



H3ABioNet

Pan African Bioinformatics Network for H3Africa

Introduction to Bioinformatics Online Course: IBT

Genomics Comparative Genomics

Learning Objectives

Session 3: Comparative Genomics

- **Part 1: (intro)** Comparative Genomics: what is it?
- **Part 2:** Genomic Variation / Comparative Genomics: WWWH?
- **Part 3:** The input
- **Part 4:** The methods
- **Part 5:** The output

Learning Outcomes

Session 1: Comparative Genomics

- Navigate through genomic resources to:
 - ▶ Retrieve information on a specific gene (sequence, variants, orthologs...)
 - ▶ View and interpret genomic alignments

Part 1

(intro)

Comparative Genomics: What is it?



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics| Fatma Guerfali

Comparative Genomics: a definition

- Comparative genomics is based on the fact that a genomic variation is happening in all organisms.
- These changes affects several features in a genome (structure, organization, functions...).
- The changes could help monitor evolution between organisms (species...)

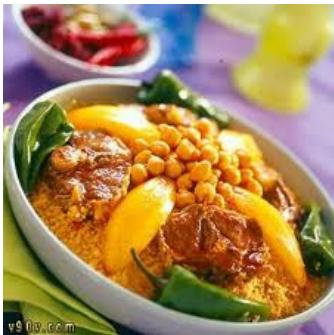
Comparative Genomics: a definition

- Comparative genomics deals with the process of comparing the sequences of **whole or parts of genomes**.
- Goal:
 - identify **similarities and differences** between **features** in these genomes
 - Identify **evolutionary relationship** between organisms

Comparative Genomics: a definition

- Comparative genomics : features + evolution

ORGANISM A1
GENOME A1
FEATURES A1



ORGANISM A2
GENOME A2
FEATURES A2



ORGANISM B
GENOME B
FEATURE B



Adapted from <http://recettes-aymen.over-blog.fr/>
<http://quebueno.be/content/6>
Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Part 2

Genomic Variation & Comparative Genomics: WWWH?

Genomic variation : WWWH

WHY ?

WHEN ?

WHAT ?

HOW ?

Genomic variation : WWWH

WHY ?

● Why would a genome evolve ?

WHEN ?

→ Genomic plasticity allows an organism to:

WHAT ?

- adapt to environmental changes
- Find the best evolution path
- Acquire virulence genes, enhanced pathogenicity
- Resistance to drugs
- Increase survival chances of members of a population
- ...

HOW ?



Genomic variation : WWWW

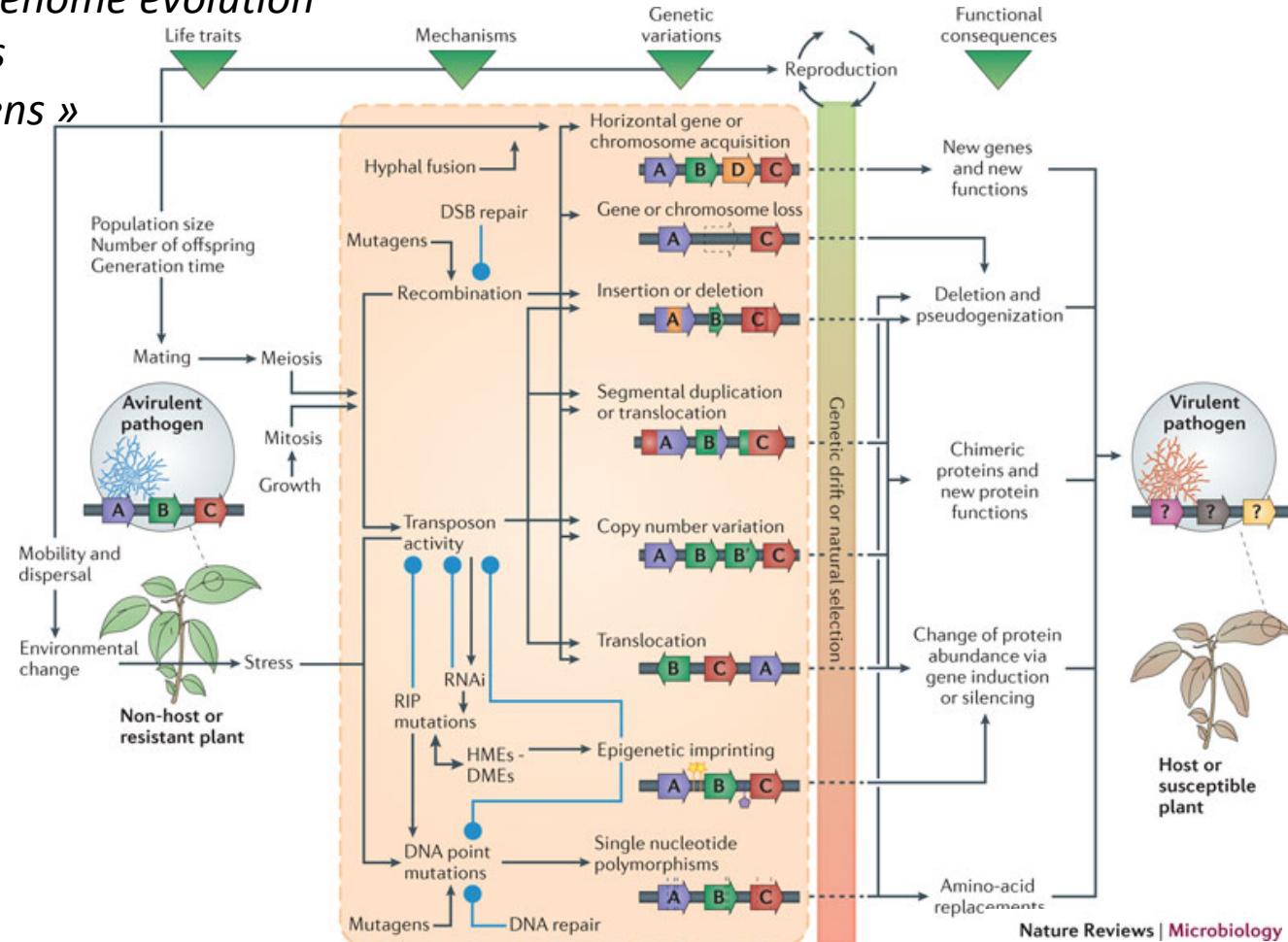
WHY ?

Example : « *Genome evolution in filamentous plant pathogens* »

WHEN ?

WHAT ?

HOW ?



Nature Reviews | Microbiology

(Raffaele & Kamoun, 2012)

Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Genomic variation : WWWH

WHY ?

● Factors/events

WHEN ?

- Gene transfer
- Environmental pressure for selection
 - pH
 - temperature
 - host
 - pathogen
- ...

HOW ?

→ A genetic variation could occur in response to such factors / events

Genomic variation : WWWH

WHY ?

● What could be affected ?

WHEN ?

- Overall genomic sequence (re-arrangements)
- DNA structure
- Regulatory elements
- Genes size, number, function, density
- Nucleotide composition
- ...

WHAT ?

HOW ?



H3ABioNet

Pan African Bioinformatics Network for H3Africa



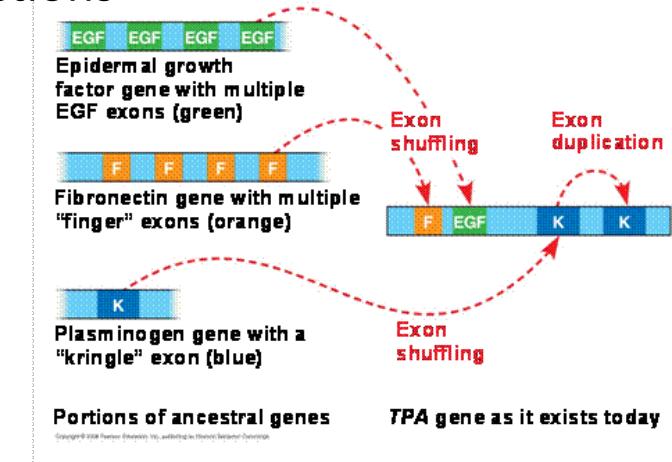
Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Genomic variation : WWWH

WHY ?

How could this happen ?

- Large genetic structural variations (duplication, recombination...)
- Transposable elements (retrotransposons...)
- Evolution of multigene families
- Evolution of genes with novel functions
- Exon shuffling
- Tandem repeats modification
- ...



Genomic variation : WWWH

WHY ?

● How to measure the changes in a genome?

WHEN ?

- Sequence variation ← Comparative Genomics
 - between 2 genomes (1 reference)
 - between several genomes
- Other variations (structure, folding...)

WHAT ?

HOW ?



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Comparative Genomics : WWWH

WHY ?

WHEN ?

WHAT ?

HOW ?

Comparative Genomics : WWWH

WHY ?

● Why would we compare genomes ?

WHEN ?

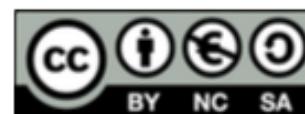
- Identify evolutionary history
- Highlight synteny
- Identify genomic rearrangements (large SV events...)
- Study convergent evolution for some organisms (e.g. virus)
- understand disease outbreak
- Identify pathogenicity markers, drug targets
- ...

HOW ?



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Comparative Genomics : WWWH

WHY ?

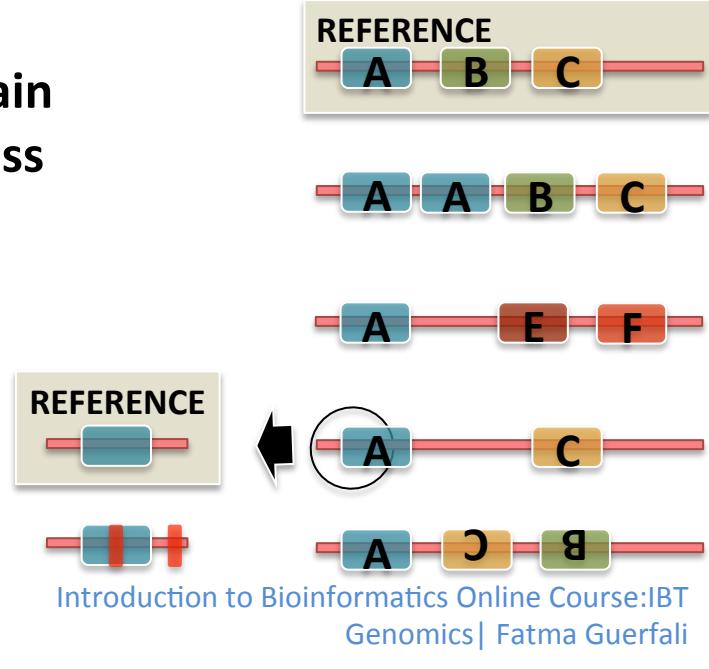
- Why would we compare “portions” of genomes ?

WHEN ?

- Comparing smaller portions of a genome allows to zoom into regions of genomic re-arrangements
- Could be **genes**:
 - Screen for functional genes **gain**
 - Screen for functional genes **loss**
 - Gain of a **new function**
 - **Exons** (length, number...)
 - Conserved **pathways**
 - **Coding / non-coding**
 - ...

WHAT ?

HOW ?



Comparative Genomics : WWWH

WHY ?

● When do we need to use comparative genomics ?

WHEN ?

→ Establish genetic and evolutionary relationship between :

WHAT ?

- Entire organisms
- Sequences

HOW ?



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Comparative Genomics : WWWH

WHY ?

- What we generally compare are **features** of 1 or more genomes to features of a another **genome (reference)**

WHEN ?

- A genome is complex and composed of different elements (regulatory, structural...)

WHAT ?

- In fact, there are different types of DNA features that can be compared between 2 genomes:

- **DNA sequences** (small, large, coding/non-coding)
- **Genes** (nature, order...)
- Regulatory elements
- ...



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Comparative Genomics : WWWH

WHY ?

- Could be classified in:

WHEN ?

- Genome **structure**

WHAT ?

