

Intégration de données hétérogènes

Bases de données non structurées

Master 2

Bioinformatique et Biologie des Systèmes

Programme de la journée “marathon sprinté”

- Matin
 - Aperçu sur
 - Big data, NoSQL, base de données de graphes, neo4j
 - Prise en main de neo4j
 - Construction d'une base de connaissances (graphe) sur les gènes/proéines d'*Escherichia coli*
 - données d'expression
 - données d'interaction protéine-protéine
 - données phylogénomique
 - annotations Gene Ontology
- Après midi
 - Concepts
 - caractérisation d'un ensemble
 - Application
 - Annotations GO sur-représentées parmi les ensembles de gènes/protéines présentant à la fois
 - coexpression,
 - interaction protéine-protéine,
 - conservation du voisinage génomique

Tendances et observations

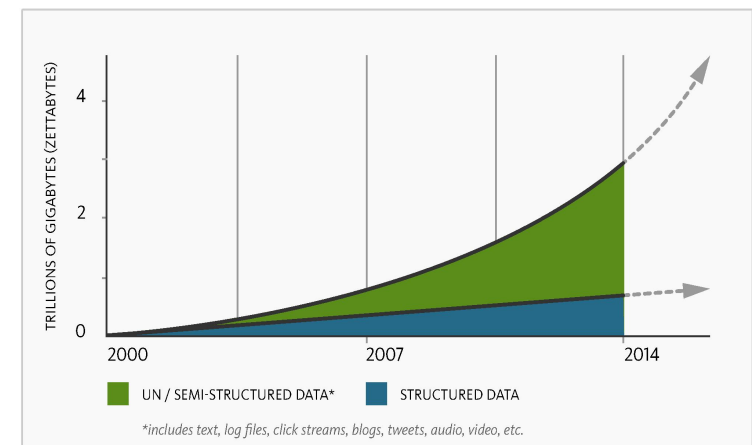
Big data

- “Big data is high **volume**, high **velocity**, and/or high **variety** information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization.” (Gartner, 2012)

Données non structurées

- Pareja-Tobes *et al.* 2015 : "Sometimes analysis can be scheduled (annotating a genome for a research project), while in other cases it is inherently impossible to even predict when you will need to run them (think of analyzing data coming from an outbreak). The resources needed or their characteristics vary wildly depending on the specific analysis. We are thus in a context where in general **we do not know in advance what we need to access**, how fast, when, and how."

NoSQL Databases



NoSQL databases

What is “NoSQL”?

- term used in late 90s for a different type of technology:
Carlo Strozzi: http://www.strozzi.it/cgi-bin/CSA/tw7//en_US/NoSQL/
- “Not Only SQL”? (but many RDBMS are also “not just SQL”)
- “NoSQL is an accidental term with no precise definition”. first used at an informal meetup in 2009 in San Francisco (presentations from Voldemort, Cassandra, Dynomite, HBase, Hypertable, CouchDB, and MongoDB)

NoSQL: Database technologies that are (mostly)

- Not using the relational model (nor the SQL language)
- Designed to run on large clusters (horizontally scalable)
- No schema - fields can be freely added to any record
- Open source
 - Based on the needs of 21st century web estates

Other characteristics (often true):

- easy replication support (fault-tolerance, query efficiency)
- simple API
- eventually consistent (not ACID)

Flexible scalability

- horizontal scalability instead of vertical

Dynamic schema of data

- different levels of flexibility for different types of DB

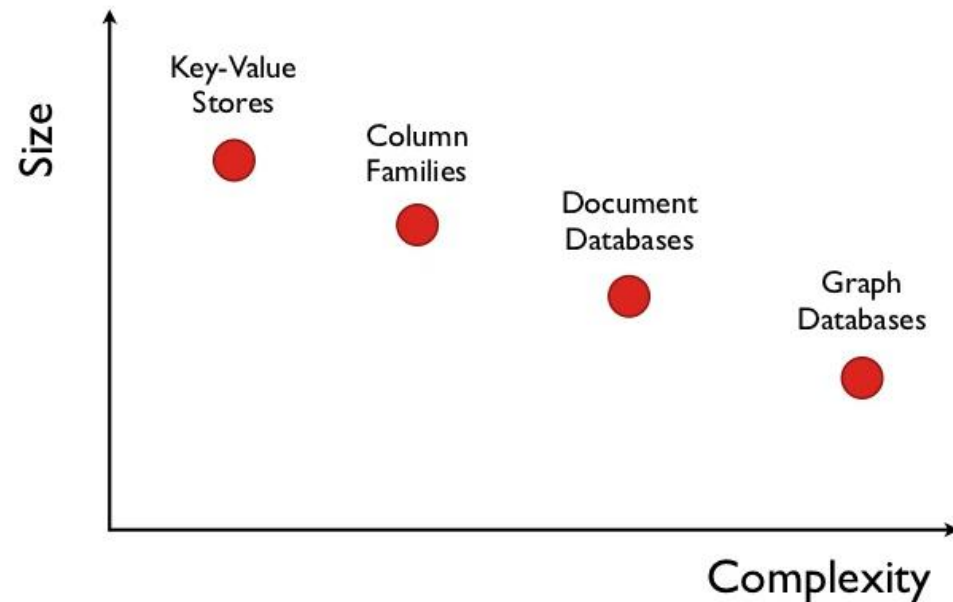
Efficient reading

- spend more time storing the data, but read fast
- keep relevant information together

