# Intégration de données hétérogènes
# Enrichissement
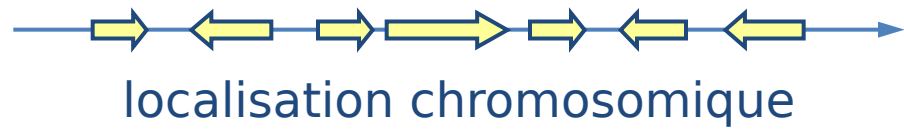
Master 2

Bioinformatique et Biologie des Systèmes
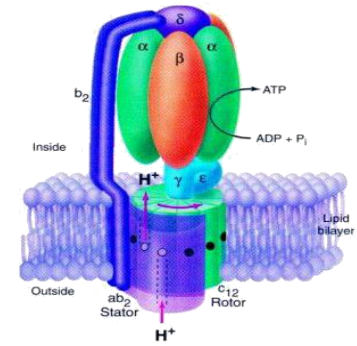
localisation chromosomique

voies
métaboliques

complexes
protéiques

ensembles de gènes

Gene
Ontology

domaines
protéiques

co-citation

# Définitions

- (Identifiants de) <u>gène</u> → ARNm → protéine

- $G$ : ensemble des gènes d'un organisme

- *Fonction de regroupement* : relation entre gènes basée sur un indice de similarité.

- *Ensemble de (gènes) voisins* : ensemble de gènes $E$ $\subseteq G$ regroupés par une fonction de regroupement.

- *Voisinage* : sous-ensemble de P($G$) formant un ensemble d'ensembles de voisins, $V \subseteq$ P($G$), regroupés par une même fonction de regroupement.



$$V = \{E_1, E_2, E_3, E_4, E_5, E_6\} \subseteq \text{P}(G)$$

- Un voisinage est un ensemble (d'ensembles de voisins) ordonné par la relation d'inclusion $\subseteq$



diagramme de Hasse de $V$

$V = \{E_1, E_2, E_3, E_4, E_5, E_6\}$

20 S Proteasome (yeast) closed state

α
β
β
α

two views

PDB 1JD2



un complexe → un ensemble de protéines

# Exemple de critère de regroupement : voies métaboliques



[Kanehisa et Goto, 2000]

une voie métabolique → un ensemble de protéines

clustering hiérarchique
des profils



Conditions expérimentales

un cluster → un ensemble de gènes

# • Recherche d'ensembles similaires

**ensemble requête**
$Q \subseteq G$



**Quels sont les
ensembles cibles
qui lui sont similaires ?**

base de données
de voisinages :
**ensembles cibles**

# Mesure de (dis)similarité



- Loi hypergéométrique : probabilité d'avoir au moins le nombre d'éléments communs observé entre 2 échantillons issus d'une même population

$$p-valeur(c,t,q,g) = \sum_{k=c}^{\min(q,t)} \frac{\binom{t}{k}\binom{g-t}{q-k}}{\binom{g}{q}}$$

*avec*
  - $g = |G|$ : taille de la population
  - $q = |Q|$ : taille de l'ensemble requête
  - $t = |T|$ : taille de l'ensemble cible
  - $c = |Q \cap T|$ : nombre d'éléments communs
- Autres mesures :
  - Loi binomiale
  - $\chi^2$
  - ratio, pourcentage

- # Recherche d'ensembles similaires

**base de données
de voisinages :
ensembles cibles**

**ensemble requête**
$Q \subseteq G$

$Q$

**Quels sont les
ensembles cibles
qui lui sont similaires ?**

- **Probabilité d'obtenir une p-valeur aussi faible par hasard : fonction de répartition des p-valeurs minimales**

- **Simulations**

RandomSet_1, minPi = M1
RandomSet_2, minPi = M2
.
.
RandomSet_n, minPi = Mn

Étant donnée une p-valeur $p$
Combien ont un meilleur score ?



y = pourcentage d'ensembles aléatoires ayant une *p-valeur* ≤ x

0,05

0,0002

x = *p-valeur*

levure *Saccharomyces cerevisiae*
n=500, q=9, g=5786, KEGG Pathways

[Barriot *et al.* 2004]

*Saccharomyces cerevisiae*
n=500, q=6-9-200-500-1000,
g=5786, KEGG Pathways

*Saccharomyces cerevisiae*
n=500, q=50, g=5786,
—— GO molecular function,
- - - Ferea *et al.*, 1999

$N = \{S_1, S_2, S_3, S_4, S_5, S_6\}$

Hasse diagram of $N$

a target set $T$ is **pertinent** if
$$Q \cap T \neq \emptyset$$
and
$$\nexists\, T' \in N \text{ such that } T' \subset T \text{ and } T' \cap Q = T \cap Q$$
and
$$\nexists\, T' \in N \text{ such that } T \subset T' \text{ and } T' - Q = T - Q$$

# Pertinence definition

- Q a non empty query set
- N a neighborhood
- a target set T ∈ N
- T pertinent if

$$Q \cap T \neq \varnothing$$

and

$$\nexists\ T' \in N \text{ such that } T' \subset T \text{ and } T' \cap Q = T \cap Q$$

and

$$\nexists\ T' \in N \text{ such that } T \subset T' \text{ and } T' - Q = T - Q$$