- Genome **function** (coding / non-coding)

HOW ?

- Genome **evolution**

Comparative Genomics : WWWH

WHY ?

- Comparative genomics uses **Sequence Alignment**
- Comparative genomics is based on **Phylogeny** that relies on several key issues:

WHAT ?

- Several genomes are sequenced and available
- Homology between genes (similar functions)
- ...

HOW ?

→ Use **complex** model genomes to **infer knowledge** (**Annotation**: function...) to **unknown or less complex** genomes



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Comparative Genomics : WWWH

WHY ?

- Algorithms/programs

WHEN ?

- *in vitro*:

WHAT ?

- Fluorescence In Situ Hybridization (FISH)
- Spectral Karyotyping (SKY) and Multiplex-FISH (M-FISH)
- Comparative Genomic Hybridization

HOW ?



H3ABioNet

Pan African Bioinformatics Network for H3Africa



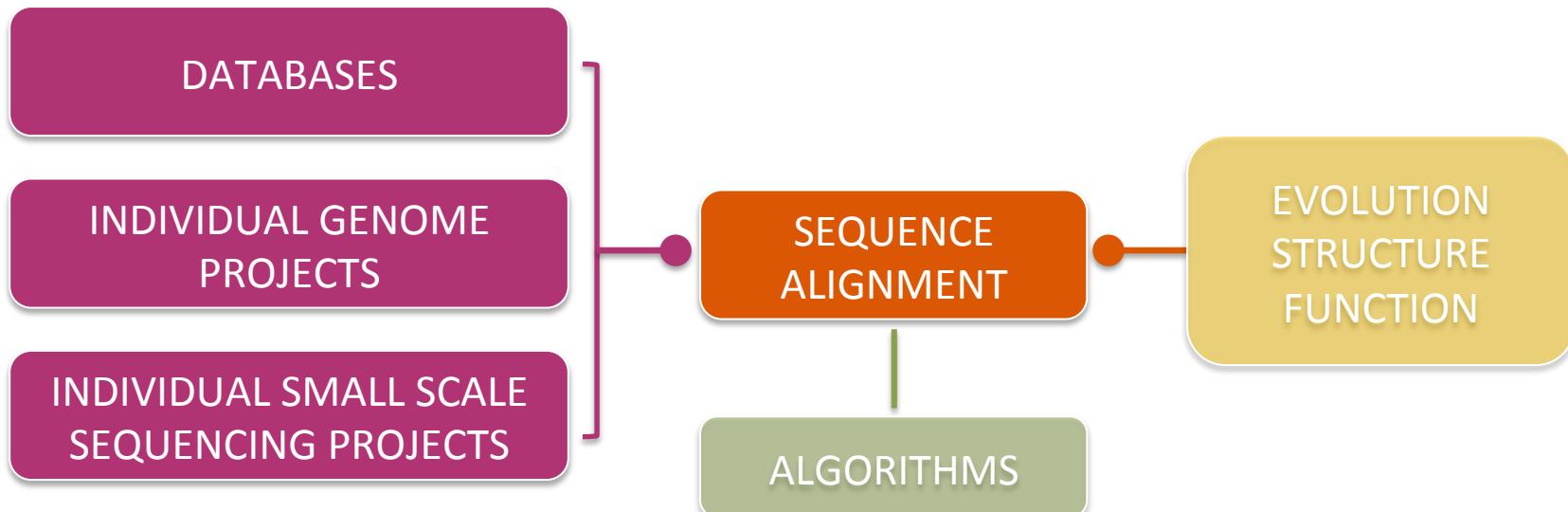
Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Keywords in Comparative Genomics

THE INPUT

THE METHODS

THE OUTPUT



Part 3

Comparative Genomics: The input



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics| Fatma Guerfali

Availability of genomes sequenced



Availability of genomes sequenced

NCBI Resources ▾ How To ▾ Sign in to NCBI

Genome

Genome Information by organism

[Overview \[16788\]](#) [Eukaryotes \[3332\]](#) [Prokaryotes \[72198\]](#) [Viruses \[5641\]](#) [Plasmids \[7519\]](#) [Organelles \[8493\]](#)

Availability of databases

OXFORD JOURNALS CONTACT US MY BASKET MY ACCOUNT

Nucleic Acids Research

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE SEARCH

OXFORD JOURNALS CONTACT US MY BASKET MY ACCOUNT

DATABASE The Journal of Biological Databases and Curation

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE SEARCH

Availability of databases



OXFORD JOURNALS

Nucleic Acids Research

CONTACT US MY BASKET MY ACCOUNT

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE SEARCH

Assembly: a resource for assembled genomes at NCBI

Paul A. Kitts*, Deanna M. Church, Françoise Thibaud-Nissen, Jinna Choi, Vichet Hem, Victor Sapojnikov, Robert G. Smith, Tatiana Tatusova, Charlie Xiang, Andrey Zherikov, Michael DiCuccio, Terence D. Murphy, Kim D. Pruitt and Avi Kimchi

Author Affiliations

Author Notes

 *To whom correspondence should be addressed. Tel: +1 301 435 6091; Fax: +1 301 480 0109; Email: kitts@ncbi.nlm.nih.gov

Received September 12, 2015.
Revision received October 21, 2015.
Accepted October 29, 2015.

« *The NCBI Assembly database (www.ncbi.nlm.nih.gov/assembly/) provides stable accessioning and data tracking for genome assembly data (...)*

The Assembly database reports **metadata** such as assembly names, simple statistical reports of the assembly (number of contigs and scaffolds, contiguity metrics such as contig N50, total sequence length and total gap length) as well as the assembly update history. *The Assembly database also tracks the relationship between an assembly submitted to the International Nucleotide Sequence Database Consortium (INSDC) and the assembly represented in the NCBI RefSeq project. Users can find assemblies of interest by querying the Assembly Resource directly or by browsing available assemblies for a particular organism. Links in the Assembly Resource allow users to easily download sequence and annotations for current versions of genome assemblies from the NCBI genomes FTP site»*

Availability of databases



OXFORD JOURNALS

DATABASE The Journal of Biological Databases and Curation

CONTACT US MY BASKET MY ACCOUNT

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE SEARCH

ANItools web: a web tool for fast genome comparison within multiple bacterial strains



Na Han^{1,2}, Yujun Qiang^{1,2} and Wen Zhang^{1,2,*}

+ Author Affiliations

✉ Corresponding author: Tel: (86)-010-61739446; Email: zhangwen@icdc.cn

Received November 25, 2015.
Revision received April 5, 2016.
Accepted April 28, 2016.

« Background: Early classification of prokaryotes was based solely on phenotypic similarities, but modern prokaryote characterization has been strongly influenced by advances in genetic methods. **With the fast development of the sequencing technology, the ever increasing number of genomic sequences per species offers the possibility for developing distance determinations based on whole-genome information. The average nucleotide identity (ANI), calculated from pair-wise comparisons of all sequences shared between two given strains, has been proposed as the new metrics for bacterial species definition and classification.**

Results: In this study, we developed the web version of ANItools (<http://ani.mypathogen.cn/>), which helps users directly get ANI values from online sources. A database covering ANI values of any two strains in a genus was also included (2773 strains, 1487 species and 668 genera) »

File formats

- Different file formats may be accepted
- Blast (NCBI) accepts:
 - Fasta sequence
 - simple sequence
 - Accession Number
 - Local file from disk (megablast)
- Whole sequence (or subsequence)

FASTA

```
>gi|343488507:5001-16300 Homo sapiens CD74 molecule (CD74), RefSeqGene on
chromosome 5
CTGCCCTGGGGAGCCCCCCCCCCCCACATCCTGCCCGCAAAGGCAGCTCACCAAAGTGGGTATTCC
AGCCTTTGTAGCTTCACTTCCACATCTACCAAGTGGCGGAGTGGCCTCTGTGGACGAATCAGATTCC
TCTCCAGCACCACCTTAAGAGGCGAGCCGGGGGTCAAGGTCCCAGATGCACAGGAGGAGAACGAGGAG
CTGTCGGGAAGATCAGAACGCCAGTCATGGATGACCAGCGCACCTTATCTCCAACAATGAGCAACTGCC
```

GenBank

```
1 ctgcctgggg agcccccccg ccccacatcc tgccccgcaa aaggcagctt caccaaagtg
61 gggtatttcc agccttgta gcttcactt ccacatctac caagtggcg gagtggcctt
121 ctgtggacga atcagattcc tctccagcac cgactttaag aggcgagccg gggggtcagg
181 gtcccagatg cacaggagga gaagcaggag ctgtcggaa gatcagaagc cagtcatgga
```

Accession

ACCESSION [NG_029730](#) REGION: 5001..16300
 VERSION NG_029730.1 GI:343488507

[>gi|343488507](#)



<http://database.oxfordjournals.org/>

Introduction to Bioinformatics Online Course:IBT
 Genomics | Fatma Guerfali

File formats

NIH U.S. National Library of Medicine NCBI National Center for Biotechnology Information Sign in to NCBI

BLAST®

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

NEWS

BLAST and HTTPS
The BLAST URL API is moving to HTTPS.
Thu, 28 Jul 2016 12:00:00 EST [More BLAST news...](#)

Web BLAST

Nucleotide BLAST
nucleotide ▶ nucleotide

blastx
translated nucleotide ▶ protein

tblastn
protein ▶ translated nucleotide

Protein BLAST
protein ▶ protein

BLAST Genomes

Search

Human Mouse Rat Microbes

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

File formats

NIH > U.S. National Library of Medicine > NCBI

BLAST® » blastn suite

Home Recent Results Saved Strategies Help

Standard Nucleotide BLAST

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

Enter Query Sequence

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

[Reset page](#) [Bookmark](#)

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#)

[Clear](#)

Query subrange [?](#)

From

To

Or, upload file [Choisir le fichier](#) aucun fichier sél. [?](#)

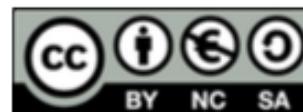
Job Title

Enter a descriptive title for your BLAST search [?](#)

Align two or more sequences [?](#)


H3ABioNet

Pan African Bioinformatics Network for H3Africa



File formats

NIH U.S. National Library of Medicine NCBI National Center for Biotechnology Information Sign in to NCBI

BLAST® » blastn suite Home Recent Results Saved Strategies Help

Microbial Nucleotide BLAST

blastn blastp blastx tblastn

Enter Query Sequence

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

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) Clear

Query subrange [?](#)

From To

Or, upload file Choisir le fichier aucun fichier sél. [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

Choose Search Set

Database Representative genomes only All genomes [?](#)

Title: Refseq prokaryote representative genomes (contains refseq assembly)
Molecule Type: mixed DNA
Update date: 2016/07/30
Number of sequences: 286522

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. [?](#)


H3ABioNet

Pan African Bioinformatics Network for H3Africa



<http://blast.ncbi.nlm.nih.gov/>
Introduction to Bioinformatics Online Course:IBT Genomics | Fatma Guerfali

File formats

NIH U.S. National Library of Medicine NCBI National Center for Biotechnology Information Sign in to NCBI

BLAST® > blastn suite Home Recent Results Saved Strategies Help

Microbial Nucleotide BLAST

blastn blastp blastx tblastn

Enter Query Sequence

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

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#) Query subrange [?](#)

From To

Or, upload file Choisir le fichier aucun fichier sél. [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

Choose Search Set

Database Representative genomes only All genomes [?](#)

Title: Prokaryote complete refseq genomes; Prokaryote complete genbank genomes without refseq
 Description: Complete prokaryote Genome Database
 Molecule Type: mixed DNA
 Update date: 2015/05/27
 Number of sequences: 9981

Organism Optional

Complete genomes
 Draft genomes
 Complete plasmids
 Complete bacteriophages

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. [?](#)

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



Part 4

Comparative Genomics: The methods



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics| Fatma Guerfali

Sequence Alignment for DNA

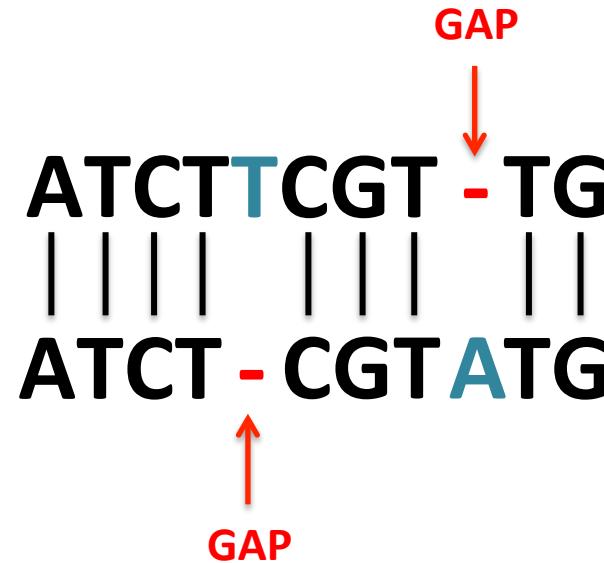
- DNA sequence alignment consists of aligning 2 DNA sequences in order to identify regions showing sequence similarity
- This highlights regions showing relationship in terms of:
 - Evolution
 - Structure
 - Function

Sequence Alignment for DNA



For simple sequences:

Compare **ATCTTCGTTG** and **ATCTCGTATG**



Sequence Alignment for DNA

- For sequences that are not as simple :
 - long sequences
 - more complex sequences (divergent...)
 - large number of sequences
 - ...

