

## Lecture 1 - An Introduction to Bioinformatics

A5ASC3.1	14	SIKLWPPSQTRLLLVERMANNLST..PSIFTRK..YGSLSKEEARENQIEEVACSTANQ.....HYEKEPDGDGGSAVQLYAKECSKLILEVLK	101
B4F917.1	13	SIKLWPPSESTRIMLVDRMTNNLST..ESIFSRK..YRLLGKQEAHENAKTIEELCFALADE.....HFREEPDGDGSSAVQLYAKETSKMMLEV	100
A9S1V2.1	23	VFKLWPPSQGTREAVRQKMALKLSS..ACFESQS..FARIELADAQEHARAIEEVAFGAAQE.....ADSGGDKTGSAVVMVYAKHASKLML	109
B9GSN7.1	13	SVKLWPPGQSTRMLVERMTKNFIT..PSFISRK..YGLLSKEEAEEDAKKIEEVAFAAANQ.....HYEKQPDGDGSSAVQIYAKESSRLM	100
Q8H056.1	30	SFSIWPPPTQRTRDAVVRRLVDTLGG..DTILCKR..YGAVPAADAEPAARGIEAEAFDAAA..SGEAAATASVEEGIKALQLYSKEVSRRLLDF	120
Q0D4Z3.2	44	SLSIWPPSQRTRDAVVRRLVQTLVA..PSILSQR..YGAVPEAEAGRAAAAVERAEAYAAVTES..SSAAAAAPASVEDGIEVLQAYSKEVSRRLL	135
B9MVW8.1	56	SFSIWPPPTQRTRDAIISRLIETLST..TSVLSKR..YGTIPKEEASEASRRIEEEAFSGAST.....VASSEKDGLLEVQLYSKEISKRMLET	141
Q0IYC5.1	29	SFAVWPPTRRTRDAVVRRLVAVLSGDTTTALRKRYRYGAVPAADAERAARAVEAQAFDAASA....SSSSSSVEDGIETLQLYSREVSNRLLA	121
A9NW46.1	13	SIKLWPPSESTRMLVERMTDNLSS..VSFFSRK..YGLLSKEEEAENAKRIETAFLAAND.....HEAKEPNLDDSSVVQFYAREASKLMLEALK	100
Q9C500.1	57	SLRIWPPTQKTRDAVLNRRIETLST..ESILSKR..YGTLSDDATTVALIEEEEAYGVASN.....AVSSDDDGKIKILELYSKEISKRMLES	142
Q2HRI7.1	25	NYSIWPPKQRTRDAVKNRLIETLST..PSVLTKR..YGTMSADEASAAAIQIEDEAFSVANA.....SSSTSNDNVTILEVYYSKEISKRMIE	110
Q9M7N3.1	28	SFKIWPPPTQRTRAVVRRLVETLTS..QSVLSKR..YGVVIPEEDATSAAARIIEEEAFSVASV..ASAASSTGGRPEDEWIEVLHIYSQEIXQ	119
Q9M7N6.1	25	SFSIWPPPTQRTRDAVINRLIESLST..PSILSKR..YGTLPQDEASETARLIEEEEAFAAAGS.....TASDADDGIEILQVYSKEISKRMID	110
Q9LE82.1	14	SVKMWPPSKSTRMLVERMTKNITT..PSIFSRK..YGLLSVEEAEQDAKRIEDLAFATANK.....HFQNEPDGDGTSAVHVYAKESSKLM	101
Q9M651.2	13	SIKLWPPSLPTRKALIERITNNFSS..KTIFTEK..YGSLTKDQATENAKRIEDIAFSTANQ.....QFEREPDGDGGS	100
B9R748.1	48	SLSIWPPPTQRTRDAVITRLIETLSS..PSVLSKR..YGTISHDEAESARRIEDEAFGVANT.....ATSAEDDGLIEILQLYSKEISRRMLD	133

Ian Simpson

[ian.simpson@ed.ac.uk](mailto:ian.simpson@ed.ac.uk)

# Learning Outcomes

By the end of our course you should be able to:

- Communicate between biological and computational domains to facilitate effective inter-disciplinary working.
- Use and/or implement Bioinformatics tools, services and software in practical research scenarios.
- Have sufficient background knowledge, skills and understanding to discover and apply additional bioinformatics techniques.

