# Bioinformatics
# A Communication/Signal Processing Perspective

Khalid Sayood
Occult Information Lab
Department of Electrical and Computer Engineering
University of Nebraska-Lincoln

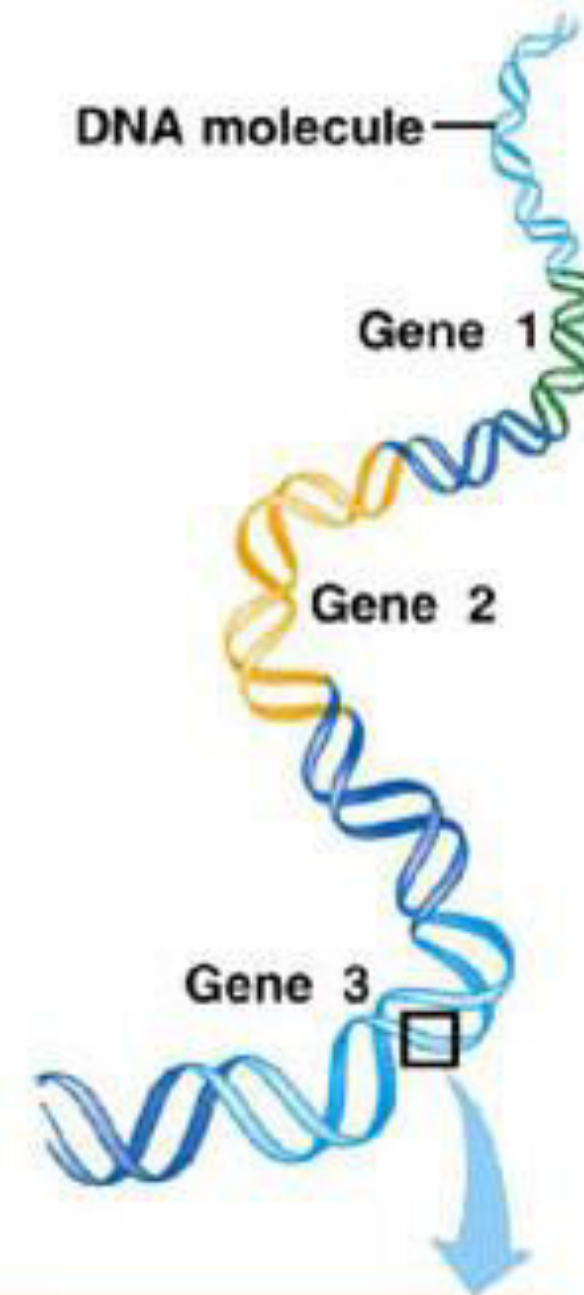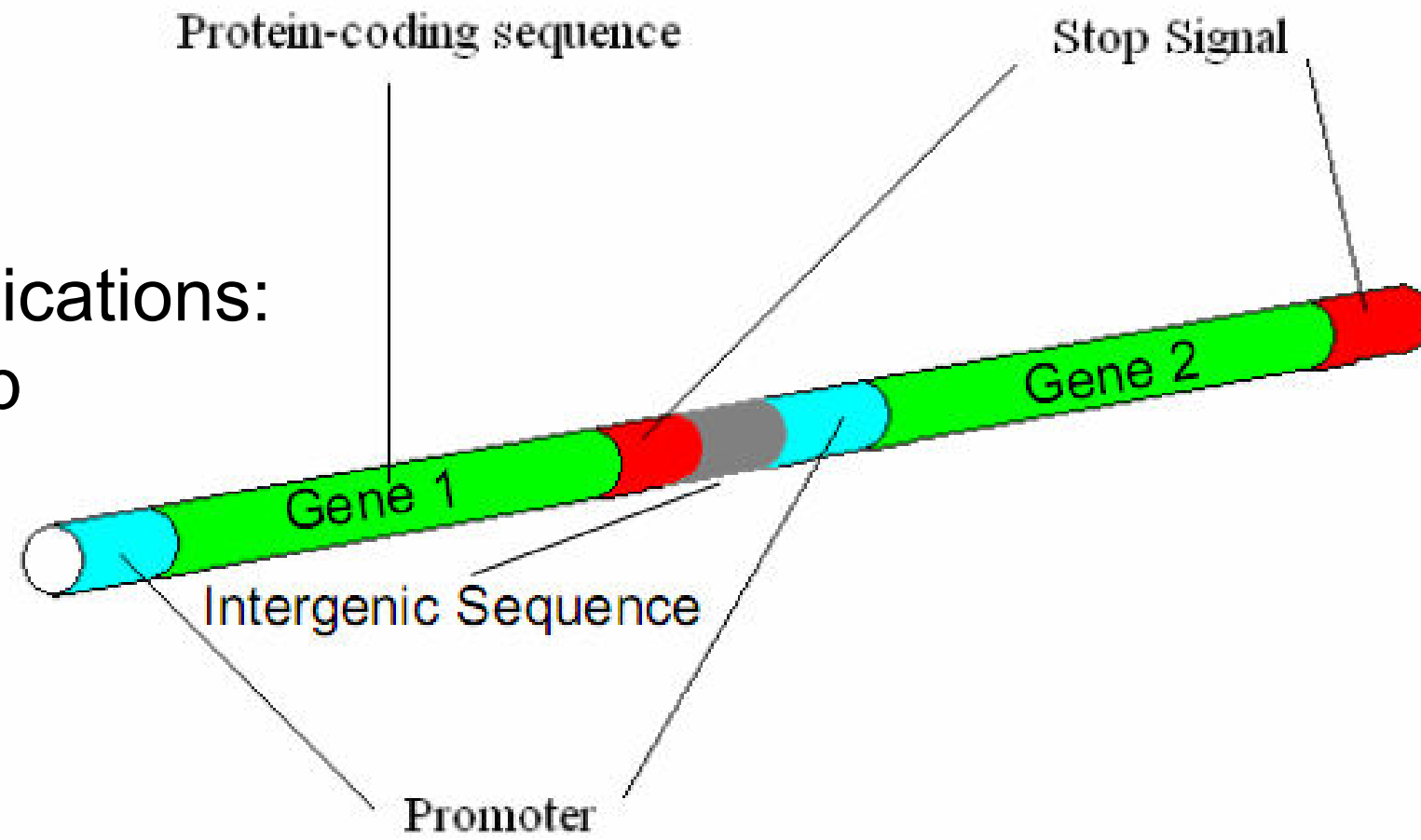# Outline

