

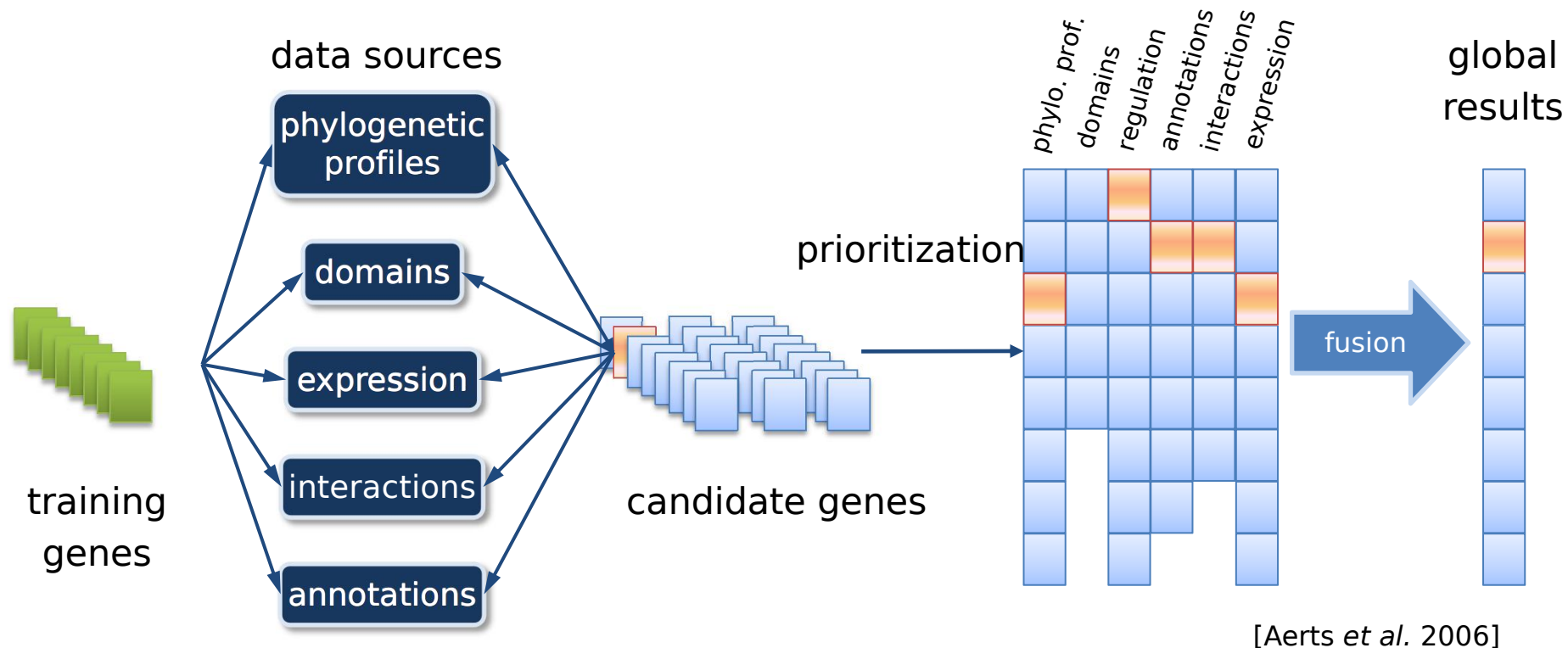
Intégration de données hétérogènes Priorisation

Master 2

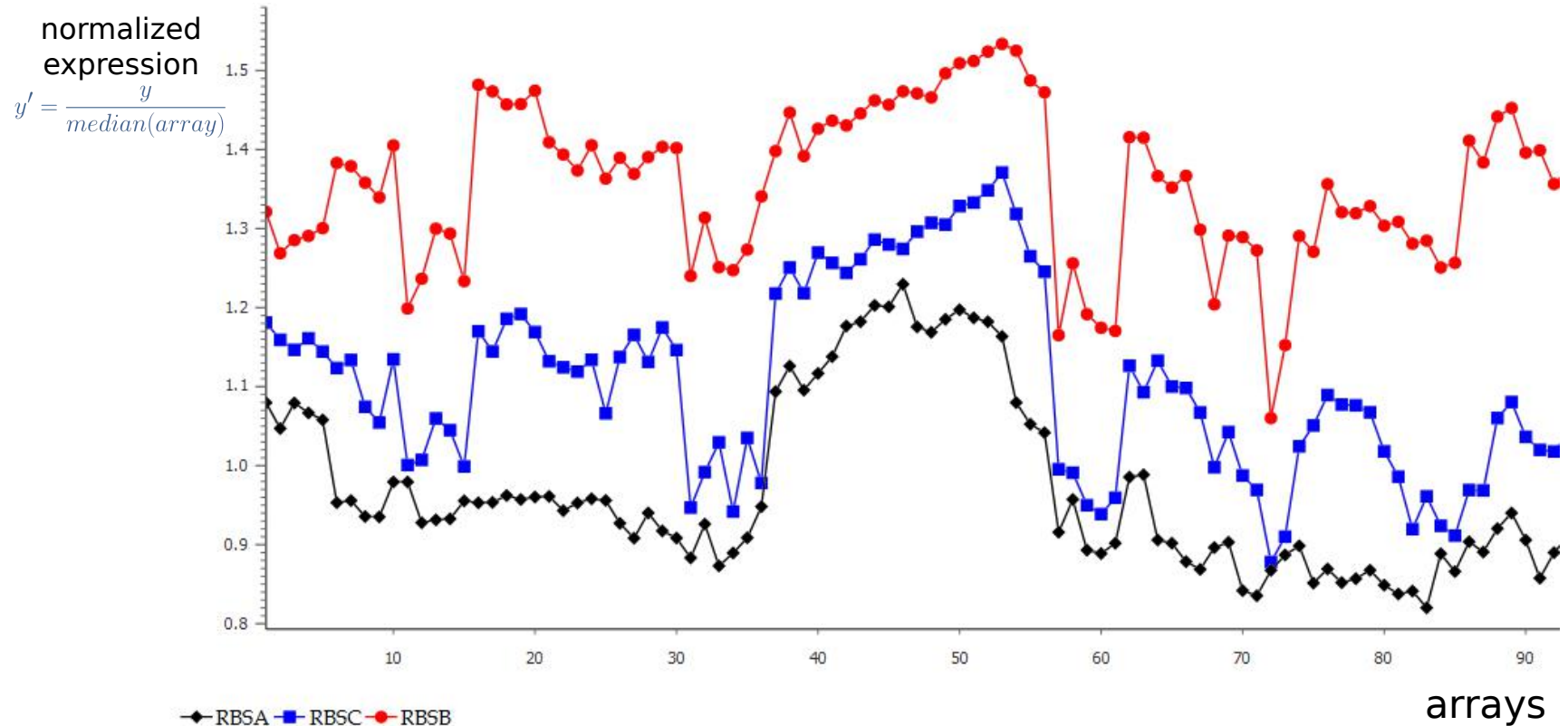
Bioinformatique et Biologie des Systèmes

Candidate gene prioritization by genomic data fusion

- Observation: more and more post-genomic data available
- Paradox: more difficult to select the best candidates
- Goal: objective and comprehensive evaluation of candidates



Gene expression



- a gene: set of expression values in various experimental conditions
- a pair of genes: dissimilarity index based on Pearson's correlation coefficient
- score : average dissimilarity