$N_1$

$T_1 = \{a,b,x\}$

$T_2 = \{a,b\}$

$T_3 = \{a\}$

$N_2$

$T_4 = \{a,b,x,y\}$

$T_5 = \{a,b,x\}$

$T_6 = \{a,x\}$

$Q = \{a,b,e\}$
$Q' = \{a,b\}$

Local decision

$|c| > 0$
$|d| < \min(\{d_{parents}\})$
$|c| > \max(\{c_{children}\})$

$c' = c$
$d' \supset d$

$c' \supset c$
$d' \supset d$

$c' \supset c$
$\mathbf{d' = d}$

$c = T \cap Q$
$d = T - Q$
$T = c \cup d$

$c' \subset c$
$d' \subset d$

$c' \subset c$
$d' = d$

$\mathbf{c' = c}$
$d' \subset d$

[Barriot, Dutour, Sherman, 2007, *BMC Bioinformatics*]

# Illustration

- Pertinence des comparaisons & redondance des résultats

# A small portion of the DAG is searched



sets having no common elements are not interesting

sets having common elements and that have bad *p-values*

up to millions of target sets in the DAG of the neighborhood

sets having common elements and that may have good *p-values*

query elements

thousands of elements (genes or proteins)

# Complex 440.30.10 mRNA splicing

| GO Term | Description | Target size | Common elements |
|---|---|---|---|
| **GO:0000398** | **nuclear mRNA splicing, via spliceosome** | **84** | **33** |
| GO:0000377 | RNA splicing, via transesterification reactions with bulged adenosine as nucleophile | 84 | 33 |
| GO:0000375 | RNA splicing, via transesterification reactions | 88 | 33 |
| GO:0008380 | RNA splicing | 99 | 33 |
| GO:0006397 | mRNA processing | 108 | 33 |
| GO:0016071 | mRNA metabolism | 132 | 33 |
| **GO:0006396** | **RNA processing** | **262** | **34** |
| GO:0016070 | RNA metabolism | 360 | 34 |
| GO:0043283 | biopolymer metabolism | 812 | 34 |
| GO:0006139 | nucleobase, nucleoside, nucleotide and nucleic acid metabolism | 1057 | 34 |
| GO:0044238 | primary metabolism | 2191 | 34 |
| GO:0044237 | cellular metabolism | 2407 | 34 |
| GO:0008152 | metabolism | 2465 | 34 |
| **GO:0000245** | **spliceosome assembly** | **10** | **5** |
| GO:0006461 | protein complex assembly | 61 | 5 |
| **GO:0006374** | **nuclear mRNA splicing via U2-type spliceosome** | **8** | **8** |
| *GO:0000391* | *U2-type spliceosome dissembly* | *2* | *2* |
| GO:0000390 | spliceosome dissembly | 2 | 2 |
| *GO:0000370* | *U2-type nuclear mRNA branch site recognition* | *2* | *2* |
| GO:0000348 | nuclear mRNA branch site recognition | 2 | 2 |
| **GO:0000393** | **spliceosomal conformational changes to generate catalytic conformation** | **3** | **3** |

- each node has only 1 parent

- Algorithm
  - parses the input with a stack of stacks at the time it is loaded
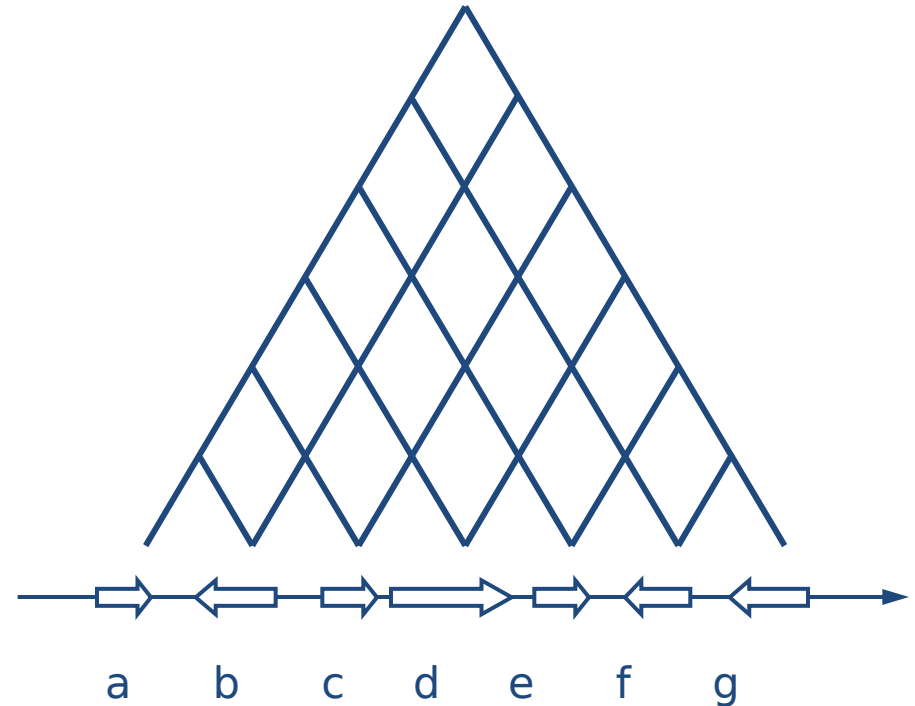  - $O(|G|)$ time