- Beginning
- Middle
- End

# My roots

- Signal Processing
- Communication
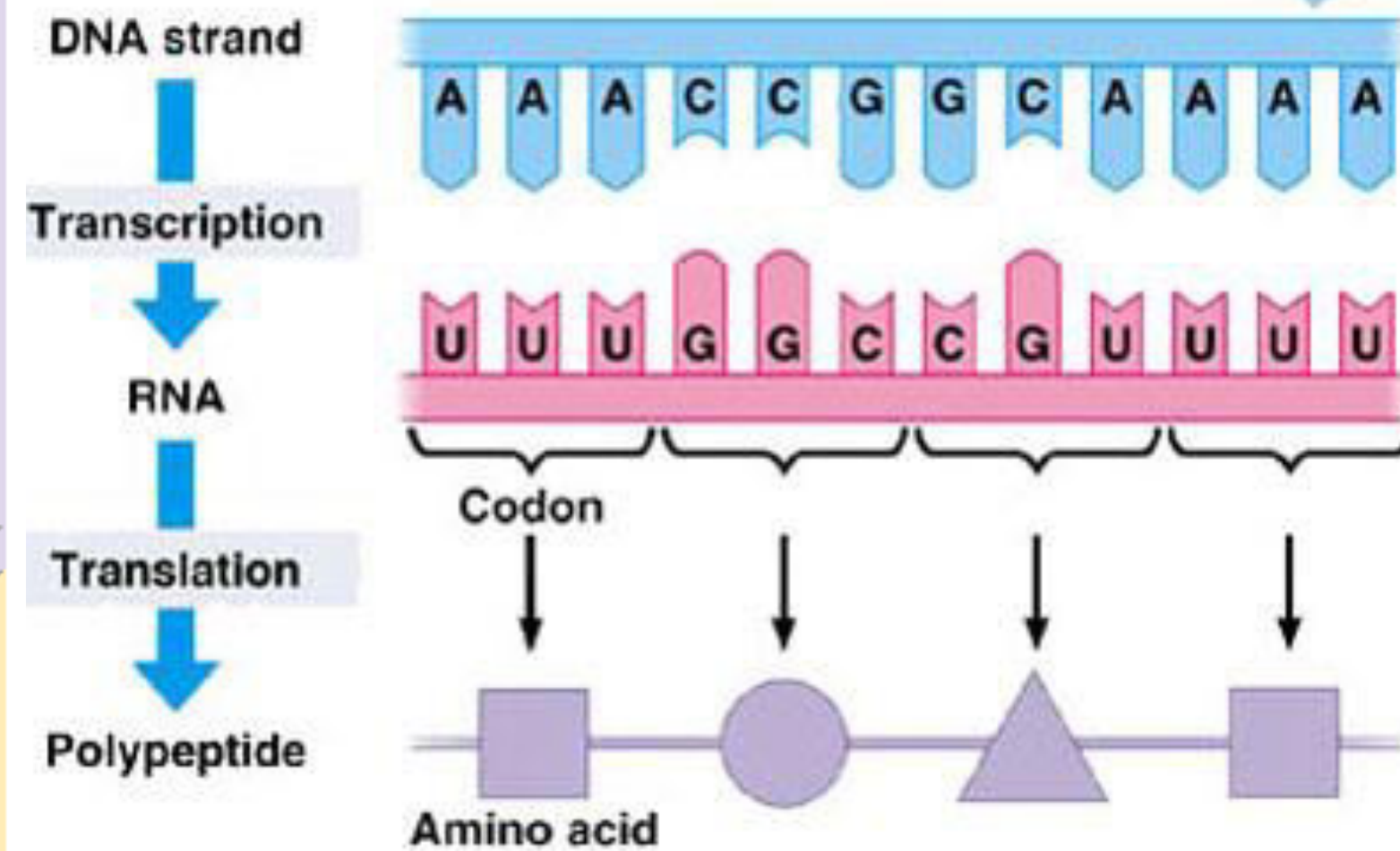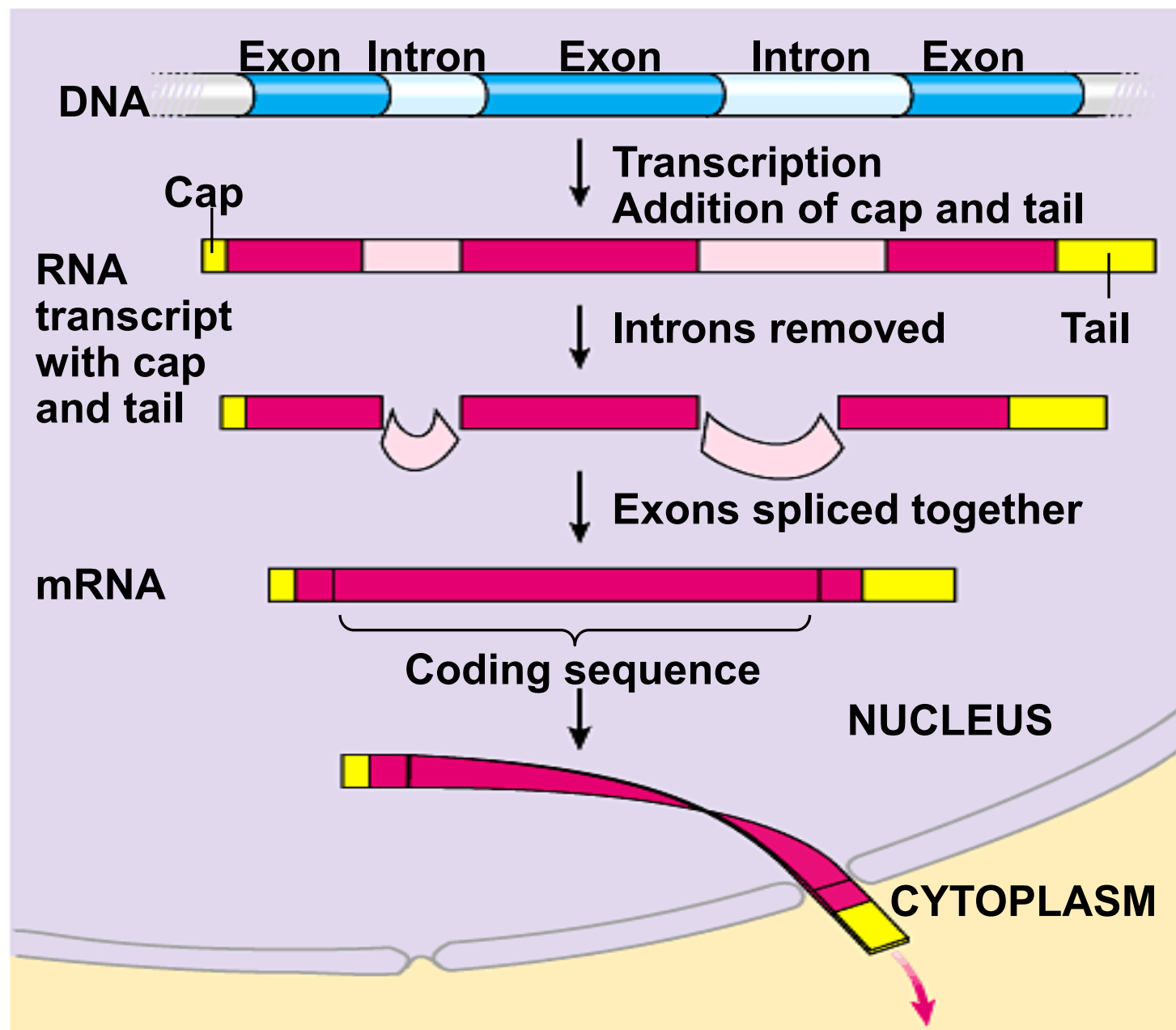- Information Theory
- Data Compression

# My roots

- Data Compression - The art or science of finding compact representations of information in data
  - Speech ⟶ Images ⟶ Video ⟶ DNA

Prokaryotes

Analogy to Communications:
Start, Message, Stop

Protein-coding sequence

Stop Signal

Gene 1

Gene 2

Intergenic Sequence

Promoter

DNA molecule

Gene 1

Gene 2

Gene 3

Eukaryotes

Exon  Intron      Exon       Intron    Exon

DNA

Transcription
Addition of cap and tail

Cap

RNA
transcript
with cap
and tail

Tail

Introns removed

Exons spliced together

mRNA

Coding sequence

NUCLEUS

CYTOPLASM

DNA strand

A A A C C G G C A A A A

Transcription

RNA

U U U G G C C G U U U U

Translation

Codon

Polypeptide

Amino acid

# Definitions of Bioinformatics

- Original definition: Study of informatic processes in biotic systems – Pauline Hogeweg and Ben Hesper

# Definitions of Bioinformatics

- Original definition: Study of informatic processes in biotic systems – Pauline Hogeweg and Ben Hesper

- Bioinformatics is conceptualizing **biology in terms of molecules** (in the sense of physical-chemistry) and then applying **"informatics"** techniques (derived from disciplines such as applied math, CS, and statistics) to understand and **organize the information associated** with these molecules,  **on a large-scale**. (Mark Gerstein, 1999)

# Bioinformatics is a management and analysis information system for life sciences.

**Data Storage and Management**

**Data Analysis**

**Interpretation of Results**

## Protein Structure Prediction

- Protein/RNA tertiary structure
- Docking
- Drug Design

## Molecular Sequence Analysis

- Homology Search
- Phylogeny Construction
- Whole Genome Sequencing
- Gene Finding

## Functional Genomics and Proteomics

- Microarrays
- Biomarker Discovery

## Systems Biology

- Pathways
- Network based wholistic approach

From RNA-seq reads to differential expression results
•Alicia Oshlack, Mark D Robinson andMatthew D Young
Genome Biology201011:220

Preprocessing:

*FastQC*

*Trimmomatic*

Indexing the Reads:

*Bowtie 2*

Aligning to hg19/KSHV

*TopHat 2*

Local alignment

Extracting Features

*CuffDiff*

P < 0.01

FPKMs values

Differentially Expressed Genes

# So what is the problem?

- What is actually going on.

# Is there hope (for me)?

**Hope springs eternal**

Digital Camera (1975)

Genome Sequencer (2018)

# What do we mean by a Communication Theory Perspective

# What do I mean by a Communication Theory Perspective

Information exists in the form of a stochastic process

# How do we deal with stochastic processes?

- Look at the signal using different basis sets – frequency domain processing.
- Look at correlation structures.
- Look at models.

All these involve averaging of some sort

All these result in the discovery of underlying structure

They can also result in dimensionality reduction.

# Realizations of a stochastic process

# Frequency profile

# Statistical Profile

# Models

But we have sequences – we don't have numbers



GAGACATTCAGTG

But we have sequences – we don't have numbers

A

C

G

T

GAGACATTCAGT

But we have sequences – we don't have numbers

|   | G | A | G | A | C | A | T | T | C | A | G | T | G |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| C | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| G | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| T | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 |

**Identification of Protein Coding Regions Using the Modified Gabor-Wavelet Transform**
Mena-Chalco et al. IEEE/ACM Trans. Comp. Bio

But we have sequences – we don't have numbers

| G | A | G | A | C | A | T | T | C | A | G | T | G |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| AA | AC | AG | AT | CA | CC | CG | CT | GA | GC | GG | GT | TA | TC | TG | TT |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0 | 1 | 2 | 1 | 2 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | 0 | 1 | 1 | 1 |

But we have sequences – we don't have numbers

$$I_k = \sum_{X \in A} \sum_{Y \in A} p_k(X,Y) \log\left( \frac{p_k(X,Y)}{p(X)p(Y)} \right)$$

