

Sequences analysis

Bioinformatics teachings

<http://bioinfomed.fr> - Olivier Croce -

Généralités

Bioinformatique => L'approche in silico de la biologie

Activités principales

- * Acquisition, organisation, stockage des données biologiques (ex. bases de données)
- * Utilisation ou conception de logiciels
- * Données => Analyse, comparaison ou modélisation
- * Production de nouvelles données biologiques

Quelques conseils

- * Méfiez-vous des résultats donnés par les logiciels :
 - => un logiciel modeste bien utilisé, donnera toujours de meilleurs résultats qu'un bon logiciel mal utilisé
 - => La qualité des résultats est parfois diminuée au profit de la rapidité
 - => Beaucoup de logiciels ne font que de la prédiction
- * Méfiez-vous des banques de données (séquences par exemple) :
 - => Les données se sont pas toujours fiables
 - => La mise à jour n'est pas toujours récente

La réalité mathématique n'est pas la réalité biologique

Les ordinateurs/logiciels sont des outils, qu'il faut apprendre à bien utiliser

Quelques liens utiles en bioinformatique



* La Société Française de BioInformatique (SFBI)
<http://sfbi.impg.prd.fr/>



* Logiciels pour la biologie de l'Institut Pasteur
<http://bioweb.pasteur.fr/>



* Le Pôle Bioinformatique Lyonnais (PBIL)
<http://pbil.univ-lyon1.fr/pbil.html> <http://npsa-pbil.ibcp.fr/>



* European Bioinformatics Institute (EBI)
<http://www.ebi.ac.uk/>

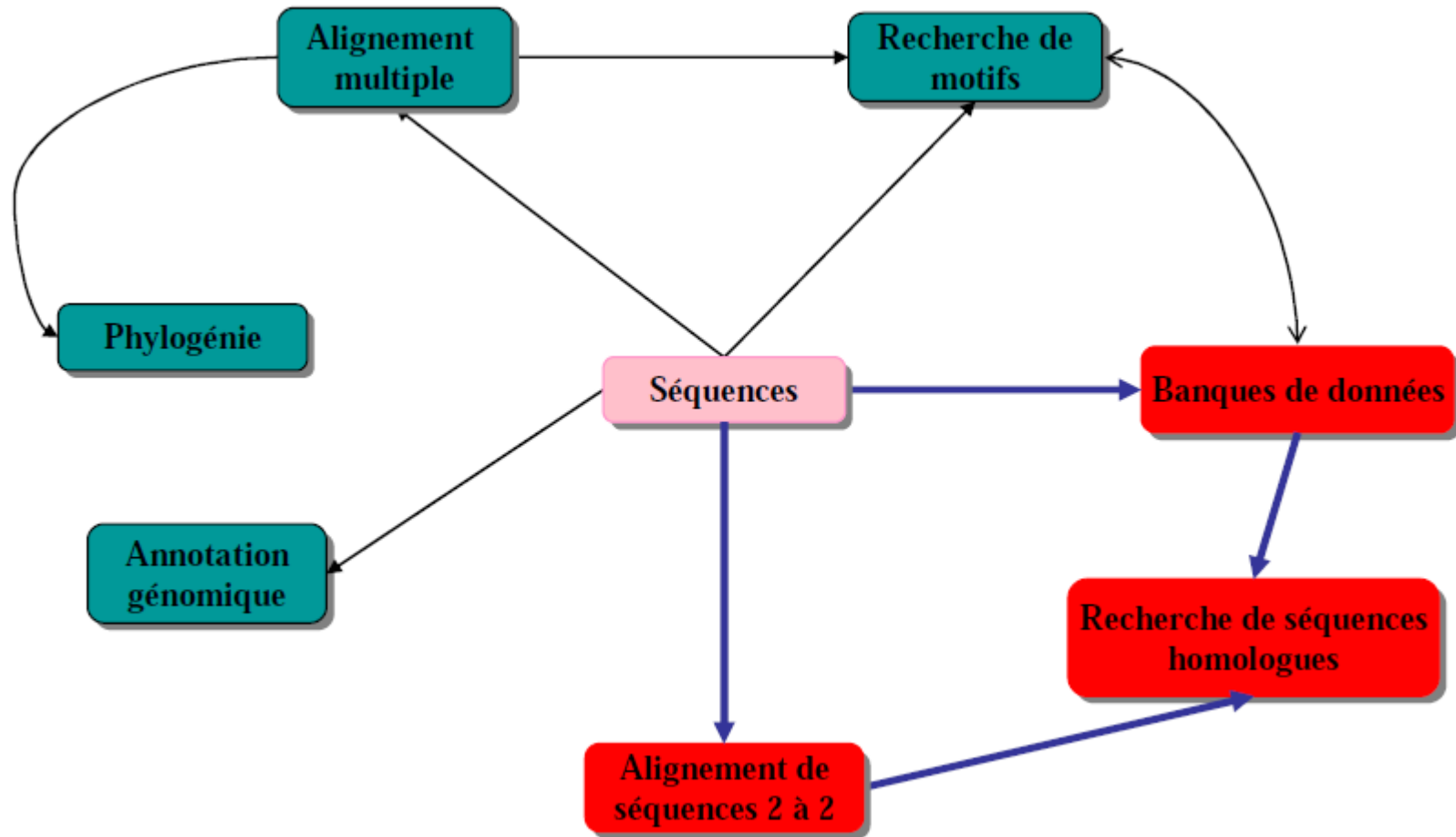


* Les outils de protéomique d'ExpASy
<http://www.expasy.org/tools/>



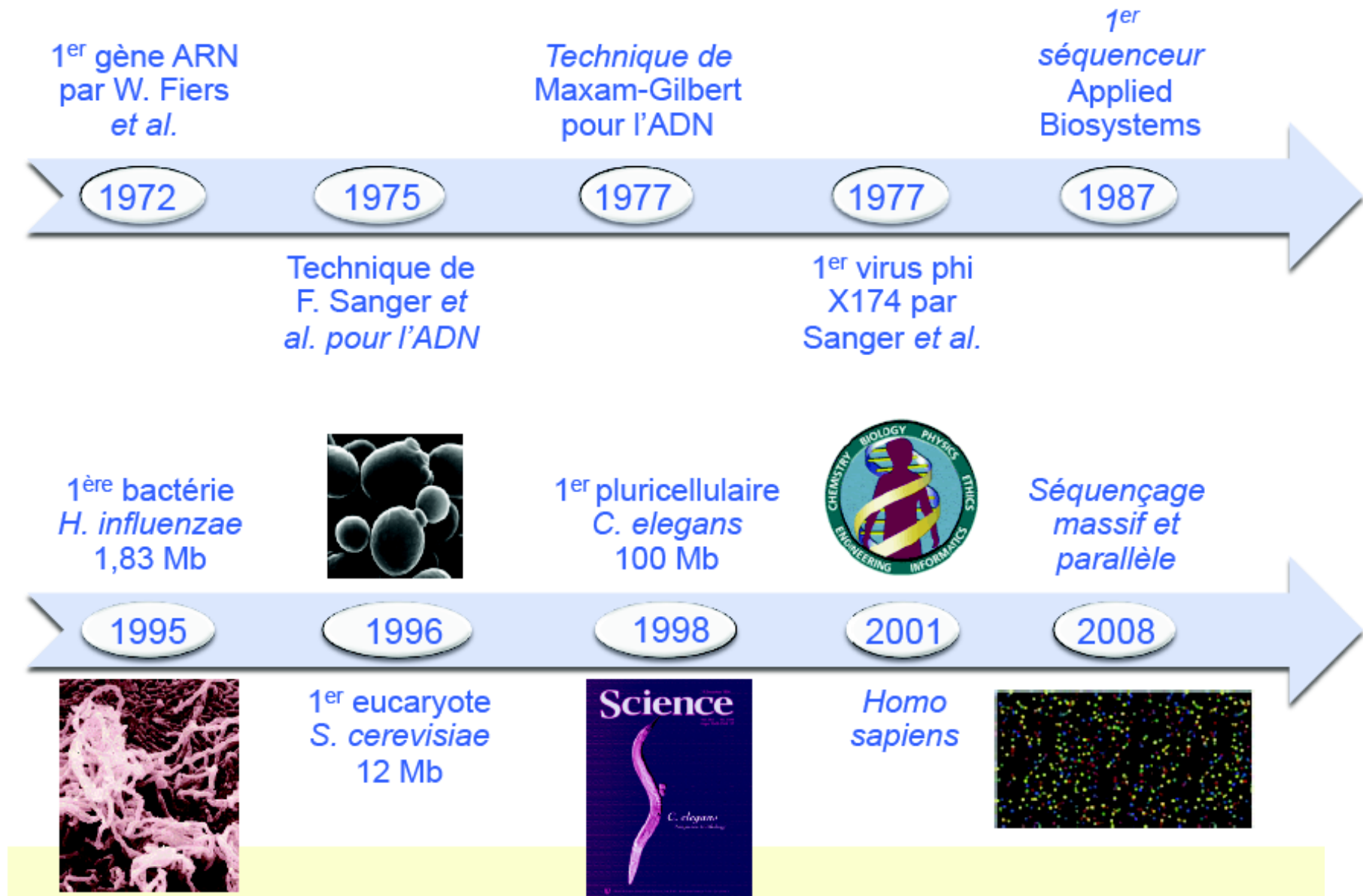
* National Center for Biotechnology Information (NCBI)
<http://www.ncbi.nlm.nih.gov/>

Analyse de séquences, pourquoi ?



