# Ontology Alignment

### **Ontology Alignment**

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Ontology alignment challenges

### **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- cvtokine 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**

#### **GENE ONTOLOGY (GO)** SIGNAL-ONTOLOGY (SigO) immune response **Immune Response** i- acute-phase response i- Allergic Response Antigen Processing and Presentation i- anaphylaxis i- antigen presentation i-B Cell Activation i- antigen processing B Cell Development i- Complement Signaling i- cellular defense response i- cytokine metabolism • synonym complement activation i- Cytokine Response i- cytokine biosynthesis synonym cytokine production i- Immune Suppression i- Inflammation p- regulation of cytokine. i- Intestinal Immunity biosynthesis i- Leukotriene Response i- Leukotriene Metabolism Natural Killer Cell Response i-B-cell activation **T** Cell Activation i- B-cell differentiation i- T Cell Development i- B-cell proliferation i- T Cell Selection in Thymus i- cellular defense response i- T-cell activation i- activation of natural killer

# **Ontologies with overlapping information**

- Use of multiple ontologies
  - custom-specific ontology + standard ontology
  - □ different views over same domain
  - overlapping domains
- Bottom-up creation of ontologies experts can focus on their domain of expertise

→ important to know the inter-ontology relationships

#### GENE ONTOLOGY (GO)

#### SIGNAL-ONTOLOGY (SigO)

Immune Response

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 differentiationi- B-cell proliferationi- cellular defense response

i- T-cell activationi- activation of natural killer cell activity

. . .

i- Allergic Response
i- Antigen Processing and Presentation
i- B Cell Activation
i- B Cell Development
i- Complement Signaling synonym complement activation
i- Cytokine Response
i- Immune Suppression
i- Inflammation
i- Intestinal Immunity
i- Leukotriene Response
i- Leukotriene Metabolism
i- Natural Killer Cell Response
i- T Cell Activation
i- T Cell Development

i- T Cell Selection in Thymus

### **Ontology Alignment**



Defining the relations between the terms in different ontologies

### **Ontology Alignment**

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Ontology alignment challenges

### An Alignment Framework



### 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

# Preprocessing

### Preprocessing

### For example,

- Selection of features
- Selection of search space

## Matchers

### **Matcher Strategies**

- Strategies based on linguistic matching
- Structure-based strategies
- Constraint-bas
- Instance-based
- Use of auxiliar



### Edit distance

- Number of deletions, insertions, substitutions required to transform one string into another
- $\Box$  aaaa  $\rightarrow$  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



# Propagation of similarity valuesAnchored matching



# Propagation of similarity valuesAnchored matching



# Propagation of similarity valuesAnchored matching



### **Matcher Strategies**

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



### **Matcher Strategies**

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



- Similarities between data types
   Similarities based on cordinalities
- Similarities based on cardinalities

## **Matcher Strategies**

- Strategies based on linguisti
- Structure-based strategies
- Constraint-based approached
- Instance-based strategies





- Instance-based
- Use life science literature as instances

### Structure-based extensions

# Learning matchers – instancebased 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
  - Naive 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

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

### **Matcher Strategies**

- Strategies based linguist
- Structure-based strategie
- Constraint-based approa
- Instance-based strategies
- Use of auxiliary information



### Use of WordNet

- Use WordNet to find synonyms
- Use WordNet to find ancestors and descendants in the isa 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

	linguistic	structure	constraints	instances	auxiliary
$\mathbf{ArtGen}$	name	parents, children		domain	WordNet
				specific	
				documents	
ASCO	name,	parents, children,			WordNet
	label	siblings,			
	description	path from root			
Chimaera	name	parents, children			
• • • • • • • • • • • • • • • • • • • •		p ======; ========			
FCA-Merge	name			domain	
8-				specific	
				documents	
FOAM	name	parents children	equivalence	documents	
FOAM	label	parents, cinturen	equivalence		
CLUE	label	111 . 1 1		• - +	
GLUE	name	neighborhood		instances	
HOONE					*** 15* .
HCONE	name	parents, children			WordNet
IF-Map				instances	a reference
					ontology
iMapper		leaf, non-leaf,	domain,	instances	WordNet
		children,	range		
		related node			
OntoMapper		parents, children		documents	
(Anchor-)	name	direct graphs			
PROMPT					
SAMBO	name,	is-a and part-of,		domain	WordNet,
	synonym	descendants		specific	UMLS
		and ancestors		documents	
S-Match	label	path from root	semantic		WordNet
		* · · ·	relations		
			codified		
			in labels		
			111 100010		