k=6

X                                        Y

GAGACAT○○○○○○○○○○○○○○○

X                                                              Y

k=14
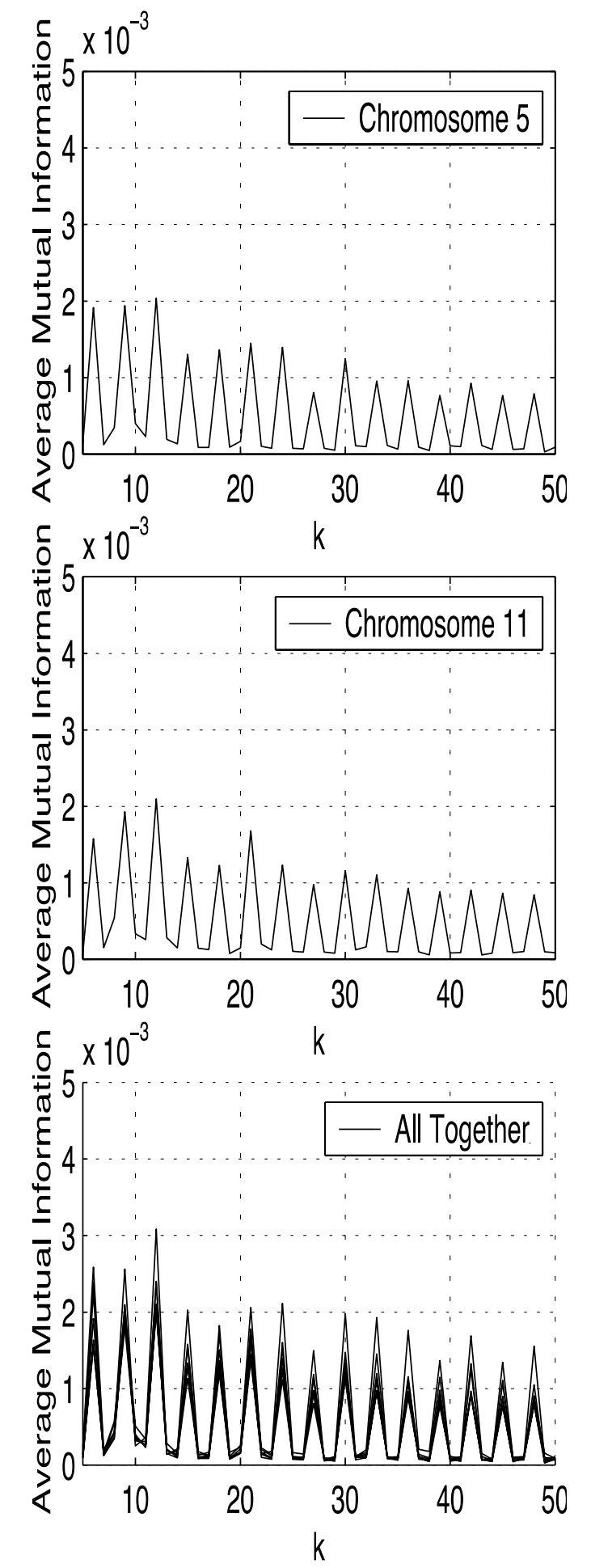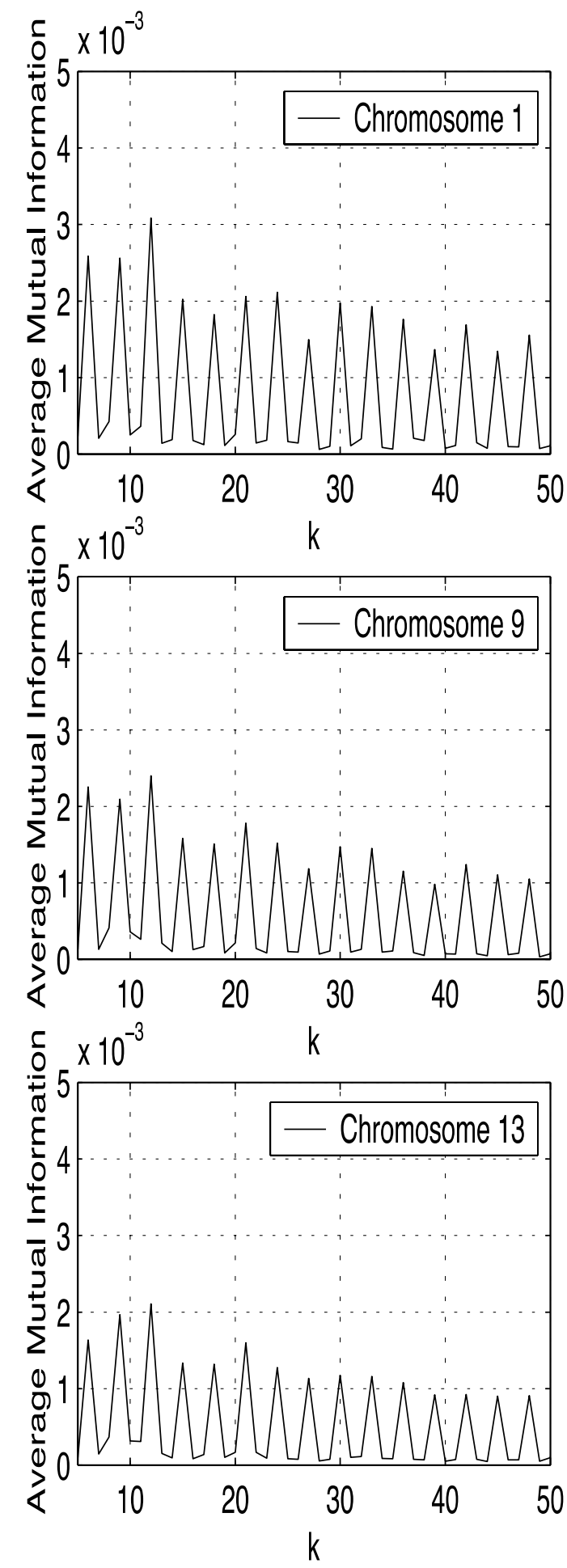
# AMI Profile for Human Chromosome 1

# Human chromosomes
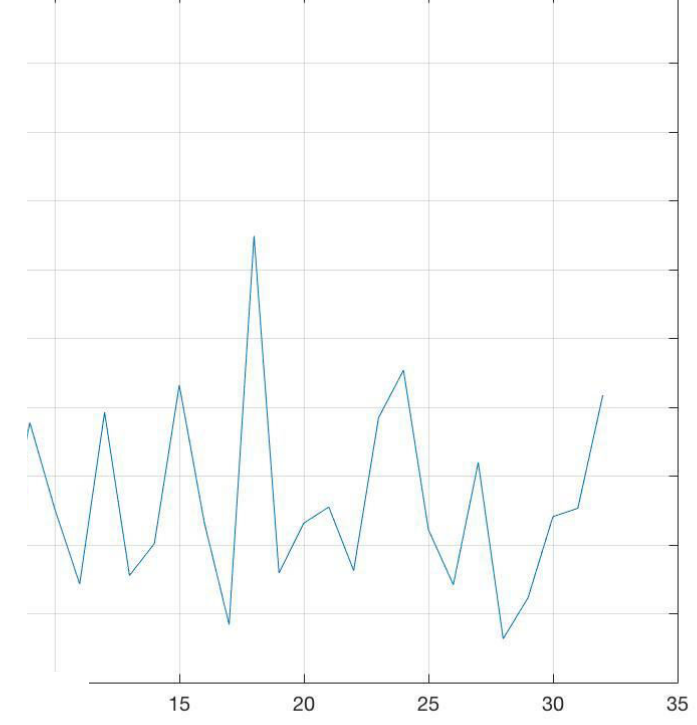
# Mouse Chromosome

a) C. Elegans Chromosomes

b) S. Cerevisae Chromosomes
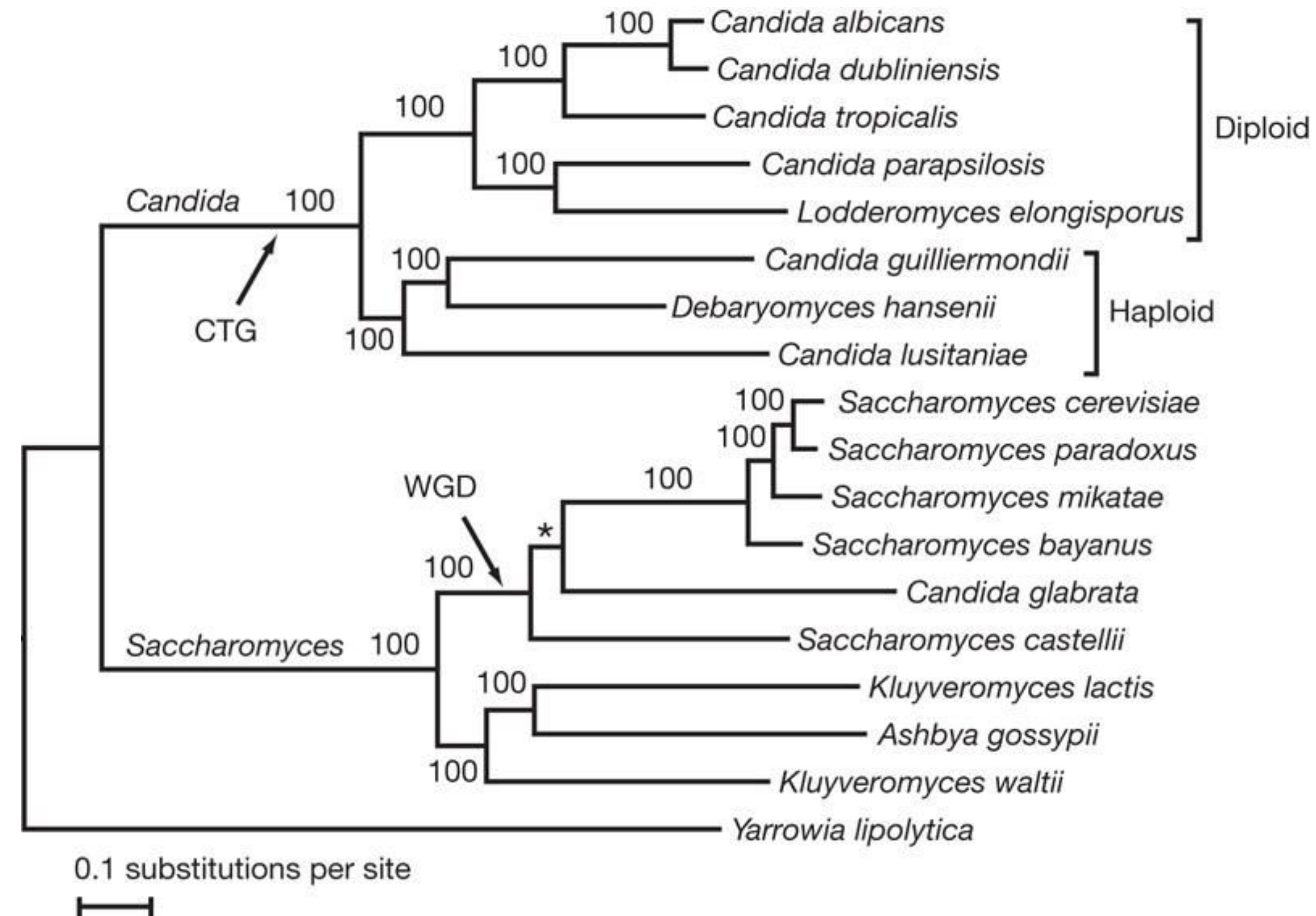
- Phylogeny for *Candida* and *Saccharomyces* clades based on multiple sequence alignment of 706 orthologous genes

- Posterior probabilities shown

- WGD: Whole Genome Duplication

- CTG: Translation of CTG codons as serine rather than leucine

- Ribosomal DNA (rDNA) is commonly used to evaluate species relatedness

- The rDNA gene complex contains 3 genes, each of which are ribosomal components once transcribed

