

Bioinformatique

Emploi du temps, groupes de TP, supports de cours et de TP disponibles sur <http://silico.biotoul.fr> sur la page Enseignement

Contrôle continu : 30%

Contrôle terminal : 70%

Objectifs pédagogiques :

- Aperçu de quelques domaines d'application de la bioinformatique
 - traitement d'image : mesures phénotypiques
 - gestion des données
 - statistiques : corrélation phénotype - génotype
 - banques de données et sites Web publiques
 - modélisation de systèmes biologiques et simulations

Intervenants

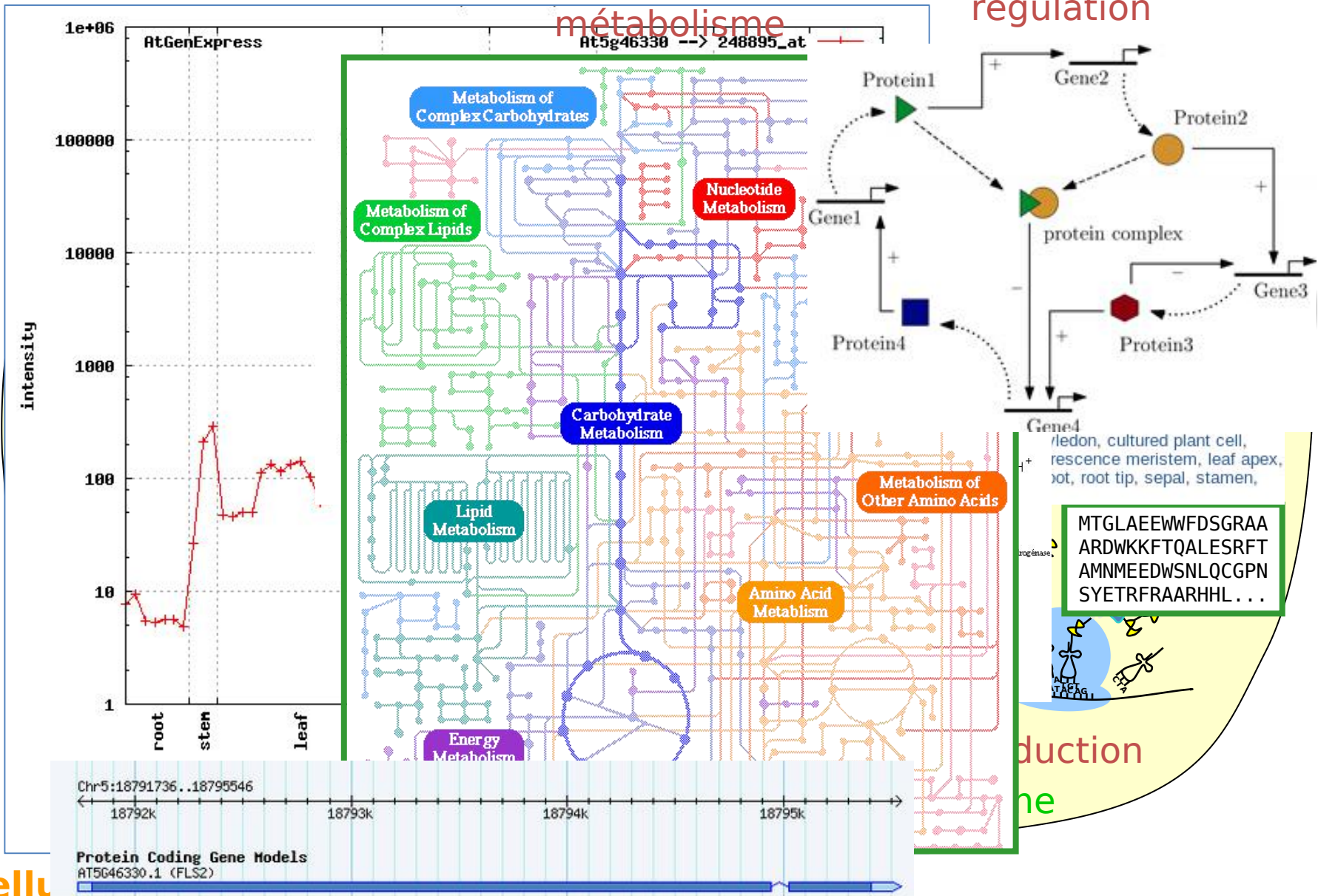
- Roland Barriot
- Maxime Bonhomme
- Franck Delavoie
- Gwennaele Fichant
- Elodie Gaulin
- Florie Gosseau
- Silvia Kocanova

- Qu'est-ce que la bioinformatique ?
- Plusieurs réponses :
 - pas de l'informatique *bio*
 - traitement de l'information biologique
 - un domaine spécialisé de la biologie
 - un domaine de recherche multidisciplinaire : biologie, santé, mathématique, informatique, physique, éthique
- Une définition : traitement des informations biologiques par des méthodes informatiques et/ou mathématiques.
- Les données et connaissances biologiques posent de nouveaux défis spécifiques pour leur gestion, représentation, et leur analyse/traitement

Historique et domaines d'application liés à la bioinformatique

- 1970
- Apparition du terme bioinformatique
 - Disponibilité de séquences de gènes (et du 1^{er} génome complet en 1977, bactériophage phi X174)
 - Comparaison de séquences
- 1980
- Alignement de séquences et recherche de séquences similaires
 - Analyse phylogénétique
 - Prédiction de fonction
 - Banques de données
- 1990
- Débuts de la génomique et génomique comparative ; analyse du contenu d'un génome
 - Débuts des analyses globales de l'expression des gènes et des protéines
 - Famille de gènes, de protéines, de domaines
 - Prédiction de la structure 3D des protéines
- 2000
- Généralisation des approches globales
 - Traitement de graphes : interactions protéine-protéine, réseau de régulation de l'expression des gènes, réseau métaboliques
 - Apprentissage automatique (*data mining*)
 - Fouille de texte (*text mining*)
 - Biologie des systèmes : modélisation de systèmes biologiques
 - Intégration de données hétérogènes, visualisation
 - Médecine personnalisée
 - NGS/Séquençage très haut débit
 - Traitement d'images liées à la microscopie
- 2010
- Ethique, confidentialité des données
 - Modélisation d'une cellule, d'un organisme, d'une population, d'un écosystème
 - Biologie de synthèse

(Quelques) données et connaissances disponibles



Données omics

- **Génome**

- séquence(s) nucléique(s) de l'ensemble des chromosomes d'un organisme
- ensemble des gènes d'un organisme

- **Transcriptome**

- ensemble des ARNm ou transcrits présents dans une cellule ou une population de cellules dans des conditions données

- **Protéome**

- ensemble des protéines présentes dans une cellule ou une population de cellules dans des conditions données

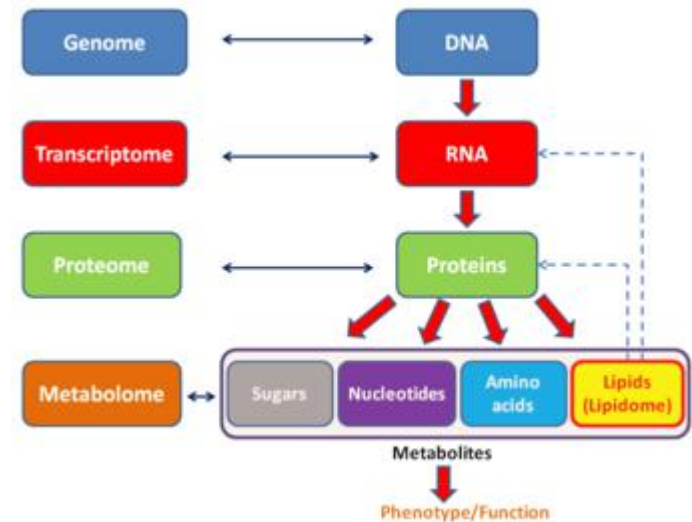
- **Interactome**

- ensemble des interactions moléculaires pouvant survenir *in vivo*
- ensemble des interactions moléculaires dans des conditions données
- ensemble des interactions au sein d'un organisme : moléculaires, physiques, génétiques, fonctionnelles, ...

- **Métabolome**

- ensemble des métabolites présents dans une cellule ou une population de cellules dans des conditions données

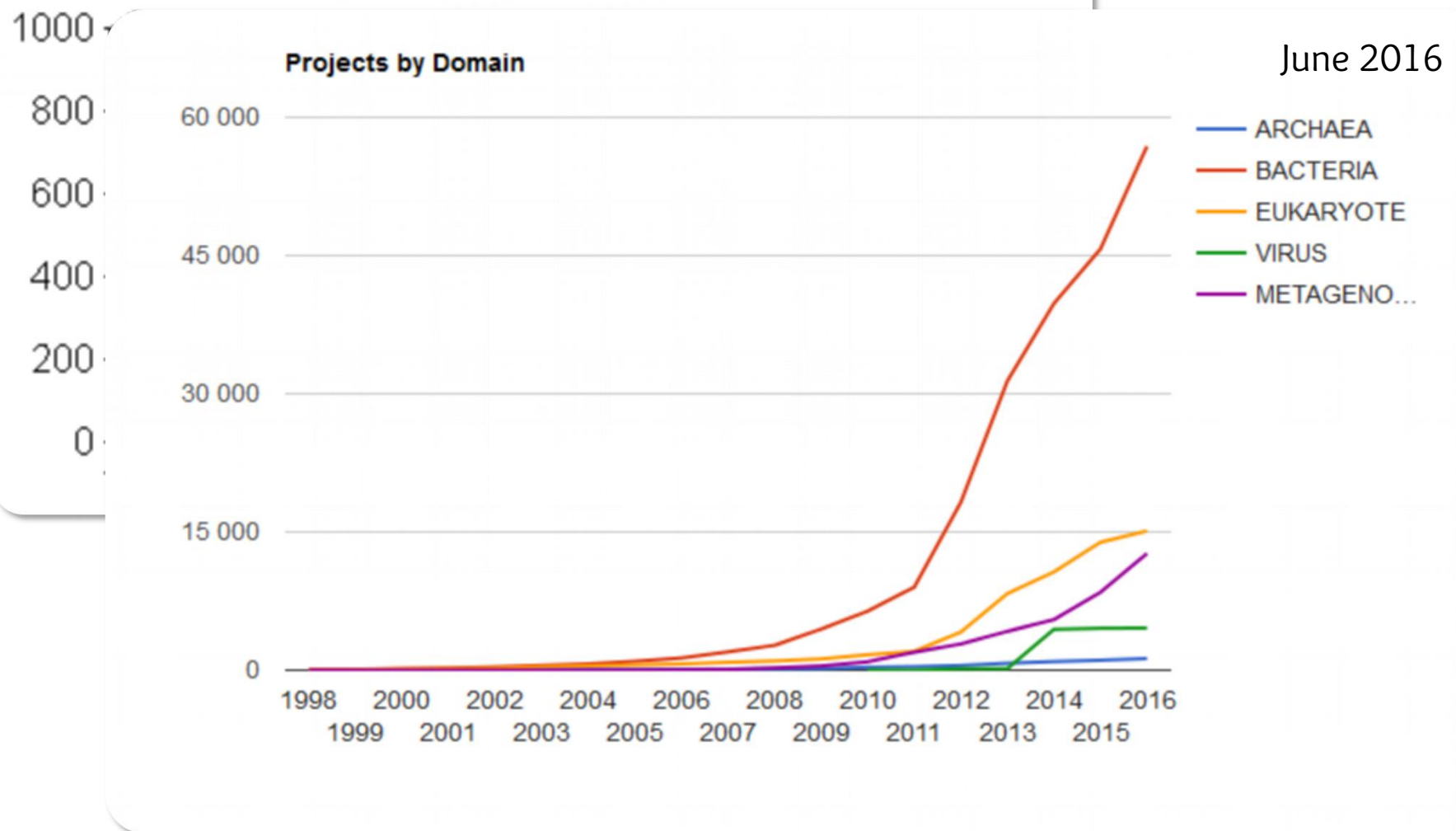
- Exome, lipidome, phénome, régulome, sécrétome, ...