Bases de données de séquences nucléotidiques

- => Origine des données : Séquençage de molécules d'ADN ou d'ARN
- => Nécessité de stocker beaucoup de données et de manière organisée.



Bases de données de séquences nucléotidiques

Banques nucléiques, collaboration => International Nucleotide Sequence Database Collaboration

* Association des 3 banques nucléiques :

EMBL-Bank (European Molecular Biology Laboratory) – EBI
<http://www.ebi.ac.uk/embl/>

GenBank (banque des Etats-Unis d'Amérique) – NCBI
<http://www.ncbi.nlm.nih.gov/Genbank/>

DDBJ (DNA DataBank of Japan) – CIB
<http://www.ddbj.nig.ac.jp/>

- Echange quotidien des données
- Répartition de la collecte des données
- Chaque banque (est supposée) collecter les données de son continent

* Les données stockées :

- 1 séquence + ses annotations = 1 entrée
- Fragments de génomes : un ou plusieurs gènes, un bout de gène, séquence intergénique
- Génomes complets
- ARNm, ARNt, ARNr, ... (fragments ou entiers)

Bases de données de séquences nucléotidiques

Banques nucléiques, format d'une entrée :

- 3 parties :

Description générale
de la séquence

« Features »

Description des objets
biologiques présents
sur la séquence

La séquence

```
ctccggcagc ccgaggtcat cctgctagac tcagacctgg atgaacctat agacttgcgc      60
tcggtcaaga gccgcagcga ggccggggag ccgccagct ccctccaggt gaagcccagag      120
acaccggcgt cggcggcggg ggcggtggcg gcggcagcgg caccaccac gacggcggag      180
```

- Chaque ligne commence par un mot-clé
 - Deux lettres pour EMBL
 - Maximum 12 lettres pour Genbank et DDBJ
- Fin d'une entrée : //

EMBL, description générale de la séquence

- ID : toujours la 1ère ligne d'une entrée

Accession	Version	Topologie	Molécule	Classe	Taxonomie	Taille seq
M71283	SV 1	linear	genomic DNA	STD	BCT	1322 BP

- AC : numéros d'accension

- Un n°acc principal pour chaque entrée, unique

- Une liste de n°acc secondaires (historique de l'entrée)

- DT : dates de création et de dernière version

- DE : description du contenu de l'entrée

- KW : mots-clés ; peu renseigné

- OS, OC : organisme contenant la séq. et sa taxonomie

- RN, RC, RX, RP, RA, RT, RL : réf. bibliographiques

- Uniquement les références données par les auteurs de l'entrée

GenBank et DDBJ, description générale

- LOCUS : toujours la première ligne d'une entrée

Locus name	Taille seq	Molécule	Topologie	Division	Date
BACCOMQP	1322 bp	DNA	linear	BCT	26-APR-1993

- DEFINITION = DE
- ACCESSION = AC
- VERSION ~ DT
- KEYWORDS = KW
- SOURCE, ORGANISM = OS, OC
- REFERENCE, AUTHORS, TITLE, JOURNAL, ... = R...

Banques nucléiques, lignes FT (Features)

Format (partagé par toutes les banques) :

- **Key** : un seul mot indiquant un groupe fonctionnel
 - Vocabulaire contrôlé, hiérarchique
 - gene : séquence complète du gène (y compris les introns)
 - CDS : séquence codante (sans les introns, entre ATG et Stop)
- **Location** : instructions pour trouver l'objet sur la séquence de l'entrée

- **Qualifiers** : description précise du groupe fonctionnel
 - Format : /qualifier="commentaires libres"
 - /gene="comQ" : nom du gène concerné
 - /note="competence regulation" : information concernant la fonction

Banques nucléiques, localisation des « keys »

467 : l'annotation ne concerne qu'une seule base

109..1105 : entre les positions 109 et 1105 (inclues)
Toujours la position la plus petite en premier

<1..21 ou 1275..>1322 : « Keys » tronqués

Commence avant le premier nt de l'entrée

Se termine après le dernier nt de l'entrée (taille seq = 1322)

<234..888 : début réel inconnu, mais avant 234

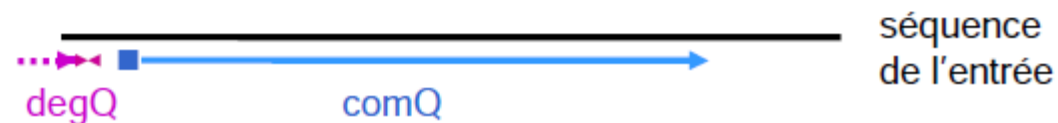
234..>888 : fin réelle inconnue, mais après 888

complement(340..565) : séquence complémentaire inversée à celle de l'entrée (brin -)

join(12..78,134..202) : fragments indiqués mis bout à bout (concaténés) ; nombre de fragments illimité

Exemple de « Feature » d'une séquence ADN

```
FT  CDS                <1..21
FT                      /codon_start=1
FT                      /db_xref="SWISS-PROT:Q99039"
FT                      /transl_table=11
FT                      /gene="degQ"
FT                      /protein_id="AAA22322.1"
FT                      /translation="YAMKIS"
FT  terminator         21..47
FT                      /gene="degQ"
FT  promoter           109..140
FT                      /gene="comQ"
FT  mRNA               146..1105
FT                      /partial
FT                      /gene="comQ"
```



Le format FASTA

- Utilisé par les logiciels d'analyse de séquences
- Fichier texte simple, créé avec n'importe quel éditeur de texte (ex. blocnote, textpad, etc.) : toto.fas, toto.fna, toto.fa, etc..
- Une ligne débutant avec « > » + un identifiant unique (obligatoire) + des commentaires (optionel)

La séquence brute (pas d'espace, ni de nombre), sur une ou plusieurs lignes

```
>AC25300 Human Polycomb 2 homolog (hPc2) mRNA, partial cds
ctccggcagcccgagggtcatcctgctagactcagacctggatgaaccac
Ctccggcagcccgagggtcatcctgctagactcagacctggatgaaccat
agacttgcgctcgggtcaagagccgcagcagggccggggagccgcccagct
ccctccaggtgaagcccgagacaccggcgtcggcggcgggtggcgggtggcg
Gcggcagcggcaccaccacgacggcggagaagcct
>ID23003 hPc2 gene
ggacgaacctgcagagtcgctgagcgcgagttcaagcccttctttgggaata
taattatcaccgacgtcaccgcgaactgcctcaccgttactttcaaggag
tacgtgacggtg
```

Banques nucléiques, inconvénients

- * Difficulté de mise à jour des données

Version plus récente d'une séquence ou d'une annotation dans d'autres banques (ex : banques dédiées à un génome complet)