→ need Algorithms !

Sequence Alignment for DNA

- 2 key approaches

- **Global Alignment**

→ Optimizes the alignment to span the full length of sequences that are aligned.

ATCATTGCGTGTGACTGTG
A - - TT - G - TGAC - - TG

- **Local Alignment**

→ Optimizes the alignment to take into account regions of the highest similarity between divergent sequences.

ATCATTGCGTGTGACTGTG
- - - ATT - G - TGACTG - -

Sequence Alignment for DNA

- Algorithms efficiency and choice depends on the number of sequences to compare
- Pairwise Alignment**
Sequence alignment of **2** sequences
→ Output: function, structure, evolutionary relationship
- Multiple sequence Alignment**
Sequence alignment of **3 or more** sequences (same length)
→ Output: homology, evolutionary relationship

Sequence Alignment for DNA

Pairwise Alignment

- A **Pairwise Alignment** is an optimized local or global alignment of **2 sequences**.
 - **3 methods:**
 - Dot-matrix
 - Dynamic programming
 - Word-based

*NB: Efficiency can be reduced in low complexity regions (repetitive sequences...) Can be evaluated by the MUM (Maximum Unique Match)
→ Long MUM sequences = more related sequences*



Sequence Alignment for DNA

Pairwise Alignment

- Dot-matrix method

- 2 sequences (A and B) are aligned using a 2-dimensional matrix
- Identity is shown with a dot
- Diagonal shows high similarity

→ Dot plot of the sequence R against sequence Q

INSERTION INTO QUERY

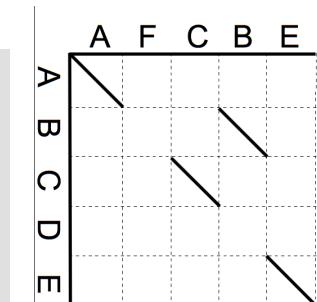
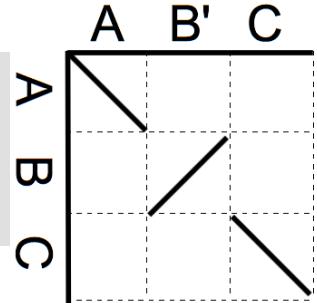
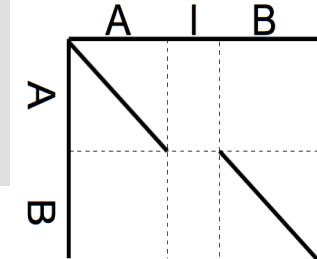
R: AB
Q: AIB

INVERSION INTO QUERY

R: ABC
Q: AB'C

REARRANGEMENT INTO QUERY

R: ABCDE
Q: AFCBE



Sequence Alignment for DNA

Pairwise Alignment

- Dynamic programming can

- Use a scoring matrix
- Assign a match score (+), a mismatch score (-), and a gap penalty (-).
- Use two different gap penalties for opening a gap and for extending a gap (gap opening >> gap extension)
 - generally results in less gaps in an alignment and gaps are grouped together = more biological relevance.

- Different algorithms for Global and Local Alignments

Sequence Alignment for DNA

Pairwise Alignment

- Dynamic programming

- **Global Alignment**

- Needleman–Wunsch algorithm.

- **Local Alignment**

- Smith–Waterman algorithm.

Sequence Alignment for DNA

Pairwise Alignment



Global Alignment

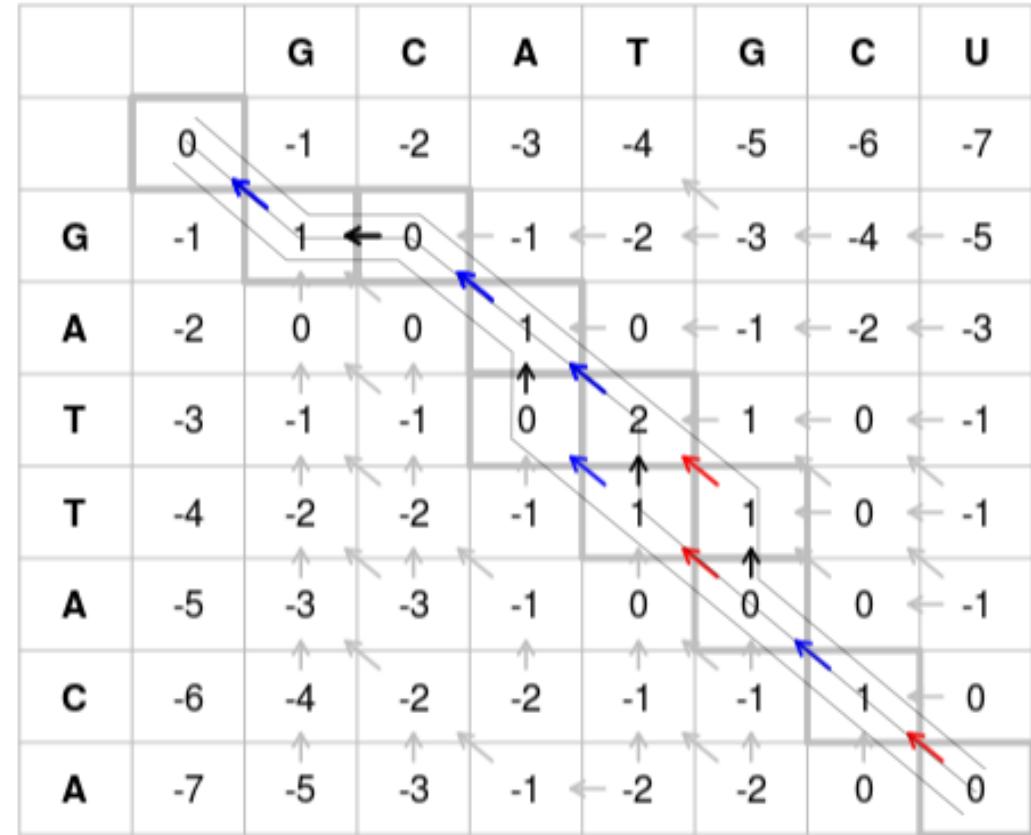
→ Needleman–Wunsch algorithm.

Needleman-Wunsch

match = 1

mismatch = -1

gap = -1



Sequence Alignment for DNA

Pairwise Alignment

- Local Alignment

→ Smith-Waterman algorithm.

| | | | | | | | | | |
|---|---|---|---|---|---|----|----|----|----|
| | - | A | C | A | C | A | C | T | A |
| - | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 2 | 1 | 2 | 1 | 2 | 1 | 1 | 2 |
| G | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | |
| C | 0 | 1 | 3 | 2 | 3 | 2 | 3 | 2 | 2 |
| A | 0 | 2 | 2 | 5 | 4 | 5 | 4 | 4 | 4 |
| C | 0 | 1 | 4 | 4 | 7 | 6 | 7 | 6 | 6 |
| A | 0 | 2 | 3 | 6 | 6 | 9 | 8 | 8 | 8 |
| C | 0 | 1 | 4 | 5 | 8 | 8 | 11 | 10 | 10 |
| A | 0 | 2 | 3 | 6 | 7 | 10 | 10 | 10 | 12 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | - | A | C | A | C | A | C | T | A |
| - | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ |
| G | 0 | ↑ | ↖ | ↑ | ↖ | ↖ | ↑ | ↖ | ↑ |
| C | 0 | ↑ | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ |
| A | 0 | ↖ | ↑ | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ |
| C | 0 | ↑ | ↖ | ↑ | ↖ | ↖ | ↖ | ↖ | ↖ |
| A | 0 | ↖ | ↑ | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ |
| C | 0 | ↑ | ↖ | ↑ | ↖ | ↖ | ↖ | ↖ | ↖ |
| A | 0 | ↖ | ↑ | ↖ | ↖ | ↖ | ↖ | ↖ | ↖ |

Sequence Alignment for DNA

Pairwise Alignment

Word-based method

- Optimal alignment not guaranteed, but efficient and faster than dynamic programming
- Useful for databases searches
- « words » are small portions (length k) of the query sequence that are used to screen the database

→ Ex: BLAST

Sequence Alignment for DNA

Pairwise Alignment

Word-based method

- BLAST (Basic Local Alignment Search Tool)
- Algorithm to compare a query sequence to a library or database of sequences
- Allows to estimate identity with a certain confidence threshold
- Popular in the scientific community (time efficiency...)

Sequence Alignment for DNA

Pairwise Alignment

- Word-based method

- BLAST

- ▶ The query sequence can be filtered to exclude low-complexity regions
- ▶ **Seeding:** list all possible words of length (DNA: default **k=11**)
- ▶ Search database for matching words using a **scoring matrix = calculate the match score**
- ▶ A **threshold score** is evaluated to top-rank the most similar sequences
- ▶ Process repeated for all words of the query

QUERY

ATCATTGCGTGTG

k=11

ATCATTGCGTGTG

TCATTGCGTGTG

CATTGCGTGTG

ATTCGTTGACTGTG

TTCGTTGACTGTG

TCGTTGACTGTG

CGTTGACTGTG



DATABASE ATCATTGCGTGTG

ATTCGTTGACTGTG

TCGTTGACTGTG

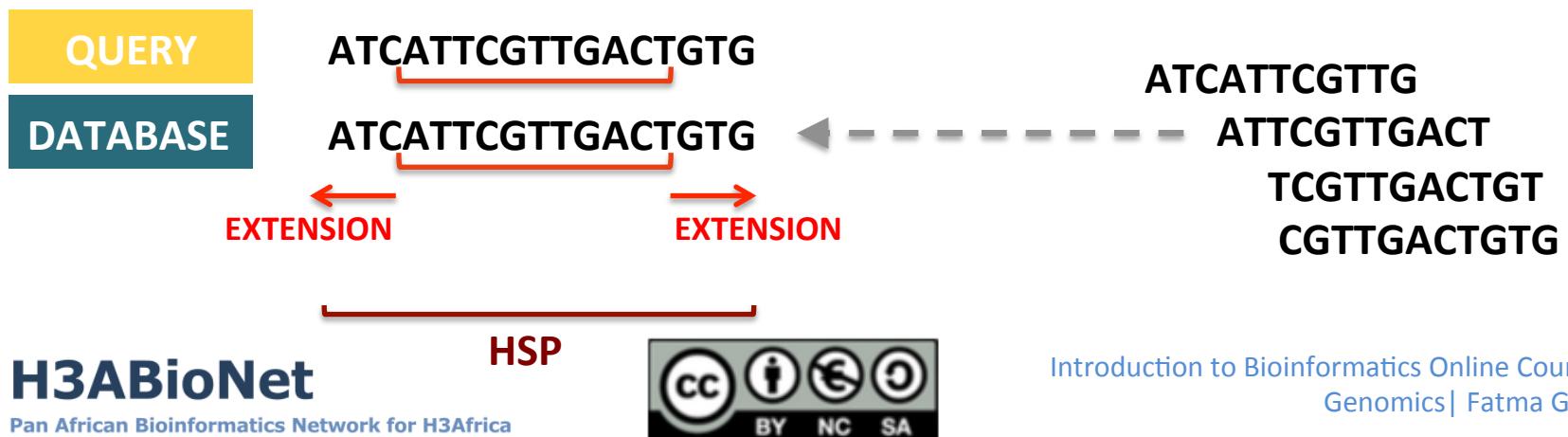
CGTTGACTGTG

Sequence Alignment for DNA

Pairwise Alignment

- Word-based method

- BLAST
 - Best match is followed by an **extension in both direction, with scoring**
 - Extension continued only if the alignment is above the threshold
 - The contiguous alignment without gaps (now possible) and a higher score is the **HSP** (High Scoring Segment Pair)