source : Wikipedia

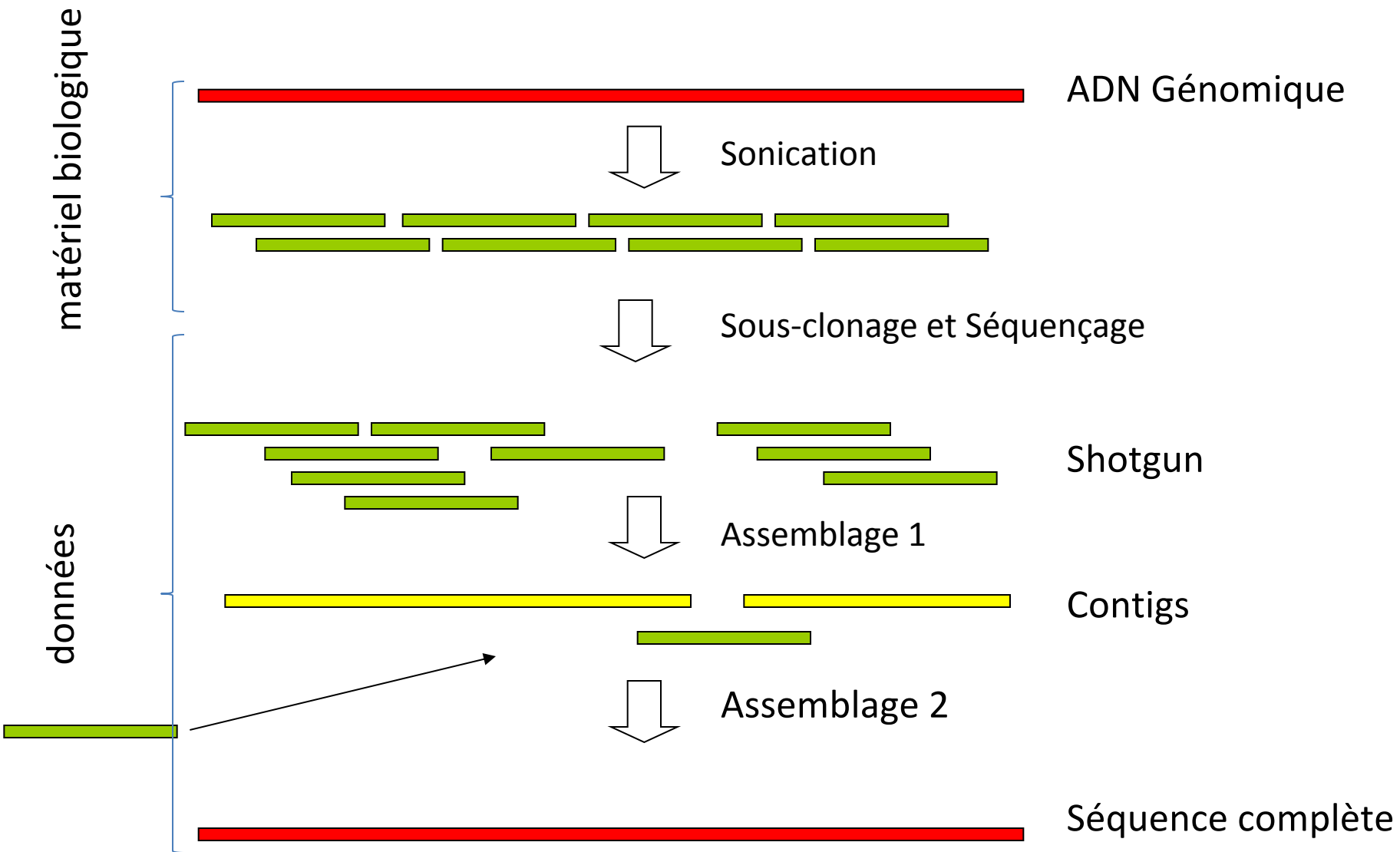
Séquences disponibles : quelques chiffres

Completely Sequenced Genomes © January 2009



- Annotation
 - régions codantes, régions régulatrice, ...
 - prédiction fonctionnelle
- Reconstruction du réseau métabolique
- Analyse des relations génotype/phénotype
- Analyses évolutives
- Conception de puces d'expression
- Identification de protéine
- Prédiction de structure

Assemblage d'un génome

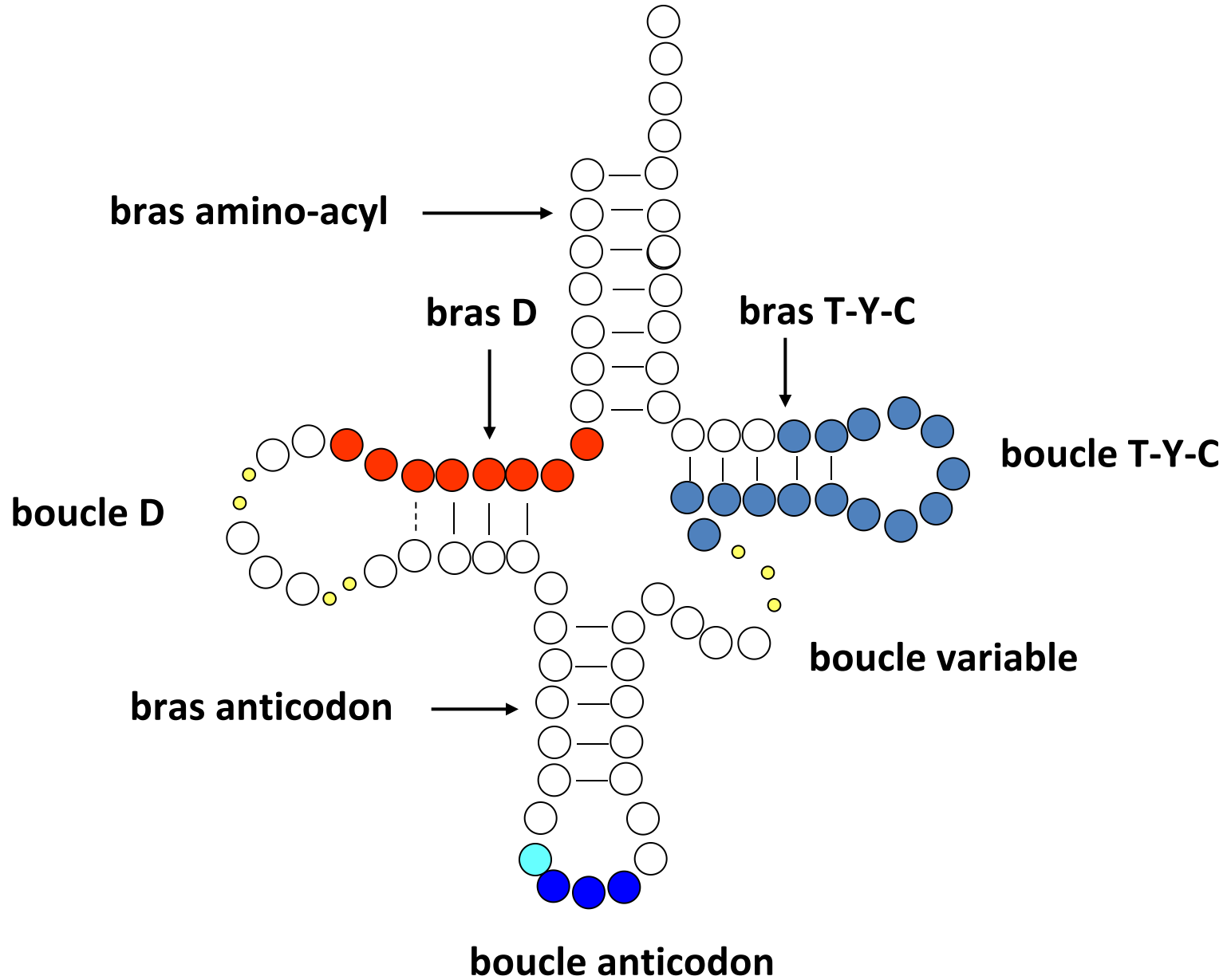


L'analyse manuelle d'une séquence peut s'avérer laborieuse

TCCTGGCCTACATGTTCTTTGGCAAAGGATCTTCAAATCAACGGCTCCCGGTGCGGCGATCATCCATTTCTTCGGAGGGATT
CACGAGATTTACTTCCCGTACATTCTGATGAAACCTGGCCCTGATTCTCGCAGCCATTGCCGGCGGAGCAAGCGGACTCTTA
ACATTACGATCTTTAATGCCGGACTIONGTGCGCGGCAGCGTCACCGGGAAGCATTATCGCATTGATGGCAATGACGCCAAGAGG
AGGCTATTTGCGCGTATTGGCGGGTGTATTGGTTCGCTGCAGCTGTATCGTTTCATCGTTTCAGCAGTGATCCTGAAATCCTCTA
AAGCTAGTGAAGAAGACCTGGCTGCCGCAACAGAAAAAATGCAGTCCATGAAGGGGAAGAAAAGCCAAGCAGCAGCTGC
TTAGAGGCGGAACAAGCCAAAGCAGAGAAGCGTCTGAGCTGTCTCCTGAAAGCGCGAACAAAATTATCTTTTTCGTGTGAT
CCGGGATGGGATCAAGTGCCATGGGGGCATCCATCTTAAGAAACAAAGTGAAAAGCGGAGCTTGACATCAGTGTGACCA
ACACGGCCATTAACAATCTGCCAAGCGATGCGGATATTGTCATCACCCACAAAGATTTAACAGACCCGCGCGAAAGCAAAGCT
GCCGAACGCGACGCACATATCAGTGGATAACTTCTTAACAGCCCGAAATACGACGAGCTGATTGAAAAGCTGAAAAGTAAT
CTTATAGAAAGAGAGTATTGTCATGCAAGTACTCGCAAAGGAAACATTAAGTCAATCAAACGGTATCATCAAAGAAGAGGC
TATCAAATTGGCAGGCCAGACGCTGATTGACAACGGCTACGTGACAGAGGATTACATTAGCAAATGTTTGACCGTGAAGAA
ACGTCTTCTACGTTTATGGGGAATTTCAATTGCCATTCCACACGGCACAGAAGAAGCGAAAAGCGAGGTGCTTCACTCAGGAA
TTTCAATCATAAGATTCCAGAGGGCGTTGAGTACGGAGAAGGCAACACGGCAAAAGTGGTATTCGGCATTGCGGGTAAAA
ATAATGAGCATTTAGACATTTTGTCTAACATCGCCATTATCTGTTTCAGAAGAAGAAACATTGAACGCCTGATCTCCGCTAAAGC
GAAGAAGATTTGATCGCCATTTCAACGAGGTGAACTGACATGATCGCCTTACATTTTCGGTTCGGGAAATATCGGGAGAGGAT
TTATCGGCGCGCTGCTTACCCTCCGGCTATGATGTGGTGTGTTGCGGATGTGAACGAAACGATGGTCAGCCTCCTCAATGA
AAAAAAGAATACACAGTGGAACTGGCGGAAGAGGGACGTTTCATCGGAGATCATTGGCCCGGTGAGCGCTATTAACAGCG
GCAGTCAGACCGAGGAGCTGTACCGGCTGATGAATGAGGCGGCGCTCATCACACAGCTGTTCGGCCCGAATGTCCTGAAG
CTGATTGCCCGTCTATCGCAGAAGGTTTAAAGACGAAGAAATACTGCAAACACACTGAATATCATTGCCTGCGAAAATATGAT
TGGCGGAAGCAGCTTCTTAAAGAAAGAAATATACAGCCATTTAACGGAAGCAGAGCAGAAATCCGTCAGTGAAACGTTAGG
TTTTCCGAATTCTGCCGTTGACCGGATCGTCCCGATTACGATCATGAAGACCCGCTGAAAGTATCGGTTGAACCATTTTTTCG
AATGGGTCATTGATGAATCAGGCTTTAAAGGGAAAACACCAGTCATAAACGGCGCACTGTTTGTGATGATTTAACGCCGTA
CATCGAACGGAAGCTGTTTACGGTCAATACCGGACACGCGGTACAGCGTATGTCGGCTATCAGCGCGGACTCAAACGGT
CAAAGAAGCAATTGATCATCCGGAATCCGCGGTGTTGTTTCATTTCGGCGCTGCTTGAAACTGGTGACTATCTCGTCAAATCGT
ATGGCTTTAAGCAAACACTGAACACGAACAATATATTAATAAATCAGCGGTGCTTTTTAAATCCTTTTCAATTCGGACGATGTGACC
CGCGTAGCGAGGTCACCTCTCAGAAAATGGGAGAAAATGTAGACTTGTAGGCCCGGCAAAGAAAATAAAAGAACCGAAT
GCACTGGCTGAAGGAATTGCCGACGACTGCGCTTCGATTTACCGGTGACCCTGAAGCGGTTGAACTGCAAGCGCTGAT
CGAAGAAAAGGATACAGCGGCGTACTTCAAGAGGTGTGCGGCATTACAGTCCCATGAACCGTTGCACGCCATCATTTTAAAG
AACTTAATCAATAACCGACCCCGTGACACAATGTCACGGGCTTTTTACTATCTCGCAATCTAGTATAATAGAAAGCGCTTA
CGATAACAGGGGAAGGAGAATGACGATGAAACAATTTGAGATTGCGGCAATACCGGGAGACGGAGTAGGAAAGAGGTTGT
AGCGGCTGCTGAGAAAGTGCTTCATACAGCGGCTGAGGTACACGGAGGTTTGTTCATTCTCATTACAGCTTTTCCATGGAGC
TGTGATTACTTGGAGCACGGCAAAAATGATGCCCGAAGATGGAATACATACGCTTACTCAATTTGAAGCAGTTTTTGGGA
GCTGTTCGGAAATCCGAAGCTGGTTCGGATCATATATCGTTATGGGGCTGCTGCTGAAATCCGGAGGGAGCTTGAGCTTTCC