- * Forte redondance

Un même fragment de séquence présent dans plusieurs entrées

- * Annotations peu normalisées

Difficulté de recherche d'une information précise

- * Annotations peu précises

Peu de descriptions sur les gènes et leurs produits

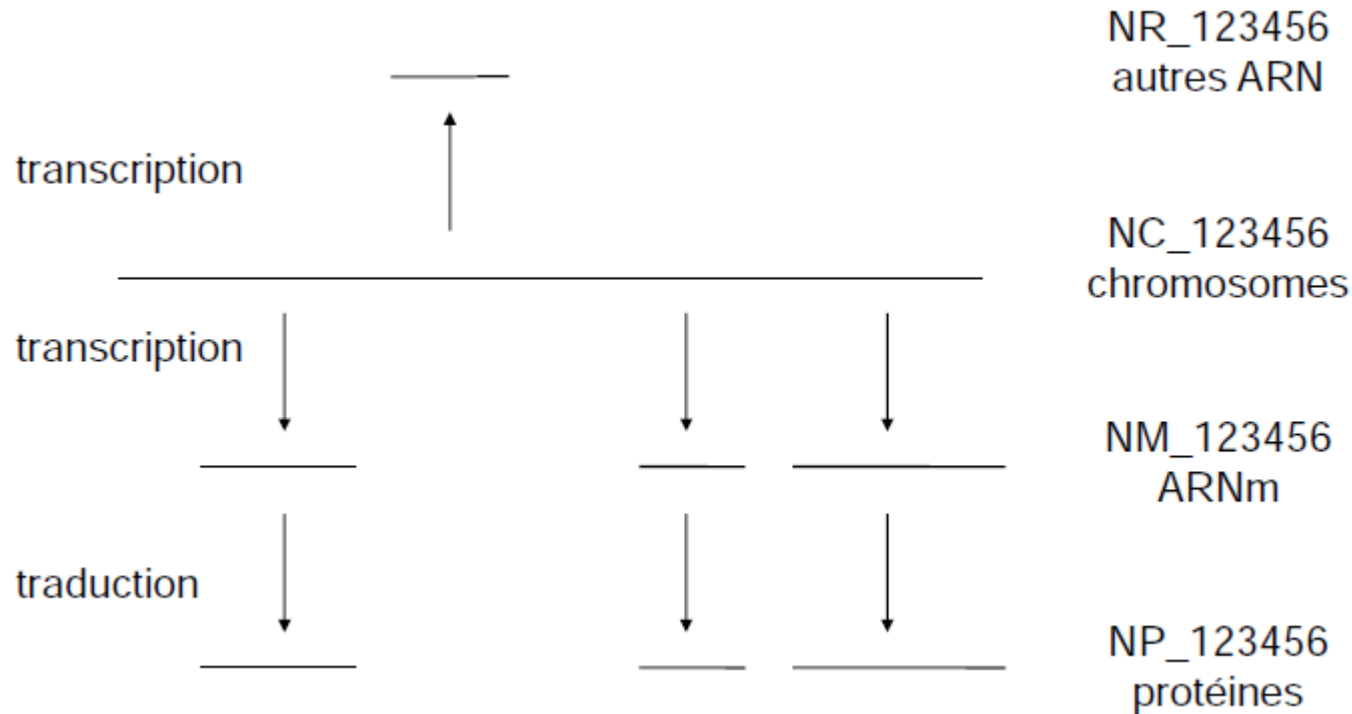
- * Erreurs dans les annotations et dans les séquences

Bases de données de séquences nucléotidiques

Bases connues : **nr / nt** et **RefSeq** (=Reference Sequence collection)

=> bases de séquences non-redondantes.

=> **nr/nt** plus complet que **RefSeq**, RefSeq + contrôlée et avec liens explicites entre les séquences nucléiques et protéiques.



Bases de données de séquences proteiques

=> Origine des données

- Traduction automatique des séquences d'ADN ou d'ARNm
- Séquençage de protéines (rare car + long et coûteux)

=> Les données stockées : séquences + annotations = une entrée

- Protéines entières
- Fragments de protéines

1956 : F. Sanger établit la séquence en aa de l'insuline

1965 : Atlas of Protein Sequences, M. Dayhoff

Version papier jusqu'en 78, puis version électronique

1984 : création de PIR-NBRF (Protein Information Resource - National Biomedical Research Foundation) Collaboration avec MIPS (Allemagne) et JIPID (Japon)

1986 : création de SwissProt

Collaboration entre SIB (Swiss Institut of Bioinformatics) et EBI

Fin 2003 : UniProt (Universal Protein Resource)

Mise en commun des informations de PIR et SwissProt/TrEMBL

<http://www.expasy.uniprot.org/>

Bases de données de séquences proteiques

UniProt est en 2 parties

SwissProt

- Données corrigées et validées par des experts
- Haut niveau d'annotations
 - * Description de la fonction (références associées)
 - * Localisation des domaines fonctionnels
 - * Modifications post-traductionnelles
 - * Existence de variants, ...
- Redondance minimale
- Nombreux liens vers d'autres banques (60 BD)

TrEMBL

- Entrées supplémentaires à SwissProt (pas encore annotées)
- Traduction automatique de l'EMBL

Bases de données de séquences proteiques

SwissProt/TrEMBL, format d'une entrée : format basé sur celui de l'EMBL

- Mot-clé de 2 lettres au début de chaque ligne
- Format différent pour les Features
- Mots-clés supplémentaires :
 - GN : les différents noms du gène qui code pour la protéine (OR) et les différents gènes qui codent pour la même protéine (AND)
 - OX : références croisées vers les banques taxonomiques
 - CC : commentaires, lignes très documentées dans SwissProt
 - KW : mots-clés issus d'un dictionnaire
 - DR : références vers d'autres banques de données
 - Vers les séquences nucléiques (EMBL/GenBank/DDBJ)
 - Vers les structures 3D
 - ...

Informations découpées en blocs pour plus de lisibilité

- CC -!- TOPIC: First line of a comment block;
- CC second and subsequent lines of a comment block.

De nombreux sujets sont abordés

- FUNCTION : description générale de la fonction de la protéine
- CATALYTIC ACTIVITY : description des réactions catalysées par les enzymes
- DEVELOPMENTAL STAGE : description du stade spécifique auquel la protéine est exprimée
- SUBUNIT : complexes dont fait partie la protéine (+ partenaires)
- ...

Pourquoi interroger une banque ?

- Obtenir des informations nouvelles et pertinentes
- Aide à la mise au point d'expériences
- Validation des résultats d'une expérience

- Trouver si des séquences sont déposées dans les banques.

- Identifier des protéines homologues (avec ancêtre commun):
 - orthologue : organismes différents
 - paralogue : organismes identiques

- Déterminer si des séquences ont une fonction similaire ou proche.

- Déterminer des familles de protéines ayant un domaine conservé.

- Localiser des régions codantes et non codantes (aligner des séquences génomiques ADN et des séquences exprimées (cDNAs, ESTs)).

- Etablir des relations entre les séquences.

=> Comment interroger une banque ?

1 - Interrogation par annotations

Interrogation des annotations d'une banque => Recherche dans les annotations

Systèmes d'interrogation de banques de données : Entrez, SRS, Acnuc

=> Recherche de mots ou expressions dans le texte des entrées

- * Obtention de données pertinentes (pas trop de résultats, mais tous ceux relatifs à notre problématique)

- * Simplicité d'utilisation (syntaxe d'interrogation intuitive)

- * Réponse rapide

- * Possibilité d'analyse des résultats (couplage à d'autres outils)

=> Systèmes d'interrogation de banques de données :

- * **Entrez**, permet aussi d'interroger des banques de séquences (<http://www.ncbi.nlm.nih.gov>)

Même fonctionnement que pour interroger PubMed

- * **SRS** : un autre système d'interrogation, + élaboré, mais plus complexe (<http://www.dkfz.de/srs/> ou liste à : <http://bioblog.instem.com/download/srs-parser-and-software-downloads/public-srs-installations/>)

- * **ACNUC** (http://doua.prabi.fr/search/query_fam)

1 - Interrogation par annotations

Exemple de serveur SRS : <http://www.dkfz.de/srs/>

Pour une recherche simple

The screenshot displays the SRS web interface. At the top, there is a navigation bar with tabs: Quick Search, Library Page, Query Form, Tools, Results, Projects, Views, Databanks, and a HELP button. The main content area is divided into several sections:

- Temporary Project:** Shows a project ID '3pXta1Ogl9x'.
- Tips:** A sidebar with helpful information, including links to the Help Center and a document about linking to the SRS server.
- Quick Text Search:** The central search area. It features a dropdown menu set to 'Nucleotide Sequences', a text input field containing 'ATPase' (circled in red), and a 'Search' button (indicated by a red arrow). Below the input field, it lists searched databanks: EMBL, EMBLCON, and EMBLDCS.
- News and Announcements:** A section titled 'Important notes to users.' containing a list of updates and changes to the database, such as '18.10.04 - Links between EMBLRELEASE, EMBLNEW, EMBLTPA and TAXONOMY are now changed to go via Species and Organism fields rather than via NCBI_TaxID.'

