

CSCI8980: Applied Machine Learning in Computational Biology

Introduction to Bioinformatics

Rui Kuang

Department of Computer Science and Engineering

University of Minnesota

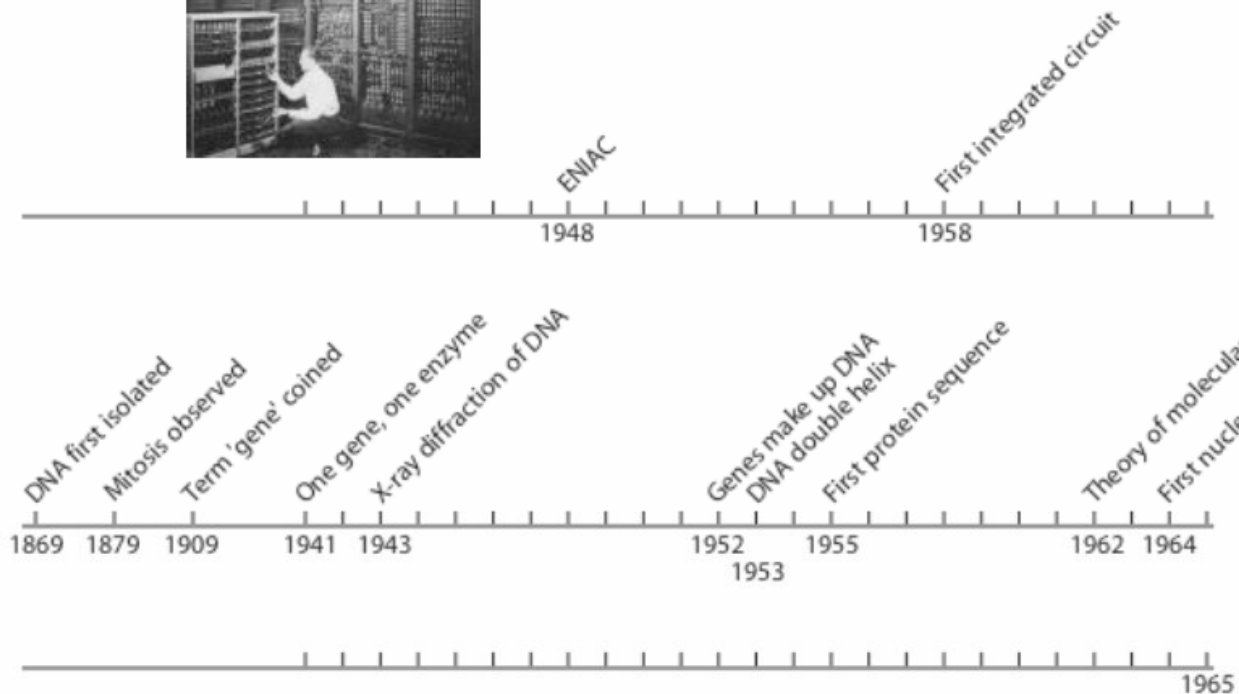
kuang@cs.umn.edu

UNIVERSITY OF MINNESOTA

Twin Cities • Duluth • Morris • Crookston • Rochester • Other Locations



History of Bioinformatics



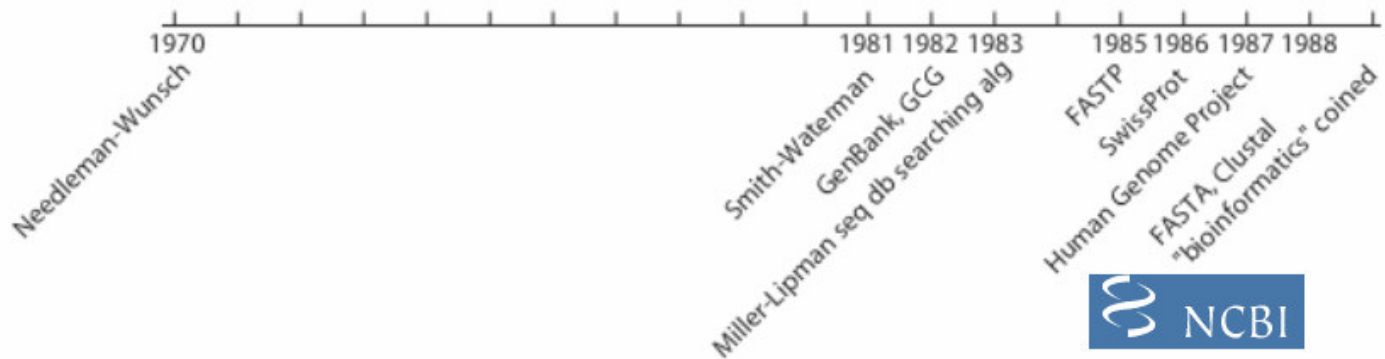
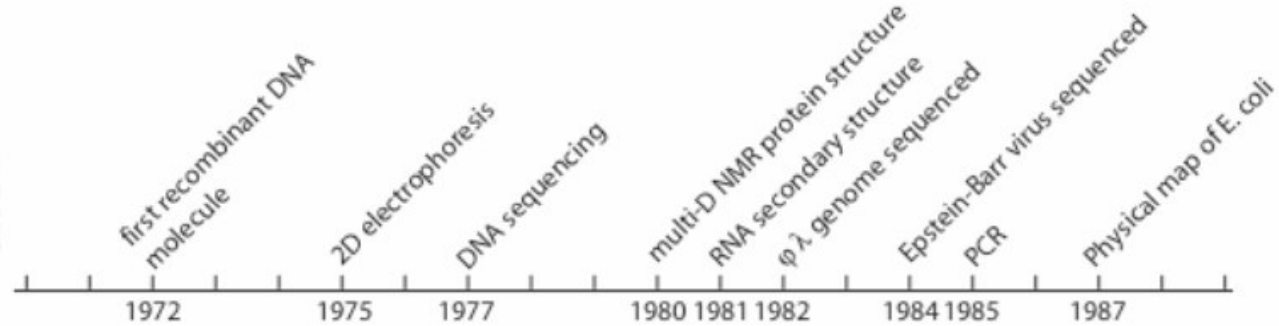
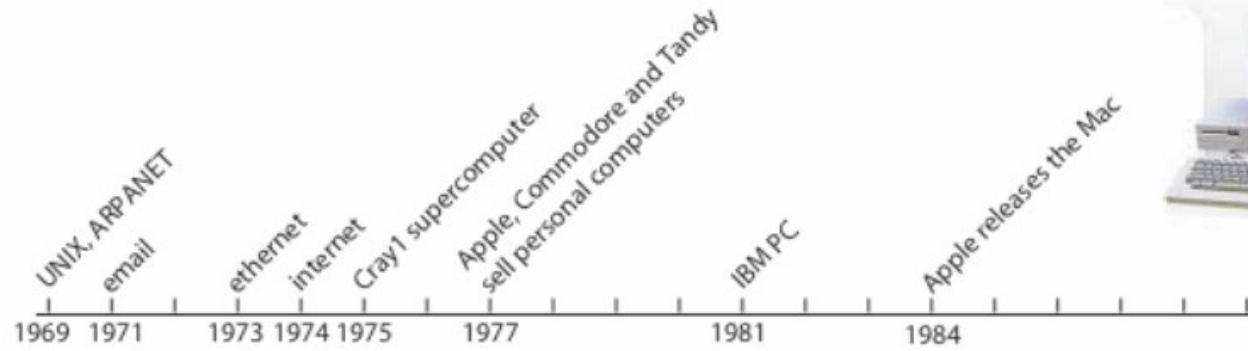
Margaret Dayhoff