- Internal transcribed spacer (ITS) 1 and ITS2 separate these genes

- ITS regions have 2 benefits:

  1. Easy to design primers (ribosome genes highly conserved, many copies)

  2. Spacers diverge more quickly than ribosome genes

- Distance matrix $D$ generated by calculating pairwise distance $d_{ij}$ between AMI profiles $\boldsymbol{x}_i$ and $\boldsymbol{x}_j$

- Distance defined in two ways:
  1. Correlation distance (angle between profiles)
  $$d_{ij} = 1 - \cos\theta = 1 - \frac{\boldsymbol{x}_i \cdot \boldsymbol{x}_j}{\|\boldsymbol{x}_i\| \|\boldsymbol{x}_j\|}$$

  2. Euclidean distance
  $$d_{ij} = \|\boldsymbol{x}_i - \boldsymbol{x}_j\|$$
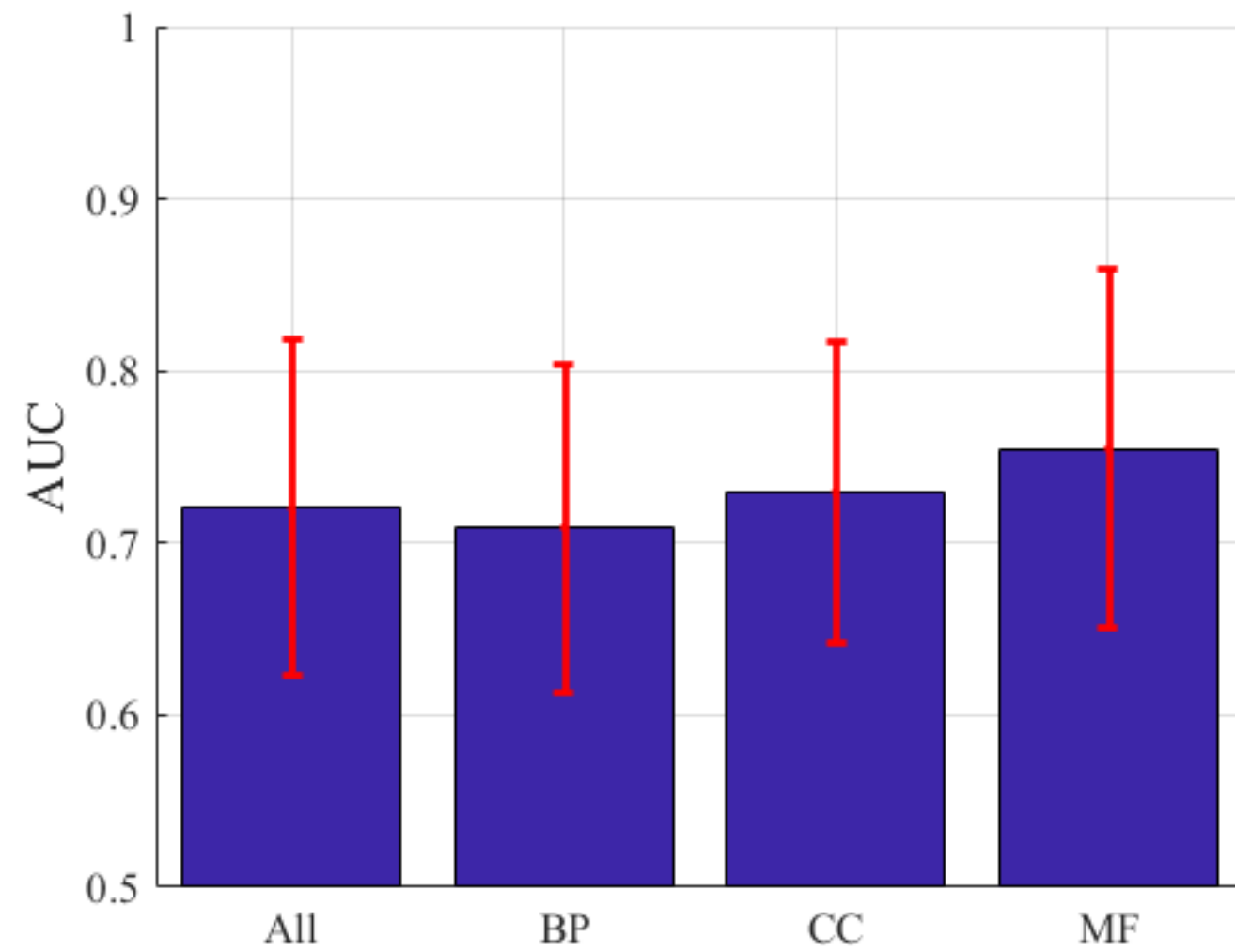
- Phylogenetic trees generated using PHYLIP (neighbor joining)

# GO Prediction



"High Abundance" GO terms

"Low Abundance" GO terms

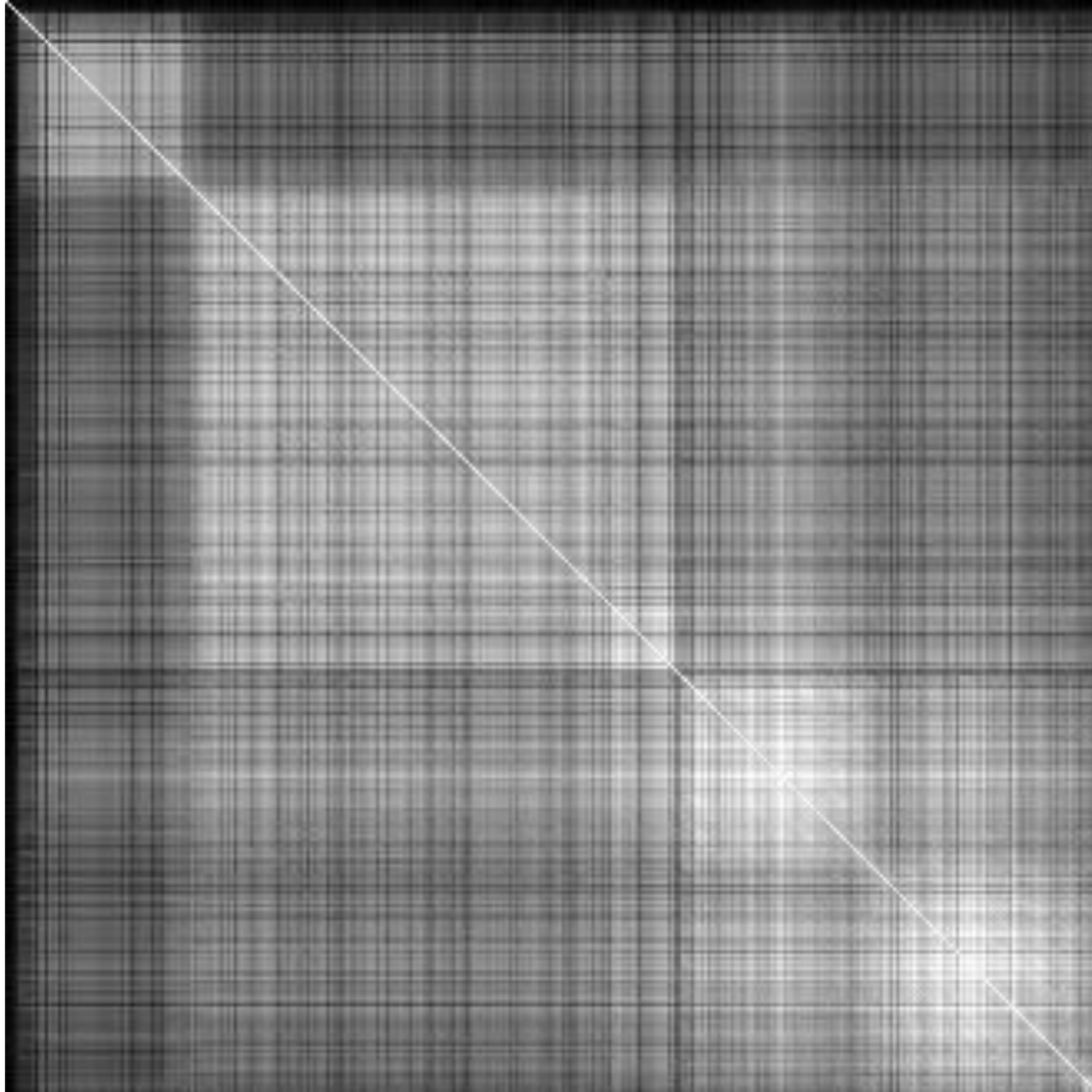BP: Biological Processes, CC: Cellular Component, MF: Molecular Function

