

# Reasoning over large-scale biological systems with heterogeneous and incomplete data

Anne SIEGEL, CNRS, Rennes

Dyliss team (Inria, Univ Rennes, CNRS)

Institut de Recherche en Informatique et Systèmes Aléatoires

# Short presentation

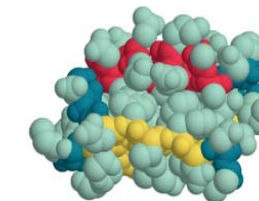
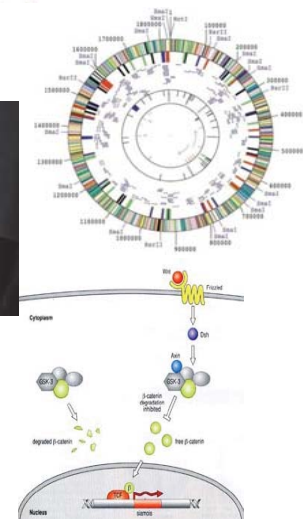
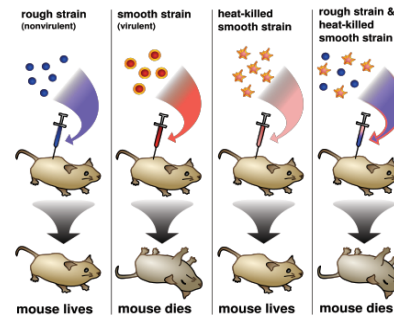
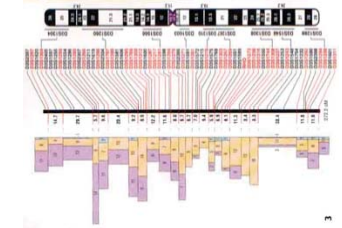
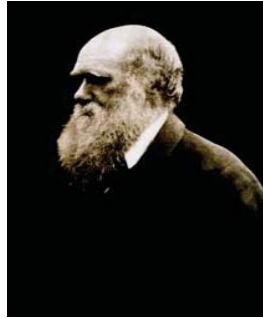
- **Research field**
  - Discrete dynamical systems & fractals
  - Systems biology
  - Knowledge representation
- **IRISA & INRIA Rennes**
  - 800 members, >40 teams
  - Univ Rennes, CNRS, Inria, etc...
- **Bioinformatics@Rennes**
  - GenOuest: platform, resource center
  - Genscale : NGS data analysis
  - **Dyliss: Integration of heterogeneous data**



# LIFE SCIENCE DATA



# From life science... to data science



## Naturalist approach

- Observing and deducing

## Experimental approach

- Perturbating and observing

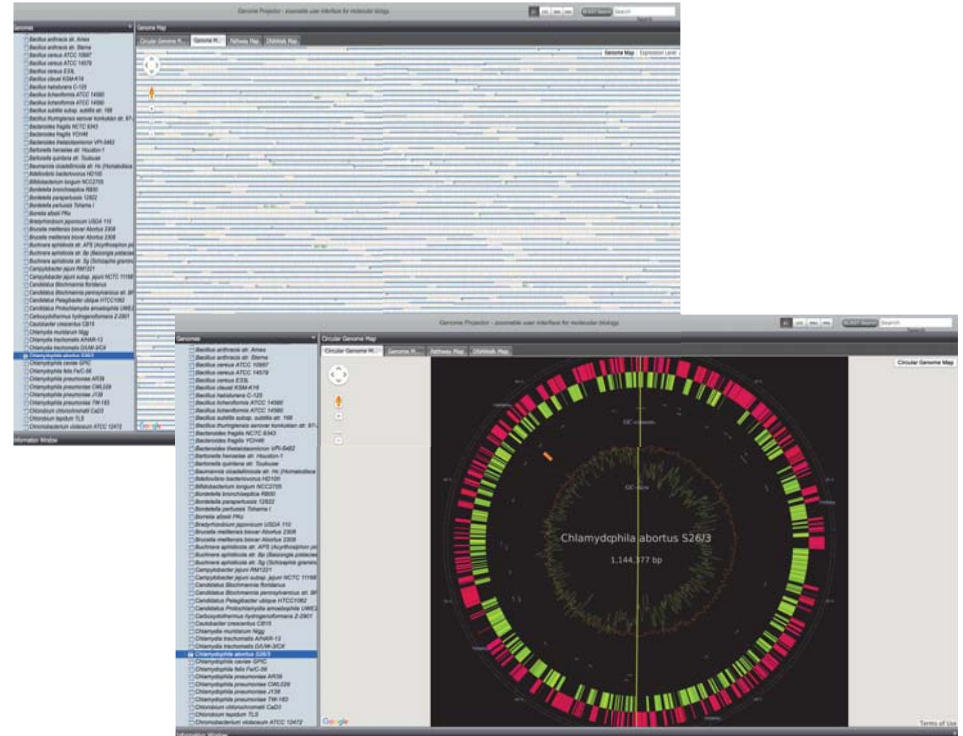
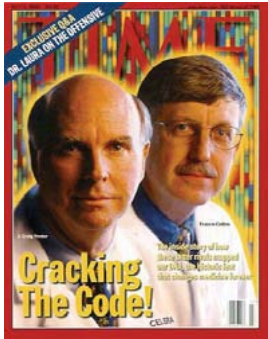
## Modern biology

- Measuring at lower scales



**Data science !**

# Biomolecular data: genomes



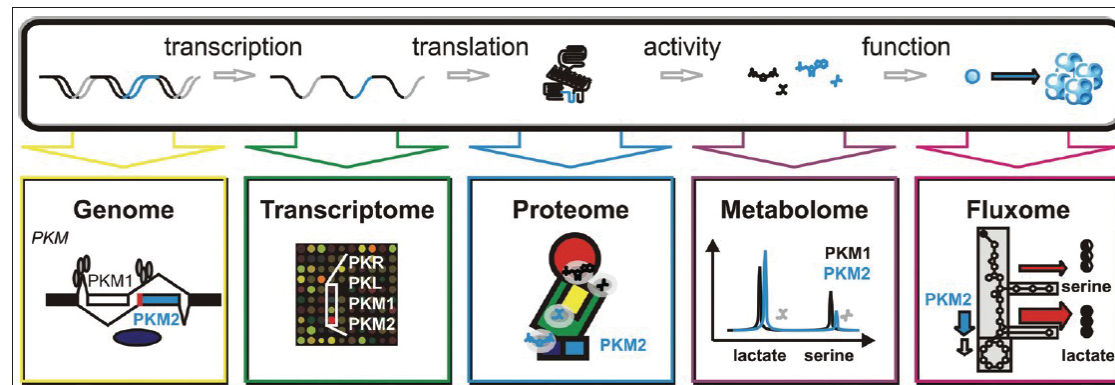
## Genome sequencing

- Very smart computational issues
- Bioinformatics

## Thousands of publicly available genomes

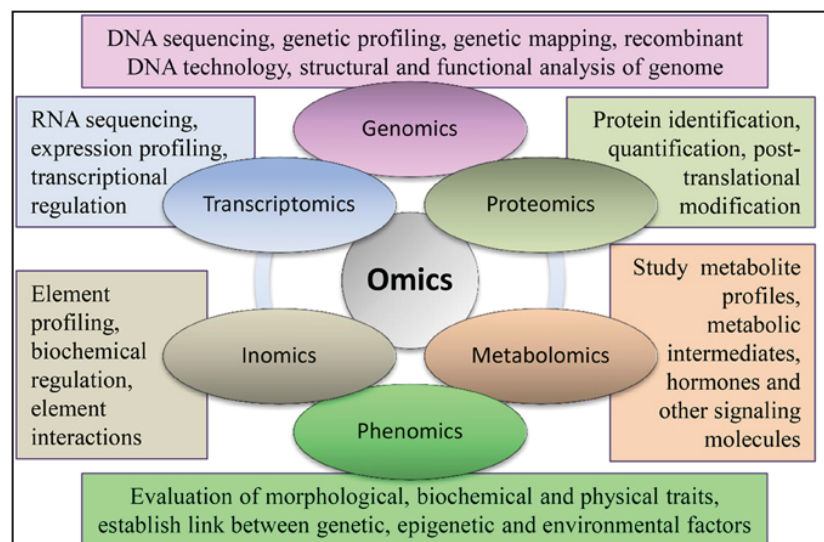
- Exploration, mapping and analysis

# What do we do with genomic data ?



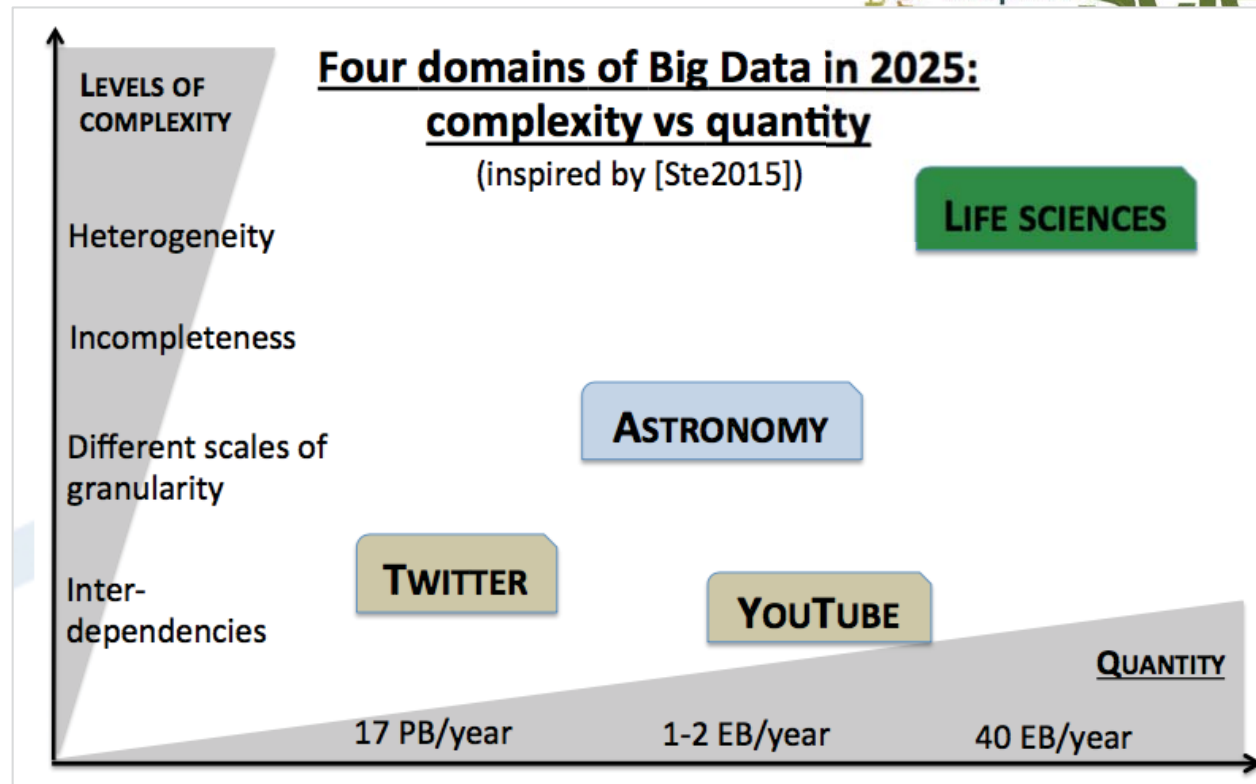
Assign a function to each DNA fragment

Develop new technologies to validate/refine the assigned functions



➔ Data deluge !

# Life science data nightmare



## Data characteristics

- Large-scale
- Incomplete
- Inter-dependent
- Heterogeneous / multi-scale

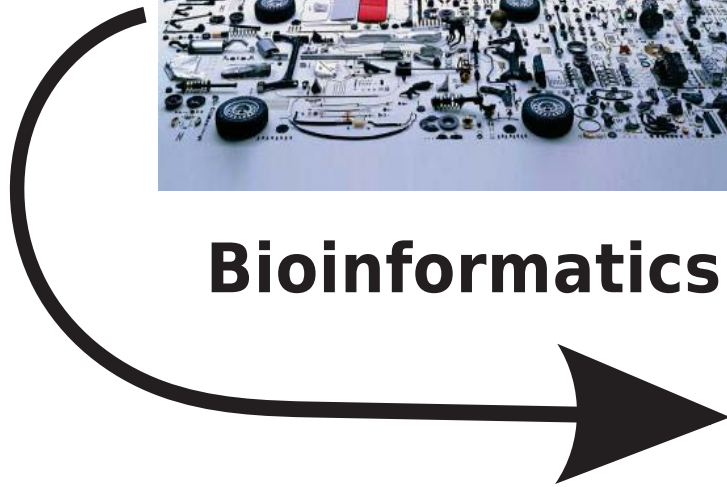


**How to integrate them?**

# SYSTEMS BIOLOGY

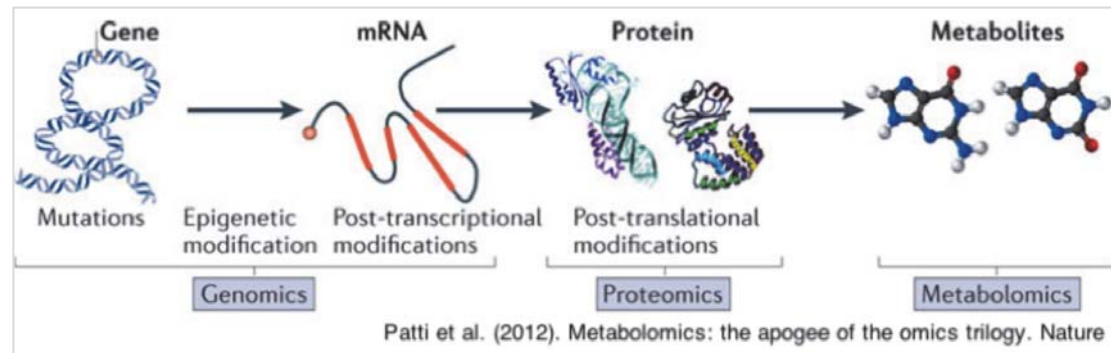


**Bioinformatics**



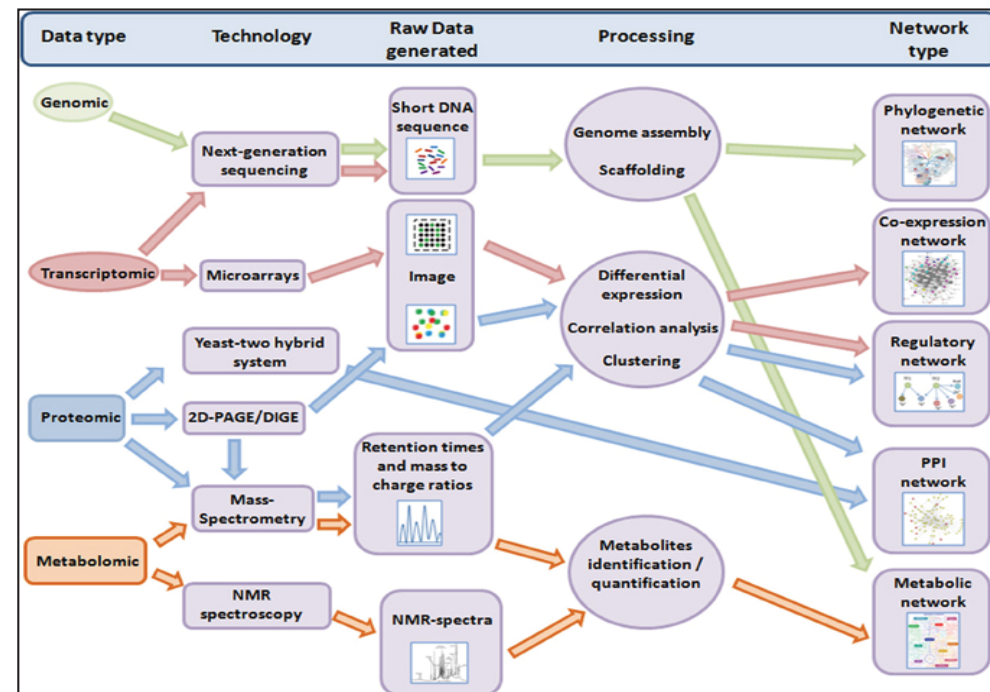
**Systems biology**

# Setting all together

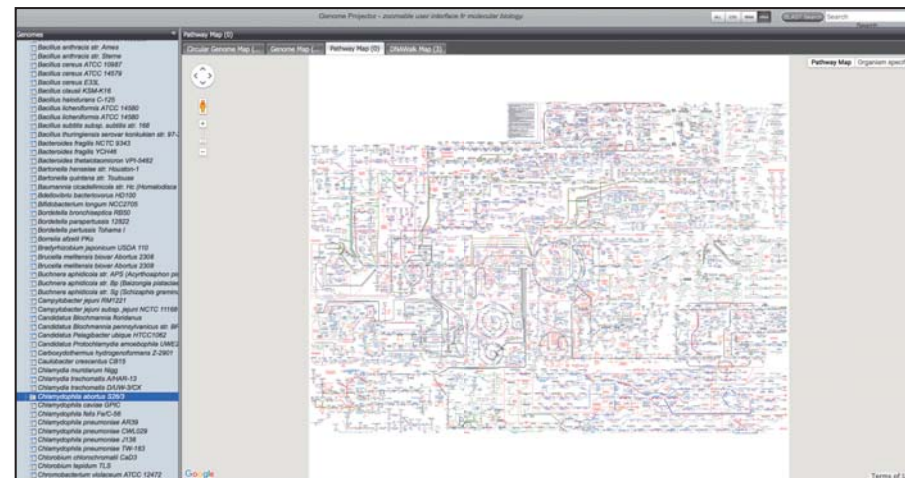
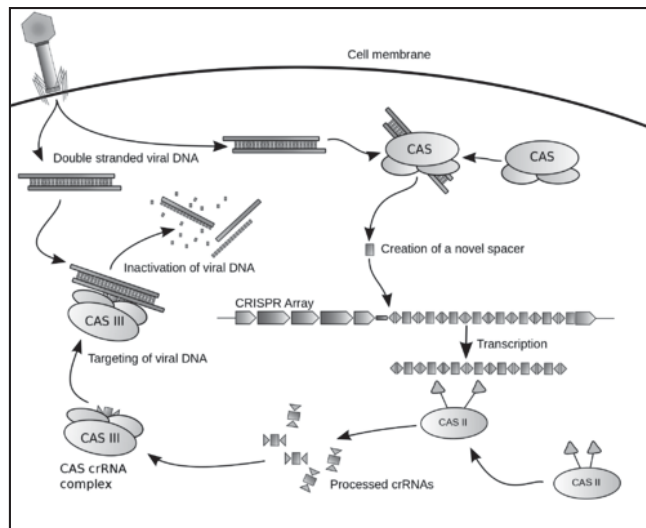
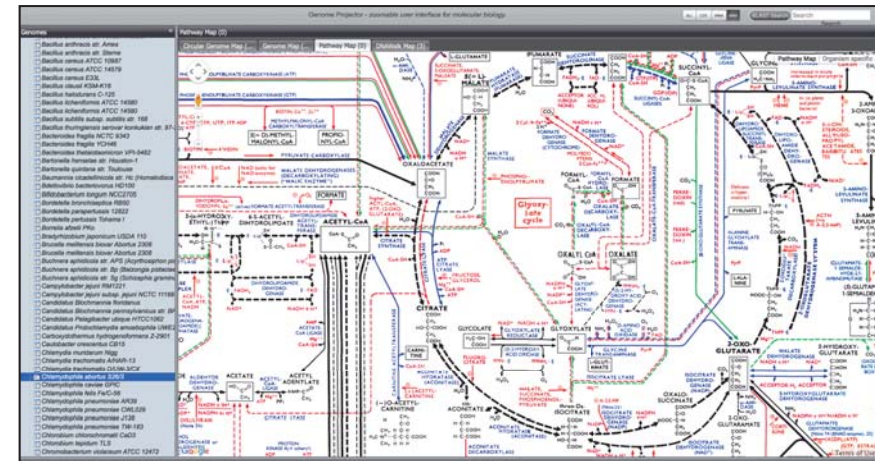
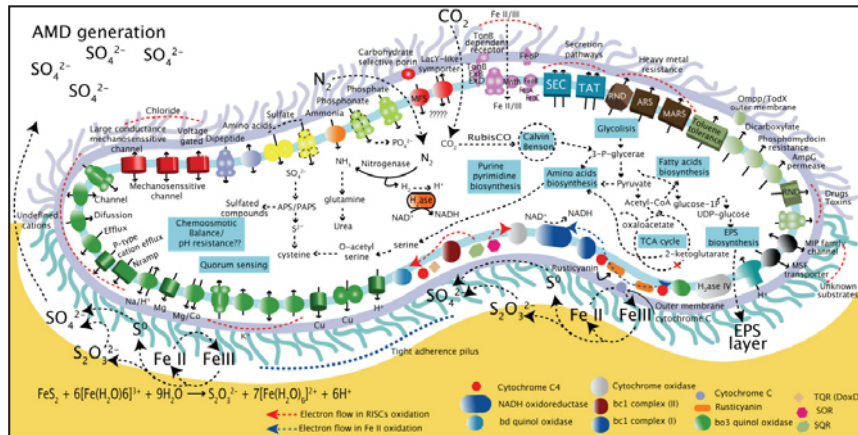