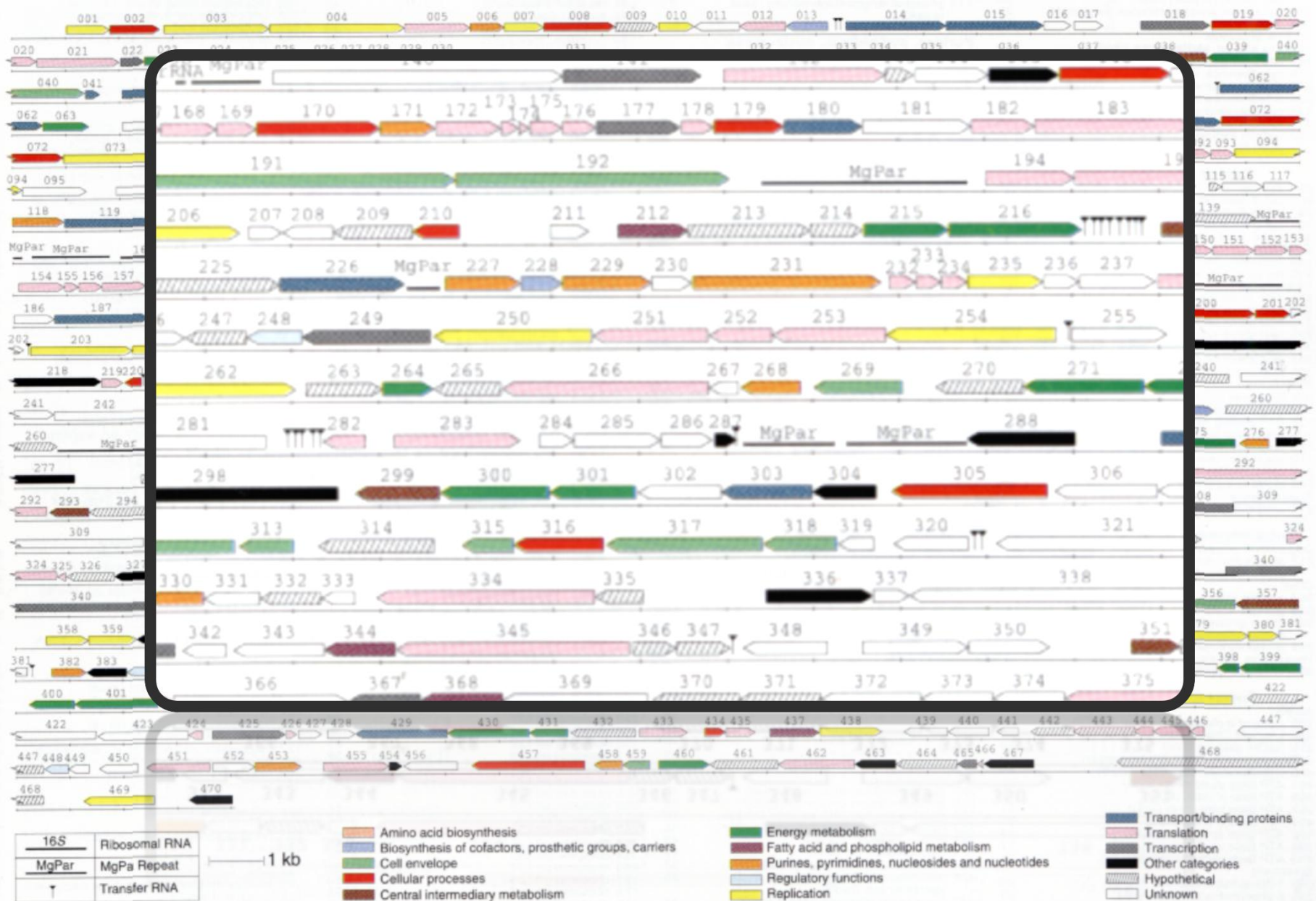
- Identification des gènes codant pour :
 - . les ARNr
 - . les ARNt
 - . les protéines
- Identification des unités de traduction
- Identification des unités de transcription (promoteur et terminateur)
- Pour les gènes codant pour les protéines, prédiction fonctionnelle par recherche de similarité de séquences (Blast) et classification en grandes classes fonctionnelles (ex: biosynthèse des acides aminés, métabolisme énergétique....)

Structure secondaire canonique d'une séquence d'ARNt



Génome de *Mycoplasma genitalium*

Distribution des unités de traduction et classification fonctionnelle



Stockage des séquences : Banques de séquences

ID Q8DPI7_STRR6 PRELIMINARY; PRT; 286 AA.
AC Q8DPI7;
DT 01-MAR-2003, integrated into UniProtKB/TrEMBL.
DT 01-MAR-2003, sequence version 1.
DT 02-MAY-2006, entry version 10.
DE DNA processing Smf protein.
GN Name=smf; OrderedLocusNames=spr1144;
OS Streptococcus pneumoniae (strain ATCC BAA-255 / R6).
OC Bacteria; Firmicutes; Lactobacillales; Streptococcaceae;
OC Streptococcus.
OX NCBI_TaxID=171101;
RN [1]
RP NUCLEOTIDE SEQUENCE [LARGE SCALE GENOMIC
RX MEDLINE=21429245; PubMed=11544234;
RX DOI=10.1128/JB.183.19.5709-5717.2001;
RA Hoskins J., Alborn W.E. Jr., Arnold J., Blaszcak L.
RA DeHoff B.S., Estrem S.T., Fritz L., Fu D.-J., Fuller V
RA Gilmour R., Glass J.S., Khoja H., Kraft A.R., Lagac
RA LeBlanc D.J., Lee L.N., Lefkowitz E.J., Lu J., Matsu
RA McAhren S.M., McHenney M., McLeaster K., Mund
RA Norris F.H., O'Gara M., Peery R.B., Robertson G.T.
RA Sun P.-M., Winkler M.E., Yang Y., Young-Bellido M
RA Zook C.A., Baltz R.H., Jaskunas S.R., Rosteck P.R.
RA Glass J.I.;
RT "Genome of the bacterium Streptococcus pneumo
RL J. Bacteriol. 183:5709-5717(2001).
CC -----
CC Copyrighted by the UniProt Consortium, see http:
CC Distributed under the Creative Commons Attribut
CC -----
DR EMBL; AE008487; AAK99947.1; -; Genomic_DNA.
DR PIR; A95147; A95147.
DR PIR; G98014; G98014.
DR GenomeReviews; AE007317_GR; spr1144.
DR BioCyc; SPNE1313:SPR1144-MONOMER; -.
DR GO; GO:0009294; P:DNA mediated transformati
DR InterPro; IPR003488; SMF.
DR Pfam; PF02481; SMF; 1.
DR TIGRFAMs; TIGR00732; dprA; 1.
KW Complete proteome.
SQ SEQUENCE 286 AA; 31583 MW; CF12DB83AE36
MELFMKITNY EYKLLKKSGL TNQILKVL EYGENVDQE
FQIDDAHLSK EFQKFPFSI LDDCYPWDL EYDAPVLI
CSKQGAKSVE KVIQGLENEL VIVSGLAKGI DTAAHMAA
NKRLLQDYIGN DHLVLSEYGP GEQPLKFHFP ARNRRIAG
AMEEGRDVFA IPGSILDGLS DGCHHLIQEG AKLVTSGQDV LAEFEF
//

Exemple d'une entrée protéique dans la banque de données UniProt

UniProt > UniProtKB Downloads · Contact · Documentation/Help

Search Blast * Align Retrieve ID Mapping *

Search in Protein Knowledgebase (UniProtKB) Query Search Advanced Search > Clear

Q8DPI7 (Q8DPI7_STRR6) ★ Unreviewed, UniProtKB/TrEMBL
 Last modified December 14, 2011. Version 39. History...

Clusters with 100%, 90%, 50% identity | Third-party data text xml rdf/xml gff fasta

Names · Attributes · Ontologies · Sequences · References · Cross-refs · Entry info Customize order

Names and origin

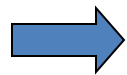
Protein names	Submitted name: DNA processing Smf protein (EMBL AAK99947.1)
Gene names	Name: smf (EMBL AAK99947.1) Ordered Locus Names: spr1144
Organism	Streptococcus pneumoniae (strain ATCC BAA-255 / R6)
Taxonomic identifier	171101 [NCBI]
Taxonomic lineage	Bacteria > Firmicutes > Lactobacillales > Streptococcaceae > Streptococcus

Protein attributes

Sequence length	286 AA.
Sequence status	Complete.
Protein existence	Predicted

Ontologies

Transcriptome : ensemble des ARNm ou transcrits présents dans une population de cellules dans des conditions données.

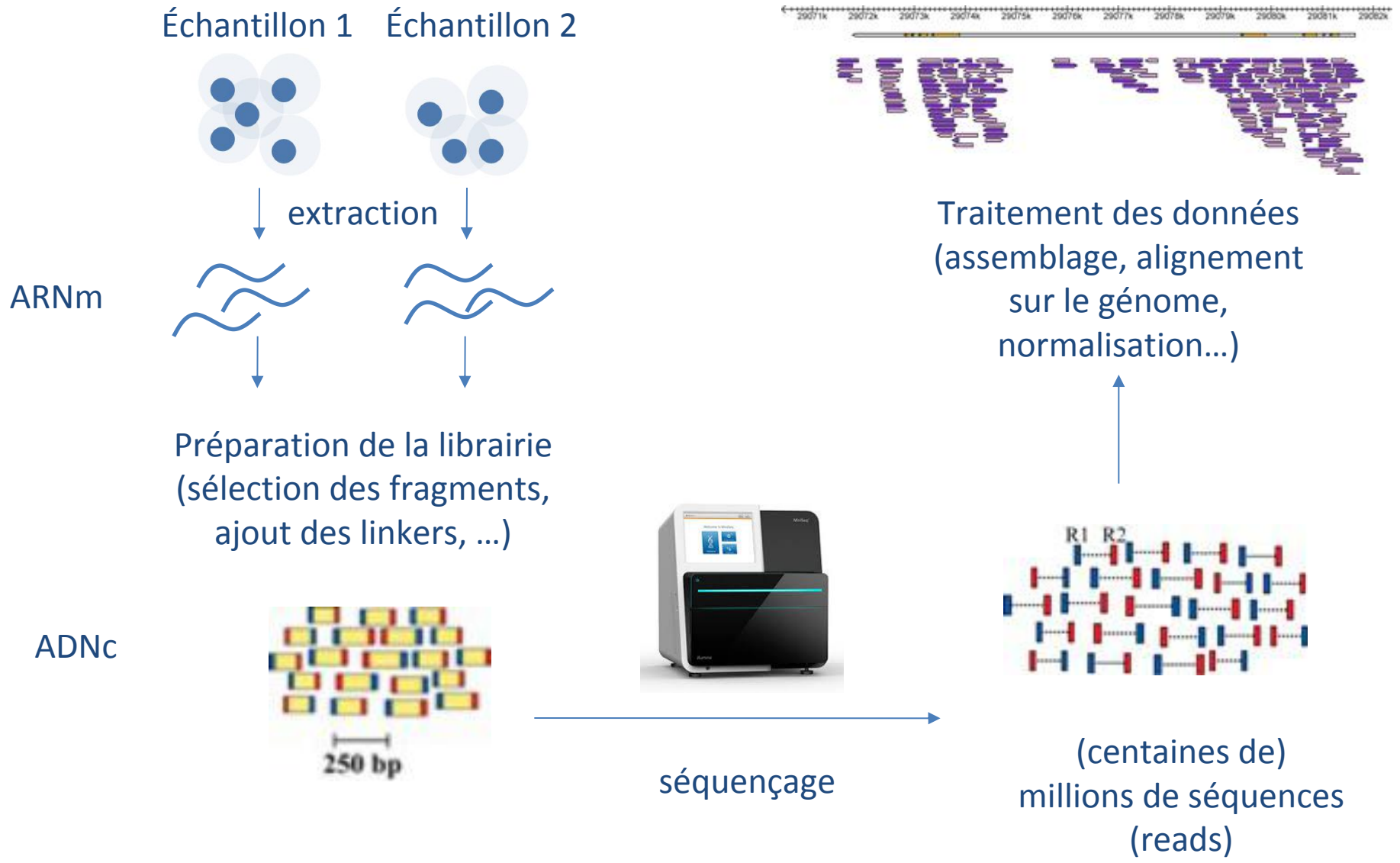


Accès au niveau d'expression de milliers de gènes simultanément (potentiellement l'ensemble des gènes d'un organisme)
= *instantané* de l'état d'une cellule ou d'une population de cellules