1 - Interrogation par annotations

Nombre d'entrées (séquences + annotations) trouvées

Entrées disponibles, téléchargeables

LION SRS Help Center ?

Quick Searches | Select Databanks | Query Form | Tools | Results | Projects | Custom Views | Information

Reset Query "[EMBL-alltext:ATPase*]" found 88018 entries next

Apply Options to:

selected results only
 unselected results only

Result Options

Link to related information: [Link](#)

Save results: [Save](#)

Display Options

View results using: SeqSimpleView

Show 30 results per page

Printer friendly view

[Apply Display Options](#)

EMBL (Release)	Accession	Description	SeqLength
<input type="checkbox"/> EMBL (Release):AA415099	AA415099	Mg0055 RCW Lambda Zap Express Library Pyricularia grisea cDNA clone RCW55 similar to Nuclear Control of ATPase mRNA Expression (NCA3), mRNA sequence.	497
<input type="checkbox"/> EMBL (Release):AF146054	AF146054	AF146054 Lentinula edodes L54 Lentinula edodes cDNA similar to plasma membrane proton ATPase; Pma, mRNA sequence.	282
<input type="checkbox"/> EMBL (Release):AF487323	AF487323	AF487323 Tuber borchii fruit body Tuber borchii cDNA clone VA5 similar to putative sodium P-type ATPase of Neurospora crassa, mRNA sequence.	592
<input type="checkbox"/> EMBL (Release):AI392066	AI392066	NC5G10T7 Conidial Neurospora crassa cDNA clone NC5G10 3' similar to V-type ATPase subunit G, N. crassa, mRNA sequence.	610
<input type="checkbox"/> EMBL (Release):AI392186	AI392186	NCSM2B6T7 Subtracted Mycelial Neurospora crassa cDNA clone SM2B6 3' similar to putative 20 kD subunit of the V-ATPase, N. crassa, mRNA sequence.	501
<input type="checkbox"/> EMBL (Release):AI392219	AI392219	NCSP1A7T7 Subtracted Perithecial Neurospora crassa cDNA clone SP1A7 3' similar to putative 20 kD subunit of the V-ATPase, N. crassa, mRNA sequence.	503
<input type="checkbox"/> EMBL (Release):AI392246	AI392246	NCSP1A7T3 Subtracted Perithecial Neurospora crassa cDNA clone SP1A7 5' similar to putative 20 kD subunit of the V-ATPase, N. crassa, mRNA sequence.	529
<input type="checkbox"/> EMBL (Release):AI392371	AI392371	NCSM2G3T7 Subtracted Mycelial Neurospora crassa cDNA clone SM2G3 3' similar to putative 20 kD subunit of the V-ATPase, N. crassa, mRNA sequence.	467
<input type="checkbox"/> EMBL (Release):AI392500	AI392500	NCSP6A4T7 Subtracted Perithecial Neurospora crassa cDNA clone SP6A4 3' similar to putative 20 kD subunit of the V-ATPase, N. crassa, mRNA sequence.	442
<input type="checkbox"/> EMBL (Release):AI392596	AI392596	NCSP4E6T7 Subtracted Perithecial Neurospora crassa cDNA clone SP4E6 3' similar to putative 20 kD subunit of the V-ATPase, N. crassa, mRNA sequence.	669
<input type="checkbox"/> EMBL (Release):AI398415	AI398415	NCW01B12T7 Westergaards Neurospora crassa cDNA clone W01B12 3' similar to vacuolar ATP synthase 14 kD subunit (V-ATPase F subunit), mRNA sequence.	595
<input type="checkbox"/> EMBL (Release):AI399483	AI399483	NCSP6A4T3 Subtracted Perithecial Neurospora crassa cDNA clone SP6A4 5' similar to putative 20 kD subunit of the V-ATPase, N. crassa, mRNA sequence.	553
<input type="checkbox"/> EMBL (Release):AI399576	AI399576	NCSP6E7T3 Subtracted Perithecial Neurospora crassa cDNA clone SP6E7 5' similar to ATPase, H+ transporting, plasma membrane, N. crassa, mRNA sequence.	524
<input type="checkbox"/> EMBL (Release):AI399577	AI399577	NCSP6E7T7 Subtracted Perithecial Neurospora crassa cDNA clone SP6E7 3' similar to ATPase, H+ transporting, plasma membrane, N. crassa, mRNA sequence.	557
<input type="checkbox"/> EMBL (Release):AJ843936	AJ843936	Pleurotus ostreatus EST POLAM106	297
<input type="checkbox"/> EMBL (Release):AW179985	AW179985	MgA0038f MgA Library Mycosphaerella graminicola cDNA clone MgA0038 5' similar to plasma membrane H(+)-ATPase, mRNA sequence.	519
<input type="checkbox"/> EMBL (Release):AW180008	AW180008	MgA0063f MgA Library Mycosphaerella graminicola cDNA clone MgA0063 5' similar to (AF036763) P-ATPase, mRNA sequence.	759
<input type="checkbox"/> EMBL (Release):AW180009	AW180009	MgA0064f MgA Library Mycosphaerella graminicola cDNA clone MgA0064 5' similar to endoplasmic reticulum-type Ca-2+-ATPase, mRNA sequence.	733
<input type="checkbox"/> EMBL (Release):AW180064	AW180064	MgA0127f MgA Library Mycosphaerella graminicola cDNA clone MgA0127 5' similar to plasma membrane H(+)-ATPase, mRNA sequence.	577

1 - Interrogation par annotations

Changer de bases de données

The screenshot shows the LION SRS web interface. At the top, there is a navigation bar with the LION logo and a 'Help Center' link. Below the navigation bar, there are several tabs: 'Quick Searches', 'Selected Databanks', 'Query Form', 'Tools', 'Results', 'Projects', 'Custom Views', and 'Information'. A red arrow points to the 'Selected Databanks' tab. Below the navigation bar, there is a search box with a 'Quick Search' button and a 'Reset' button. On the left side, there is a 'Search Options' panel with instructions on how to use the search. Below the search options, there is a 'Browse Entries' button. On the right side, there is a 'Available Databanks' section with a list of databases and checkboxes to select them. A red arrow points to the 'EMBL (Release)' checkbox. The list of databases includes:

- Combined Sequence Databases
 - all EMBLALL GENBANK (Release+updates) UNIPROT
 - ENSEMBL ENSEMBL_CDNA ENSEMBL_RNA
 - REFSEQ_DNA_ALL ENSEMBL_PROT REFSEQ_PROT_ALL
- General DNA Databases
 - all EMBL (Release) EMBL (Updates) EMBL (Whole Genome Shotgun Sequences) GENBANK (Release)
 - GENBANK (Updates) NRNUC RefSeq_DNA (release) RefSeq_DNA (updates)
- General Protein Databases
 - all SWISSPROT (Release) SPTR EMBL NRPEP
 - IPI GENPEPT RefSeq_Prot (release)
 - RefSeq_Prot (updates) UNIREF100 UNIREF90
 - UNIREF50 COGS_SEQUENCES EXPROT
 - SWISSPROT (Updates) SWISSPROT_SPLICEVAR
- Outdated Protein Databases
- cDNA Sequences
- EST&STS
- Eukaryotic Genomes - NCBI - Genomes
- Eukaryotic Genomes - NCBI - Proteins
- Eukaryotic Genomes - NCBI - mRNAs
- Eukaryotic Genomes - NCBI - Contig Layouts
- Eukaryotic Genomes - Ensembl - Genomes
- Eukaryotic Genomes - Ensembl - Proteins
- Eukaryotic Genomes - Ensembl - cDNAs
- Eukaryotic Genomes - Ensembl - RNAs
- Eukaryotic Genomes - Other Libs
- Microbial Genomes
- Genome Annotations
- Unigene representative sequences
- Unigene Clusters - taxonomic subsections
- Organism specific tissue distribution
- Non-Coding RNAs
- Contig Databases

1 - Interrogation par annotations

SRS@EMBL-EBI Quick Search Library Page **Query Form** Tools Results Projects Views Databanks HELP

Reset search MEDLINE Patent Abstracts EMBL
IMGTHLA UniProt

Search Options

Combine search terms with:

Use wildcards

Get results of type:

Result Display Options

View results using:

or

Create a view

Show
results per page

Fields you can search | **Your search terms**

In a single field, you can separate multiple values by &, |, !