Gene function = regulation of a intra-cellular transformation procedure

- Biological interactions !
- Graphs / networks



# What we get...

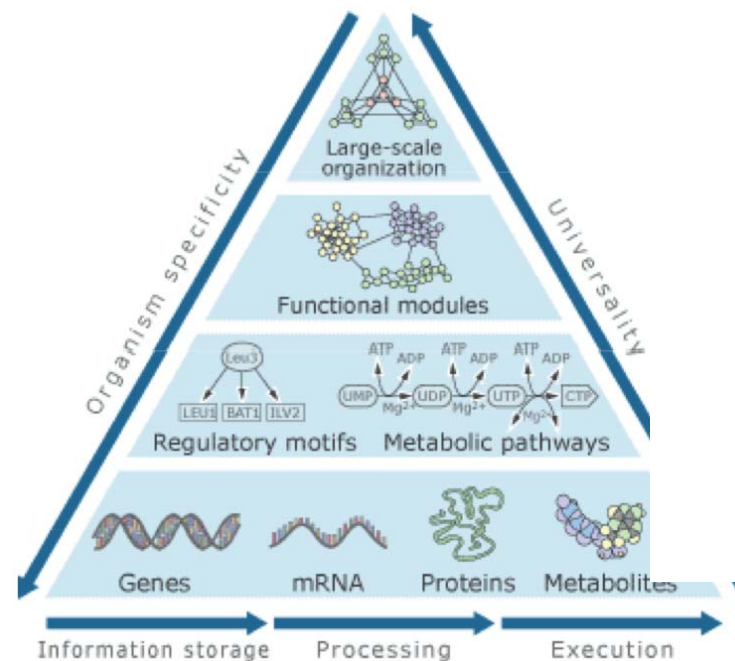


Large-scale graph description of interactions between compounds

# Systems biology

## Statement : **biology is a complex system**

- « Requires to examine the structure and dynamics of a cellular function rather than the characteristics of isolated parts of a cell » (Kitano, 2002)

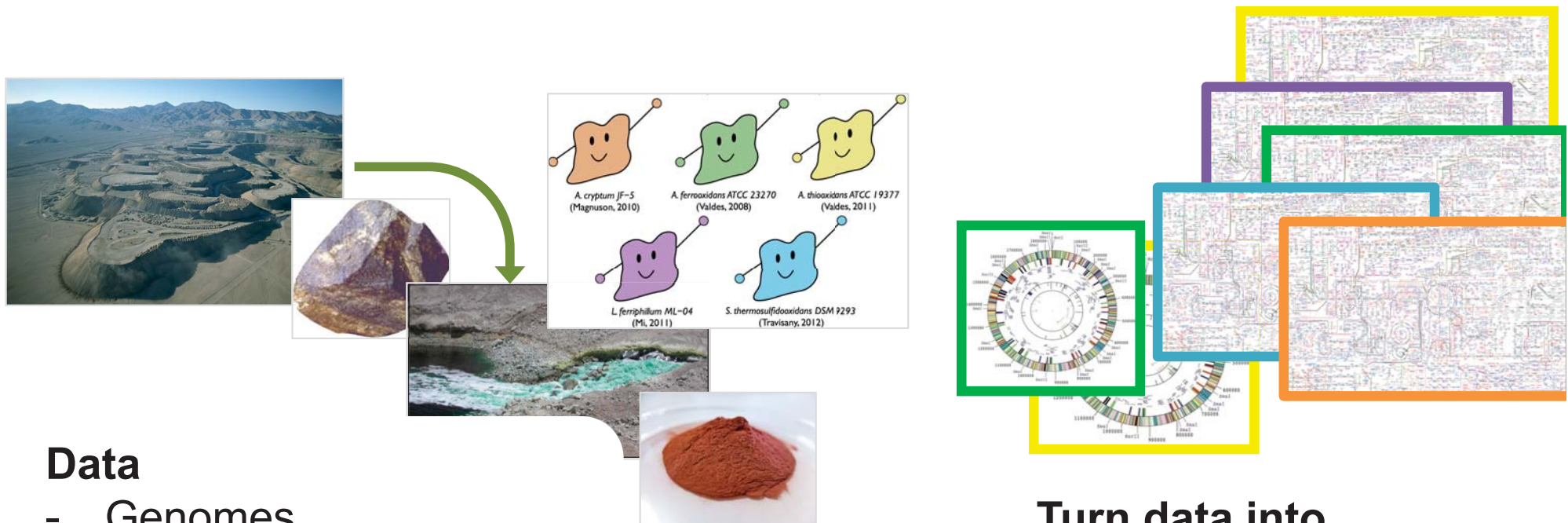


## Systems biology: **Interpreting multi-layer data and graphs**

- Produce predictive statements that can be experimentally validated

# Case-study: extremophile mining consortium

*Role of an **empirical taylor-made consortium** of bacteria in copper extraction from ore ?*



## Data

- Genomes
- Expression data
- Metabolic compounds

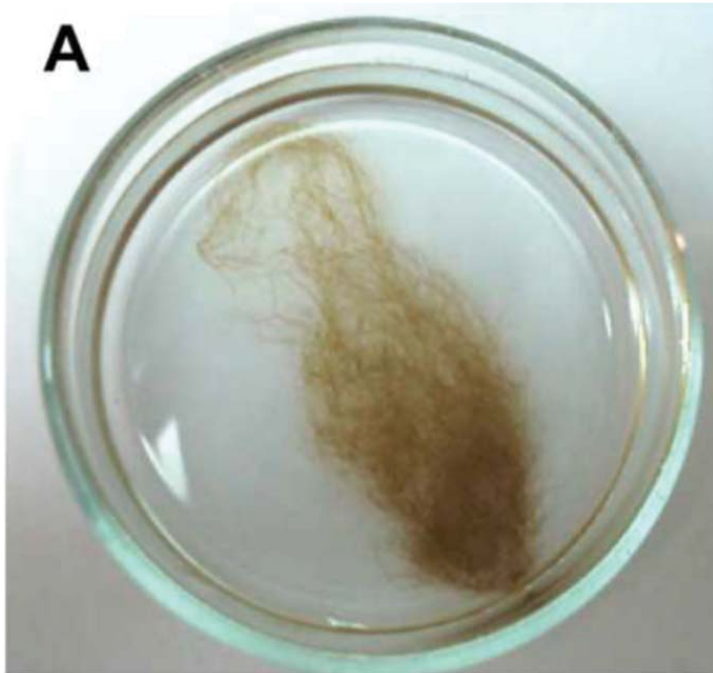
## Turn data into

- genomics maps
- interaction maps

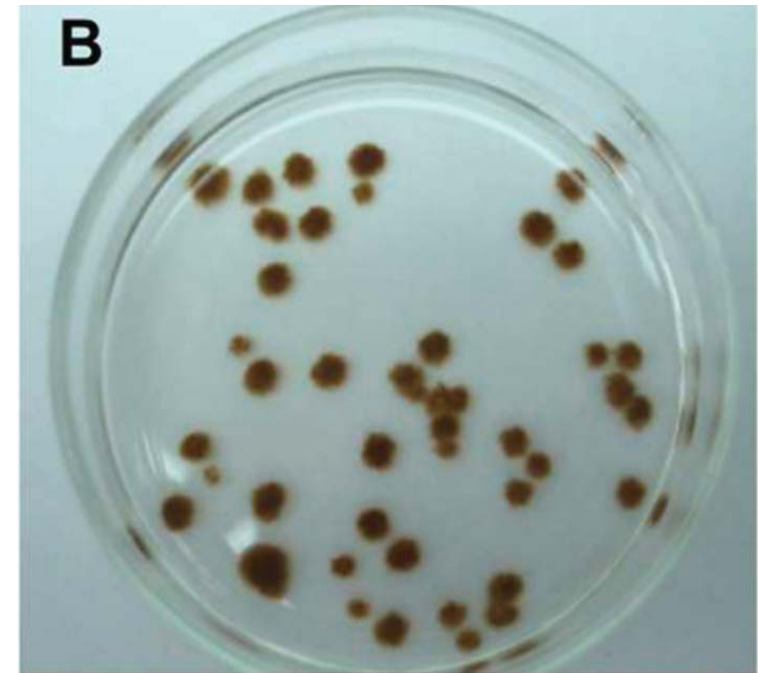
**Understand the contribution of each bacteria to the complete system ?**

➤ **integrative and systems biology**

## A second case-study : algal metabolism



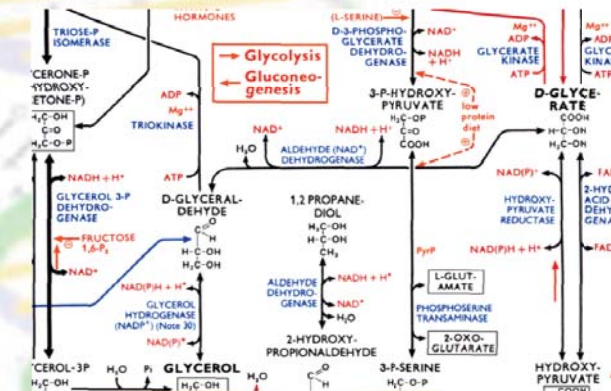
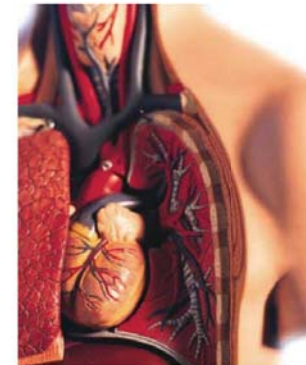
*E. siliculosus*



In axenic condition....

*Ectocarpus*  
[Dittami2014, Tapia2016]

**What is the role of environmental bacteria ?**



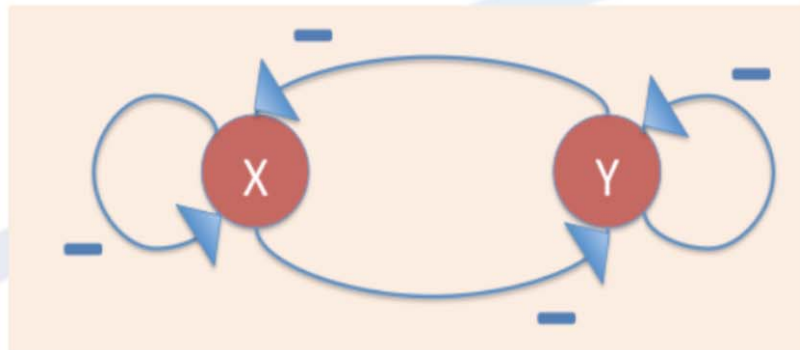
## Are molecular/cellular different than others ?

# Dynamical systems

## Historical motivation

Model the evolution of the set of components in a system according to time.

$$F : \begin{array}{ccc} \mathbb{T} & \times & \mathbb{S} \\ (t & , & \mathbf{z}) \\ \text{(time} & , & \text{state)} \end{array} \begin{array}{c} \rightarrow \mathbb{S} \\ \mapsto F(t, \mathbf{z}) \\ \text{new state at time } t \end{array}$$



$$\begin{aligned} \frac{dX}{dt} &= \frac{k}{K + Y^n} - aX \\ \frac{dY}{dt} &= \frac{l}{L + X^n} - bY \end{aligned}$$

Parameterized  
numerical system

$$f(X) \leftarrow 1 - Y$$

$$f(Y) \leftarrow 1 - X$$

Boolean model with  
asynchronous update  
scheme

## Identification of a dynamical system

Find the **best function F** which parcimounously explains and describes the observed responses of a system.

# Model identification since the 18th century

4

## What has always allowed a model identification

- **A priori knowledge about the laws governing the system**
  - Predetermined shape for the function  $F$
- **Limited number of components**
  - Reduction of the search space
- **Wide panel of sensors and perturbations**
  - Discriminate parameters

$$F : \begin{array}{ccc} T & \times & S \\ (t & , & \mathbf{z}) \\ \text{(time} & , & \text{state)} \end{array} \begin{array}{c} \rightarrow S \\ \mapsto F(t, \mathbf{z}) \\ \text{new state at time } t \end{array}$$

## Where is the complexity ?

- The search space grows exponentially with the number of measured compounds



**The more compounds we measure,  
the less identifiable a system is.**

# Differences between application domains

5

## Physical sciences

- **Knowledge.**  
Fundamental laws of physics.
- **Sensors.**  
Numerous.
- **Perturbations.**  
Various protocols in controlled frameworks.
- **System description.**  
Independent components

## Biological sciences

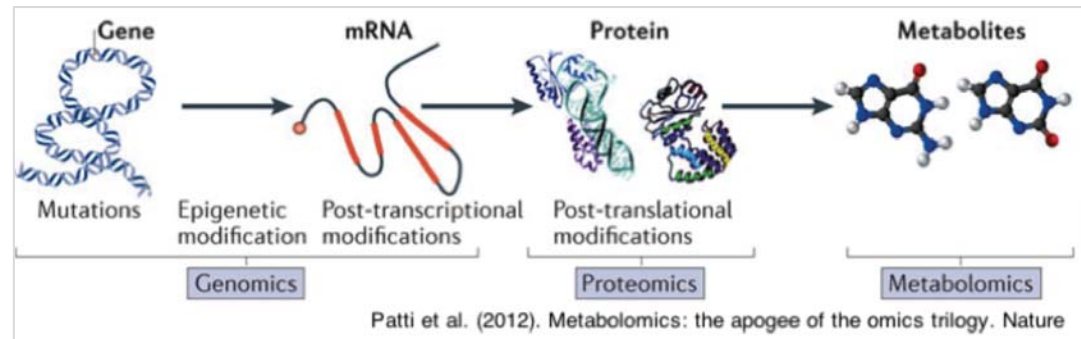
- **Knowledge.**  
Empirical laws
- **Sensors.**  
Low quality although numerous.
- **Perturbations.**  
Quite few according to sensors
- **System description.**  
Hidden dependencies

# Today's molecular/cellular biological systems

6

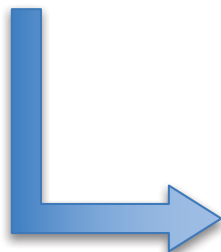
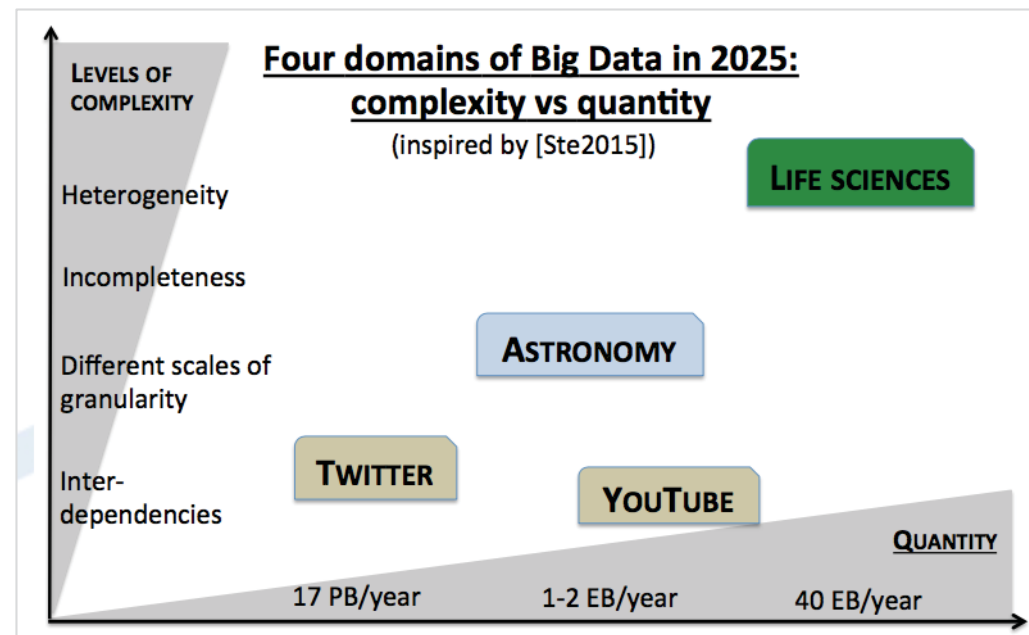
## Omics data.

- Large-scale
- Noisy
- Heterogeneous.



## Biological systems characteristics

- Large-scale
- Empirical laws
- Few data wrt the search space size



**Biological systems observed with omics data are not uniquely identifiable**

## Strategy: combine dynamical systems and constraints programming

### Describe a system by a family of abstract models

- Reason over a family of models instead of selecting a single one

### (Logical) knowledge representation

- Search space **description**
- Structured knowledge (link open data)

### Discrete dynamical systems

- **Links** between multi-scale observations.
- **Invariants** of model families.