## Foundations

Bioinformatics  
Algorithms

Biological  
Concepts

Data Types  
& Sources

## Skills

Python  
Scripting

Using  
Web Tools &  
Services

Building Analytic  
Workflows

## Application

Case  
Studies

Coursework  
Mini-projects

Inter-disciplinary  
communication in  
practice

# Course Organisation

## Teaching

### Lectures

Weeks 1-10 on Friday 10:00-10:55

Lecture Theatre 1 Appleton Tower

Also streamed live and recorded

Slides released >24h prior to lecture

### Computing Lab

Weeks 2, 4, 6, and 8 on Wednesday 14:00-15:55

Appleton Tower Lab 6.06

### Reading List

(called Resource List with all reading material)

### Materials

Everything available on courseware site:-

<https://opencourse.inf.ed.ac.uk/bio1>

## Assessment

- 2x Coursework Assignments, 50% each
- Submission via Gradescope
- CW1 released 06/10/23, submit 27/10/23
- CW2 released 30/10/23, submit 17/11/23
- dedicated discussion channel on Piazza

## Optional Mini-Course

### Introduction to Python

- Optional catch-up mini-course for those who don't have prior experience of scripting and/or Python
- practice computational notebooks
- Small challenge coding puzzles to aid reading list
- dedicated discussion channel on Piazza

# Bioinformatics 1

## Communication

Discussion Board on Piazza and Course Announcements

# Course Topics

Week	Labs	CW	Date Commencing	Week Topic Title
1			18th September	Introduction to Bioinformatics
2			25th September	Pairwise Sequence Alignment
3*		6th	2nd October	Basic Local Alignment and Search Tool (BLAST)
4			9th October	Multiple Sequence Alignment
5			16th October	Exploring Biological Databases
6		27th	23rd October	Working with Biological Databases
7		30th	30th November	Ontologies & Functional Enrichment Analysis
8			6th November	Biological Network Analysis
9		17th	13th November	Introduction to Next Generation Sequencing
10			20th November	Course Wrap-Up, Q&A

\*NB pre-recorded lecture

# Objectives of the Coursework

## Aim

- to apply methods
- to develop research skills
- gain experience using bioinformatics sources and literature
- learn scientific writing and reporting

## Task

- questions to rehearse key concepts
- self-directed research
- strong focus on precise description of methods and sources
- experimentation
- analytical, showing insight and evidence of evidence gathering to support conclusions

# The Extended Common Marking Scheme

<https://web.inf.ed.ac.uk/infweb/student-services/ito/students/common-marking-scheme>

Grade	Mark	Degree Award	Description
A1	90-100	1st class or MSc with distinction	<b>Excellent.</b> Outstanding in every respect, the work is well beyond the level expected of a competent student at their level of study. The work should meet the criteria for an A2 grade and should also evidence a clear understanding of the limits of the state of knowledge, and their consequences, for the topic at hand.
A2	80-89	1st class or MSc with distinction	<b>Excellent.</b> Outstanding in some respects, the work is often beyond what is expected of a competent student at their level of study. Demonstrates that the student is actively extending their knowledge and capacity well beyond required materials and making new connections independently: for example, by showing a strong grasp of a range of related materials that are optional or not directly provided, or by demonstrating unusual creativity, depth of analysis, or synthesis with other areas of study.
A3	70-79	1st class or MSc with distinction	<b>Excellent.</b> Very good or excellent in most respects, the work is what might be expected of a very competent student. It indicates that the student has an excellent grasp of the required materials for the course, and may have demonstrated some limited knowledge of or fluency with additional optional materials, if provided.
B	60-69	2(i) or MSc with merit	<b>Very Good.</b> Good or very good in most respects, the work displays thorough mastery of the relevant learning outcomes.
C	50-59	2(ii) or MSc	<b>Good.</b> The work clearly meets requirements for demonstrating the relevant learning outcomes.
D	40-49	3rd class or PG Diploma/Cert	<b>Pass (undergraduate or Diploma level).</b> The work meets minimum requirements for demonstrating the relevant learning outcomes. A satisfactory performance for undergraduate
E	30-39	Fail	<b>Marginal Fail.</b> The work fails to meet minimum requirements for demonstrating the relevant learning outcomes.
F	20-29	Fail	<b>Clear Fail.</b> The work is very weak and/or incomplete in important respects.
G	10-19	Fail	<b>Bad Fail.</b> The work is extremely weak or mostly incomplete/absent.
H	0-9	Fail	<b>Bad Fail.</b> The work is absent or of very little, if any, consequence to the area in question.

