

# ODLS 2016

Workshop on Ontologies and Data in Life Sciences,  
Sep 29-30, Halle (Saale), Germany

## **Ontological interpretation of biomedical database annotations**

Filipe Santana da Silva<sup>1</sup>, Ludger Jansen<sup>2</sup>,  
Fred Freitas<sup>1</sup>, Stefan Schulz<sup>3</sup>

<sup>1</sup>Centro de Informática (CIn), Universidade Federal de Pernambuco (UFPE), Recife, Brazil

<sup>2</sup>Institut für Philosophie, Universität Rostock, Germany

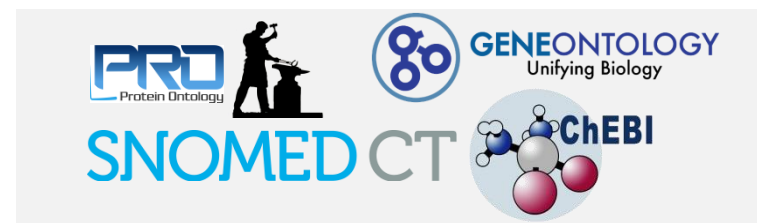
<sup>3</sup>Institut für Medizinische Informatik, Statistik und Dokumentation, Medizinische Universität  
Graz, Austria

# Biological Databases and Bio-Ontologies

- Two worlds: Bio-DBs



- Bio-ontologies



- How are they related to each other?
- Can their content be expressed by a unified model of meaning?
- Is database content of ontological nature?
- Can OWL be used as a language to express
  - bio-ontology content
  - bio-database structure
  - bio-database content
- Which is the added value?

# Biological Databases and Bio-Ontologies

## ■ Two worlds: Bio-DBs



## Bio-ontologies



- Store summarized results of laboratory experiments
- Classical database structure
- Values:
  - Numeric
  - Textual
  - Symbolic (codes from ontologies)

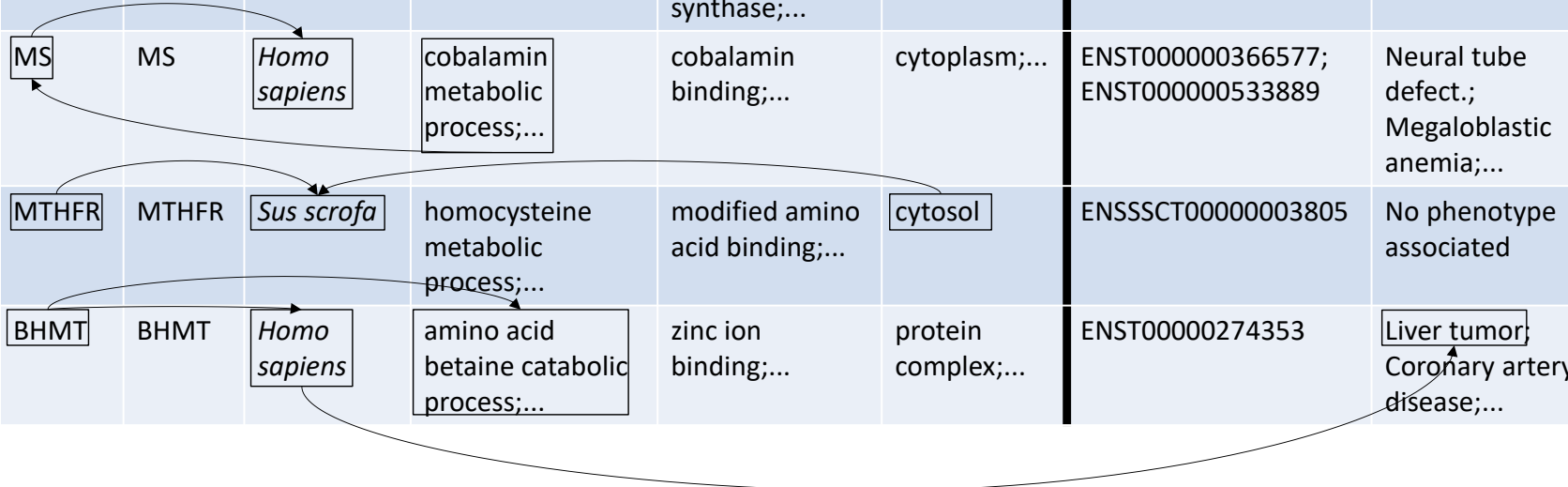


- Provide definitions
- Provide axioms that are universally true
- Obey formal semantics
- Main use case:
  - annotation of biological database entries

