



***mig***

Unité Mathématique, Informatique et Génomique

# Ontology design

**Transys training session**

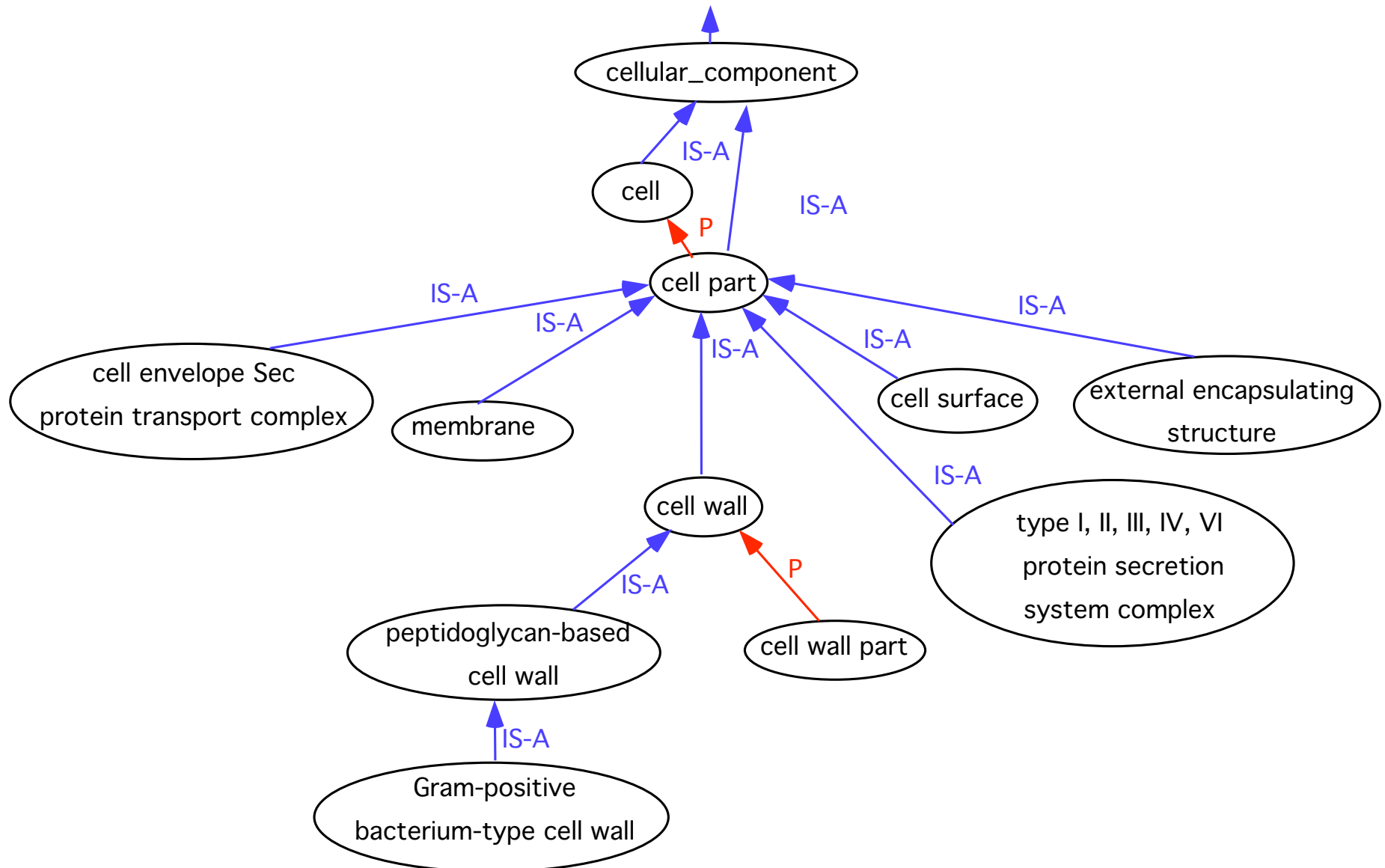
**Jouy-en-Josas – 3rd June 2009**



ALIMENTATION  
AGRICULTURE  
ENVIRONNEMENT

**INRA**


## Example of Ontology from Gene Ontology



AmiGO! Your friend in the Gene Ontology.