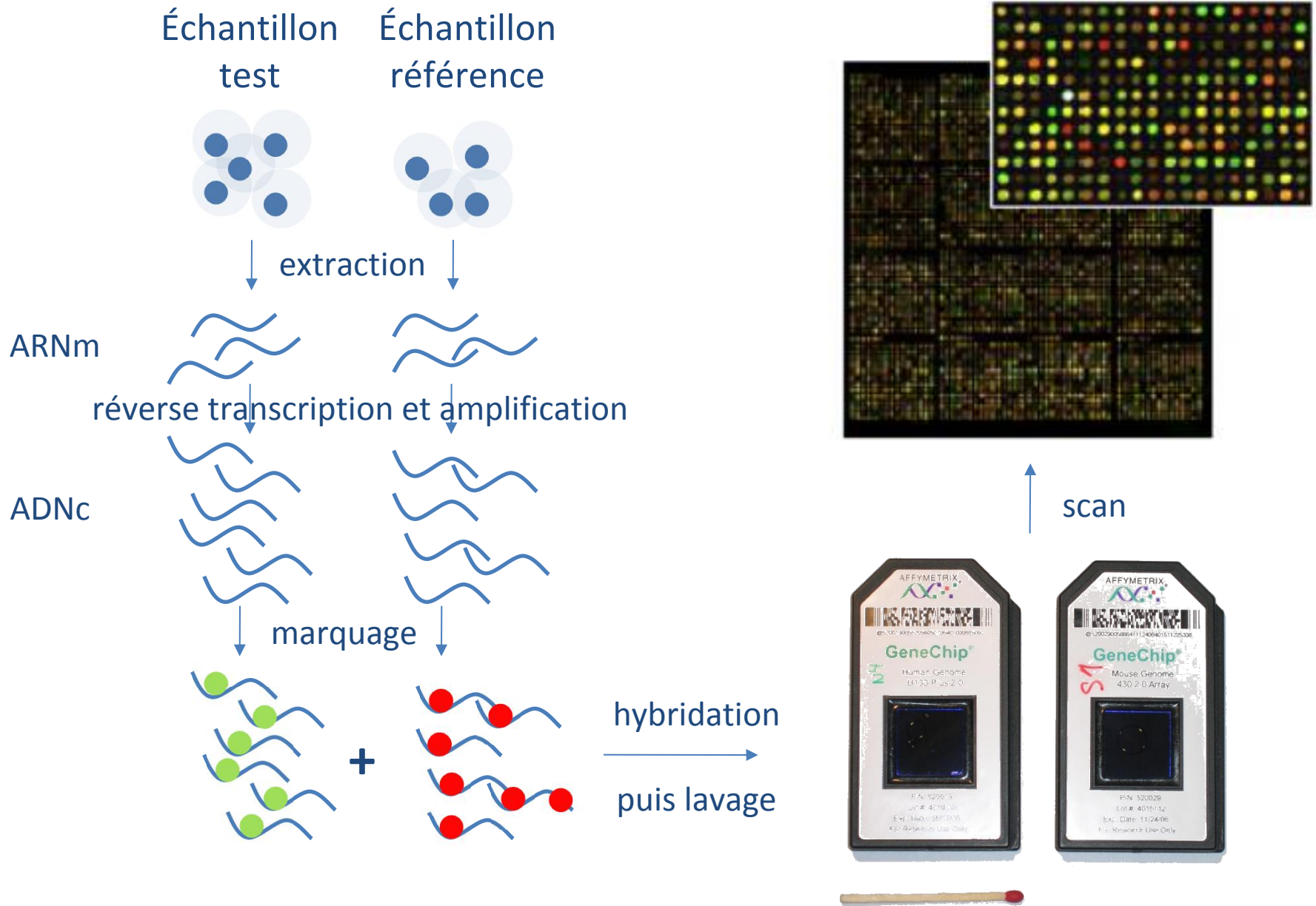
Données d'expression des gènes obtenues par :

- qPCR
- Puces à ADN
- Séquençage ultra-haut débit

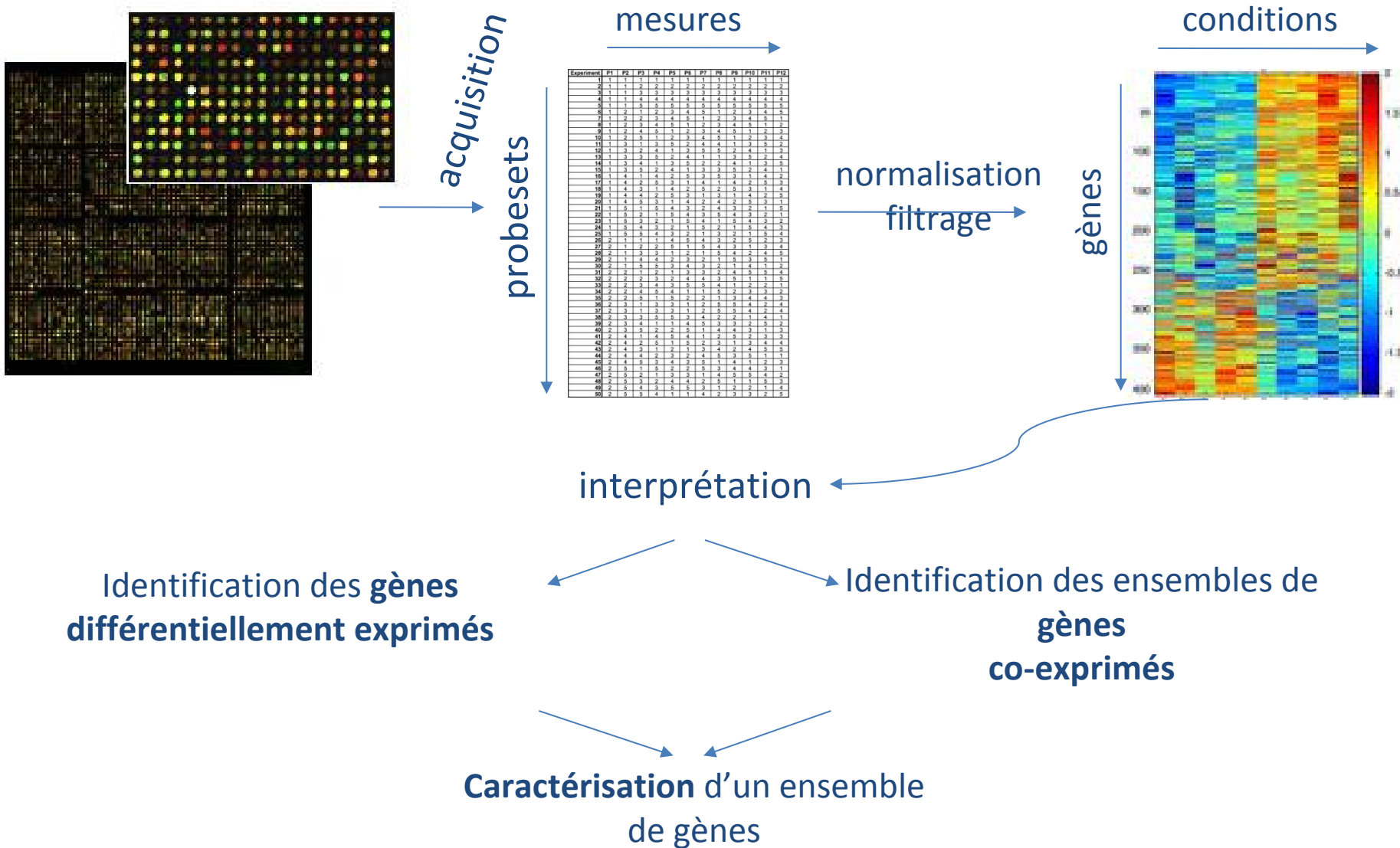
Acquisition des données (RNAseq)



Transcriptome : acquisition des données

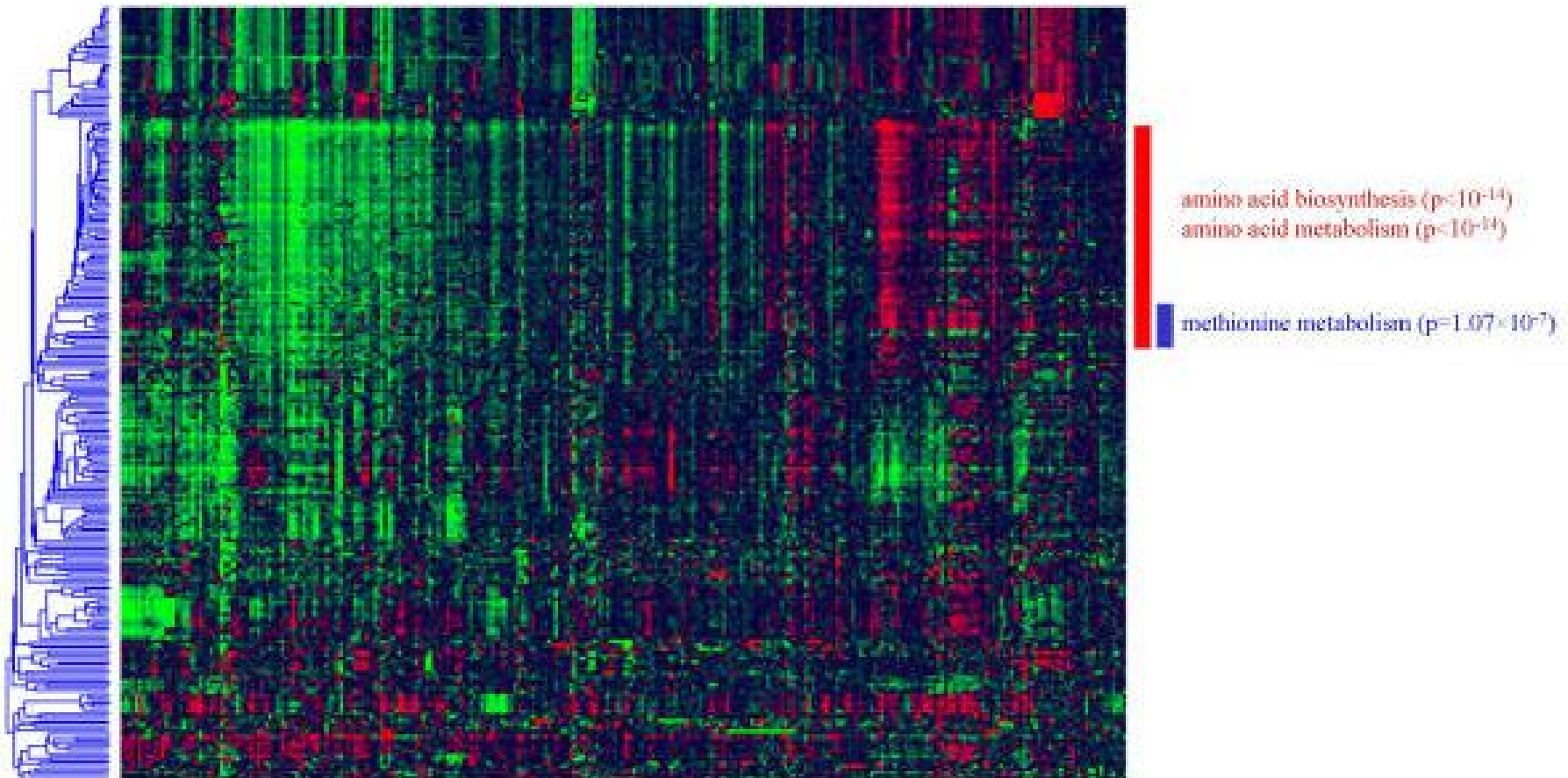


Transcriptome : analyse et interprétation des données



Transcriptome : gènes co-exprimés

- Motivation : les gènes ayant des profils d'expression similaires sont potentiellement co-régulés et participeraient donc à un même processus biologique
- But : regrouper les gènes impliqués dans un même processus biologique

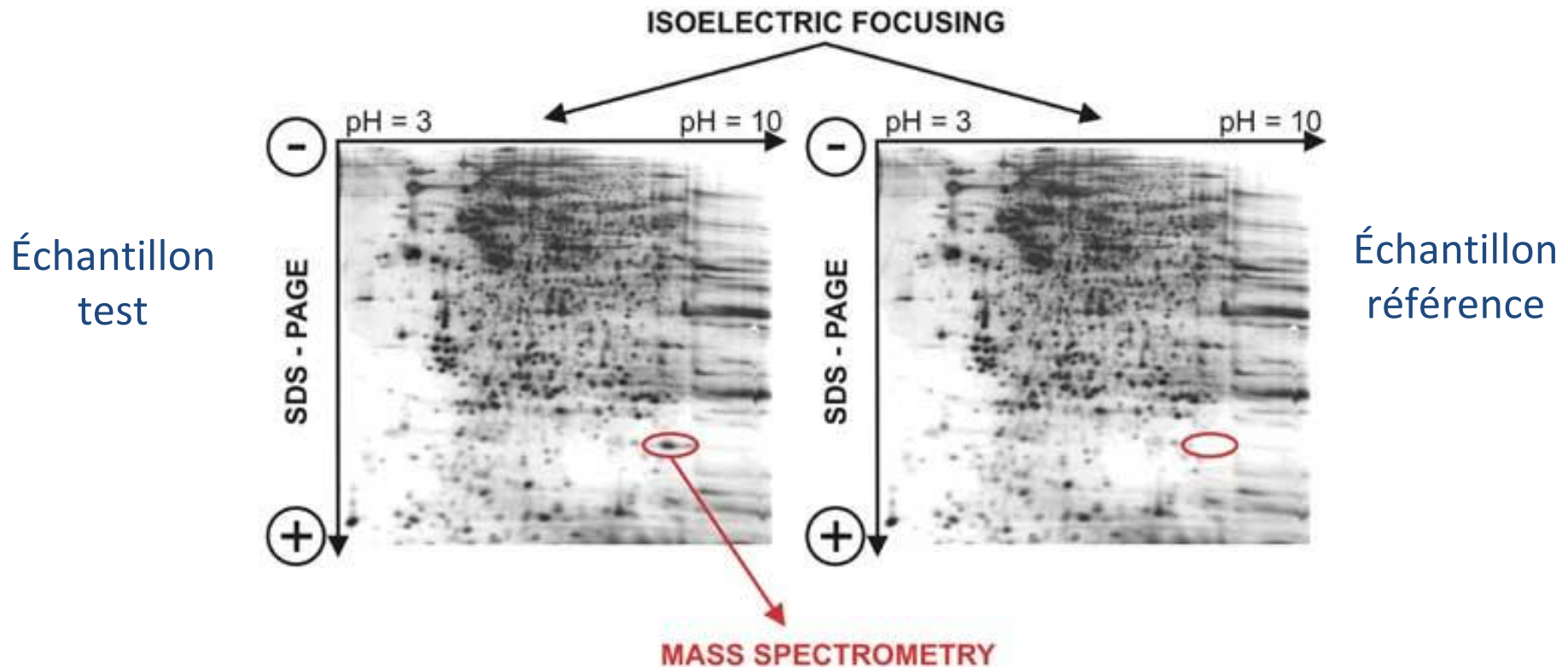


Protéomique

Protéome : ensemble des protéines exprimées dans une cellule, une partie d'une cellule (membranes, organites) ou un groupe de cellules (organe, organisme, groupe d'organismes) dans des conditions données et à un moment donné.

