

**Stefan Schulz**

University Medical Center, Freiburg, Germany

**Elena Beisswanger**

Language and Information Engineering Lab,  
Jena, Germany

**Olivier Bodenreider**

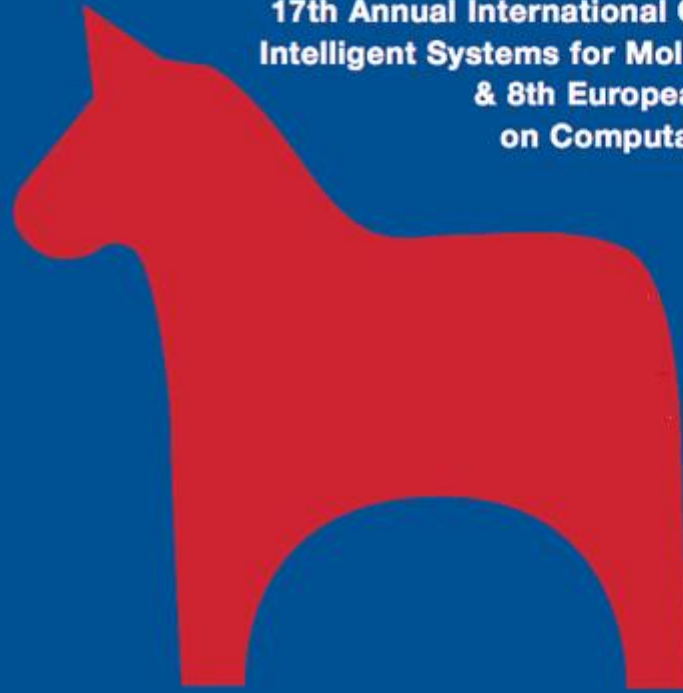
National Library of Medicine,  
Bethesda, MD, USA

**László van den Hoek**

**Erik M. van Mulligen**

Erasmus Medical Center,  
Rotterdam, The Netherlands

17th Annual International Conference on  
Intelligent Systems for Molecular Biology  
& 8th European Conference  
on Computational Biology



# **Alignment of the UMLS Semantic Network with BioTop *Methodology and Assessment***

# Ontology Alignment

- Linking two ontologies by detecting semantic correspondences between their representational units
- Types of correspondences: equivalence, subsumption, others
- Purpose of ontology alignment:
  - Creating interoperability between semantically annotated data
  - Enriching semantics
  - Cross-Validation of ontologies
- Requirements of ontology alignment:
  - comparable scope
  - comparable context
  - comparable semantic foundations

