# IBT 432 Aplikasi Bioinformatika
## Analisis komparasi genomika II: Pengenalan Genomika Komparatif

**Riza Arief Putranto**

# Rencana Perkuliahan

1. ~~Kontrak belajar dan pengenalan bioinformatika aplikatif~~
2. ~~Database sekuen dan analisis genomika~~
3. ~~Anotasi sekuen ke genom - Praktik~~
4. ~~Analisis komparasi genomika I~~
5. Analisis komparasi genomika II
6. Analisis komparasi genomika III
7. Analisis komparasi genomika – Praktik
8. Protein modelling I
9. Protein modelling II
10. Protein modelling III
11. Protein modelling - Praktik
12. Visualisasi protein modelling
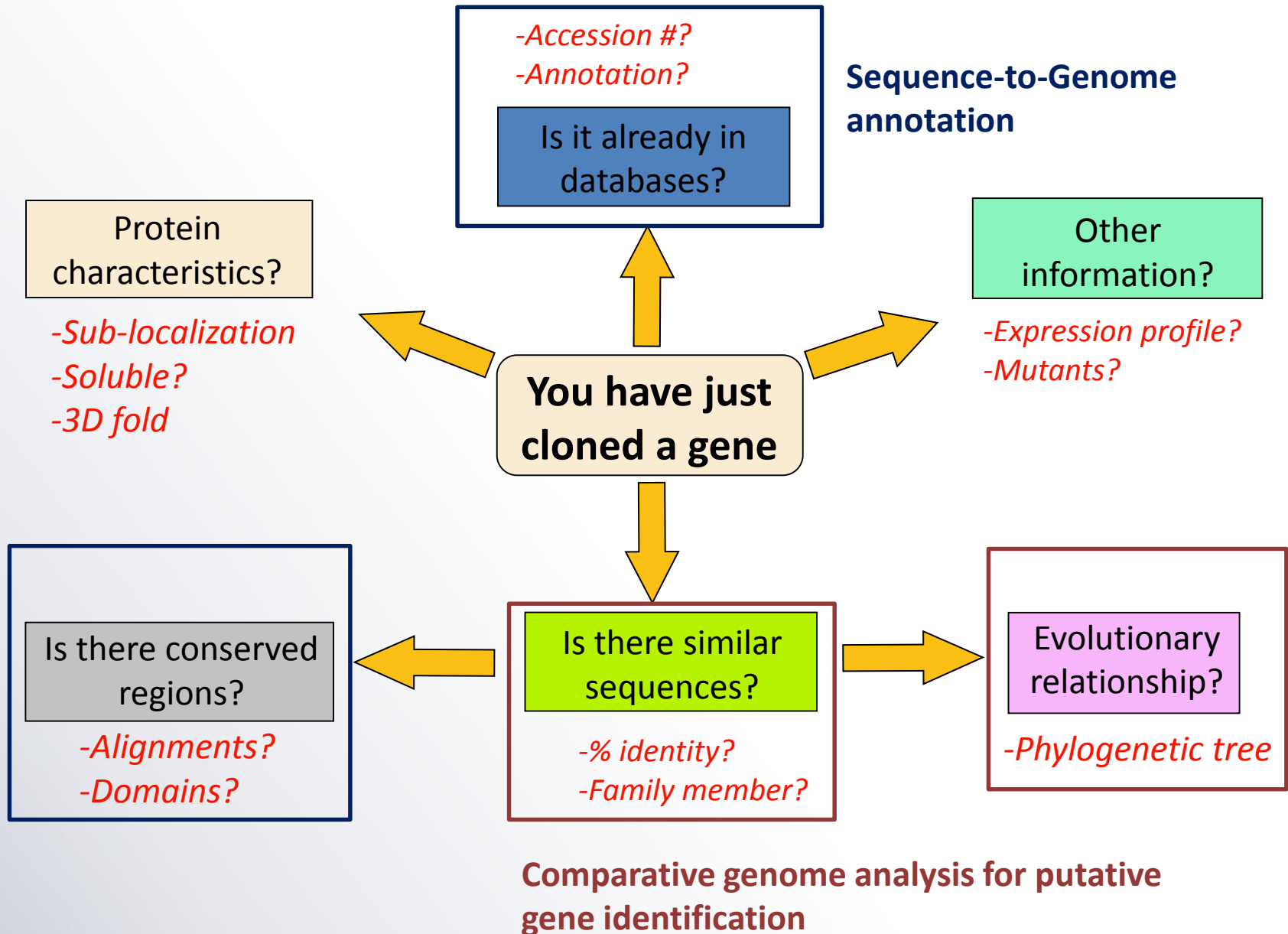13. Visualisasi protein modelling - Praktik
14. Presentasi mahasiswa