# Combinations

### **Combination Strategies**

- Usually weighted sum of similarity values of different matchers
- Maximum of similarity values of different matchers



### **Filtering techniques**

Threshold filtering

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



# **Filtering techniques**

### Double threshold filtering

- (1) Pairs of concepts with similarity higher than or equal to **upper** threshold are alignment suggestions
- (2) 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

			Align Concept in mouse and h	iuman	
matchers:	1.0NGram1.0TermBasic1.0TermWN1.0UMLSM1.0Naive Bayes	single threshold: double threshold:	0.6 • upper 0.6 lower 0.4 •	weighted-sum combination maximum-based combination	use preprocessed data
Start Computation	Finish Computation	Interrupt Computation	interrupt at: 1000	J	

comments to sambo@ida.liu.se

# **Example alignment system SAMBO** – suggestion mode

nose_MA	nose_MeSH		
nasal_cavity_epithelium definition: MA:0001324 synonym: nasal mucosa part-of: nasal_cavity	nasal_mucosa definition: MESH:A.04.531.520 synonym: nasal epithelium part-of:		
nasal_cavity_epithelium nasal_mucosa			
= Equiv. Concepts ≤ Sub-Concept ≥	Super-Concept << Undo >> Skip to Next		

# Example alignment system SAMBO – manual mode



# **Ontology Alignment**

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Ontology alignment challenges

# **Evaluation measures**

Precision:

# correct mapping suggestions

# mapping suggestions

• Recall:

# correct mapping suggestions

# correct mappings

F-measure: combination of precision and recall

## Ontology Alignment Evaluation Initiative

http://oaei.ontologymatching.org/

- Since 2004
- Evaluation of systems
- Different tracks (2017)
  - anatomy, conference, large biomedical ontologies, disease and phenotype
  - multilingual: multifarm (9 languages)
  - □ process model
  - □ interactive
  - □ instance
  - □ link discovery for spatial data

Evaluation measures
 Precision/recall/f-measure
 recall of non-trivial mappings

□ full / partial golden standard

- 17 systems participated
  - □ benchmark (13)
    - ASMOV: p = 0.95, r = 0.90
  - $\square$  anatomy (11)
    - AOAS: f = 0.86, r+ = 0.50
    - SAMBO: f =0.81, r+ = 0.58
  - $\Box$  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

 $\Box$  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 2008 – anatomy track**

## Align

- □ Mouse anatomy: 2744 terms
- NCI-anatomy: 3304 terms
- □ Mappings: 1544 (of which 934 'trivial')

#### Tasks

- □ 1. Align and optimize f
- □ 2-3. Align and optimize p / r
- 4. Align when partial reference alignment is given and optimize f

# **OAEI 2008 – anatomy track#1**

- 9 systems participated
- SAMBO
  - □ p=0.869, r=0.836, r+=0.586, f=0.852
- SAMBOdtf
  - □ p=0.831, r=0.833, r+=0.579, f=0.832
- Use of TermWN and UMLS

# **OAEI 2008 – anatomy track#4**

Can we use given mappings when computing suggestions?