# Outline

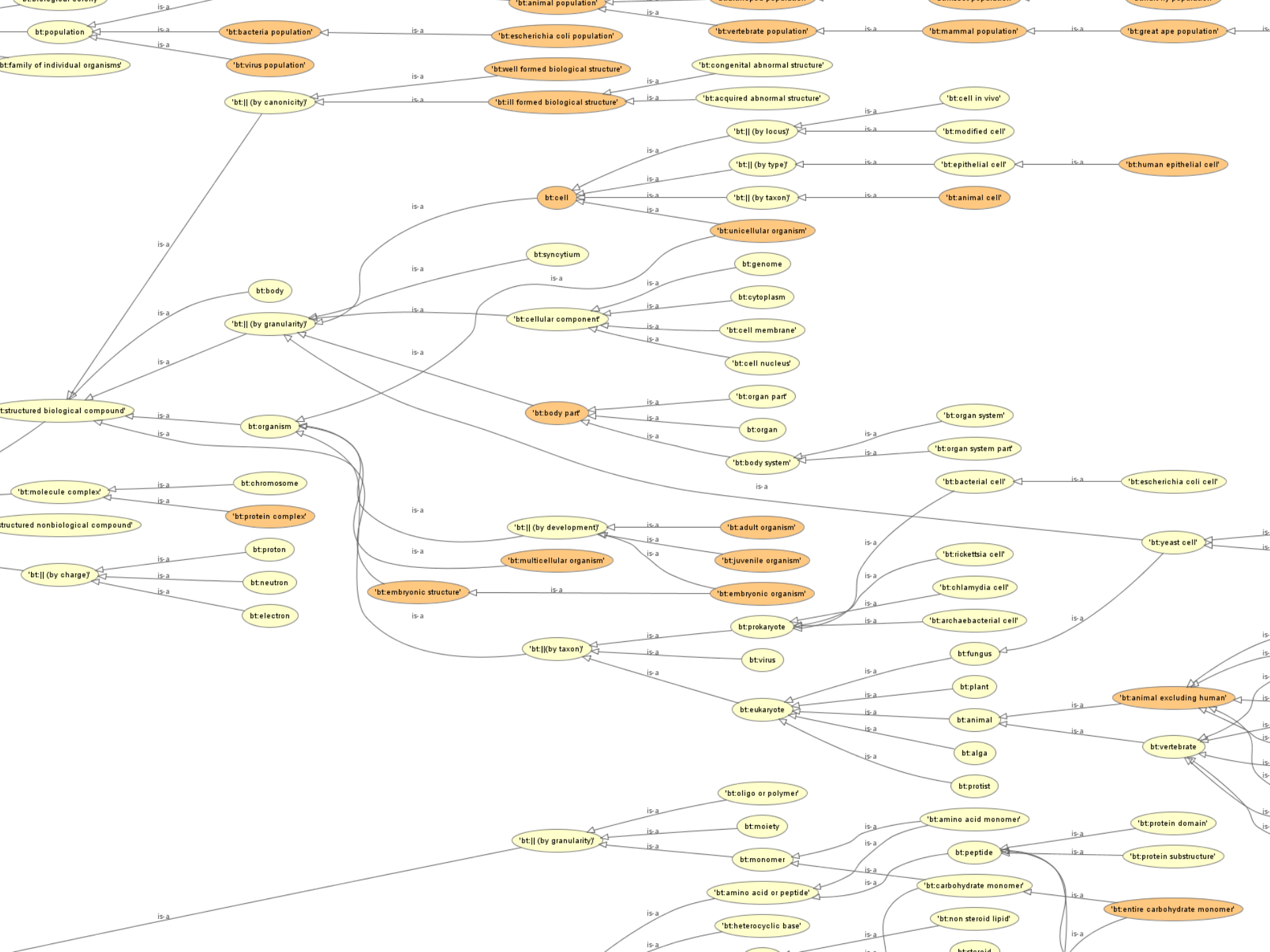
- Introduction
  - BioTop
  - UMLS SN
- Methodology
  - UMLS SN: formal redefinition
  - Interactive Mapping
- Assessment
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- Conclusion

# Outline

- **Introduction**
  - BioTop
  - UMLS SN
- Methodology
  - UMLS SN: formal redefinition
  - Interactive Mapping
- Assessment
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- Conclusion

# BioTop – a Life Science Upper Ontology

- Recent development (starting 2006, Freiburg & Jena)
- Goal: to provide formal definitions of upper-level types and relations for the biomedical domain
- Uses description logics (OWL-DL)
  - 339 classes, 60 relation types
  - 373 subclass axioms
  - 80 equivalent class axioms, 66 disjoint class axioms
- Compatible with BFO and DOLCE lite
- links to OBO ontologies
- downloadable from: <http://purl.org/biotop>