Thanks to Luce Skrabanek

History of Bioinformatics

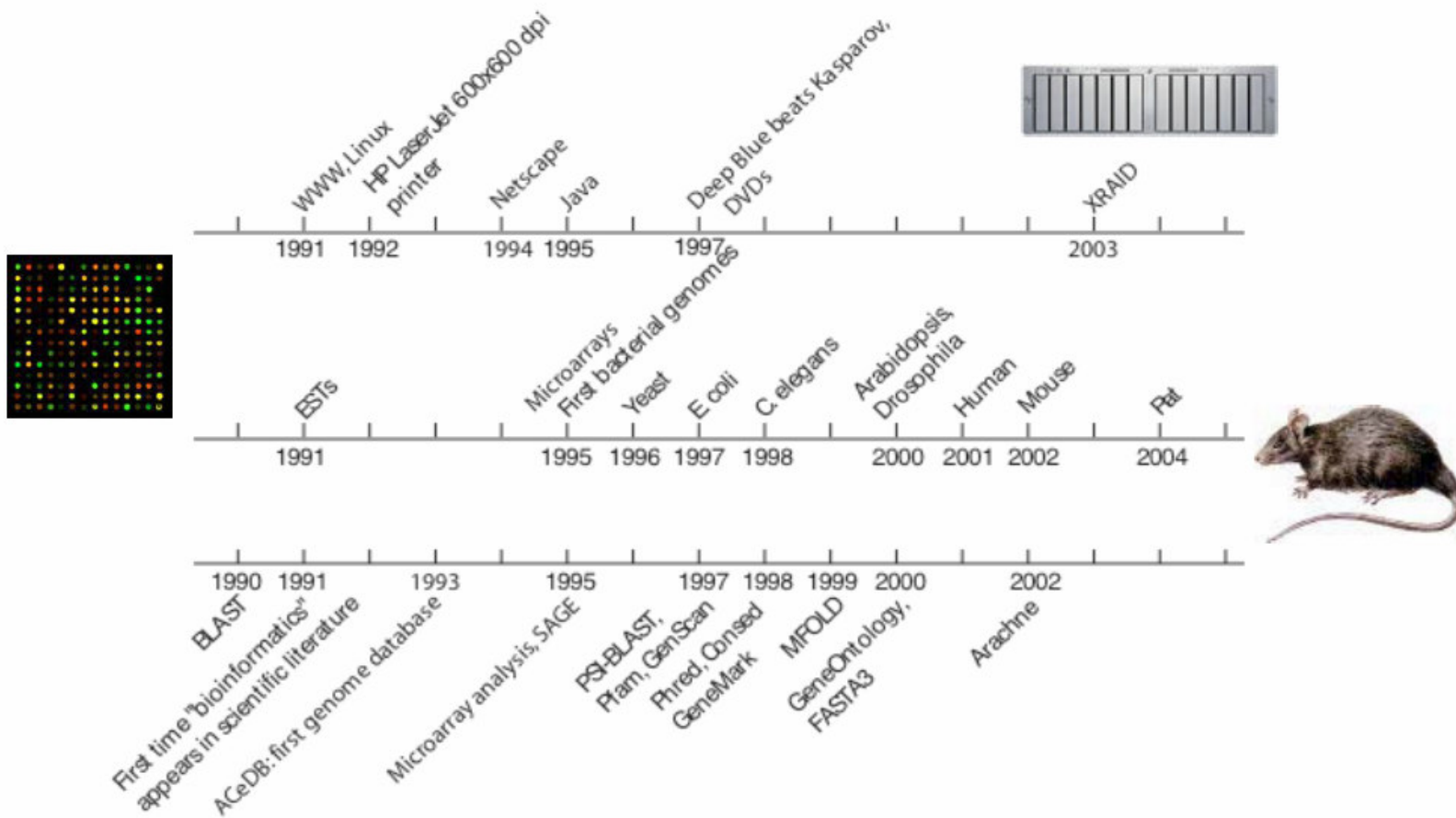


Paul Berg



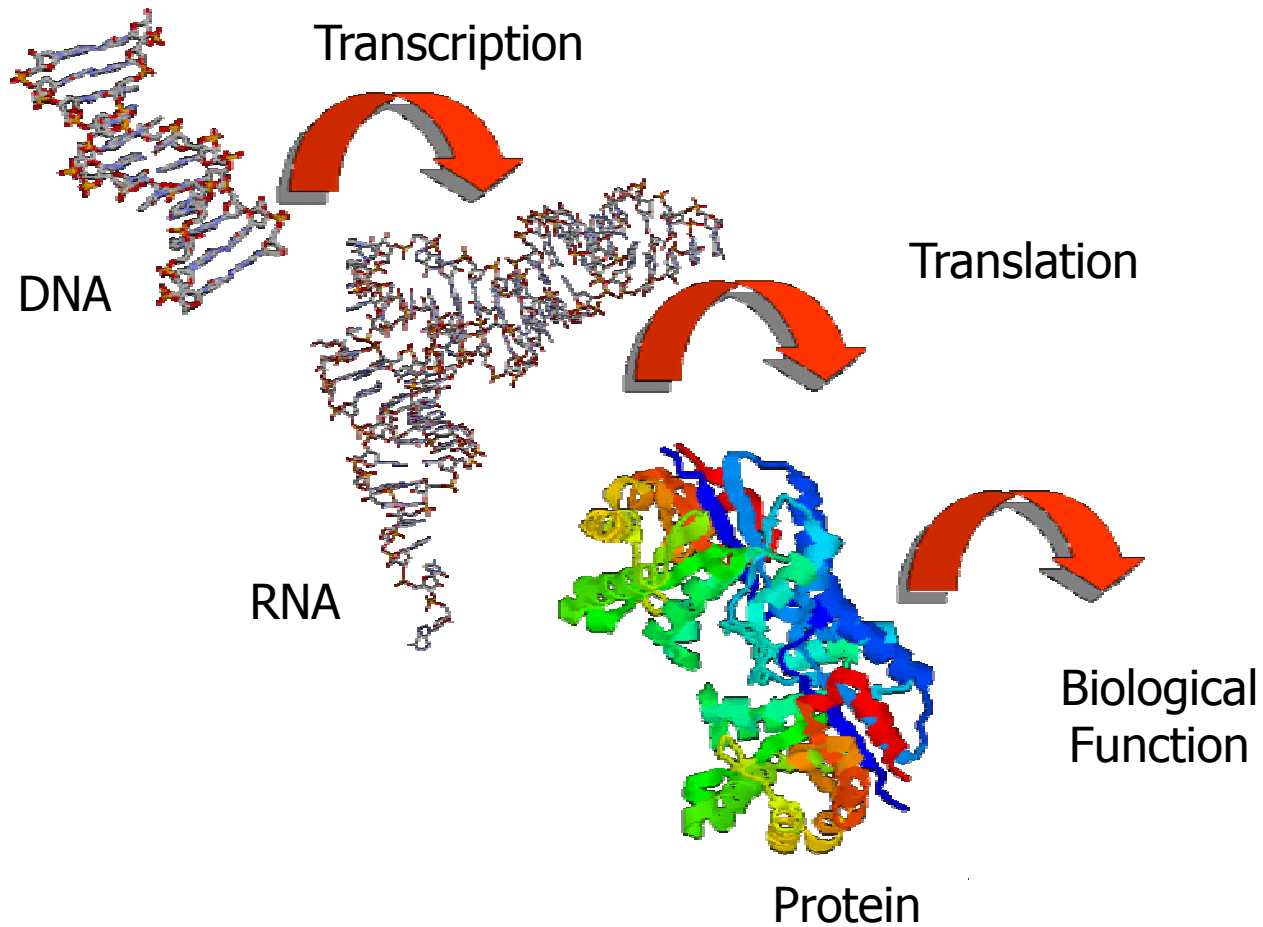
Thanks to Luce Skrabanek

History of Bioinformatics



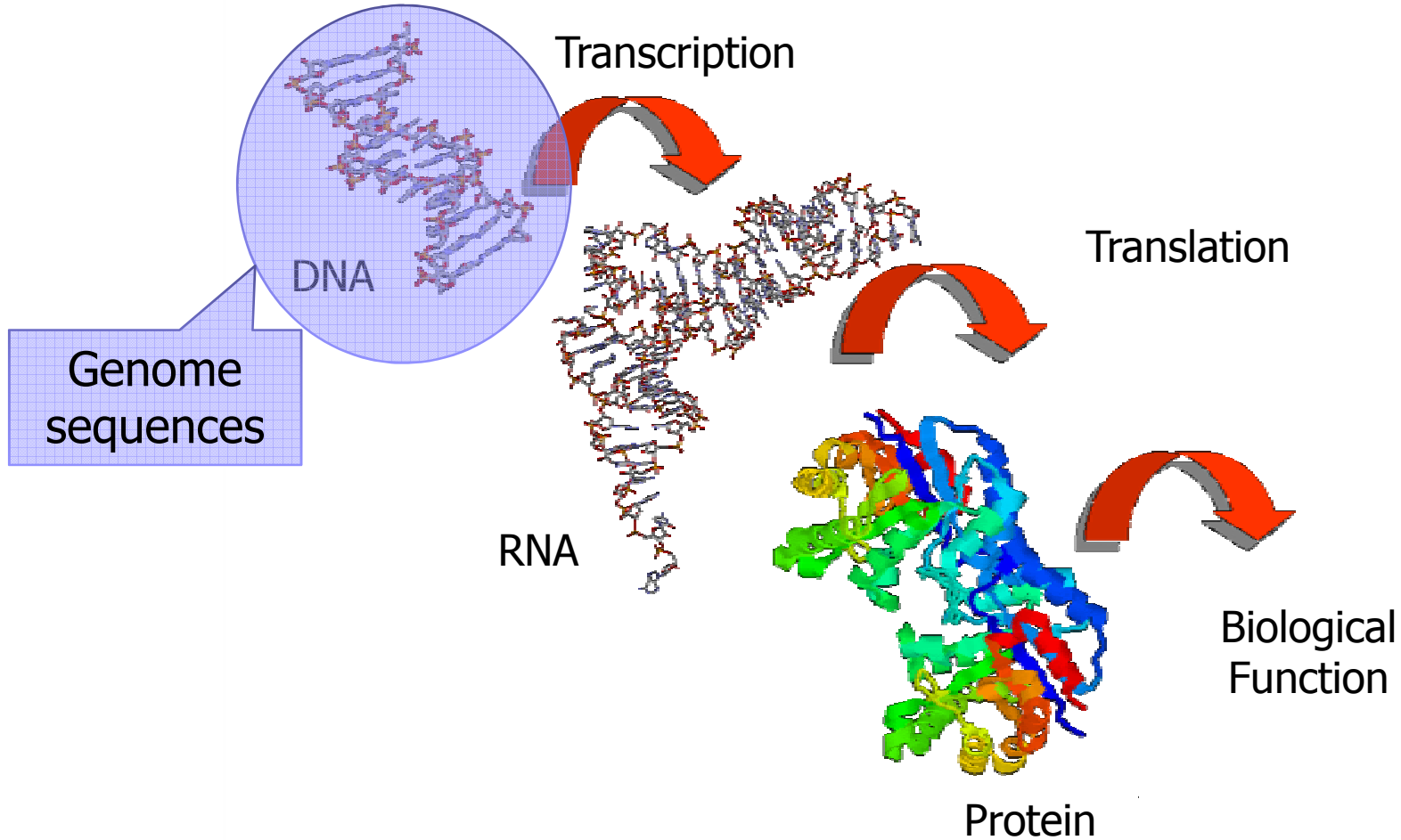
Thanks to Luce Skrabanek

Biological Data

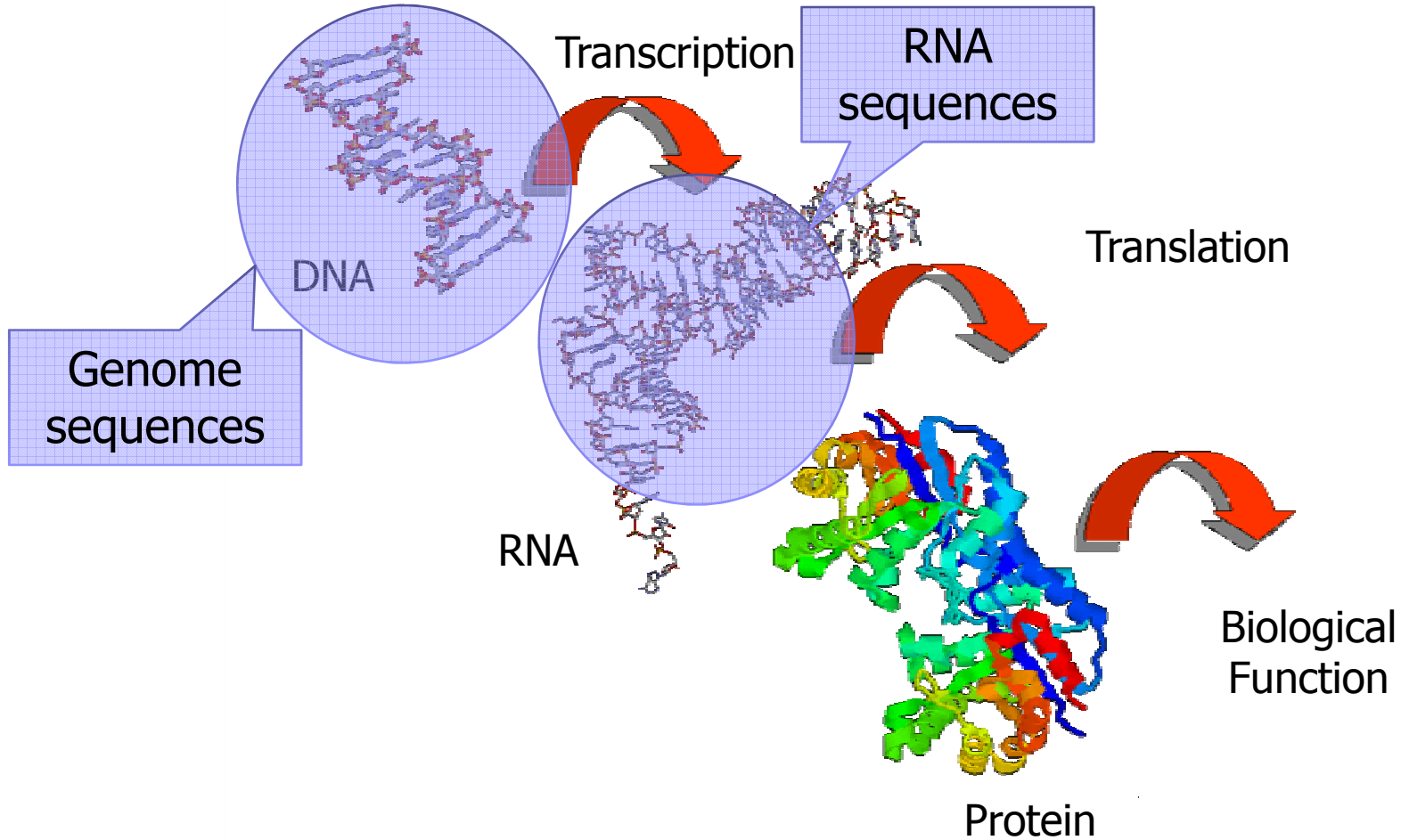


Jacques van Helden, David Gilbert and A.C. Tan, 2003

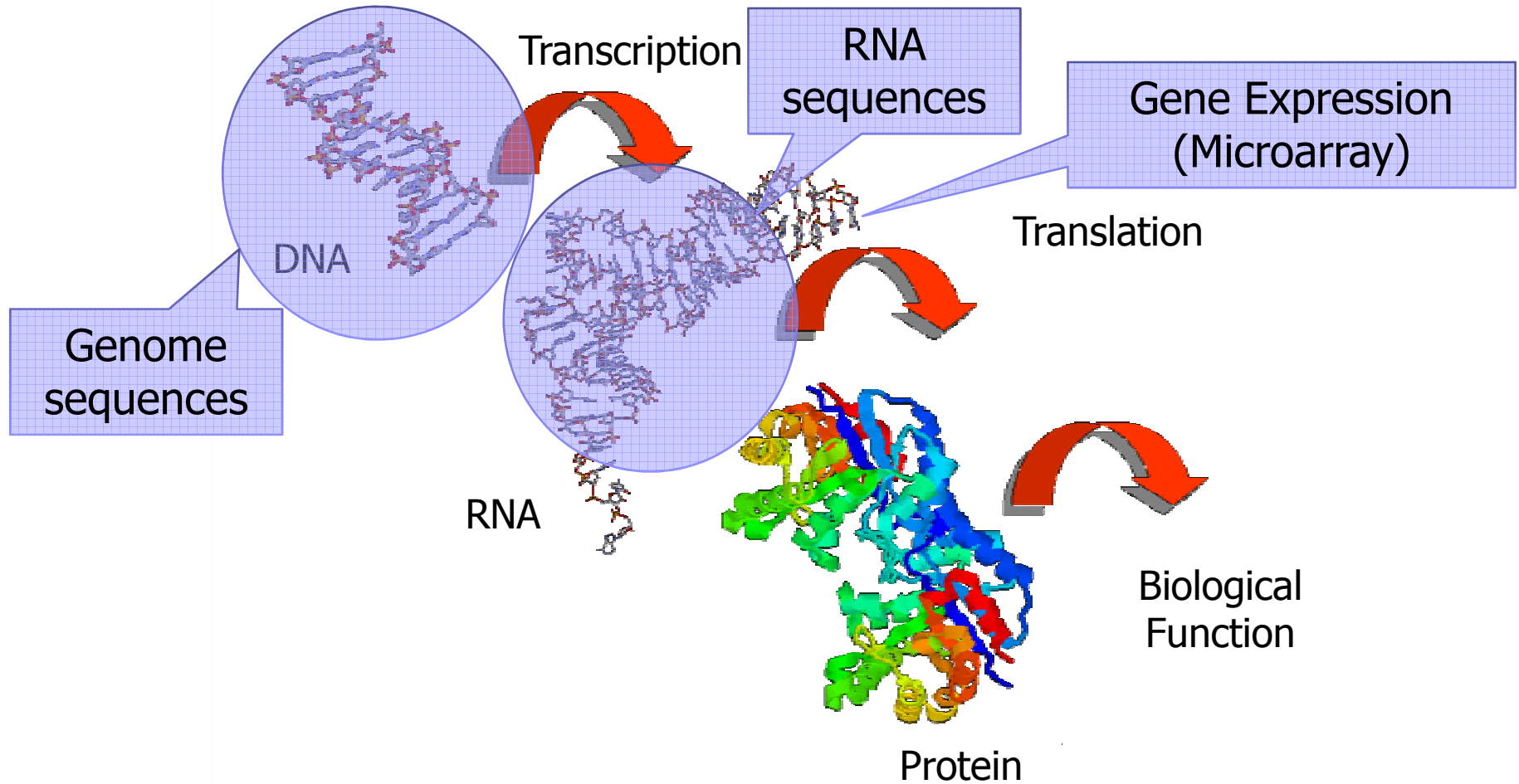