Cost saving

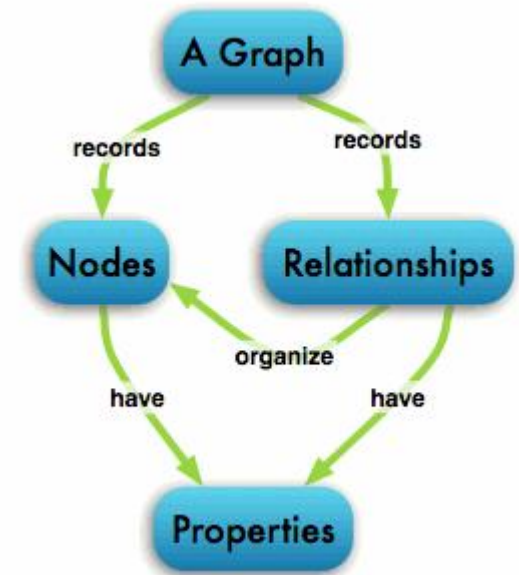
- designed to run on commodity hardware
- typically open-source (with a support from a company)

- **MapReduce** programming model
 - running over a distributed file system
- **Key-value stores**
- **Document databases**
- **Column-family stores**
- **Graph databases**



source: <http://www.slideshare.net/emileifrem/nosql-east-a-nosql-overview-and-the-benefits-of-graph-databases>

- **Open source** graph database
 - The most **popular**
- Initial release: 2007
- Written in: **Java**
- OS: cross-platform
- Stores data as **nodes** connected by directed, typed **relationships**
 - With properties on both
 - Called the “property graph”

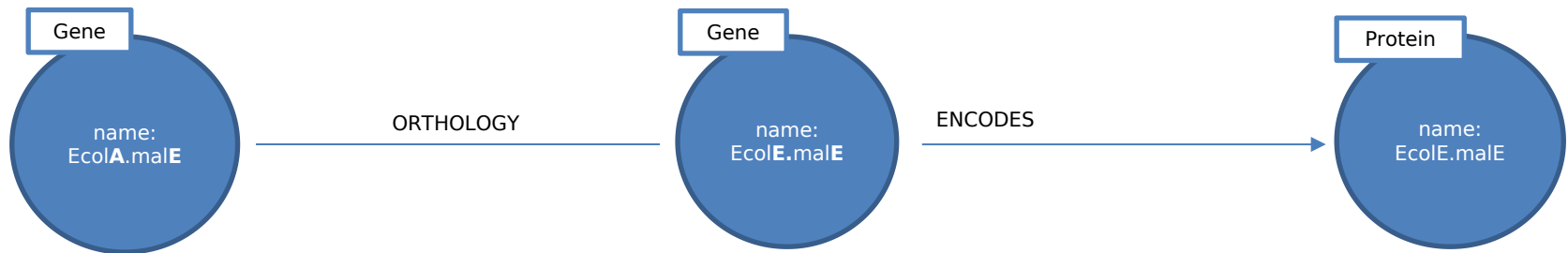


- **reliable** – with full ACID transactions
- durable and fast – disk-based, native storage engine
- **scalable** – up to several billion nodes/relationships/properties
- highly-available – when distributed (replicated)
- **expressive** – powerful, human readable graph query language
- fast – powerful traversal framework
- **embeddable** - in Java program
- simple – accessible by REST interface & Java API