= *instantané* de l'état d'une cellule ou d'une population de cellules

Séparation des protéines par gels d'électrophorèse (1D, 2D) puis identification des spots par spectrométrie de masse

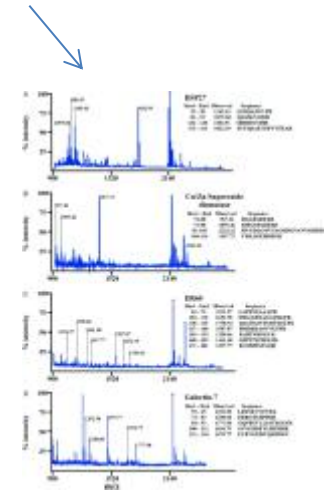
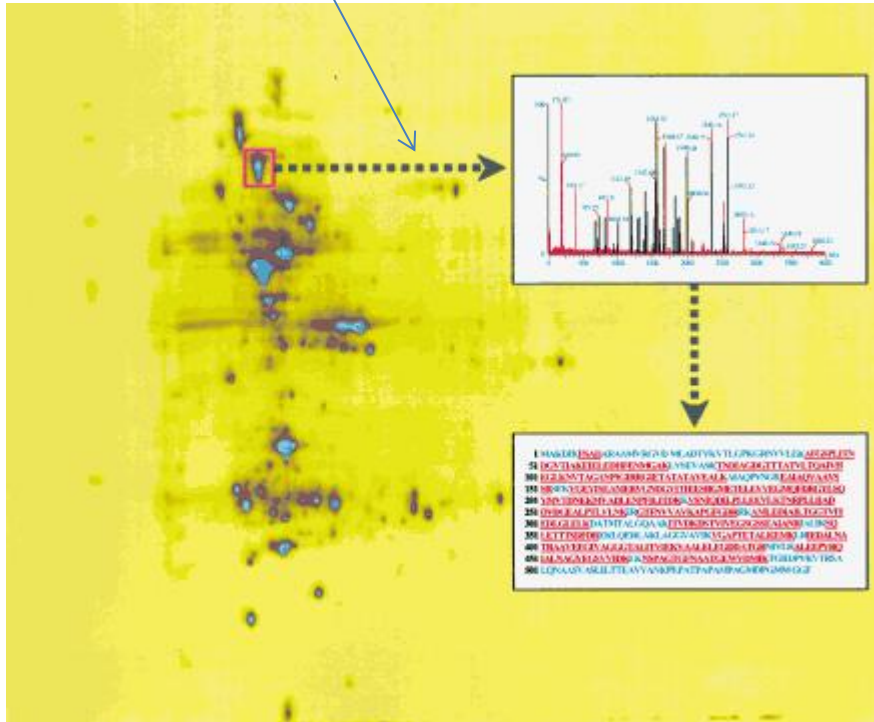


Protéomique : Identification de protéine

Digestion du spot par une enzyme (ex: trypsine) et mesure du poids des peptides obtenus

Digestion *in silico* du protéome

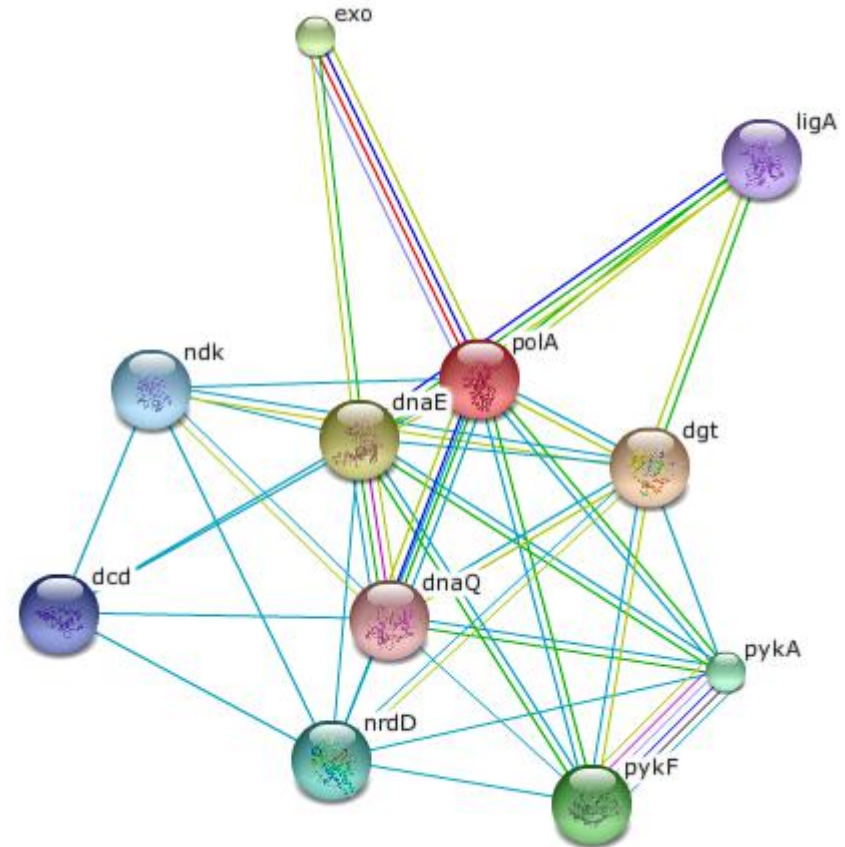
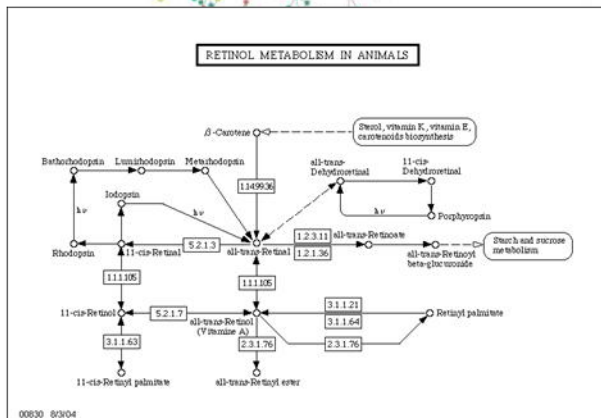
Recherche des protéines correspondant au profil observé



Réseaux de gènes et de protéines

Réseaux :

- d'interactions protéine - protéine, génétiques, fonctionnelles, ...
- de régulation des gènes
- métabolisme (enzymes – substrats)
- transduction du signal



Prédiction de structure

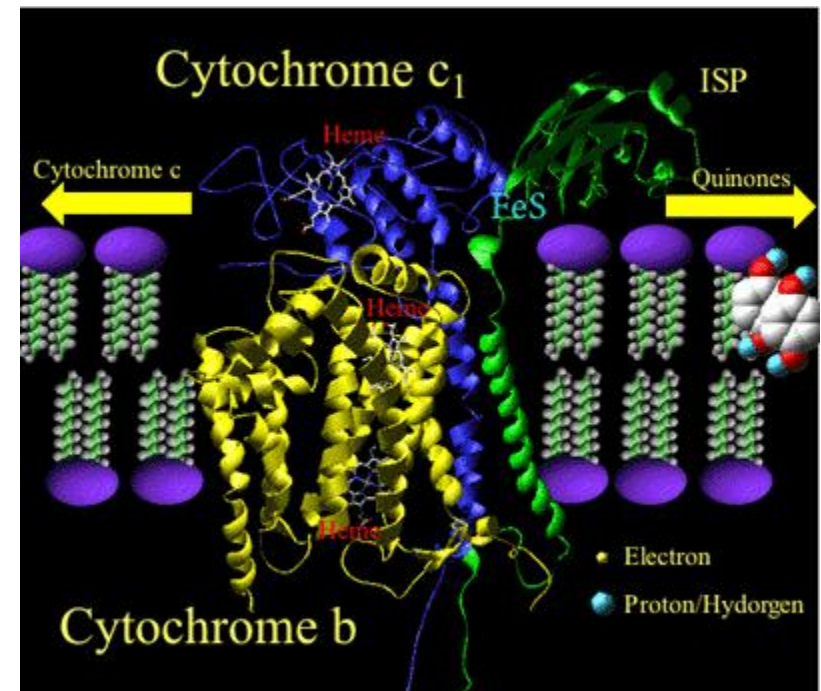
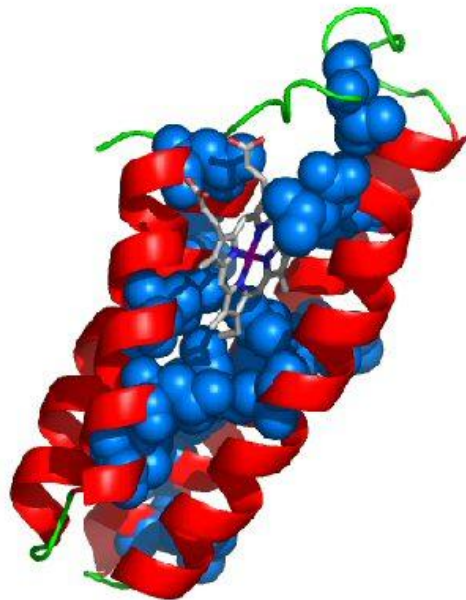
Séquence protéique

>gi|5524211|gb|AAD44166.1| cytochrome b

LCLYTHIGRNIYYGSYLYSETWNTGIMLLITMATAFMGYVLPWGQMSFWGATVITNLFSAIPYIGTNLV
 EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFHPYTIKDFLG
 LLILLLLLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGGVLAFLSIVL
 GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFPLIAGX
 IENY

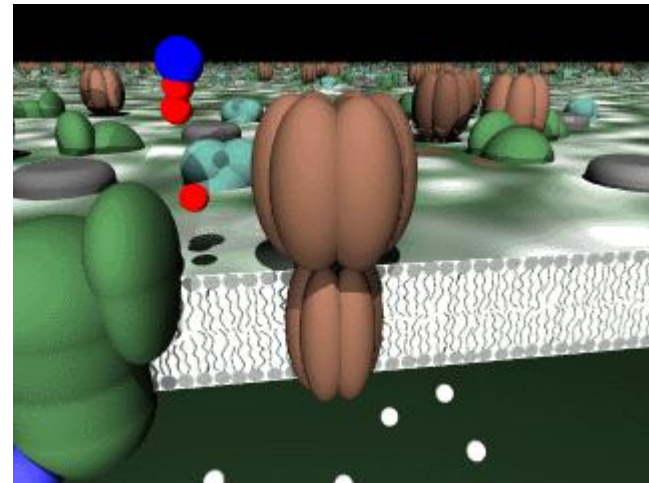
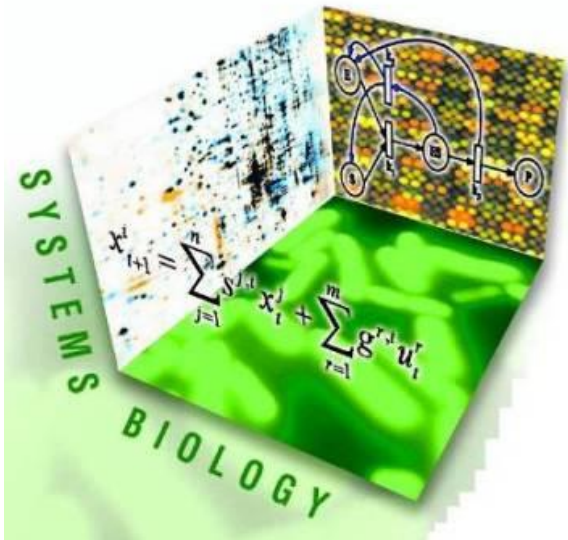


Prédiction ou résolution
de la structure tridimensionnelle

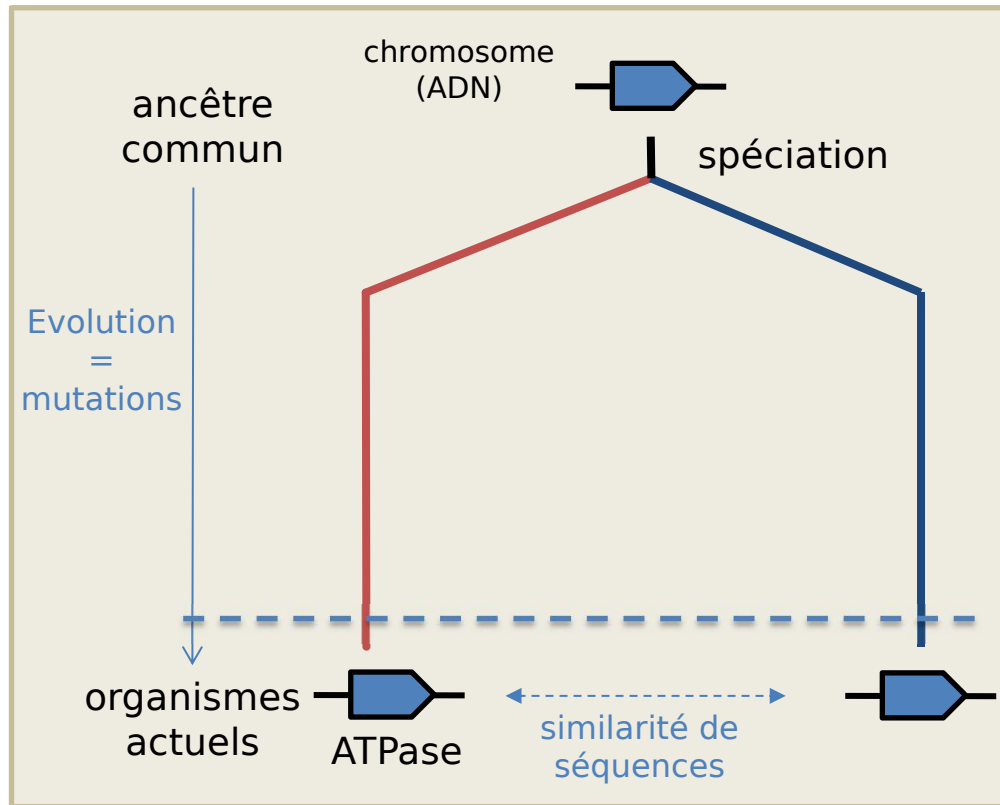


Intégration et synthèse des connaissances

- modélisation d'un système
 - processus biologique (respiration)
 - organite (mitochondrie)
 - cellule
 - population
 - écosystème



À terme : simulation d'une cellule virtuelle et prédiction de son comportement



Du temps du séquençage des premiers gènes et génomes

- la conservation/similarité des séquences impliquait des fonctions similaires
- les annotations des gènes caractérisés expérimentalement étaient transférés aux nouveaux gènes/génomes séquencés

Homologie, orthologie et orthologie 1:1

Homologie : deux gènes sont homologues s'ils ont divergé à partir d'une même séquence ancêtre

Paralogie

- définition : deux gènes sont paralogues si leur divergence a commencé après la **duplication** du gène ancêtre.
- hypothèse : mécanisme évolutif d'acquisition de nouvelles fonctions par accumulation de mutations sur une des deux "copies" du gène ancêtre non soumis à une pression de sélection ?
- remarque : les deux gènes paralogues sont au départ dans le même génome.