Sequence Alignment for DNA

Pairwise Alignment

- Word-based method

- BLAST

- A scoring matrix is used to evaluate the quality of the alignment
- A scoring matrix is a predefined substitution matrix (match = 1, mismatch = 0...)
- ex: BLOSUM

Sequence Alignment for DNA

NIH > U.S. National Library of Medicine > NCBI National Center for Biotechnology Information Sign in to NCBI

BLAST® » blastn suite

Home Recent Results Saved Strategies Help

Standard Nucleotide BLAST

blastn blastp blastx tblastn tblastx

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) BLASTN programs search nucleotide databases using a nucleotide query. [more...](#)

Query subrange

Or, upload file aucun fichier sél.

Job Title NG_029730:Homo sapiens CD74 molecule (CD74),... 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.): Human genomic plus transcript (Human G+T) (176974 sequences)

Exclude Models (XM/XP) Uncultured/environmental sample sequences

Limit to Sequences from type material [YouTube](#) [Create custom database](#)

Entrez Query 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

BLAST Search database Human G+T using Megablast (Optimize for highly similar sequences) Show results in a new window

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


H3ABioNet

Pan African Bioinformatics Network for H3Africa


 Introduction to Bioinformatics Online Course:IBT Genomics | Fatma Guerfali
<http://blast.ncbi.nlm.nih.gov/>

Sequence Alignment for DNA

Algorithm parameters

marked with ♦ sign

Note: Parameter values that differ from the default are highlighted in yellow and

[Restore default search parameters](#)

General Parameters

| | | |
|-------------------------------------|---|---|
| Max target sequences | <input type="text" value="100"/> ♦ | Select the maximum number of aligned sequences to display  |
| Short queries | <input checked="" type="checkbox"/> Automatically adjust parameters for short input sequences  | |
| Expect threshold | <input type="text" value="10"/>  | |
| Word size | <input type="text" value="28"/>  | |
| Max matches in a query range | <input type="text" value="0"/>  | |

Scoring Parameters

| | |
|------------------------------|---|
| Match/Mismatch Scores | <input type="text" value="1,-2"/>  |
| Gap Costs | <input type="text" value="Linear"/>  |

Filters and Masking

| | |
|---------------|--|
| Filter | <input checked="" type="checkbox"/> Low complexity regions  |
| | <input type="checkbox"/> Species-specific repeats for: <input type="text" value="Homo sapiens (Human)"/>  |
| Mask | <input checked="" type="checkbox"/> Mask for lookup table only  |
| | <input type="checkbox"/> Mask lower case letters  |


H3ABioNet

Pan African Bioinformatics Network for H3Africa

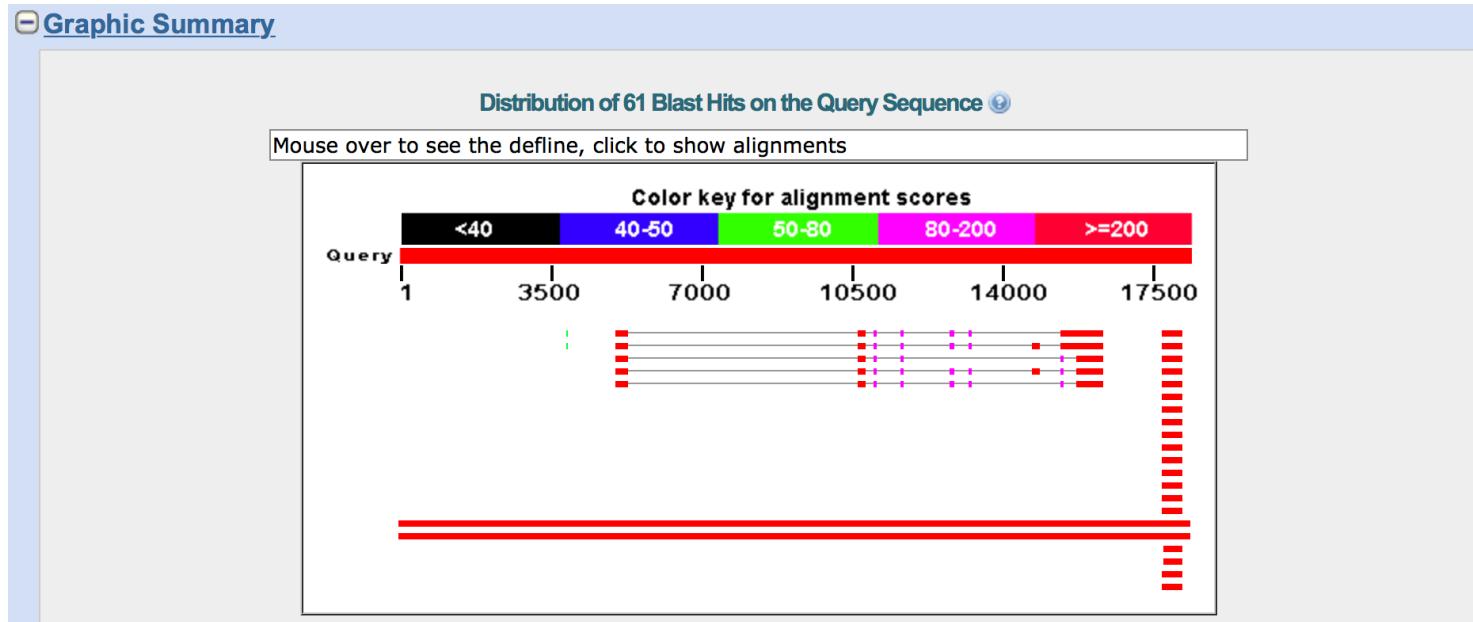

 Introduction to Bioinformatics Online Course:IBT
 Genomics | Fatma Guerfali
<http://blast.ncbi.nlm.nih.gov/>

Sequence Alignment for DNA

Pairwise Alignment

- Word-based method

- BLAST output
- A list of sequences that have the best match to the query

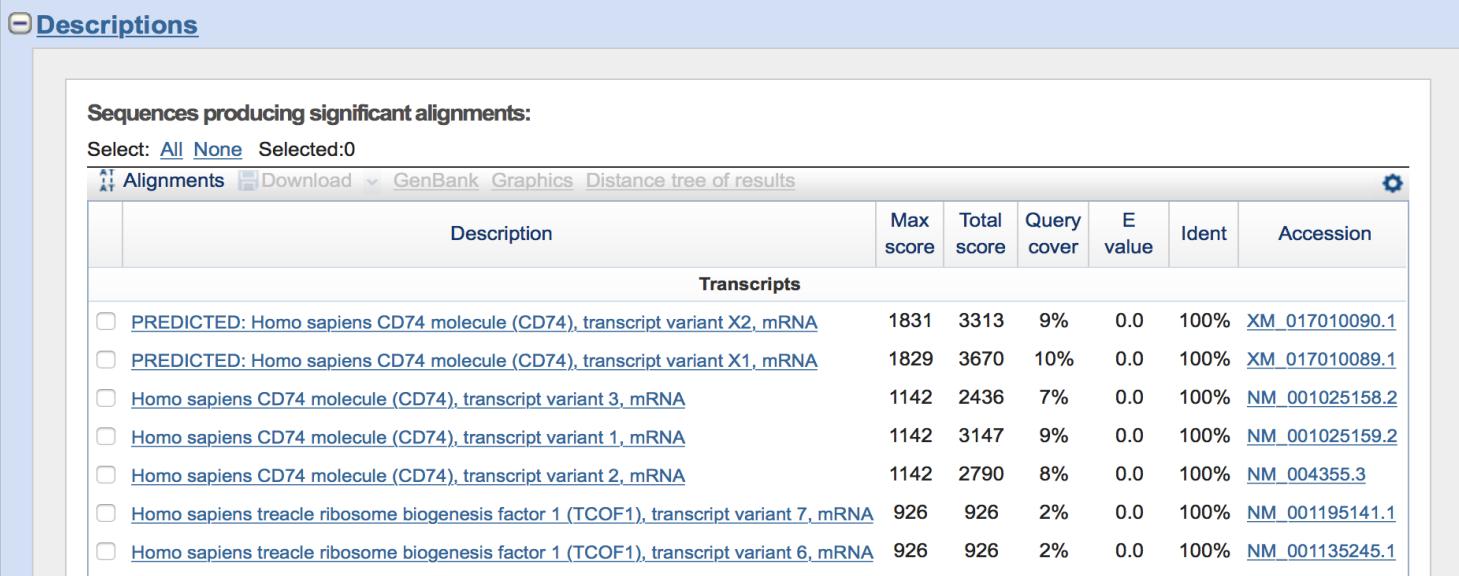


Sequence Alignment for DNA

Pairwise Alignment

- Word-based method

- BLAST output
- e-value: probability that the alignment is found by chance
(the lower the e-value, the more interesting the match)



The screenshot shows a 'Descriptions' tab selected in a BLAST search interface. The main panel displays a table of sequences producing significant alignments. The columns include Description, Max score, Total score, Query cover, E value, Ident, and Accession. The table lists various transcripts for the Homo sapiens CD74 gene, along with some entries for TCOF1.

| | Description | Max score | Total score | Query cover | E value | Ident | Accession |
|--------------------------|---|-----------|-------------|-------------|---------|-------|----------------|
| Transcripts | | | | | | | |
| <input type="checkbox"/> | PREDICTED: Homo sapiens CD74 molecule (CD74), transcript variant X2, mRNA | 1831 | 3313 | 9% | 0.0 | 100% | XM_017010090.1 |
| <input type="checkbox"/> | PREDICTED: Homo sapiens CD74 molecule (CD74), transcript variant X1, mRNA | 1829 | 3670 | 10% | 0.0 | 100% | XM_017010089.1 |
| <input type="checkbox"/> | Homo sapiens CD74 molecule (CD74), transcript variant 3, mRNA | 1142 | 2436 | 7% | 0.0 | 100% | NM_001025158.2 |
| <input type="checkbox"/> | Homo sapiens CD74 molecule (CD74), transcript variant 1, mRNA | 1142 | 3147 | 9% | 0.0 | 100% | NM_001025159.2 |
| <input type="checkbox"/> | Homo sapiens CD74 molecule (CD74), transcript variant 2, mRNA | 1142 | 2790 | 8% | 0.0 | 100% | NM_004355.3 |
| <input type="checkbox"/> | Homo sapiens treacle ribosome biogenesis factor 1 (TCOF1), transcript variant 7, mRNA | 926 | 926 | 2% | 0.0 | 100% | NM_001195141.1 |
| <input type="checkbox"/> | Homo sapiens treacle ribosome biogenesis factor 1 (TCOF1), transcript variant 6, mRNA | 926 | 926 | 2% | 0.0 | 100% | NM_001135245.1 |

Sequence Alignment for DNA

Pairwise Alignment

- Word-based method

- BLAST output
- Alignment details : sequences (query and database) aligned with % identity...)