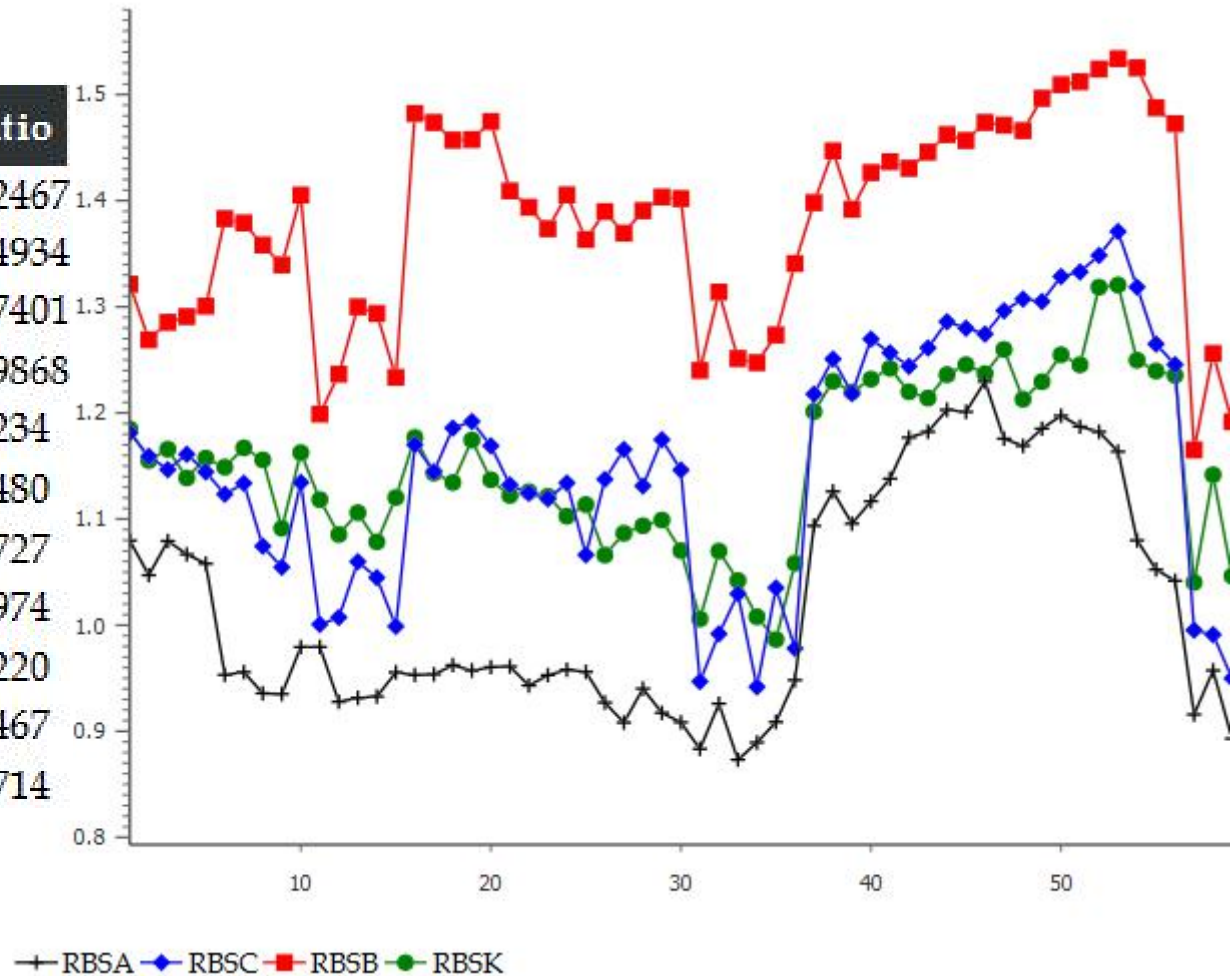


gene pairwise dissimilarity matrix

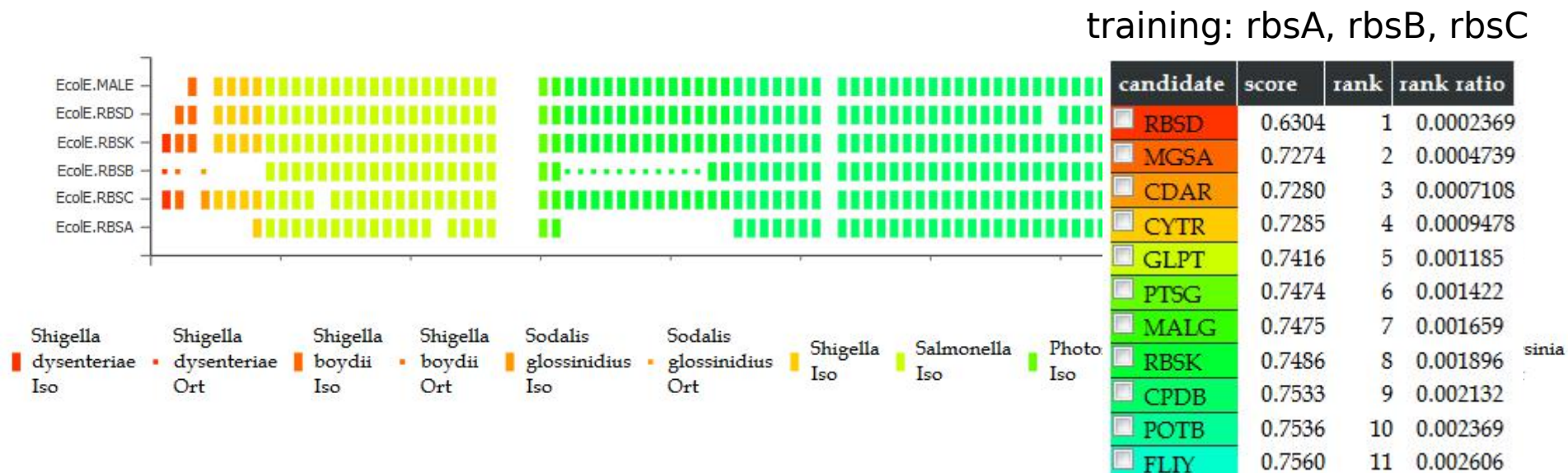
Gene expression illustration

- training: rbsA, rbsB, rbsC in *E. coli* K-12

candidate	score	rank	rank ratio
RBSK	0.1870	1	0.0002467
RBSD	0.2695	2	0.0004934
FDOI	0.3288	3	0.0007401
MALE	0.3514	4	0.0009868
MALK	0.3537	5	0.001234
FDOG	0.3551	6	0.001480
FDOH	0.3670	7	0.001727
TREB	0.3679	8	0.001974
NUPG	0.3841	9	0.002220
LAMB	0.3850	10	0.002467
MALF	0.3933	11	0.002714



Phylogenetic profiles



- a gene: presence/absence of orthologs 1:1 in other genomes
- pair of genes: dissimilarity index based on the Jaccard index
- score: average dissimilarity



gene pairwise dissimilarity matrix

Phylogenetic data

- Phylogenetic profiles



- gene pairwise distance matrix computation
 - Hypothesis: genes located near each other in a set of genomes are likely to be functionally related
 - g_1' and g_2' orthologs 1:1 of $gene_1$ and $gene_2$ in another genome i
 - Probability that the distance D_i is smaller than the observed distance d_i