http://amigo.geneontology.org/cgi-bin/amigo/go.cgi Google

EPIPAGRI Portail:Micr... - Wikipédia Wipo QuaeroWS AlloCiné iHOP Cocitations GO browser AnnuaireINRA Epipagri Alvis Doodle Alvis Search Engine >>

 the Gene Ontology AmiGO

Search Browse GOOSE Other Tools Help

Search the Gene Ontology database

cell wall part

☒ GO terms ☐ genes or proteins ☐ exact match

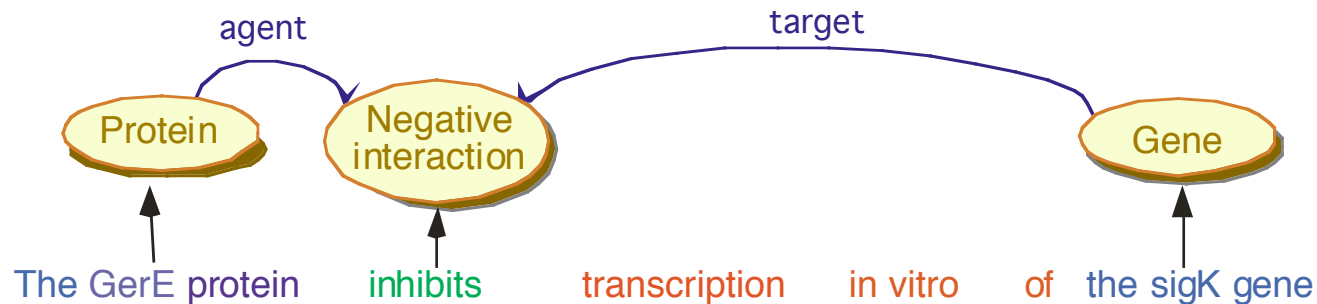
Soumettre

[Try AmiGO Labs](#)

GO database release 2009-06-02  
[Cite this data](#) • [Terms of use](#) • [GO helpdesk](#)  
Copyright © 1999-2009 [the Gene Ontology](#)