[Alignments](#)

Download ▾ GenBank Graphics Sort by: E value ▾ Next ▾ Previous ▾ Descriptions

PREDICTED: Homo sapiens CD74 molecule (CD74), transcript variant X2, mRNA
 Sequence ID: [ref|XM_017010090.1](#) Length: 1777 Number of Matches: 7

| Range 1: 787 to 1777 GenBank Graphics | | | | | ▼ Next Match | ▲ Previous Match | Related Information |
|---------------------------------------|--|---------------|-----------|-----------|--------------|------------------|---------------------|
| Score | Expect | Identities | Gaps | Strand | | | |
| 1831 bits(991) | 0.0 | 991/991(100%) | 0/991(0%) | Plus/Plus | | | |
| Query 15310 | AGAGTCACTGGAACTGGAGGACCCGTCTTCTGGGCTGGGTGTGACCAAGCAGGATCTGGG | | | | 15369 | | |
| Sbjct 787 | AGAGTCACTGGAACTGGAGGACCCGTCTTCTGGGCTGGGTGTGACCAAGCAGGATCTGGG | | | | 846 | | |
| Query 15370 | CCCAGGTAAGGGCCTTGCAAGAGGGCATCTGGTCACCAGCAGCTCATCCCCAGCAGGGCC | | | | 15429 | | |
| Sbjct 847 | CCCAGGTAAGGGCCTTGCAAGAGGGCATCTGGTCACCAGCAGCTCATCCCCAGCAGGGCC | | | | 906 | | |
| Query 15430 | AGCTCCTTGTGGGCAGGTGAAGGAGTGTGACGCTGGGCCACTCTCAAACATTCTGGGA | | | | 15489 | | |
| Sbjct 907 | AGCTCCTTGTGGGCAGGTGAAGGAGTGTGACGCTGGGCCACTCTCAAACATTCTGGGA | | | | 966 | | |

Sequence Alignment for DNA

Pairwise Alignment

Word-based method

- BLAST have different variant queries according to the type of query sequence (**Q**) and type of sequence in the database (**R**):

| | Q | | R |
|---------|-------------------------|---|-------------------------|
| BLASTN | Nucleic Acid | → | Nucleic Acid |
| BLASTX | Translated Nucleic Acid | → | Protein |
| TBLASTX | Translated Nucleic Acid | → | Translated Nucleic Acid |
| TBLASTN | Protein | → | Translated Nucleic Acid |
| BLASTP | Protein | → | Protein |

Sequence Alignment for DNA

Multiple Sequence Alignment

- **Multiple Sequence Alignment** have been developed to handle more than 2 sequences at a time.
- Align all queried sequences to form a query group.
- Allows to identify **conserved sequences** portions among a group of queried sequences that are:
 - known to be evolutionarily related
 - of unknown/supposed evolutionar relationship → this multiple alignment helps to establish their evolutionar relationships (phylogenetic trees)



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Adapted from different web resources

Introduction to Bioinformatics Online Course:IBT
Genomics| Fatma Guerfali

Sequence Alignment for DNA

Multiple Sequence Alignment

- Different methods:

- Dynamic programming
- Progressive method
- Iterative method : HMMs (Hidden Markov Models)
- ...

→ Evaluation of Conservation across sequences

Sequence Alignment for DNA

Multiple Sequence Alignment

● Dynamic programming

- Optimized for 2 sequences, so computationally expensive here
- Extends the sequence matrix from 2 sequences to the number of sequences in the query : alignment between pairs of sequences

→ MSA



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Adapted from different web resources

Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Sequence Alignment for DNA

Multiple Sequence Alignment

● Progressive method

- Aligns sequences to identify the most similar ones
- Then progressively adds all other related sequences of the group

→ Clustal (clustal-Omega: medium-large alignments)

→ T-Coffee (small alignments)

Adapted from different web resources

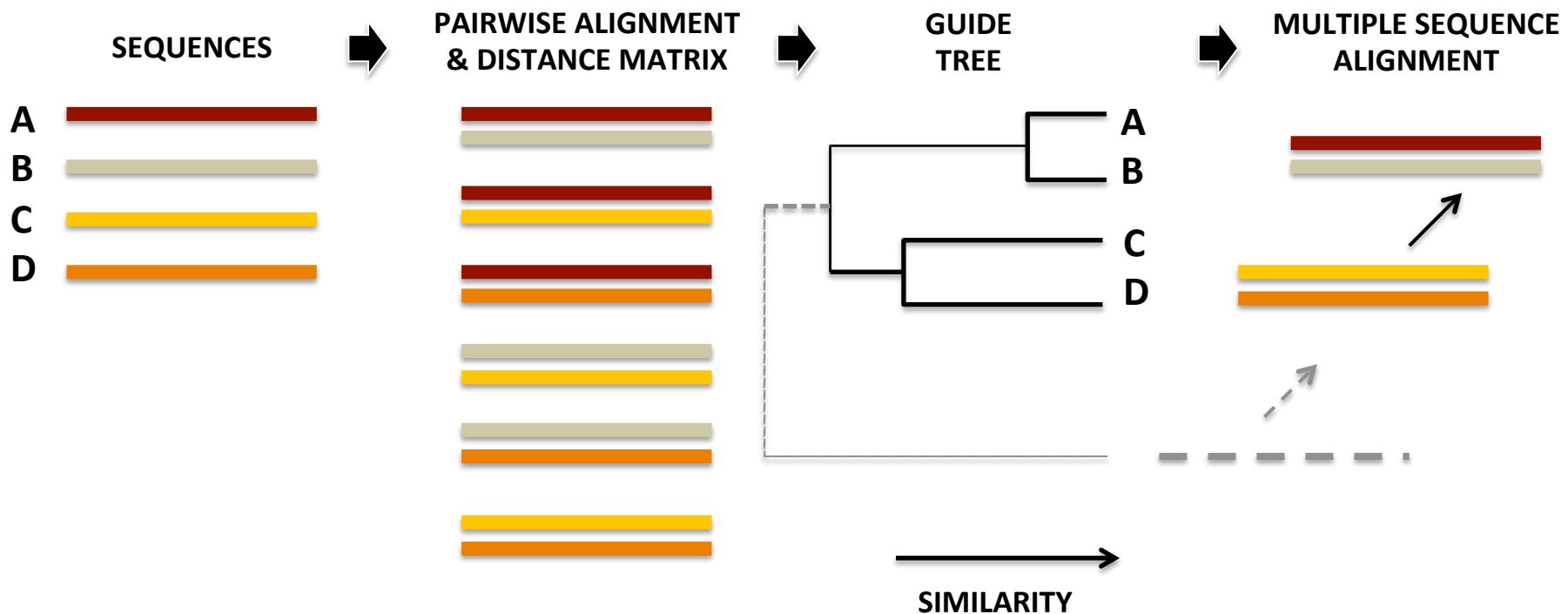
Introduction to Bioinformatics Online Course:IBT
Genomics| Fatma Guerfali

Sequence Alignment for DNA

Multiple Sequence Alignment

- Progressive method

Clustal

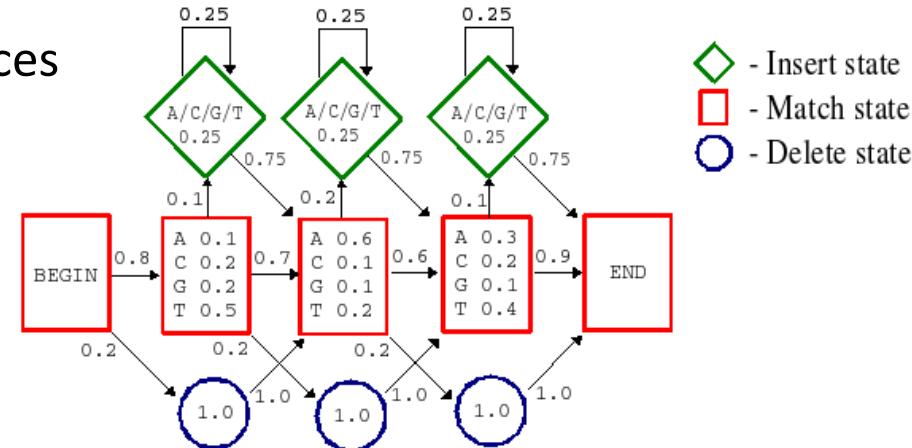


Sequence Alignment for DNA

Multiple Sequence Alignment

HMMs

- Determine Probability scores for multiple sequence alignments
- The aligned sequences serve as a group, no need to previously order the sequences
- Can build an HMM profile
- Improved for more distant sequences



<http://www.cbs.dtu.dk>

Introduction to Bioinformatics Online Course:IBT
Genomics| Fatma Guerfali

Sequence Alignment for DNA

Multiple Sequence Alignment

● Conservation scores

- Based on the fact that the highest conservation is maintained through evolution for the most important functions (promoters, essential enzymes, exons...)
- Regulatory regions might be generally evolving « faster »
- Multiple alignments → identify what **elements reject substitutions** (substitutions occur in neutral DNA, do not occur if an **element is functionally constrained**)

Sequence Alignment for DNA

Multiple Sequence Alignment

● Conservation scores

- **PhyloP** (Phylogenetic p-values)

Measures Base Conservation from non-coding regions

- **PhastCons** (part of PHAST: PHylogenetic Analysis

with Space/Time models)

Measures Base Conservation based on HMM model

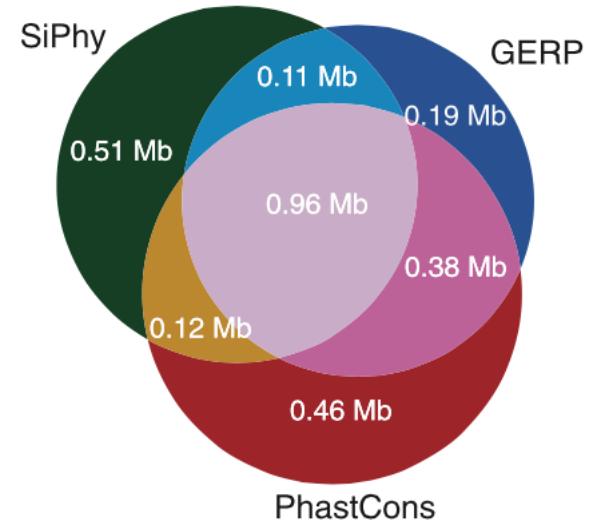
- **GERP** (Genomic Evolutionary Rate Profiling)

Measures Base Conservation to estimate the neutral evolution rate in genomes

- **SiPhy** (SIte-specific PHYlogenetic analysis)

Models the pattern of substitution (biased nucleotide substitution, HMM)

