with minimal alignment work
 - Complementary to the other approaches

(Tan, Lambrix 2007)

Idea

- Select small segments of the ontologies
- Generate alignments for the segments (expert/oracle)
- Use and evaluate available alignment algorithms on the segments
- Recommend alignment algorithm based on evaluation on the segments

Framework



Experiment case - Ontologies



- NCI thesaurus
 - National Cancer Institute, Center for Bioinformatics
 - Anatomy: 3495 terms
- MeSH
 - National Library of Medicine
 - Anatomy: 1391 terms

Experiment case - Oracle

- UMLS
 - Library of Medicine
 - Metathesaurus contains > 100 vocabularies
 - NCI thesaurus and MeSH included in UMLS
 - Used as approximation for expert knowledge
 - 919 expected alignments according to UMLS

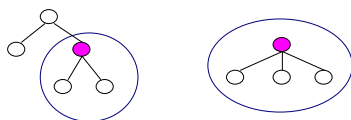
Experiment case – alignment strategies



- Matchers and combinations
 - N-gram (NG)
 - Edit Distance (ED)
 - Word List + stemming (WL)
 - Word List + stemming + WordNet (WN)
 - NG+ED+WL, weights 1/3 (C1)
 - NG+ED+WN, weights 1/3 (C2)
- Threshold filter
 - thresholds 0.4, 0.5, 0.6, 0.7, 0.8

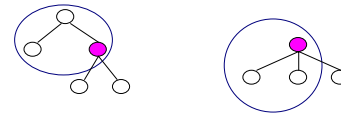
Segment pair selection algorithms

- SubG
 - Candidate segment pair = sub-graphs according to is-a/part-of with roots with same name; between 1 and 60 terms in segment
 - Segment pairs randomly chosen from candidate segment pairs such that segment pairs are disjoint



Segment pair selection algorithms

- Clust - Cluster terms in ontology
 - Candidate segment pair is pair of clusters containing terms with the same name; at least 5 terms in clusters
 - Segment pairs randomly chosen from candidate segment pairs

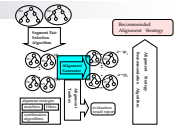


Segment pair selection algorithms

- For each trial, 3 segment pair sets with 5 segment pairs were generated
- SubG: A1, A2, A3
 - 2 to 34 terms in segment
 - level of is-a/part-of ranges from 2 to 6
 - max expected alignments in segment pair is 23
- Clust: B1, B2, B3
 - 5 to 14 terms in segment
 - level of is-a/part-of is 2 or 3
 - max expected alignments in segment pair is 4

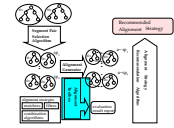
Segment pair alignment generator

- Used UMLS as oracle



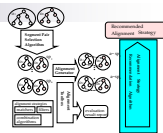
Alignment toolbox

- Used KitAMO as toolbox
- Generates reports on similarity values produced by different matchers, execution times, number of correct, wrong, redundant suggestions



Recommendation algorithm

- Recommendation scores: F, F+E, 10F+E
- F: quality of the alignment suggestions
 - average f-measure value for the segment pairs
- E: average execution time over segment pairs, normalized with respect to number of term pairs
- Algorithm gives ranking of alignment strategies based on recommendation scores on segment pairs

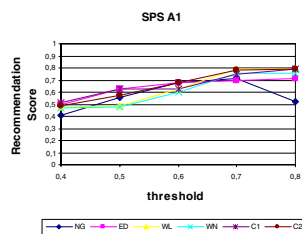


Expected recommendations for F

- Best strategies for the whole ontologies and measure F:
 1. (WL,0.8)
 2. (C1,0.8)
 3. (C2,0.8)