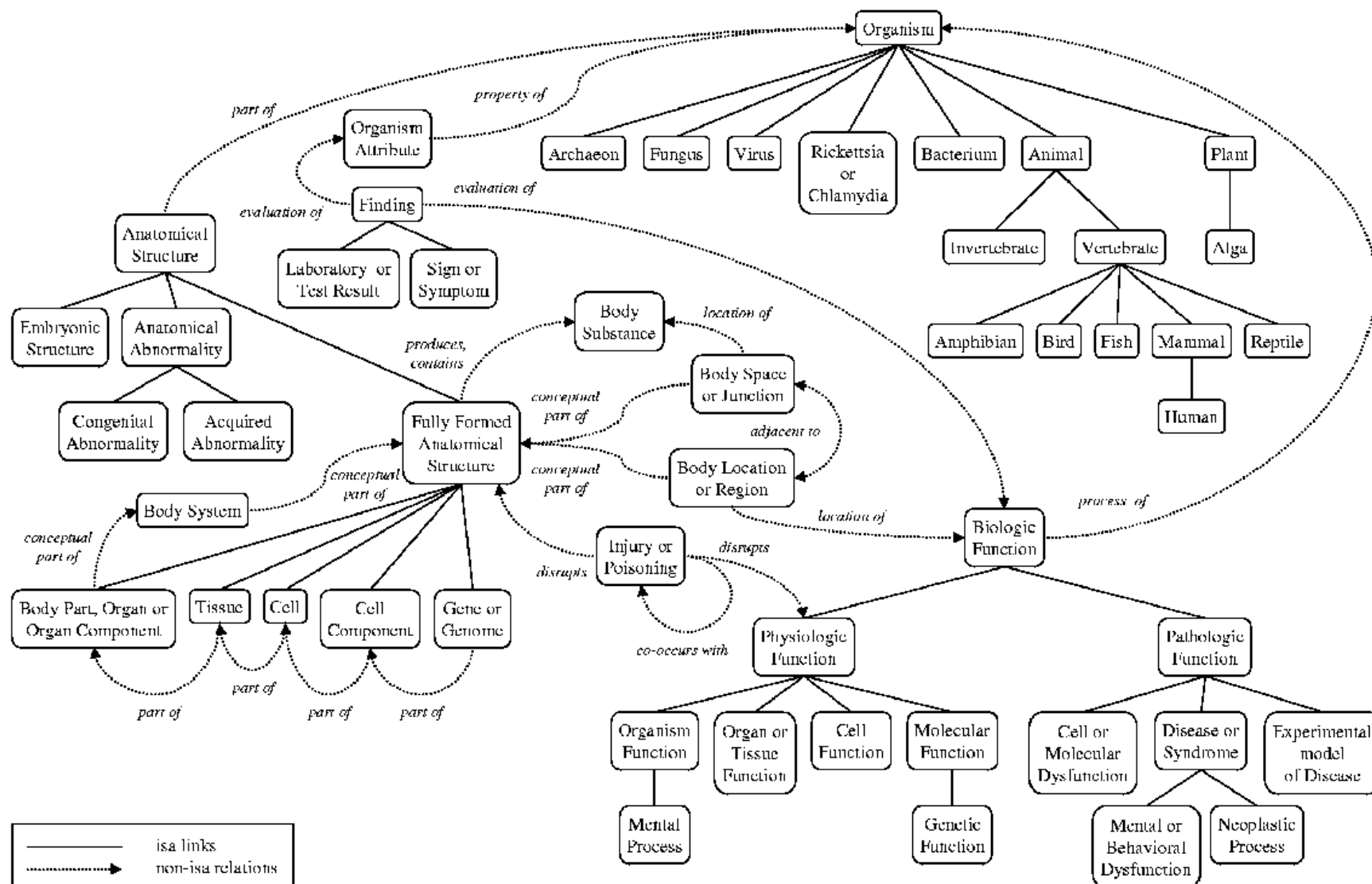


# UMLS Semantic Network (SN)

- Upper-level semantic categorization framework for all (~1 M) concepts of the **UMLS Metathesaurus**
- Tree of 135 semantic types (e.g. *Tissue*, *Diagnostic\_Procedure*)
- 53 associative relationships (e.g., *treats*, *location\_of*)
- 612 relational assertions (triples), sanctioning the domain and range of relations  
{*Tissue*; *location\_of*; *Diagnostic\_Procedure*}
- mainly unchanged in the last 20 years

Unified Medical Language System (UMLS):  
Metathesaurus links over 100  
biomedical vocabularies

# UMLS Semantic Network (SN)





# Comparison UMLS-SN - BioTop

	UMLS-SN	BioTop
Types / Classes	135	339
Relation Types	53	60 (object properties)
Axioms	612	509
Semantics	Implicit Frame-like Closed-world (?)	Explicit (description logics) Set-theoretic Open-world
Class subsumption $\sqsubseteq$	+	+
Relation subsumption $\sqsubseteq$	+	+
Domain / Range Restrictions	+	+
Relation Inheritance blocking	+	—
Full Definitions $\equiv$	—	+
Disjoint Partitions	—	+
Negations $\neg$	—	+
Existential Restrictions $\exists$	—	+
Value Restrictions $\forall$	—	+

# Outline

- Introduction
  - BioTop
  - UMLS SN
- Methodology
  - UMLS SN: formal redefinition
  - Interactive Mapping
- Assessment
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- Conclusion

# Outline

- Introduction
  - BioTop
  - UMLS SN
- **Methodology**
  - UMLS SN: formal redefinition
  - Interactive Mapping
- Assessment
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- Conclusion