Biological Data



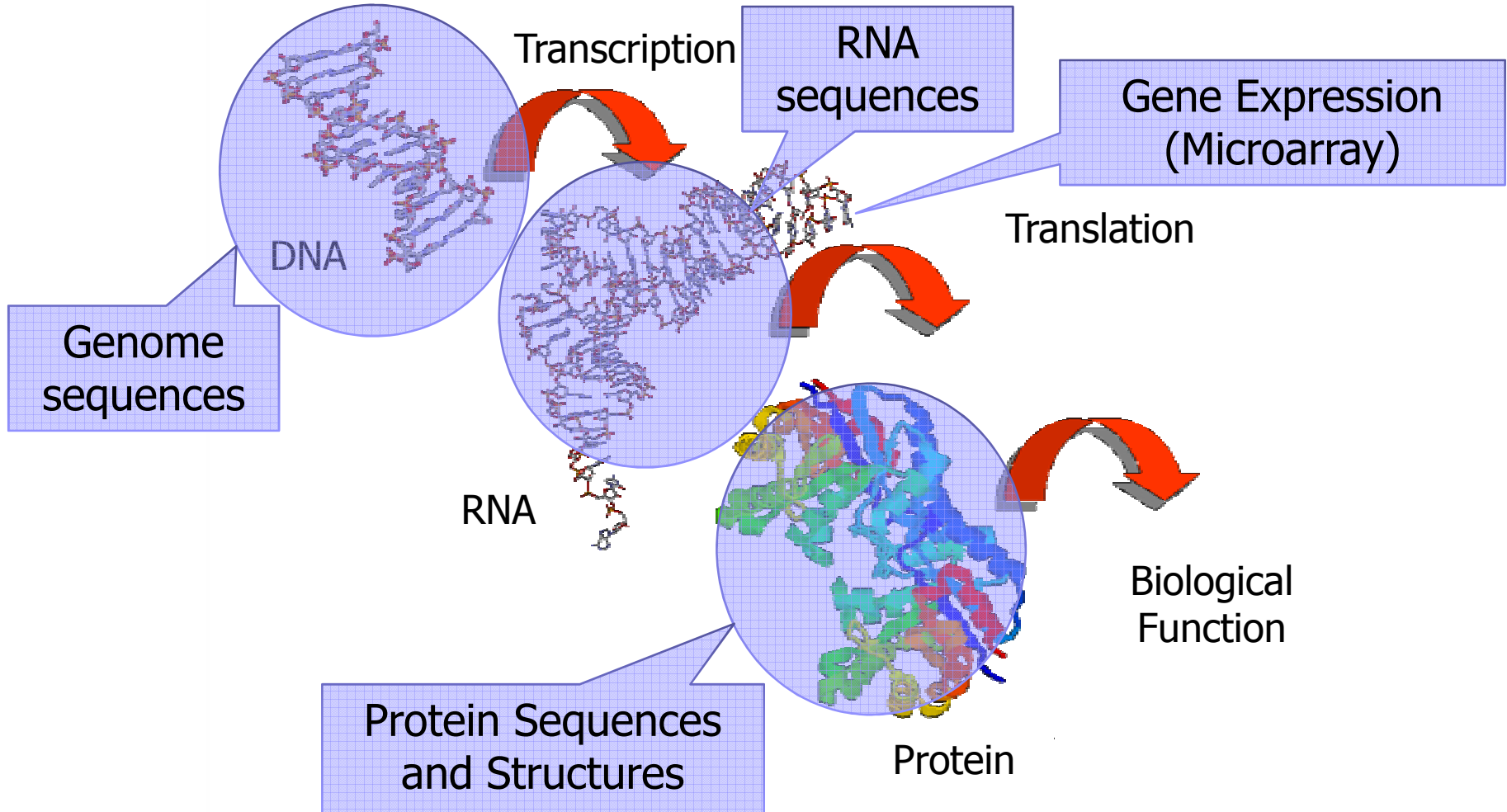
Biological Data



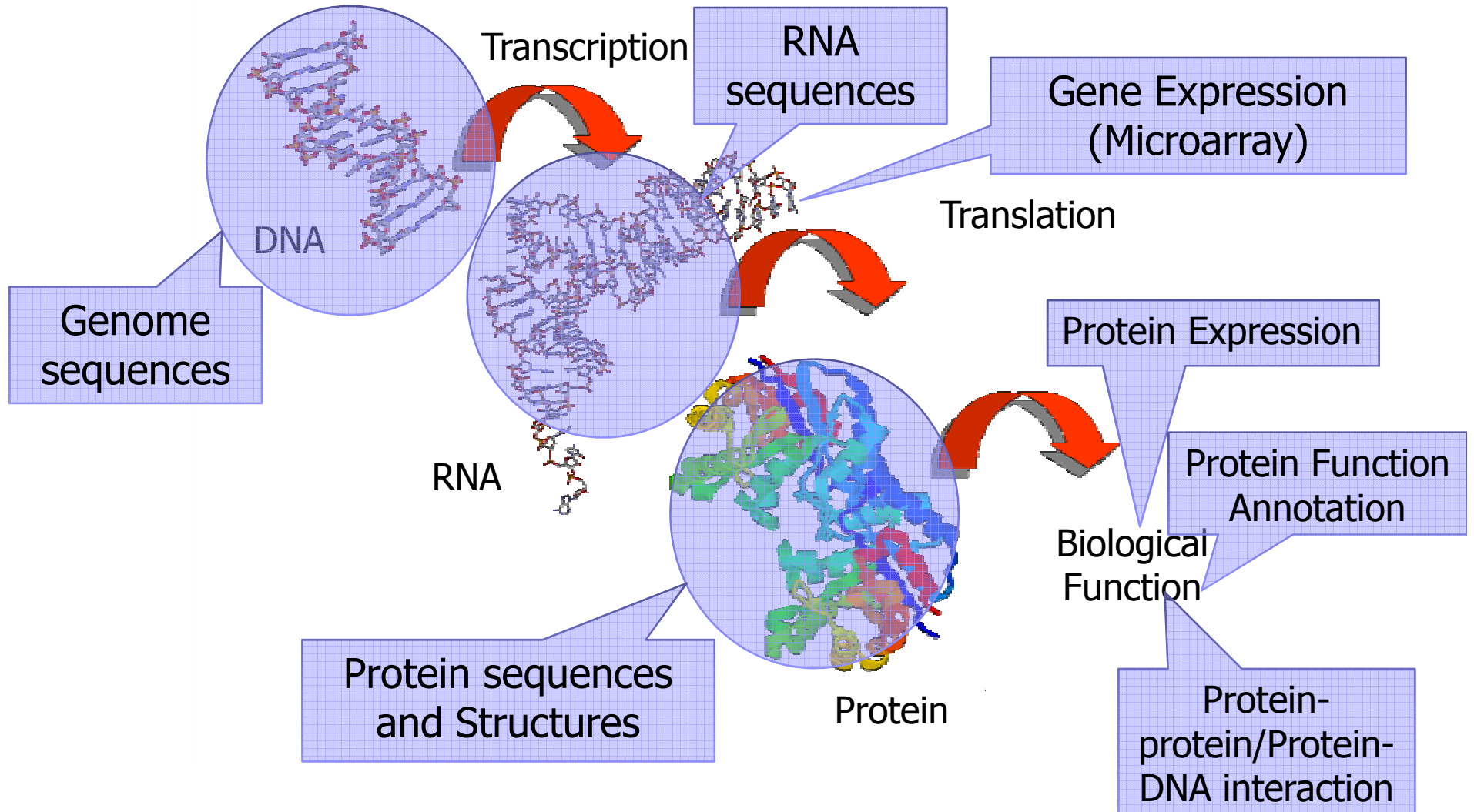
Biological Data



Biological Data



Biological Data

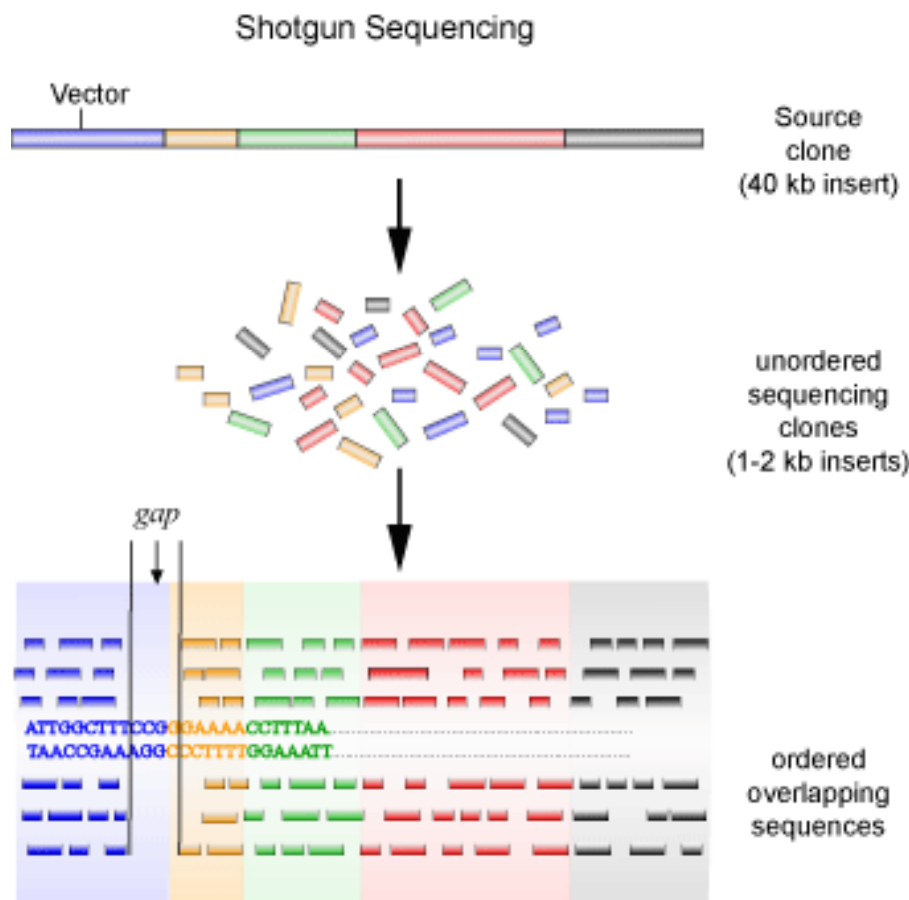




Other Data

- SNPs
- Organism-specific databases
- Genomes
- Molecular pathways
- Scientific literature
- Disease information
-

Combinatory Algorithms



- Get multiple copies of DNA segments.