# Example – tabular Bio-DB structure



ID	PR Protein	PR Gene	Organism (NCBI Tax)	GO Biological Process	GO Molecular Function	GO Cellular component	Ensembl ID	Ensembl Phenotype
F1MEW4	CBS	CBS	<i>Bos taurus</i>	blood vessel remodelling;...	Cysthationine beta-synthase;...	cytoplasm;...	ENSBTAT00000000184; ...	No phenotype associated
Q99707	MS	MS	<i>Homo sapiens</i>	cobalamin metabolic process;...	cobalamin binding;...	cytoplasm;...	ENST000000366577; ENST000000533889	Neural tube defect; Megaloblastic anemia;...
F1RF82	MTHFR	MTHFR	<i>Sus scrofa</i>	homocysteine metabolic process;...	modified amino acid binding;...	cytosol	ENSSSCT00000003805	No phenotype associated
Q93088	BHMT	BHMT	<i>Homo sapiens</i>	amino acid betaine catabolic process;...	zinc ion binding;...	protein complex;...	ENST00000274353	Liver tumor; Coronary artery disease;...



# Example – tabular Bio-DB structure



ID	PR Protein	PR Gene	Organism (NCBI Tax)	GO Biological Process	GO Molecular Function	GO Cellular component	Ensembl ID	Ensembl Phenotype
F1MEW4	CBS	CBS	<i>Bos taurus</i>	blood vessel remodelling;...	Cysthationine beta-synthase;...	cytoplasm;...	ENSBTAT00000000184; ...	No phenotype associated
Q99707	MS	MS	<i>Homo sapiens</i>	cobalamin metabolic process;...	cobalamin binding;...	cytoplasm;...	ENST000000366577; ENST000000533889	Neural tube defect; Megaloblastic anemia;...
F1RF82	MTHFR	MTHFR	<i>Sus scrofa</i>	homocysteine metabolic process;...	modified amino acid binding;...	cytosol	ENSSSCT00000003805	No phenotype associated
Q93088	BHMT	BHMT	<i>Homo sapiens</i>	amino acid betaine catabolic process;...	zinc ion binding;...	protein complex;...	ENST00000274353	Liver tumor; Coronary artery disease;...

'is included in'

'is included in'

'is included in'

'is included in'

'is participant in'

'is part of'

'includes'

# Example – tabular Bio-DB structure

more abstract:

- a database record informs about experimental evidence that:
  - Proteins of the type *Prot1*
    - participate in Processes of type *BProc<sub>1</sub>... BProc<sub>k</sub>* within organisms of type *Org<sub>1</sub>*
    - are active in cellular components of type *CComp<sub>1</sub>* or *CComp<sub>2</sub>* or ... *CComp<sub>x</sub>* within organisms of type *Org<sub>1</sub>*
    - participate in Processes that have small molecules of type *Mol1, Mol<sub>2</sub>...Mol<sub>y</sub>* as outcome
    - within organisms of type *Org<sub>1</sub>* – if dysfunctional – *Org<sub>1</sub>* has dispositions to develop the phenotypes (disorders) *Phen<sub>1</sub>...Phen<sub>z</sub>*

This information is not explicitly contained in the database – it is implicitly shared by database users and curators