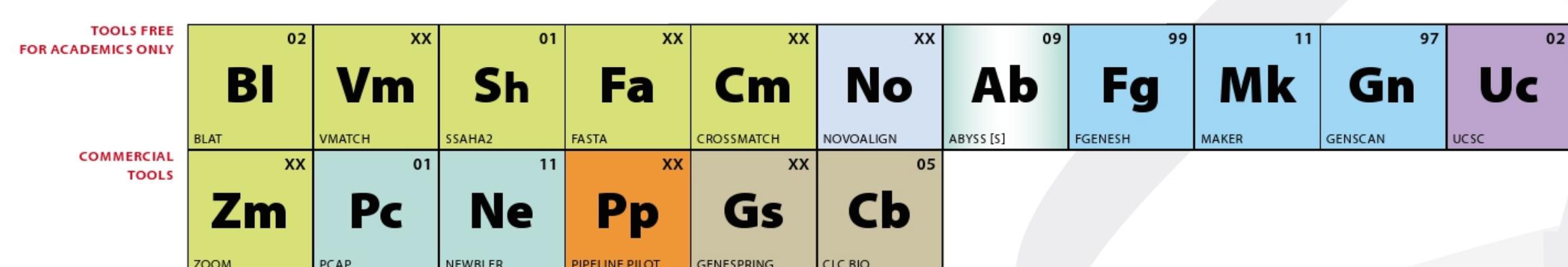
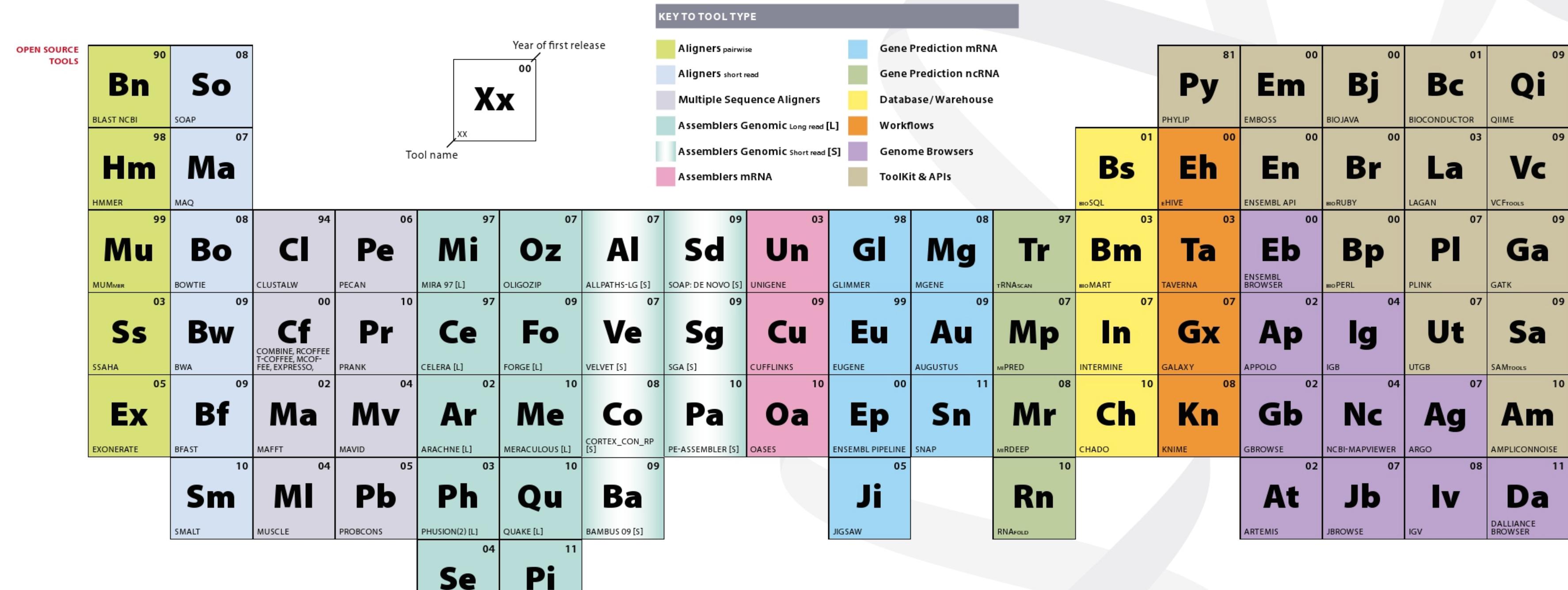
# Computing Labs