→ partial reference alignment given with all trivial and 50 non-trivial mappings

#### SAMBO

□ p=0.636 $\rightarrow$ 0.660, r=0.626 $\rightarrow$ 0.624, f=0.631 $\rightarrow$ 0.642

- SAMBOdtf
  - □  $p=0.563 \rightarrow 0.603$ ,  $r=0.622 \rightarrow 0.630$ ,  $f=0.591 \rightarrow 0.616$

(measures computed on non-given part of the reference alignment)

- 11 systems
- Anatomy:
  - □ best system f=0.943, p=0.95, r=0.936, r+=0.832, 47 seconds
  - □ 5 systems produce coherent mappings

## OAEI Anatomy Track 2007-2016\*

#### Components

- Almost all systems implement preprocessing, matchers, combination, filtering components
- Debugging component and GUI rarely implemented
- Matching strategies
  - Variety of string-based strategies
  - □ Most often string and structured-based strategies
- Use of background knowledge
  - □ Almost all systems use sources of background knowledge

<sup>\*</sup> Dragisic Z, Ivanova V, Li H, Lambrix P, <u>Experiences from the Anatomy track in the</u> <u>Ontology Alignment Evaluation Initiative</u>, *Journal of Biomedical Semantics* 8:56, 2017.

# Evaluation of algorithms

# Cases

#### □ GO vs. SigO

GO: 70 terms	SigO: 15 terms
GO-immune defense	SigO-immune defense

GO: 60 terms	
GO-behavior	

SigO: 10 terms SigO-behavior



## **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)

Terminological matchers



Basic learning matcher (Naïve Bayes)



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

Domain matcher (using UMLS)



- 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 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



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



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

# Complementary evaluation

Alignment cubes

- Interactive visualization of alignments
- Region-level, mapping level
- Missing mappings
- Often found mappings
- http://www.ida.liu.se/~patla00/research/AlignmentCubes/

#### Alignment cubes



# **Ontology Alignment**

- Ontology alignment
- Ontology alignment strategies
- Evaluation of ontology alignment strategies
- Ontology alignment challenges

# Challenges

- Large-scale matching evaluation
- Efficiency of matching techniques
  - parallellization
  - distribution of computation
  - approximation of matching results (not complete)
  - modularization of ontologies
  - optimization of matching methods

# Challenges

Matching with background knowledge
 partial alignments
 reuse of previous matches
 use of domain-specific corpora
 use of domain-specific ontologies

Matcher selection, combination and tuning
 recommendation of algorithms and settings

# Challenges

User involvement
 visualization
 user feedback

- Explanation of matching results
- Social and collaborative matching
- Alignment management: infrastructure and support

# **Further reading**

Starting points for further studies

<u>http://www.ontologymatching.org</u>
 (plenty of references to articles and systems)

 Ontology alignment evaluation initiative: <u>http://oaei.ontologymatching.org</u> (home page of the initiative)

- Euzenat, Shvaiko, *Ontology Matching*, Springer, 2007.
- Shvaiko, Euzenat, Ontology Matching: state of the art and future challenges, *IEEE Transactions on Knowledge and Data Engineering* 25(1):158-176, 2013.
- Dragisic Z, Ivanova V, Li H, Lambrix P, <u>Experiences from the Anatomy track in the Ontology Alignment Evaluation Initiative</u>, *Journal of Biomedical Semantics* 8:56, 2017.

Systems at LiU / IDA / ADIT

 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)

- Lambrix P, Kaliyaperumal R, <u>A Session-based Ontology Alignment Approach</u> <u>enabling User Involvement</u>, *Semantic Web Journal* 8(2):225-251, 2017.
- Ivanova V, Bach B, Pietriga E, Lambrix P, <u>Alignment Cubes: Towards Interactive Visual Exploration and Evaluation of Multiple Ontology Alignments</u>, 16th International Semantic Web Conference, 400-417, 2017.

 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, Lambrix, A method for recommending ontology alignment strategies, *International Semantic Web Conference*, 494-507, 2007.
   Ehrig, Staab, Sure, Bootstrapping ontology alignment methods with APFEL, *International Semantic Web Conference*, 186-200, 2005.
   Mochol, Jentzsch, Euzenat, Applying an analytic method for matching approach selection, *International Workshop on Ontology Matching*, 2006.
   (recommendation of alignment strategies)
- Lambrix, Liu, Using partial reference alignments to align ontologies, *European Semantic Web Conference*, 188-202, 2009.
  (use of partial alignments in ontology alignment)

Lambrix, Strömbäck, Tan, Information integration in bioinformatics with ontologies and standards, chapter 8 in Bry, Maluszynski (eds), *Semantic Techniques for the Web*, Springer, 2009. ISBN: 978-3-642-04580-6.