# Ontology for text interpretation

## Simple example of text annotation by an ontology

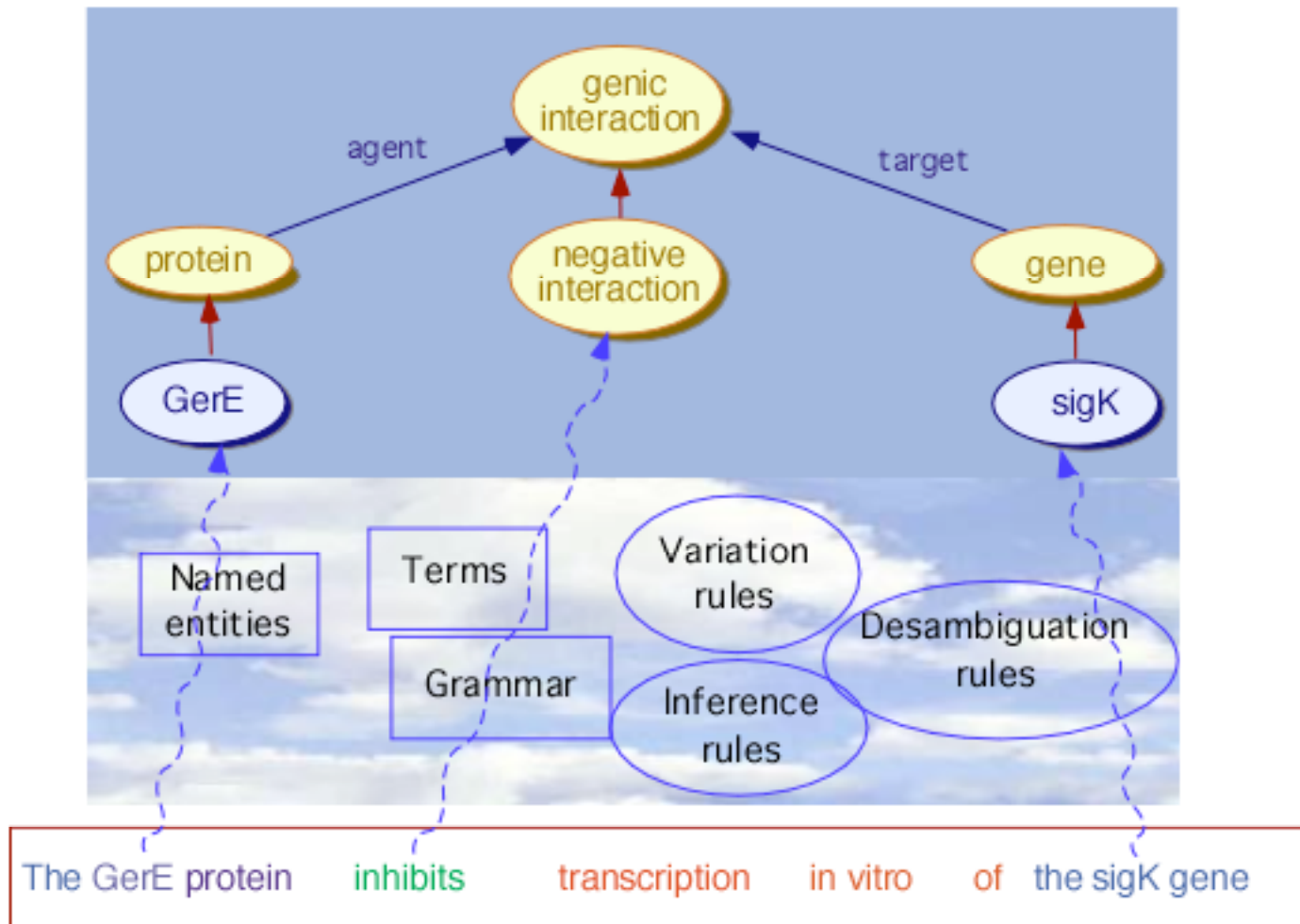


Associate text *fragments to concepts and relations of the ontology*

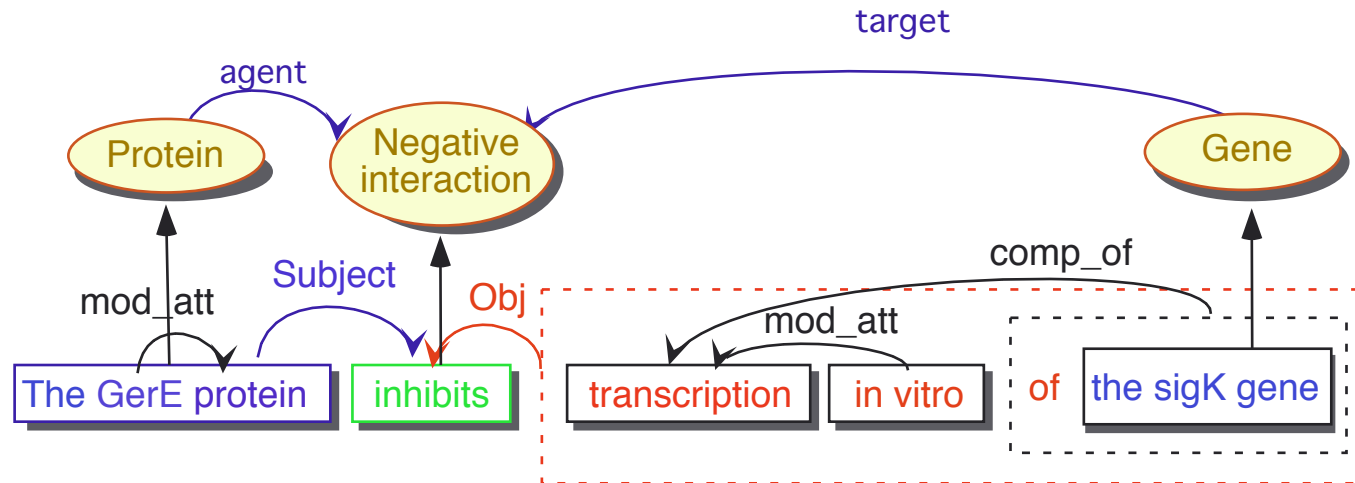
### In the simplest case

- The text terms coorespond to labels of ontology concepts (*GerE*, *inhibit*, *sigK*)
- The relations are deried from the order of the words (protein and gene surround and interaction verb)

# Text to ontology



# From text to semantic annotation



**Semantic relation**

agent(Ger\_protein, inhibit), target(sigK gene, inhibit), ...