# Contents

- ❑ **History of HG Project**
- ❑ **Definition of CG**
- ❑ **Concept of synteny**
- ❑ **Objectives of CG**
- ❑ **Concept of homology**
- ❑ **Visualization of CG**
- ❑ **Conclusion**

# History

- **Human Genome Project (2001-2003)** decided to use smaller genomes as warm-up for human genome
- Resulted in sequencing the following:
  - Many bacteria
  - Model-organism genomes
    - ✓ Yeast, *C. elegans*, *Arabidopsis*, *Drosophila*
    - ✓ Comparison of these genome sequences provided basis for field of comparative genomics

# Definition of comparative genomics (CG)

- **Comparative genomics** is a field of biological research in which the **genomic features of different organisms are compared**.

- The genomic features may include the **DNA sequence, genes, gene order, regulatory sequences**, and other genomic structural landmarks.

- In this branch of genomics, **whole or large parts** of genomes resulting from genome projects are compared to study **basic biological similarities** and **differences** as well as evolutionary relationships between organisms.

# Why are we doing CGs?



Remember the tree of life and how they hosted **2.3 million named species** of animals, plants, fungi and microbes. The tree of life traces **the origin of life** through **3.5 billion years of evolution**.

The tree of life was built using comparative genomics

http://www.tolweb.org/tree/

# Sizes of genomes and numbers of genes

# CG built on the concept of synteny

- Synteny: **genes** that are in the **same relative position** on **two different chromosomes**

- Genetic and physical maps compared between species

  – Or between chromosomes of the same species

- **Closely related species** generally have **similar order of genes on chromosomes**

- Synteny can be used to **identify genes in one species** based on **map position in another**

# Mouse vs Human synteny

- When sequences from mouse and human genomes are compared, **regions of remarkable synteny were found**

- Genes are in **almost identical order** for long stretches along the chromosome

**Human**
Chr 14

**Mouse**
Chr 14

Mouse vs Human synteny

"It is all about structure"

Chromosomes

# Human vs Rat synteny



Vitt *et al.* 2004 Genome Research, 14(4), 640-650.

# Human vs many species synteny



Froenicke *et al.* 2006 Genome Research, 16(3), 306-310.

# Rice vs Maize synteny

Soderlund *et al.* 2011 *Nucleic Acids Research*, 39(10), e68-e68.

# The objectives of CG

- **Comparison** of **genomic sequences** from **different species** can help identify the following:
  - **Gene structure** (Exon, Intron, 5'UTR, 3'UTR)
  - **Gene function** (Metabolism, Binding, etc)
  - **Regulatory sequences** (Promoters, Enhancers, etc)

# How to create a CGs analysis – molecular phylogeny

➤ The use of **molecular data** to establish the **relationship** between **species**, **organisms** or **gene families**

❑ **Homology**

Sequences that **share common ancestry**

Homologous genes can be similar in sequence, but **similar sequences are not necessarily homologous**

❑ **Orthologs**

Homologs in **different species** derived by a **speciation** event

❑ **Paralogs**

Homologs in the **same or different species** derived by a **duplication** event

Jensen, R. A. 2001 *Genome Biology,* 2, 1-3

# Understanding the concept of homologs

# Understanding the concept of homologs



1. A1 and B1 – ...
2. A1 and B2 – ...
3. A2 and B1 – ...
4. A2 and B2 – ...
5. A1 and A2 – ...
6. B1 and B2 – ...

Jensen, R. A. 2001 *Genome Biology,* 2, 1-3

# Understanding the concept of homologs



1. A1 and B1 – **paralog**
2. A1 and B2 – **paralog**
3. A2 and B1 – **paralog**
4. A2 and B2 – **paralog**
5. A1 and A2 – **ortholog**
6. B1 and B2 – **ortholog**

Jensen, R. A. 2001 *Genome Biology,* 2, 1-3

# Understand the concept of conserved region

Conserved sequences are **similar** or **identical** sequences in nucleic acids (DNA and RNA) or proteins **across species** (**orthologous** sequences) or within a genome (**paralogous** sequences). **Conservation** indicates that a sequence has been **maintained by natural selection**.