- 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

- Neo4j : Labeled property graph

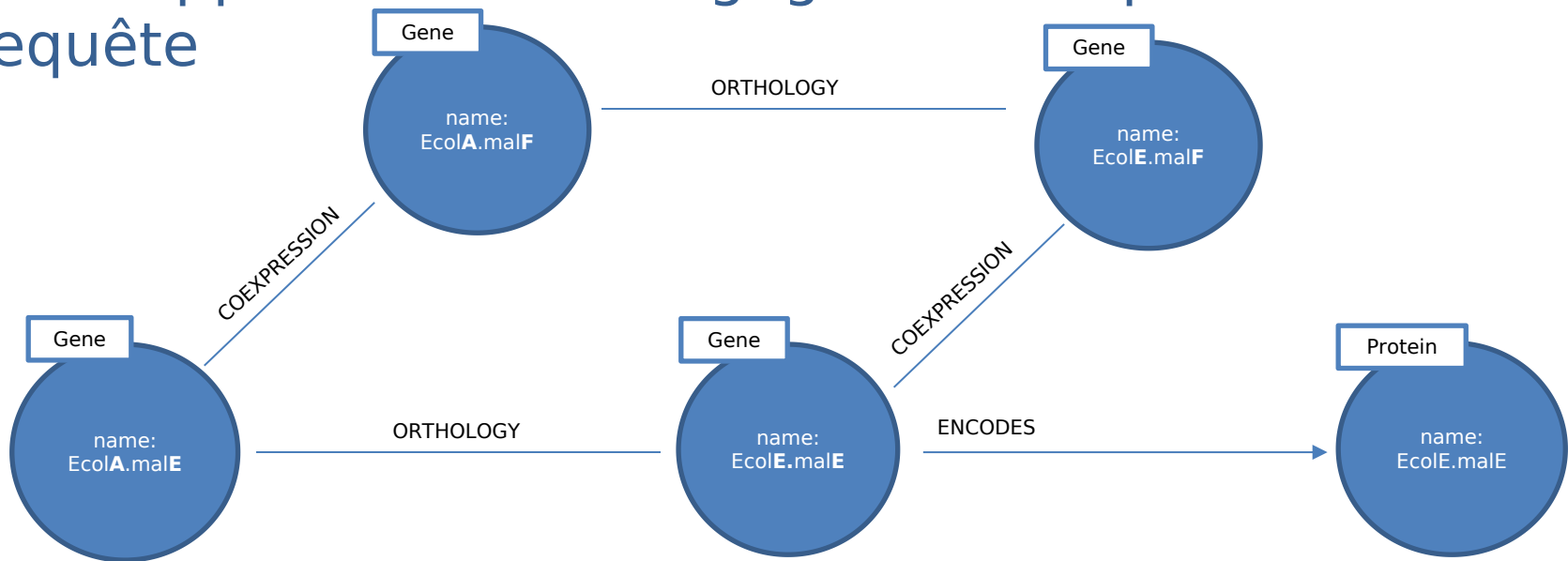


Labeled Property Graph

`(EcolA.malE:Gene) - [:ORTHOLOGY] - (EcolE.malE:Gene) - [:ENCODES] ->(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



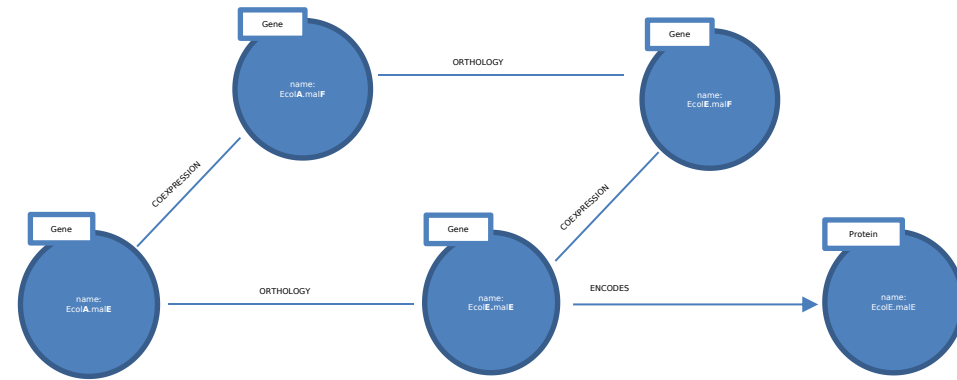
Labeled Property Graph

```

(EcoIE.maIE:Gene) -[:COEXPRESSION]- (EcoIE.maIF) -[:ORTHOLOGY]-
(EcoIA.maIF:Gene) -[:COEXPRESSION]- (EcoIA.maIE:Gene) -[:ORTHOLOGY]- (EcoIE.maIE:Gene)
-[:ENCODES]->(EcoIE.maIE:Protein)
  