### Solving optimisation problems

- **Replace laws by constraints**
- Extract robust information

# KNOWLEDGE REPRESENTATION

```

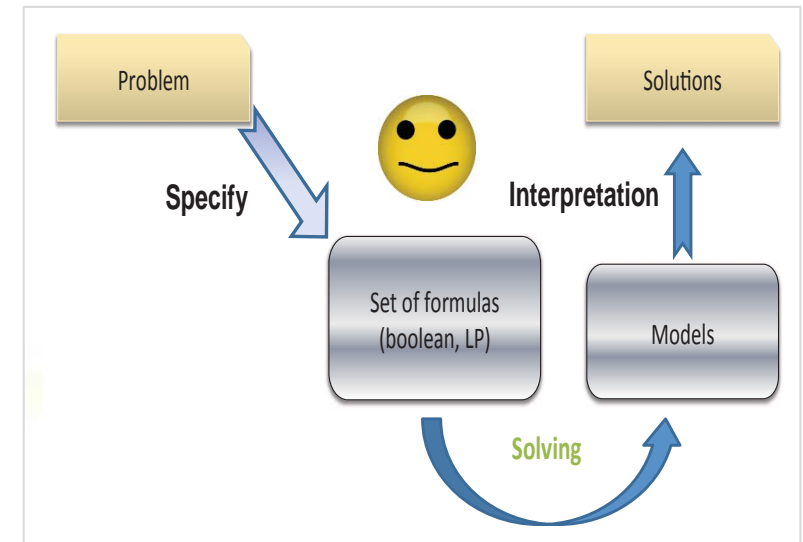
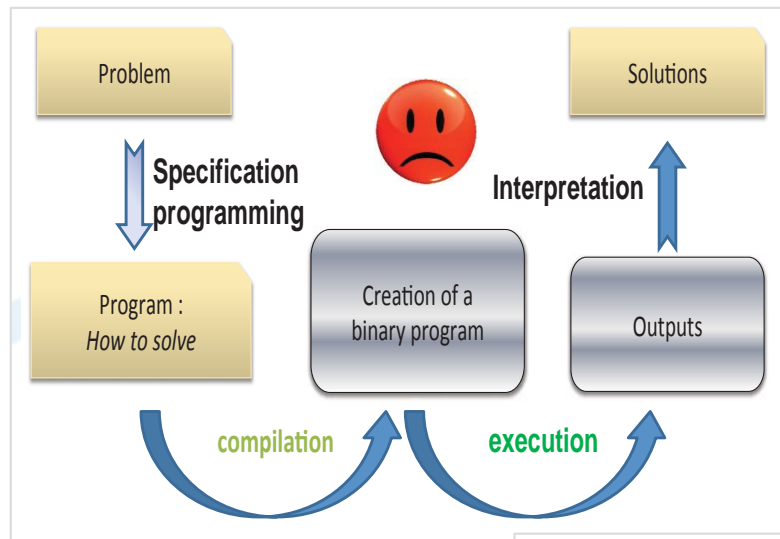
1{murderer(ms_Scarlet); murderer(colonel_Mustard)}1.
1{weapon_of_crime(revolver); weapon_of_crime(candlestick)}1.
1{place_of_crime(kitchen); place_of_crime(hall);
    place_of_crime(dining_room)}1.

crime_record(ms_Scarlet, 7). crime_record(colonel_Mustard, 4).

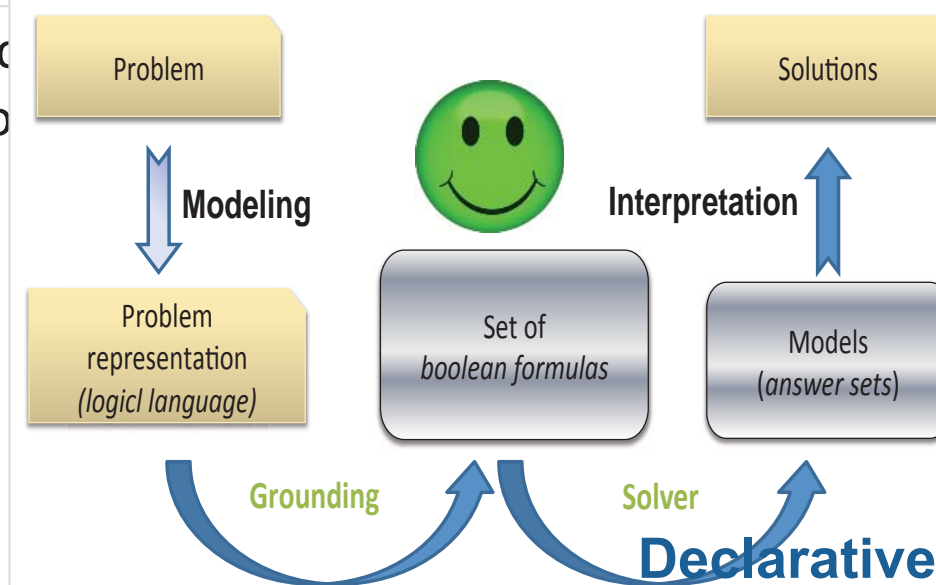
weapon_of_crime(candlestick).
:- place_of_crime(kitchen).
place_of_crime(hall) :- murderer(colonel_Mustard), not
    weapon_of_crime(revolver).

sol (X, Y, Z) :- murderer(X), weapon_of_crime(Y), place_of_crime(Z).
#maximize{W , sol : sol (X, Y, Z) , crime_record(X, W) , murdered(W)}.
#show sol /3.
    
```

# Solving combinatorial problems



Write a program which  
how the problem should be solved



Write (boolean, linear)  
constraints (*SAT, ILP, ...*)

**Declarative programming**

**Answer set programming.**

**Describe what you want to solve**

➤ **Problem = axioms & rules**

➤ **No need of algorithm**

# ASP logical rules : declarative programming

$$\underbrace{K \{ atom_1; \dots; atom_n \} L}_{\text{head}} \underbrace{:-}_{\text{"smiley"}} \underbrace{atom_{n+1} ; \dots; atom_r; not atom_{r+1}; \dots; not atom_s.}_{\text{body}}$$

**If** all terms on the **right side** are true,  
**then** at least **K** and at most **L** terms are true  
 on the **left side**.

**If** nothing on the **left side**,  
**then** always false.

**If** nothing on the **right side**,  
**then** always true.

$:- K\{atom_1, \dots atom_N\}L.$

$K\{atom_1, \dots atom_N\}L.$

**Optimisation rule**

**#maximize{W, atom(X): condition(X), W}.**

## High-level model language

- Propositional logics
- Model for negation

## Highly performant solving technics

- SAT-based and deductive-DB technics
- Decidable: no infinite loop

# Link with systems biology ?

**Integrative and systems biology is a very relevant field to challenge ASP technologies**

- Repair large-scale interaction graph with **branch and bound** solving heuristics (KR 2010)
- Scale metabolic network completion problem with **unsatisfiable core** solving strategy (LPNMR 2013)
- Design experiments with **incremental solving** (Frontiers 2015)
- Implement and benchmark **constrains propagators** (TPLP 2018)

**Linear constrains atoms**

$$\&sum\{a1*x1; \dots ; aI *xI \} \leq k$$

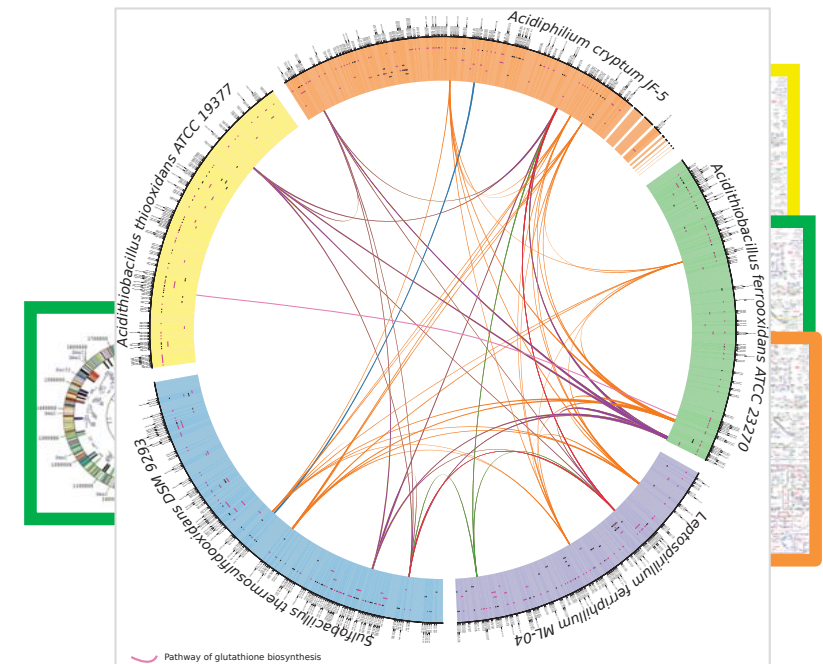
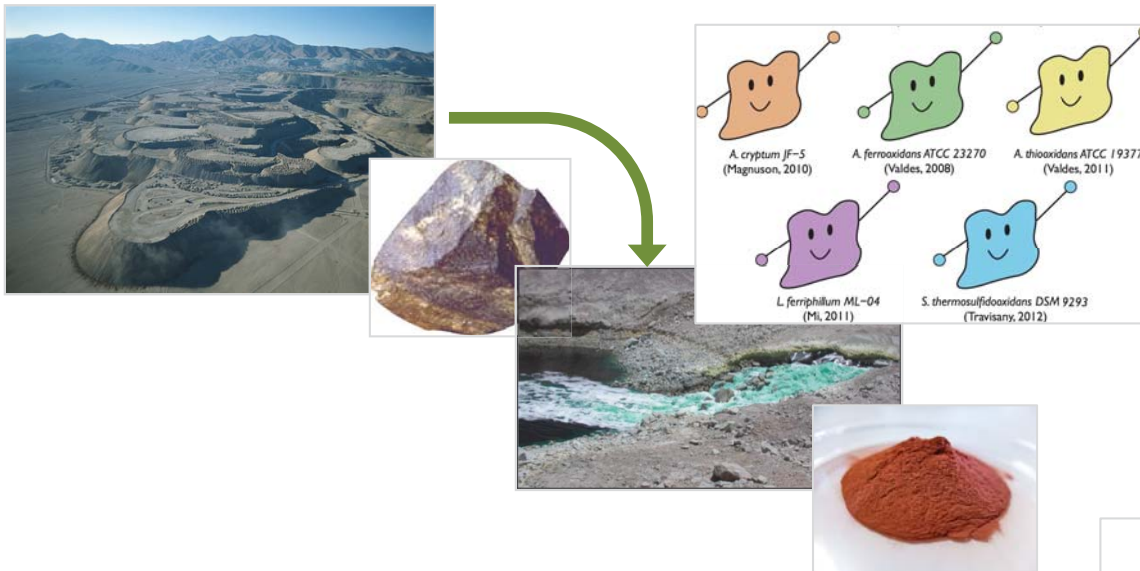
Problem statement  
& modelling



Solving heuristics  
& problem reformulation

# Application: extremophile mining consortium

Role of an **empirical taylor-made consortium** of bacteria in copper extraction from ore ?



« **NAD(H) biosynthesis** metabolic pathways of **A. Cryptum** complements metabolic functions spread between the five strains »

## ASP program

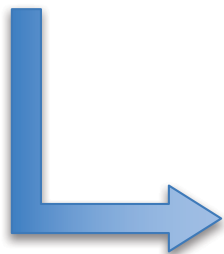
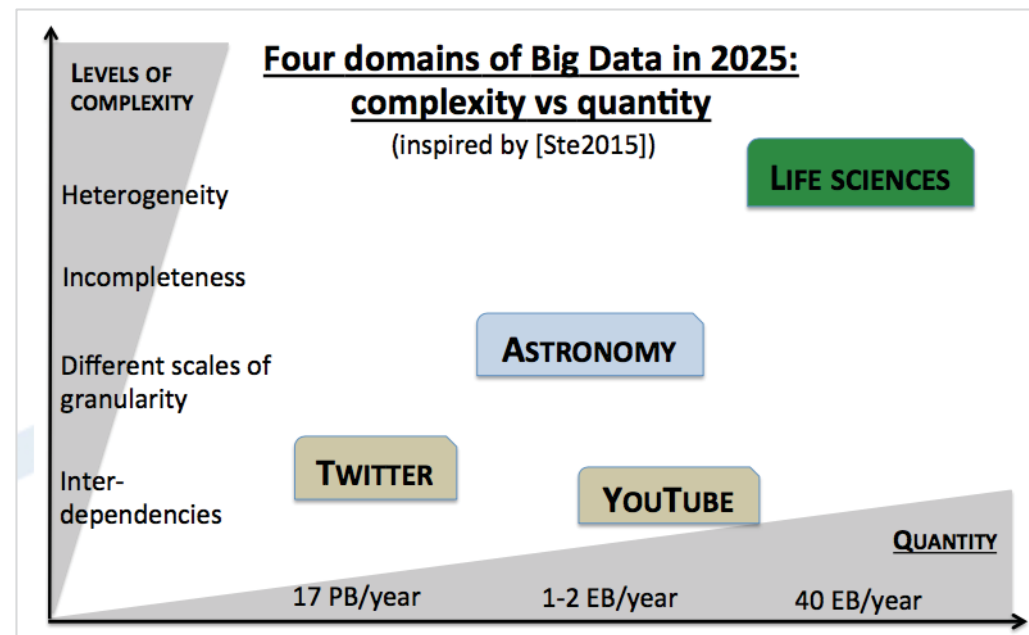
→ graph alignment / static modeling

→ chains of reactions explaining the capability of the consortium to produce the compounds (LPNRM'13, Microbiology open'15)

# BACK TO DYNAMICAL SYSTEMS

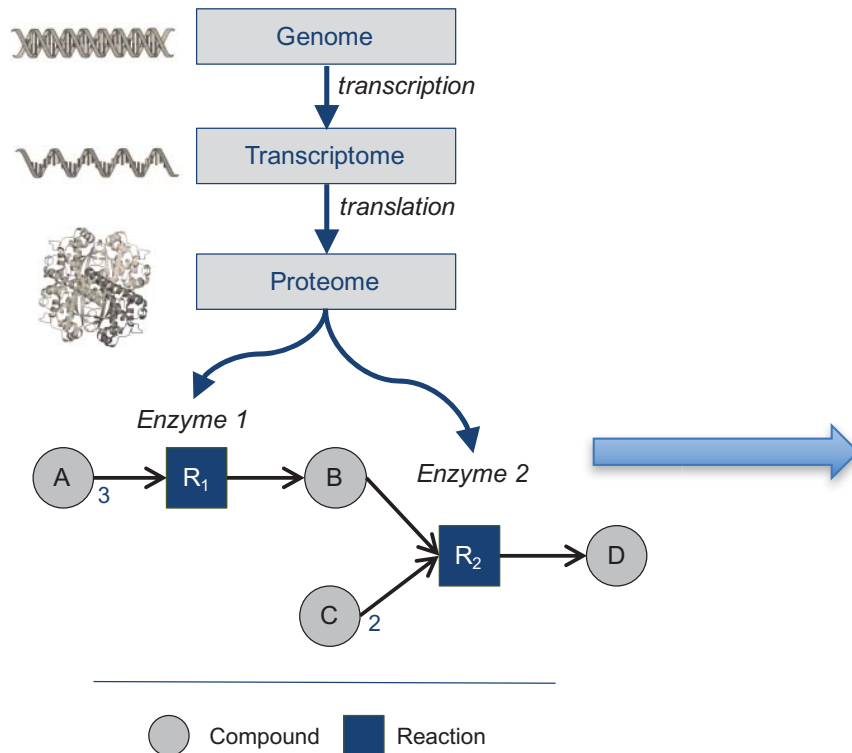
## Biological systems characteristics

- Large-scale
- Empirical laws
- Few data wrt the search space size

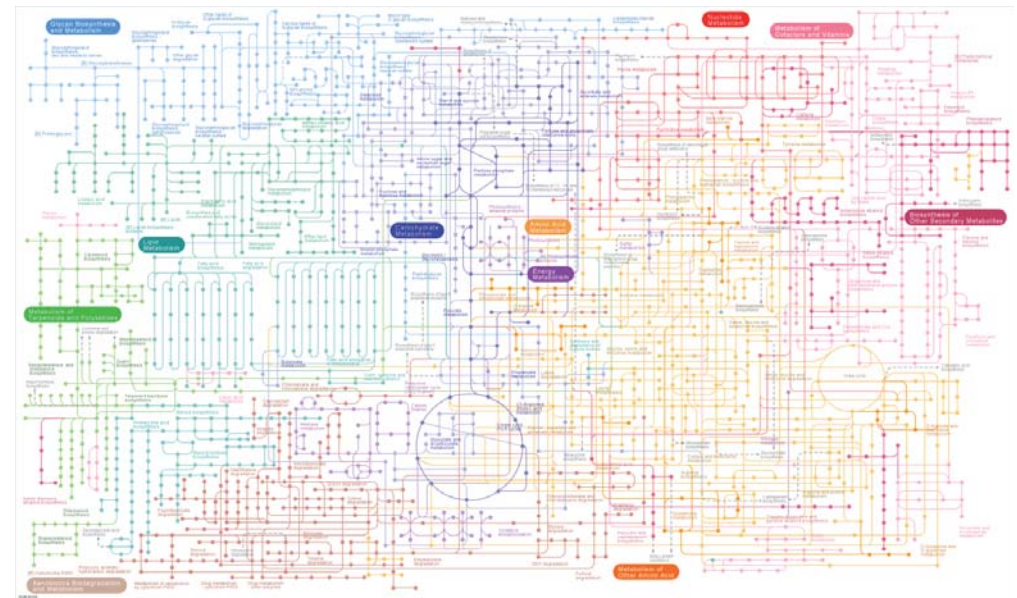


**Biological systems observed with omics data  
are not uniquely identifiable**

# Underlying tool : from genes to dynamical systems



1 genome  
 ⇒ 1 metabolic network  
 = bipartite directed graph

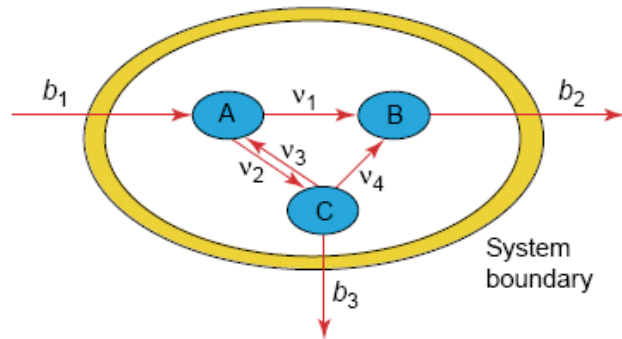


Link between genes  
 and functions

Large scale metabolic network

**All expected metabolic capabilities of an organism**

# How to model fluxes ?



$$\frac{dA}{dt} = -v_1 - v_2 + v_3 + b_1$$

$$\frac{dB}{dt} = v_1 + v_4 - b_2$$

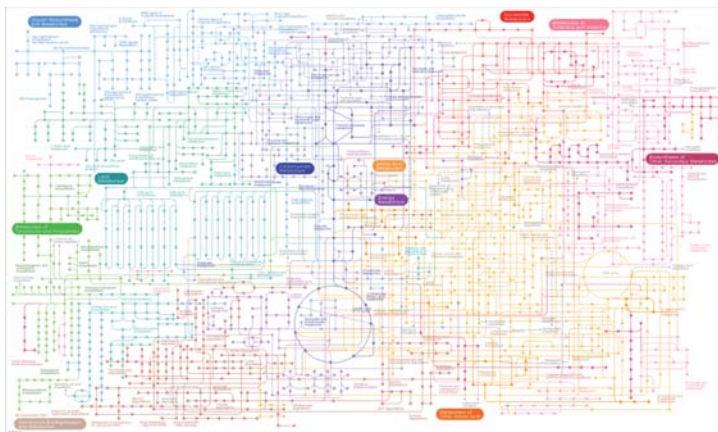
$$\frac{dC}{dt} = v_2 - v_3 - v_4 - b_3$$

$$\frac{dx}{dt} = S \cdot v(x)$$

$$v([substrat]) = Vm[Substrat] / (Km + [Substrat])$$

## Back to high school chemistry

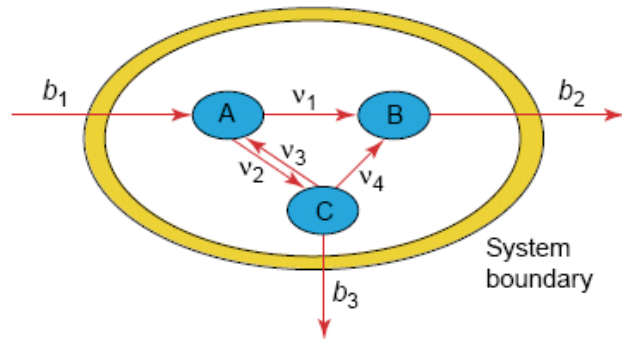
- Two parameters have to be estimated for each reaction