**Semantic category**

is\_a(Ger\_protein, protein), is\_a(inhibit, negative\_interaction), ...

**Syntactic dependency**

subject(Ger\_protein, inhibit), obj(transcription, inhibit), ...

**Term**

term(GerE protein), term(in vitro) term(sigK gene), ...

**Grammatical category**

cat(the, det), cat(GerE, nom), cat(inhibit, verbe), ...

**Named entity**

entity(GerE, protein), entity(sigK, gene), entity(sigma K, protein), ...

**Segmented text**

word(the), word(Ger\_protein), word(inhibit), word(transcription), ...

## Ontology as a formal knowledge model

New facts can be automatically derived from the formal representation of knowledge in a model with guarantees of validity and completeness.

### *Simple example*

A first order logic model on ability of birds to fly.

Specific classes inherit from properties of more general classes.

Ontology =

$\{\text{Bird}(X) \rightarrow \text{Fly}(X),$	All birds fly
$\text{Canary}(X) \rightarrow \text{Bird}(X)\}$	All canaries are birds

---

$\text{Canary}(X) \rightarrow \text{Fly}(X)$	Then, all canaries fly
--	------------------------

Observation:  $\text{Canary}(\text{titi})$

Derived fact:  $\text{Bird}(\text{titi}), \text{Fly}(\text{titi})$

# Ontology and model

The ontology is the *science of being, of categories, and structures of the objects, of properties, events, processes and relations.*

Concept used in *Artificial Intelligence*,

Because of the need for **declarative representations** as reusable and general as possible,  
Corresponding the objects and processes that they represent.

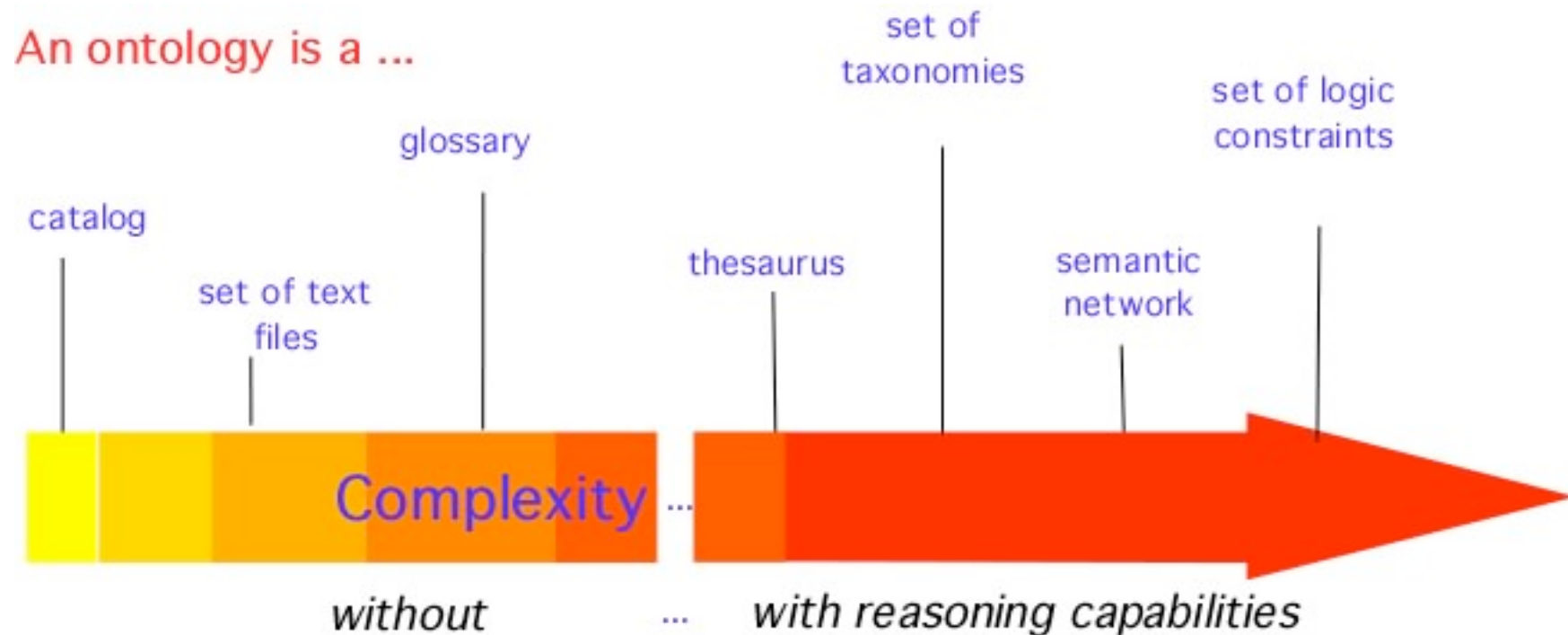
**Formal ontology** is an active research field in particular in *logic models* that are able to simulate reasoning knowledge by automatic reasoning mechanisms. (ex. *Owl* language).