# Ontological framework

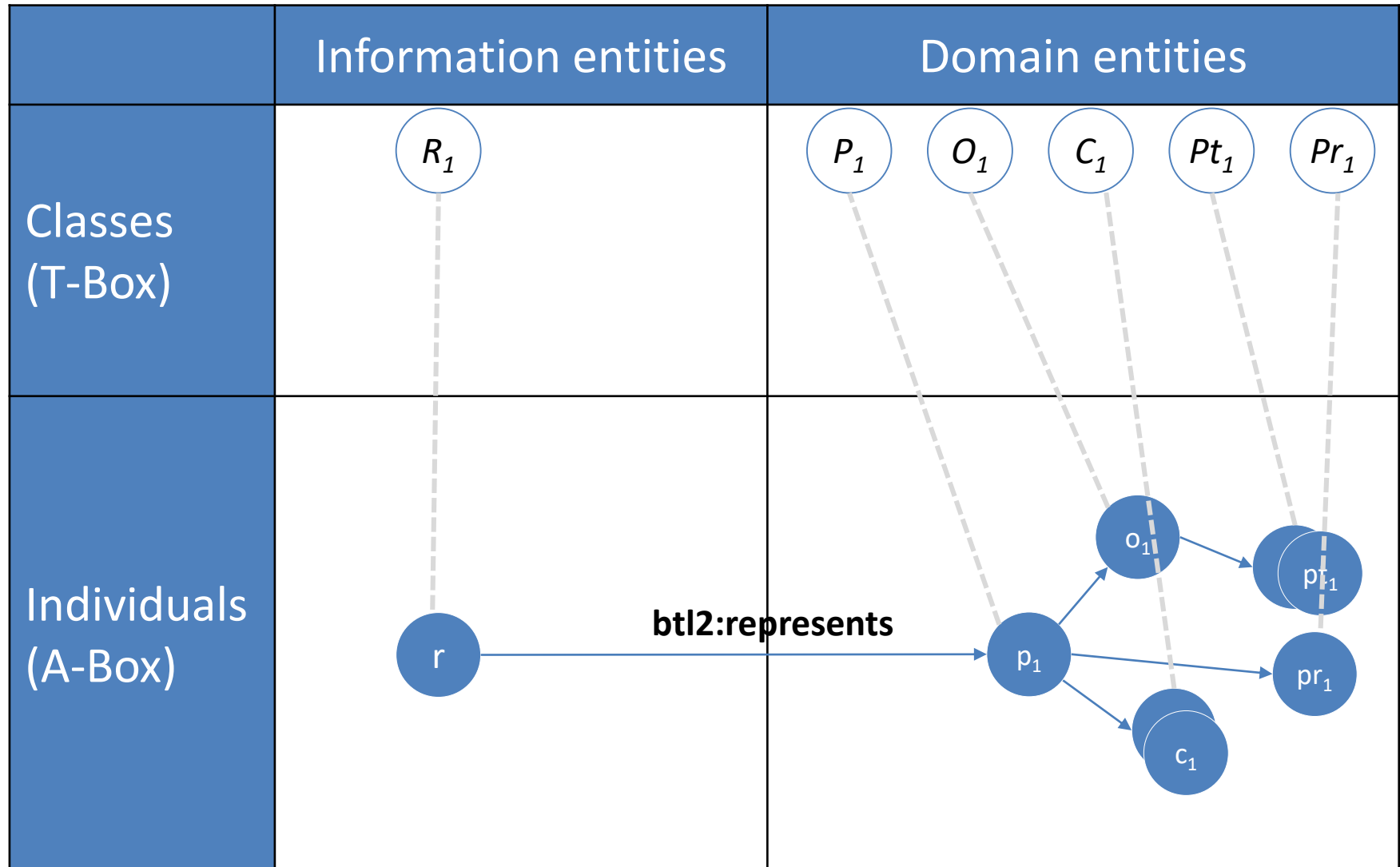
	Information entities	Domain entities
Classes (T-Box)	<ul style="list-style-type: none"><li>• Bio-DB</li><li>• Database record</li><li>• Data item</li></ul>	<ul style="list-style-type: none"><li>• <i>Homo sapiens</i></li><li>• <i>Megaloblastic anemia</i></li><li>• <i>Cobalamin binding</i></li><li>• <i>Methionin synthase</i></li><li>• (...)</li></ul>
Individuals (A-Box)	<ul style="list-style-type: none"><li>• Uniprot</li><li>• Ensembl</li><li>• Database record about Methionin Synthase in Homo Sapiens</li><li>• Data item, such as "cobalamin binding" in this record</li></ul>	<ul style="list-style-type: none"><li>• John Doe, of which tissue is stored in a biobank and analysed in a lab</li><li>• John's megaloblastic anemia</li><li>• A cobalamin binding process observed in the lab within a tissue sample from John</li><li>• a dysfunctional Methionin synthase protein molecule in John's tissue</li></ul>