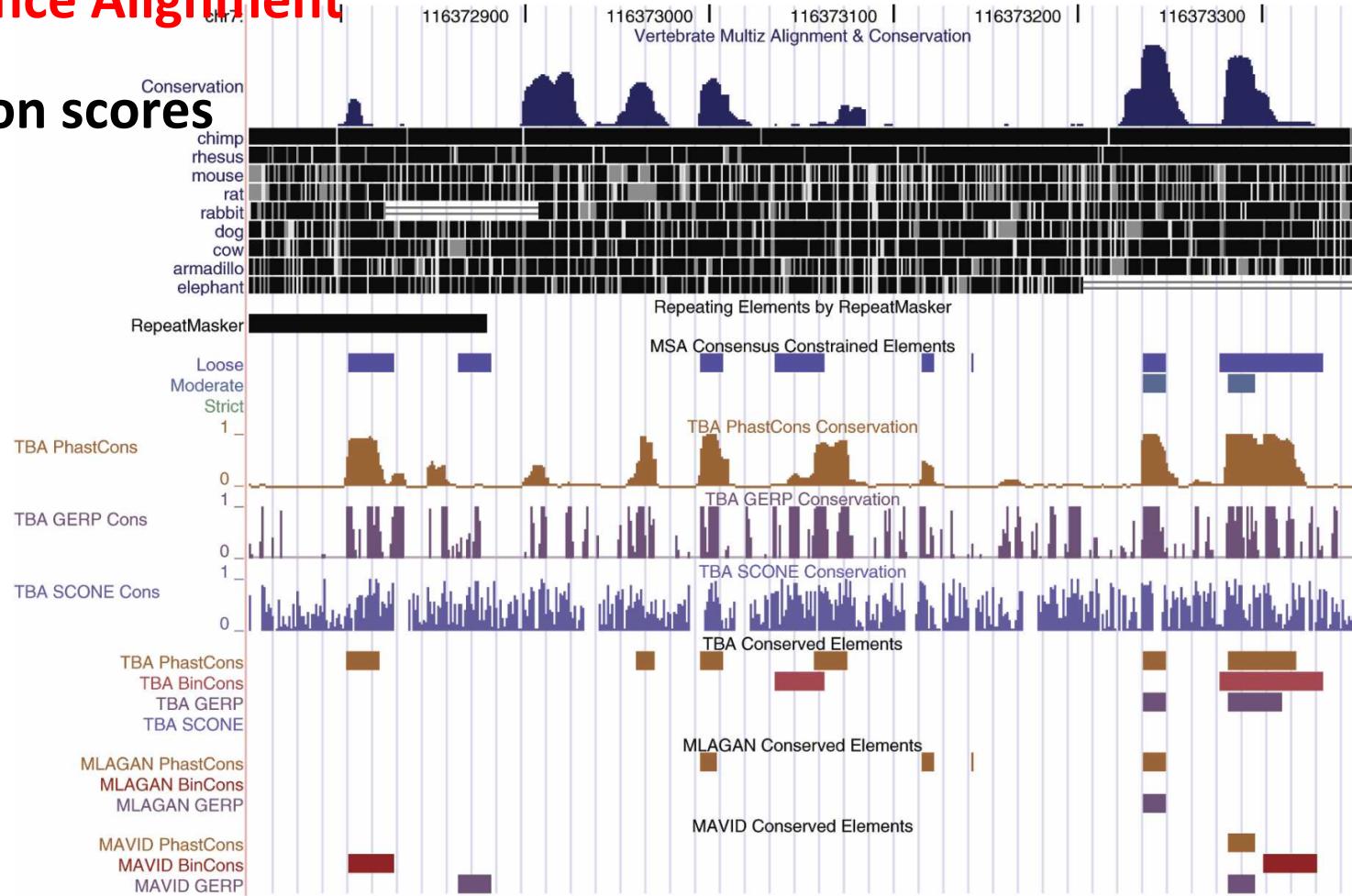
Base Overlap



Sequence Alignment for DNA

Multiple Sequence Alignment

● Conservation scores



Sequence Alignment for DNA

Whole Genome Alignment

- **MUMmer Ultra-fast alignment of large-scale DNA**

(<http://mummer.sourceforge.net>)

- Alignment of entire genomes (complete or draft)
- **Maximal Unique Matcher:** Find the MUMs = subsequences that occur only once in both genomes compared and not extendable anymore
- Suffix-tree based approach

- **WebACT Artemis Comparison Tool (ACT)**

(www.webact.org)

Visualize the alignment of publically available prokaryotic genomes

- **Databases: Ensembl, VISTA**



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Part 5

Comparative Genomics: The output

Genomic Structure

- Analyzing a genome structure means to analyze:

- ▶ **At the Genome level**

- Base composition (%GC, codon bias, nucleotide distribution...)
- Genome organization (SV events, genomic rearrangements...)
- Sequence conservation (regulatory elements, repetitive regions...)
- Synteny (conserved or not)

- ▶ **At the Gene level**

- Gene order

Genomic Function

- Analyzing functions in a genome means to analyze:

- ▶ At the non-coding sequence level

- Regulatory functions...

- ▶ At the coding sequence level

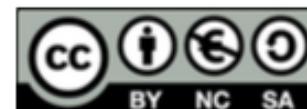
- Insights into functions

- Compare gene sequences
 - Compare protein-coding portions

- How ?

Different algorithms help identify portions of the genome coding for proteins

- *Ab initio* approaches
 - Using homology ...



<http://www.cbs.dtu.dk>

Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Genomic Evolution

- From the Multiple sequence alignments we can infer homology and estimate the evolution distance between sequences/organisms.

- Analysis could be:
 - Based on entire **Genome** comparisons
 - Based on **Gene** comparisons
 - ...

Genomic Evolution

● Based on entire genome comparison

- Phylogenetic relationship between organisms
- Previous dogma:

"Anything found to be true of E. coli must also be true of elephants"
(Jacques Monod, 1954)

- Need to be related to its phylogenetic context !!
- Different outputs expected depending if :
 - Comparing closely related species
 - Comparing evolutionarily distant species



Genomic Evolution: in easy words

COMPARING CLOSELY
RELATED SPECIES

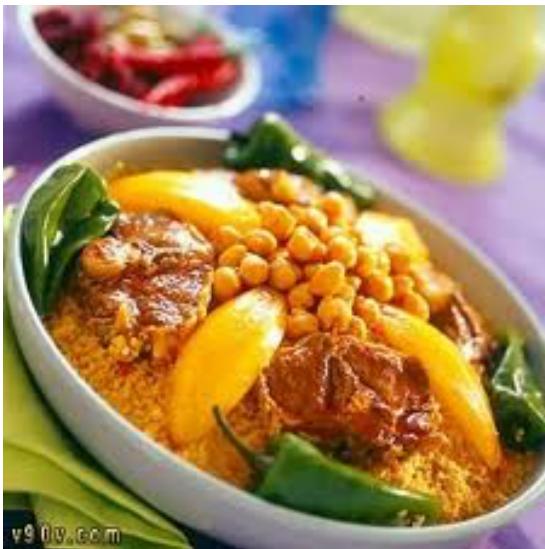


COMPARING EVOLUTIONARILY
DISTANT SPECIES



Comparative Genomics and Genome Evolution in easy words

COMPARING CLOSELY RELATED SPECIES



A: NOT SPICY



B: SPICY

COMPARING EVOLUTIONARILY DISTANT SPECIES



C: NOT SPICY

Adapted from <http://recettes-aymen.over-blog.fr/>
<http://quebueno.be/content/6>

Introduction to Bioinformatics Online Course:IBT
Genomics| Fatma Guerfali

Genomic Evolution

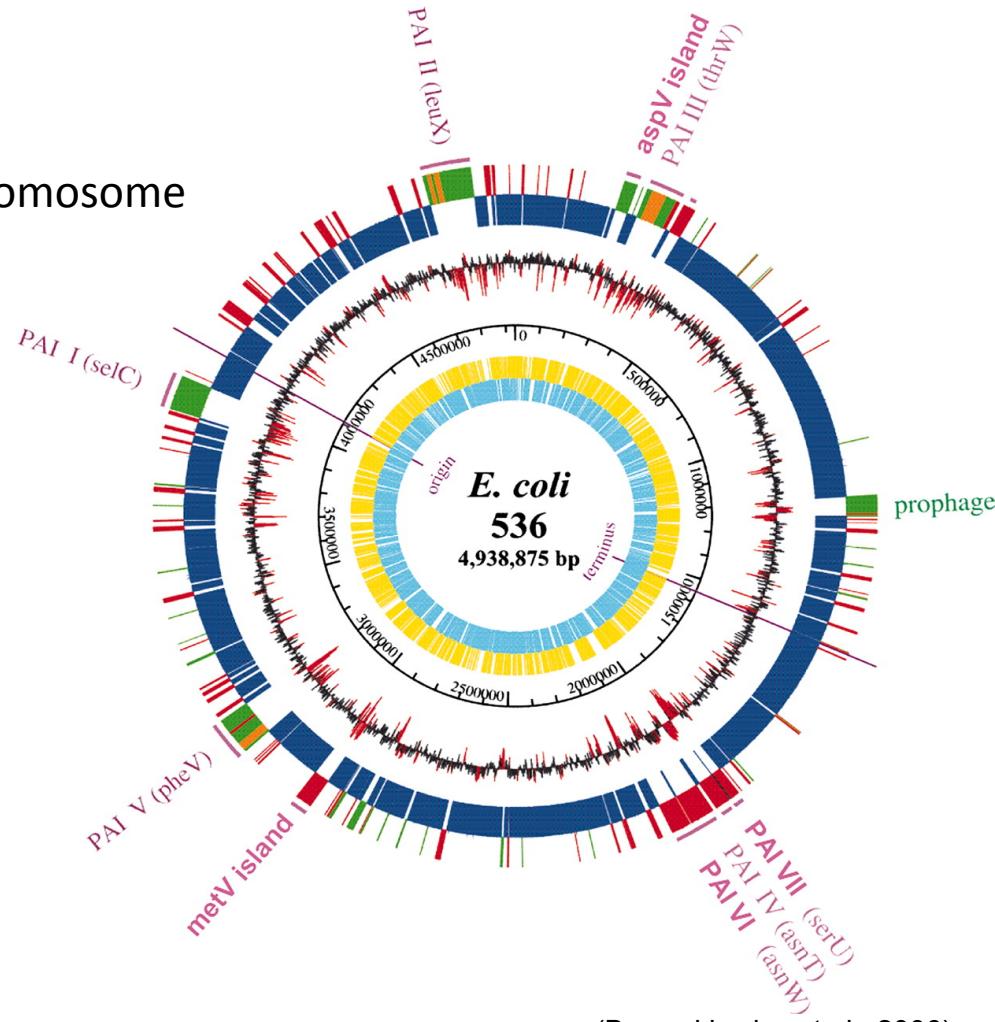
Based on entire genome comparison

Examining the dynamics of closely related genomes helps to build therapeutic strategies for Ebola virus:

- Ebola virus largest outbreak (2014)
 - Comparison of 100 available Ebolavirus (Filoviridae) genomes to each other + to other viral genomes.
 - Filoviridae are different from all other viral genomes
 - Filovirus genomes : sequence diversity but proteins with similar functions and gene order
- Ebolavirus genomes very similar but different in intergenic regions and genes of specific function = potential vaccine candidates.

Genomic Structure

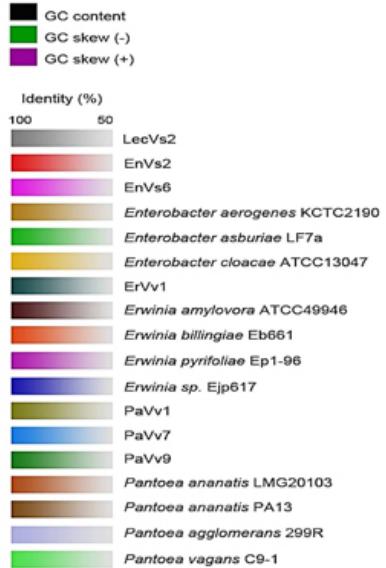
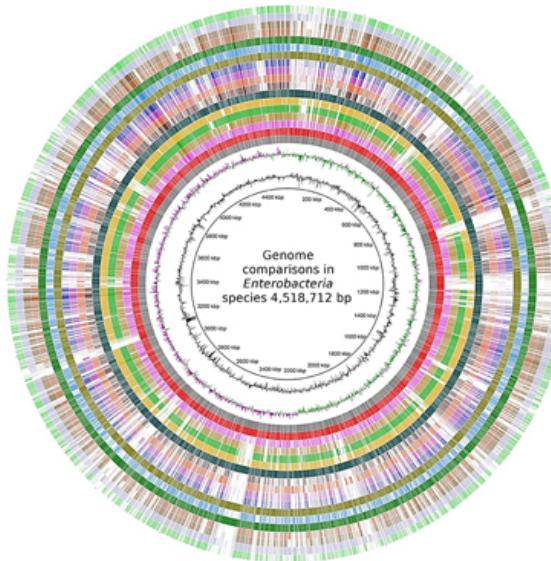
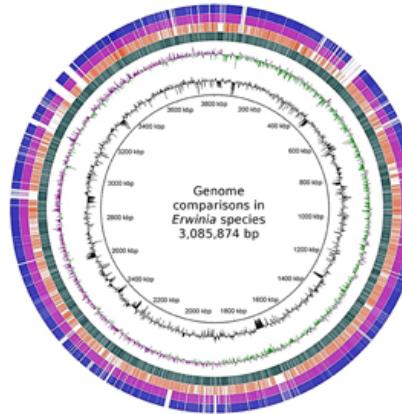
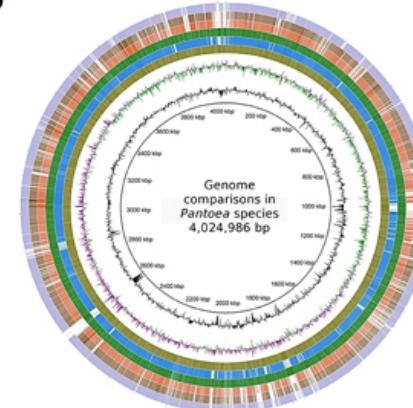
- Genetic map of the UPEC strain 536 chromosome