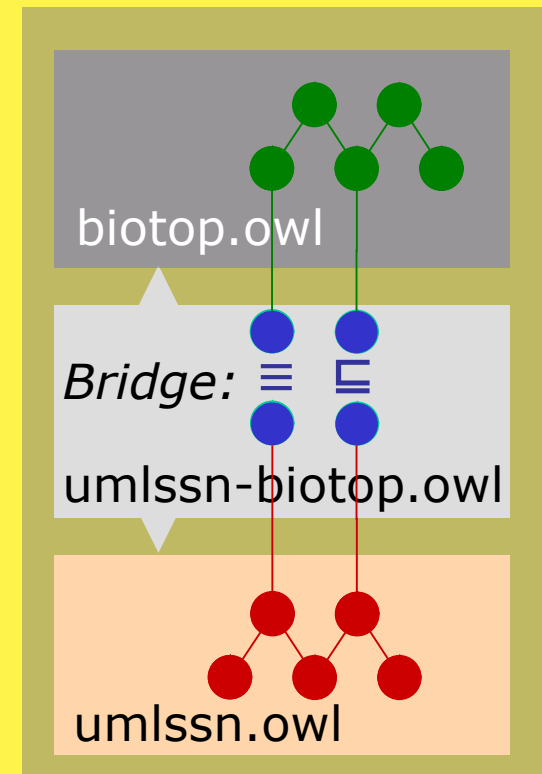
# Methodology

1. Prerequisite: provide description logics semantics to the UMLS SN:

umlssn.owl

2. Building a bridging ontology

- Subsumption  $\sqsubseteq$
- Equivalence  $\equiv$



# **Redefinition of UMLS SN semantics**

# Redefinition of UMLS SN semantics

- **Semantic Types**, e.g.: *Tissue*, *Diagnostic\_Procedure*:
  - Types extend to classes of individuals
  - subsumption hierarchies = is-a hierarchies (every instance of a child is also an instance of each parent)
  - no explicit disjoint partitions
- **Semantic Relations**, e.g.: *treats*, *location\_of*:
  - Reified as classes, **not** represented as OWL object properties
- **Triples**, e.g.:  $\{Tissue; location\_of; Diagnostic\_Procedure\}$ 
  - domain and range restrictions = value restrictions on the roles *has-domain* and *has-range*

# UMLS SN: Why SRs as classes ...

## and not OWL object properties? (I)

treats	Range	Disease	Person
		Domain	
Drug		allowed	disallowed
Physician		disallowed	allowed

$$\textit{TreatingPerson} \equiv \textit{Action} \sqcap \exists \textit{has\_domain}.\textit{Physician} \sqcap \exists \textit{has\_range}.\textit{Person} \sqcap \forall \textit{has\_domain}.\textit{Physician} \sqcap \forall \textit{has\_range}.\textit{Person}$$

$$\textit{TreatingDisease} \equiv \textit{Action} \sqcap \exists \textit{has\_domain}.\textit{Drug} \sqcap \exists \textit{has\_range}.\textit{Disease} \sqcap \forall \textit{has\_domain}.\textit{Drug} \sqcap \forall \textit{has\_range}.\textit{Disease}$$

$$\textit{Treating} \equiv \textit{TreatingPerson} \sqcup \textit{TreatingDisease}$$

# UMLS SN: Why SRs as classes ..

## and not OWL object properties? (II)

- Source Representation

“Defined not  
Inherited”

*Idea\_or\_Concept conceptual\_part\_of Behavior*

- Target Representation

*Conceptual\_part\_of\_Domain\_Idea\_Or\_Concept\_Range\_Behavior\_Rest\_Class*  $\sqsubseteq$   
*Conceptual\_part\_of*  $\sqcap$   
 $\forall$  *has\_domain*. *Idea\_Or\_Concept\_Rest\_Class*  $\sqcap$   
 $\forall$  *has\_range*. *Behavior\_Rest\_Class*

*Idea\_Or\_Concept\_Rest\_Class*  $\equiv$  *Idea\_Or\_Concept*  $\sqcap$   $\neg$  *Temporal\_Concept*  $\sqcap$   
 $\neg$  *Qualitative\_Concept*  $\sqcap$   $\neg$  *Quantitative\_Concept*  $\sqcap$   
 $\neg$  *Spatial\_Concept*  $\sqcap$   $\neg$  *Functional\_Concept*