(largest overview of systems)

# Ontology Debugging
## Defects in ontologies

- Syntactic defects
  - E.g. wrong tags or incorrect format
- Semantic defects
  - E.g. unsatisfiable concepts, incoherent and inconsistent ontologies
- Modeling defects
  - □ E.g. wrong or missing relations

## Example - incoherent ontology

#### Example: DICE ontology

 Brain ⊑ CentralNervousSystem п BodyPart п ∃systempart.NervousSystem п ∃ region.HeadAndNeck п ∀region.HeadAndNeck

A brain is a central nervous system and a body part which has a system part that is a nervous system and that is in the head and neck region.

CentralNervousSystem ⊑ NervousSystem

A central nervous system is a nervous system.

BodyPart ⊑¬NervousSystem

Nothing can be at the same time a body part and a nervous system.

#### Slide from G. Qi

## Example - inconsistent ontology

#### Example from Foaf:

- Person(timbl)
- Homepage(timbl, <u>http://w3.org/</u>)
- Homepage(w3c, <u>http://w3.org/</u>)
- Organization(w3c)
- InverseFunctionalProperty(Homepage)
- DisjointWith(Organization, Person)
- Example from OpenCyc:
  - ArtifactualFeatureType(PopulatedPlace)
  - ExistingStuffType(PopulatedPlace)
  - DisjointWith(ExistingObjectType,ExistingStuffType)
  - ArtifactualFeatureType 
    ExistingObjectType

## Example - missing is-a relations

- In 2008 Ontology Alignment Evaluation Initiative (OAEI) Anatomy track, task 4
  - □ Ontology MA : Adult Mouse Anatomy Dictionary (2744 concepts)
  - Ontology NCI-A : NCI Thesaurus anatomy (3304 concepts)
  - □ 988 mappings between MA and NCI-A
    - 121 missing is-a relations in MA
    - 83 missing is-a relations in NCI-A

## Influence of missing structure

Ontology-based querying.



## Influence of missing structure

#### Incomplete results from ontology-based queries

Publed.gov U.S. National Library of Medicine National Institutes of Health	Search: PubMed  Limits  Scleral Diseases" [MeSH]	Advanced search Help Search Clear
Medical Subject Headings (MeSH) All MeSH Categories Diseases Category		return 1617 articles return 695 articles 57% results are missed
<ul> <li>Eye Diseases</li> <li>Scleral Diseas</li> <li>Scleritis</li> </ul>	<u>es</u>	

## Defects in ontologies and ontology networks

- Ontologies and ontology networks with defects, although often useful, also lead to problems when used in semantically-enabled applications.
- → Wrong conclusions may be derived or valid conclusions may be missed.

#### **Overview of debugging approach**



## Debugging semantic defects

#### Example : an Incoherent Ontology

Consider the following TBox  $\mathcal{T}^*$ , where A, B and C are primitive and  $A_1, \ldots, A_7$  defined concept names:

 $\begin{array}{ll} ax_1:A_1 \stackrel{.}{\sqsubseteq} \neg A \sqcap A_2 \sqcap A_3 & ax_2:A_2 \stackrel{.}{\sqsubseteq} A \sqcap A_4 \\ ax_3:A_3 \stackrel{.}{\sqsubseteq} A_4 \sqcap A_5 & ax_4:A_4 \stackrel{.}{\boxminus} \forall s.B \sqcap C \\ ax_5:A_5 \stackrel{.}{\sqsubseteq} \exists s. \neg B & ax_6:A_6 \stackrel{.}{\sqsubseteq} A_1 \sqcup \exists r.(A_3 \sqcap \neg C \sqcap A_4) \\ ax_7:A_7 \stackrel{.}{\sqsubseteq} A_4 \sqcap \exists s. \neg B \end{array}$ 



The ontology is incoherent!

The set of unsatisfiable concepts are :  $\{A_1, A_3, A_6, A_7\}$ .