<input type="button" value="i"/>	<input type="text" value="AllText"/> <input type="button" value="v"/>	<input type="text" value="ATPase"/>
<input type="button" value="i"/>	<input type="text" value="AllText"/> <input type="button" value="v"/>	
<input type="button" value="i"/>	<input type="text" value="AllText"/> <input type="button" value="v"/>	
<input type="button" value="i"/>	<input type="text" value="AllText"/> <input type="button" value="v"/>	

Create a view

Select the fields you want displayed in your view and choose the format

Choose 1 or more fields:

Display As: Table List

Sequence Format:

1 - Interrogation par annotations

The screenshot shows the EMBL SRS Query Form interface. The 'Query Form' tab is highlighted with a red circle. The search bar contains the text 'search EMBL UniProt UniProt/TrEMBL'. The 'Search Options' panel on the left is also circled in red, showing 'Combine search terms with: & (AND)', 'Use wildcards' checked, and 'Get results of type: Entry'. The 'Result Display Options' panel shows 'View results using: * Names only *' and 'Show 30 results per page'. The 'Fields you can search' dropdown menu is open, listing various fields like 'AllText', 'Organism Name', 'Sequence Length', etc., and is circled in red. The 'Your search terms' table contains 'ATPase' and 'homo sapiens'. A red arrow points to the 'Search' button.

SRS@ EMBL-EBI Quick Search Library Page **Query Form** Tools Results Projects Views Databanks **HELP**

Reset search EMBL UniProt UniProt/TrEMBL

Search Options

Combine search terms with: & (AND)

Use wildcards

Get results of type: Entry

Result Display Options

View results using: * Names only *

or

Create a view

Show 30 results per page

Fields you can search	Your search terms
In a single field, you can separate multiple values by &, , ! <input type="button" value="v"/> <input type="button" value="Search"/>	
<input type="button" value="i"/> AllText <input type="button" value="v"/>	ATPase
<input type="button" value="i"/> Organism Name <input type="button" value="v"/>	homo sapiens
<input type="button" value="i"/> AllText	
<input type="button" value="i"/> ID	
<input type="button" value="i"/> Sequence Length	
<input type="button" value="i"/> Accession Number	
<input type="button" value="i"/> Entry Creation Date	
<input type="button" value="i"/> Description	
<input type="button" value="i"/> Keywords	
<input type="button" value="i"/> Organism Name	
<input type="button" value="i"/> NCBITaxid	
<input type="button" value="i"/> Organelle	
<input type="button" value="i"/> References: RefPosition	
<input type="button" value="i"/> References: MedlineID	
<input type="button" value="i"/> References: RefGroup	
<input type="button" value="i"/> References: Authors	
<input type="button" value="i"/> References: Title	
<input type="button" value="i"/> References: Journal	
<input type="button" value="i"/> References: VolumeNo	
<input type="button" value="i"/> References: FirstPage	
<input type="button" value="i"/> References: Year	

Displayed in your view and choose the format

Display As: Table List

Sequence Format: embl

1 - Interrogation par annotations

Query Results - Microsoft Internet Explorer

Fichier Edition Affichage Favoris Outils ?

Précédente Rechercher Favoris Média

Adresse <http://srs.ebi.ac.uk/srsbin/cgi-bin/vqgetz> OK

SRS@EMBL-EBI

Quick Search Library Page Query Form Tools Results Projects Views Databanks HELP

Reset Query "[([libs={embl uniprot sptrembl}-AllText:ATPase*) & (([libs-Organism:homo*] & [libs-Organism:sapiens*]) | [libs-Organism:homo sapiens*])]" found 3981 entries next

Apply Options to:

- selected results only
- unselected results only

Result Options

Launch analysis tool:
BlastN

Show tools relevant to these results:

Link to related information:

Save results:

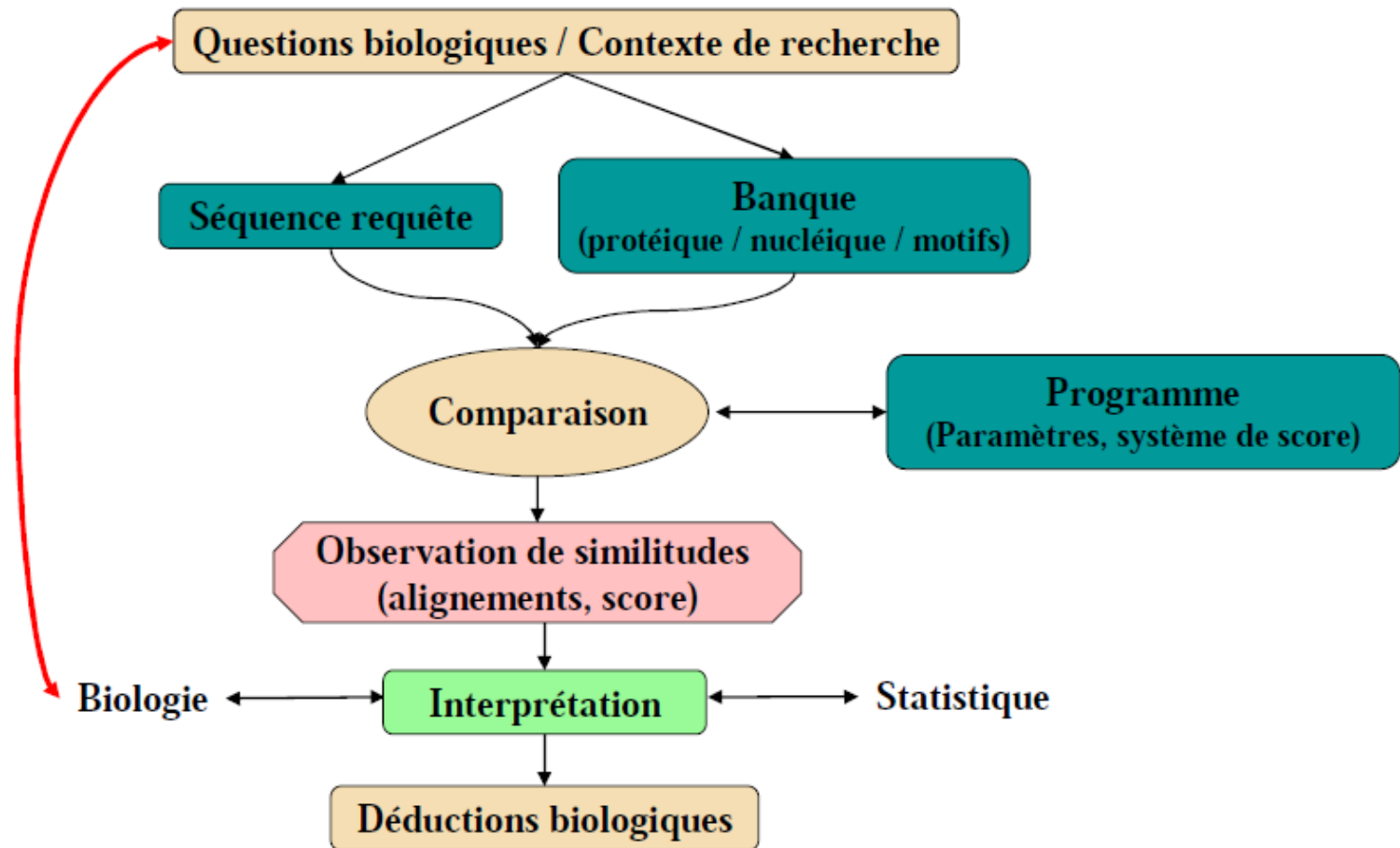
Display Options

EMBL UniProt/TrEMBL UniProt	Accession	Description	SeqLength
<input type="checkbox"/> EMBL:HSM800272	AL049929	Homo sapiens mRNA; cDNA DKFZp564C0582 (from clone DKFZp564C0582)	1884
<input type="checkbox"/> EMBL:HSM802065	AL137377	Homo sapiens mRNA; cDNA DKFZp434K0126 (from clone DKFZp434K0126)	1862
<input type="checkbox"/> EMBL:HSM802623	AL161996	Homo sapiens mRNA; cDNA DKFZp434P0831 (from clone DKFZp434P0831)	5785
<input type="checkbox"/> EMBL:HSM804282	AL832971	Homo sapiens mRNA; cDNA DKFZp666G172 (from clone DKFZp666G172)	3613
<input type="checkbox"/> EMBL:HSM804295	AL832984	Homo sapiens mRNA; cDNA DKFZp666D103 (from clone DKFZp666D103)	2400
<input type="checkbox"/> EMBL:HSM800103	AL833813	Homo sapiens mRNA; cDNA DKFZp564C236 (from clone DKFZp564C236)	2161
<input type="checkbox"/> EMBL:HSM805474	AL834179	Homo sapiens mRNA; cDNA DKFZp761L1023 (from clone DKFZp761L1023)	3387
<input type="checkbox"/> EMBL:HSM805442	AL834370	Homo sapiens mRNA; cDNA DKFZp762C1113 (from clone DKFZp762C1113)	2543
<input type="checkbox"/> EMBL:BC018811	BC018811	Homo sapiens proteasome (prosome, macropain) 26S subunit, ATPase, 4, mRNA (cDNA clone IMAGE:2961362).	1426
<input type="checkbox"/> EMBL:BC018859	BC018859	Homo sapiens proteasome (prosome, macropain) 26S subunit, ATPase, 1, mRNA (cDNA clone IMAGE:2111012)	1559