( ( (   a   (   b       c ) )  d )   e     )

**tree**

# Further optimizations (2/2)

- DAG is implicit, e.g. adjacent genes on the chromosome:
  - store the genes order
  - $\Theta(|G|)$ space instead of $\Theta(|G|^2)$
  - each pair of genes defines an interval which defines a set
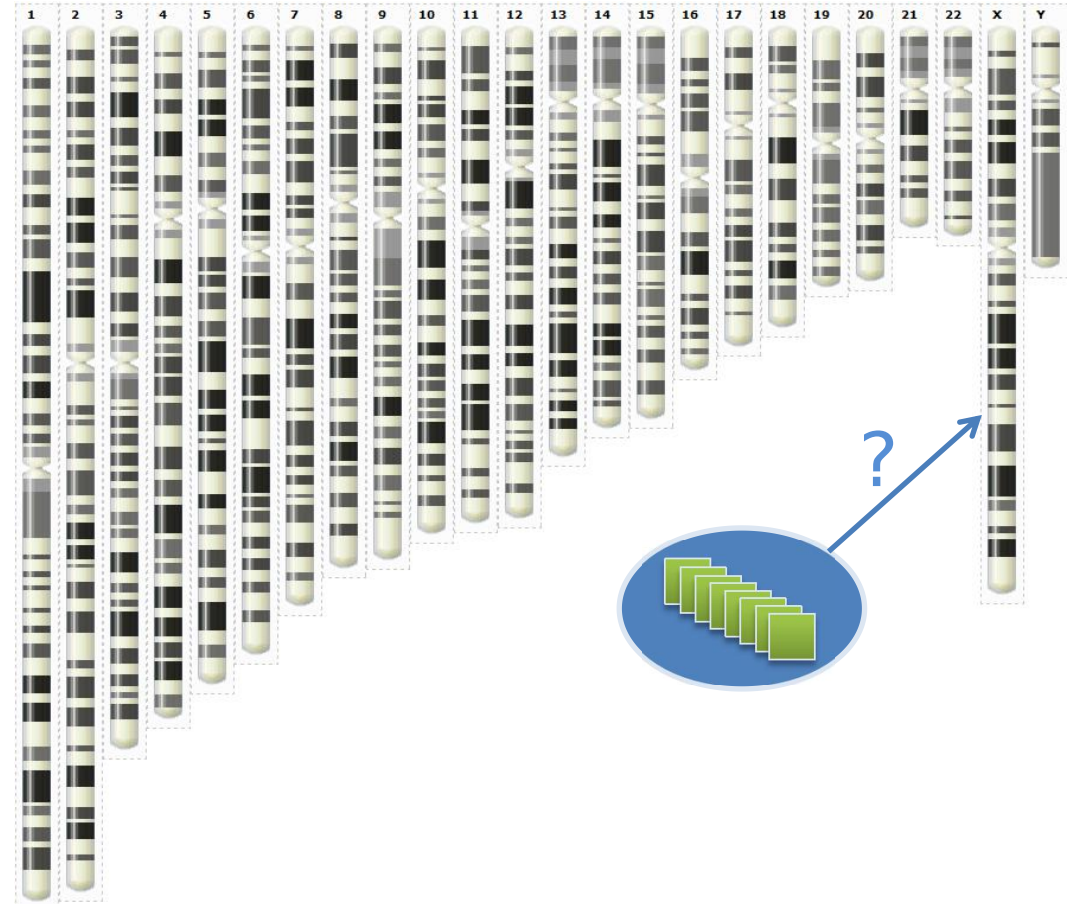- requires a specific algorithm
  - $O(|Q|^2)$ time