*Behavior\_Rest\_Class*  $\equiv$  *Behavior*  $\sqcap$   $\neg$  *Individual\_Behavior*  $\sqcap$   
 $\neg$  *Social\_Behavior*



# Representation of SRs and triples

- All triples including R are defined as subclasses of R

*Affects\_Domain\_Cell\_Component\_Range\_Physiologic\_Function*  $\sqsubseteq$

*Affects*  $\sqcap \forall \text{has\_domain. Cell\_Component} \sqcap$   
 $\forall \text{has\_range. Physiologic\_Function}$

- All parents are fully defined by the union of their children

*Brings\_About*  $\equiv$  *Produces*  $\sqcup$  *Causes*

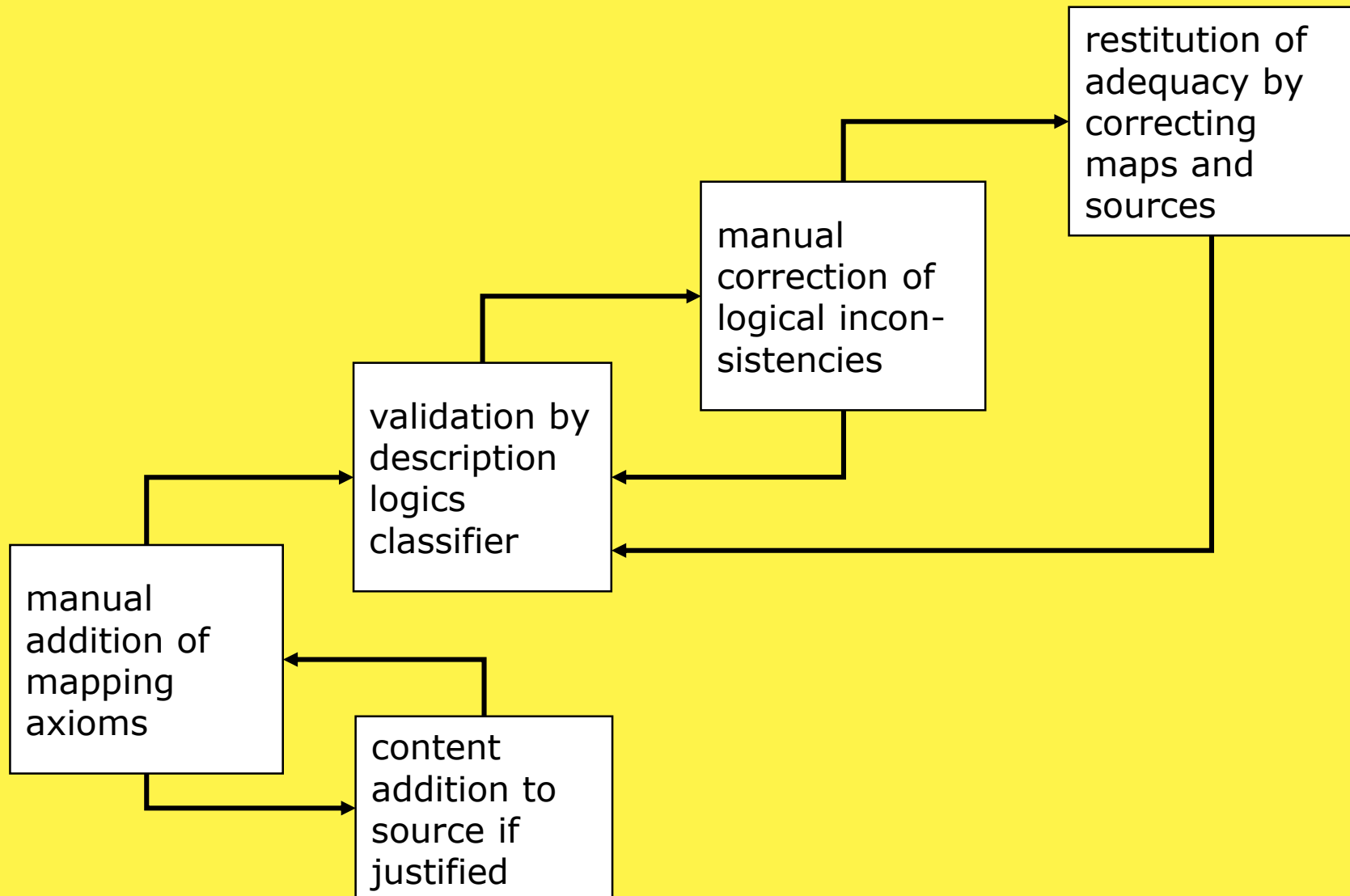
# Mapping

# Mapping

- Fully manually, using Protégé 4, consistency check with Fact++ and Pellet 1.5, supported by explanation plugin\*
- Analyzing
  - UMLS SN hierarchies and free-text definitions
  - BioTop formal and free-text definitions
- Iterative check of
  - logic consistency (DL classifier)
  - domain adequacy (analysis of new entailments)

\*(Horridge ISWC 2008)

# Mapping workflow



# Mapping of UMLS Types

- Direct Match (often after content addition to BioTop):

```
sn:Plant ≡ bt:Plant
```

- Restriction mapping:

```
sn:AnatomicalAbnormality ≡ bt:OrganismPart ⊓ ∃□  
bt:bearerOf.bt:PathologicalCondition
```

- Union:

```
sn:Gene_Or_Genome ≡ bt:Gene ⊔ bt:Genome.
```

- Out of scope

```
sn:Daily_Or_Recreational_Activity ⊆ bt:Action ⊓ ∃□ bt:hasParticipant.bt:Human
```

- No mapping

```
sn:Idea_or_concept
```

# Mapping of UMLS Relations

- Mapping of domain and range

```
sn:hasDomain ≡ bt:hasAgent  
sn:hasRange ≡ bt:hasPatient
```

- Mapping of (reified) SN relations

```
sn:Affects ≡ bt:Affecting
```

- Linkage of (reified) SN relations to BioTop relations by augmented restrictions:

```
sn:hasDomain ∨ (bt:physicalPartOf ∨ (ImmaterialPhysicalEntity ⊔ MaterialEntity)) ⊓  
sn:hasRange ∨ (bt:hasPhysicalPart ∨ (ImmaterialPhysicalEntity ⊔ MaterialEntity))