# Denotation

	Information entities	Domain entities
Classes (T-Box)	<ul style="list-style-type: none"> <li>• Bio-DB</li> <li>• Database record</li> <li>• Data item</li> </ul>	<ul style="list-style-type: none"> <li>• <i>Homo sapiens</i></li> <li>• <i>Megaloblastic anemia</i></li> <li>• <i>Cobalamin binding</i></li> <li>• <i>Methionin synthase</i></li> <li>• (...)</li> </ul>
Individuals (A-Box)	<ul style="list-style-type: none"> <li>• UniProt</li> <li>• Ensembl</li> <li>• Database record about Methionin Synthase in Homo Sapiens</li> <li>• Data item, such as "cobalamin binding" in</li> </ul>	<ul style="list-style-type: none"> <li>• John Doe, of which tissue is stored in a biobank and analysed in a lab</li> <li>• John's megaloblastic anemia</li> <li>• A cobalamin binding process observed in the lab within a tissue sample from John</li> <li>• a dysfunctional Methionin synthase protein molecule in John's tissue</li> </ul>



# Case 1: database entry represents individuals



# Multiple defined subclasses

'Prot<sub>i</sub> Dysf in Org<sub>i1</sub> with Phen<sub>i1,...,il</sub>' equivalentTo  
 'Prot<sub>i</sub> Dysf in Org<sub>i1</sub> and 'is part of' some (Org<sub>i1</sub> and includes some Phen<sub>i1,...,im</sub>)  
 'Prot<sub>i</sub> in Org<sub>i1</sub> in BProc<sub>i1,...,il</sub>' equivalentTo 'Prot<sub>i</sub> in Org<sub>i1</sub>' and 'is participant in' some BProc<sub>i1,...,im</sub>  
 'Prot<sub>i</sub> in Org<sub>i1</sub> in CComp<sub>i1,...,il</sub>' equivalentTo 'Prot<sub>i</sub> in Org<sub>i1</sub> and 'is included in' some CComp<sub>i1,...,im</sub>  
 'Prot<sub>i</sub> in Org<sub>i1</sub> with Mol<sub>i1,...,il</sub>' equivalentTo  
 'Prot<sub>i</sub> in Org<sub>i1</sub>' and 'is participant in' some (Process and 'has participant' some Mol<sub>i1,...,im</sub>)

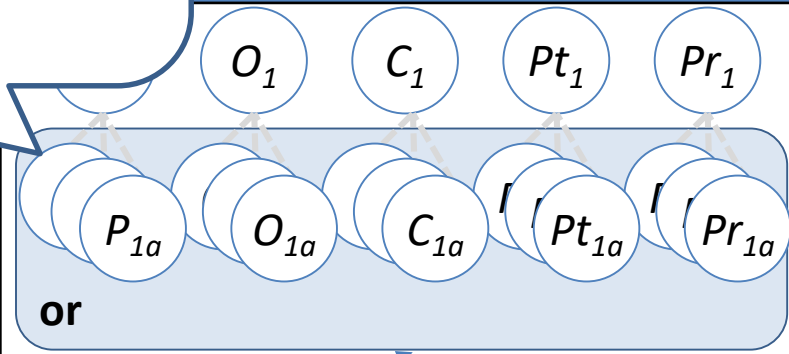
'Org<sub>i1</sub> with Prot<sub>i</sub>' equivalentTo Org<sub>i1</sub> and 'has part' some Prot<sub>i</sub>  
 'Org<sub>i1</sub> with Prot<sub>i</sub> Dysf' equivalentTo Org<sub>i1</sub> and 'has part' some 'Prot<sub>i</sub> Dysf'  
 'Org<sub>i1</sub> with Phen<sub>i1,...,im</sub> and Prot<sub>i</sub> Dysf' equivalentTo  
 'Org<sub>i1</sub> with Prot<sub>i</sub> Dysf' and includes some Phen<sub>i1,...,il</sub>