**implicit**

## Set of genes of interest

Examples

- Differentially expressed genes
- Co-expressed genes
- Tissue specific genes
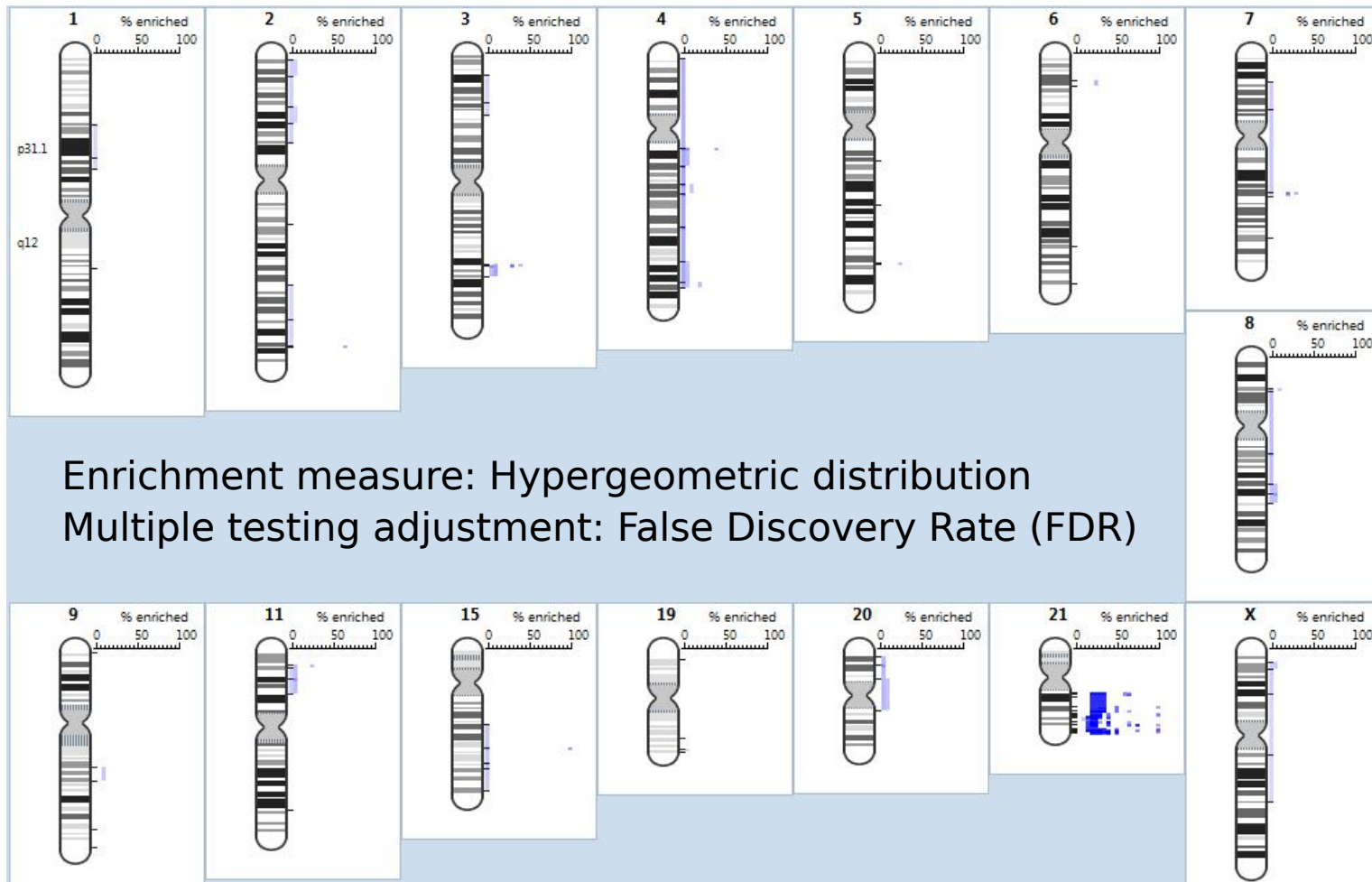- Partners of a protein complex
- Imprinted genes
- ...

→ Question: Do those genes surprisingly cluster in the genome?

Goal: consider every possible region for enrichment

# Down Syndrome differentially expressed genes

## Experiment:

Published list of **differentially expressed genes** in **Down syndrome patients** from Mao, R., C.L. Zielke, H.R. Zielke, and J. Pevsner, Global up-regulation of chromosome 21 gene expression in the developing Down syndrome brain (2003) *Genomics* **81:** 457-467.



Issues:

- Number of regions to test
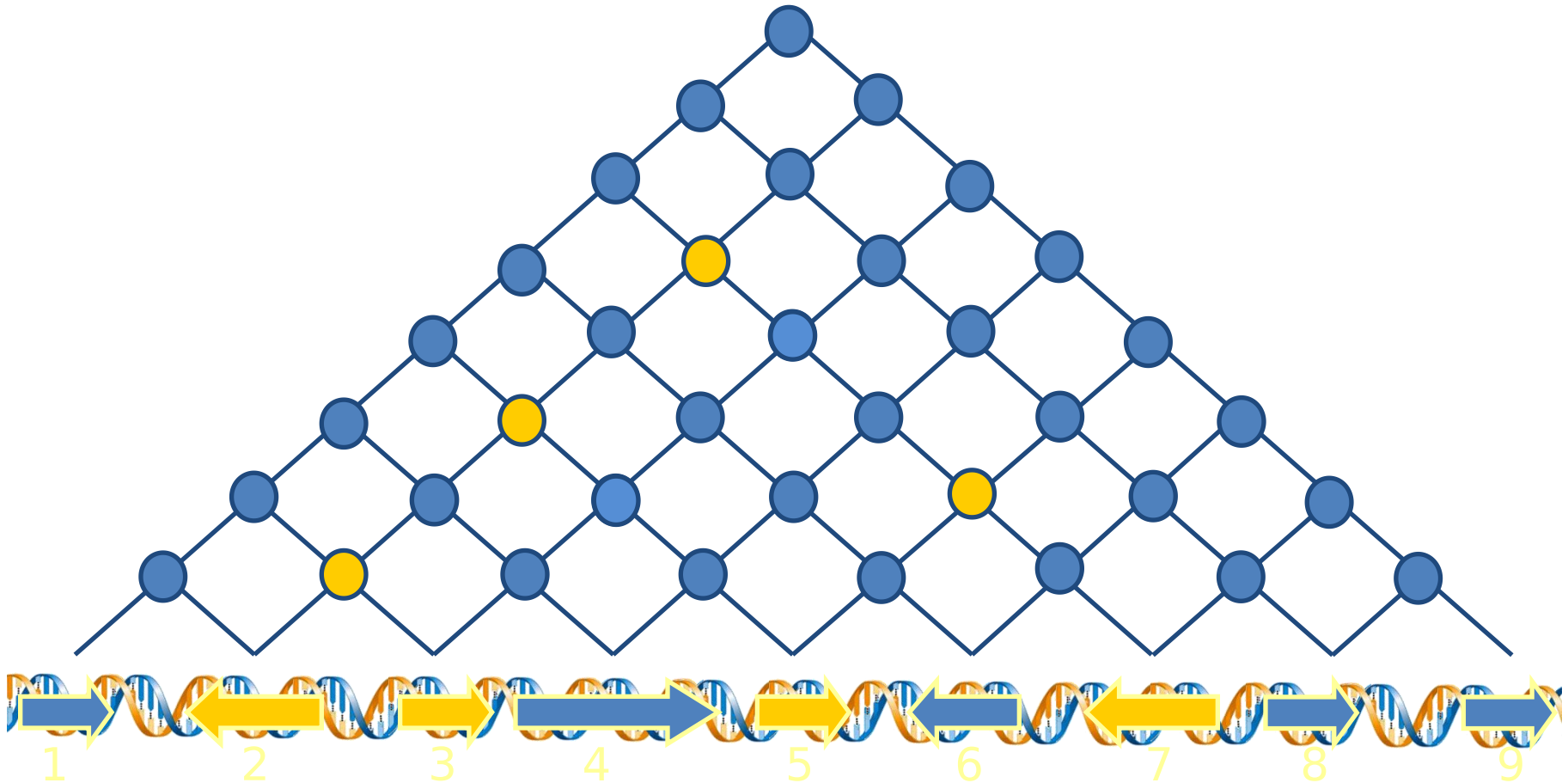
- False positives