## Intractable in practice !

- Overapproximation of the dynamics

# Quasi-steady state hypothesis



$$\frac{dx}{dt} = S \cdot v(x) = 0$$

$$v([substrat]) = V_m [Substrat] / (K_m + [Substrat])$$

**= constant**

## Metabolic compounds do not accumulate

- Fluxes have constant values
- Fluxes are constrained by linear values
- The system optimises a global objective

*r* is active if

$$v_r > 0 \text{ and}$$

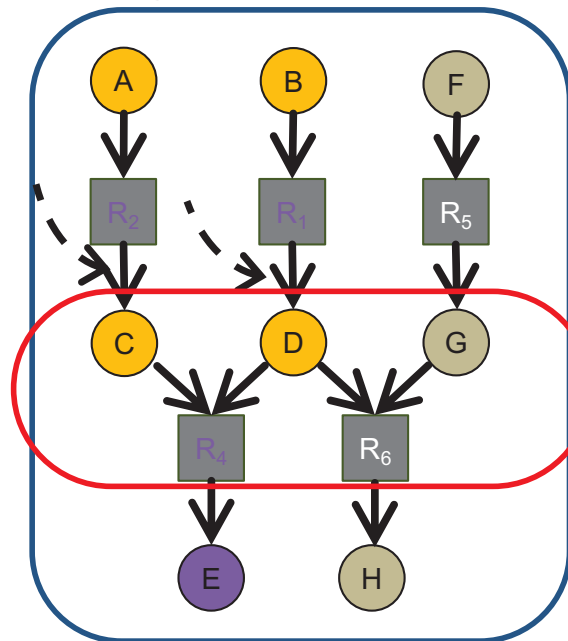
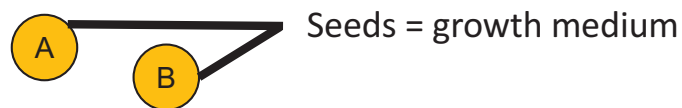
$$s.v = 0 \text{ and}$$

$$lb < v < ub$$




Replace kinetic constants by global optimisation hypotheses

# Growing phase hypothesis

**Functionality:** recursive graph-based semantics



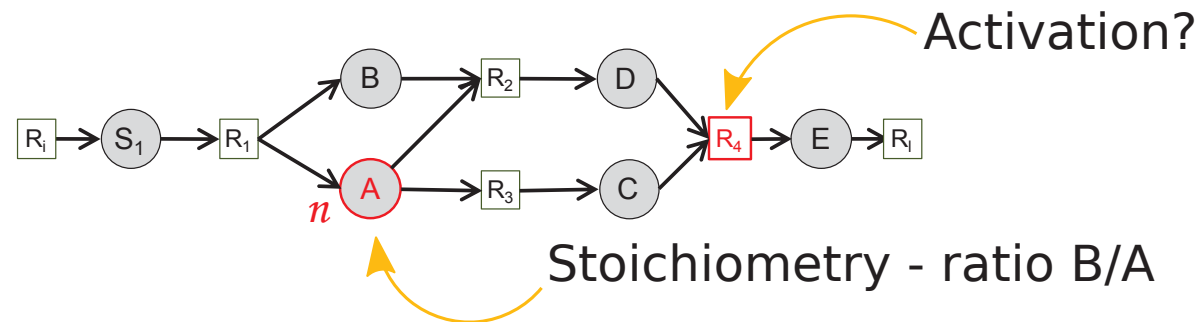
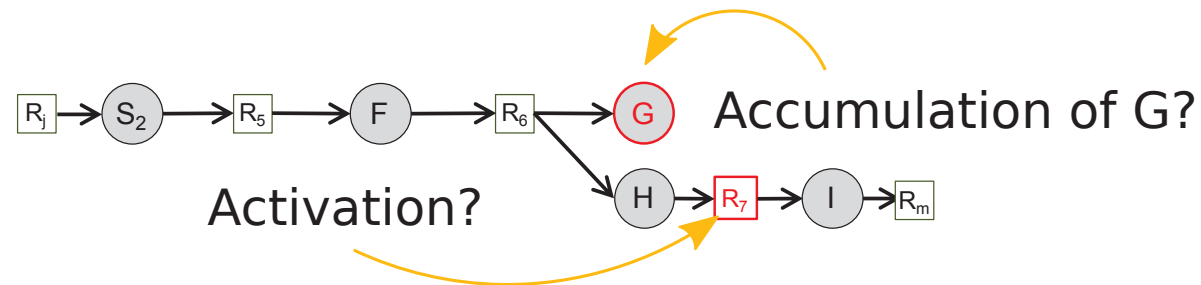
**“and” condition checked recursively**

-  Non-producible metabolite
-  Metabolite reachable from the seeds
-  Reaction

```
scope(M): - seed(M).
scope(M): - product(M, R), reaction(R), scope(Mi): reactant(Mi, R).
```

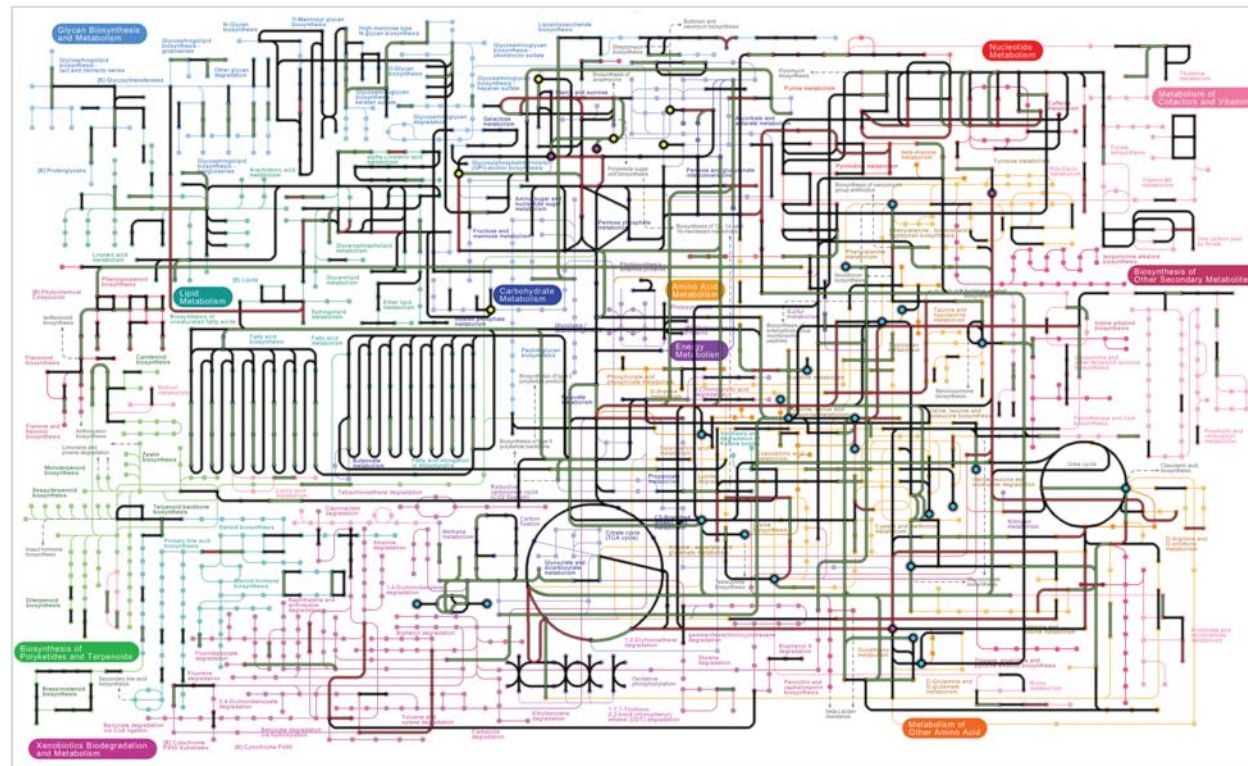
**Study paths in hypergraphs**

# Everything is a matter of choices

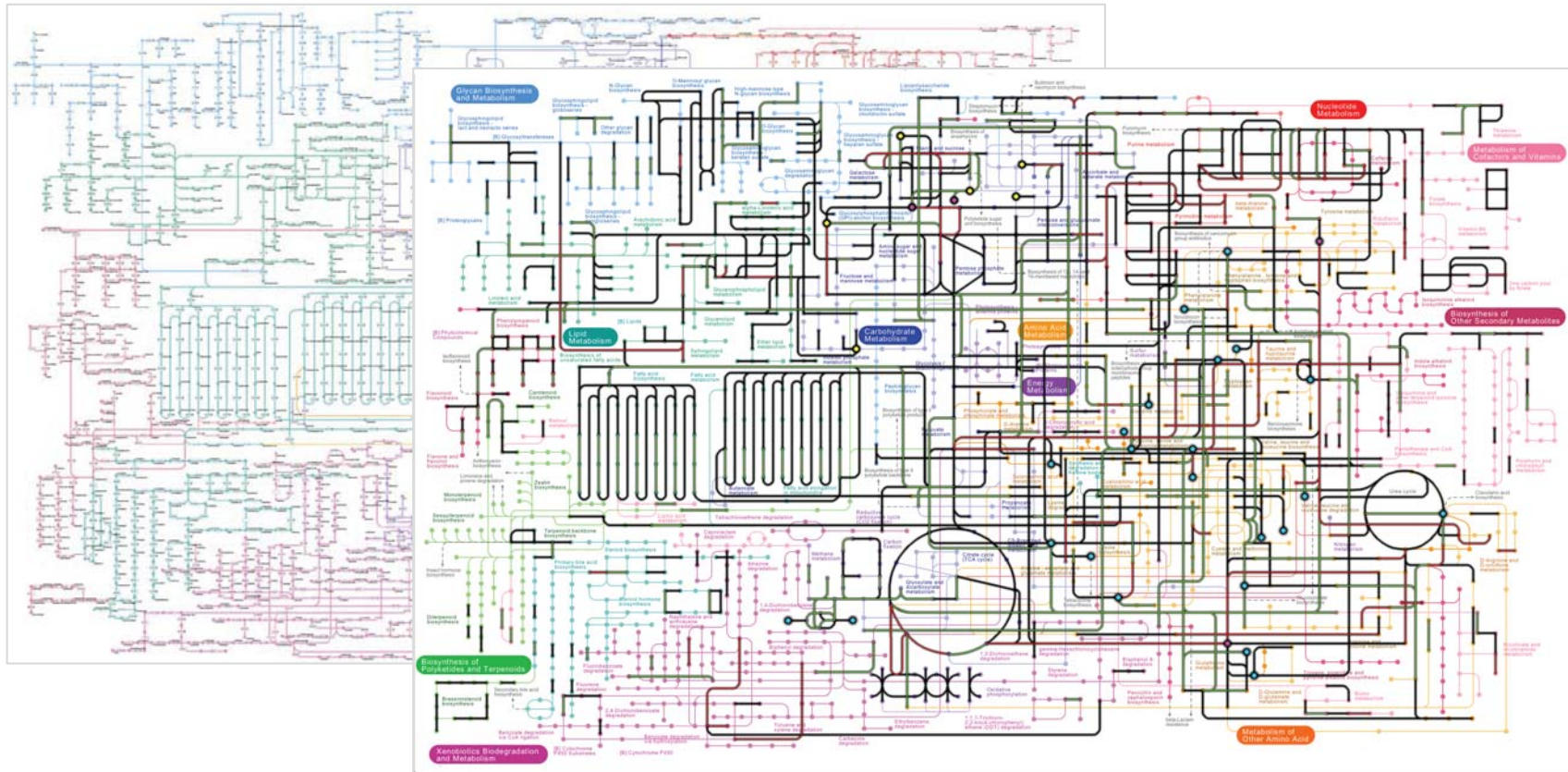


**The reaction status of the reactions is different according to the approximation**

- No choice but dealing with such overapproximation !
- Use the flexibility of ASP language to handle these questions



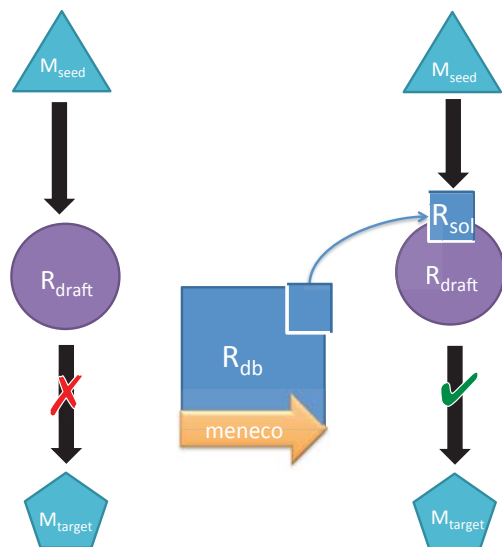
# Data incompleteness



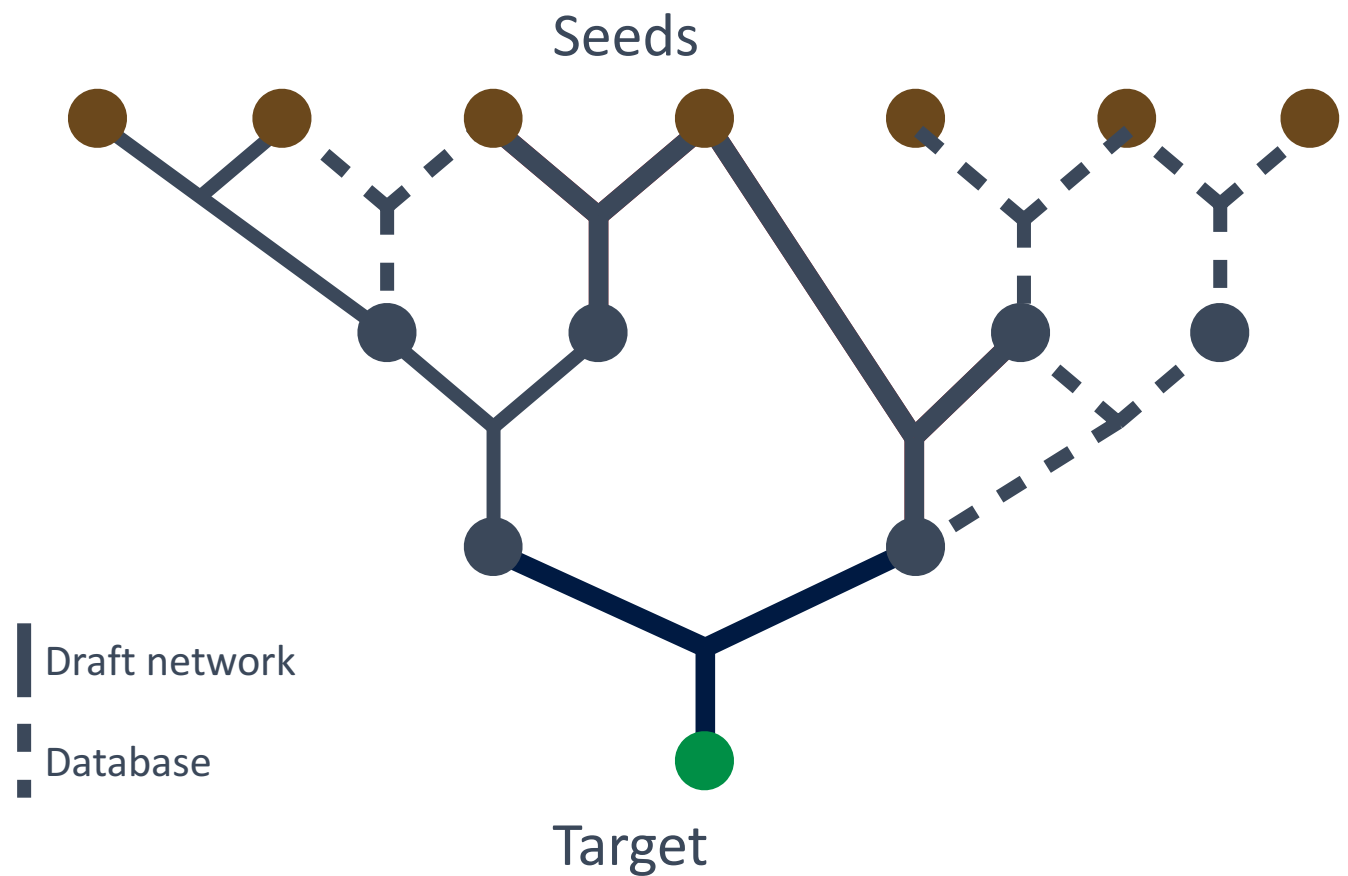
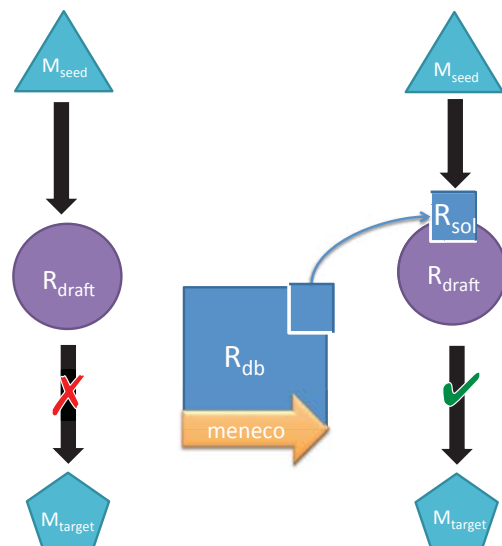
**Metabolic networks built from NGS sequencing**

➤ no possible biomass production.