- Structured sessions putting into practice topics we've covered in lectures
- Weeks 2, 4, 6, & 8, @14:00 AT6.06
- Python - Jupyter Notebooks & Worksheets
- Q&As in lab and on site discussion boards
- Opportunity to meet and discuss the topics with your peers and demonstrators
- Asynchronous working possible if you cannot make labs
- University Noteable System available 24/7



# The “Elements” of Bioinformatics

Econometrics "The Elements of Biostatistics"



# What is Bioinformatics ?

## bioinformatics, *n.*

View as: Outline | [Full entry](#)

OED | Oxford English Dictionary  
*The definitive record of the English language*

Pronunciation: Brit. /bɪənɪfɔ'mætrɪks/, U.S. /'baɪən̩,ɪnfər'mædɪks/

Etymology: < [BIO-](#) [comb. form](#) + [INFORMATICS n.](#). Compare Dutch *bioinformatica* (1978 or earlier).... ([Show More](#))

With *sing.* concord. The branch of science concerned with information and information flow in biological systems, esp. the use of computational methods in genetics and genomics.

[Categories »](#)

1976 *Acta Biotheoretica* 25 3 We publish papers on the philosophy of biology, biomathematics and bioinformatics.

1987 *Science* 4 Sept. 1108/3 A new research program in bioinformatics. This is intended to bring together research in computing science, structural biology, and molecular genetics.

2001 *N.Y. Times* 4 Jan. B6/2 A leader in cutting-edge fields like bioinformatics, in which computers are used to decipher genes and proteins.

2007 *Wired* May 125/1 Delwart chops all the remaining genetic material into little pieces... He uses specially designed bioinformatics software to check them.

Luscombe NM, Greenbaum D, Gerstein M. What is bioinformatics? A proposed definition and overview of the field. *Methods Inf Med.* 2001;40(4):346-58.

# Branches of Bioinformatics

## "Statistics & Machine Learning"

Mathematical & Statistical Modelling

Dynamical modelling (simulation)

Data integration methods

## "Computational Biology"

High Performance Computing (HPC)