2 - Interrogation par similarité - Recherche dans les séquences

Pourquoi ? => Savoir si ma séquence ressemble à d'autres déjà connues ; Trouver toutes les séquences d'une même famille ; Rechercher toutes les séquences qui contiennent un motif donné

Comment ? => Comparaison d'une séquence aux séquences de la banque : BLAST, FASTA, Blat, YASS, ...



Programme = alignement + statistique

→ analyse biologique

2 - Interrogation par similarité – les alignements

Notion clé : l'alignement – 2 types d'alignements :



→ utilisé en phylogénie

Alignement global: alignement de toute la séquence A avec toute la séquence B

Méthode employée pour aligner des séquences dont on soupçonne l'homologie.

L'alignement est optimisé sur toute la longueur des séquences.

L'algorithme de référence est celui de **Needleman & Wunsch**.



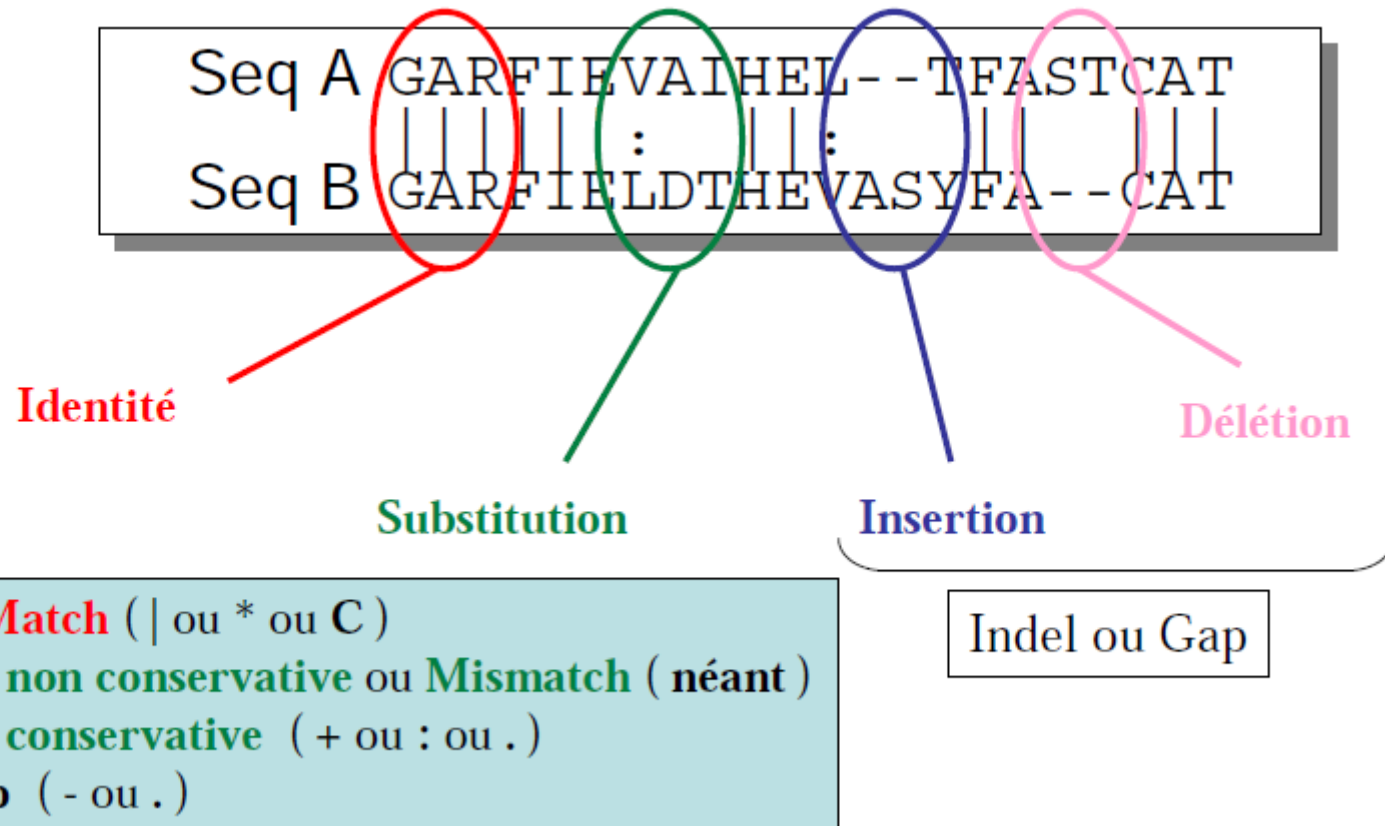
→ utilisé pour la recherche de similarité



Alignement local: alignement de sous-séquences de A avec des sous-séquences de B

Aligne seulement les régions dont le score est supérieur à un seuil donné. Utilisé lorsque l'on veut aligner deux séquences de taille très différente. (par ex. dans une recherche de sous-séquence). Beaucoup plus rapide que l'alignement global. Algorithme de **Smith & Waterman**.

2 - Interrogation par similarité – types d'alignements



3 situations sont possibles pour une position donnée de l'alignement :

- les caractères sont les mêmes: Identité ou match
- les caractères ne sont pas les mêmes: Substitution
- les caractères existent dans un cas et pas dans l'autre : Insertion ou Délétion

2 - Interrogation par similarité - Définitions

Identité (globale) :

Proportion des paires de résidus identiques entre deux séquences alignées. (Exprimé généralement en %).

Similitude (=similarité) :

Mesure de la ressemblance entre deux séquences. Le degré de similitude est quantifié par un score basé sur le % de similarité (% identité + % substitutions conservatives) des séquences. De 100% à quelques nucléotides/aminoacides en commun.

Gaps ou InDels :

Proportion d'Indels entre deux séquences alignées (Exprimé en %).

Homologie :

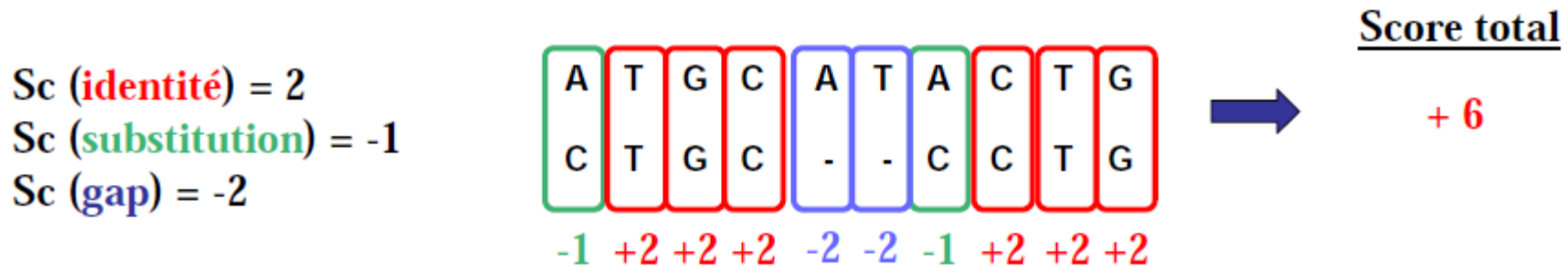
Deux séquences sont homologues si elles ont un ancêtre commun. Il n'y a pas vraiment de degré d'homologie (remarque : on ne dit pas: tres homologue, faible homologie, etc.) Il n'y a pas vraiment de limite, mais en dessous de 20-25% (twilight zone), il devient très difficile de distinguer une homologie d'une ressemblance fortuite.