What are the root causes of these defects?

#### **Explain the Semantic Defects**

• We need to identify the sets of axioms which are necessary for causing the logic contradictions.



• For example, for the unsatisfiable concept "*A*<sub>1</sub>", there are two sets of axioms.

 $ax_1:A_1 \sqsubseteq \neg A \sqcap A_2 \sqcap A_3$  $ax_2:A_2 \sqsubseteq A \sqcap A_4$ 

$$ax_{1}:A_{1} \stackrel{\models}{=} \neg A \sqcap A_{2} \sqcap A_{3}$$
$$ax_{3}:A_{3} \stackrel{\models}{=} A_{4} \sqcap A_{5}$$
$$ax_{4}:A_{4} \stackrel{\models}{=} \forall s.B \sqcap C$$
$$ax_{5}:A_{5} \stackrel{\models}{=} \exists s. \neg B$$

#### Minimal Unsatisfiability Preserving Sub-TBoxes (MUPS)

**Definition 1** Let A be a concept which is unsatisfiable in a TBox  $\mathcal{T}$ . A set  $\mathcal{T}' \subseteq \mathcal{T}$  is a *minimal unsatisfiability-preserving sub-TBox (MUPS)* of  $\mathcal{T}$  if

- A is unsatisfiable in  $\mathcal{T}'$ , and
- A is satisfiable in every sub-TBox  $\mathcal{T}'' \subset \mathcal{T}'$ .

We will abbreviate the set of MUPS of  $\mathcal{T}$  and A by  $mups(\mathcal{T}, A)$ .  $mups(\mathcal{T}^*, A_1) = \{\{ax_1, ax_2\}, \{ax_1, ax_3, ax_4, ax_5\}\}$ 

• The MUPS of an unsatisfiable concept imply the solutions for repairing.

 $\rightarrow$  Remove at least one axiom from each axiom set in the MUPS

#### Example

$$mups(\mathcal{T}^*, A_1) = \{\{ax_1, ax_2\}, \{ax_1, ax_3, ax_4, ax_5\}\}$$
  

$$mups(\mathcal{T}^*, A_3) = \{\{ax_3, ax_4, ax_5\}\}$$
  

$$mups(\mathcal{T}^*, A_6) = \{\{ax_1, ax_2, ax_4, ax_6\},$$
  

$$\{ax_1, ax_3, ax_4, ax_5, ax_6\}\}$$
  

$$mups(\mathcal{T}^*, A_7) = \{\{ax_4, ax_7\}\}$$

• Possible ways of repairing all the unsatisfiable concepts in the ontology:

$$\{ax_1, ax_3, ax_4\}$$



How to represent all these possibilities?

#### Minimal Incoherence Preserving Sub-TBox (MIPS)

**Definition 2** Let  $\mathcal{T}$  be an incoherent TBox. A TBox  $\mathcal{T}' \subseteq \mathcal{T}$  is a minimal incoherencepreserving sub-TBox (MIPS) of  $\mathcal{T}$  if

- $\mathcal{T}'$  is incoherent, and
- every sub-TBox  $\mathcal{T}'' \subset \mathcal{T}'$  is coherent.

$$\begin{split} mups(\mathcal{T}^*, A_1) &= \{ \{ax_1, ax_2\}, \{ax_1, ax_3, ax_4, ax_5\} \} \\ mups(\mathcal{T}^*, A_3) &= \{ \{ax_3, ax_4, ax_5\} \} \\ mups(\mathcal{T}^*, A_6) &= \{ \{ax_1, ax_2, ax_4, ax_6\}, \\ \{ax_1, ax_3, ax_4, ax_5, ax_6\} \} \\ mups(\mathcal{T}^*, A_7) &= \{ \{ax_4, ax_7\} \} \end{split}$$

We will abbreviate the set of MIPS of  $\mathcal{T}$  by  $mips(\mathcal{T})$ . For  $\mathcal{T}^*$  we get three MIPS:

$$mips(\mathcal{T}^*) = \{ \{ax_1, ax_2\}, \{ax_3, ax_4, ax_5\}, \{ax_4, ax_7\} \}$$