```

# Outline

- Introduction
  - BioTop
  - UMLS SN
- Methodology
  - UMLS SN: formal redefinition
  - Interactive Mapping
- Assessment
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- Conclusion

# Outline

- Introduction
  - BioTop
  - UMLS SN
- Methodology
  - UMLS SN: formal redefinition
  - Interactive Mapping
- **Assessment**
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- Conclusion



# Assessment: Cross-evaluation

- Formative evaluation of BioTop: Mapping and subsequent classification unveils hidden problems in BioTop:
  - Faulty disjointness axioms (e.g. *bt:Organic Chemical* was disjoint from *bt:Carbohydrate*)
  - ambiguities: Sequence as information entity vs. sequence as molecular structure
  - granularity mismatches:  
e.g. Chromosome as molecule

# Assessment: NE co-occurrences

- Named Entity tagging, UMLS concept pairs identified in 15 M PubMed abstracts

Semantic Type 1: UMLS ID	NE 1	Semantic Type 2: UMLS ID	NE 2
<b>Enzyme:</b> C0916840	superoxide reductase	<b>Organic_Chemical:</b> C0001992	aldehyde
<b>Finding:</b> C0883391	free testosterone index	<b>Laboratory_Procedure:</b> C0020980	immunoassay
<b>Food:</b> C1145642	sorghum	<b>Invertebrate:</b> C0009276	beetles
<b>Functional_Concept:</b> C0332240	idiopathic	<b>Pharmacologic_Substance:</b> C0011685	desipramine
<b>Functional_Concept:</b> C1510670	feeds	<b>Intellectual_Product:</b> C0023683	life table
<b>Gene_or_Genome:</b> C0087142	v-Jun	<b>Mammal:</b> C0025920	C3H
<b>Gene_or_Genome:</b> C0600449	essential gene	<b>Hazardous_or_Poisonous_Substance:</b> C0000511	4-nitroquinolone-1-oxide
<b>Geographic_Area:</b> C0027978	New Zealand	<b>Idea_or_Concept:</b> C0018741	health resources
<b>Hazardous_or_Poisonous_Substance:</b> C0036248	stx	<b>Organic_Chemical:</b> C0000967	acetal

- Expert rating with sample of co-occurrences: which are semantically related?