$$p_i = Pr(D_i \leq d_i) = \frac{2d_i}{N_i - 1}$$

- For a set of M other genomes

$$d = Pr(D_1 \leq d_1, \dots, D_M \leq d_M) = \prod_{i=1}^M p_i$$

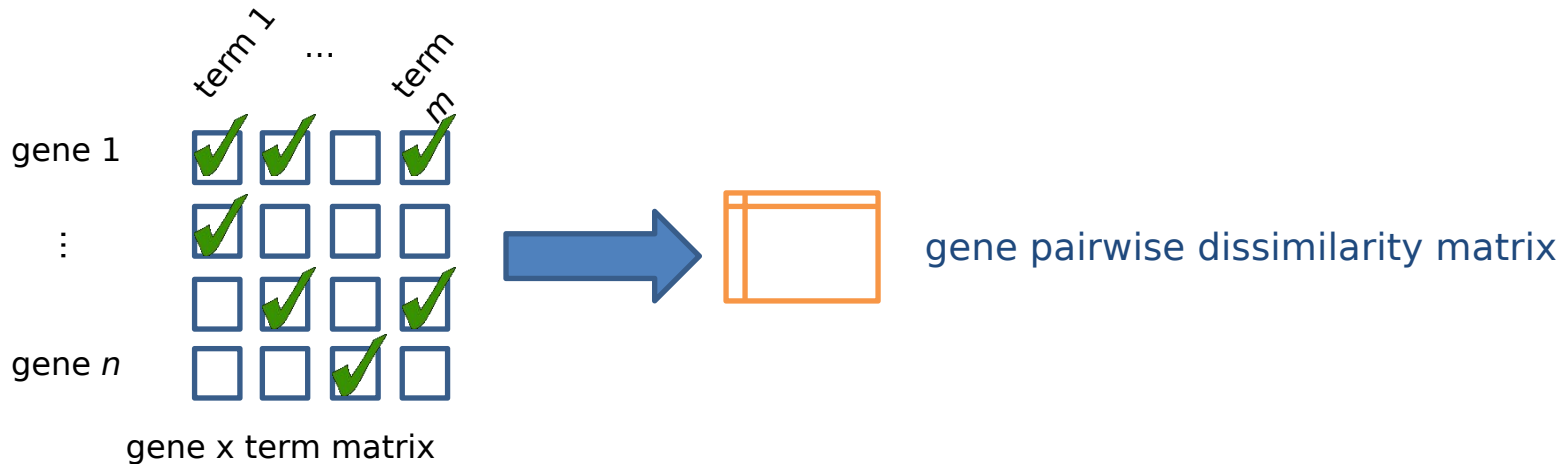
- M depends on the pair of genes considered
- d not comparable between genes (e.g. $0.1^6 = 10^{-6}$ vs. $0.5^{20} = 9.5 \cdot 10^{-7}$)
- normalization: log transformation, z-score, average of distance matrix and its transpose

Phylogenetic data: genome selection

- Reference genomes should not be too evolutionary close to the genome of interest
- Reference genomes should not be redundant in order not to introduce biases
- Need to estimate the relevance of a genome with respect to
 - A genome of interest: is it not too closely related? is it informative?
 - A set of already selected reference genomes: redundancy vs. additional signal
- Parameters
 - Rearrangements
 - Significance of genes proximity on the chromosome
 - Core genome size
 - Maximize the coverage of the genome of interest

Approaches:

- gene-term matrix: distance between rows
 - manhattan/euclidean, Jaccard, ...
- **but:** same weight for each GO-term
- based on GO-term similarity
 - adapt weight to information content



- Node/term information content

$$IC(term) = -\log p(term) \quad \text{with } p(term) = \text{freq}(term)$$

- $MICA(t_1, t_2)$: Maximum Information Common Ancestor

$$MICA(t_1, t_2) = \arg \max IC(t_i), t_i \in \text{ancestors}(t_1, t_2)$$

- $sim_{res}(t_1, t_2) = IC(MICA(t_1, t_2))$ [Resnik, 1995]

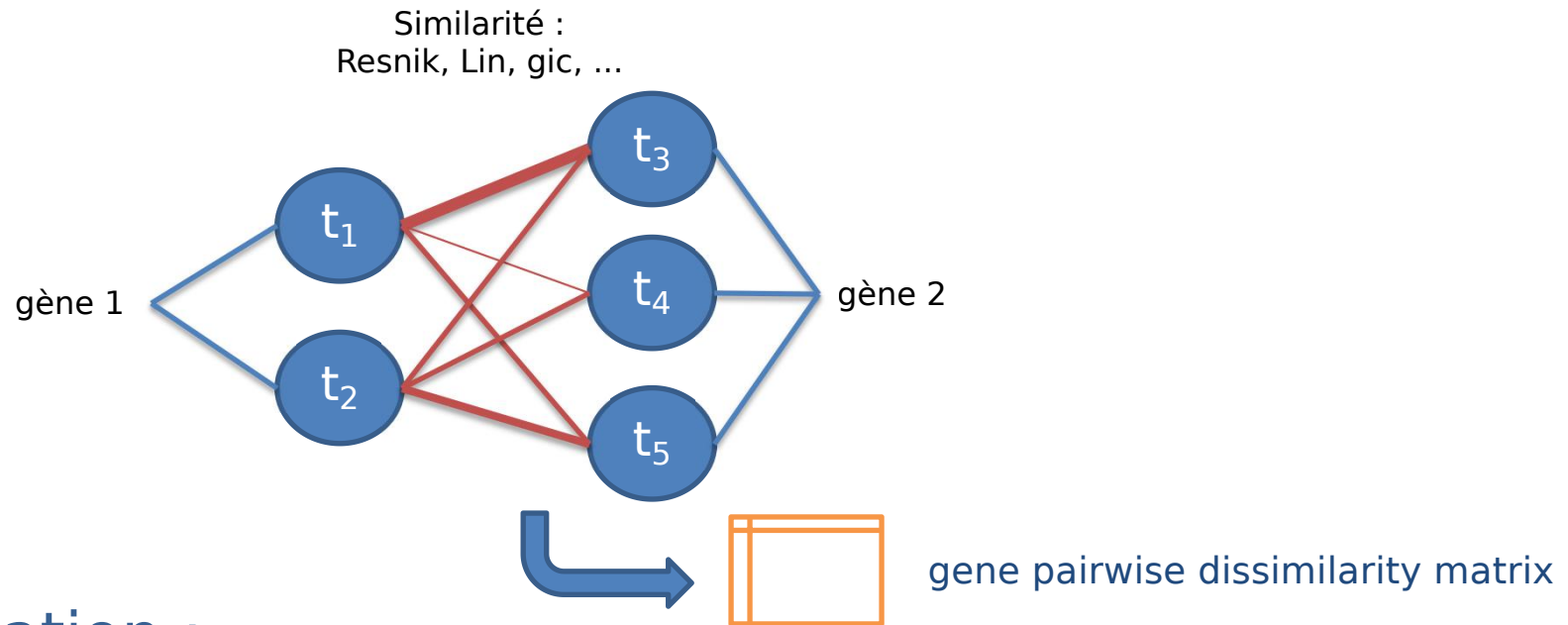
- $sim_{lin}(t_1, t_2) = IC(MICA(t_1, t_2)) / (IC(t_1) + IC(t_2))$ [Lin, 1998]

- $sim_{gic}(t_1, t_2) = \frac{\sum_{t \in \{GO(t_1) \cap GO(t_2)\}} IC(t)}{\sum_{t \in \{GO(t_1) \cup GO(t_2)\}} IC(t)}$ [Pesquita et al., 2008]

Simimilarité entre gènes basée sur la similarité entre termes GO

Possibilités :

- Similarité moyenne des termes communs au 2 gènes
- Similarité maximale, ex : t_1-t_3
- Best Match Average (bma), ex : $\text{ave}(t_1-t_3, t_2-t_5)$