```

Labeled property graph

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



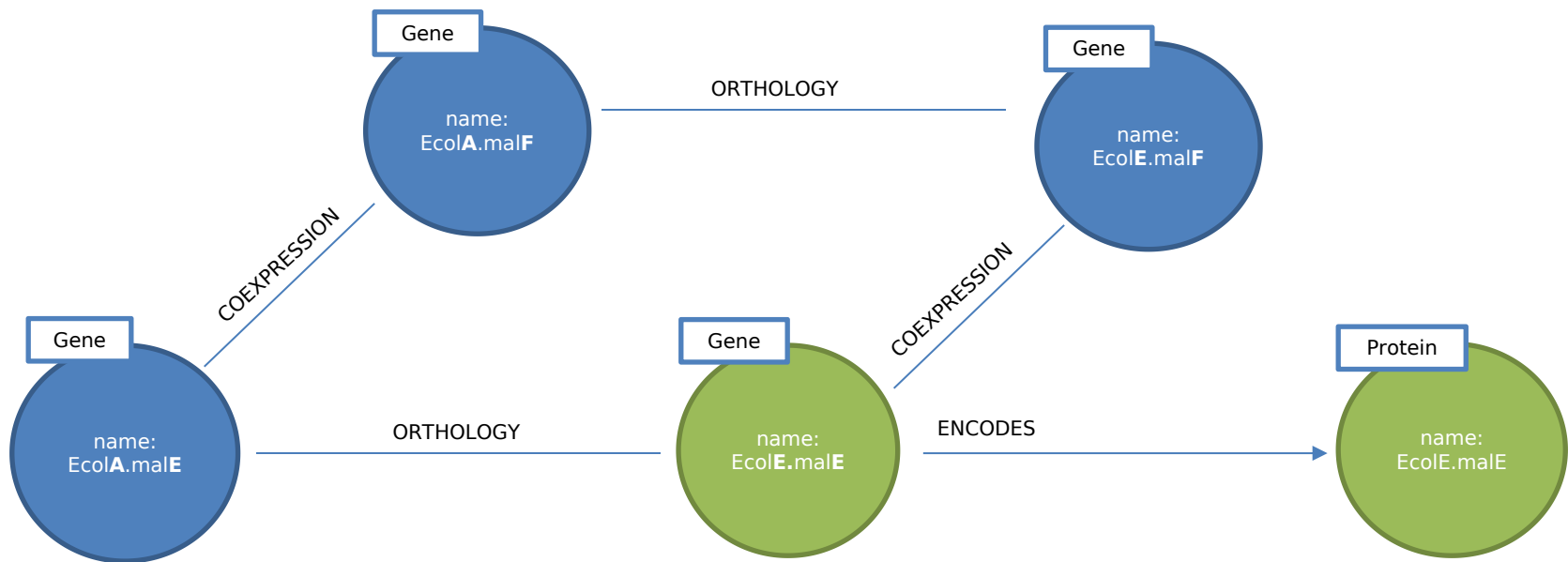
- Propriétés des sommets : de type clé/valeur
- Étiquettes des sommets : une ou plusieurs afin de les regrouper (Gene, Protein)
- Relations : orientées, peuvent avoir des propriétés comme les sommets.

Langage de requête, exemple : Cypher

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

OU

```
MATCH (g:Gene {name: 'EcoLE.malE'}) -[:ENCODES]->(p:Protein)
RETURN g,p
```



- Neo4j graph **query language**
 - For querying and updating
- Declarative – we say **what** we want
 - **Not how** to get it
 - **Not** necessary to express **traversals**
- Human-readable
- Inspired by SQL and SPARQL
- Still growing = syntax changes are often

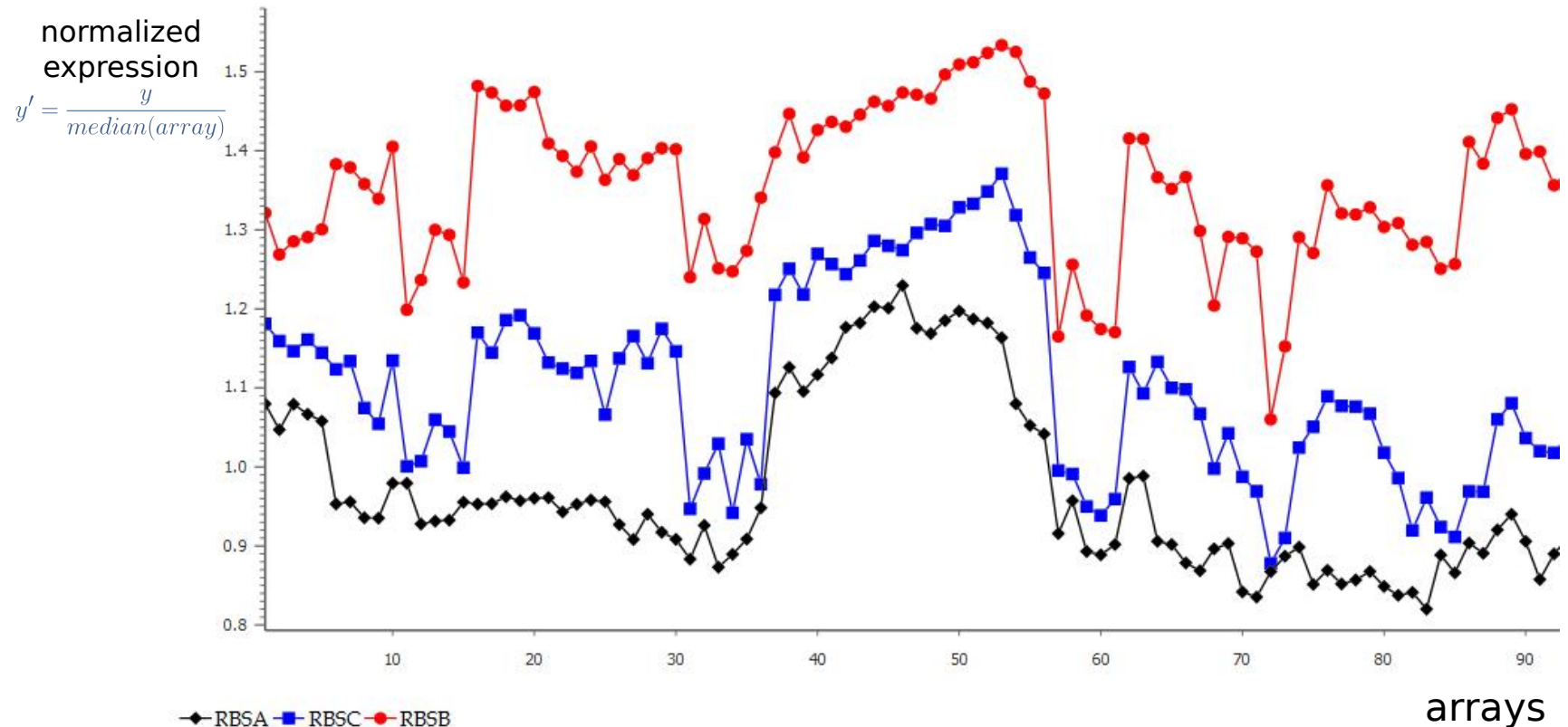
<http://neo4j.com/docs/stable/cypher-query-lang.html>