A noter aussi que des séquence sans ressemblance apparente peuvent aussi être homologues (on le retrouve par ex. au niveau 3D) =>Des séquences homologues ne sont pas nécessairement similaires.

2 - Interrogation par similarité – Qualité alignements : Le score

Le score de l'alignement est la somme de toutes les positions 2 à 2.

➤ **Exemple:** On peut associer une récompense (positive) à des symboles alignés identiques et une pénalité (négative) à un substitution ou à un gap.



=> bon alignement donnera le score maximum entre 2 sequences

2 - Interrogation par similarité – alignement sur Bases de seq.

Comparaison d'une séquence aux séquences de la banque

=> Identifier les [fragments de] séquences de la banque ayant une forte ressemblance avec la séquence requête

* Requête : une ou plusieurs séquences (ADN/ARN ou protéine)

* Résultat : une liste de séquences ressemblant à la séquence entrée (classées par score, identité, couverture, e-value)

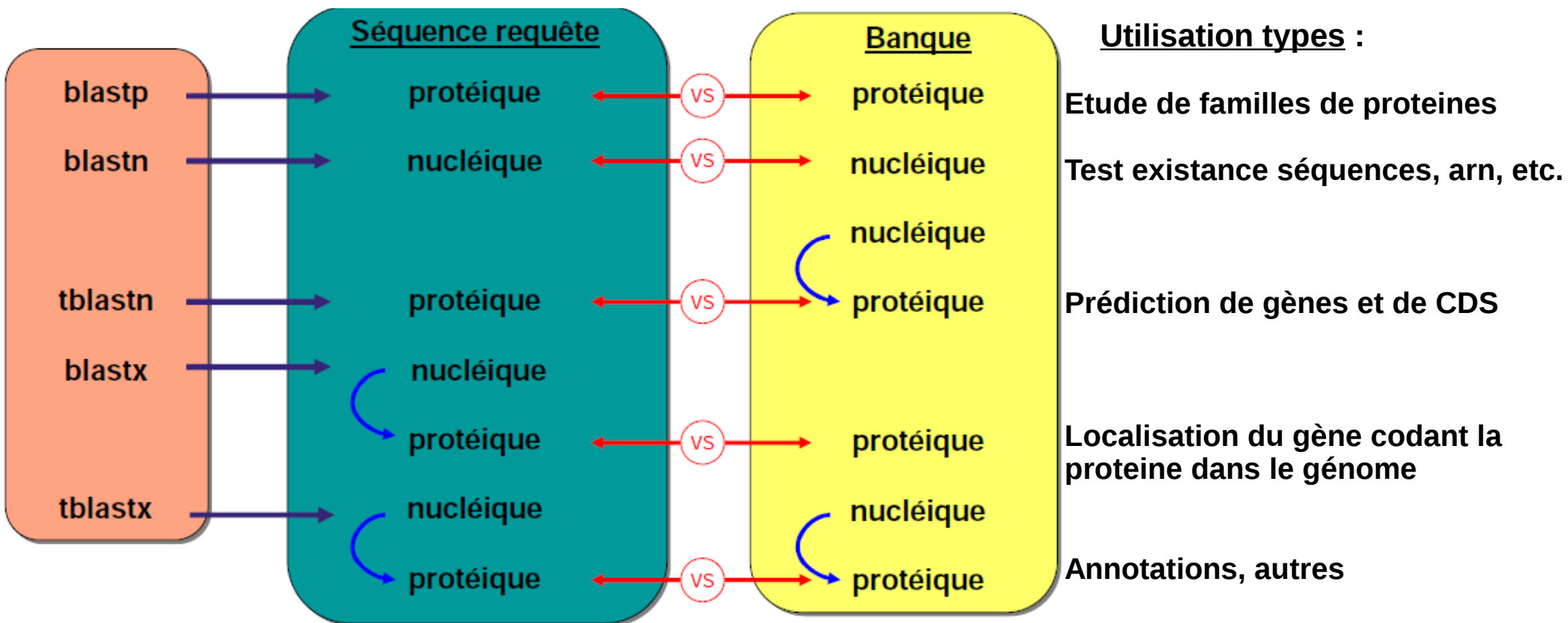
=> exemples suivants avec "Blast" (avec interface graphique du NCBI)

<http://blast.ncbi.nlm.nih.gov/Blast.cgi>

2 - Interrogation par similarité – Blast

L'algorithme est basée sur un modèle statistique (Karlin et Altschul, 1990) qui s'applique aux comparaisons de séquences sans insertion / délétion.

Différents programmes :



[blastn](#)[blastp](#)[blastx](#)[tblastn](#)[tblastx](#)

Enter Query Sequence

BLASTN programs search nucleotide databases using a nucleotide query. [more.](#)

Enter accession number(s), gi(s), or FASTA sequence(s) ⓘ

[Clear](#)

Query subrange ⓘ

From To

Or, upload file

[Chosir...](#) ⓘ

Job Title

Enter a descriptive title for your BLAST search ⓘ

 Align two or more sequences ⓘ

Choose Search Set

Database

 Human genomic + transcript Mouse genomic + transcript Others (nr etc.):

Nucleotide collection (nr/nt) ⓘ

Organism

Optional

 Enter organism name or id--completions will be suggested Exclude ⓘ

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown ⓘ

Exclude

Optional

 Models (XM/XP) Uncultured/environmental sample sequences

Entrez Query

Optional

Enter an Entrez query to limit search ⓘ

Program Selection

Optimize for

- Highly similar sequences (megablast)
- More dissimilar sequences (discontiguous megablast)
- Somewhat similar sequences (blastn)

Choose a BLAST algorithm ⓘ

BLASTSearch **database Nucleotide collection (nr/nt)** using **Blastn (Optimize for somewhat similar sequences)** Show results in a new window [Algorithm parameters](#)[General Parameters](#)

Algorithm parameters

Note: Parameter values that differ from the default are highlighted in yellow and marked with ♦ sign

General Parameters

- Max target sequences: 100 (Select the maximum number of aligned sequences to display)
- Short queries: Automatically adjust parameters for short input sequences
- Expect threshold: 10
- Word size: ♦ 16 (highlighted in yellow)
- Max matches in a query range: 0

Scoring Parameters

- Match/Mismatch Scores: 1,-2
- Gap Costs: Linear

Filters and Masking

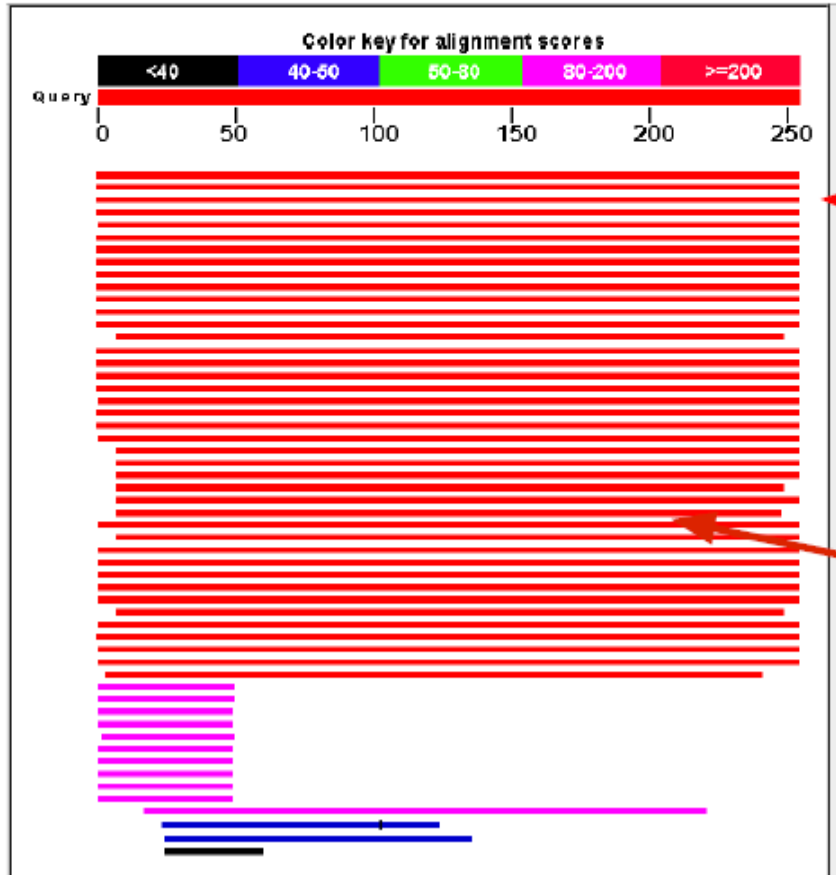
- Filter: Low complexity regions; Species-specific repeats for: Homo sapiens (Human)
- Mask: Mask for lookup table only; Mask lower case letters