(Brzuszkiewicz et al., 2006)

Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Genomic Structure

**A****B****C****D**

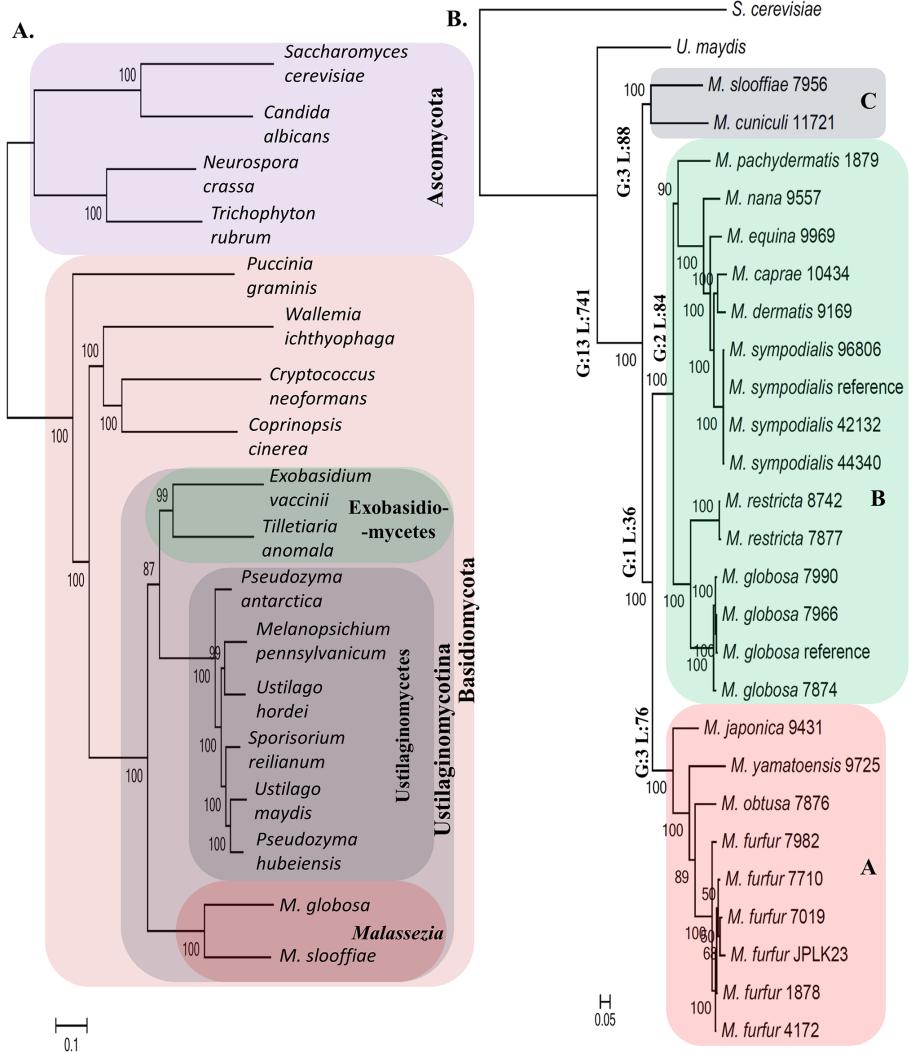
Whole Genome Comparison of
(A) all strains
(B) Enterobacter
(C) Erwinia
(D) Pantoea

Genomic Evolution

Phylogenetic relationships

« relationship of the *Malassezia* genus with respect to other fungi with sequenced genomes. »

G: gene family gain; L: gene family loss



(Wu et al., 2015)

https://en.wikipedia.org/wiki/Sequence_alignment#cite_note-mount-1
Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

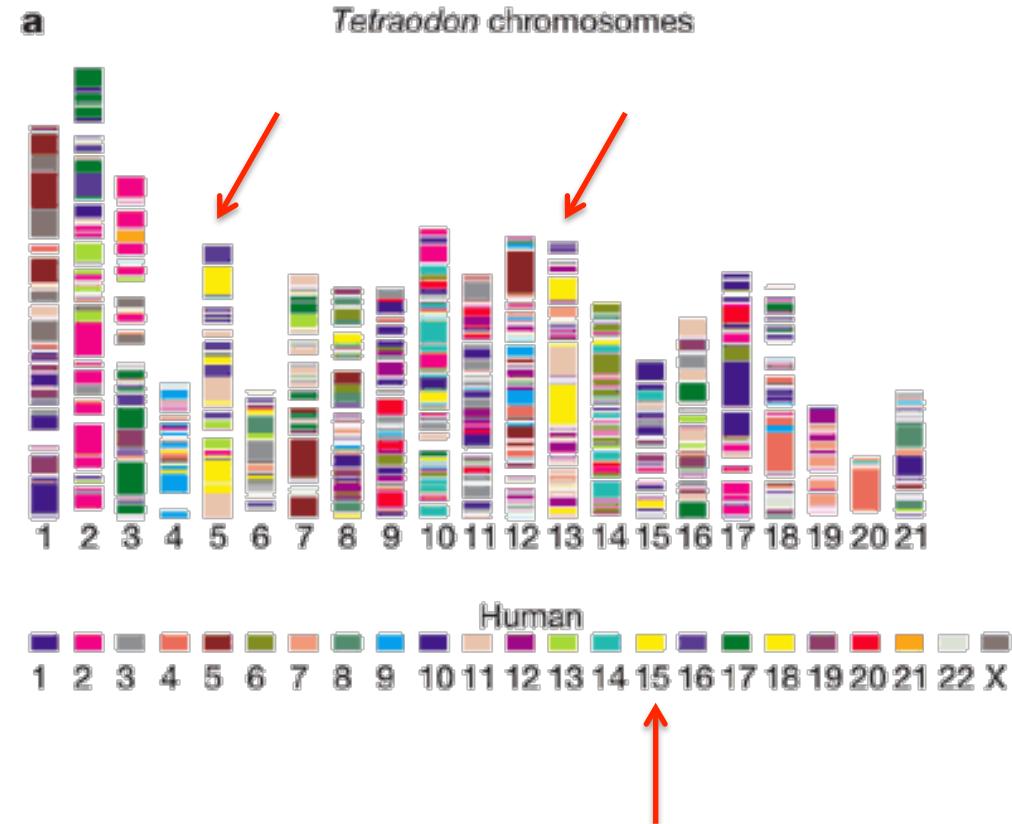
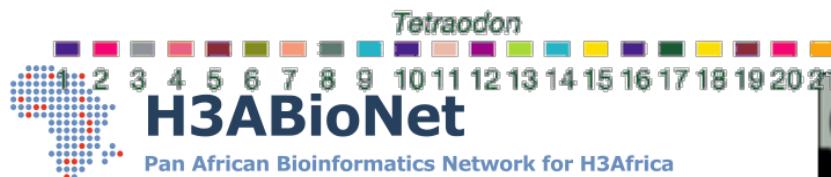
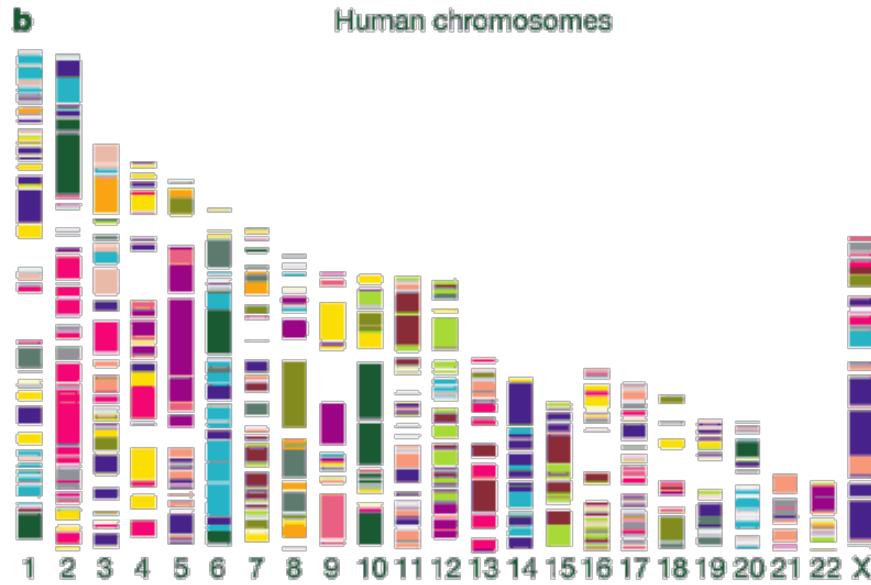
Genomic Evolution

- Based on entire genome comparison

- Synteny
 - ▶ Defined as the overall conservation of (gene/blocks) order in chromosomes between different genomes.
 - ▶ Evaluated in whole genomes, blocks could include large portions of genomes.
 - ▶ Recombination / crossing over affects groups of adjacent genes in a chromosome → **linkage group**.