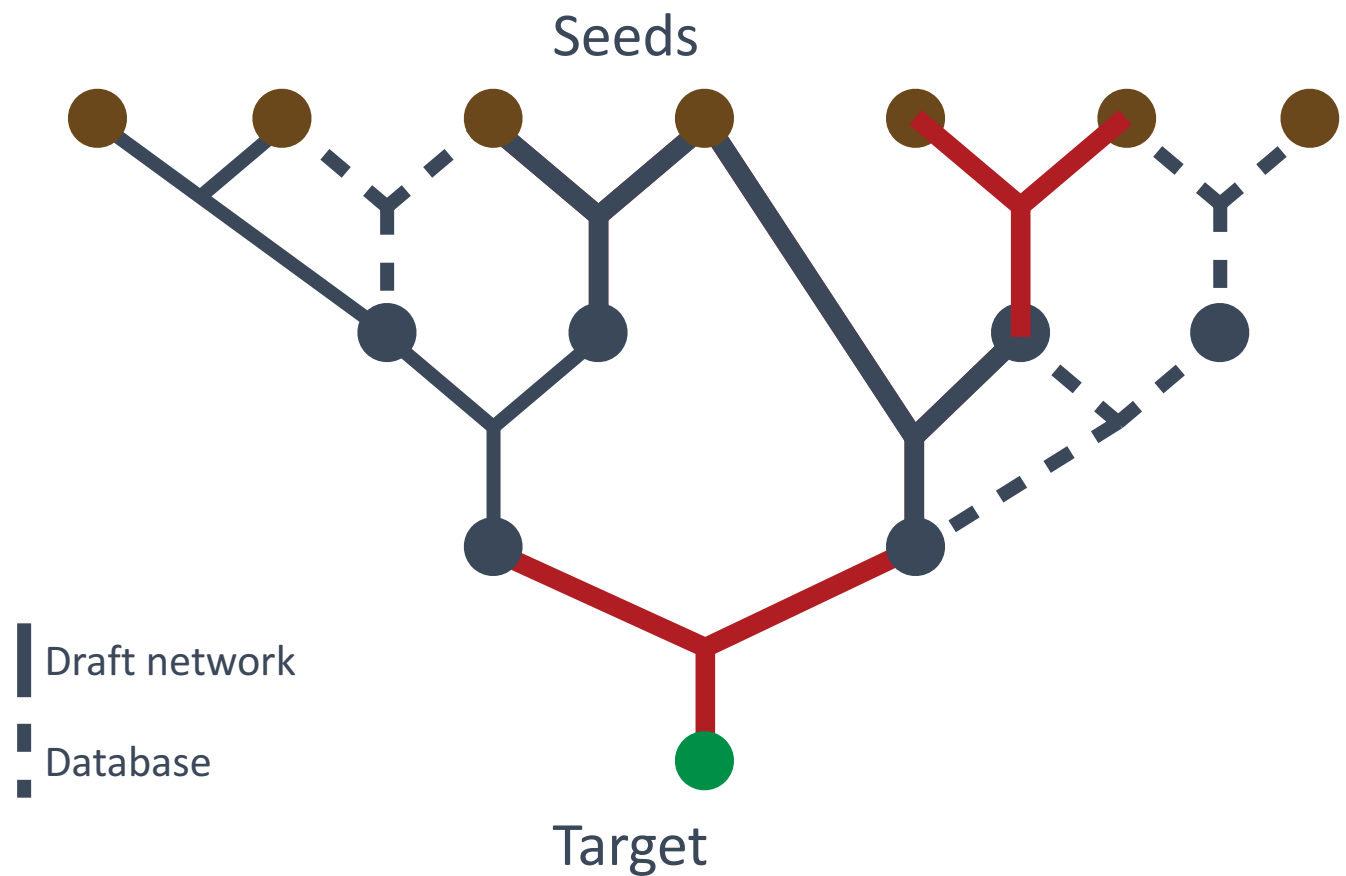
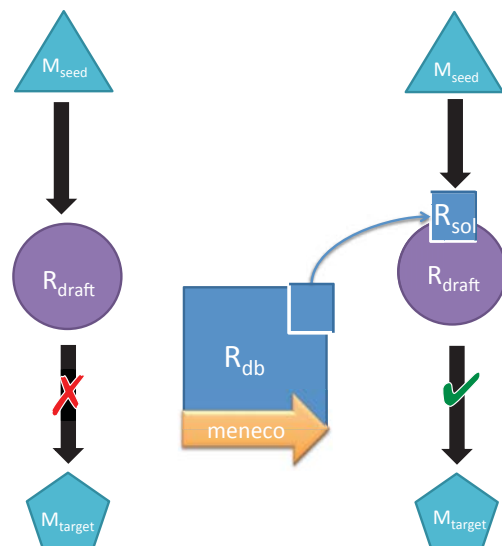
# Gapfilling a metabolic network (nutshell)



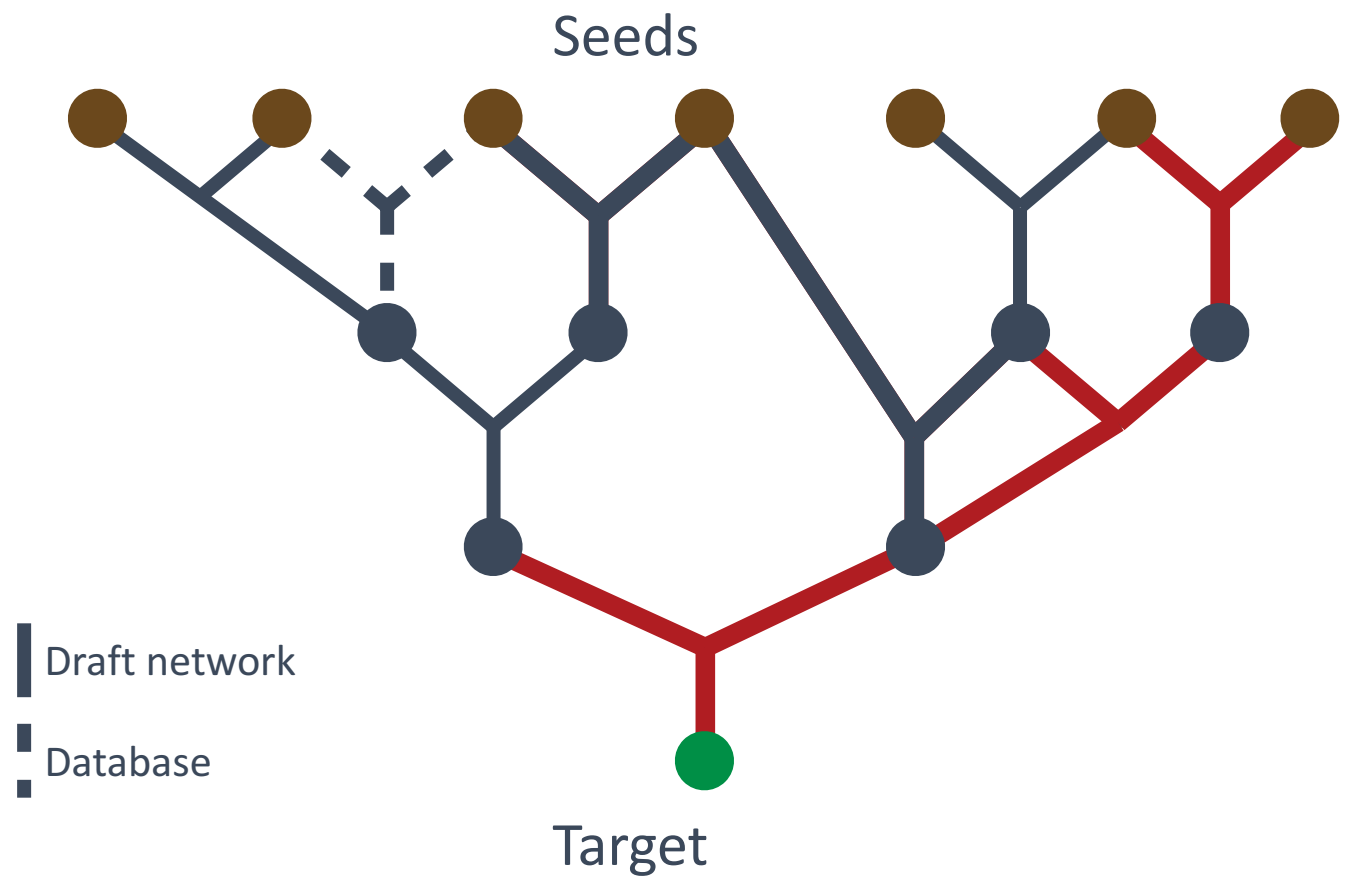
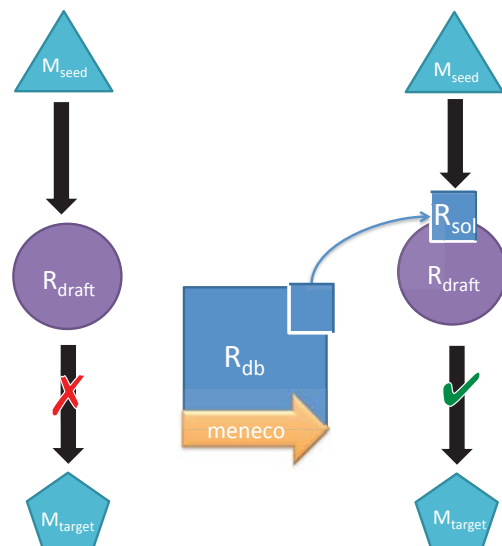
# Gapfilling a metabolic network (nutshell)



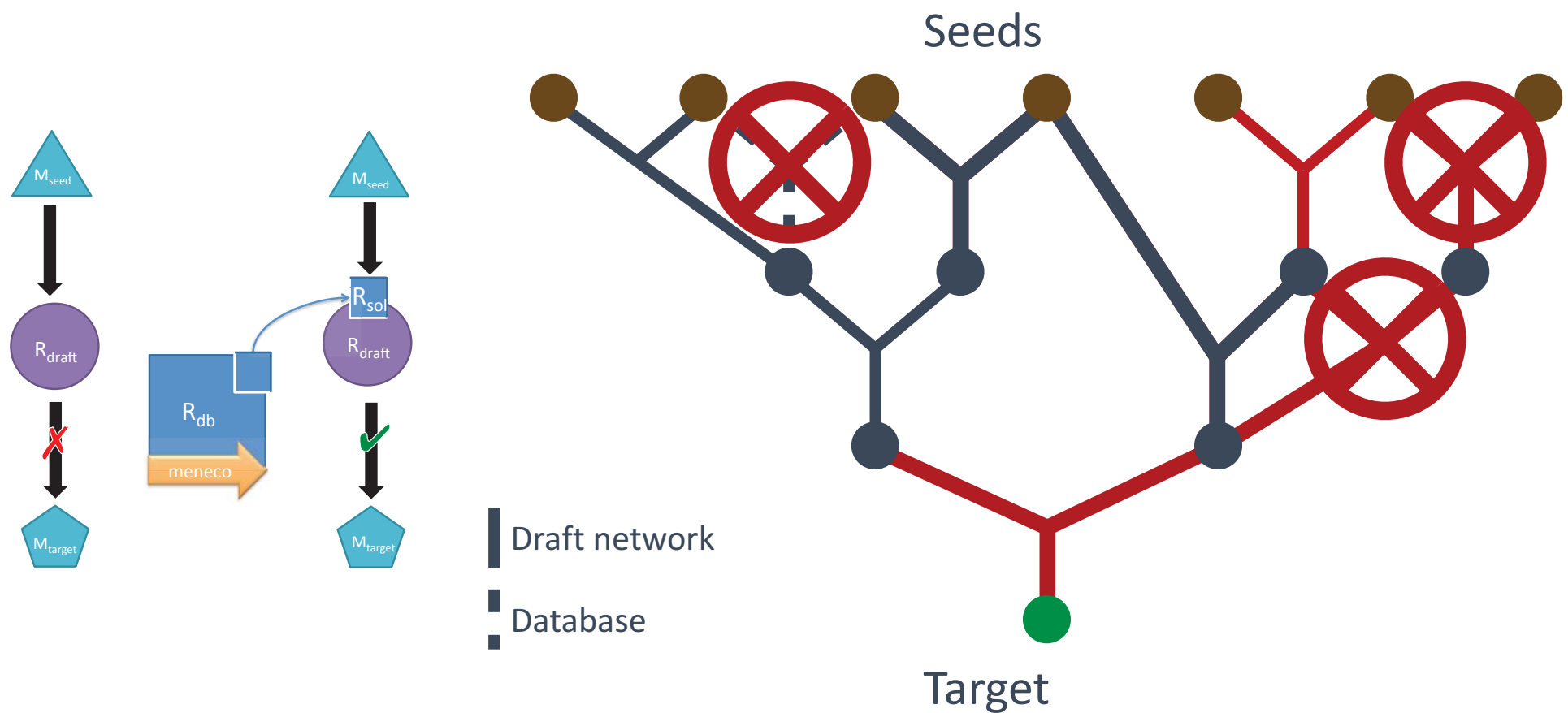
# Gapfilling a metabolic network (nutshell)



# Gapfilling a metabolic network (nutshell)



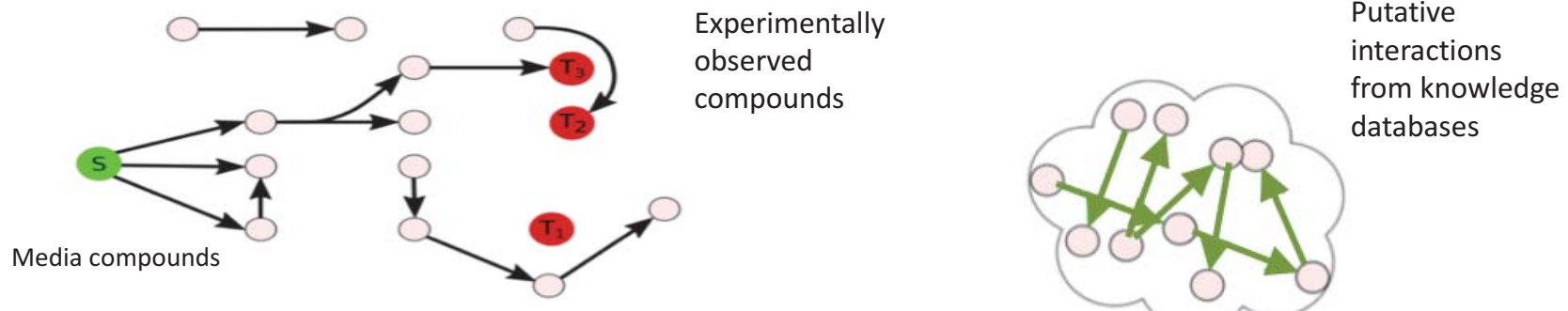
# Gapfilling a metabolic network (nutshell)



# Gapfilling a metabolic network

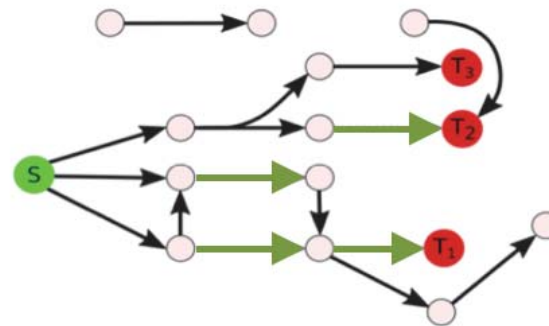
## What we have

- Graph with **non-accessible target components**
- **Knowledge database** of possible edges



## Gap-filling problem:

- Restore target accessibility
- Minimal number of reactions



$$\text{gapfilling}(S, R_T, G_1, G_{DB}) = \arg \min_{\{R_i..R_m\} \subset G_{DB}} \left( \frac{\text{size}(\text{reactants}(R_T) \setminus \text{scope}(G_1 \cup \{R_i..R_m\}))}{\text{size}\{R_i..R_m\}} \right)$$

# Meneco: ASP-based gap-filling for non-model organisms

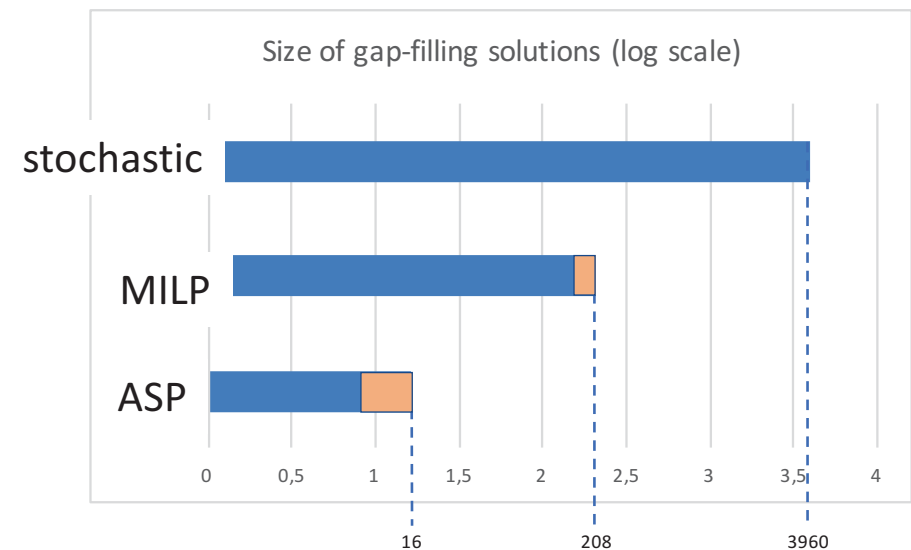
$$\text{Hybgapfilling}(S, R_T, G_1, G_{DB})$$

$$\arg \min_{\{R_i..R_m\} \subset G_{DB}} \left( \frac{\text{size}(\text{reactants}(R_T) \setminus \text{scope}(G_1 \cup \{R_i..R_m\}))}{\text{size}\{R_i..R_m\}} \right)$$

```
{reaction(r)}.
scope(M): - seed(M).
scope(M): - product(M, R), reaction(R), scope(M') : reactant(M', R).
:- target(T), not scope(T).
#minimize{ reaction(r) }.
```

**16 reactions in average are sufficient to restore degraded bacterial networks** (PLOS CB 2017)

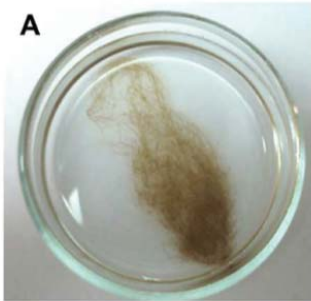
- MILP-based approaches required from 200 to 4000 reactions.



Benchmark of 10,800 bacterial networks

# Example of application

➤ **Genome: 1785 reactions, 1981 compounds**



Ectocarpus  
siliculosus

[Tapia2016]

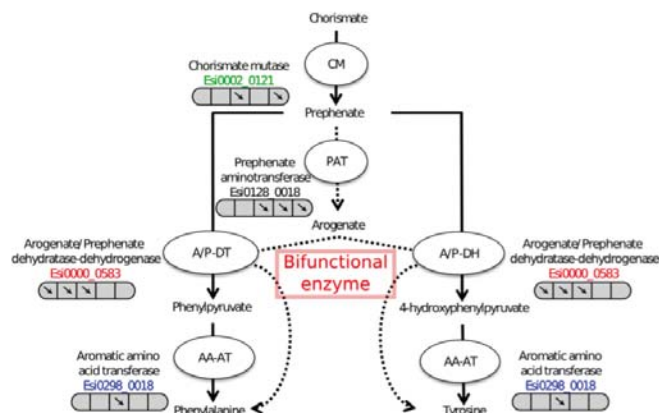
➤ **54 metabolites to produce:**

- 25 are graph-based producible
- None is FBA-based producible.