# Assessment: NE co-occurrences

		Expert judgment: should be related (52)	Expert judgment: Should not be related (93)
matching against SN triplets	SN: sanctioned	31	22
	SN: unsanctioned	21	71
Description logics classification	SN-BioTop: accepted	52	90
	SN-BioTop: rejected	0	3

- Using SN alone: very low agreement with expert rating
- Using SN+BioTop: very few rejections (only 3)
- Reasons:
  - false-positive rate: Expert rating done on NE (e.g. *Superoxide reductase unrelated with Aldehyde*), but system judgments at type level: *sn:Enzyme* related to *sn:Organic Chemical*
  - few rejections: DL's open world semantics

# Assessment: finding incompatible semantic types

- Each UMLS concept is categorized by one or more UMLS SN types
- 397 different SN type combinations
- Using UMLS-SN BioTop Bridge: 133 combinations inconsistent, affecting 6116 UMLS concepts
- Main reason: hidden ambiguities, e.g.

*sn:Manufactured Object*  $\cap$  *sn:HealthCareRelatedOrganization*

(e.g. *Hospital* as building vs. organization).

# Outline

- Introduction
  - BioTop
  - UMLS SN
- Methodology
  - UMLS SN: formal redefinition
  - Interactive Mapping
- Assessment
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- Conclusion

# Outline

- Introduction
  - BioTop
  - UMLS SN
- Methodology
  - UMLS SN: formal redefinition
  - Interactive Mapping
- Assessment
  - Ontology Cross-Validation
  - NE co-occurrence validation
  - UMLS SN cluster consistency
- **Conclusion**

# Conclusion

- Successful alignment between the (legacy) SN and the (novel) BioTop ontology
- Necessary: formal re-interpretation of SN
- Prospect: join large amount of data annotated by the SN with formal rigor of BioTop
- Strength: machine inference, consistency checking
- Challenge: Antagonize unwarranted effects of the open world semantics by making exhaustive use of disjoint partitions
- More use cases !

# Acknowledgements

- EC STREP project “BOOTStrep” (FP6 – 028099)
- Intramural Research Program of the National Institutes of Health (NIH), US National Library of Medicine
- Martin Boeker (Freiburg)
- Holger Stenzhorn (Freiburg)
- Anonymous Reviewers



## **Stefan Schulz**

University Medical Center, Freiburg, Germany

## **Elena Beisswanger**

Language and Information Engineering Lab,  
Jena, Germany

## **Olivier Bodenreider**

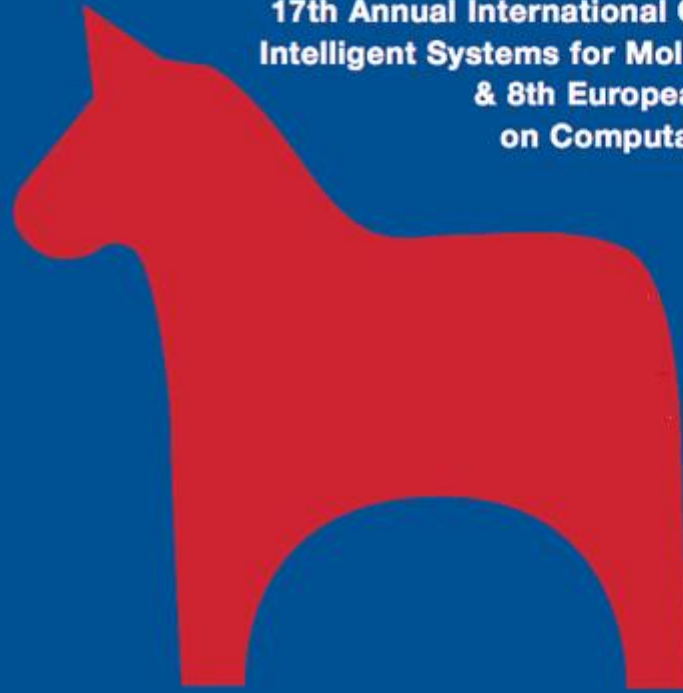
National Library of Medicine,  
Bethesda, MD, USA