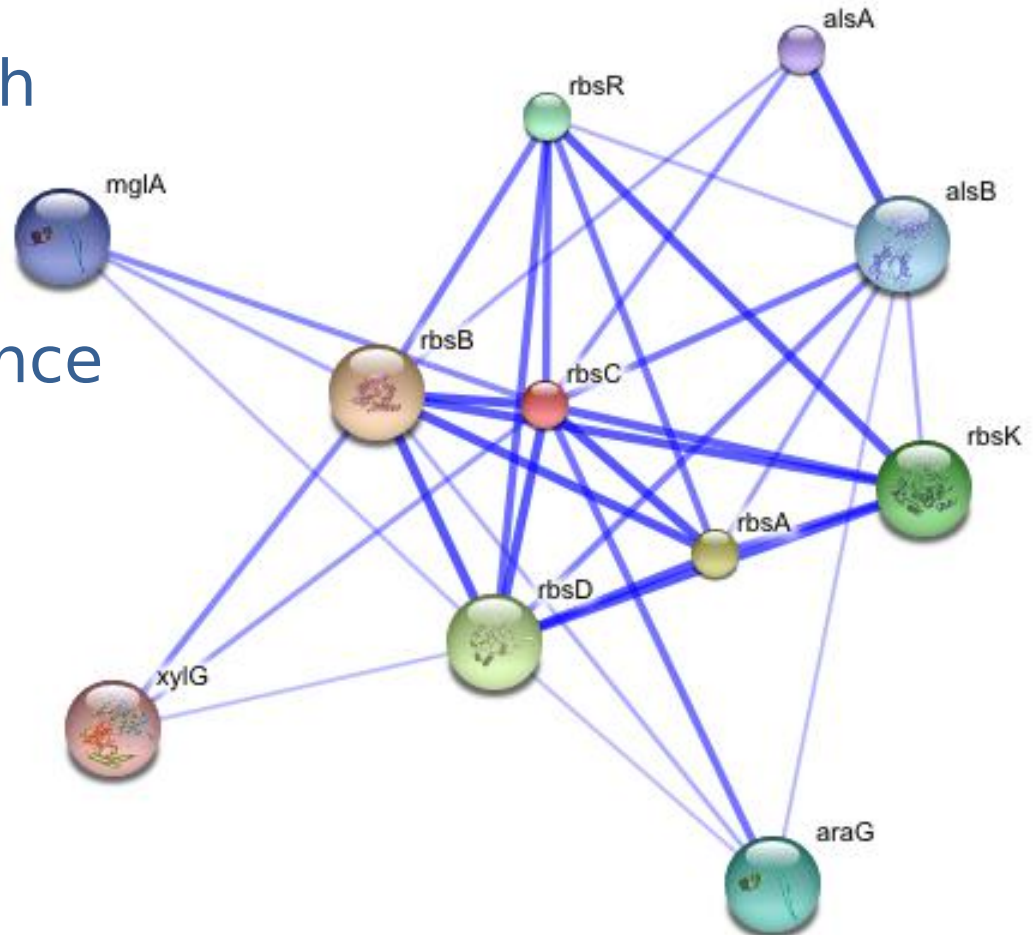


Application :

- Performances légèrement meilleures obtenues que les autres avec la combinaison Resnik + similarité maximale
- à confirmer sur d'autres jeux de données ou d'autres contextes

Interactions

- all pairs shortest path
- a pair of gene:
shortest path length
- score: average distance



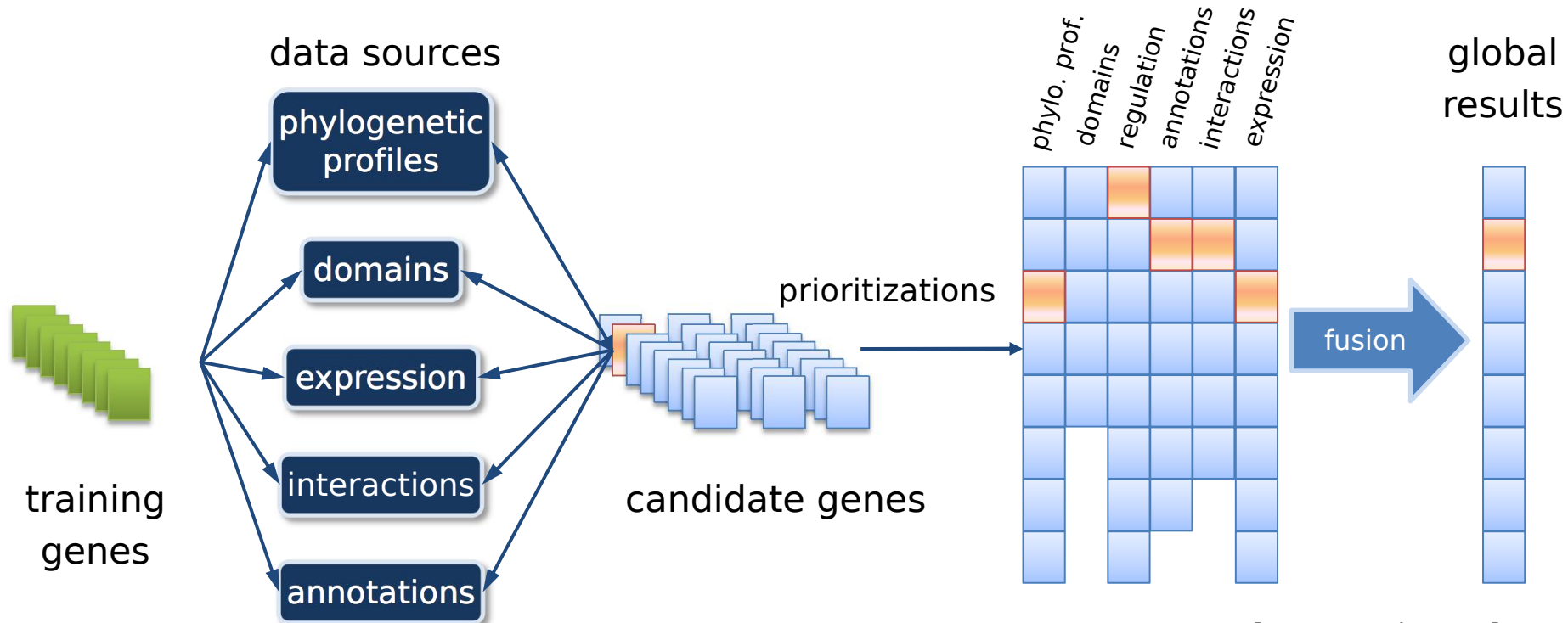
training: rbsA, rbsB, rbsC

candidate	score	rank	rank ratio
<input type="checkbox"/> RBSK	1.000	2	0.0005136
<input type="checkbox"/> RBSD	1.000	2	0.0005136
<input type="checkbox"/> RBSR	1.000	2	0.0005136
<input type="checkbox"/> ALSB	1.333	5	0.001284
<input type="checkbox"/> ALSC	1.333	5	0.001284
<input type="checkbox"/> YPHD	1.333	5	0.001284
<input type="checkbox"/> MGLC	1.667	10.5	0.002696
<input type="checkbox"/> XYLG	1.667	10.5	0.002696
<input type="checkbox"/> ALSA	1.667	10.5	0.002696
<input type="checkbox"/> YTFT	1.667	10.5	0.002696

from STRING
<http://string-db.org>

Candidate gene prioritization by genomic data fusion

- Observation: more and more post-genomic data available
- Paradox: more difficult to select the best candidates
- Goal: objective and comprehensive evaluation of candidates

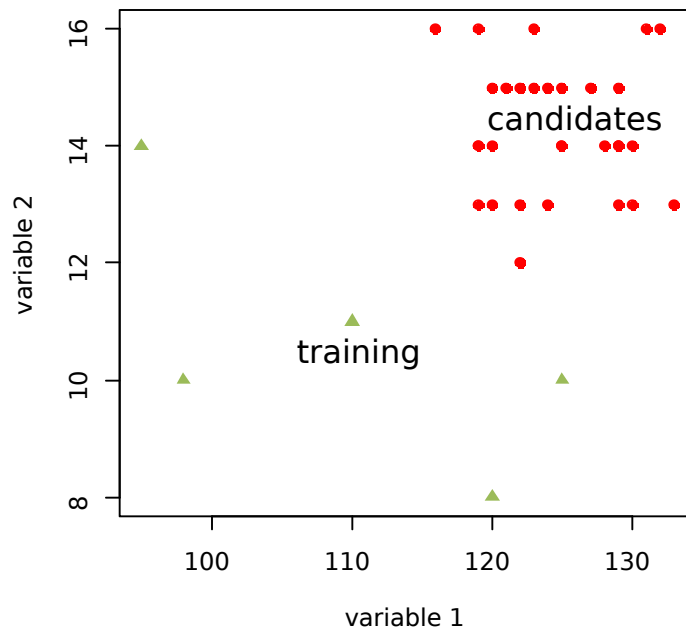


[Aerts et al. 2006]