➤ **Gapfilling**

- MILP : 500 reactions (untractable)
- ASP: 50 reactions added to the network
  - Sufficient for fluxes
  - Manual curation

Proposed after manual curation



**New bifunctional role of a specific enzym**  
(Plant Journal 2015)

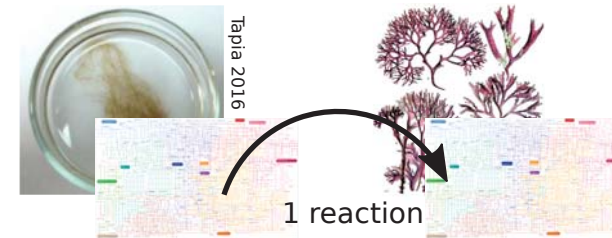
# Counter-example of application



Chondrus crispus

## Network analysis (G. Markov, SBR)

- 1943 reactions
- 149 reactions added by ASP
- **No way to produce biomass**



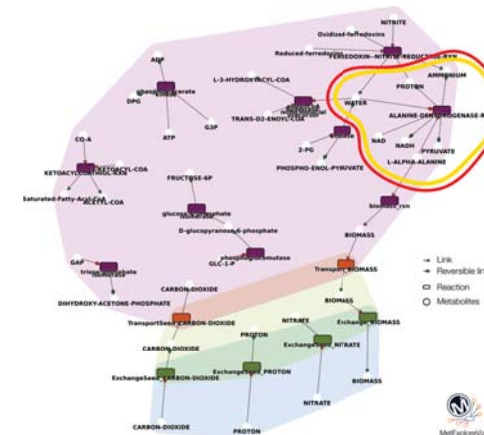
## New problem to be solved

- **Hybrid problem** (TPLP 2018)
- Constraint propagator
- Reduce the database

$\text{Hybgapfilling}(S, R_T, G_1, G_{DB}) =$

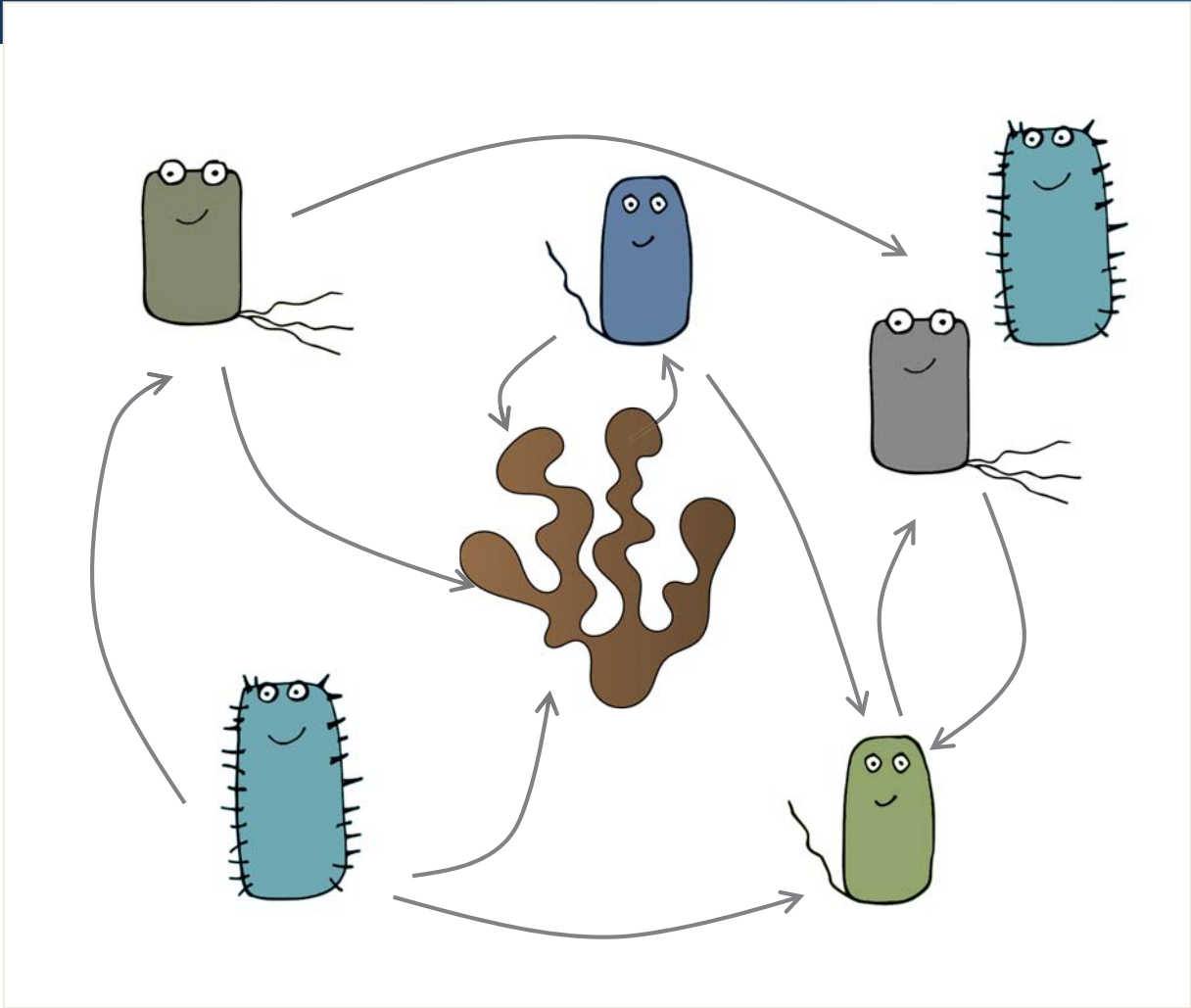
$$\arg \min_{\{R_i..R_m\} \subset G_{DB}} \left( \frac{\text{size}(\text{reactants}(R_T) \setminus \text{scope}(G_1 \cup \{R_i..R_m\}))}{\text{size}\{R_i..R_m\}} \right)$$

s.t.  $s.v = 0, v_{R_T} > 0, lb < v < ub$

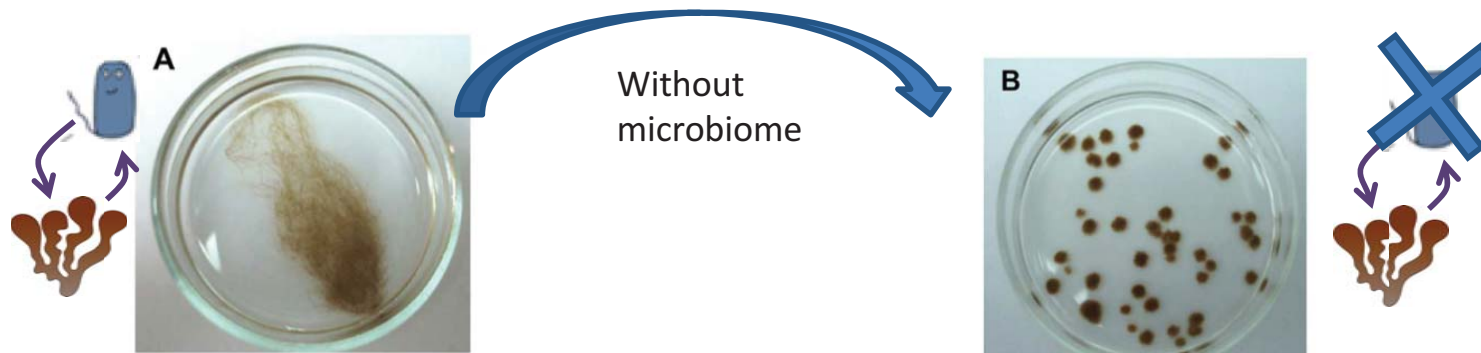


Essential reactions for alanine production in *CcrGem*

# STILL MORE COMPLEXITY



# Role of environmental bacteria ?

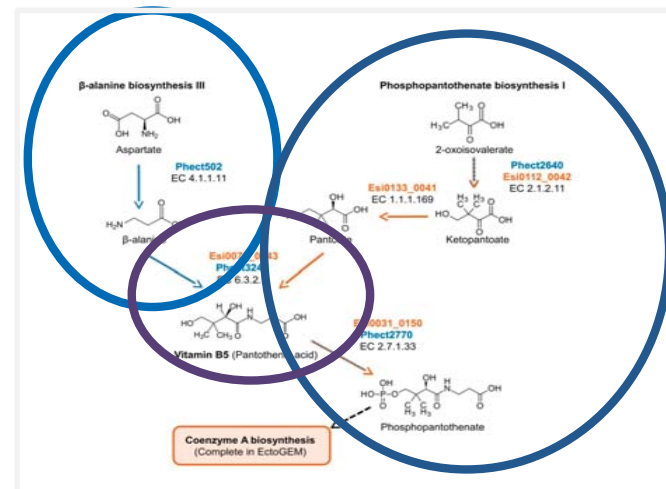
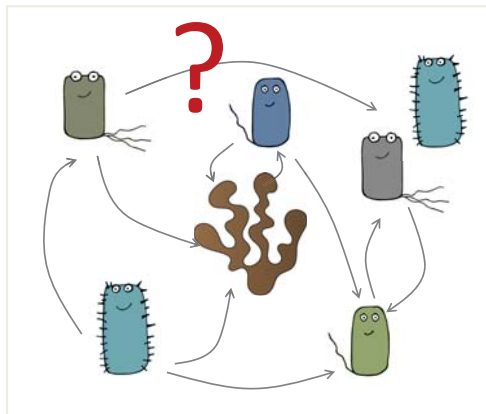


*Ectocarpus*

[Dittami2014, Tapia2016, Prigent2015]

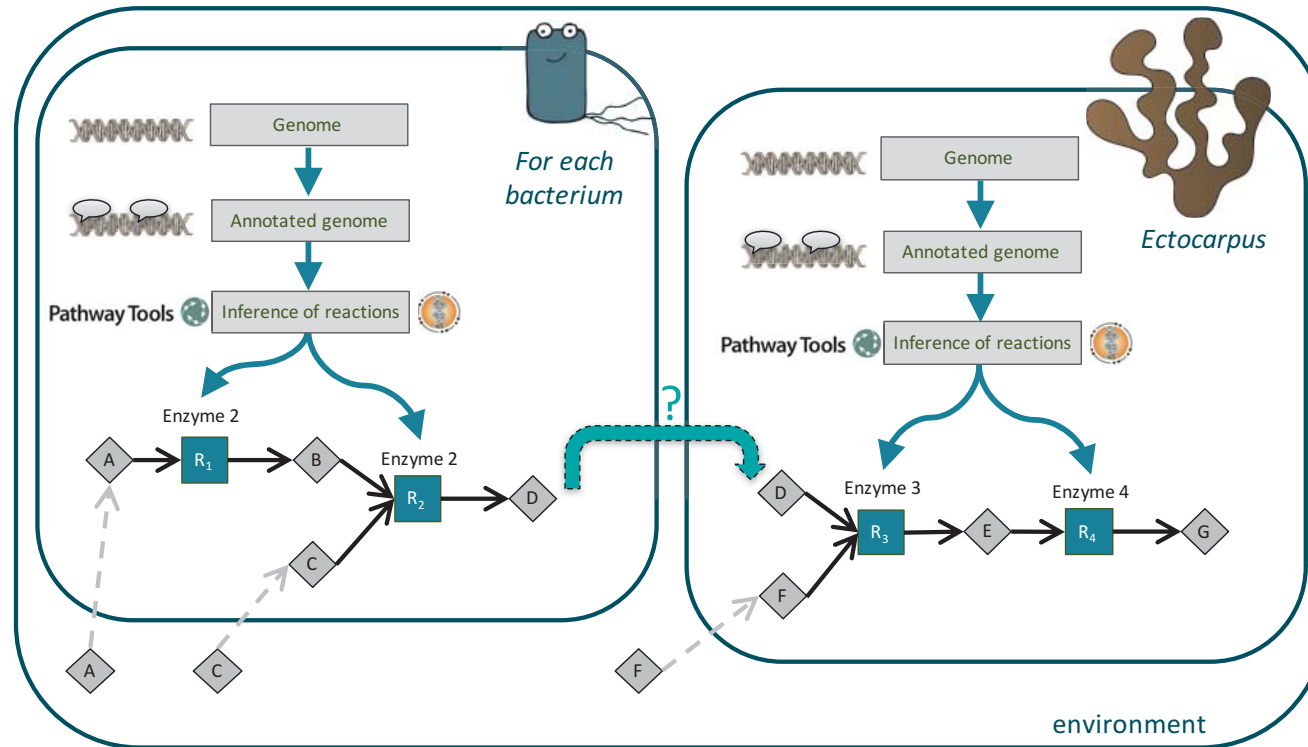


CNRS • SORBONNE UNIVERSITÉ  
Station Biologique  
de Roscoff



**Metabolism may be an explanation**  
(PLOS CB 2017)

# Systems ecology question



**Can we suggest compound exchanges that could restore the production of targeted compounds ?**

- New gap-filling problem !
- Steiner graph approach (Sagot team, 2017) or ASP implementation

# Scalability...

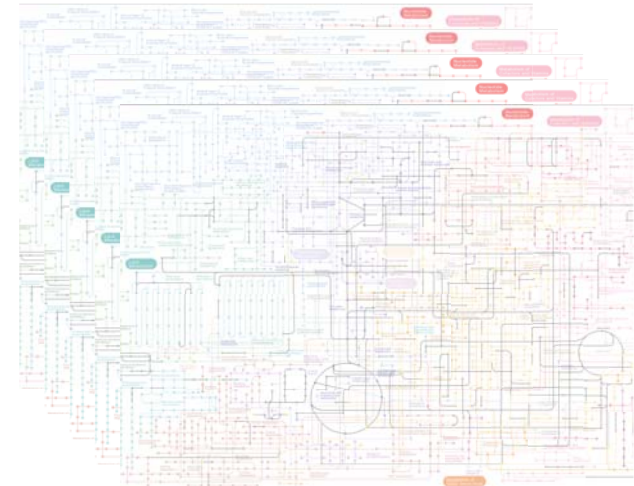
But... There are hundreds of bacteria in the environment



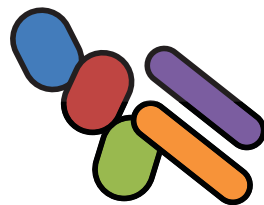
Marine biology



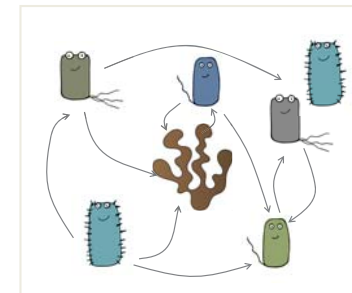
Hundreds of bacteria



Hundreds of Genome-scale models (GSMs)

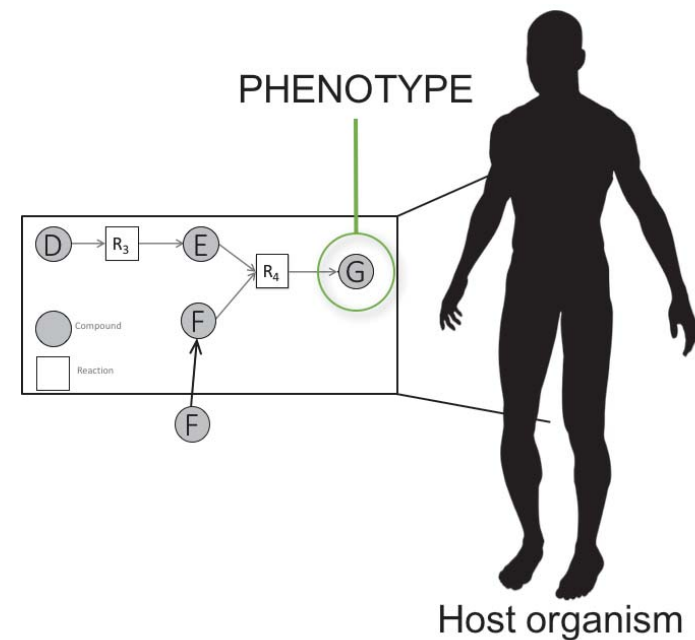


Happy few bacteria interact with the algae



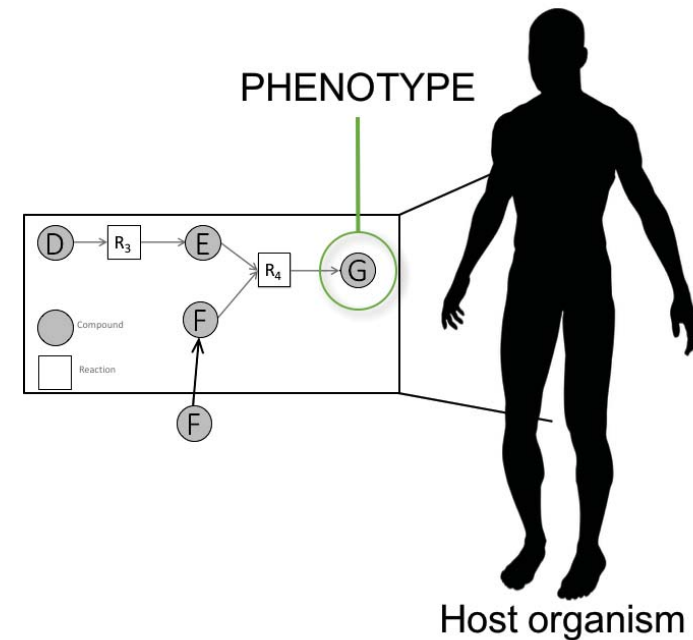
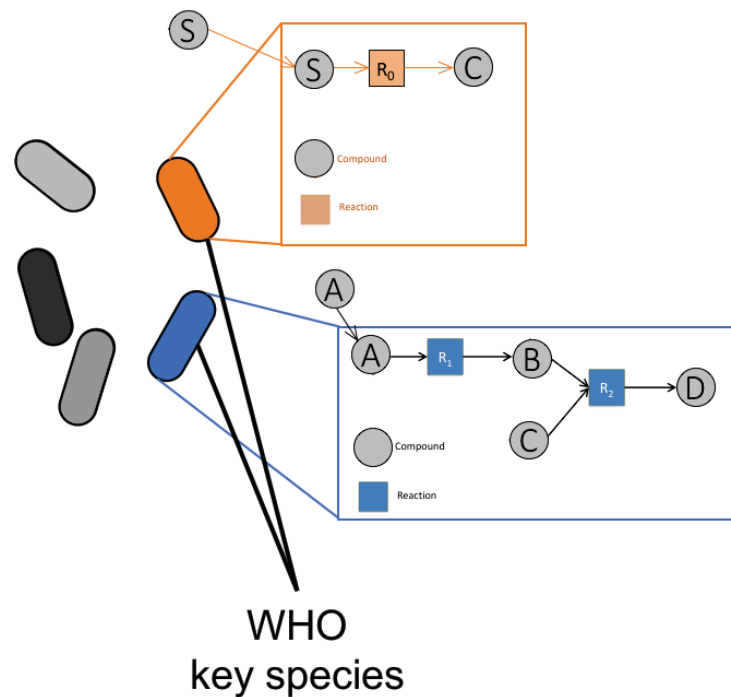
How to select communities within large microbiotas which explain the algal response to stress ?