# represents classes

## Domain entities

Classes  
(T-Box)

Individuals  
(A-Box)

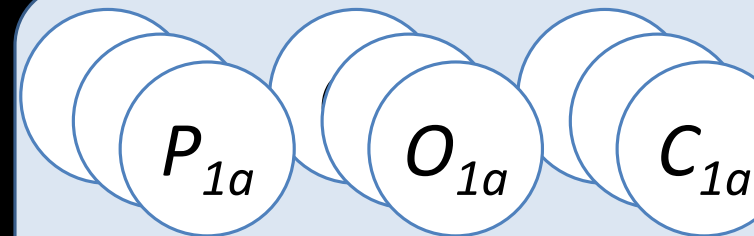


btl2:represents rdf:Type only

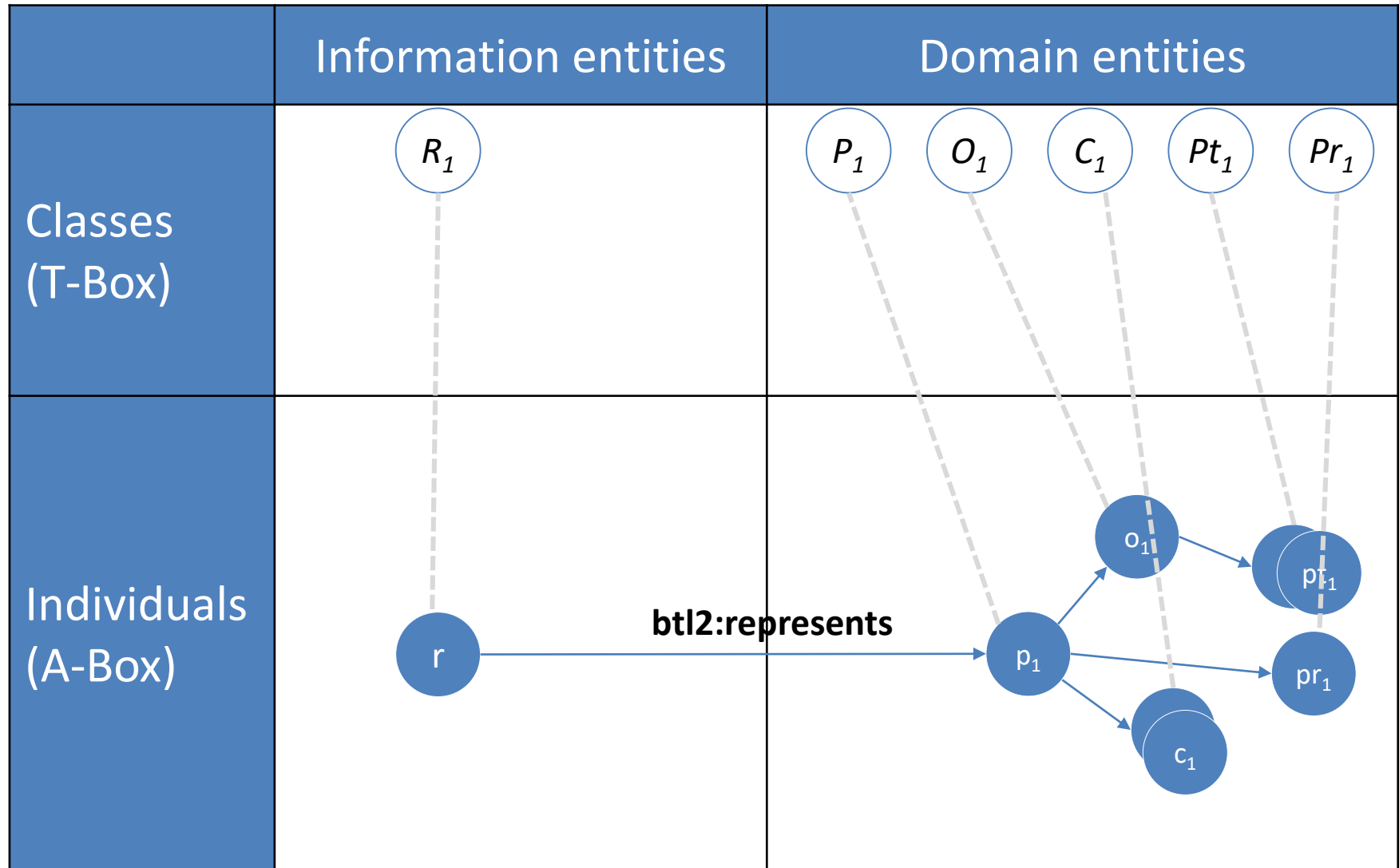


# Multiple defined subclasses

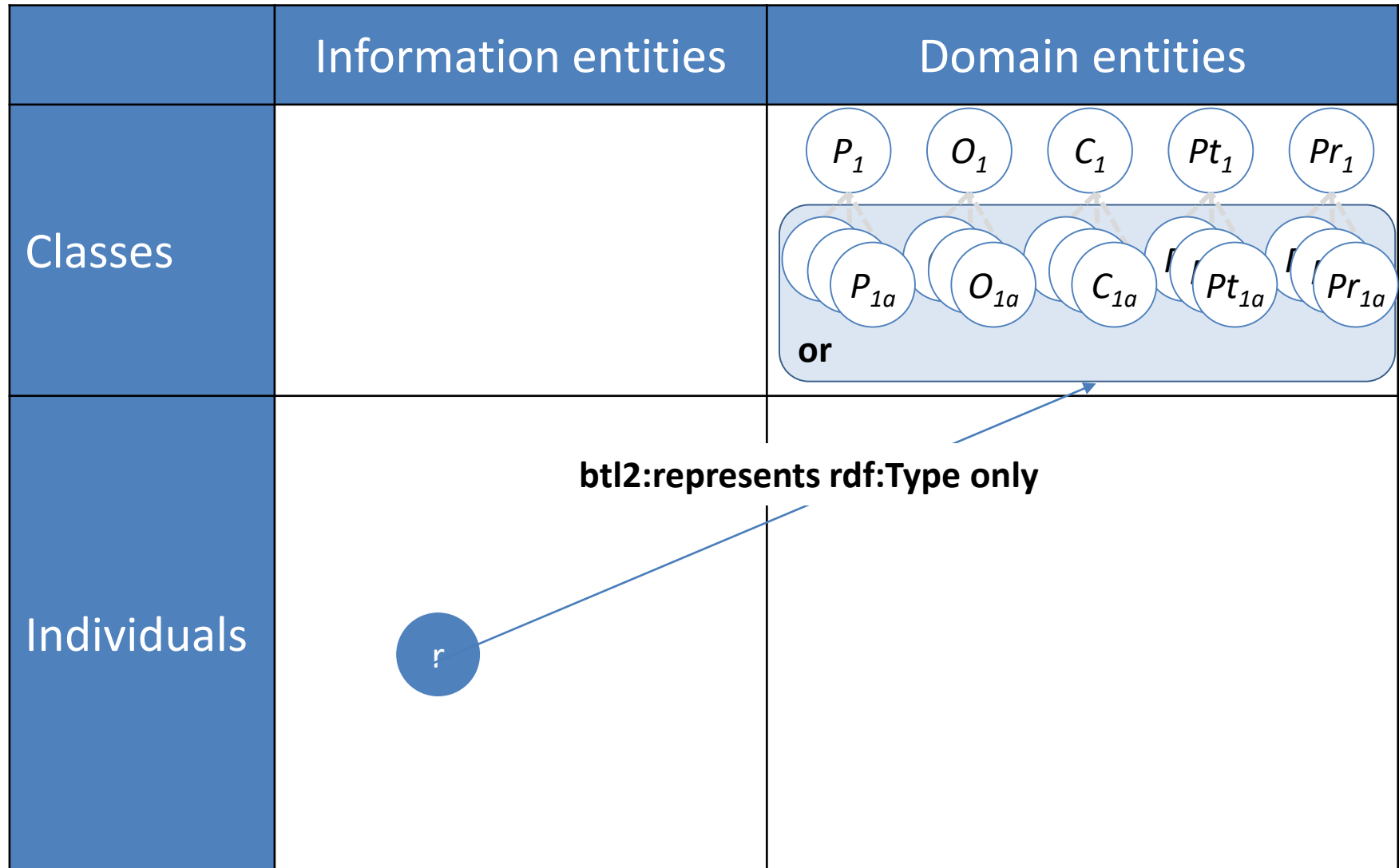
- ' $Prot_i$  Dysf in  $Org_{i1}$  with  $Phen_{i1, \dots, il}$ ' equivalentTo
  - ' $Prot_i$  Dysf in  $Org_{i1}$  and **'is part of'** some ( $Org_{i1}$  and **includes** some  $Phen_{i1, \dots, im}$ )
- ' $Prot_i$  in  $Org_{i1}$  in  $BProc_{i1, \dots, il}$ ' equivalentTo ' $Prot_i$  in  $Org_{i1}$ ' and **'is participant in'** some  $BProc_{i1, \dots, im}$
- ' $Prot_i$  in  $Org_{i1}$  in  $CComp_{i1, \dots, il}$ ' equivalentTo ' $Prot_i$  in  $Org_{i1}$  and **'is included in'** some  $CComp_{i1, \dots, im}$
- ' $Prot_i$  in  $Org_{i1}$  with  $Mol_{i1, \dots, il}$ ' equivalentTo
  - ' $Prot_i$  in  $Org_{i1}$ ' and **'is participant in'** some ( $Process$  and **'has participant'** some  $Mol_{i1, \dots, im}$ )
  