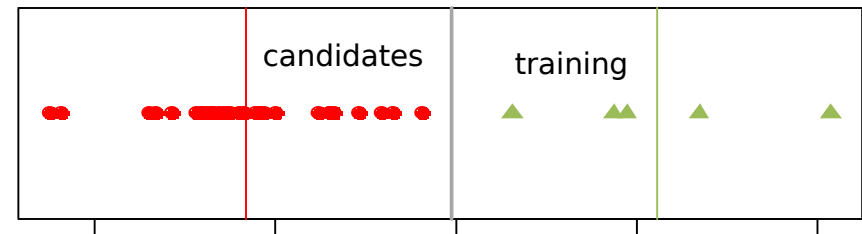
Weighted fusion through linear discriminant analysis

• Principles

- prioritize the candidate genes and including the training genes
- consider each data source as a measure for classification with classes: training/candidate
- perform discriminant analysis to weigh and separate training genes from background (candidates)

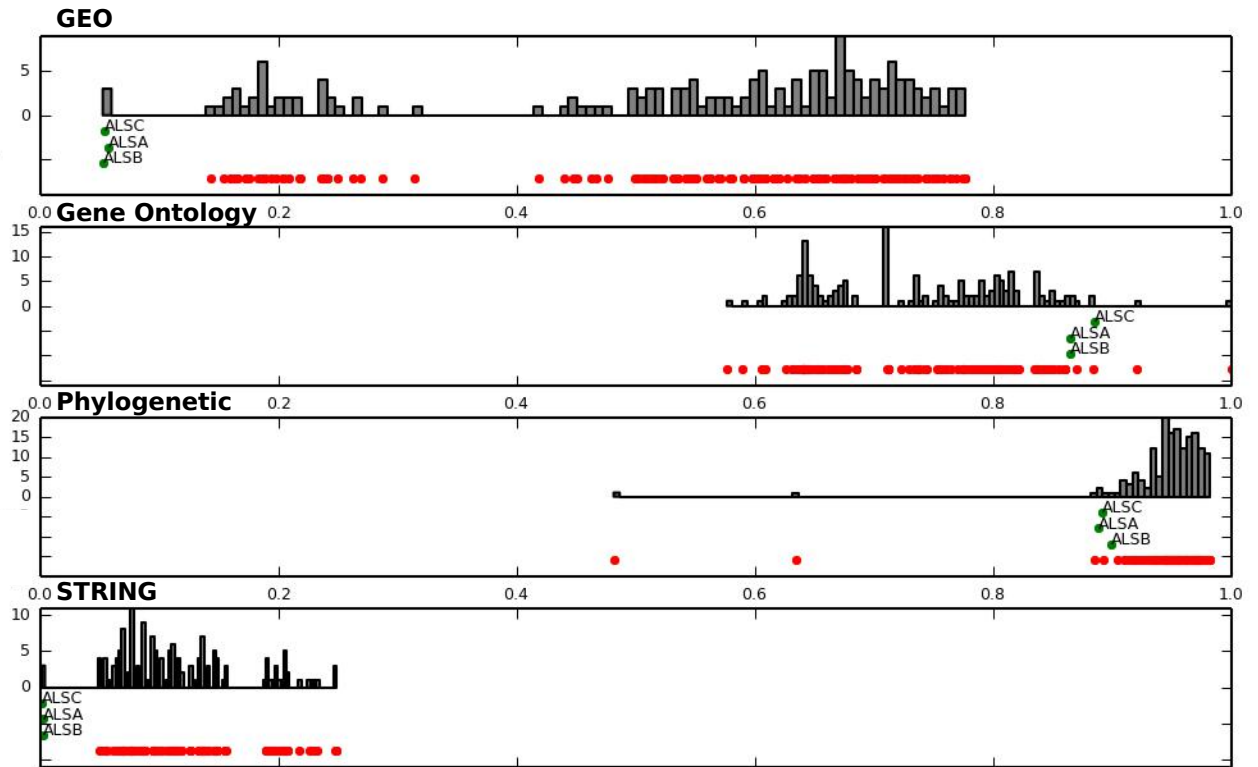


dimension	weight
variable 1	-0.1307346
variable 2	-0.7031850



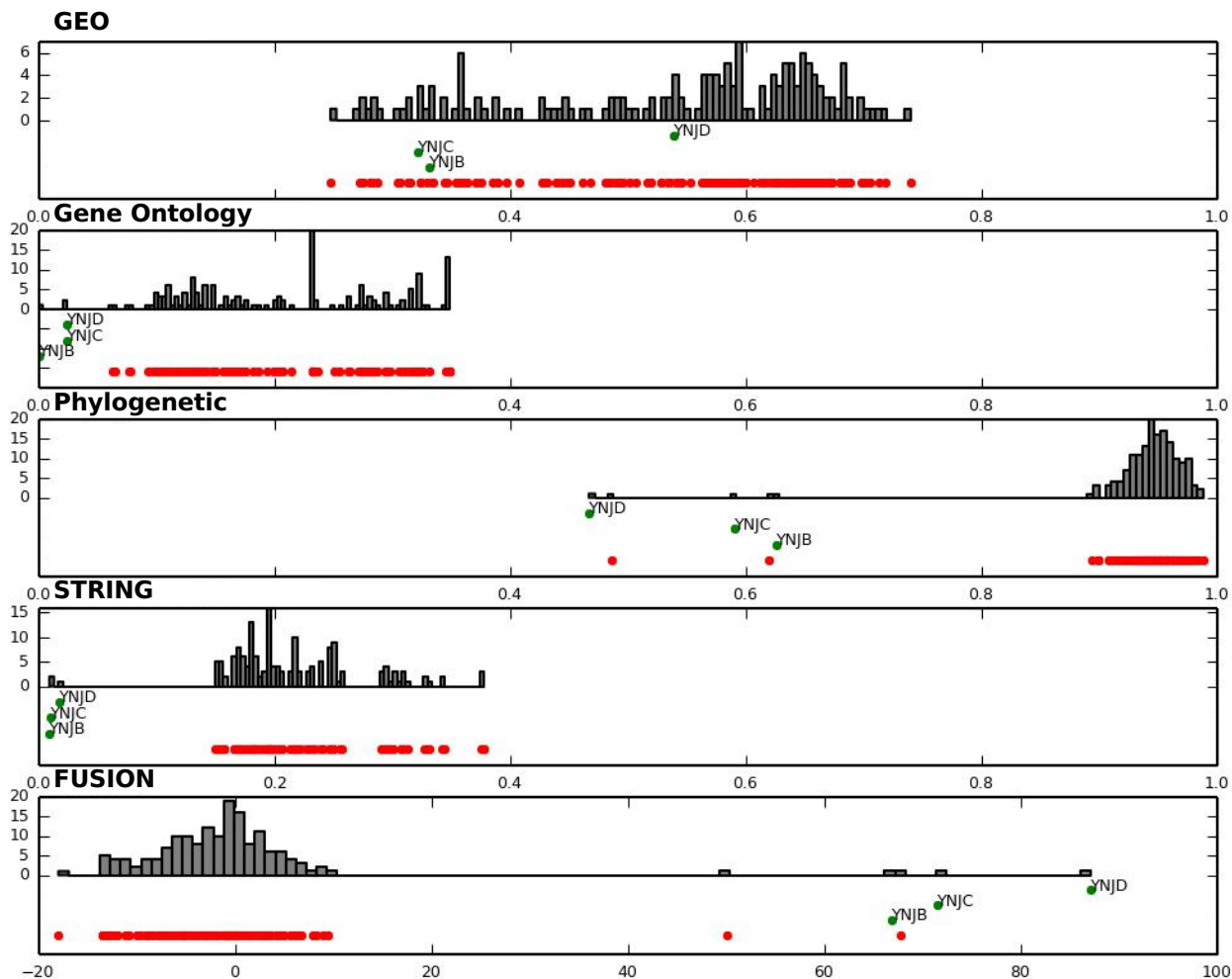
Application to *E. coli* *alsA* system: *alsA*, *alsB*, *alsC*

Data source	Weight
Expression (GEO)	3.5
Annotations (Gene Ontology)	-4.7
Phylogenetic	4.0
Interactions (STRING)	12.3