# Selecting communities of interest within [large] microbiotas



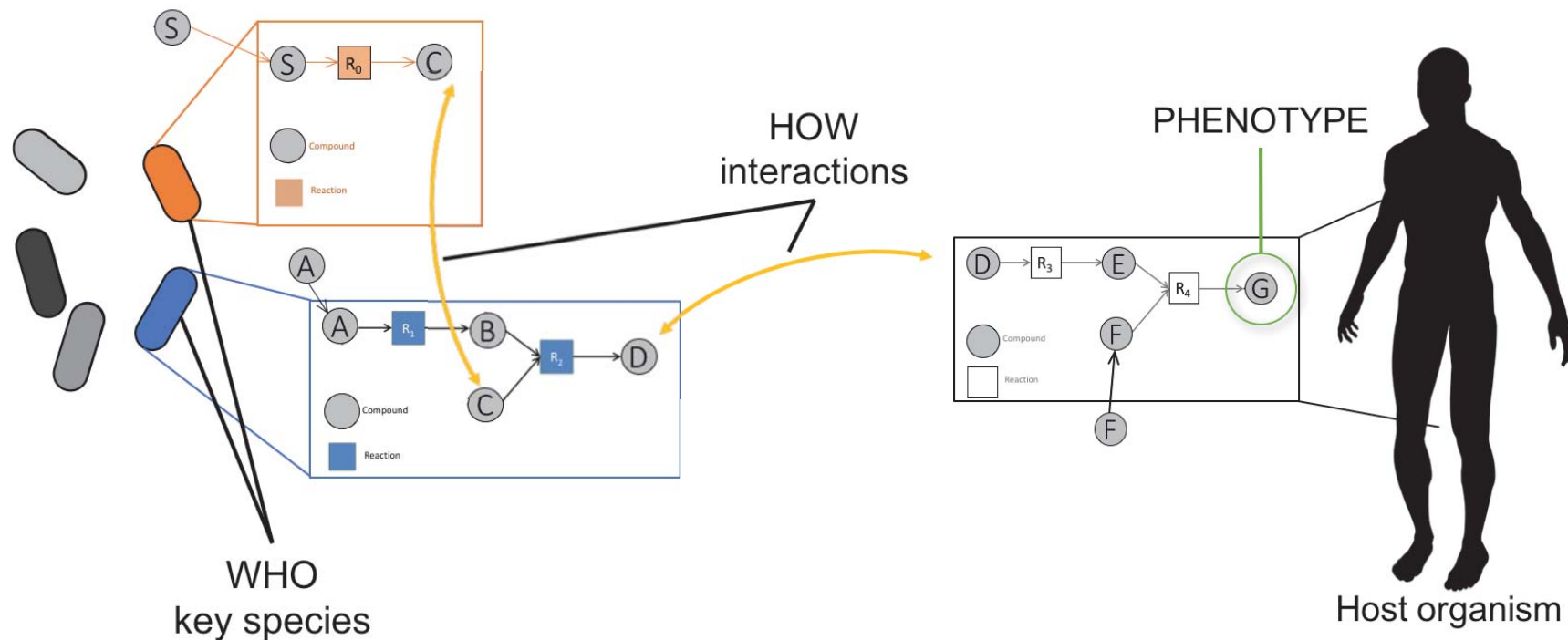
The “who”, “how” challenges of community selection

# Selecting communities of interest within [large] microbiotas



The “who”, “how” challenges of community selection

# Selecting communities of interest within [large] microbiotas

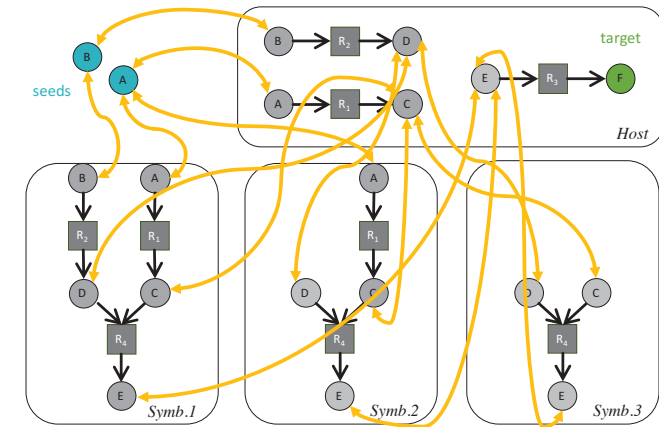


The “who”, “how” challenges of community selection

# Complexity

## Community selection problem

- Switch from hundreds of symbiots to 3 or 4
- Pinpoint a few number of putative cross-feedings



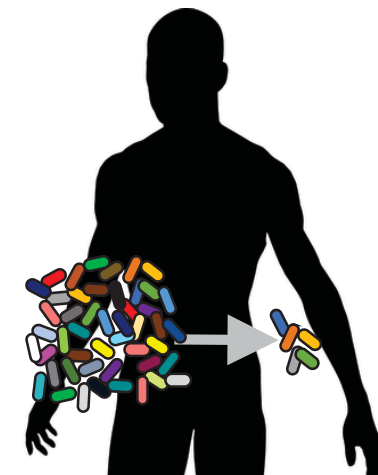
499,177 combinations of  
<6 exchanges

$$\text{Comsel}(S, T, G_1 \dots G_n) = \arg \min_{\{exchg(G_{i_1} \dots G_{i_L}) \subset \{G_1 \dots G_n\}\}} \left( \begin{array}{l} size(T \setminus MBscope(G_{i_1} \dots G_{i_L})) \\ size\{\varepsilon \subset exchg(G_{i_1} \dots G_{i_L}) \mid \\ T \cap CPscope(G_{i_1} \dots G_{i_L}, \varepsilon, S) = \\ T \cap MBscope(G_{i_1} \dots G_{i_L}, S)\} \end{array} \right)$$

- depends on the number of hyperarcs

## Size of the search space

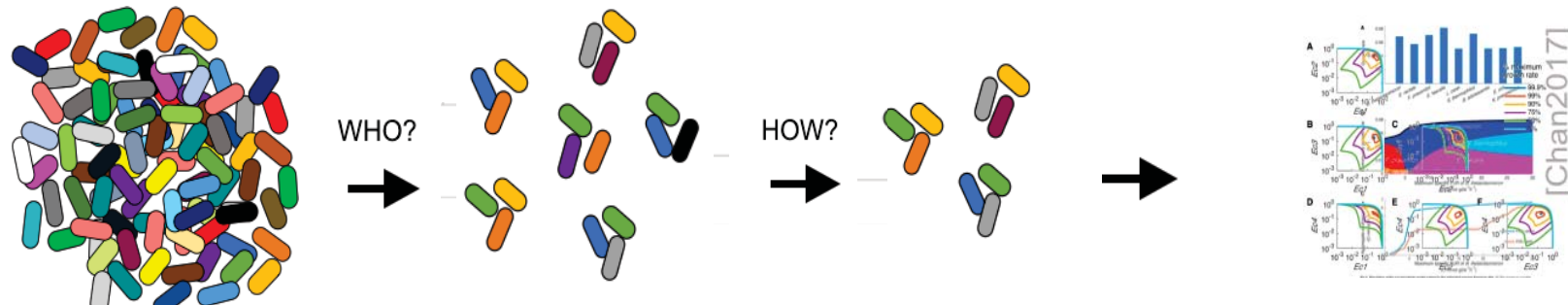
- depends on the number of symbionts



$1.62 \cdot 10^{81}$  combinations of  
<10 exchanges

**Highly combinatorial problem**

# Two-step optimization procedure



## Heuristics for the community selection problem

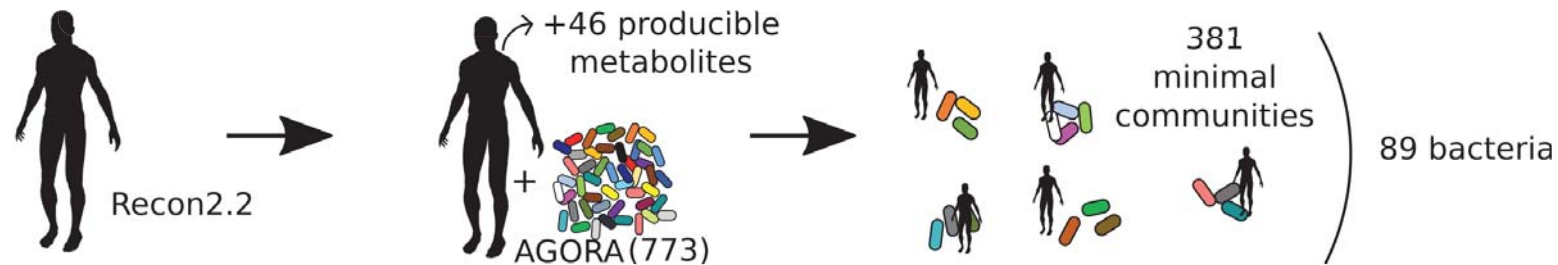
- **Who problem.**
  - Get rid of boundaries and select all minimal symbiot families
- **How problem.**
  - Sort the selected families according to the number of exchanges
- **Manual curation.**
  - Ask your favorite biologist to select the final one

$$\begin{aligned} & \text{mxdbagCnity}(S, T, G_1..G_N) \\ &= \arg \min_{\{G_{i_1}..G_{i_L}\} \subset \{G_1..G_N\}} \left( \frac{\text{size}(T \setminus \text{mxdbagScope}(G_{i_1}..G_{i_L}, S))}{\text{size}\{G_{i_1}..G_{i_L}\}} \right) \end{aligned}$$

$$\begin{aligned} & \text{cptCnity}(S, T, G_1..G_N) \\ &= \arg \min_{\substack{\{G_{i_1}..G_{i_L}\} \\ \subset \{G_1..G_N\}}} \left( \begin{aligned} & \text{size}(T \setminus \text{mxdbagScope}(G_{i_1}..G_{i_L}, S)), \\ & \text{size}\{G_{i_1}..G_{i_L}\}, \\ & \text{size}\{\mathcal{E} \subset \text{exchg}(G_{i_1}..G_{i_L})| \\ & \quad T \cap \text{cptScope}(G_{i_1}..G_{i_L}, \mathcal{E}, S) \\ & \quad = T \cap \text{mxdbagScope}(G_{i_1}..G_{i_L}, S)\} \end{aligned} \right) \end{aligned}$$

# Validation/benchmarking on human microbiome project

Context of the study [Swainston et al., 2016] [Magnúsdóttir et al., 2016]



Recon2.2

+46 producible metabolites

+ AGORA (773)

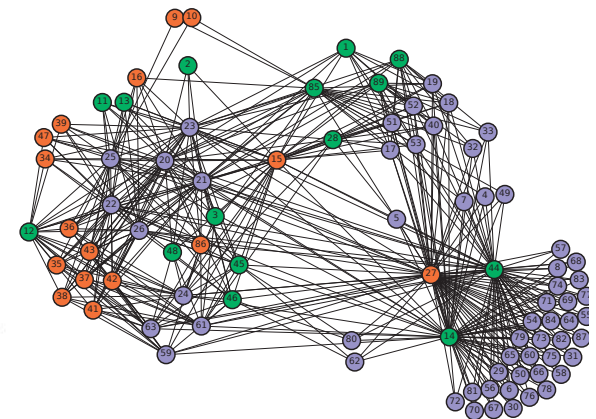
381 minimal communities

89 bacteria

<i>Altiplex fragilis</i> DSM 17242	4	<i>Chlorobacterium</i> sp. ATCC 35061	40
<i>Altiplex industrialis</i> ATCC 12060	5	<i>Pseudomonas putidus</i> DSM 18416	40
<i>Altiplex mediterranea</i> DSM 19147	6	<i>Pseudomonas putidus</i> DSM 23871	51
<i>Altiplex shufui</i> WAIL 4001	7	<i>Pseudomonas</i> sp. ATCC 11840	52
<i>Altiplex</i> sp. ATCC 61784	8	<i>Pseudomonas</i> sp. ATCC 27077	52
<i>Bacteroides</i> sp. ATCC 25411	9	<i>Staphylococcus aureus</i> ATCC 25922	53
<i>Bacteroides</i> clausii DSM 19493	10	<i>Staphylococcus aureus</i> ATCC 30066	53
<i>Bacteroides</i> fecis DSM 19494	11	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fecis DSM 19494	12	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	13	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	14	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	15	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	16	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	17	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	18	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	19	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	20	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	21	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	22	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	23	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	24	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	25	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	26	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	27	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	28	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	29	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	30	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	31	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	32	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	33	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	34	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	35	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	36	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	37	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	38	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	39	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	40	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	41	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	42	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	43	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	44	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	45	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	46	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	47	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	48	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	49	<i>Staphylococcus aureus</i> ATCC 23344	54
<i>Bacteroides</i> fragilis DSM 19494	50	<i>Staphylococcus aureus</i> ATCC 23344	54

<i>Bacillus cereus</i> AH187 F4810 72	9
<i>Bacillus cereus</i> G942	10
<i>Bacillus mycoides</i> DSM 2048	15
<i>Bacillus timonensis</i> JC401	16
<i>Brevibacillus brevis</i> DSM 100599	27
<i>Caldwellia formosa</i> DSM 4568	34
<i>Citrobacter freundii</i> Y19	35
<i>Citrobacter</i>	36
<i>Clostridium</i>	37
<i>Clostridium</i>	38
<i>Clostridium</i>	39
<i>Clostridium</i>	40
<i>Clostridium</i>	41
<i>Clostridium</i>	42
<i>Clostridium</i>	43
<i>Clostridium</i>	44
<i>Clostridium</i>	45
<i>Clostridium</i>	46
<i>Clostridium</i>	47
<i>Clostridium</i>	48
<i>Clostridium</i>	49
<i>Clostridium</i>	50
<i>Clostridium</i>	51
<i>Clostridium</i>	52
<i>Clostridium</i>	53
<i>Clostridium</i>	54
<i>Clostridium</i>	55
<i>Clostridium</i>	56
<i>Clostridium</i>	57
<i>Clostridium</i>	58
<i>Clostridium</i>	59
<i>Clostridium</i>	60
<i>Clostridium</i>	61
<i>Clostridium</i>	62
<i>Clostridium</i>	63
<i>Clostridium</i>	64
<i>Clostridium</i>	65
<i>Clostridium</i>	66
<i>Clostridium</i>	67
<i>Clostridium</i>	68
<i>Clostridium</i>	69
<i>Clostridium</i>	70
<i>Clostridium</i>	71
<i>Clostridium</i>	72
<i>Clostridium</i>	73
<i>Clostridium</i>	74
<i>Clostridium</i>	75
<i>Clostridium</i>	76
<i>Clostridium</i>	77
<i>Clostridium</i>	78
<i>Clostridium</i>	79
<i>Clostridium</i>	80
<i>Clostridium</i>	81
<i>Clostridium</i>	82
<i>Clostridium</i>	83
<i>Clostridium</i>	84
<i>Clostridium</i>	85
<i>Clostridium</i>	86
<i>Clostridium</i>	87
<i>Clostridium</i>	88
<i>Clostridium</i>	89
<i>Clostridium</i>	90
<i>Clostridium</i>	91
<i>Clostridium</i>	92
<i>Clostridium</i>	93
<i>Clostridium</i>	94
<i>Clostridium</i>	95
<i>Clostridium</i>	96
<i>Clostridium</i>	97
<i>Clostridium</i>	98
<i>Clostridium</i>	99
<i>Clostridium</i>	100

<i>Acinetobacter haemolyticus</i> NIPH 261	1
<i>Acinetobacter junii</i> SH205	2
<i>Acinetobacter pittii</i> ANC 4052	3
<i>Bacillus clausii</i> KSM K16	13
<i>Bacillus endophyticus</i> 2102	1
<i>Bacillus halodurans</i> C 125	1
<i>Bacillus mojavicus</i>	1
<b>Cluster 3:</b>	
<b>Firmicutes &amp; Proteobacteria</b>	
<i>Abiotriton granum</i> granum NCG 173078	40
<i>Brachyspira</i> 1 1 55	40
<i>Methylobacterium radiotolerans</i> JCM 2831	40
<i>Serratia liquefaciens</i> ATCC 27592	80
<i>Vibrio fluvialis</i> 560	80
<i>Vibrio furnissii</i> NCTC 11218	80