- **MATCH:** The graph **pattern** to match
- **WHERE:** **Filtering** criteria
- **RETURN:** What to return
- **CREATE/MERGE:** Creates nodes and relationships.
- **DELETE:** Remove nodes, relationships, properties
- **SET:** Set values to **properties**
- **WITH:** Divides a query into multiple parts

- **Embedded** database in Java system
- **Language**-specific connectors
 - **Libraries** to connect to a running Neo4j server (python, R, ...)
- **Cypher** query language
 - Standard language to **query** graph data
- **HTTP REST** API
- etc.



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

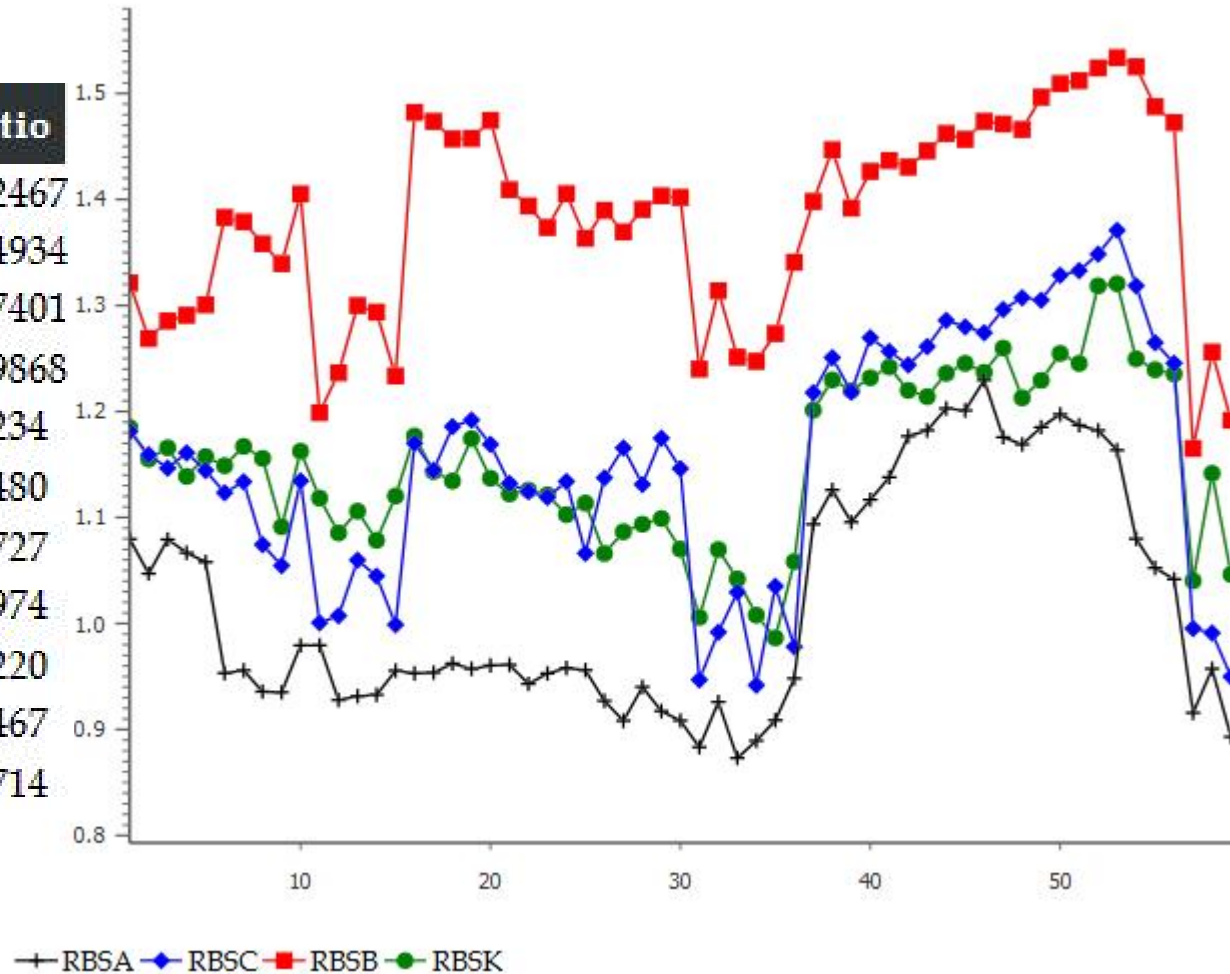


gene pairwise dissimilarity matrix

Gene expression illustration

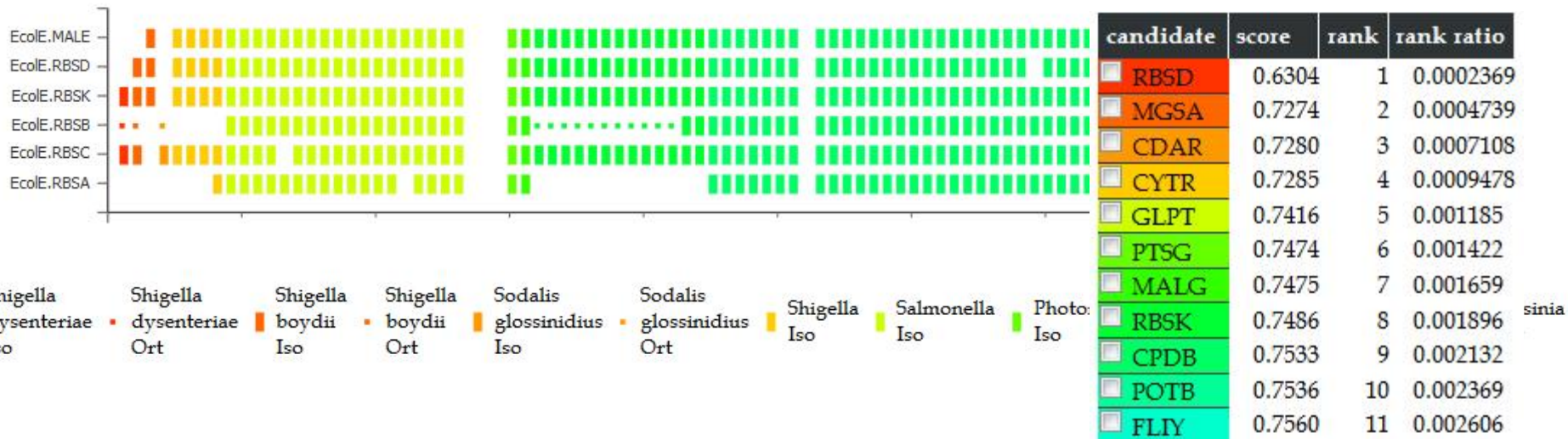
- training: rbsA, rbsB, rbsC in *E. coli* K-12

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

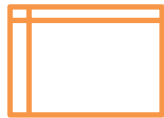


Phylogenetic profiles

training: rbsA, rbsB, rbsC



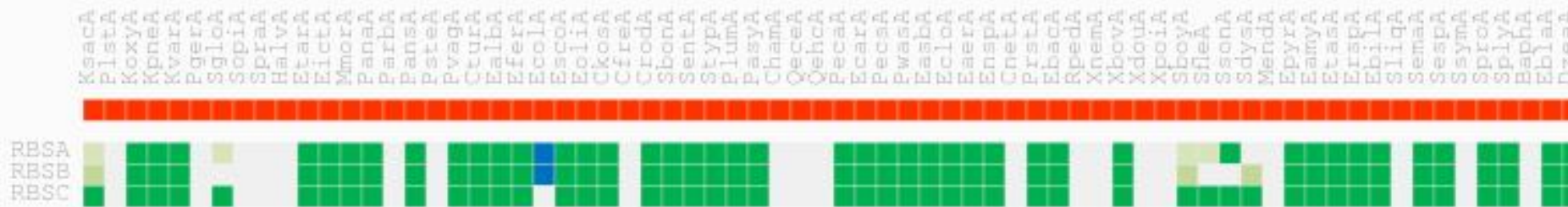
- a gene: presence/absence of orthologs 1:1 in other genomes
- pair of genes: dissimilarity index based on the Jaccard index
- score: average dissimilarity