A possible repairing is  $\{ax_i\} \cup \{ax_j\} \cup \{ax_k\}$ , where

- $ax_i \in \{ax_1, ax_2\}$
- $ax_j \in \{ax_3, ax_4, ax_5\}$
- $ax_k \in \{ax_4, ax_7\}$

# Completing the is-a structure of ontologies



#### **Repairing actions:**

- {Endocarditis  $\sqsubseteq$  PathologicalPhenomenon, GranulomaProcess  $\doteq$  NonNormalProcess}
- {Carditis  $\doteq$  CardioVascularDisease, GranulomaProcess  $\doteq$  PathologicalProcess}
- {Carditis  $\sqsubseteq$  Fracture, GranulomaProcess  $\doteq$  NonNormalProcess}

## **Description logic EL**

#### Concepts

Atomic concept	Α
Universal concept	T
Intersection of concepts	СпD
Existential restriction	∃r.C

#### Terminological axioms: equivalence and subsumption

## Generalized Tbox Abduction Problem – GTAP(**T**,**C**,Or,M)

- Given
  - □**T** a Tbox in EL
  - □C- a set of atomic concepts in T
  - $\Box M = \{Ai \subseteq Bi\}_{i=1..n} and \forall i:1..n: Ai, Bi \in \boldsymbol{C}$
  - $\Box \text{ Or: } \{\text{Ci} \subseteq \text{Di} \mid \text{Ci, Di} \in \textbf{C}\} \rightarrow \{\text{true, false}\}$

Find

□ S = {E<sub>i</sub> ⊆ F<sub>i</sub>}<sub>i=1..k</sub> such that  $\forall$  i:1..k: E<sub>i</sub>, F<sub>i</sub> ∈ C and Or(E<sub>i</sub> ⊆ F<sub>i</sub>) = true and T U S is consistent and T U S |= M

## GTAP - example



 $C = \{$  GranulomaProcess, CardioVascularDisease, PathologicalPhenomenon, Fracture, Endocarditis, Carditis, InflammationProcess, PathologicalProcess, NonNormalProcess} $\}$ 

 $T = \{ \text{ GranulomaProcess } \sqsubseteq \top, \text{ hasAssociatedProcess } \trianglerighteq \top \times \top, \\ \text{CardioVascularDisease } \trianglerighteq \text{ PathologicalPhenomenon, Fracture } \textcircled{} \text{ PathologicalPhenomenon,} \\ \exists \text{hasAssociatedProcess.PathologicalProcess } \trianglerighteq \text{ PathologicalPhenomenon,} \\ \text{Endocarditis } \underrightarrow{} \text{Carditis, Endocarditis } \underrightarrow{} \exists \text{hasAssociatedProcess.InflammationProcess,} \\ \text{PathologicalProcess } \end{bmatrix}$ 

 $M = \{$  Endocarditis  $\doteq$  PathologicalPhenomenon, GranulomaProcess  $\doteq$  NonNormalProcess  $\}$ 

#### Preference criteria

#### There can be many solutions for GTAP



#### Preference criteria

There can be many solutions for GTAP



Not all are equally interesting.

#### More informative

- Let S and S' be two solutions to GTAP(T,C,Or,M). Then,
- S is more informative than S' iff  $\mathbf{T} \cup S \models S'$  but not  $\mathbf{T} \cup S' \models S$
- S is equally informative as S' iff  $\mathbf{T} \cup S \models S'$  and  $\mathbf{T} \cup S' \models S$

## More informative

#### Blue' solution is more informative than 'green' solution



## Semantic maximality

A solution S to GTAP(T,C,Or,M) is semantically maximal iff there is no solution S' which is more informative than S.



## Subset minimality

A solution S to GTAP(T,C,Or,M) is subset minimal iff there is no proper subset S' of S that is a solution.



Combining with priority for semantic maximality

A solution S to GTAP(T,C,Or,M) is maxmin optimal iff S is semantically maximal and there is no other semantically maximal solution that is a proper subset of S.