- Redundancy

Enrichment measure: Hypergeometric distribution
Multiple testing adjustment: False Discovery Rate (FDR)

# Pertinent Regions
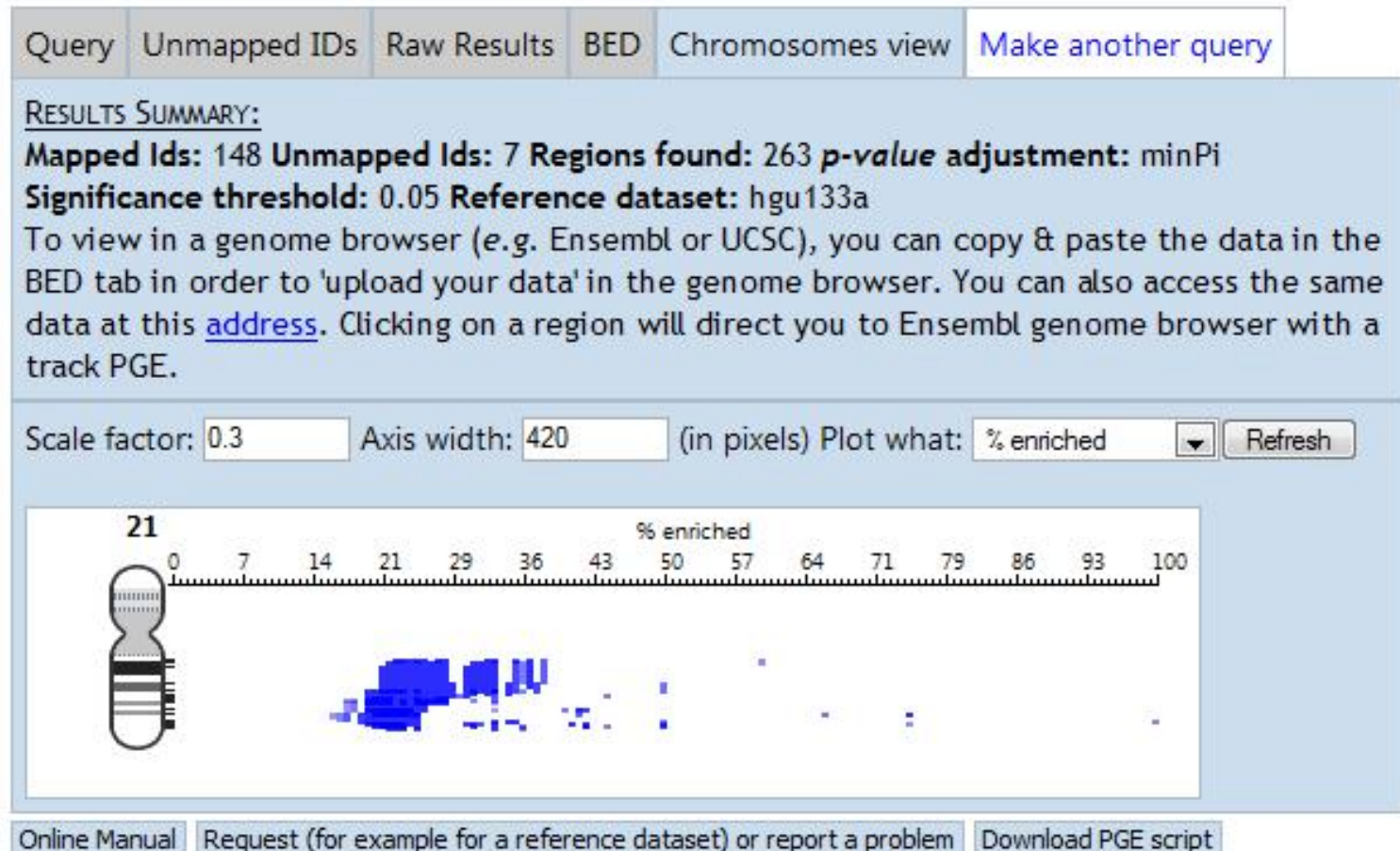
A region is pertinent if it is:
- bounded by genes of interests
- the largest, when genes of interest are consecutive