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

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

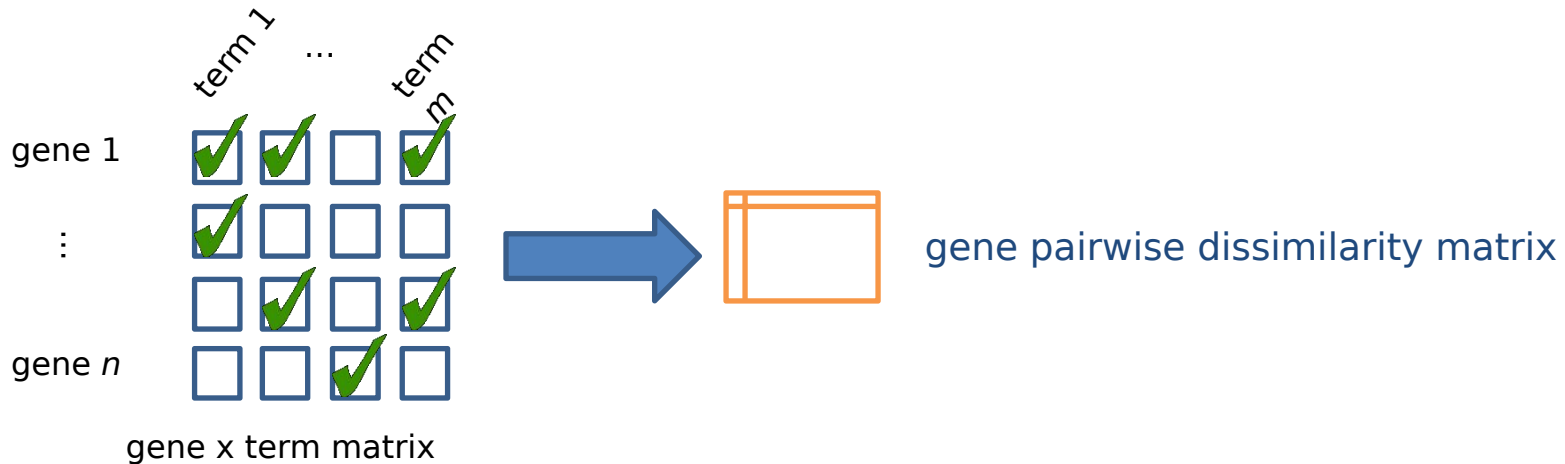
- M depends on the pair of genes considered
- d not comparable between genes (e.g. $0.1^6 = 10^{-6}$ vs. $0.5^{20} = 9.5 \cdot 10^{-7}$)
- normalization: log transformation, z-score, average of distance matrix and its transpose

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]

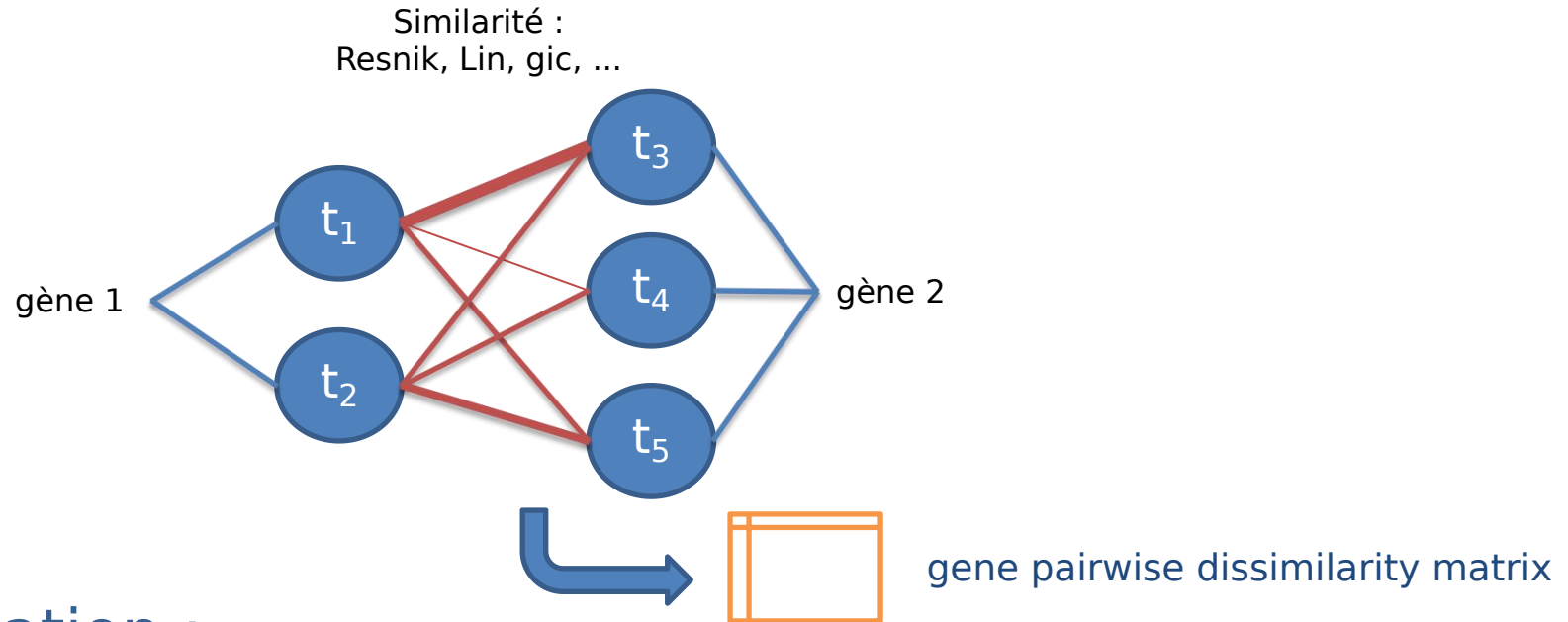
- $sim_{lin}(t_1, t_2) = IC(MICA(t_1, t_2)) / (IC(t_1) + IC(t_2))$ [Lin, 1998]