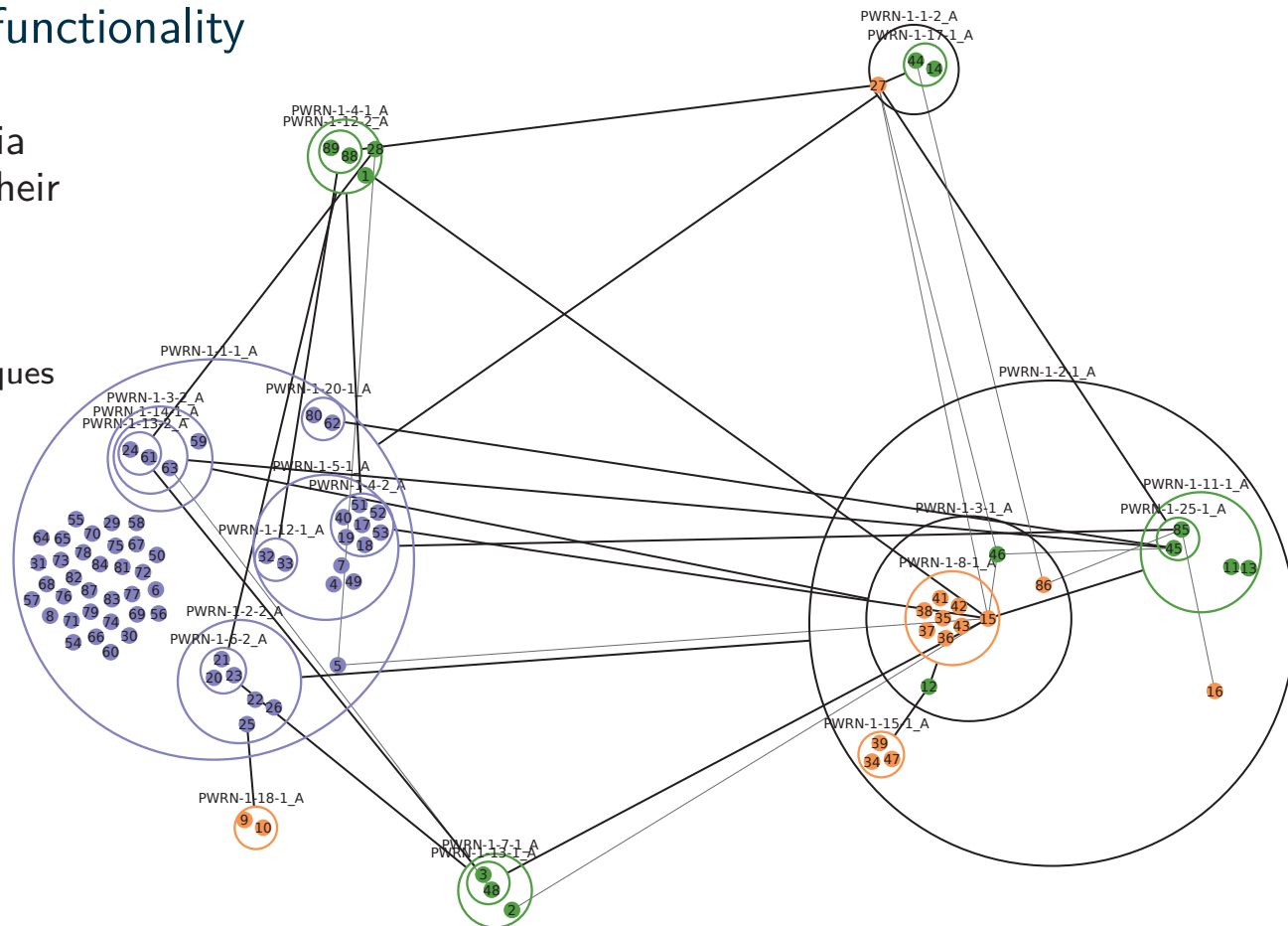
C. Frioux's Thesis. ECCB 2018

# Validation/benchmarking on human microbiome project

## Association of bacteria & functionality

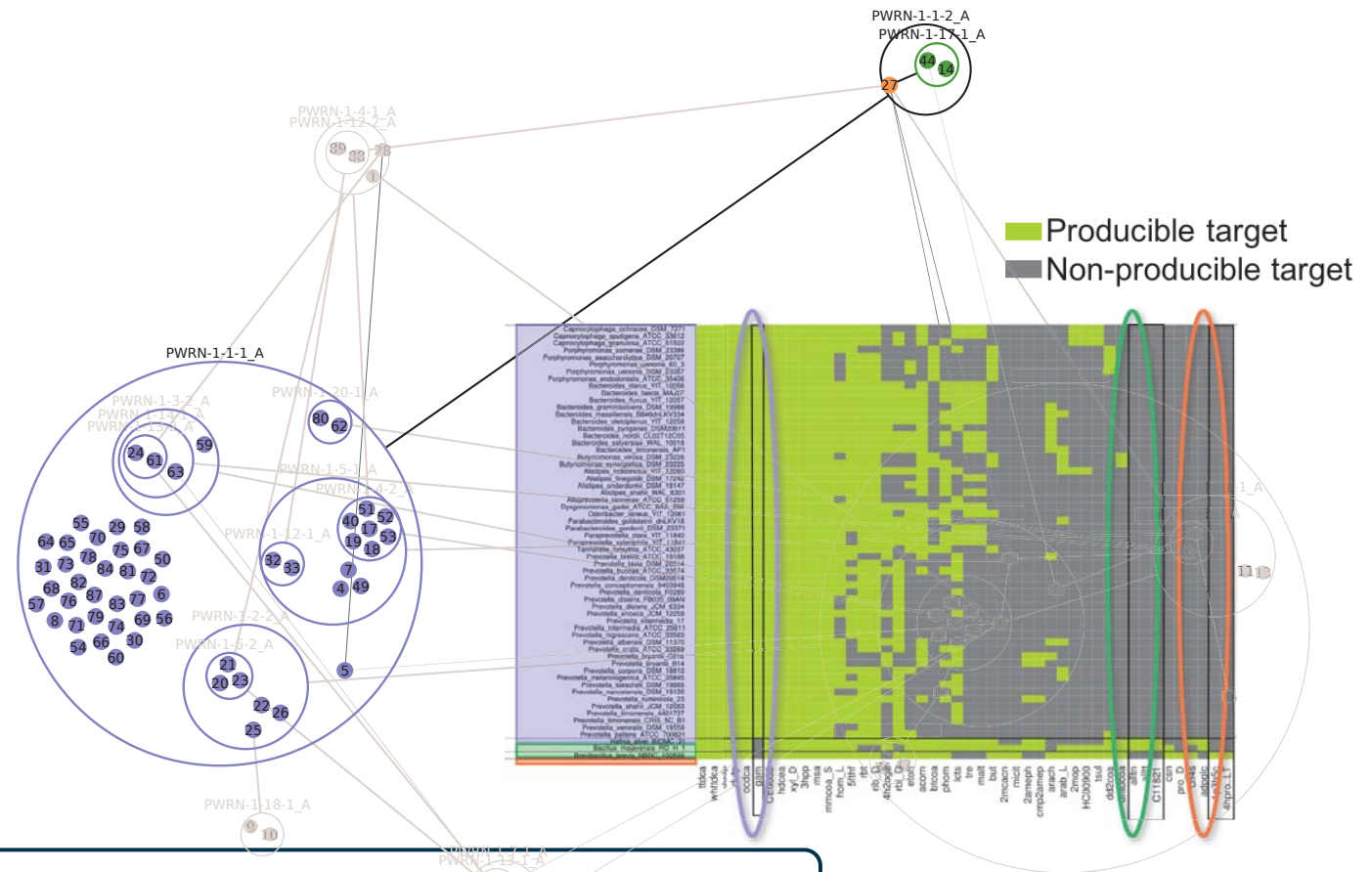
- Groups of equivalent bacteria in clusters with respect to their associations [Bourneuf et al., 2017]

- **Powernodes:** groups of bacteria, parts of bicliques
- **Poweredges:** connect bicliques



# Validation/benchmarking on human microbiome project

- Producibility of individual targets explains the communities → screening

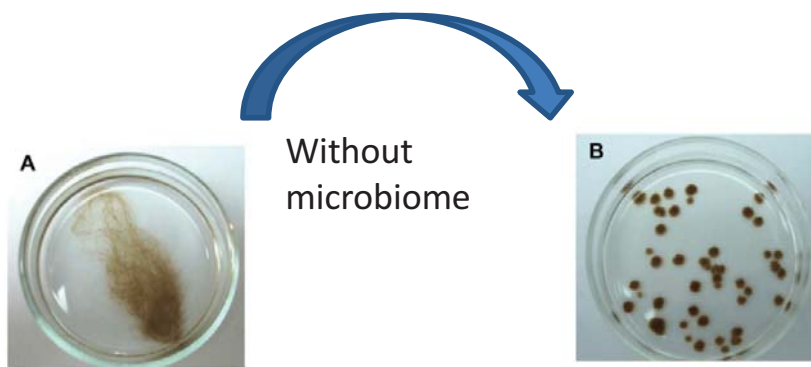
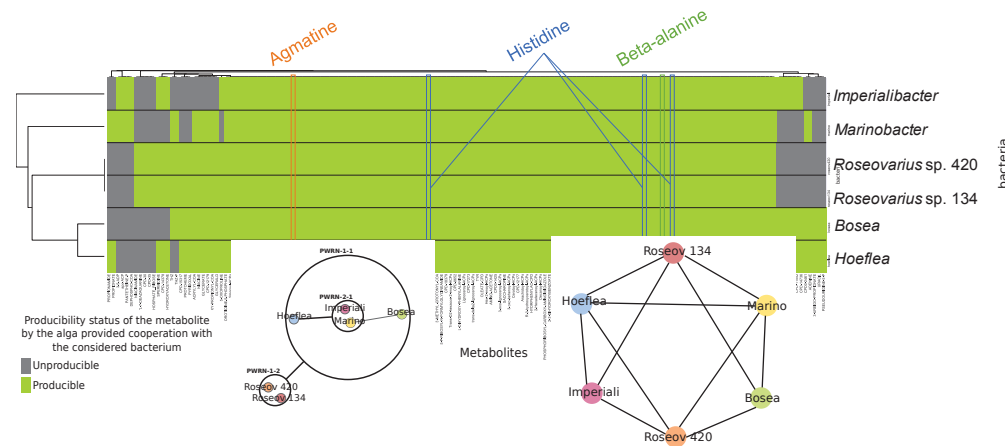


Community composition can be explained by the functional dependencies of the targets towards specific groups of bacteria

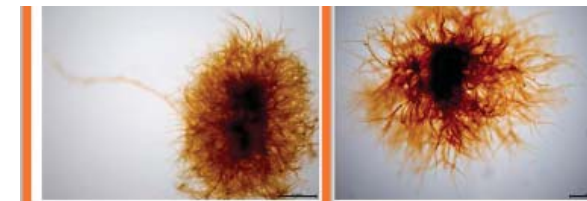
# Validation/benchmarking on human microbiome project

- *Ca. P. ectocarpi* not culturable
- 10 culturable bacteria → functional redundancy
- 6 equivalent communities of 3 bacteria

Joint work with Enora Fremy, Bertille Burgunter-Delamare & Simon Dittami



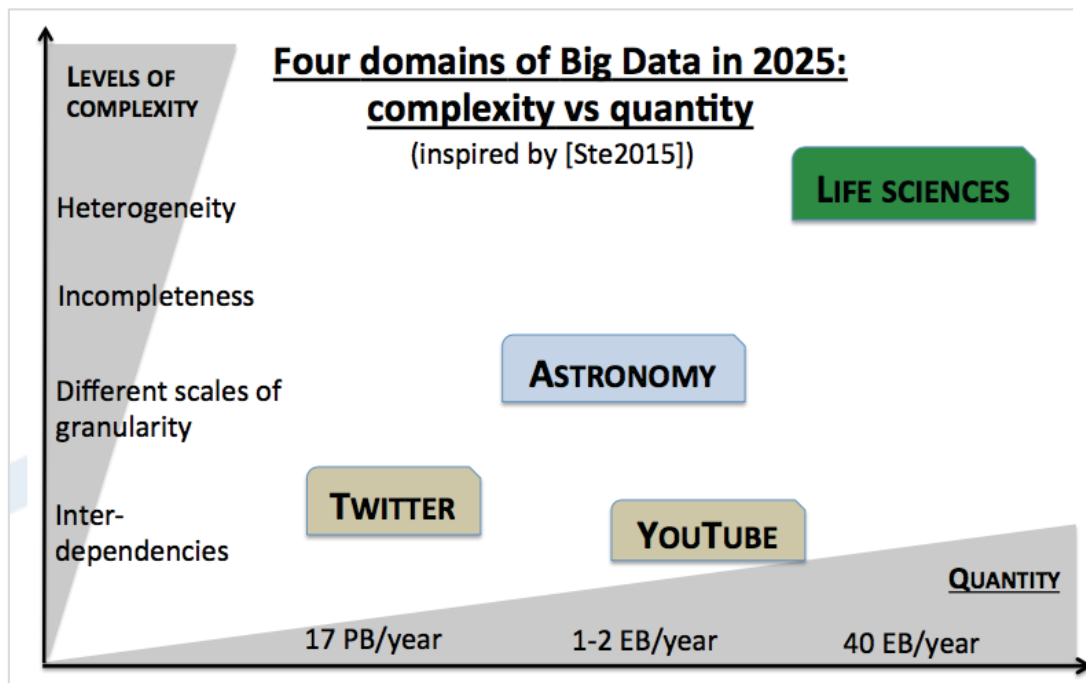
+ 3 selected  
bacteria among  
30 cultivable  
bacteria



*S. Dittami,  
Bertille Burgunter-Delamare*

The algae grew again... But with strange behaviors

## TOWARDS CONCLUSION

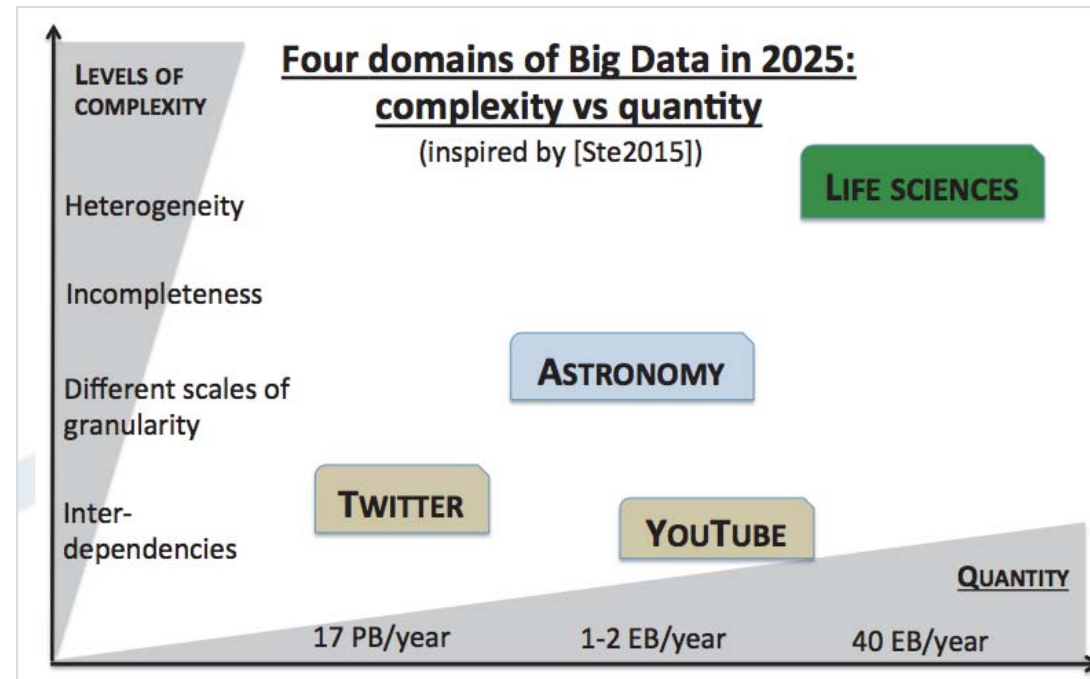
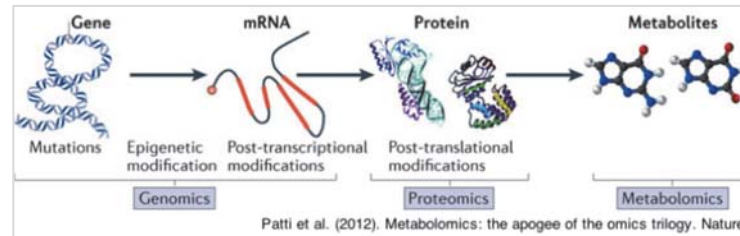


# Take home messages: life science data integration ?

- **Life science data are multi-scale and heterogeneous**
  - Linked by underlying regulatory processes
- **Systems biology ?**
  - study of complex systems which cannot be uniquely identified
- **Handling complexity for**
  - Make (dynamical) hypotheses
  - Solve optimization problems instead of identify parameters
  - Win-win collaboration with your BBF ASP-tech developers
- **We will never replace biologists**

**Molecular and cellular life science analysis is a user-assisted data science rather than a modeling system science**

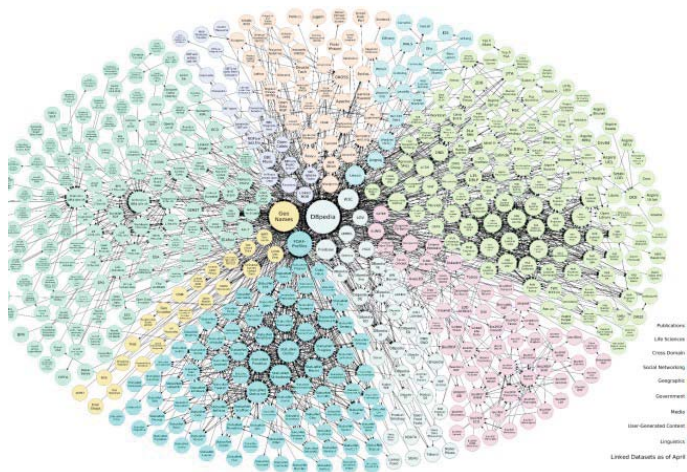
# What about the future



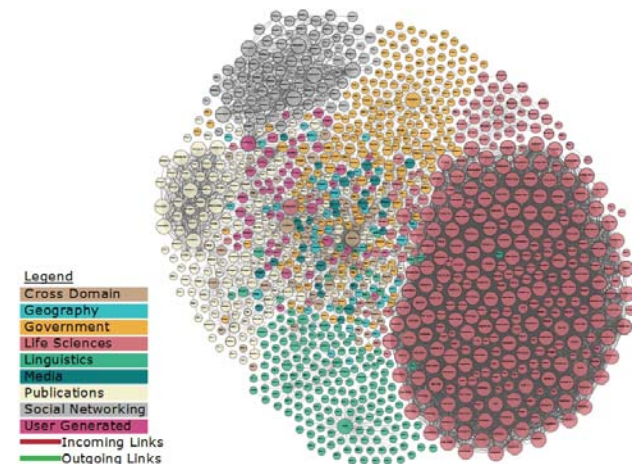
- **Size complexity**
  - Towards deep-learning ?
- **Heterogeneity complexity ?**
  - Knowledge-based methods

# Linked open data

- **More than 1500 life science databases**
  - Gene Ontology
  - Chebi
  - KEGG
  - Swissprot...
  
- **Many of these DB are being linked and can be queried**
  - Huge knowledge repositories to support reasoning



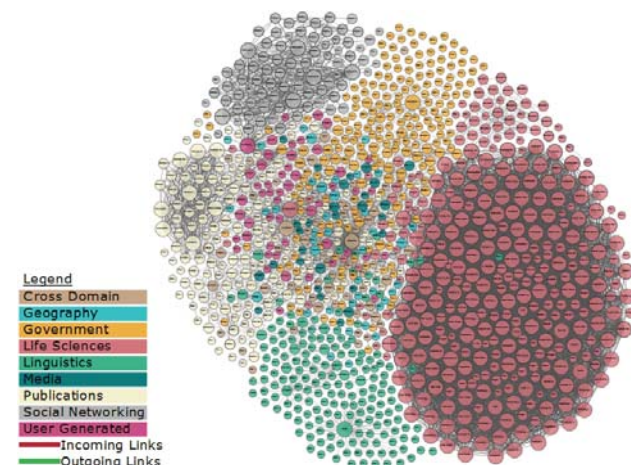
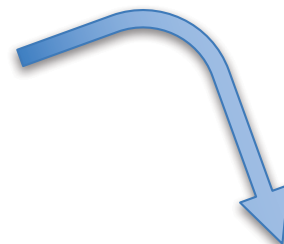
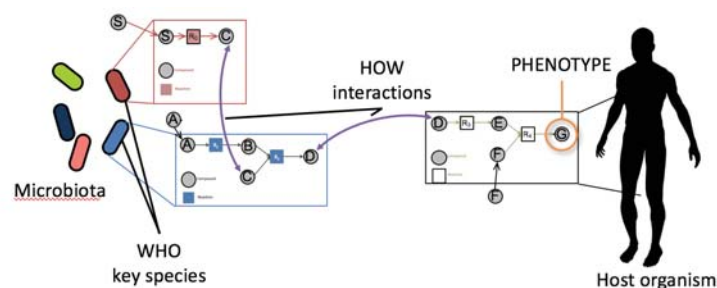
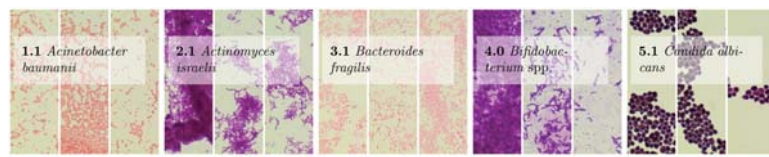
*Linked Open Data initiative (2014)*



*Linked Open Data initiative (2017)*

# The futur of life-science data analysis ?

Machine learning : compound,  
function and species identification



Knowledge representation :

Connect data

- Performant queries
- User-friendly interfaces

Formal approaches : explain

- Automatic reasoning
- Assist biologists and never replace them



# Prospective

- **Our future role : facilitate and scale life science data analysis**
  - Easy exploration of search spaces
  - **Extract dynamical features as constraints** (temporal ?)
  - Use knowledge DBs
  
- **Always explain the results**
  - Give choices to experimentalists
  - **According to all the hypotheses that we make, biologists have to double-check our predictions.**

# Acknowledgment

- **Equipe Dyliss@IRISA**

- **Métabolisme@IRISA**

- C. Frioux
- S. Prigent (INRA)
- M. Aite
- M. Chevallier
- J. Got

- **Algues@SBR Roscoff**

- S. Dittami
- H. Klijean
- B. Burgunter
- T. Tonon

- **ASP tech@Potsdam university**

- T. Schaub's team
- S. Thiele



- **CMM@Univ Chile**

- A. Maass
- M. P. Cortes
- P. Bourdon

- **Modélisation métabolisme@Nantes**

- D. Eveillard (LS2N)