Slow progressor populations
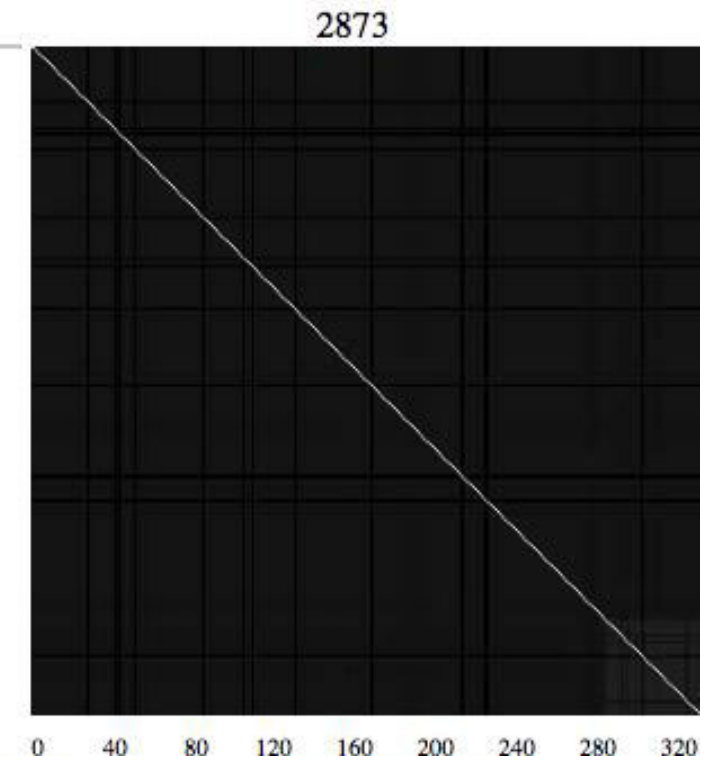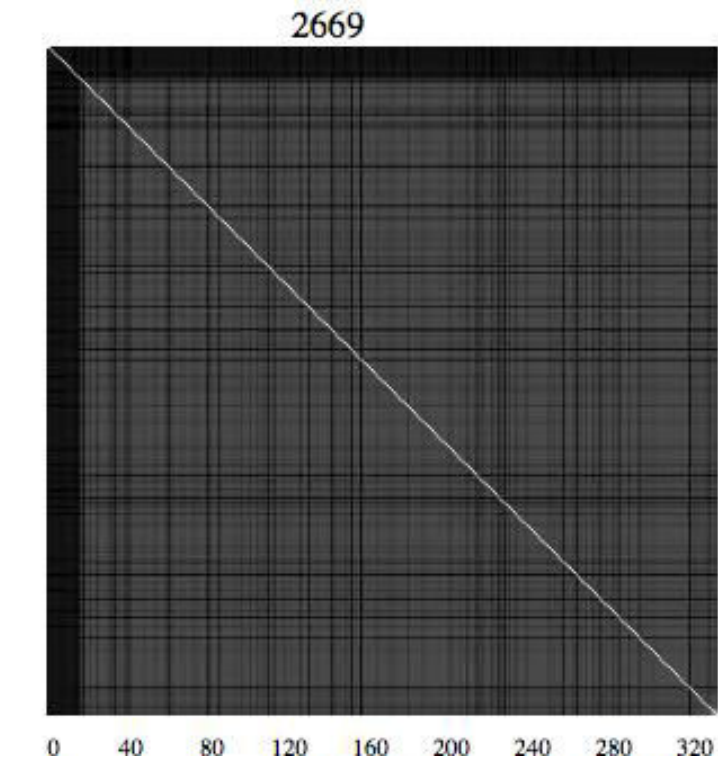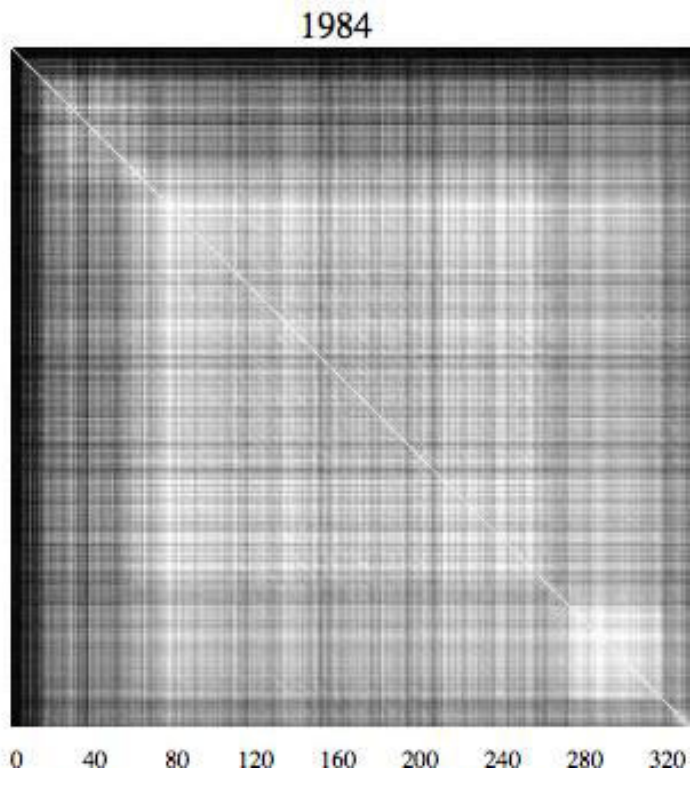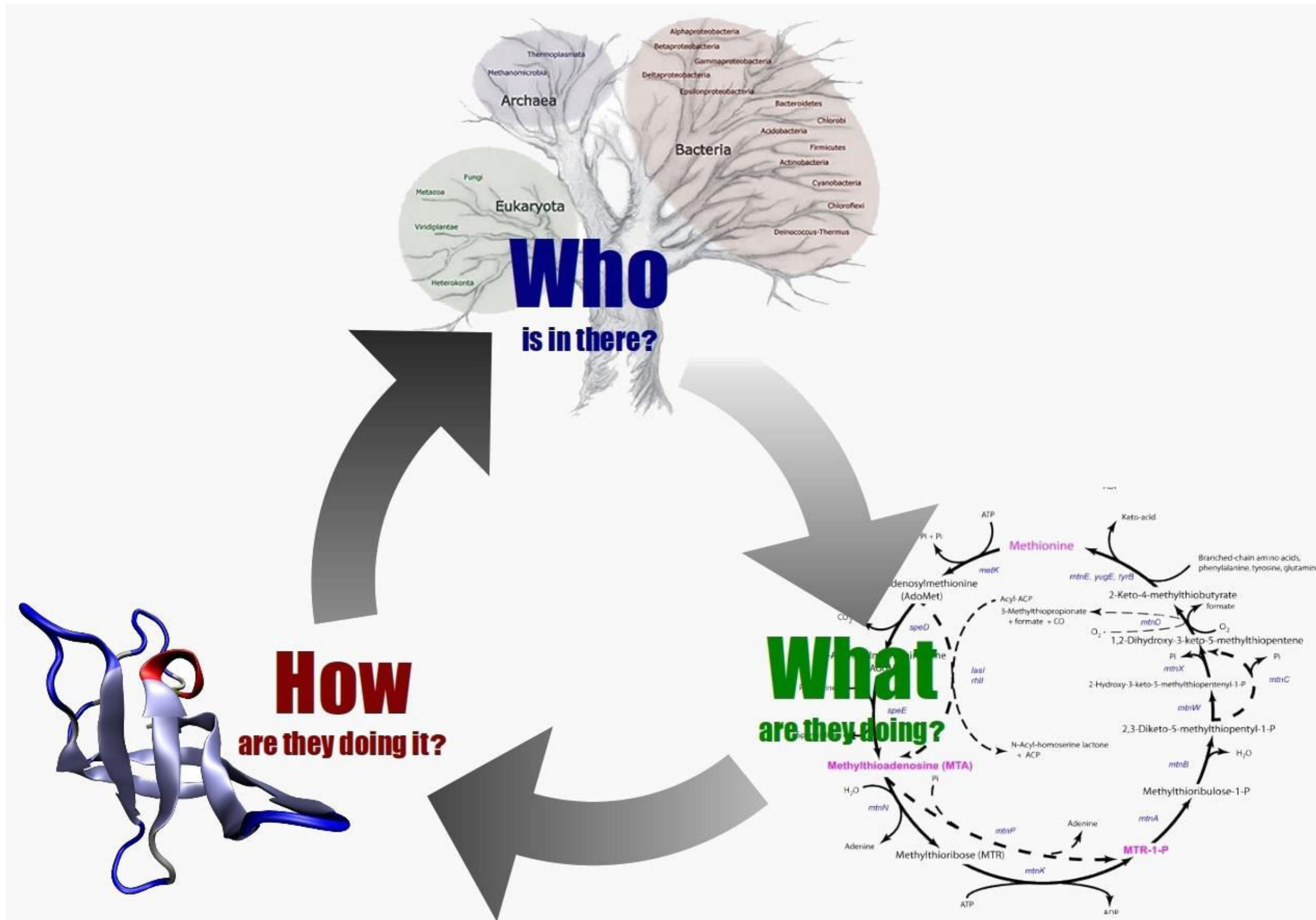