# Down Syndrome (min$P_i$)



Large regions tend to have smaller *p-values* while small regions tend to have higher percentage of enrichment
→ A smaller region included in a more significant one is pertinent if it has a much higher percentage of genes of interests (>50%)
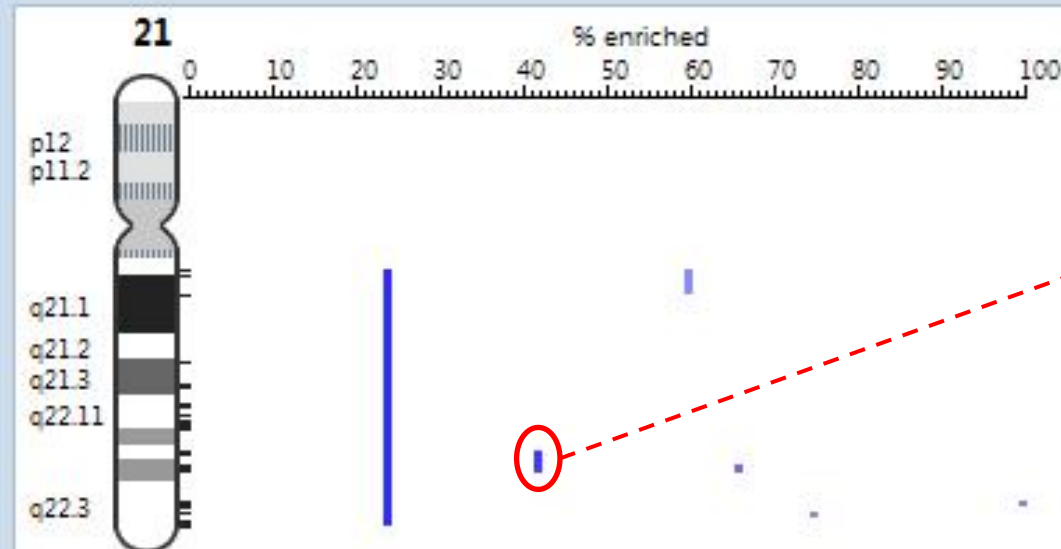
# Down Syndrome d.e.g. Final Results

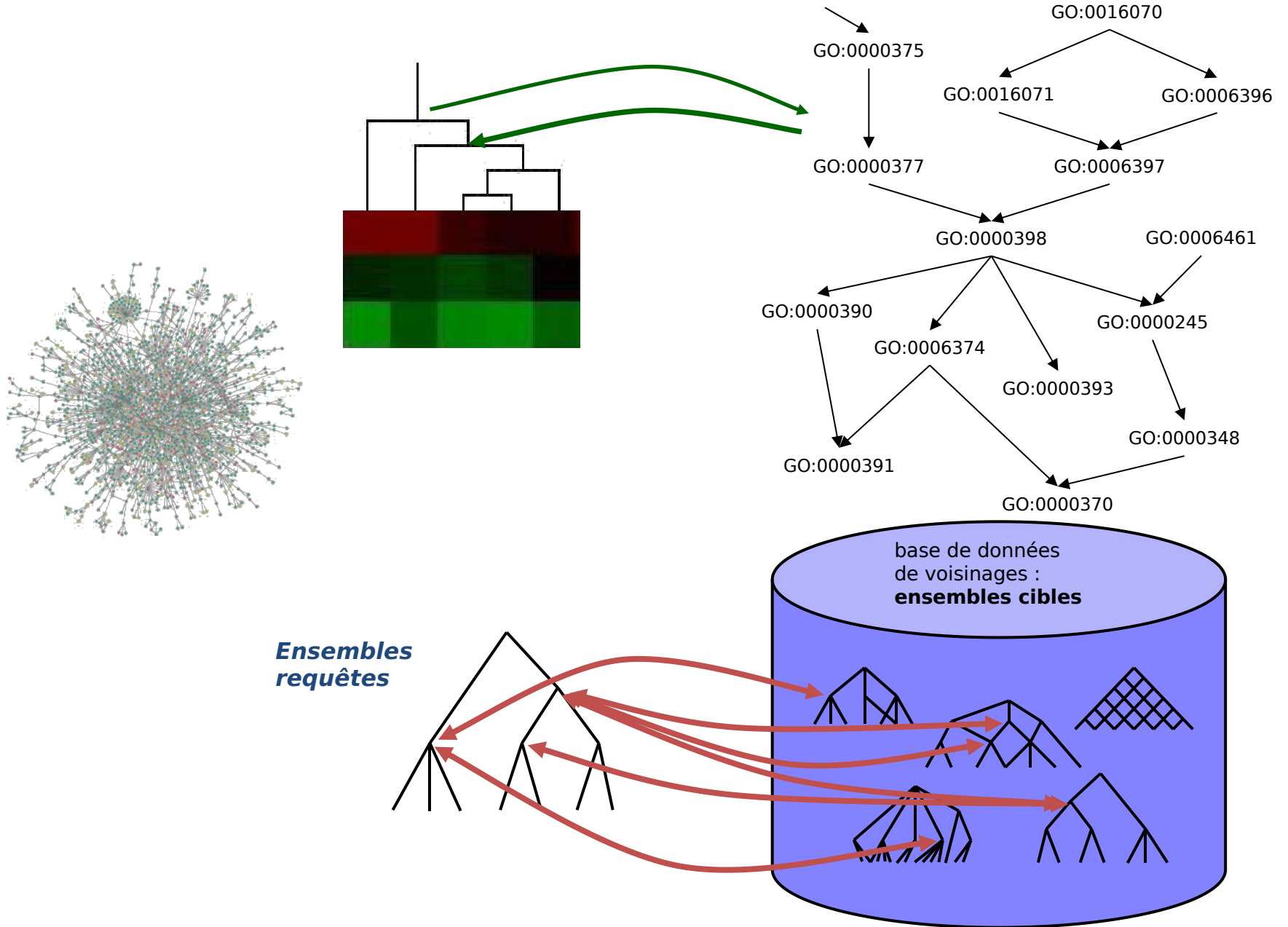| Query | Unmapped IDs | Raw Results | BED | Chromosomes view | Make another query |