Database, software and algorithm development

Web Services

## "Classical Bioinformatics"

Structural Bioinformatics

Gene Prediction

Motif Prediction & Pattern matching

Gene expression analysis

Genome Annotation

Molecular phylogenetics

## "Systems Bioinformatics"

Systems Biology

Network Biology

Drug Design

Image Analysis

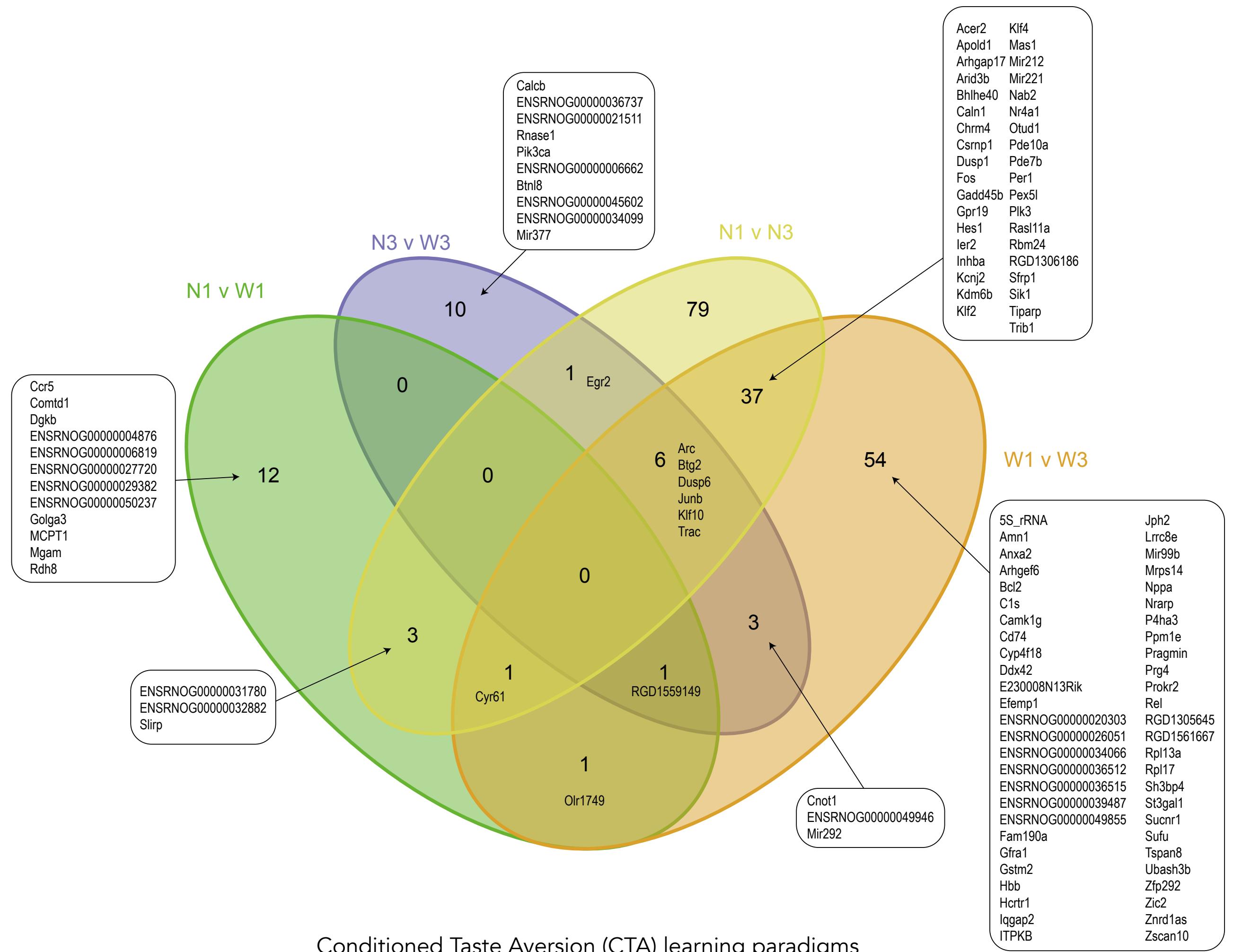
## "Statistical Genetics"