Formal approaches and empirical knowledge bases converge in particular in the biomedical domain.



# Applied Ontologies

As opposed to formal ontologie, many kinds of representation are called ontologies (*C. Welty et al.*)



# Unstructured ontologies

**Catalogues, nomenclature** : item lists with unique identifiers,

*Example* :

[SwissProt](#)

[Hugo](#) (Human Gene Nomenclature) and links to d'[Ingenuity](#) bibliographic database

**Glossaries, terminologies or dictionnaires** : « keywords » with, or without, definitions and synonyms. For document indexing and search by « string-matching ».

*Examples*

## Gene Interactions Database

[CadInteract](#)

(cocitation of gene names)

Sentence	PMID
A low concentration of <b>gerE</b> activated cotB transcription by final <b>sigma(K)</b> RNA polymerase, whereas a higher concentration was needed to activate transcription of cotX or cotC.	10788508 <a href="#">Abstract</a> <a href="#">Annotation</a>
Transcription of the cotB, cotC, and cotX genes by final <b>sigma(K)</b> RNA polymerase is activated by a small, DNA-binding protein called <b>gerE</b> .	10788508 <a href="#">Abstract</a> <a href="#">Annotation</a>
In extension of recent work on the transcription of cot gene A and the mother-cell regulatory genes <b>gerE</b> , <b>sigK</b> and spoIIID, we show that genes involved in coat formation are turned on in a regulatory cascade of at least four co-ordinately controlled gene sets.	1691789 <a href="#">Abstract</a>

# Ontology and reasoning

**Thesaurus** : terms are organized in generality hierarchies. The properties of more general classes are supposed to be true for more specific classes.

*Exemples*

MeSH

[Cells \[A11\]](#)

[Cellular Structures \[A11.284\]](#)

[Cell Membrane \[A11.284.149\]](#) +

[Cell Surface Extensions \[A11.284.180\]](#) +

► [Cell Wall \[A11.284.183\]](#)

[Cell Wall Skeleton \[A11.284.183.200\]](#)

The taxinomic relations in thesaurii are often unformal, weak, and may produce invalid inferences.



Example from MeSH : mutation → mutagenesis

Instead of *mutation is\_a\_result\_of mutagenesis*

☹ Unexpected result: PubMed query: « *mutagenesis* » retrieves documents on *Expanded DNA Repeats*

# Formal ontology and inheritance

## Gene Ontology extract

- ☒ **I** GO:0005575 : cellular\_component [166204 gene products]
- ☒ **I** GO:0005623 : cell [103566 gene products]
- ☒ **P** GO:0044464 : cell part [103534 gene products]
- ☒ **I** GO:0030312 : external encapsulating structure [1459 gene products]
- ☐ **I** **GO:0005618 : cell wall [958 gene products]** 
  - ☒ **P** GO:0044426 : cell wall part [25 gene products]
  - ☒ **I** GO:0009277 : fungal-type cell wall [228 gene products]
  - ☐ **I** GO:0009274 : peptidoglycan-based cell wall [116 gene products] 
    - ☐ **I** GO:0009276 : Gram-negative-bacterium-type cell wall [114 gene products]
    - ☐ **I** GO:0009275 : Gram-positive-bacterium-type cell wall [0 gene products]
  - ☒ **I** GO:0009505 : plant-type cell wall [281 gene products]
  - ☒ **I** GO:0031160 : spore wall [25 gene products]