Orthologie

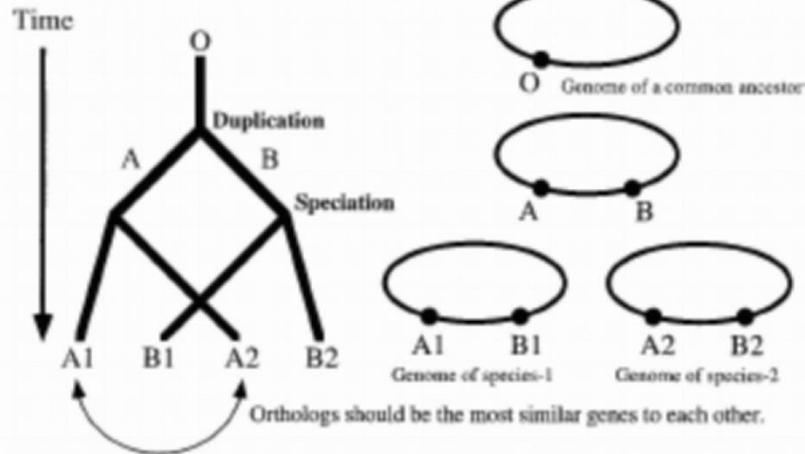
- définition : deux gènes sont orthologues si leur divergence a commencé après un événement de **spéciation** (le gène ancêtre se trouvait dans l'organisme ancêtre).
- remarque : deux gènes orthologues ne peuvent pas être présents dans le même génome.
- La relation d'orthologie est souvent abusivement utilisée pour déterminer la présence ou l'absence d'un gène dans un génome.

Orthologie 1:1

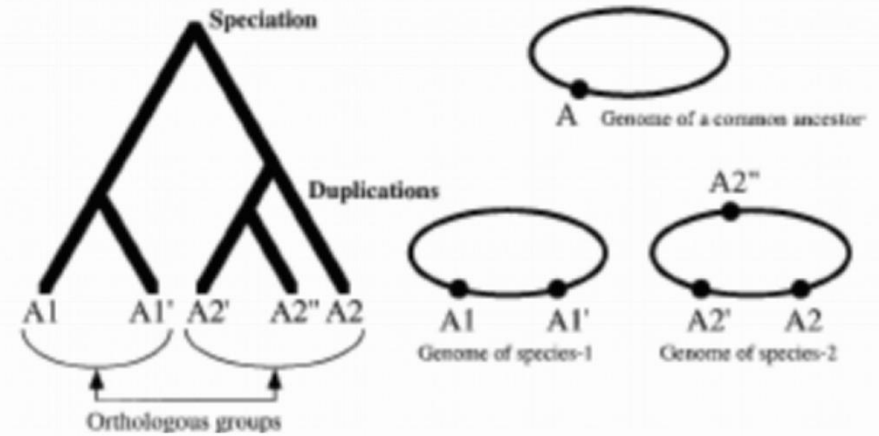
- définition : deux gènes orthologues sont orthologues 1:1 s'il n'y a pas eu de duplication (apparition de paralogue) après l'évènement de spéciation.
- remarque : la chronologie des événements de spéciation et de duplication est importante.
- Deux orthologues 1:1 ont une forte probabilité d'avoir conservé la même fonction.

Chronologie des évènements

(a) Orthologous Gene Pair

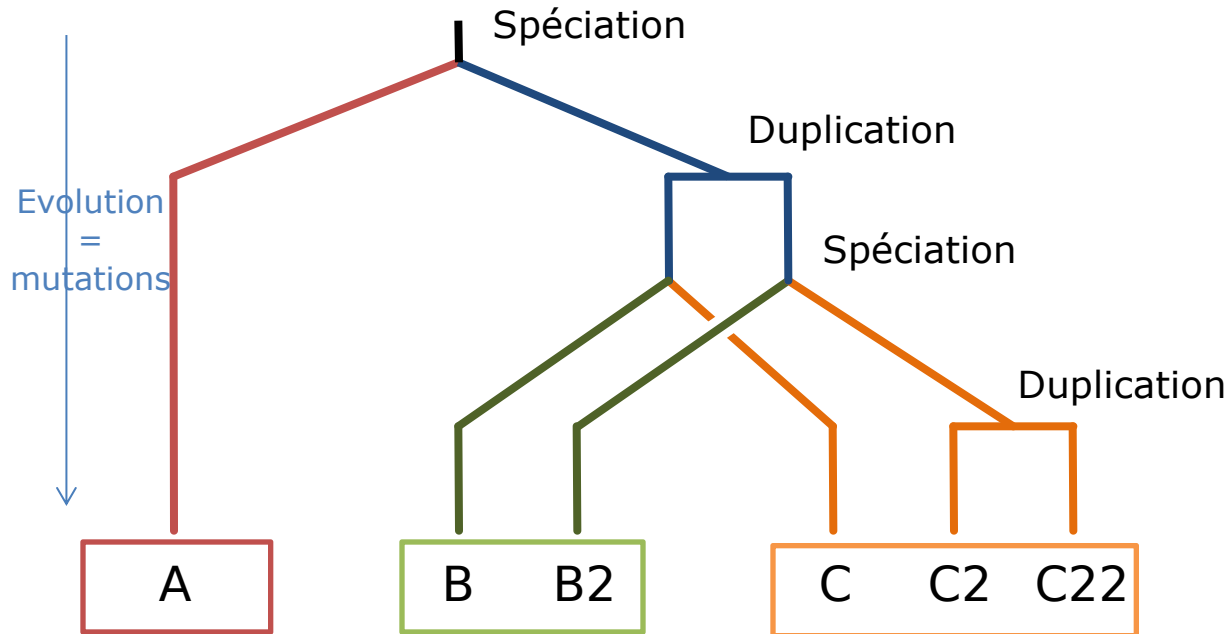


(b) Orthologous Gene Clusters

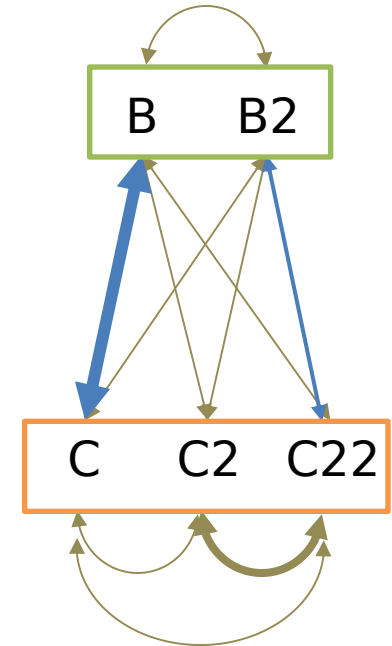


- L'orthologie 1:1 est essentielle pour la prédiction de la fonction des gènes (par transfert d'annotations)
- Problèmes :
 - la chronologie est nécessaire pour différencier paralogues, orthologues et orthologues 1:1
 - la chronologie n'est pas connue
- Deux principales approches pour identifier les orthologues 1:1
 - analyse et reconstruction de l'histoire évolutive
 - représentation sous forme de graphe des relations d'orthologie

Prédiction des orthologues 1:1

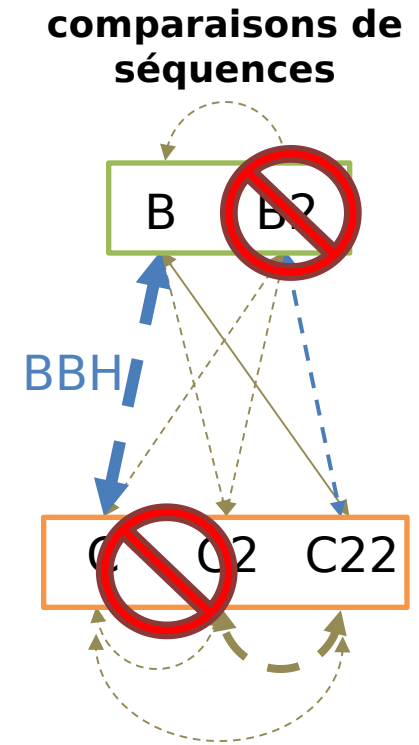
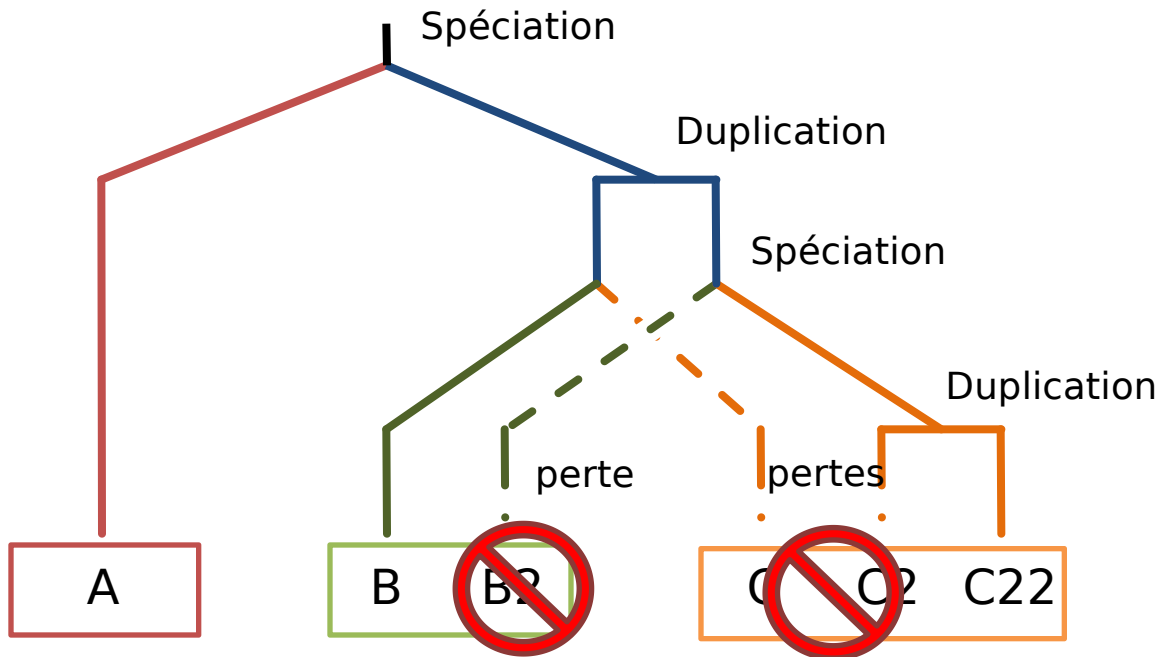