## **László van den Hoek**

## **Erik M. van Mulligen**

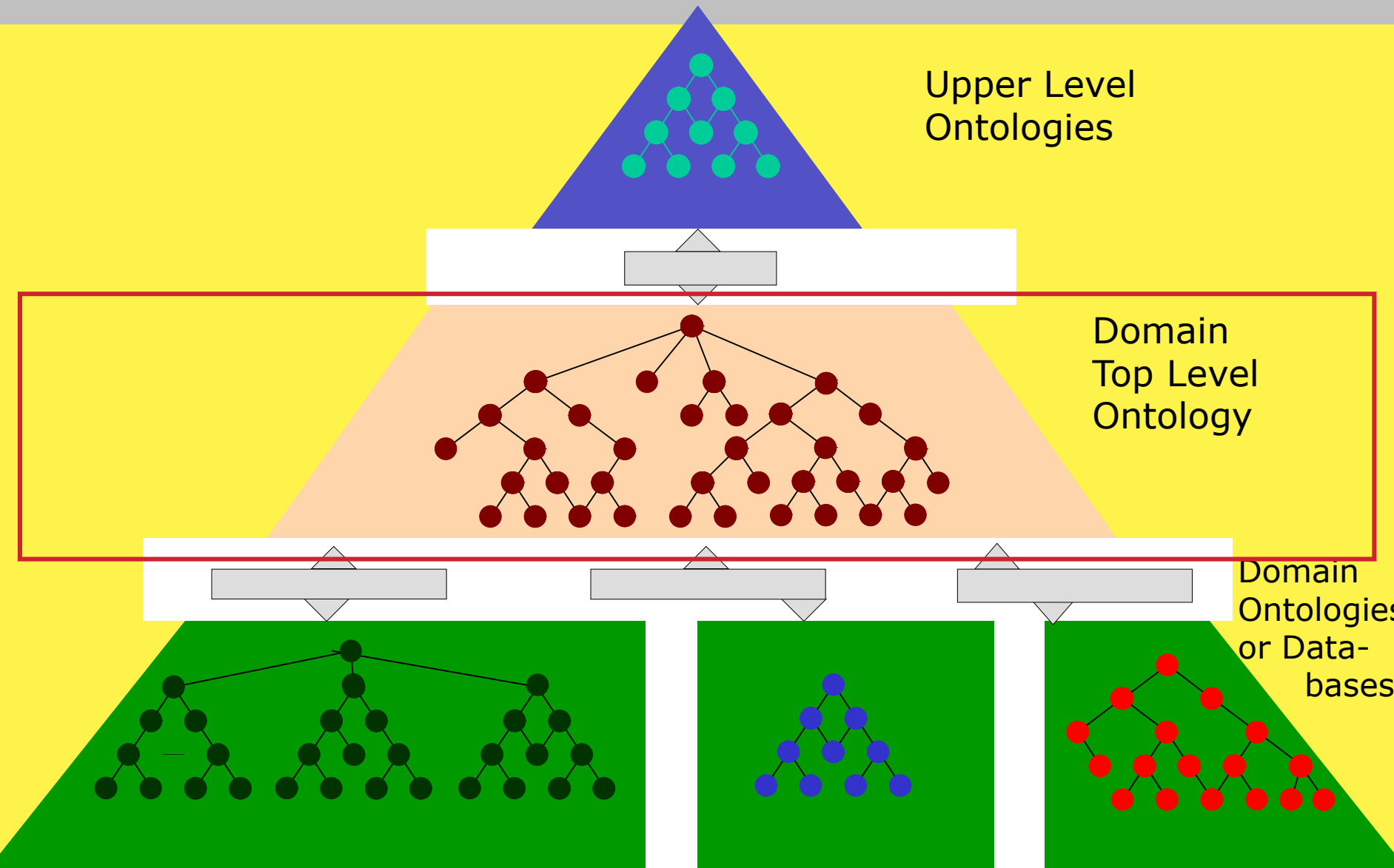
Erasmus Medical Center,  
Rotterdam, The Netherlands

17th Annual International Conference on  
Intelligent Systems for Molecular Biology  
& 8th European Conference  
on Computational Biology



# **Alignment of the UMLS Semantic Network with BioTop *Methodology and Assessment***

# Ontology Stack



# The Semantic Network of the UMLS