Genome wide linkage and association Studies

Genotype Analysis

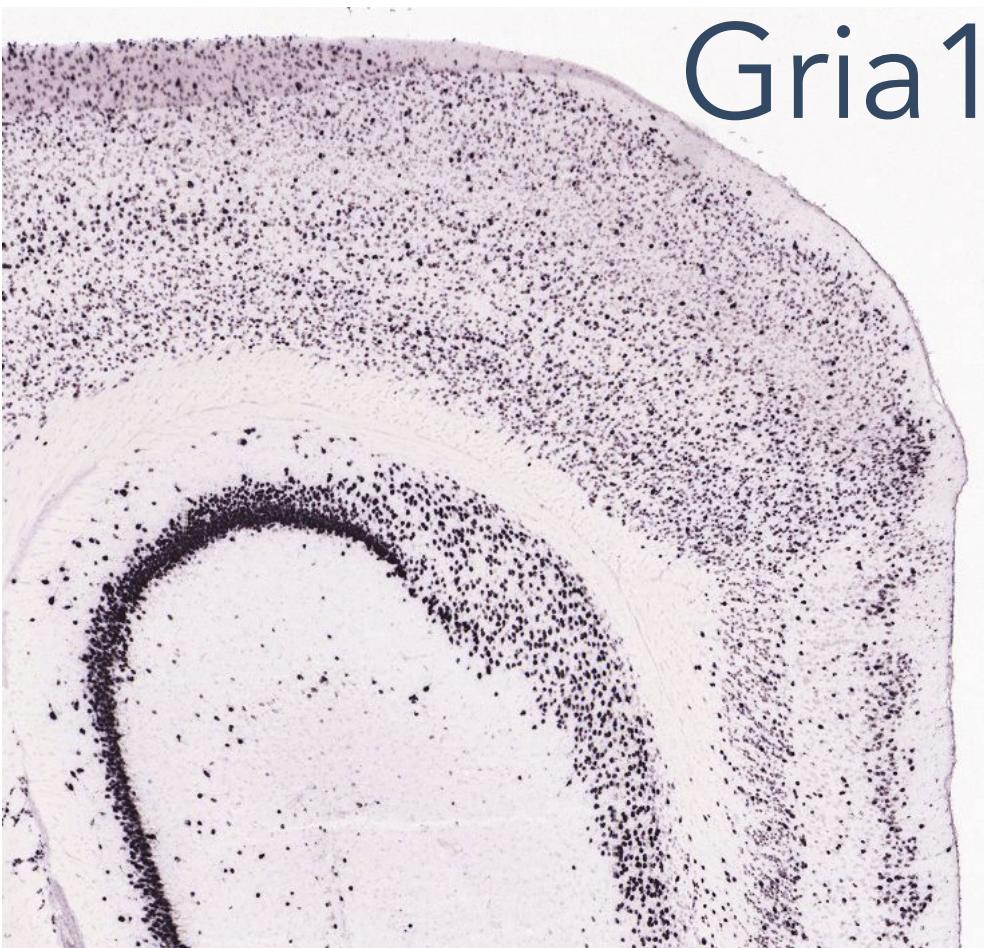
Population Genetics

# Integration of Gene Expression Data

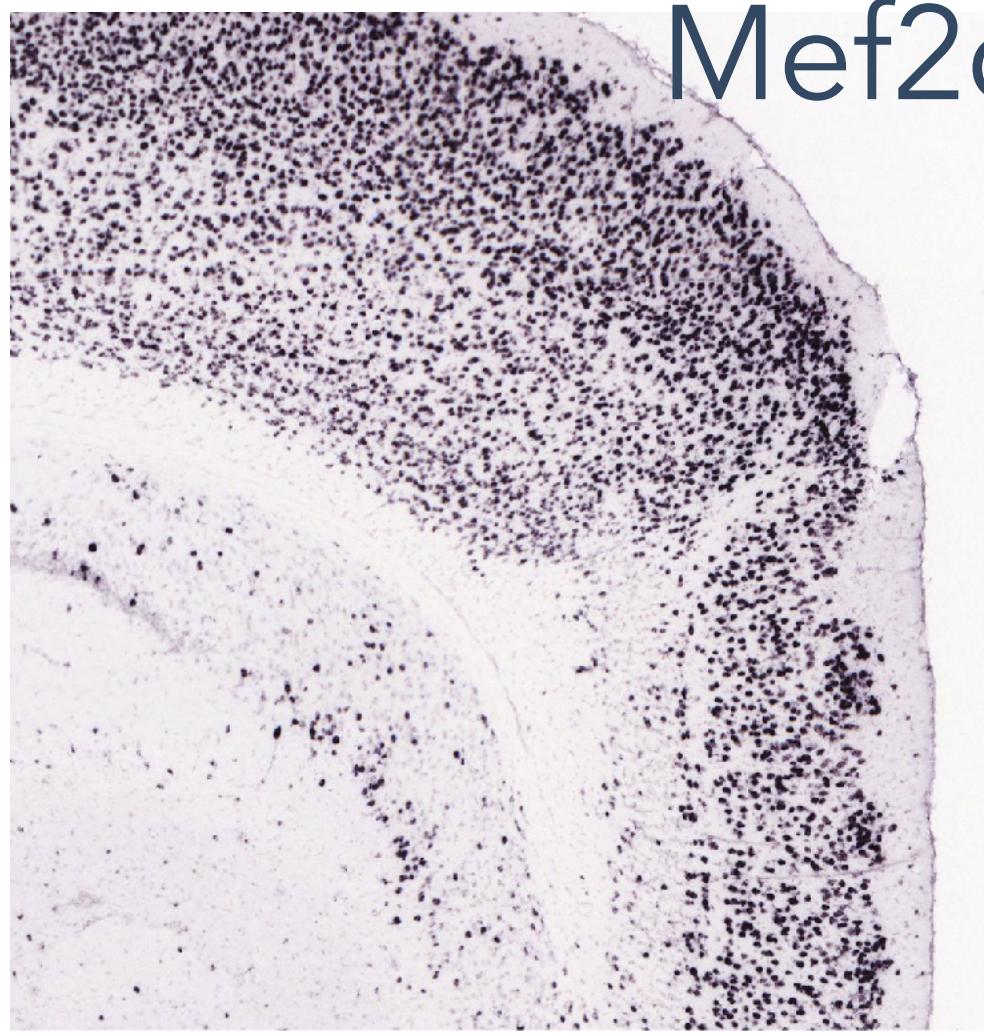


- compared gene expression in the gustatory cortex between rats subject to CTA with or without novel taste
- arrays were sampled 1 and 3 hours post-experience
- we have identified over 200 genes that are differentially expressed in the learning process
- key genes involved in the regulation of synaptic plasticity are amongst the lists and are being validated *in vivo*