comparaisons de séquences

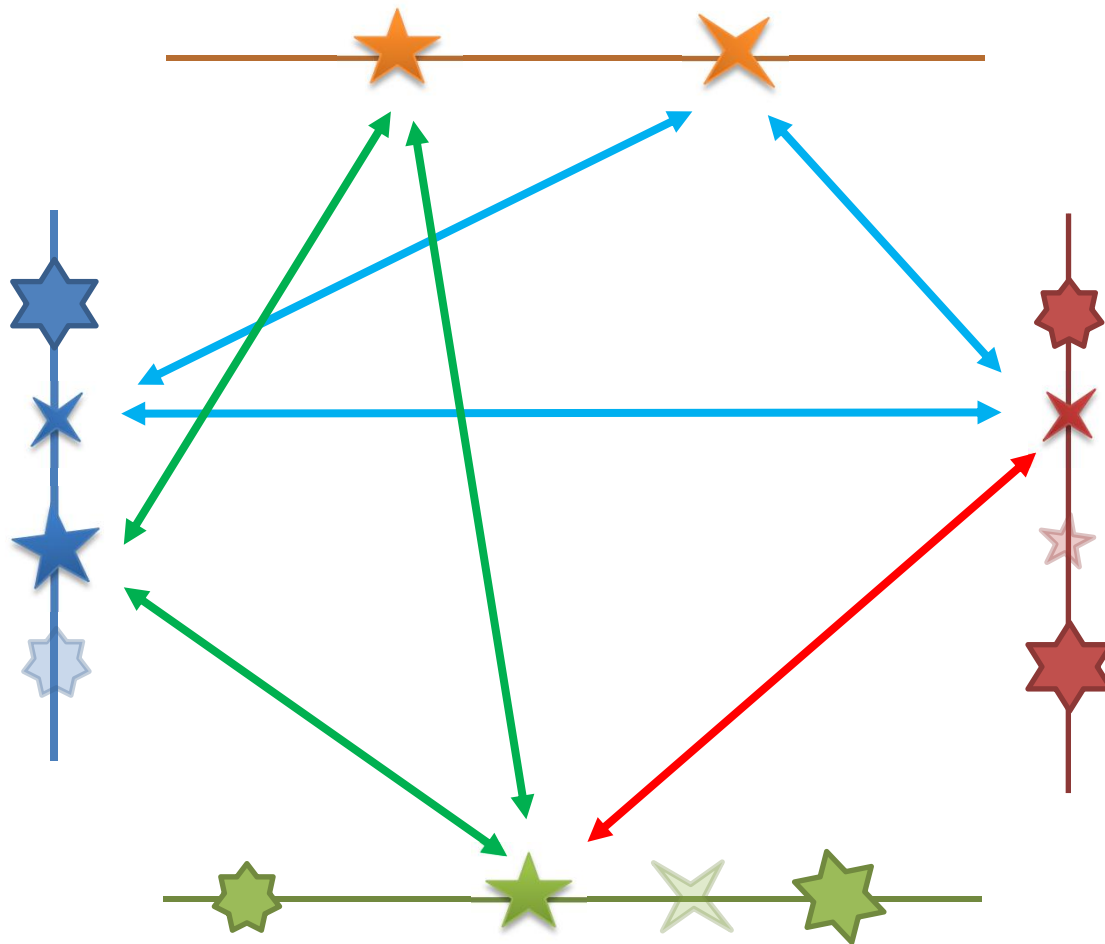


- comparaisons de toutes les séquences 2 à 2
- **B↔C prédit** $\text{sim}(B,C) > \max(\text{sim}(B,B2), \text{sim}(C, C2), \text{sim}(C,C22))$
- **B2↔C22 non prédit** $\text{sim}(B2,C22) < \text{sim}(C2,C22)$
- Remarque : calcul avec 7,223,104 séquences de 2 355 génomes =>52,000.10⁹ comparaisons

Autre évènement évolutif : perte de gènes



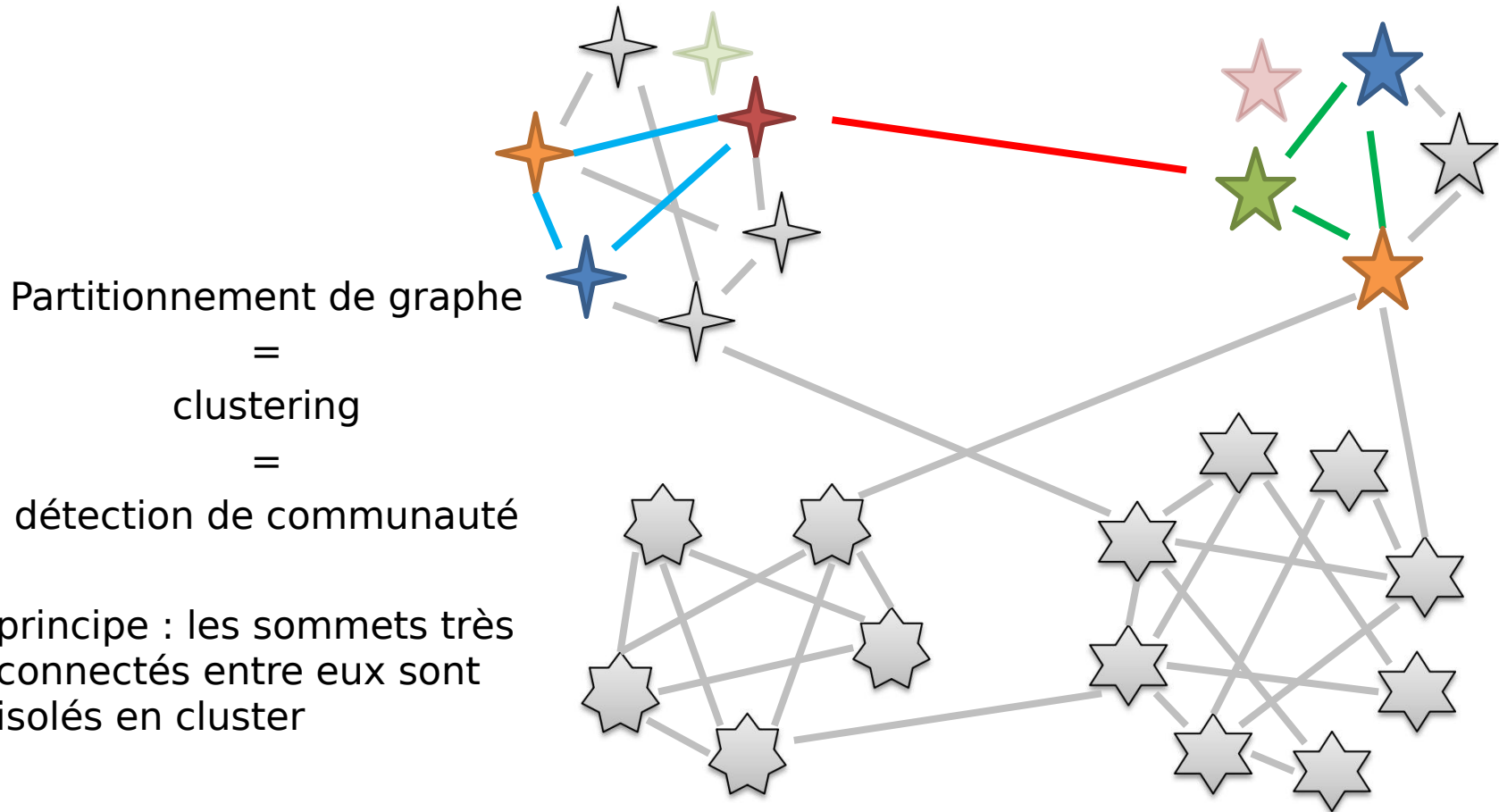
des séries de duplications et de pertes de gènes induisent des erreurs de prédiction



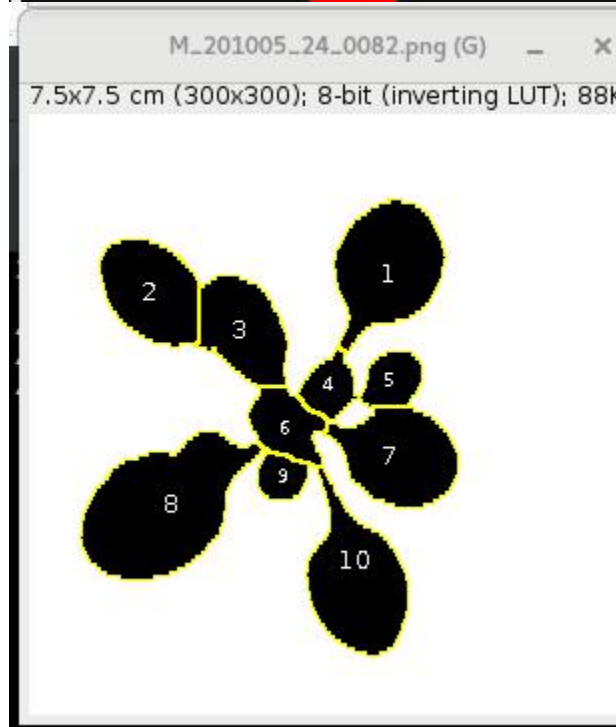
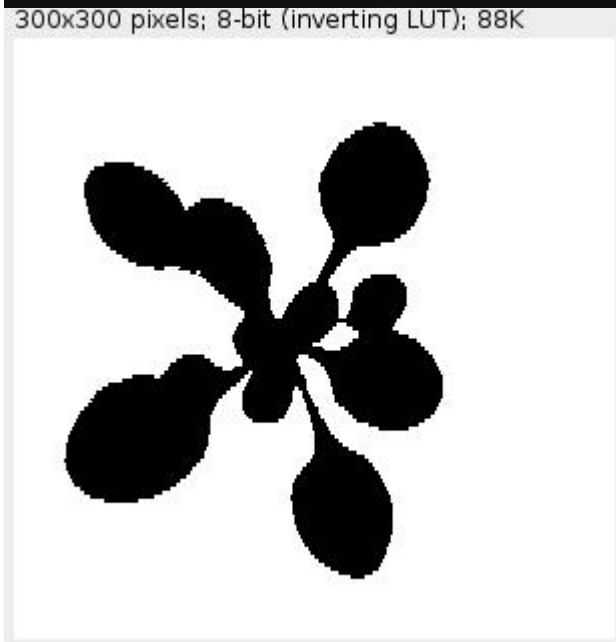
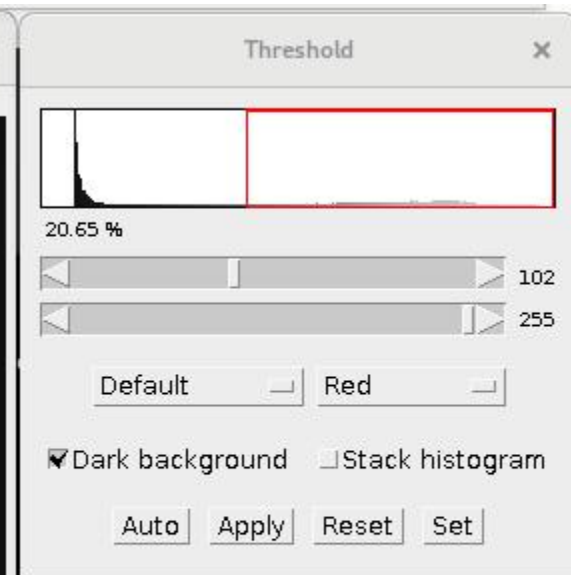
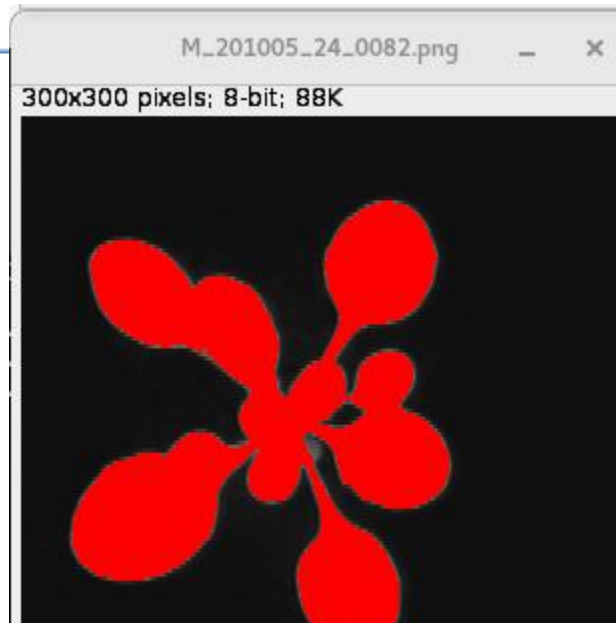
Elagage des erreurs de prédiction