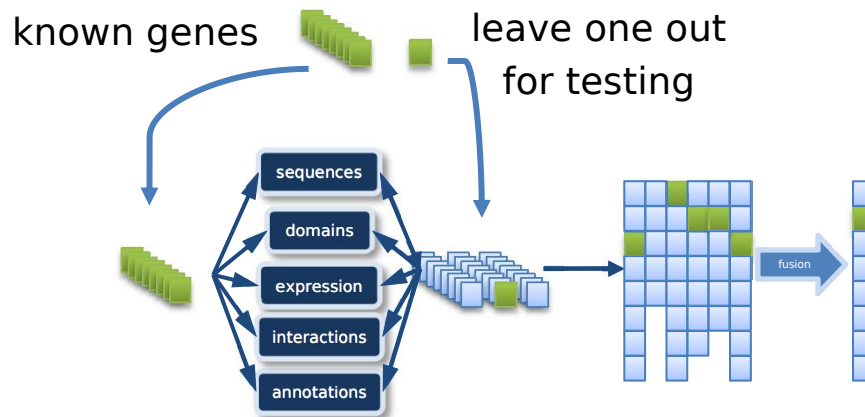
Application to *E. coli* *ynjD* system: *ynjB*, *ynjC*, *ynjD*

Data source	Weight
Expression (geo)	1.3
Annotations (Gene Ontology)	3.4
Phylogenetic	17.4
Interactions (string)	6.6



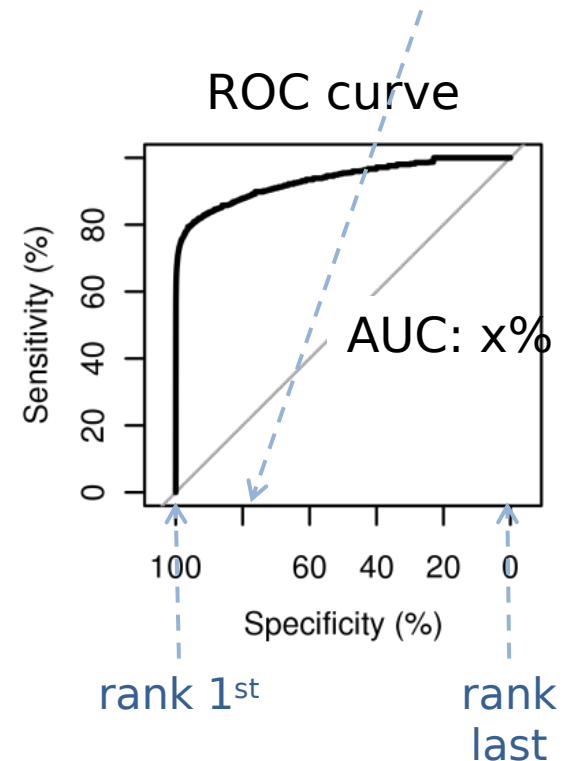
Evaluation methodology

- Leave-one-out cross validation (LOOCV)



How well does it rank?
e.g. rank ratio = $2/8 = 0.25$

- for each manually curated ABC system
 - ♦ perform LOOCV on each gene: rank ratio
 - ♦ plot Receiver Operating Characteristic (ROC) curve and consider Area Under the Curve (AUC)



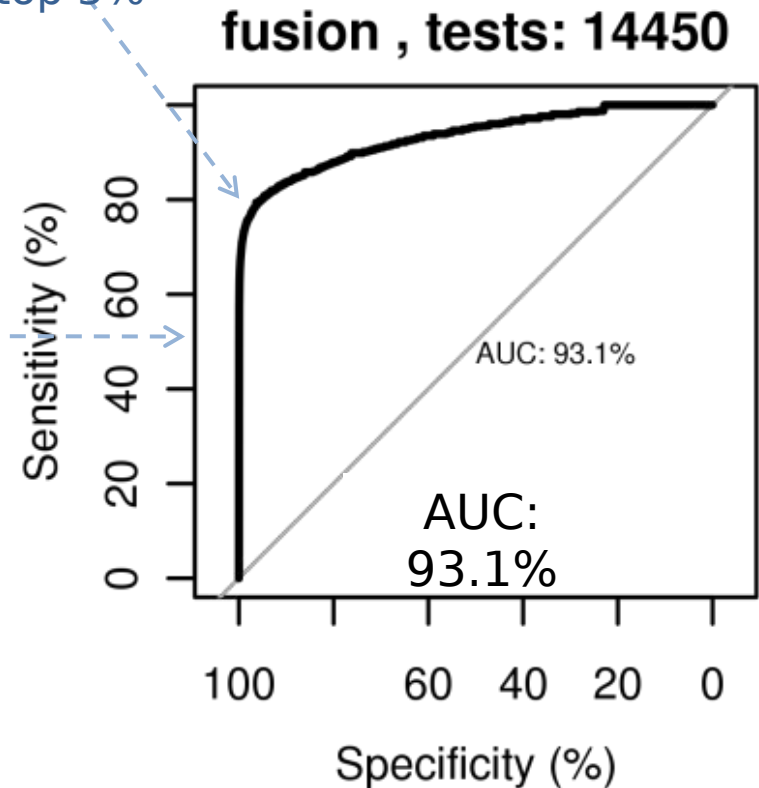
Evaluation: global results

Gold standard

- ABCdb, manually curated ABC systems:
 - ♦ 135 genomes
 - ♦ 14,450 genes
 - ♦ 4,586 ABC systems

53% of the left out genes rank 1st

80% of the left out genes rank in the top 5%



Prioritization for functional inference

Organism	Hide	Hide
Organism	Escherichia coli (strain K12)	
External Links	[UNIPROT] [NCBI]	
Taxonomic Lineage	> Bacteria > Proteobacteria > Gammaproteobacteria > Enterobacteriales > Enterobacteriaceae > Escherichia > Escherichia coli > EcolE	
Strain Name	K12	
ABCdb identifier	EcolE	
Chromosomes	EcolE01	

Assembly	Hide	Hide		
Assembly	NBD	MSD	SBP	Class
EcolE01.RBSB	★ EcolE01.RBSA	★ EcolE01.RBSC	★ EcolE01.RBSB	A_1a

Proteins	Hide	Hide	
Protein	Domain	Subfamily	TCdb
★ EcolE01.RBSB	SBP	S_1aa	3.A.1.2.1 Ribose porter (RbsC has 10 TMSs with N- and C-termini in the cytoplasm (Stewart and Hermodson, 2003))
★ EcolE01.RBSC	MSD	M_1aa	3.A.1.2.1 Ribose porter (RbsC has 10 TMSs with N- and C-termini in the cytoplasm (Stewart and Hermodson, 2003))
★ EcolE01.RBSA	NBD-NBD	N_1aN&N_1aC	3.A.1.2.1 Ribose porter (RbsC has 10 TMSs with N- and C-termini in the cytoplasm (Stewart and Hermodson, 2003))