**RESULTS SUMMARY:**

**Mapped Ids:** 148 **Unmapped Ids:** 7 **Regions found:** 6 *p-value* **adjustment:** minPi
**Significance threshold:** 0.05 **Reference dataset:** hgu133a

Scale factor: 1     Axis width: 300     (in pixels) Plot what: % enriched ▼
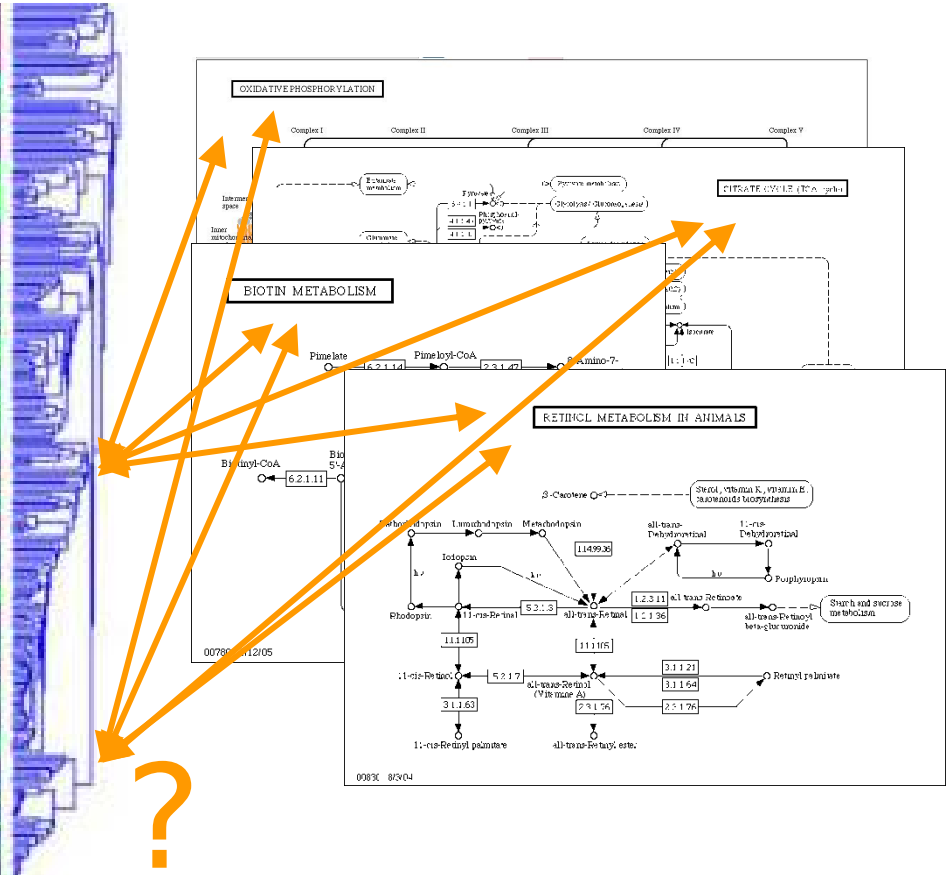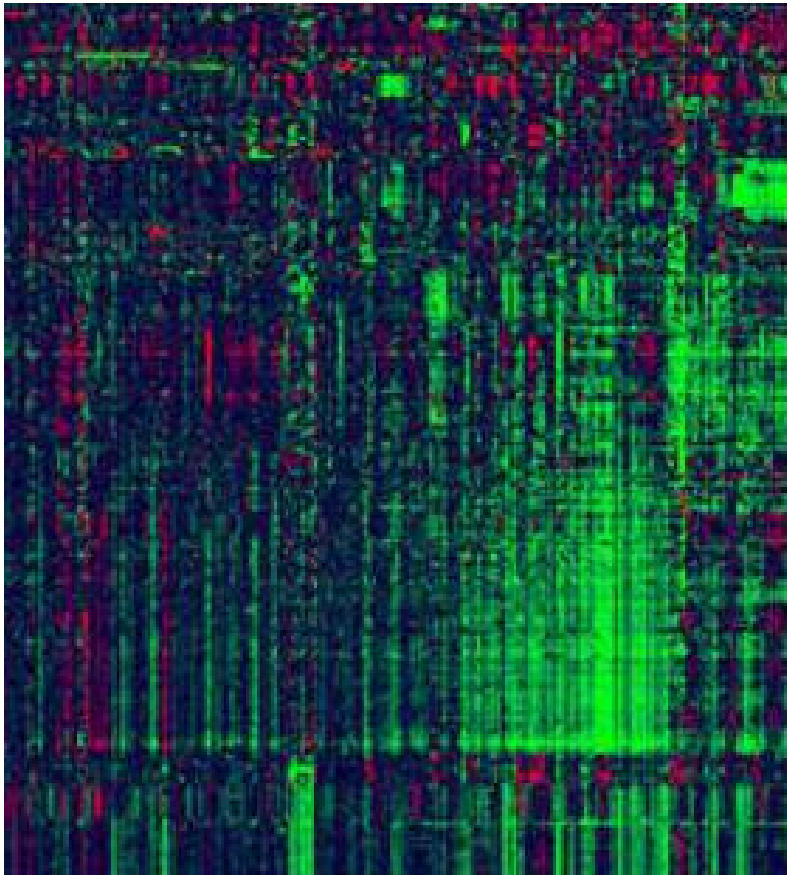
Refresh



```
p-value:            7.87E-10
p-value_adj:        <0.002
Score:              91.039
Score_adj:          INF
Common elements:    6 / 14 (43%)
Overlaping regions: 1
Start of region:    37 359 546 bp
End of region:      40 223 183 bp
Genes:
  DSCR2, DYRK1A, PCP4, PIGP, SH3BGR
  WRB
```