Part-of relation is distinguished from IS-A relation

## Gene Ontology extract

- ☒ **I** GO:0005575 : cellular\_component [166204 gene products]
- ☐ **I** GO:0005623 : cell [103566 gene products]
- ☒ **P** GO:0044464 : cell part [103534 gene products]
- ☐ **I** **GO:0016020 : membrane [30811 gene products]**
  - ☐ **I** GO:0030673 : axolemma [11 gene products]
  - ☒ **I** GO:0048475 : coated membrane [236 gene products]
  - ☒ **P** GO:0044425 : membrane part [13647 gene products]
  - ☒ **I** GO:0042175 : nuclear envelope-endoplasmic reticulum network [904 gene products]
  - ☒ **I** GO:0031090 : organelle membrane [4382 gene products]
  - ☒ **I** GO:0019867 : outer membrane [635 gene products]

## Subtilist classification extract

### 1 Cell envelope and cellular processes

#### 1.1 Cell wall

#### 1.2 Transport/binding proteins and lipoproteins

#### 1.3 Sensors (signal transduction)

#### 1.4 Membrane bioenergetics (electron transport chain and ATP synthase)

#### 1.5 Mobility and chemotaxis

#### 1.6 Protein secretion

#### 1.7 Cell division

#### 1.8 Sporulation

#### 1.9 Germination

#### 1.10 Transformation/competence

The taxinomic link does not denote here an *is-a*, nor a *kind-of* relation.

It means, « is a function related to ».

☺ Useful for gene annotation.

☹ Unexpected side-effect if used for automatic derivation.

## MultiFun (E. Coli) classification extract

### 6 Cell structure

#### 6.1 Membrane

#### 6.2 Peptidoglycan (murein)

#### 6.3 Surface antigens (ECA, O antigen of LPS)

#### 6.4 Flagellum

#### 6.5 Pilus

#### 6.6 Ribosome

#### 6.7 Capsule (M and K antigens)

☹ Membrane is defined as a cell structure instead of *part-of* the cell structure

### Basic Ontology Design Rules:

No conjunction of coordination (or, and)

No *other* (Other functions) or *Miscellaneous*

Use labels that are at most as possible self-sufficient, or at least comprehensible in the hierarchy context. (*Similar to unknown proteins / From B. subtilis ; No similarity*)