- $sim_{gic}(t_1, t_2) = \frac{\sum_{t \in \{GO(t_1) \cap GO(t_2)\}} IC(t)}{\sum_{t \in \{GO(t_1) \cup GO(t_2)\}} IC(t)}$ [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 : $\text{ave}(t_1-t_3, t_2-t_5)$

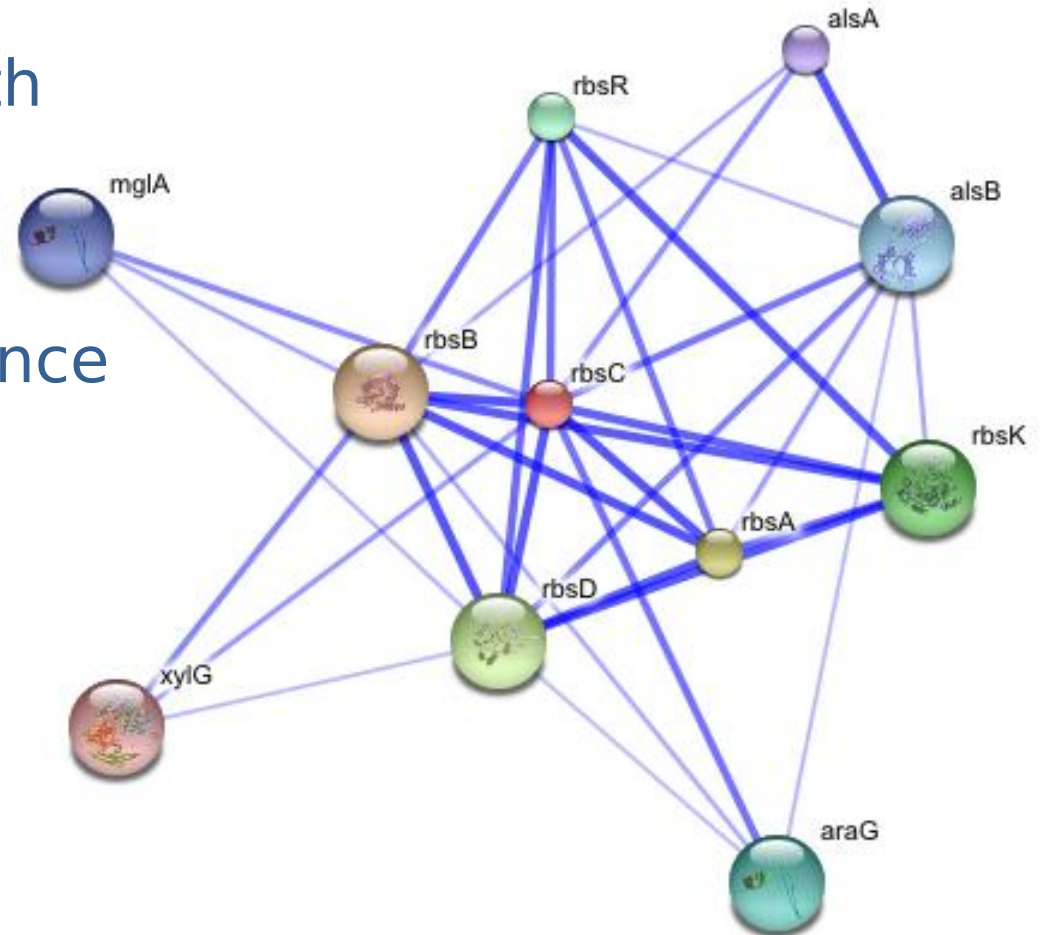


Application :

- 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>