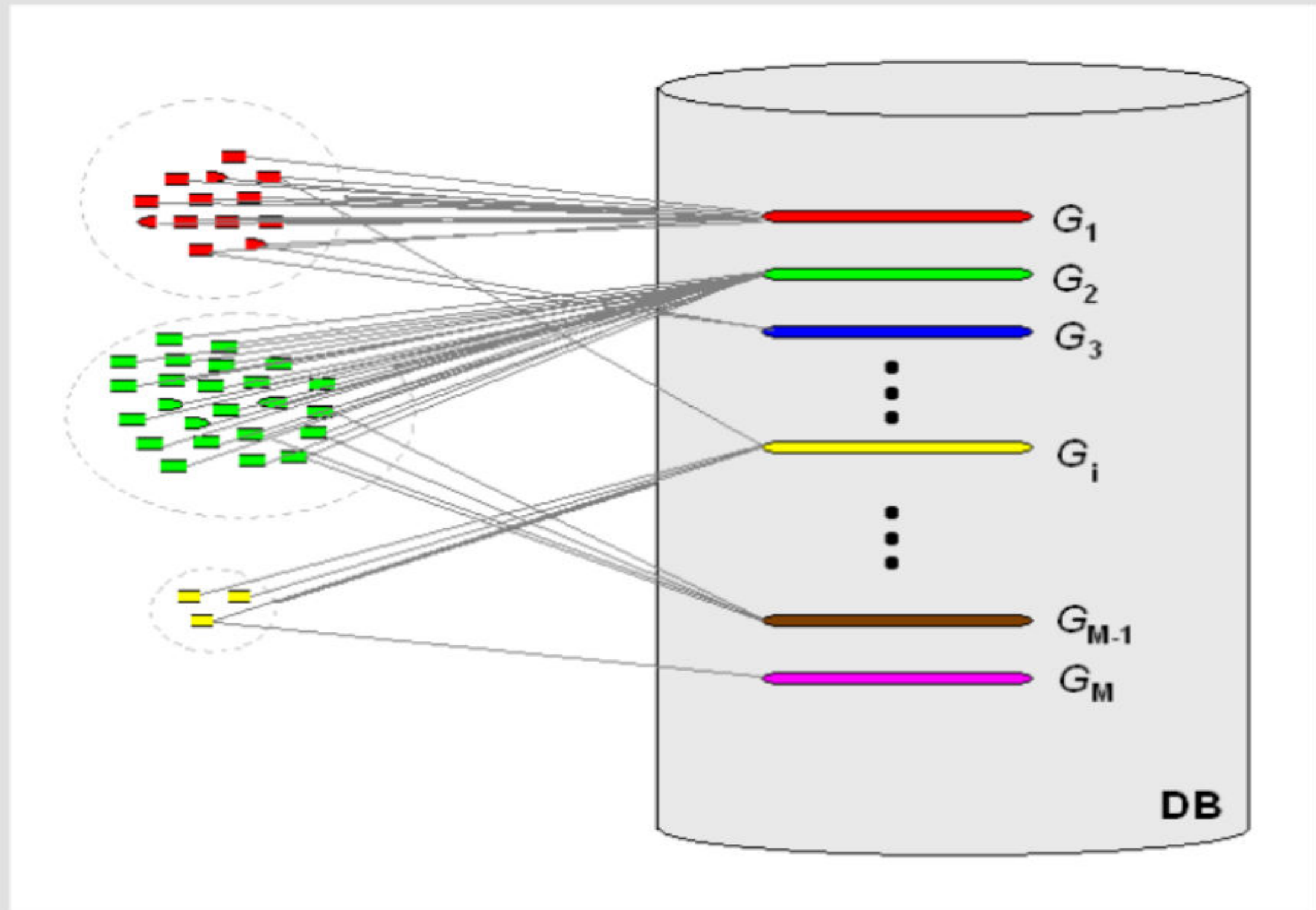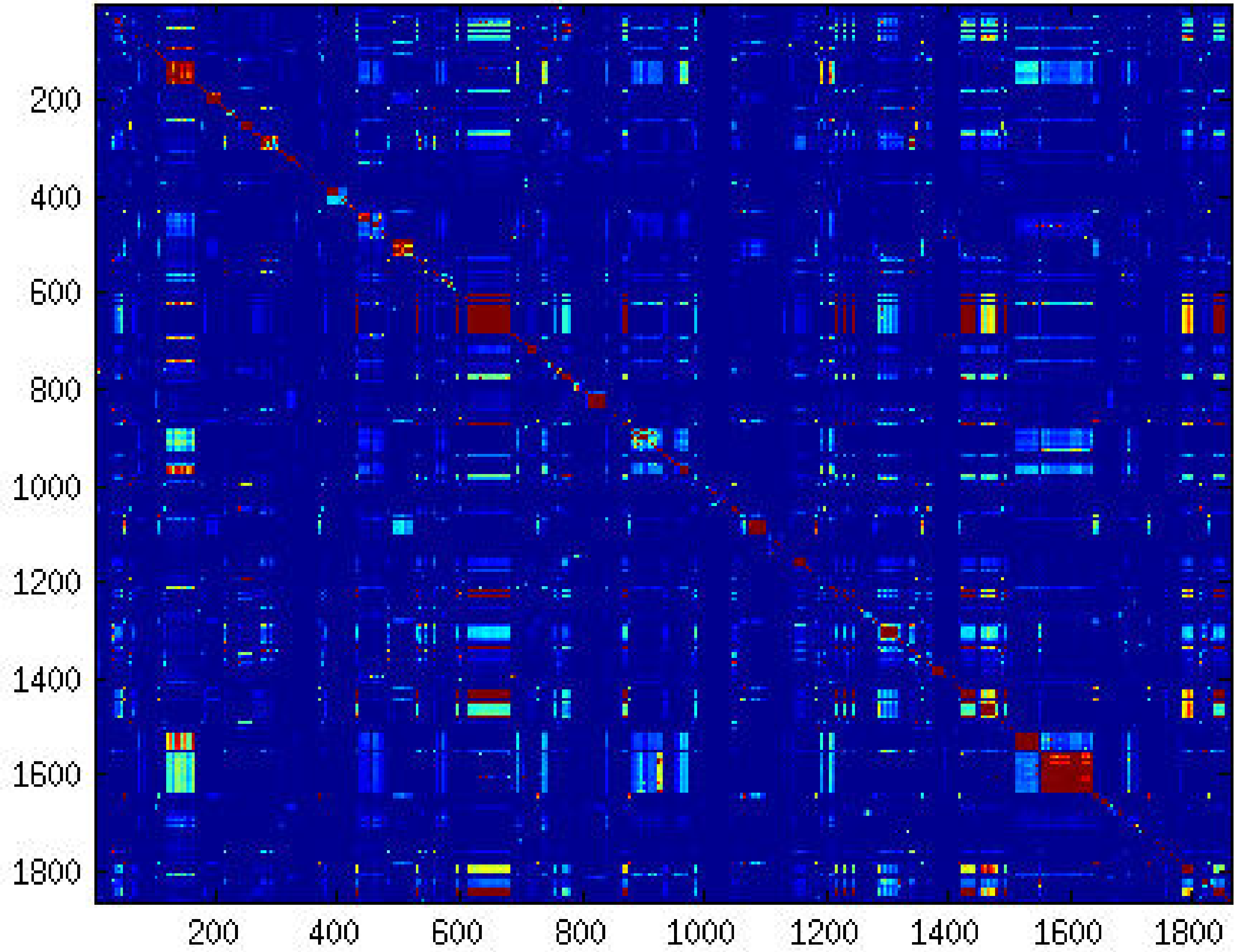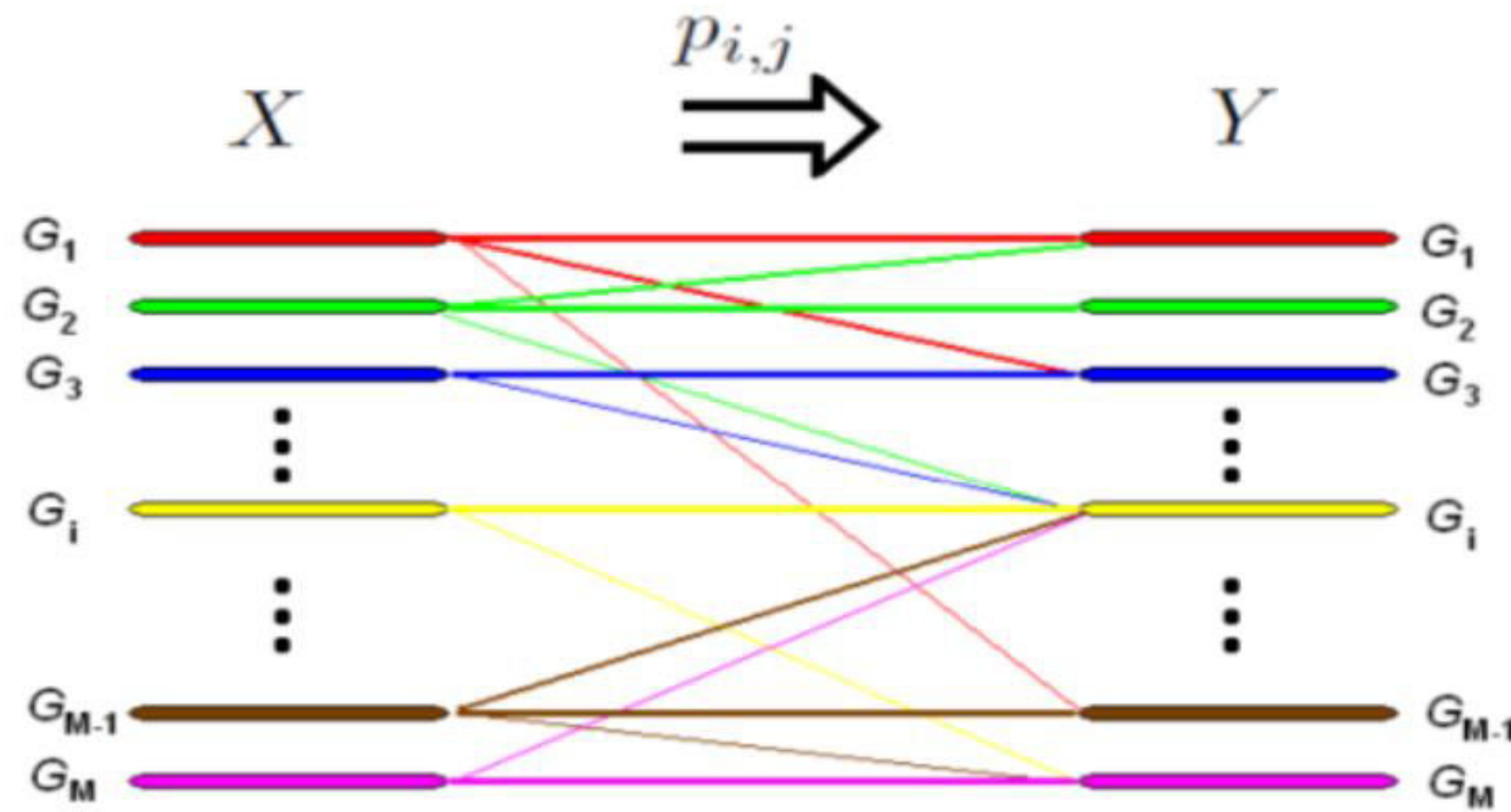
Rapid progressor populations