from ABCdb
<http://www-abcdb.biotoul.fr>

Prioritization for functional inference

Prioritization **Hide**

Hide

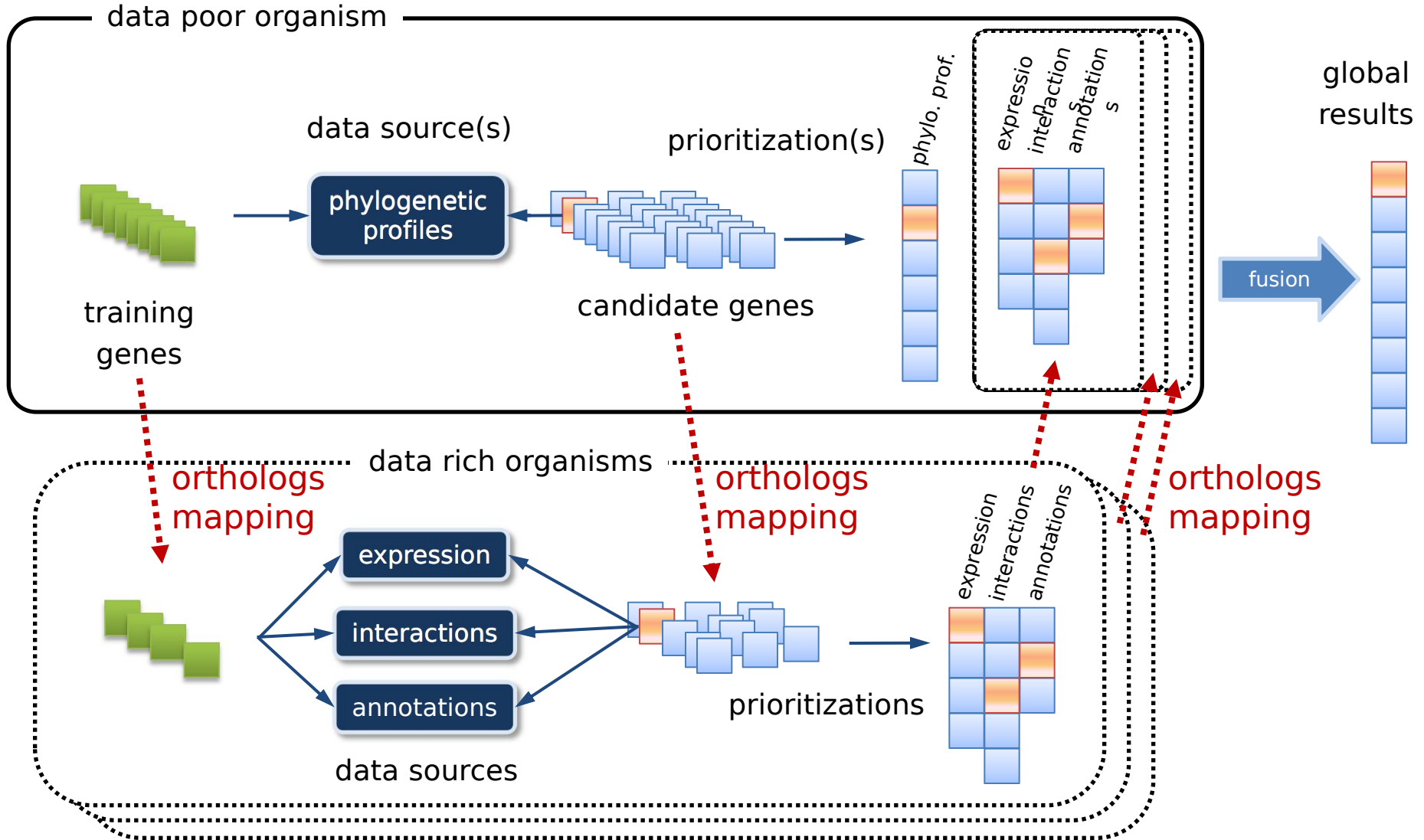
Run prioritization.

Show entries

Search:

rank	Global results	pathways (fusion)	string (fusion)	transcriptome (fusion)	phylogenetic_profiles EcolE	go (fusion)	interactome EcolE
1	RBSD (1) S: 0, RR: 0	D-ribose pyranase					
2	RBSK (2) S: 0, RR: 0	Ribokinase					
3	MALE (3) S: 0, RR: 0.001	SBP of maltose/maltodextrin/maltoologisaccharide ABC transporter					
4	DEOC (4) S: 0, RR: 0.001	Deoxyribose-phosphate aldolase					
5	RBSR (5) S: 0.001, RR: 0.001	Ribose operon repressor					
6	UDP (6) S: 0.001, RR: 0.001	Uridine phosphorylase					
7	MGLA (7) S: 0.001, RR: 0.002	NBD of galactose/glucose (methyl galactoside) ABC transporter (same subfamily)					
8	MUKF (8) S: 0.002, RR: 0.002	Chromosome partition protein mukF					
9	GAPA (9) S: 0.002, RR: 0.002	<i>CITT (2056)</i> S: 1, RR: 1	XYLF (9) S: 0, RR: 0.002	UCPA (9) S: 0.003, RR: 0.002	CPDB (9) S: 0.753, RR: 0.002	RPLN (16) S: 0.004, RR: 0.004	UDP (34) S: 1.5, RR: 0.009

Extension of the method to data poor organisms

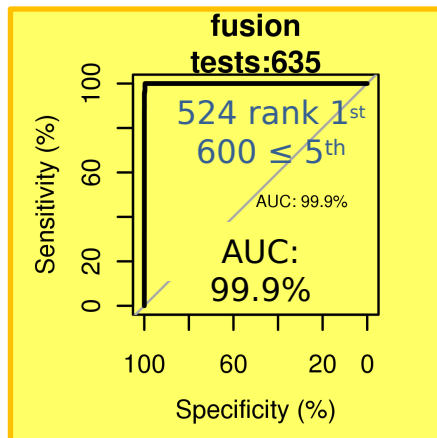
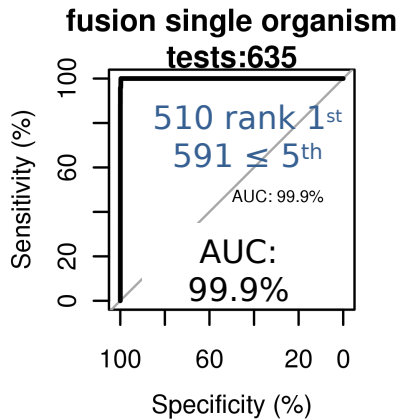


Performances: using other organisms data through orthology

Organisms:

B. subtilis, *E. coli*, *P. aeruginosa*

192 ABC systems, 635 genes



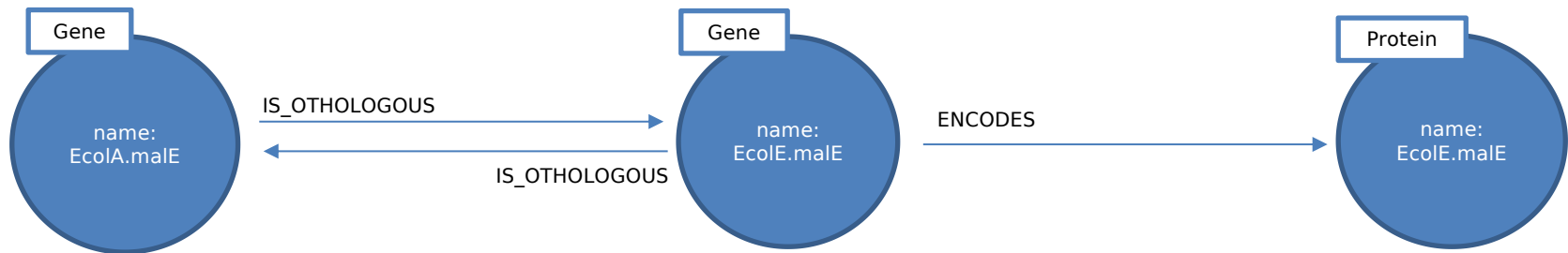
		# ABC genes	AUC (%)	data sources	
Bacteria	Actinobacteria				
	Thermotogales	Streptomyces coelicolor	434	93	
	Chlamydiales	Thermotoga maritima	170	95.3	
	Mycoplasmatales	Chlamydia trachomatis	31	89.8	
	Mycoplasmatales	Mycoplasma gallisepticum	47	82.4	
	Mycoplasmatales	Mycoplasma genitalium	36	74.8	
	Epsilonproteobacteria	Helicobacter pylori	42	90.7	
	Betaproteobacteria	Nitrosomonas europaea	76	98	
	Proteobacteria	Gammaproteobacteria	Pseudomonas aeruginosa	285	99.9 ★★★★★
			Coxiella burnetii	37	97.8
			Escherichia coli	216	99.9 ★★★★★
			Salmonella enterica	198	98.4 ★★
			Shigella flexneri	192	99.3 ★
		Enterobacteriaceae			
		Alphaproteobacteria	Bradyrhizobium japonicum	619	99.9 ★
		Alphaproteobacteria	Anaplasma marginale	18	96.2
		Clostridiales	Clostridium perfringens	141	92.7
		Streptococcaceae	Streptococcus pneumoniae	163	96.7
	Firmicutes				
	Bacillales	Staphylococcus aureus	145	98.3 ★★	
Bacillales	Bacillus subtilis	208	99.9 ★★★★★		
Cyanobacteria					
Cyanobacteria	Nostoc sp.	234	96 ★		
Cyanobacteria	Synechocystis sp.	132	96		
Archaea	Crenarchaeota				
	Crenarchaeota	Thermophilum pendens	141	89.9	
	Crenarchaeota	Metallosphaera sedula	60	81.2	
	Thaumarchaeota	Aeropyrum pernix	104	90.4	
	Thaumarchaeota	Nitrosopumilus maritimus	33	90.2	
	Euryarchaeota	Euryarchaeota	Methanocaldococcus jannaschii	40	95
		Euryarchaeota	Thermococcus onnurineus	65	90.6
		Euryarchaeota	Halobacterium sp.	83	93 ★★
		Euryarchaeota	Methanosphaera stadtmanae	35	96.5
		Euryarchaeota	Candidatus Methanoregula boonei	83	93.4
Euryarchaeota		Methanosarcina mazei	119	91.9	
Methanomicrobia					