Genomic Evolution

- Based on entire genome comparison
- Synteny

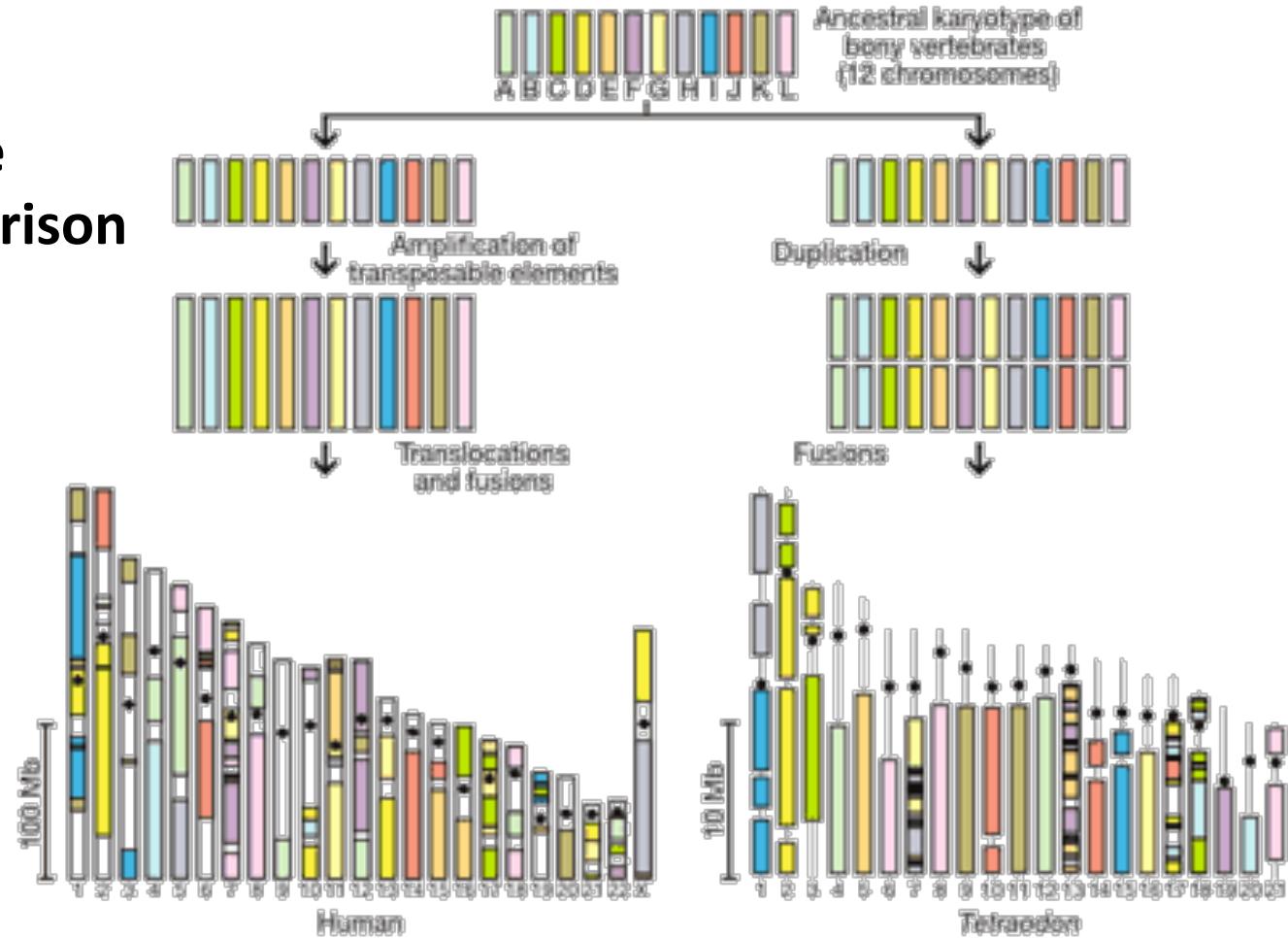


(Jaillon et al., 2004)
Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali



Genomic Evolution

- Based on entire genome comparison
- Synteny



(Jaiillon et al., 2004)

Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Genomic Evolution

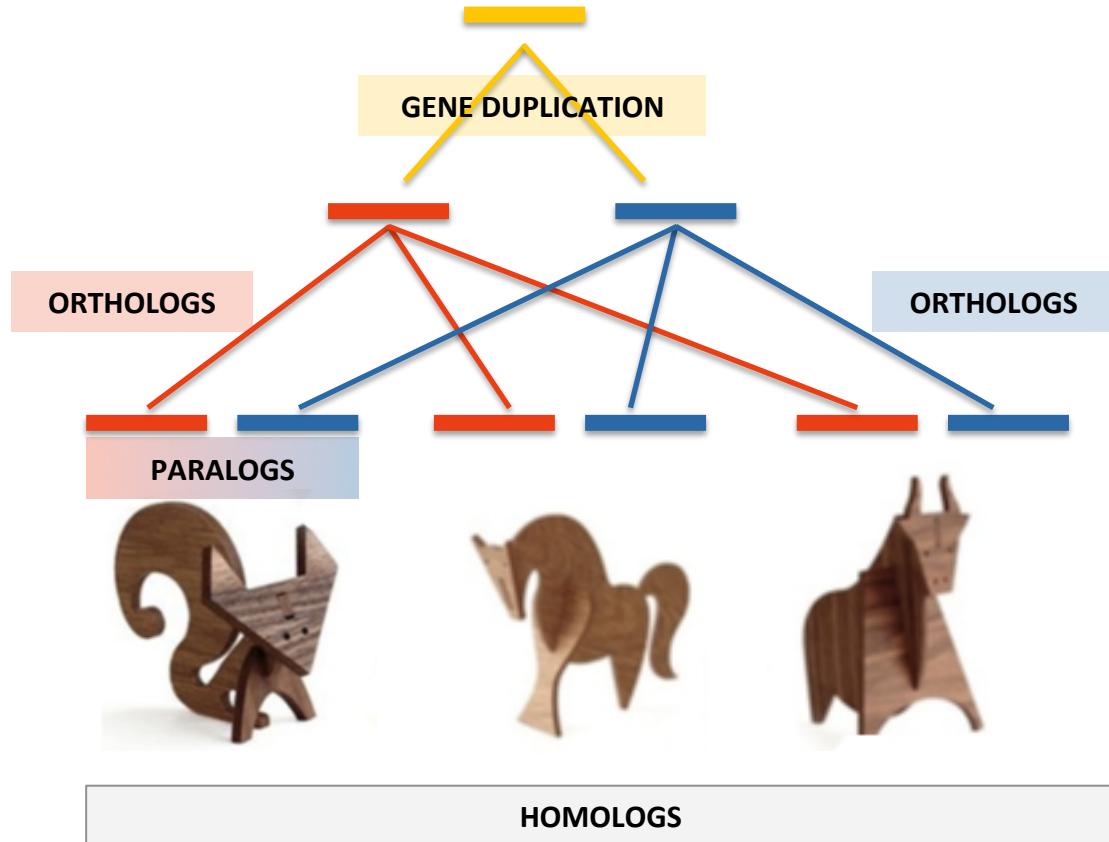
- Based on individual gene comparison

Homology

2 genes are homologs if they have a common ancestor

They can be classified in orthologs and paralogs:

- As a consequence of speciation = **Orthology**
- As a consequence of duplication = **Paralogy**



Genomic Evolution

- Based on individual gene comparison

Orthology

Finding orthologs can be the first step in whole genome alignment

- BLAST Reciprocal Best Hit (best pairs of orthologs)
- OrthoMCL (possible predictions for several species)
- EnsemblCompara (precomputed data): orthology and paralogy predictions based on phylogenies.
- eggNOG (evolutionary genealogy of genes: Non-supervised Orthologous Groups): OGs of proteins across different taxonomic levels



Genomic Evolution

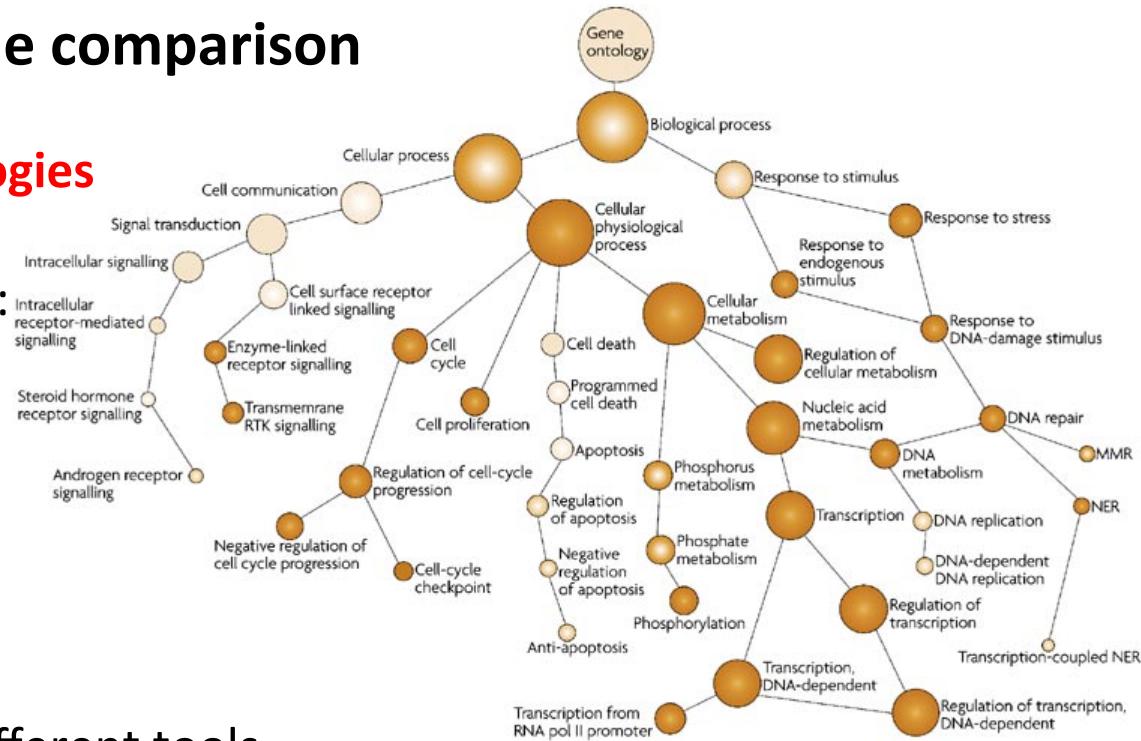
- Based on individual gene comparison

Inferring Annotation: Ontologies

Assigning functions to Genes:
GO (Gene Ontology)

- Biological Process
- Molecular Function
- Cellular Component

The GO database contains different tools to retrieve these informations.



Nature Reviews | Cancer

(Hu et al., 2007)

<http://geneontology.org/>

Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Genomic Evolution

● Functional/Structural Predictions

► impact of a mutation on the function

Analysis of the impact of aa substitution

- Structural and/or Functional effect of single point mutations SNPs
 - PolyPhen-2 ([http://genetics.bwh.harvard.edu/pph2/...](http://genetics.bwh.harvard.edu/pph2/))
 - SIFT (<http://sift.jcvi.org>)
 - VEP (http://www.ensembl.org/Homo_sapiens/Tools/VEP)

► impact of the mutation on the Struture

Impact on gene coding portions (gain/loss) or non-coding portions

Browsers

VISTA

(<http://genome.lbl.gov/vista/index.shtml>)

Collection of resources for comparative genomics

- VISTA browsers can be used to analyze pre-computed alignments or user generated or queried sequences

VISTA servers

- **mVISTA** (query sequences vs multi-species sequences)
- **rVISTA** (identification of regulatory TF binding sites)
- **gVISTA** (query sequences vs whole-genome assemblies)
- **wgVISTA** (alignment of 10Mb sequences (finished/draft): microbes...)
- ...

VISTA tools: Rviewer region viewer to compare genomic intervals



The screenshot shows the VISTA homepage. At the top, there's a banner with the text "Tools for Comparative Genomics". Below the banner, there's a navigation bar with links: "VISTA Home", "Custom Alignment", "Browser", "Enhancer DB", "Downloads", and "Publications". The main content area has two main sections. On the left, under "Submit Your Sequences", there's a link to "mVISTA" with a blue button containing the sequence "GACAC" followed by four short vertical lines and "GACAT". On the right, under "Precomputed Alignments", there's a link to "VISTA Browser" with a red button featuring a wavy line graph.

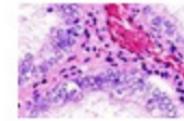
Browsers

Ensembl Browser

(<http://www.ensembl.org>)

- ▶ Comparative analyses at the genome and gene levels
- ▶ Genome sequences compared using pairwise and multiple whole-genome alignments
- ▶ These alignments help to determine
 - Synteny
 - Sequence conservation scores
 - Gene homology relationships (GeneTrees)

Gene expression in different tissues



Find SNPs and other variants for my gene

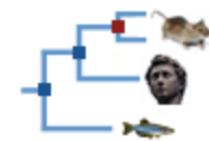


Retrieve gene sequence

```

GCCTGACTTCCGGGTGG
GGGCTTGTGGCGCGAGC
GCCCTCTGCTGCGCCTN
AOGGGACAGATTCTGA
CACCTCTGGAGCGGTTT
CCCAGTCCAGCGTGGCG
  
```

Compare genes across species



Use my own data in Ensembl



ENCODE data in Ensembl



(Herrero et al., 2015)

Introduction to Bioinformatics Online Course:IBT
Genomics | Fatma Guerfali

Comparative genomics

Take-home messages

● Input / Output

- ▶ DNA Sequences (genome, gene...)
- ▶ Homology, similarity, evolutionary distance

● Alignment

- ▶ Whole genome : MUMmer...
- ▶ Multiple genomes : MGA...
- ▶ Multiple Sequence Alignment : Clustal...
- ▶ Global/Local Sequence Alignment : BLAST...

● Input / Output files

- ▶ Fasta/GenBank to alignment or phylogenetic distances