# Metagenomics

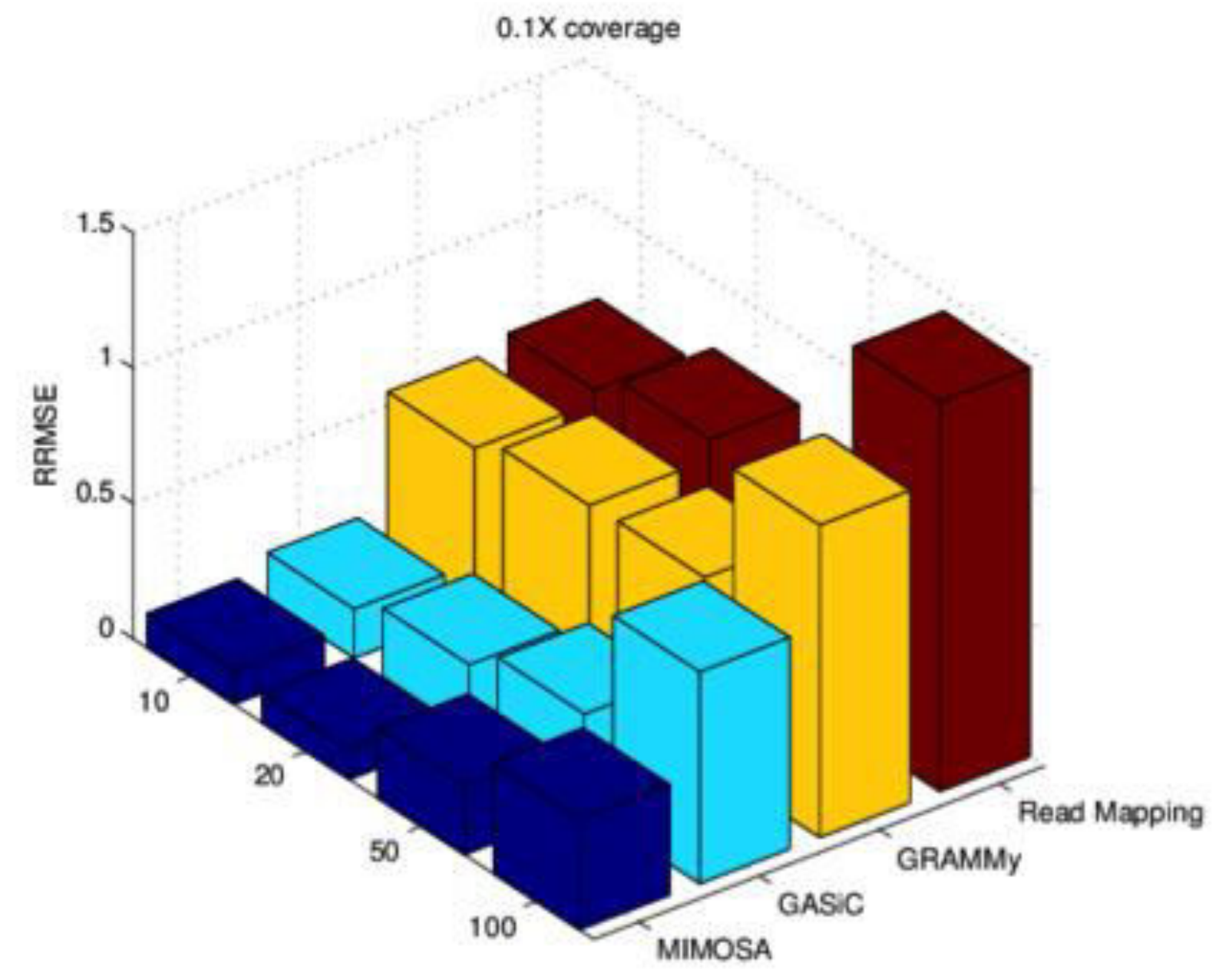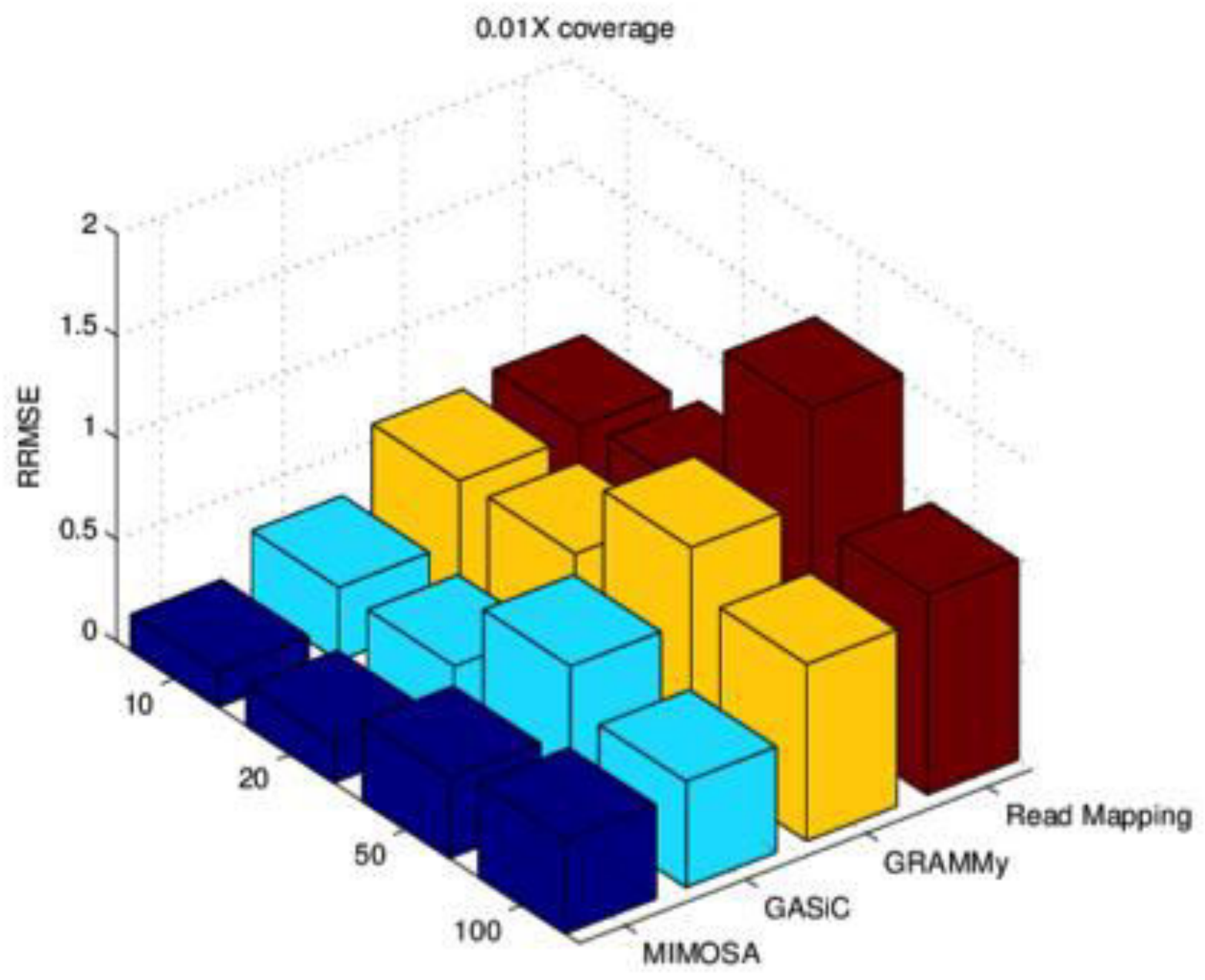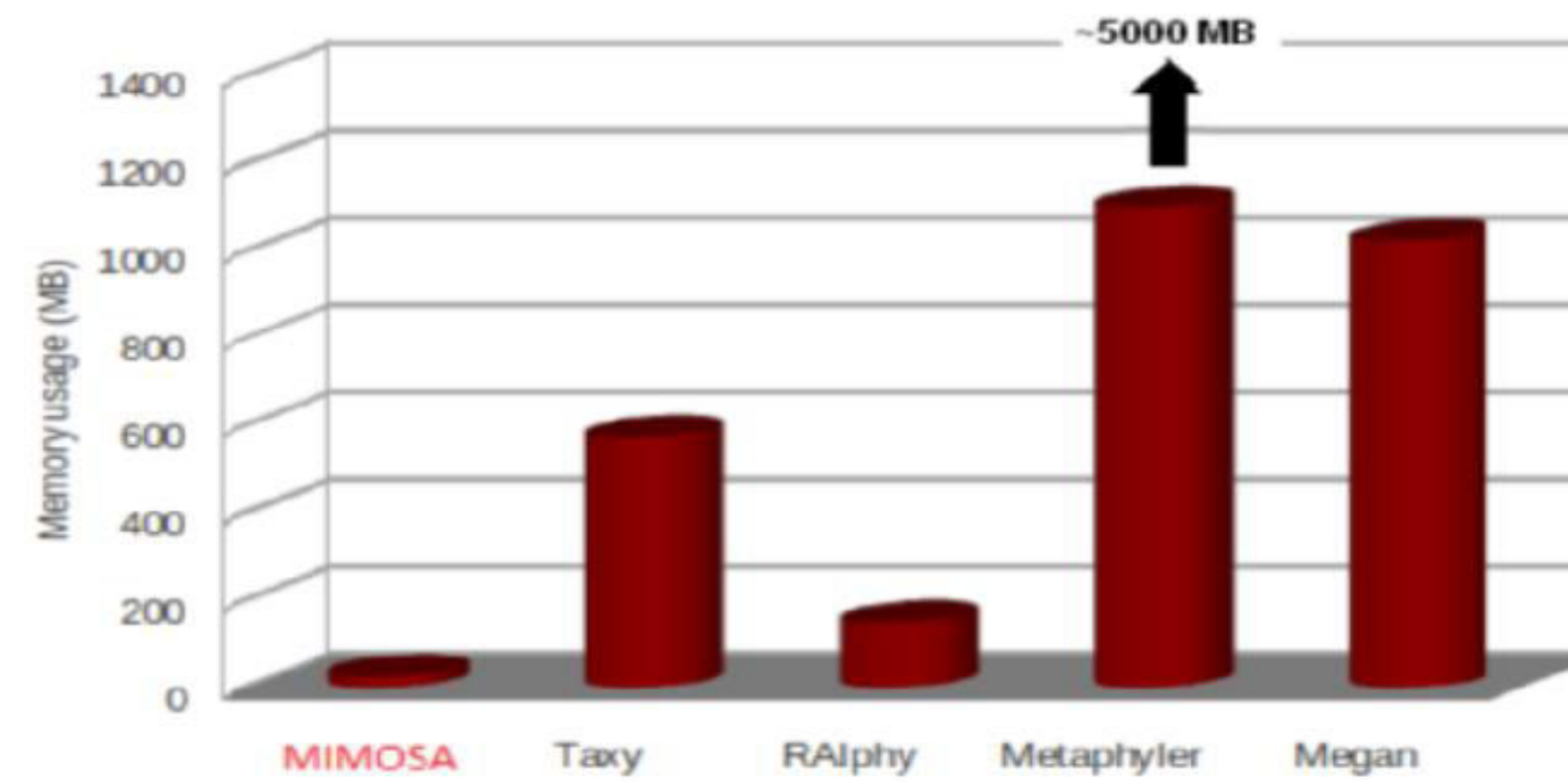# Similarities between genomes + single read processing = detection errors
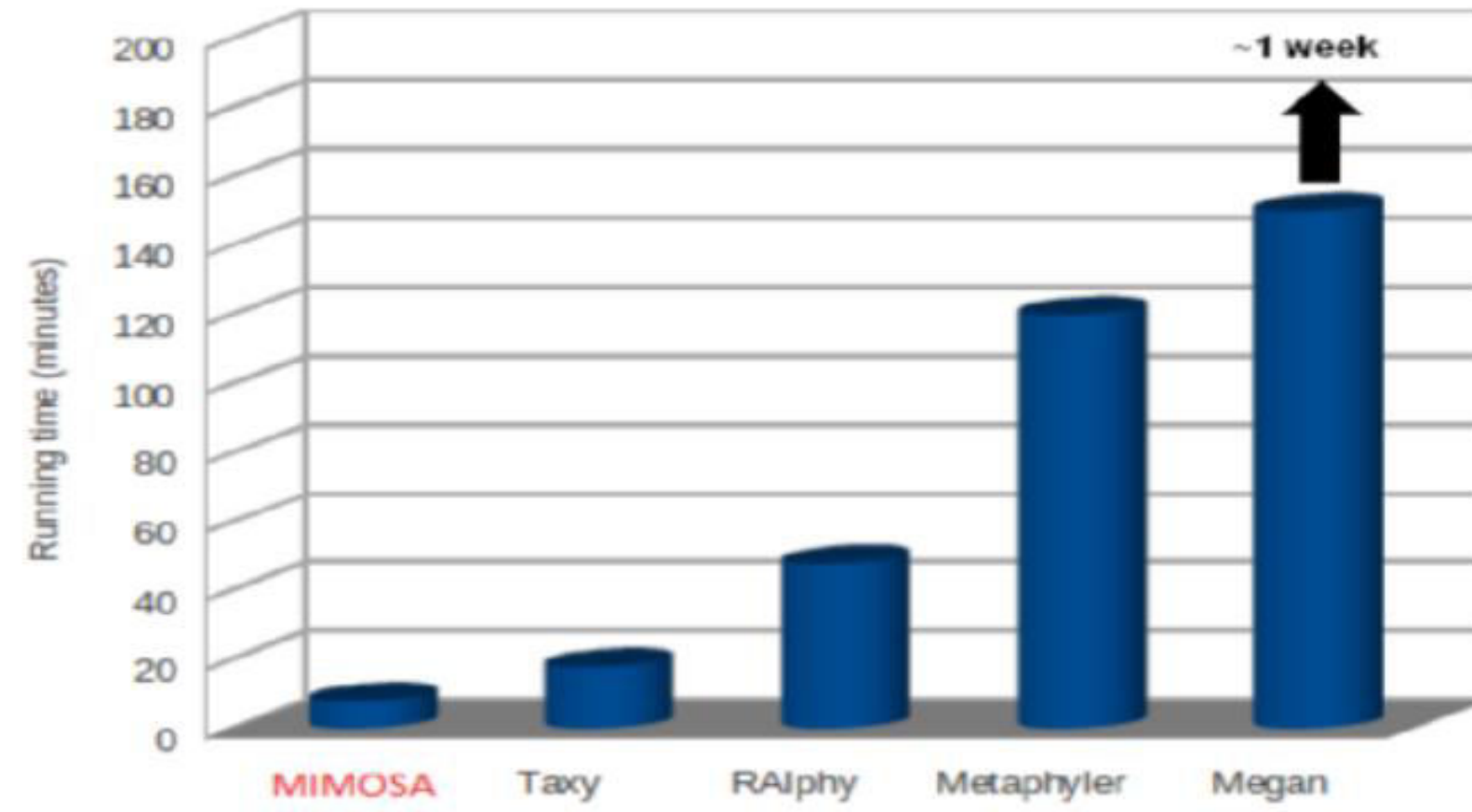
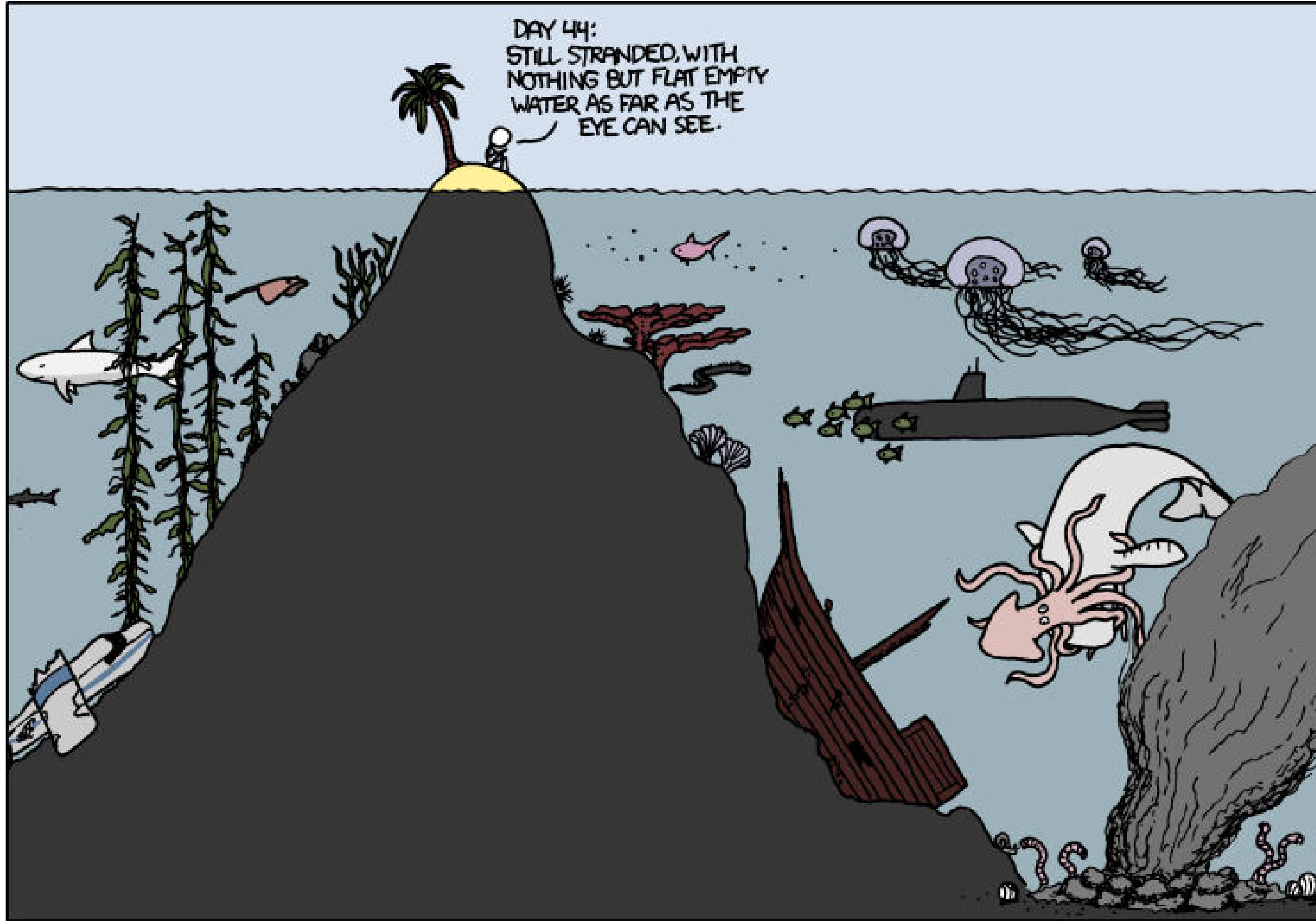Learn channel statistics to estimate true abundance given the observed detection.