- Alignment the segments to reconstruct the sequence.
- Closing the GAP with slow and expensive experiments.
- Combinatory algorithms for closing the gap with minimal number of pool tests.

CSCI8980: Applied Machine Learning in Computational Biology

Inferring Gene Regulatory Network with Bayesian Networks

Rui Kuang

Department of Computer Science and Engineering

University of Minnesota

kuang@cs.umn.edu

UNIVERSITY OF MINNESOTA

Twin Cities • Duluth • Morris • Crookston • Rochester • Other Locations





Gene Regulatory Networks

- *Gene regulatory networks*: switching on and off of genes by regulation of transcriptional machinery
- *Learning problem*: Model gene regulatory behavior using genome-wide data, extract hypotheses for wet lab testing
- *Descriptive models*, such as probabilistic graphical models, linear network models, clustering, are interpretable models to training data.
- Can check if local components of model reflect known biological mechanisms.

Gene Regulation

- *Regulatory proteins (transcription factors)* bind to non-coding *regulatory sequence (promoter)* of a gene to control rate of transcription

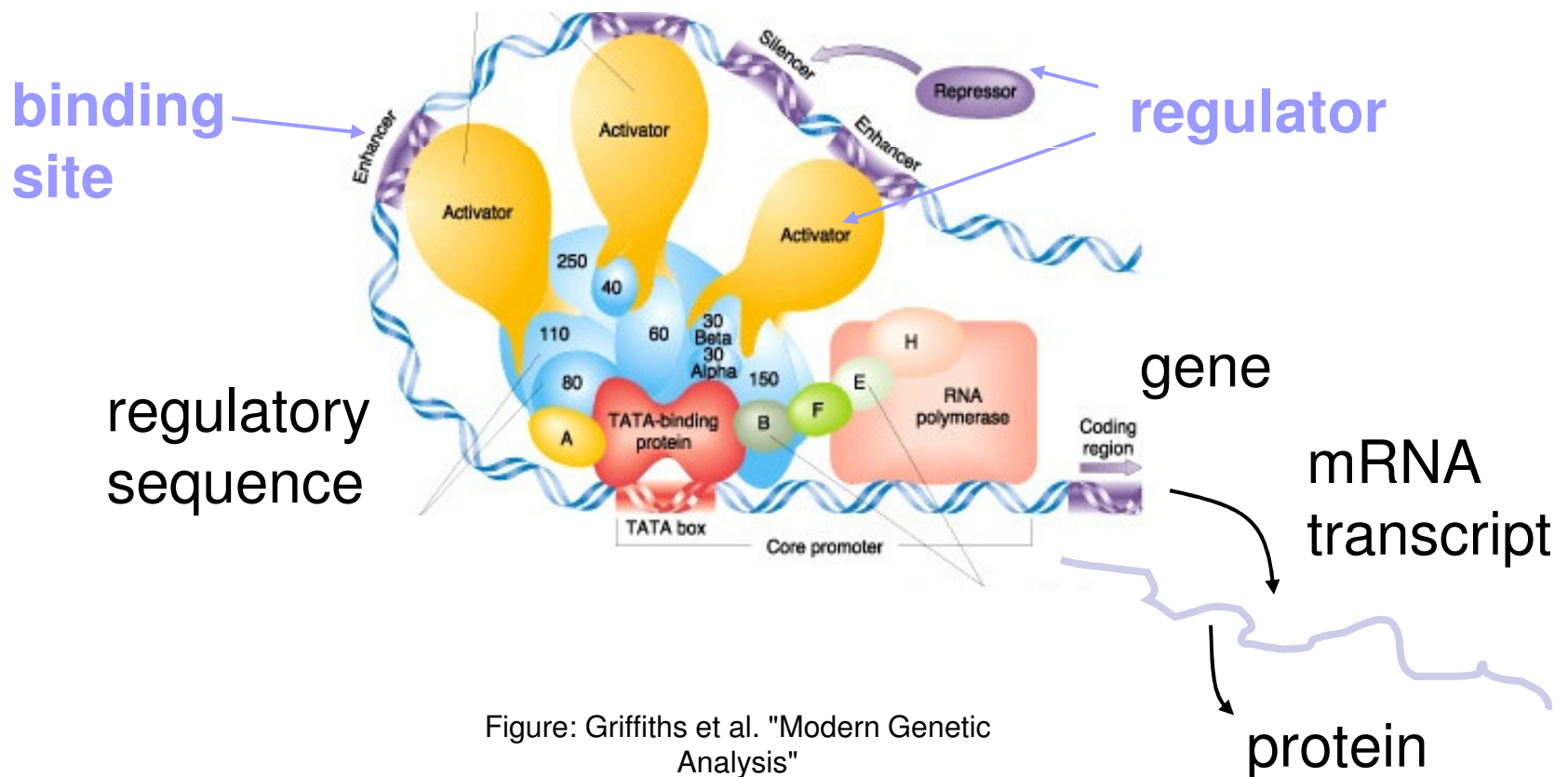


Figure: Griffiths et al. "Modern Genetic Analysis"