Performances: using other organisms data through orthology

			# ABC genes	AUC (%)	data sources	
Bacteria	Actinobacteria					
		Thermotogales	Streptomyces coelicolor	434	93	
			Thermotoga maritima	170	95.3	
		Chlamydiales	Chlamydia trachomatis	31	89.8	
		Mycoplasma	Mycoplasma gallisepticum	47	82.4	
			Mycoplasma genitalium	36	74.8	
		Epsilonproteobacteria	Helicobacter pylori	42	90.7	
		Betaproteobacteria	Nitrosomonas europaea	76	98	
			Pseudomonas aeruginosa	285	99.9	★ ★ ★ ★
			Coxiella burnetii	37	97.8	
		Gamma proteobacteria	Escherichia coli	216	99.9	★ ★ ★ ★
			Salmonella enterica	198	98.4	★ ★
			Shigella flexneri	192	99.3	★
		Alphaproteobacteria	Bradyrhizobium japonicum	619	99.9	★
			Anaplasma marginale	18	96.2	
		Clostridiales	Clostridium perfringens	141	92.7	
		Streptococcaceae	Streptococcus pneumoniae	163	96.7	
		Bacillales	Staphylococcus aureus	145	98.3	★ ★
		Bacillus subtilis	208	99.9	★ ★ ★ ★	
	Cyanobacteria	Nostoc sp.	234	96	★	
		Synechocystis sp.	132	96		
		Thermophilum pendens	141	89.9		

• Principes

- Représenter les données en tant qu'objets reliés par des relations
- Chaque objet ou relation peut avoir des attributs qui lui sont propres
- Développement d'un langage de manipulation et de requête

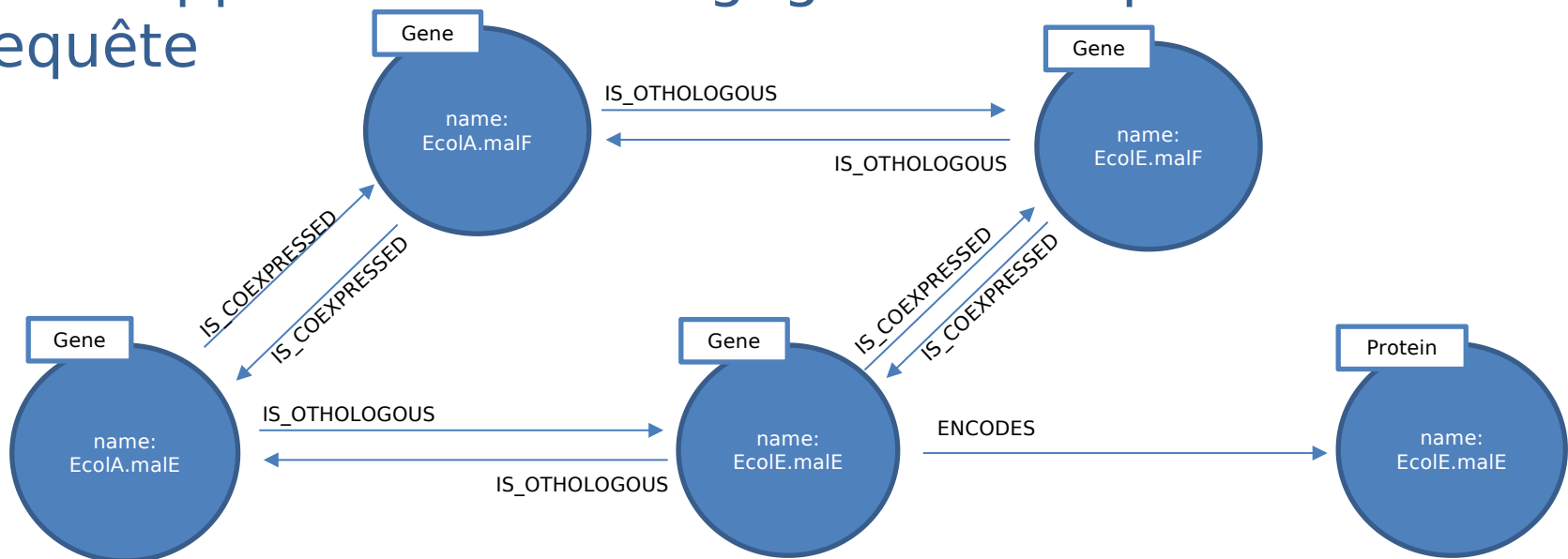


Labeled Property Graph

$(EcolA.malE:Gene) \leftarrow [:IS_ORTHOLOGOUS] \rightarrow (EcolE.malE:Gene) - [:ENCODES] \rightarrow (EcolE.malE:Protein)$

• Principes

- Représenter les données en tant qu'objets reliés par des relations
- Chaque objet ou relation peut avoir des attributs qui lui sont propres
- Développement d'un langage de manipulation et de requête



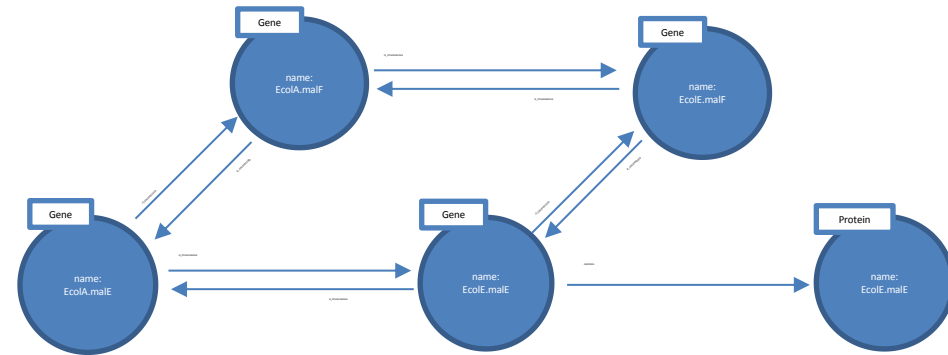
Labeled Property Graph

```

(EcoI.E.maIE:Gene) <- [:IS_COEXPRESSED] -> (EcoI.E.maIF) <- [:IS_ORTHOLOGOUS] ->
(EcoI.A.maIF:Gene) <- [:IS_COEXPRESSED] -> (EcoI.A.maIE:Gene) <- [:IS_ORTHOLOGOUS] -> (EcoI.E.maIE:Gene)
- [:ENCODES] -> (EcoI.E.maIE:Protein)
  
```


Labeled property graph

Un graphe avec propriétés étiquetées est constitué de sommets, relations, propriétés et étiquettes :



- Propriétés des sommets : de type clé/valeur
- Étiquettes des sommets : une ou plusieurs afin de les regrouper (Gene, Protein)
- Relations : orientées, peuvent avoir des propriétés comme les sommets.

Langage de requête, exemple : Cypher

```
MATCH (g:Gene) -[:ENCODES]->(p:Protein)
WHERE g.name='EcoI.E.malE'
RETURN g,p
```