$$p_{i,j} = \frac{|G_i \cap G_j|}{|G_j|}$$

RRMSE for the simulated metagenomes corresponding to a mixture of 10, 20, 50, and 100 randomly selected organisms for 0.01X and 0.1X average genome coverages.

# OCCULT INFORMATION LAB
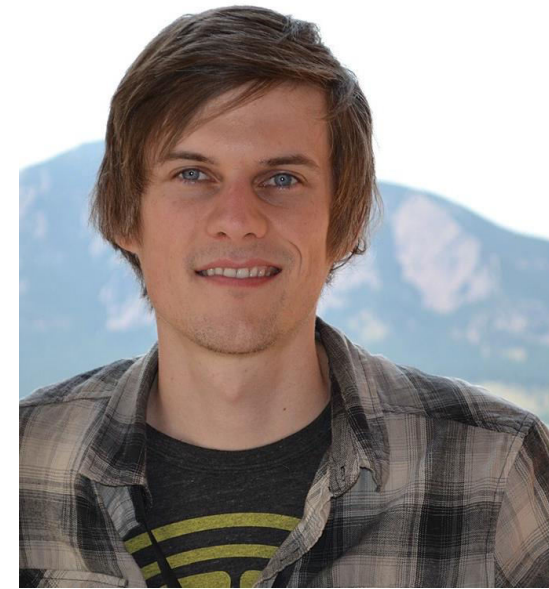
## BIOINFORMATICS

Mark Bauer

Hasan Otu

David Russell

Ufuk Nalbantoglu

Sam Way

Garin Newcomb

Amirsalar Mansouri

Dicle Yalcin

Keith Murray

Jacob Bohac