Gene Regulation

- *Regulatory proteins (transcription factors)* bind to non-coding *regulatory sequence (promoter)* of a gene to control rate of transcription

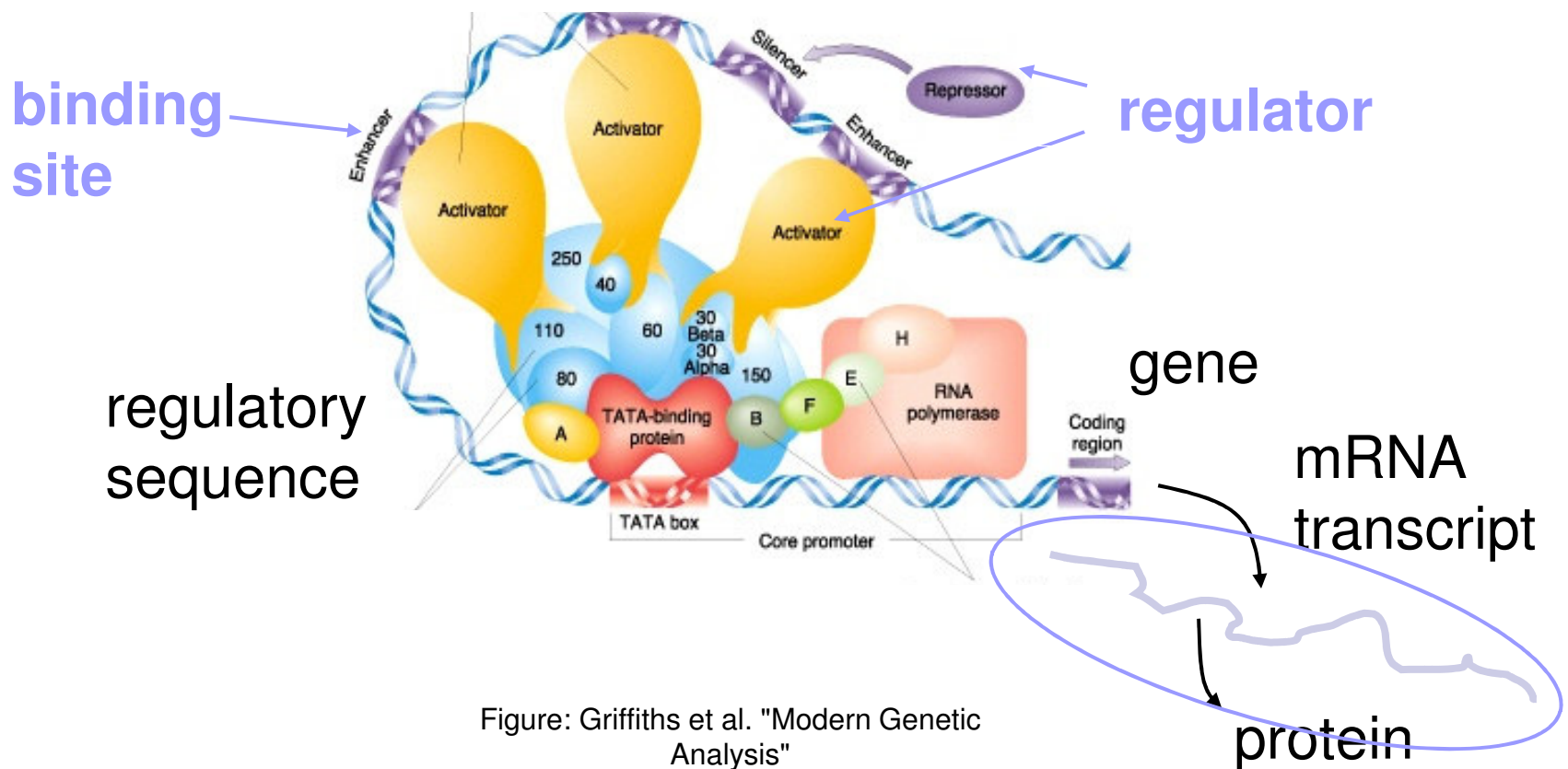
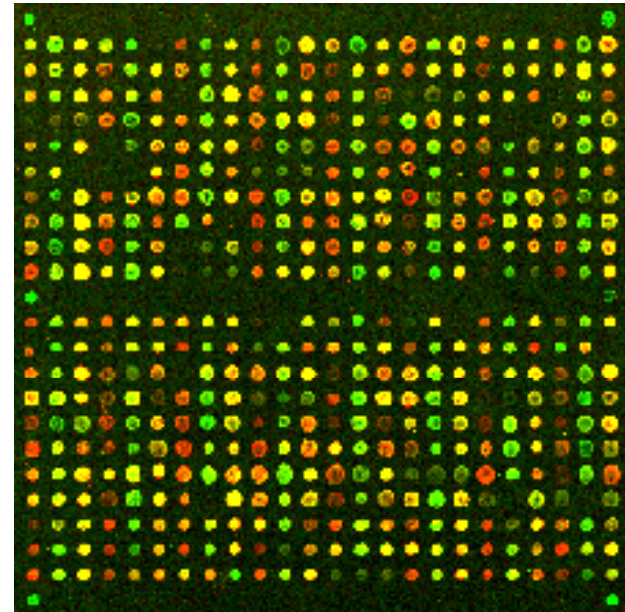


Figure: Griffiths et al. "Modern Genetic Analysis"

Genome-wide Expression Data

- *Microarray* (and other high-throughput) technologies measure *mRNA transcript expression levels* for 1000s of genes at once
- Noisy and sparse data
- Snapshot of the cellular system: *transcriptome*, i.e. protein expression not observed
- Difficult to infer regulatory relation between genes.



Regulatory Components in yeast

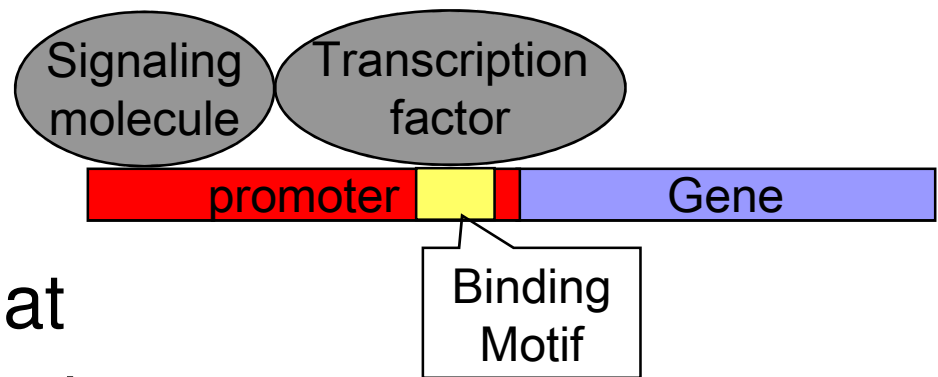
- For simple organisms like *yeast* (*S. cerevisiae*), previous studies and data sources the components needed in model:

- Known and putative *transcription factors*

- Signaling molecules* that activate transcription factors

- Known and putative *binding site “motifs”* in promoter regions

- In yeast, regulatory sequence = 500 bp upstream region





Analyze Gene Expression Data

■ Clustering

- Groups genes with similar expression patterns
- The gene clusters do not reveal the regulatory structure of the genes

■ Boolean Networks

- Deterministic models of the logical interactions between genes
- Gene is in either on state or off state
- Not feasible to learn from microarray data

■ Bayesian Networks

- Measure expression level of each gene
- Gene as random variables affecting on others
- Can possibly include other random variables, such as external stimuli, environment parameters, and biological factors



Model Validation of Genetic Regulatory Networks

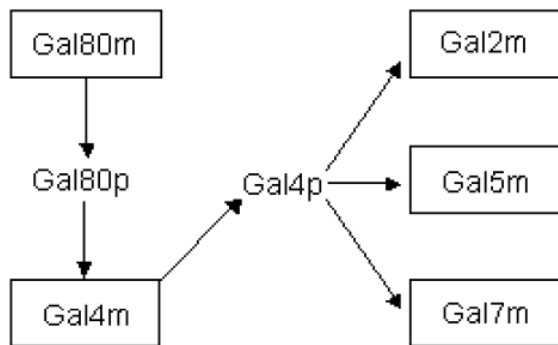
- Using *Bayesian scoring* metric to choose the right network structure

$$\begin{aligned} \text{BayesianScore}(S) &= \log p(S | D) \\ &= \log p(S) + \log p(D|S) + c, \end{aligned}$$

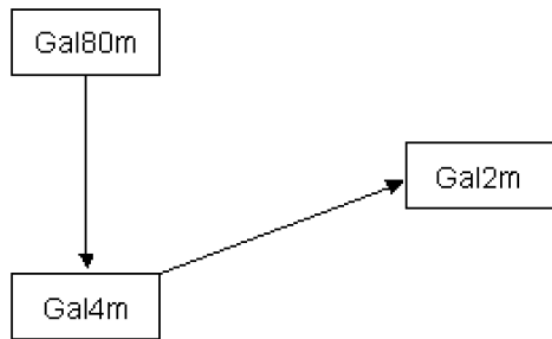
where $p(D|S)$ is the likelihood function and $P(S)$ is a prior on the model S .

- Validated on the *galactose system* in *S. cerevisiae*
- Expression data: 52 genomes worth of Affymetrix GeneChip expression data

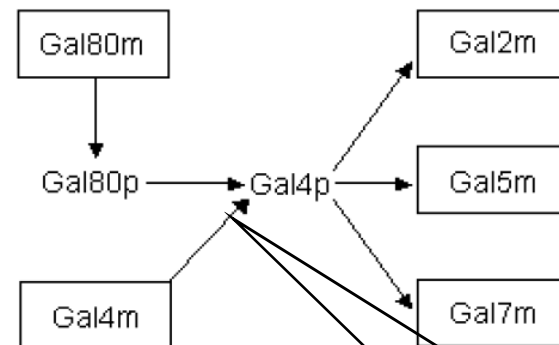
Hypothesis of Galactose System



M1

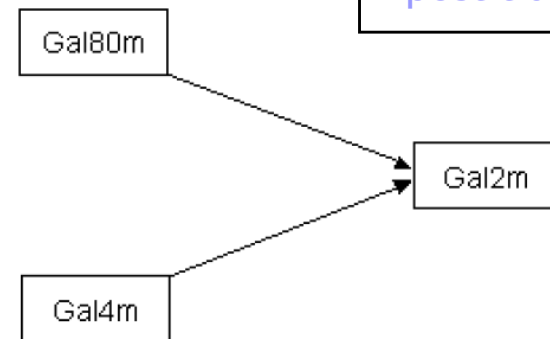


M1



M2

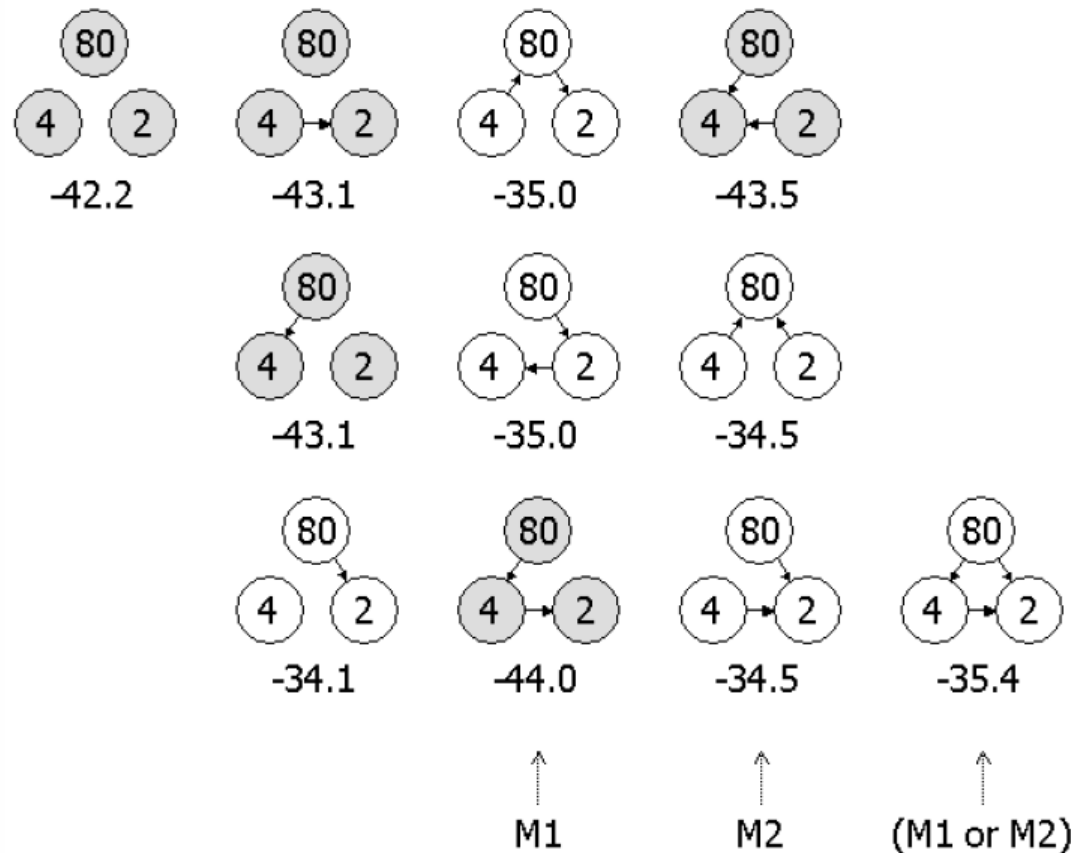
Gal80p inhibits Gal4p post-translationally



M2

Scoring Possible Structures

- Binary quantization of gene expression into up/down (3 binary random variables)



Scoring Possible Structures

- Binary quantization of gene expression into up/down (3 binary random variables)