# Ontology design in Transys

## Two short-term main goals

Information retrieval = Bibliographic search with BioAlvis

Domain modeling for sharing a common representation among Transys partners

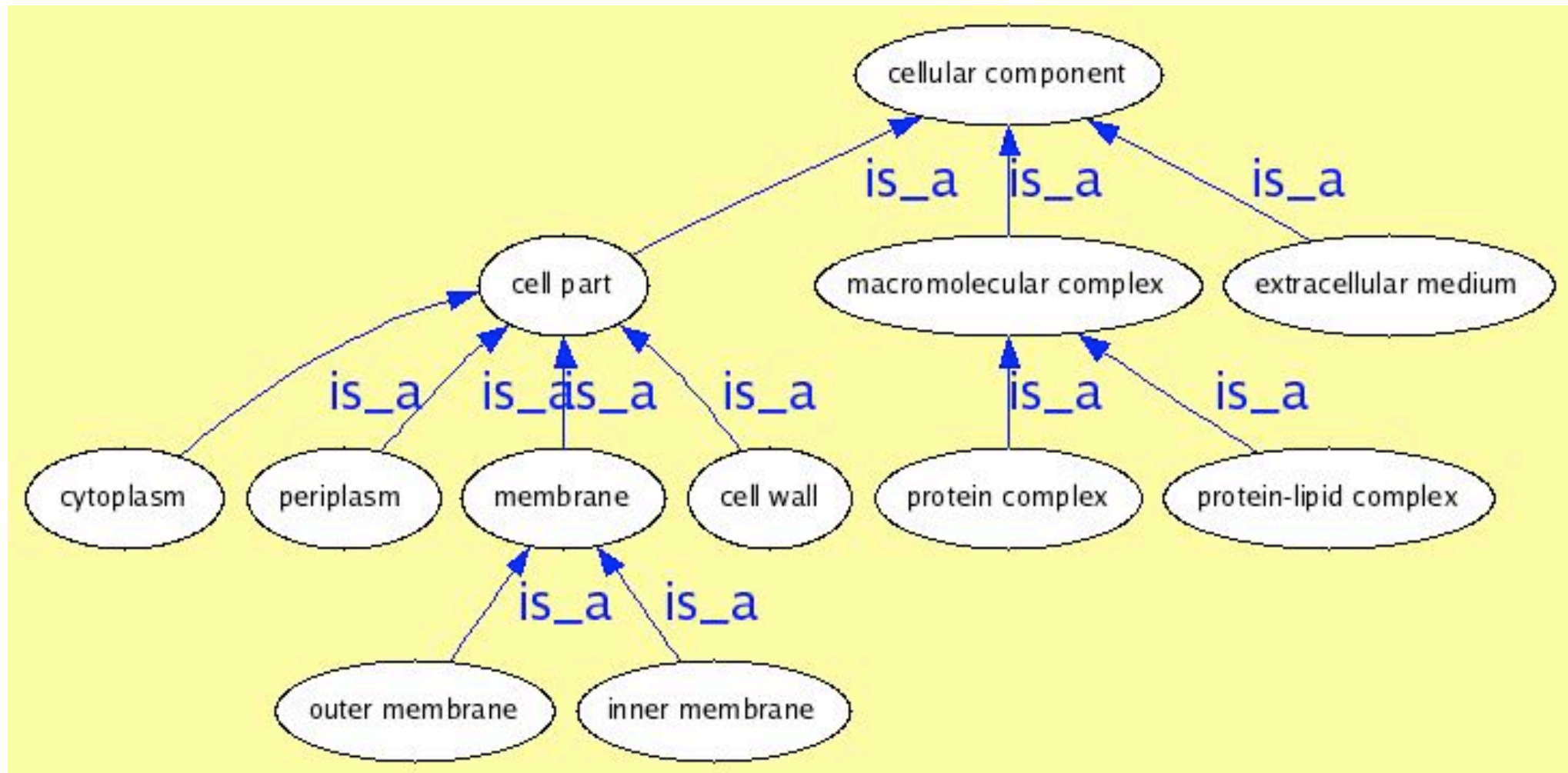
## Méthodology

1. Specify the domain and the role of the ontology
2. Acquire the knowledge from biologist expertise, scientific papers, databanks, thesaurus
3. Conceptualize
  - identify the key concept of the ontology, their properties and relations
  - Identify the terms to be used as concept labels
  - Build the structure
4. Integrate / align existing thesaurus / ontology (GO, MeSH, BactriOnto)

Non linear process.



## Transys core ontology on cell component



## Examples of term extracted from scientific papers

aminoacyl-tRNA synthetase cell division cycle elongation of protein synthesis elongation of RNA synthesis enzyme synthesis initiation of protein synthesis initiation of RNA synthesis late sporulation main glycolytic pathway phosphorylation- dephosphorylation proteolytic processing regulation of RNA synthesis ribosomal protein TCA cycle termination of protein synthesis termination of RNA synthesis adaptation to atypical condition	antibiotic production cell division cell-wall detoxification DNA recombination DNA restriction-modification and repair DNA-packaging and segregation germination membrane bioenergetics metabolism of amino-acid and related molecule metabolism of carbohydrate and metabolism of coenzyme and prosthetic group metabolism of lipid metabolism of nucleotide and nucleic acid	metabolism of sulfur motility and chemotaxis phosphate metabolism protein folding protein modification protein secretion protein synthesis RNA modification RNA synthesis sporulation transformation-competence transport-binding protein and lipoprotein transposon and IS cell envelope and cellular process information pathway intermediary metabolism
---	---	---

## **Ontology population ?**