# Défis actuels



GO:0016070

GO:0000375

GO:0016071    GO:0006396

GO:0000377    GO:0006397

GO:0000398    GO:0006461

GO:0000390    GO:0000245

GO:0006374

GO:0000393

GO:0000391    GO:0000348

GO:0000370

base de données
de voisinages :
**ensembles cibles**

*Ensembles
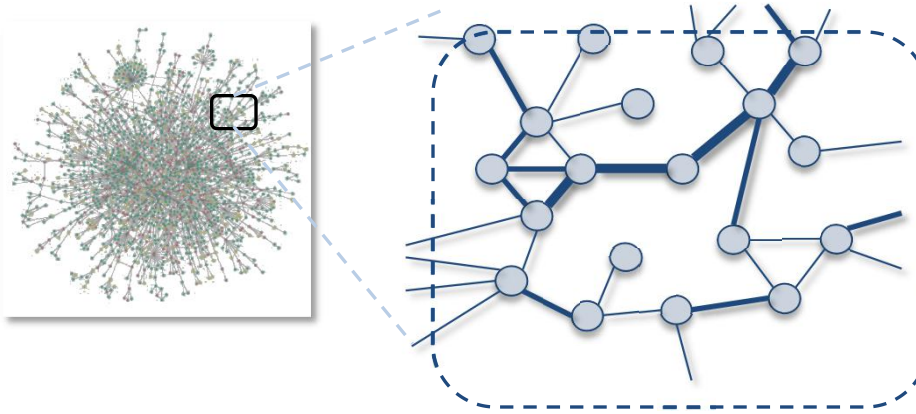requêtes*

# Analyse de données d'expression



[Ferea *et al.*, 1999]                    [Kanehisa & Goto, 2000]

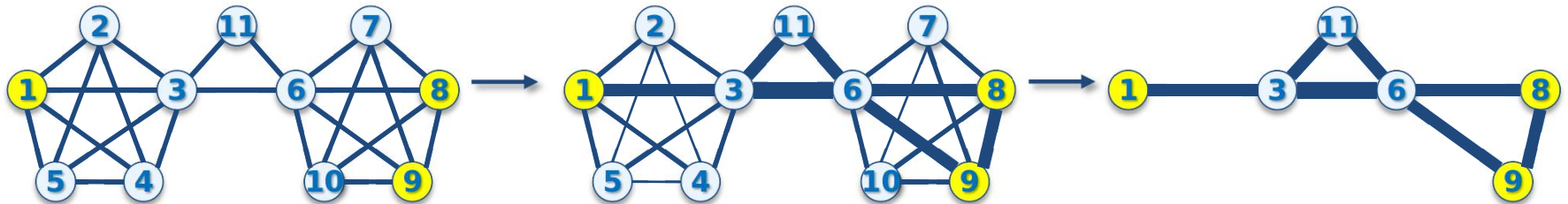# Extraction de sous-graphe pertinent & visualisation



**Idée :**
- Grands graphes d'interactions physiques et/ou fonctionnelles

- Visualiser les relations entre gènes
  d'intérêt

Gènes ayant la même annotation
ex : interaction with host

Marche aléatoire :
pondération des arcs

Surreprésentation :
sous-graphe pertinent



Visualisation du sous graphe
expliquant le mieux ce qui lie
les gènes d'intérêt