Méthodes de traitement de graphes

Répresentation: sommets = gènes, liens = orthologues 1:1

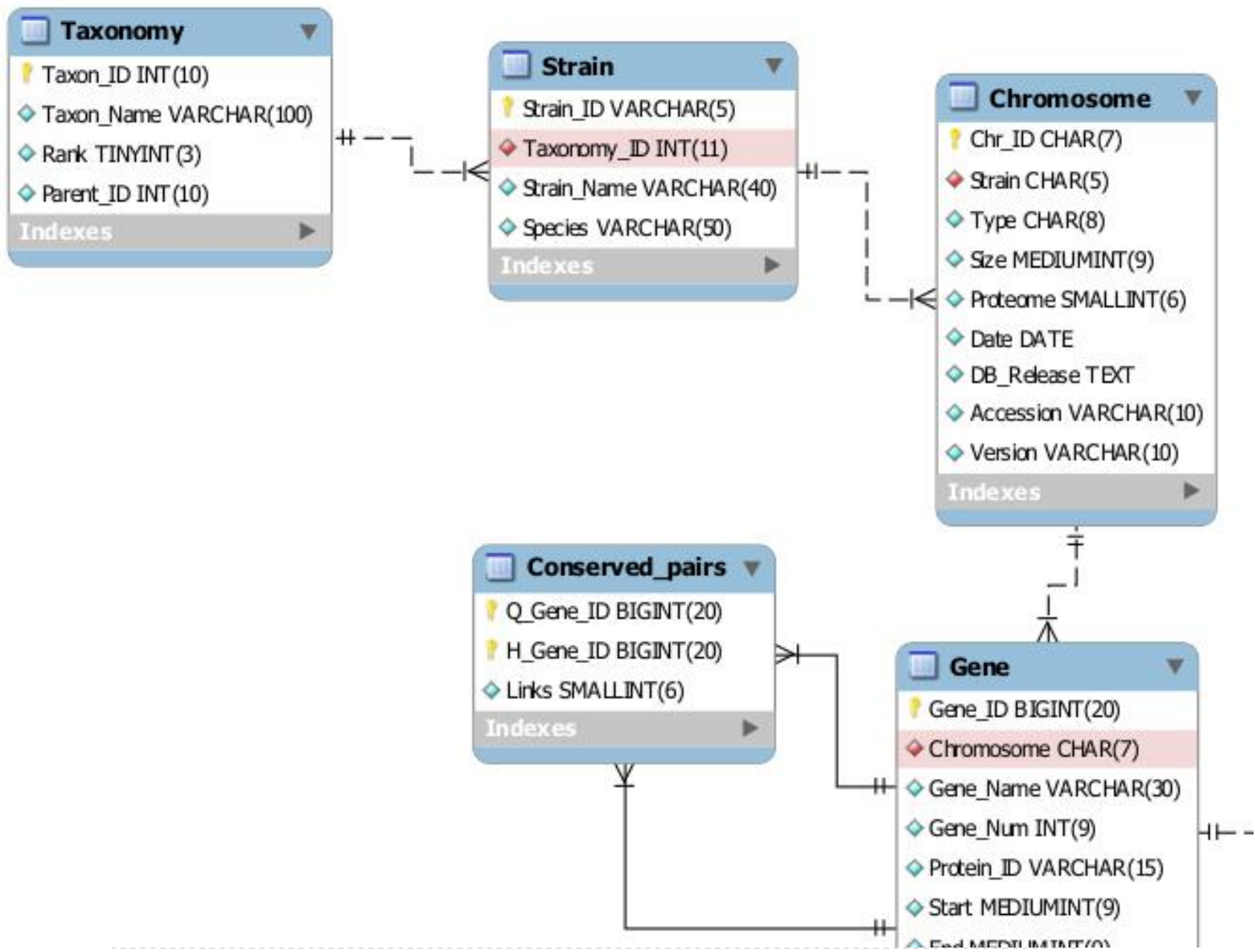



Traitements d'image



Results

	Area
1	1.719
2	1.259
3	1.126
4	0.387
5	0.415
6	0.608
7	1.437
8	2.579
9	0.289
10	1.709



 Home Help Contact About Us Subscribe Login Register

Gene Search

Search Browse Tools Portals Download Submit News ABRC Stocks

Locus: AT5G46330 [Add a Comment](#)


Representative Gene Model [AT5G46330.1](#)

Gene Model Type protein_coding

Other names: FLAGELLIN-SENSITIVE 2, FLS2, MPL12.8

Description Encodes a leucine-rich repeat serine/threonine protein kinase that is expressed ubiquitously. FLS2 is involved in MAP kinase signalling relay involved in innate immunity. Essential in the perception of flagellin, a potent elicitor of the defense response. FLS2 is directed for degradation by the bacterial ubiquitin ligase AvrPtoB. The mRNA is cell-to-cell mobile.

Map Detail Image



Chr5:18791736..18795546

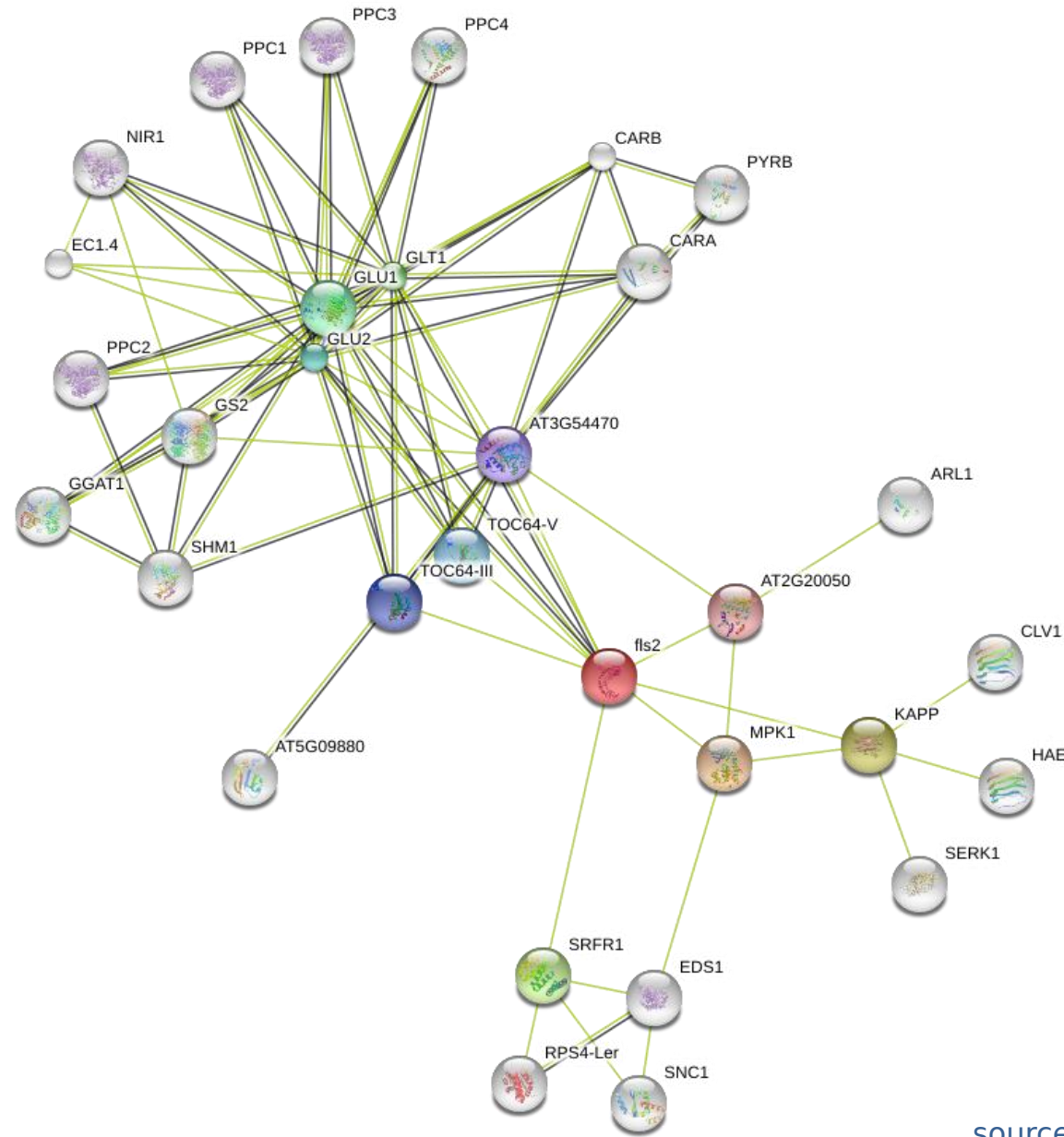
18792k 18793k 18794k 18795k

Protein Coding Gene Models
AT5G46330.1 (FLS2)

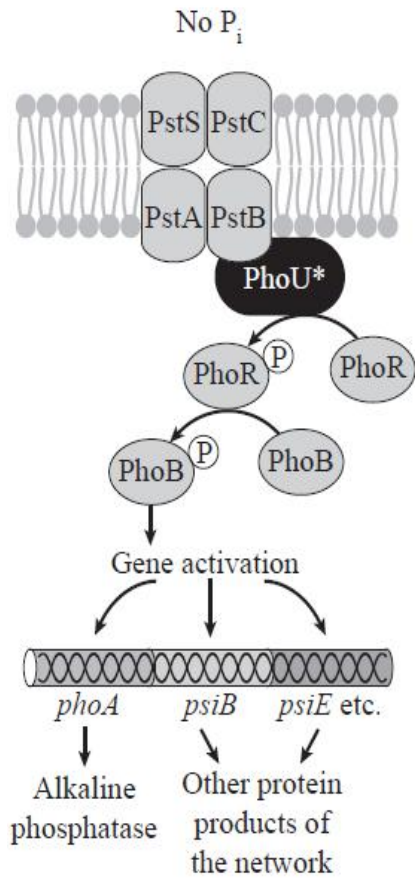
Annotations

category	relationship type	keyword
GO Biological Process	involved in	defense response by callose deposition in cell wall, defense response to bacterium, detection of bacterium, protein phosphorylation, receptor-mediated endocytosis, regulation of anion channel activity, transmembrane receptor protein tyrosine kinase signaling pathway
GO Cellular Component	located in	endosome, endosome membrane, integral component of membrane, membrane, plasma membrane
GO Molecular Function	functions in	ATP binding
GO Molecular Function	has	ATP binding, kinase activity, protein binding, protein serine/threonine kinase activity, transmembrane receptor protein serine/threonine kinase activity
Growth and Developmental Stages	expressed during	LP.02 two leaves visible stage, LP.04 four leaves visible stage, LP.06 six leaves visible stage, LP.08 eight leaves visible stage, LP.10 ten leaves visible stage, LP.12 twelve leaves visible stage, flowering stage, petal differentiation and expansion stage, plant embryo globular stage, vascular leaf senescent stage
Plant structure	expressed in	carpel, cauline leaf, collective leaf structure, cotyledon, cultured plant cell, epidermal cell, flower, guard cell, hypocotyl, inflorescence meristem, leaf apex, leaf lamina base, petiole, plant embryo, pollen, root, root tip, sepal, stamen, stem, vascular leaf
user-defined	has gene product	cell-to-cell mobile RNA

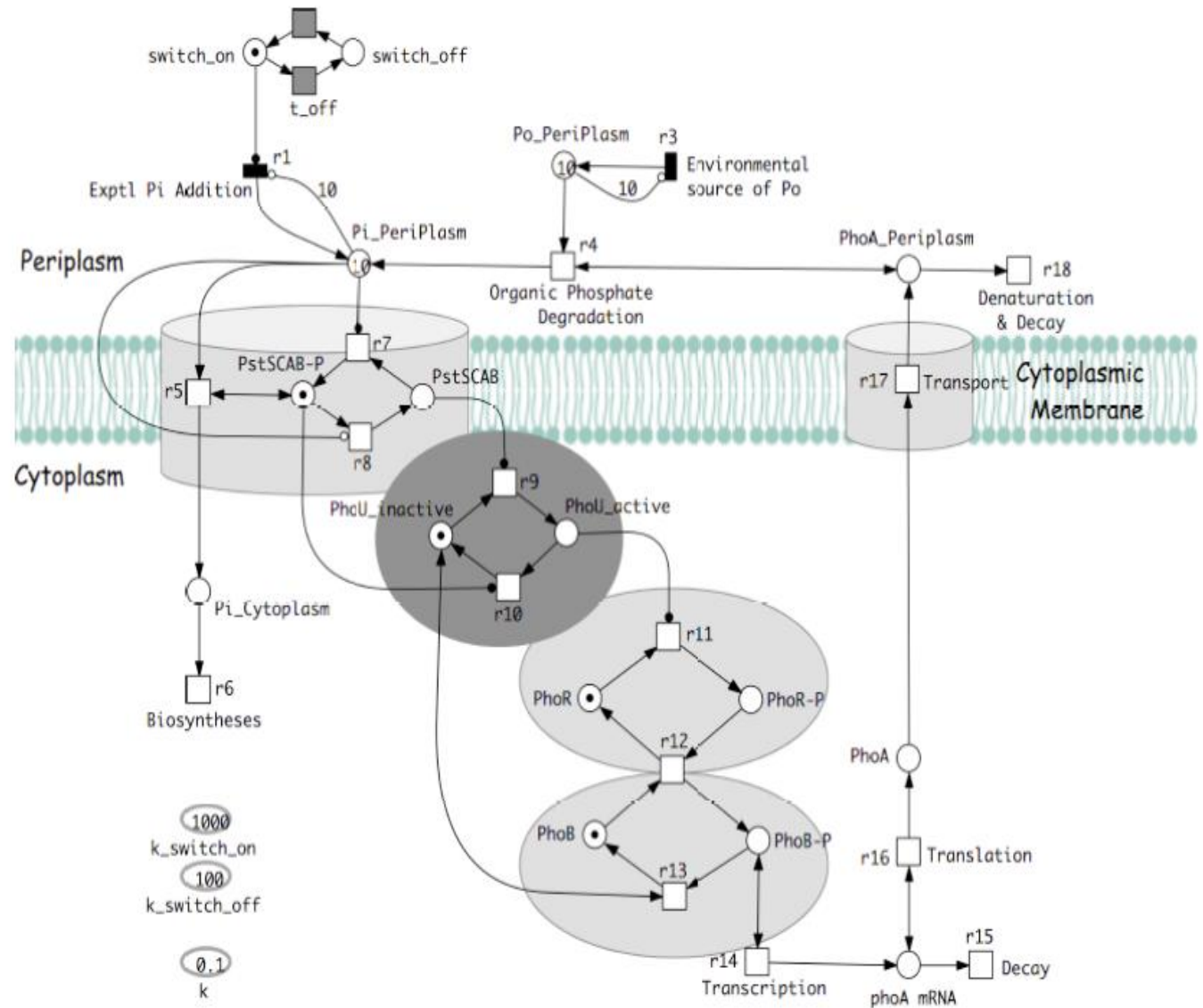
[Annotation Detail](#)



Modélisation de systèmes biologiques



Neidhardt *et al.* 1990



Durzinsky *et al.*, 2011