Combining with priority for subset minimality

A solution S to GTAP(T,C,Or,M) is minmax optimal iff S is subset minimal and there is no other subset minimal solution that is more informative than S.



Combining with equal preferences

- A solution S to GTAP(T,C,Or,M) is skyline optimal iff there is no other solution that is a proper subset of S and that is equally informative than S.
  - All subset minimal, minmax optimal and maxmin optimal solutions are also skyline optimal solutions.
  - Semantically maximal solutions may or may not be skyline optimal.

#### Preference criteria - conclusions

- In practice it is not clear how to generate maxmin or semantically maximal solutions (the preferred solutions)
- Skyline optimal solutions are the next best thing and are easy to generate

## Approach

#### Input

- □ Normalized EL TBox
- Set of missing is-a relations (correct according to the domain)
- Output a skyline-optimal solution to GTAP
- Iteration of three main steps:
  - Creating solutions for individual missing is-a relations
  - Combining individual solutions
  - Trying to improve the result by finding a solution which introduces additional new knowledge (more informative)

## Intuition 1



#### Intuitions 2/3



#### Example – repairing single is–a relation



#### Example – repairing single is–a relation



GranulomaProcess GranulomaProcess GranulomaProcess

## Algorithm - Repairing multiple is-a relations

- Combine solutions for individual missing is-a relations
- Remove redundant relations while keeping the same level of informativness
- Resulting solution is a skyline optimal solution

{InflammationProcess  $\sqsubseteq$  PathologicalProcess, Carditis  $\doteq$  CardioVascularDisease, GranulomaProcess  $\doteq$  PathologicalProcess}

## Algorithm – improving solution

- Solution S from previous step may contain relations which are not derivable from the ontology.
- These can be seen as new missing is-a relations.
- We can solve a new GTAP problem: GTAP(T U S, C, Or, S)
### Example – improving solutions



#### $GranulomaProcess \stackrel{.}{\sqsubseteq} InflammationProcess$

{InflammationProcess  $\sqsubseteq$  PathologicalProcess, Carditis  $\doteq$  CardioVascularDisease, GranulomaProcess  $\doteq$  InflammationProcess}

# Algorithm properties

#### Sound

Skyline optimal solutions

# Experiments

Two use-cases

Case 1: given missing is-a relations
 AMA and a fragment of NCI-A ontology – OAEI 2013

- AMA (2744 concepts) 94 missing is-a relations
  → 3 iterations, 101 in repairing (47 additional new knowledge)
- NCI-A (3304 concepts) 58 missing is-a relations
  → 3 iterations, 54 in repairing (10 additional new knowledge)
- Case 2: no given missing is-a relations Modified BioTop ontology
  - Biotop (280 concepts, 42 object properties) randomly choose is-a relations and remove them: 47 'missing' → 4 iterations, 41 in repairing (40 additional new knowledge)

# **Further reading**

Starting points for further studies

### Further reading ontology debugging

Semantic defects

- Schlobach S, Cornet R. Non-Standard Reasoning Services for the Debugging of Description Logic Terminologies. 18th International Joint Conference on Artificial Intelligence - IJCAI03, 355-362, 2003.
- Schlobach S. <u>Debugging and Semantic Clarification by Pinpointing</u>. 2nd European Semantic Web Conference - ESWC05, LNCS 3532, 226-240, 2005.

## Further reading ontology debugging

Completing ontologies

- Fang Wei-Kleiner, Zlatan Dragisic, Patrick Lambrix. <u>Abduction Framework</u> for Repairing Incomplete EL Ontologies: Complexity Results and <u>Algorithms</u>. 28th AAAI Conference on Artificial Intelligence - AAAI 2014, 1120-1127, 2014.
- Lambrix P, Ivanova V, <u>A unified approach for debugging is-a structure and mappings in networked taxonomies</u>, *Journal of Biomedical Semantics* 4:10, 2013.
- Lambrix P, Liu Q, <u>Debugging the missing is-a structure within taxonomies</u> <u>networked by partial reference alignments</u>, *Data & Knowledge Engineering* 86:179-205, 2013.