- ' $Org_{i1}$  with  $Prot_i$ ' equivalentTo  $Org_{i1}$  and **'has part'** some  $Prot_i$
- ' $Org_{i1}$  with  $Prot_i$  Dysf' equivalentTo  $Org_{i1}$  and **'has part'** some ' $Prot_i$  Dysf'
- ' $Org_{i1}$  with  $Phen_{i1, \dots, im}$  and  $Prot_i$  Dysf' equivalentTo
  - ' $Org_{i1}$  with  $Prot_i$  Dysf' and **includes** some  $Phen_{i1, \dots, il}$



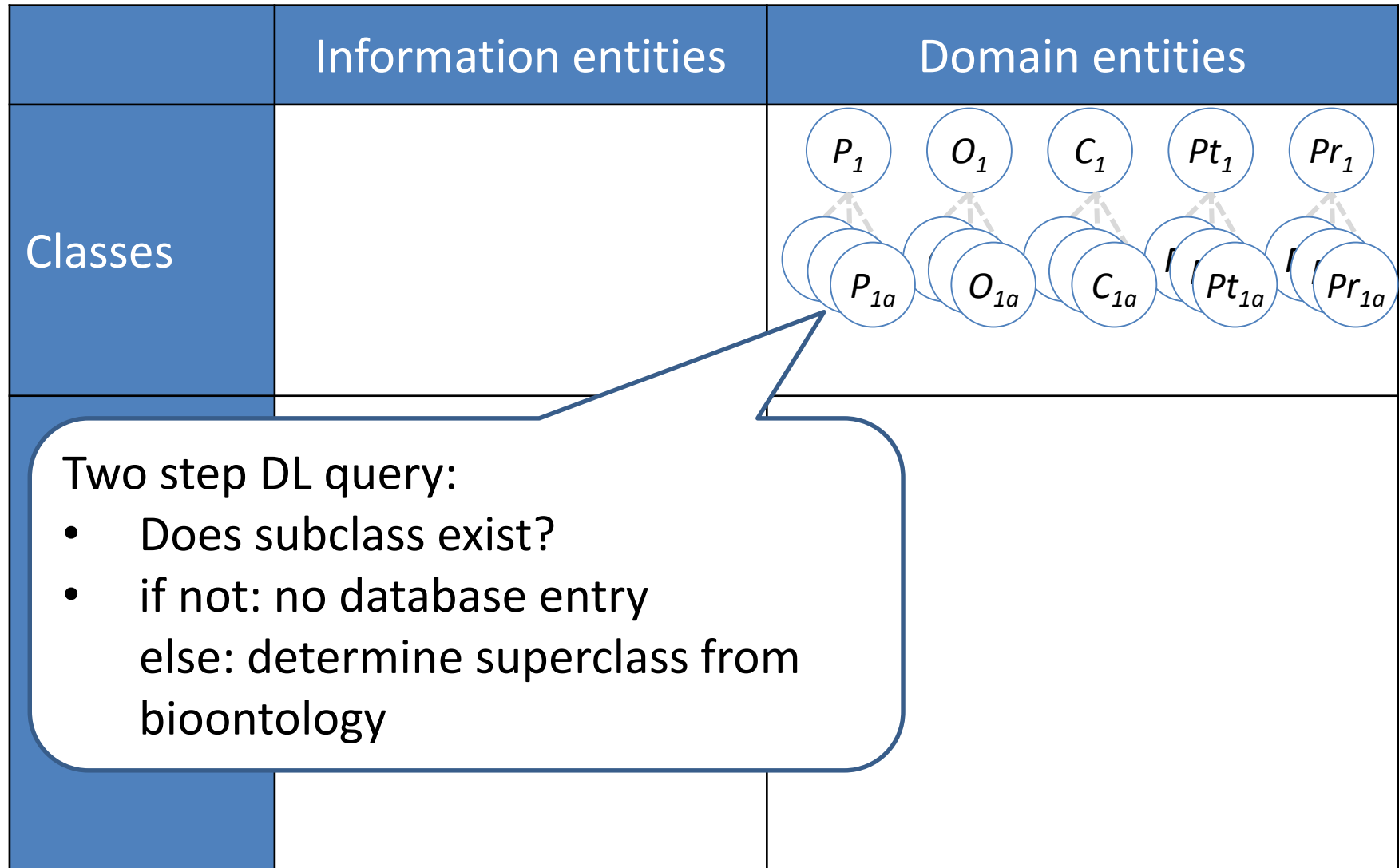
# Querying: A-box query for individuals



# Querying: T-box query for subclasses



# Querying: T-box query for subclasses



# Competency questions

(Q1) Which biological processes have proteins of the kind  $Prot_i$  as participants?

$BProc_1$  and ('has participant' some  $Prot_i$ )

(Q2) In which cellular locations is  $Prot_i$  active in organisms of the type  $Org_1$ ?

'Cellular component' and ('is included in' some  $Org_1$ )  
and (includes some  $Prot_i$ )

(Q3) Which proteins are involved in processes of the type  $BProc$  in organisms of the type  $Org_1$ ?

*Protein* and (' is participant in' some  $BProc_1$ )  
and ('is included in' some  $Org_1$ )

# Evaluation (one database record)

Model	Q1	Q2	Q3	Classes	Individuals	Axioms/ Assertions
A- Box	bp1001, bp2001, bp3001	cc1001, cc2001, cc3001	p1004	24	51	207
T- Box	<i>BProc<sub>1</sub></i>	<i>CComp<sub>1</sub></i>	<i>Prot<sub>i</sub></i>	68	0	149



# Discussion (I)

- Both modelling solutions:
  - highly productive
  - scaling problems to be expected
- A-Box solution (prototypical individuals):
  - A-box reasoning more costly
  - Makes existential assumptions
- T-Box solution (multiple subclasses)
  - Theoretically allows non-referential entries
  - Simplified model: EL++

# Discussion (II)

- Do biological database refer to ontological content?
  - No "real" universal statements on biological entities
  - Even no existential assumption
  - Dispositional statements to be discussed (see paper)
- Exercise best described as ontological representation of referring individuals
- Possible use case: non-disruptive querying of Bio-DBs where axioms of the annotation ontologies need to be explored

# Conclusion

- Four ontological approaches - IND, SUBC, DISP and HYB
  - Structure and content of BIO-DBs
- Solution:
  - Expressiveness, DB retrieval and retrieval based on DL queries
  - Interpretation:
    - Denoted entities as prototypical individuals
    - Creation of defined subclasses
    - Database content as reporting dispositions

Funding: This work was funded by *Conselho Nacional de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) 3914/2014-03* and *Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) 140698/2012-4*.