# Visualization of CGs – phylogenetic tree

❑ Visualize **evolutionary relationships** between **species** and **genes/proteins**

❑ **Rooted tree**

   – Order of evolutionary events

❑ **Unrooted tree**

   – Evolutionary relationships between descendants

# Visualization of CGs – Dot Plot

❑ A **graphical method** for **comparing two biological sequences/genomes** and identifying **regions of close similarity**

❑ **Synteny**
Gene loci are on the same chromosome

❑ **Conserved synteny**
Gene loci are on the same chromosome in different species

❑ **Collinearity**
The order of the gene loci is preserved across species

# Visualization of CGs – Dot Plot



**Match chromosome sequence from species A to species B**

# Visualization of CGs – Dot Plot

**Different pattern in the genome to genome comparison**



canonical          weak forward similarity          off-diagonal

inverted duplication          inverted repeat          palindrome          inverted transposition

Chaisson *et al.* 2006 *PNAS* **103** 19824-9
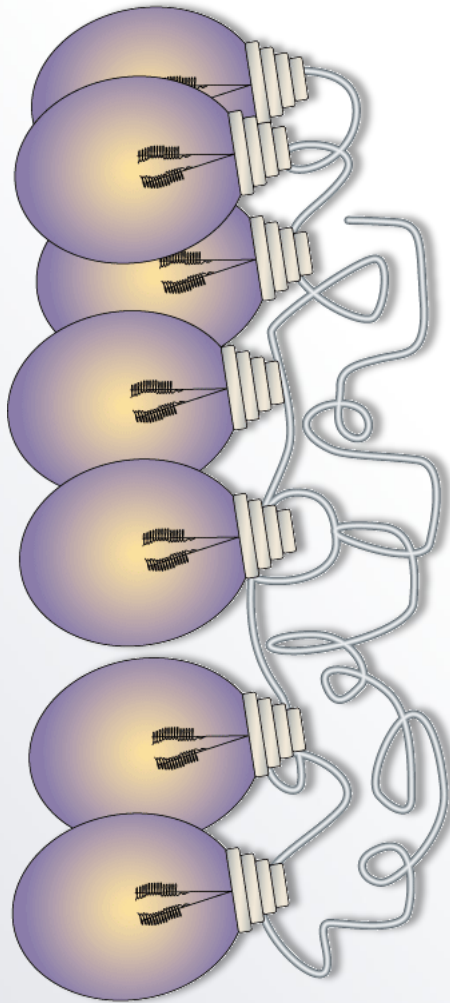
# Principals of CGs

# What do we need to do a CGs?

➤ Genome/transcriptome of **a query/target species**

➤ Genome/transcriptome of **a reference species**

➤ **Bioinformatics tools**:

    ❑ MEGA-BLAST
    ❑ Multiple Sequence Alignment
    ❑ Annotator
    ❑ Tree builder

# Conclusion: the goal of CG
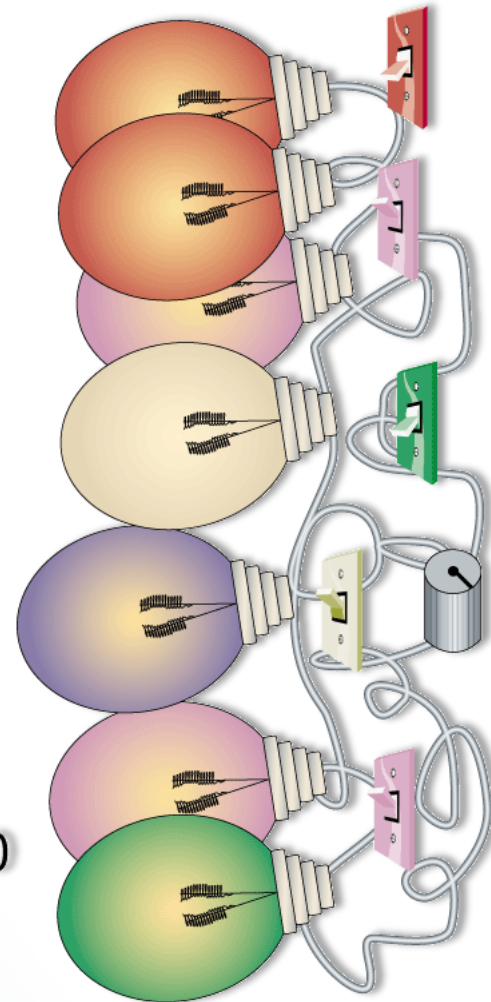


Gene 1

Gene 2

Gene 3

Gene 4

Gene 40,000

**Unknown species**          **Model species**

# It was still the fifth course, don't get dizzy yet