Results

SubG, F, SPS A1



Results

- Top 3 strategies for SubG and measure F:
 - A1: 1. (WL,0.8) (WL, 0.7) (C1,0.8) (C2,0.8)
 - A2: 1. (WL,0.8) 2. (WL,0.7) 3. (WN,0.7)
 - A3: 1. (WL,0.8) (WL, 0.7) (C1,0.8) (C2,0.8)
- Best strategy always recommended first
- Top 3 strategies often recommended
- (WL,0.7) has rank 4 for whole ontologies

Results

- Top 3 strategies for Clust and measure F:
B1: 1. (C2,0.7) 2. (ED,0.6) 3. (C2,0.6)
B2: 1. (WL,0.8) (WL, 0.7) (C1,0.8) (C2,0.8)
B3: 1. (C1,0.8) (ED,0.7) 3. (C1,0.7) (C2,0.7) (WL,0.7) (WN,0.7)
- Top strategies often recommended, but not always
- (WL,0.7) (C1,0.7) (C2,0.7) ranked 4,5,6 for whole ontologies

Results

- SubG gives better results than Clust
- Results improve when number of segments is increased
- 10F+E similar results as F
- F+E
 - WordNet gives lower ranking
 - Runtime environment has influence

Part II – Ontology Alignment

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Recommending ontology alignment strategies
- Current Issues

Current issues

- Systems and algorithms
 - Complex ontologies
 - Use of instance-based techniques
 - Alignment types (equivalence, is-a, ...)
 - Complex alignments (1-n, m-n)
 - Connection ontology types – alignment strategies

Current issues

- Evaluations
 - Need for Golden standards
 - Systems available, but not always the alignment algorithms
 - Evaluation measures
- Recommending 'best' alignment strategies

Further reading

Starting points for further studies

Further reading ontologies

- KnowledgeWeb (<http://knowledgeweb.semanticweb.org/>) and its predecessor OntoWeb (<http://ontoweb.aifb.uni-karlsruhe.de/>)
- Lambrix, Tan, Jakoniene, Strömbäck. Biological Ontologies, chapter 4 in Baker, Cheung, (eds), *Semantic Web: Revolutionizing Knowledge Discovery in the Life Sciences*, 85-99, Springer, 2007. ISBN: 978-0-387-48436-5.

(general about ontologies)

- Lambrix, Towards a Semantic Web for Bioinformatics using Ontology-based Annotation, *Proceedings of the 14th IEEE International Workshops on Enabling Technologies: Infrastructures for Collaborative Enterprises*, 3-7, 2005. Invited talk.

(ontologies for semantic web)

- OWL, <http://www.w3.org/TR/owl-features/> , <http://www.w3.org/2004/OWL/>

Further reading ontology alignment

- <http://www.ontologymatching.org>
(plenty of references to articles and systems)
- Ontology alignment evaluation initiative: <http://oaei.ontologymatching.org>
(home page of the initiative)
- Euzenat, Shvaiko, *Ontology Matching*, Springer, 2007.
- Lambrix, Tan, SAMBO – a system for aligning and merging biomedical ontologies, *Journal of Web Semantics*, 4(3):196-206, 2006.
(description of the SAMBO tool and overview of evaluations of different matchers)
- Lambrix, Tan, A tool for evaluating ontology alignment strategies, *Journal on Data Semantics*, VIII:182-202, 2007.
(description of the KitAMO tool for evaluating matchers)

Further reading ontology alignment

- Chen, Tan, Lambrix, Structure-based filtering for ontology alignment, *IEEE WETICE workshop on semantic technologies in collaborative applications*, 364-369, 2006.

(double threshold filtering technique)

- Tan H, Lambrix P, 'A method for recommending ontology alignment strategies', *International Semantic Web Conference*, 494-507, 2007.

Ehrig M, Staab S, Sure Y, 'Bootstrapping ontology alignment methods with APFEL', *International Semantic Web Conference*, 186-200, 2005.

Mochol M, Jentzsch A, Euzenat J, 'Applying an analytic method for matching approach selection', *International Workshop on Ontology Matching*, 2006.

(recommendation of alignment strategies)