BLAST

Search database Nucleotide collection (nr/nt) using Megablast (Optimize for highly similar sequences)

Show results in a new window

er to show define and scores, click to show alignments



2. représentation graphique des résultats

ce trait représente
la séquence soumise (long. 253 AA)

chaque trait de couleur représente
un alignement entre la séquence de
départ et une séquence de la
banque de donnée sélectionnée
couleur → score
longueur → taille de l'alignement

= HSP ("high scoring pair")

3. Résumé des résultats

Descriptions

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

[Alignments](#)
[Download](#)
[GenBank](#)
[Graphics](#)
[Distance tree of results](#)

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	Uncultured bacterium clone Z206 16S ribosomal RNA gene, partial sequence	2676	2676	98%	0.0	99%	GQ388829.1
<input type="checkbox"/>	Uncultured Afipia sp. clone NW-12 16S ribosomal RNA gene, partial sequence	2663	2663	99%	0.0	99%	AY568510.1
<input type="checkbox"/>	Uncultured Afipia sp. 16S ribosomal RNA gene, partial sequence	2658	2658	98%	0.0	99%	FJ572667.1
<input type="checkbox"/>	Uncultured bacterium clone MA-78-I98C 16S ribosomal RNA gene, partial sequence	2656	2656	97%	0.0	99%	HM141889.1
<input type="checkbox"/>	Uncultured bacterium clone BL 16S ribosomal RNA gene, partial sequence	2643	2643	97%	0.0	99%	FJ154970.1
<input type="checkbox"/>	Afipia sp. BAC308 16S ribosomal RNA gene, partial sequence	2643	2643	98%	0.0	99%	EU130950.1
<input type="checkbox"/>	Afipia sp. OHSU_I 16 ribosomal RNA, tRNA-Ile, tRNA-Ala, 23S ribosomal RNA, and 5S riboso	2615	2615	100%	0.0	99%	KC677616.1
<input type="checkbox"/>	Afipia sp. OHSU_II 16 ribosomal RNA, tRNA-Ile, tRNA-Ala, 23S ribosomal RNA, and 5S riboso	2615	2615	100%	0.0	99%	KC677615.1
<input type="checkbox"/>	Uncultured Afipia sp. clone AV_6J-C05 16S ribosomal RNA gene, partial sequence	2606	2606	96%	0.0	99%	EU341227.1
<input type="checkbox"/>	Uncultured bacterium clone SPCUL1A1 16S small subunit ribosomal RNA gene, partial seque	2601	2601	95%	0.0	99%	AY186080.1
<input type="checkbox"/>	Uncultured Afipia sp. clone C-12 16S ribosomal RNA gene, partial sequence	2595	2595	99%	0.0	98%	AY568503.2

4. Présentation des alignements obtenus pour chaque séquence sélectionnée de la banque

```

Sbjct 1321 |CGGGCCTTGTACACACCGCCCGTCACACCATGGGAGTTGGTTCTACCTGAAGGCAGTGCG| 1380
Query 1382 |CTAACCCGCAAGGGAGGCAGCTGACCACGGTAGGGTCAGCGACTGGGGTGAAGTCGTAAC| 1441
Sbjct 1381 |CTAACCCGCAAGGGAGGCAGCTGACCACGGTAGGGTCAGCGACTGGGGTGAAGTCGTAAC| 1440
Query 1442 |AAGGTAGCCGTA| 1453
Sbjct 1441 |AAGGTAACCGTA| 1452

```

[Download](#) [GenBank](#) [Graphics](#)

Uncultured Afipia sp. clone NW-12 16S ribosomal RNA gene, partial sequence

Sequence ID: [gb|AY568510.1](#) Length: 1487 Number of Matches: 1

Range 1: 2 to 1483 [GenBank](#) [Graphics](#)

[Next Match](#) [Previous Match](#)

Score	Expect	Identities	Gaps	Strand
2663 bits(1442)	0.0	1470/1482(99%)	8/1482(0%)	Plus/Plus
Query 4	GTTTGATCCTGGCTCAGAGCGAACGCTGGCGGCAGGCTTAACACATGCAAGTCGAGCGGG	63		
Sbjct 2	GTTTGATCCTGGCTCAGAGCGAACGCTGGCGGCAGGCTTAACACATGCAAGTCGAGCGGG	61		
Query 64	CGTAGCAATACGTCAGCGGCAGACGGGTGAGTAACGCGTGGGAACGTACCTTTTGGTTTCG	123		
Sbjct 62	CGTAGCAATACGTCAGCGGCAGACGGGTGAGTAACGCGTGGGAACGTACCTTTTGGTTTCG	121		
Query 124	GAACAACGAGGGAACTTCAGCTAATACCGGATAAGCCCTTACGGGGAAAGATTTATCG	183		
Sbjct 122	GAACAACGAGGGAACTTCAGCTAATACCGGATAAGCCCTTACGGGGAAAGATTTATCG	181		
Query 184	CCGAAAGATCGGCCCGCTCTGATTAGCTAGTTGGTGAGGTAATGGCTCACCAAGGCGAC	243		
Sbjct 182	CCGAAAGATCGGCCCGCTCTGATTAGCTAGTTGGTGAGGTAACGGCTCACCAAGGCGAC	241		
Query 244	GATCAGTAGCTGGTCTGAGAGGATGATCAGCCACATTGGGACTGAGACACGGCCCAA	303		
Sbjct 242	GATCAGTAGCTGGTCTGAGAGGATGATCAGCCACATTGGGACTGAGACACGGCCCAA	301		
Query 304	CCTACGGGAGGCAGCAGTGGGGAATATTGGACAATGGGCGAAAGCCTGATCCAGCCATGC	363		
Sbjct 302	CCTACGGGAGGCAGCAGTGGGGAATATTGGACAATGGGCGAAAGCCTGATCCAGCCATGC	361		
Query 364	CGCGTGAGTGATGAAGGCCCTAGGGTTGTAAGCTCTTTTGTGCGGGAAGATAATGACGG	423		
Sbjct 362	CGCGTGAGTGATGAAGGCCCTAGGGTTGTAAGCTCTTTTGTGCGGGAAGATAATGACGG	421		
Query 424	TACCGCAAGAATAAGCCCCGGCTAACTTCGTGCCAGCAGCCCGGTAATACGAAGGGGGC	483		

Exercice

1 – “séquences cibles” : récupération des séquences (homologues) de « *Bacillus anthracis* », *gene* « *cya* »

2 - ”séquences non-cibles” : récupération des séquences proches, mais n'étant pas de *Bacillus anthracis* »

Methodes à utiliser :

1.1 En effectuant une recherche par mots-clefs grâce aux applications suivantes : SRS (<http://www.dkfz.de/srs/>), Entrez (<http://www.ncbi.nlm.nih.gov>) ou WWWQuery (=ACNUC) (http://doua.prabi.fr/search/query_fam), récupérer sur votre disque dur toutes les séquences codantes (CDS) relatives au gène “*cya*” de *Bacillus anthracis*. Sauvegarder dans un fichier fasta l'ensemble de ces séquences cibles.

1.2 Pour compléter (ou vérifier) le jeu de séquences cibles, à l'aide d'une séquence représentative de votre résultat précédent, effectuer une recherche par similarité des séquences du gène, en utilisant le programme BLASTn (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>).

Ce programme renvoie pour chaque séquence trouvée un score, une e-value et une identité. Servez-vous de ces valeurs pour sélectionner les séquences qui pourraient être ajoutées à votre jeu de séquences cibles.

2 Toujours par similarité (BLASTn donc), récupérer les séquences proches, mais n'étant pas de « *Bacillus anthracis* », afin de constituer un jeu de séquences proches «non-cibles ». Sauvegarder ces séquences dans un autre fichier fasta.

Pour ceux qui sont à l'aise avec les notions précédentes (ou qui ont terminé), effectuer la même recherche par Blast, mais en ligne de commande en utilisant le serveur du laboratoire
=> http://www.bioinfomed.fr/__Teachings/bioinfo_urmite/blast_URMITE_tuto/blast_URMITE_tuto.html