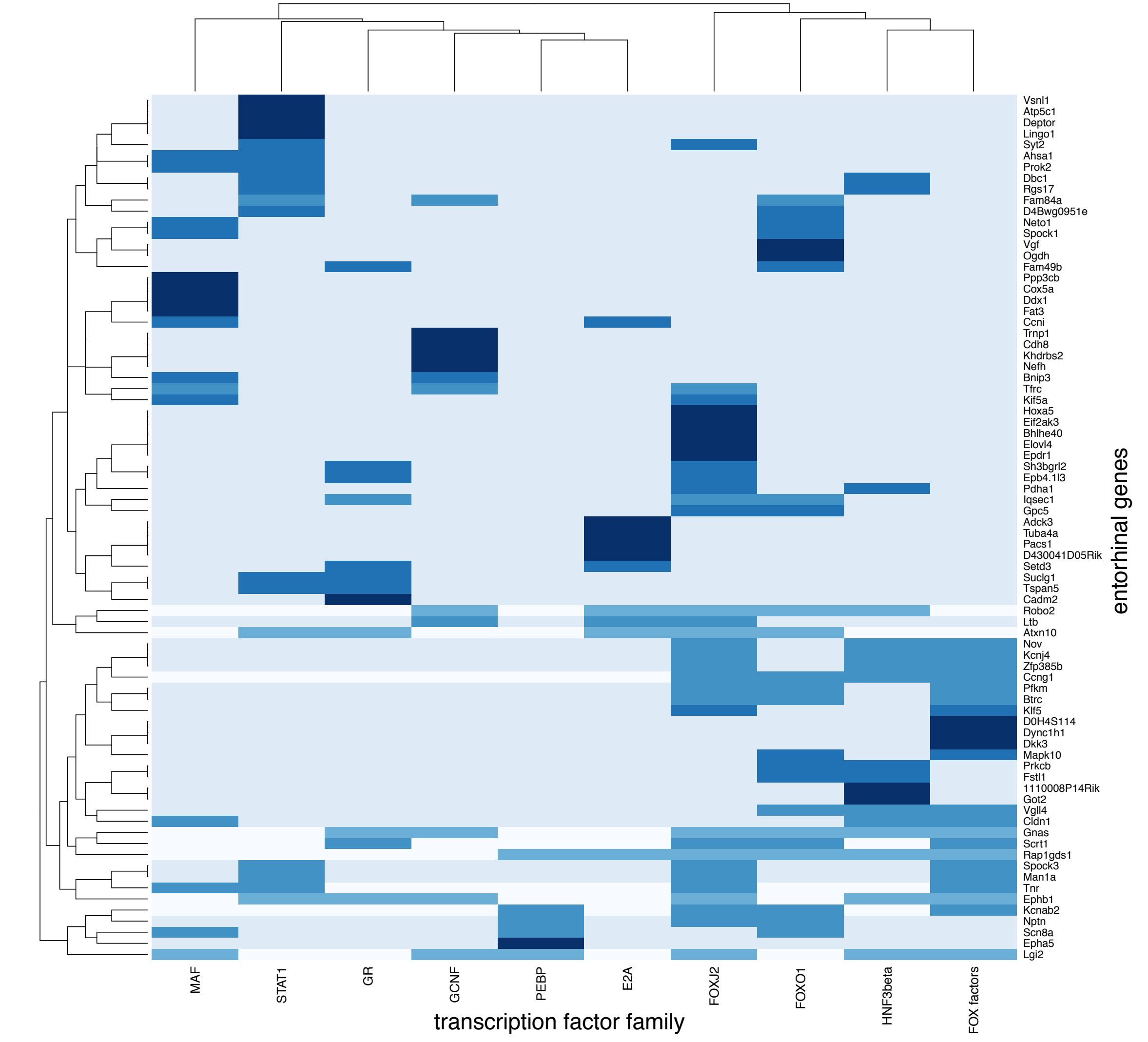
# Gene Expression Regulation in the Brain



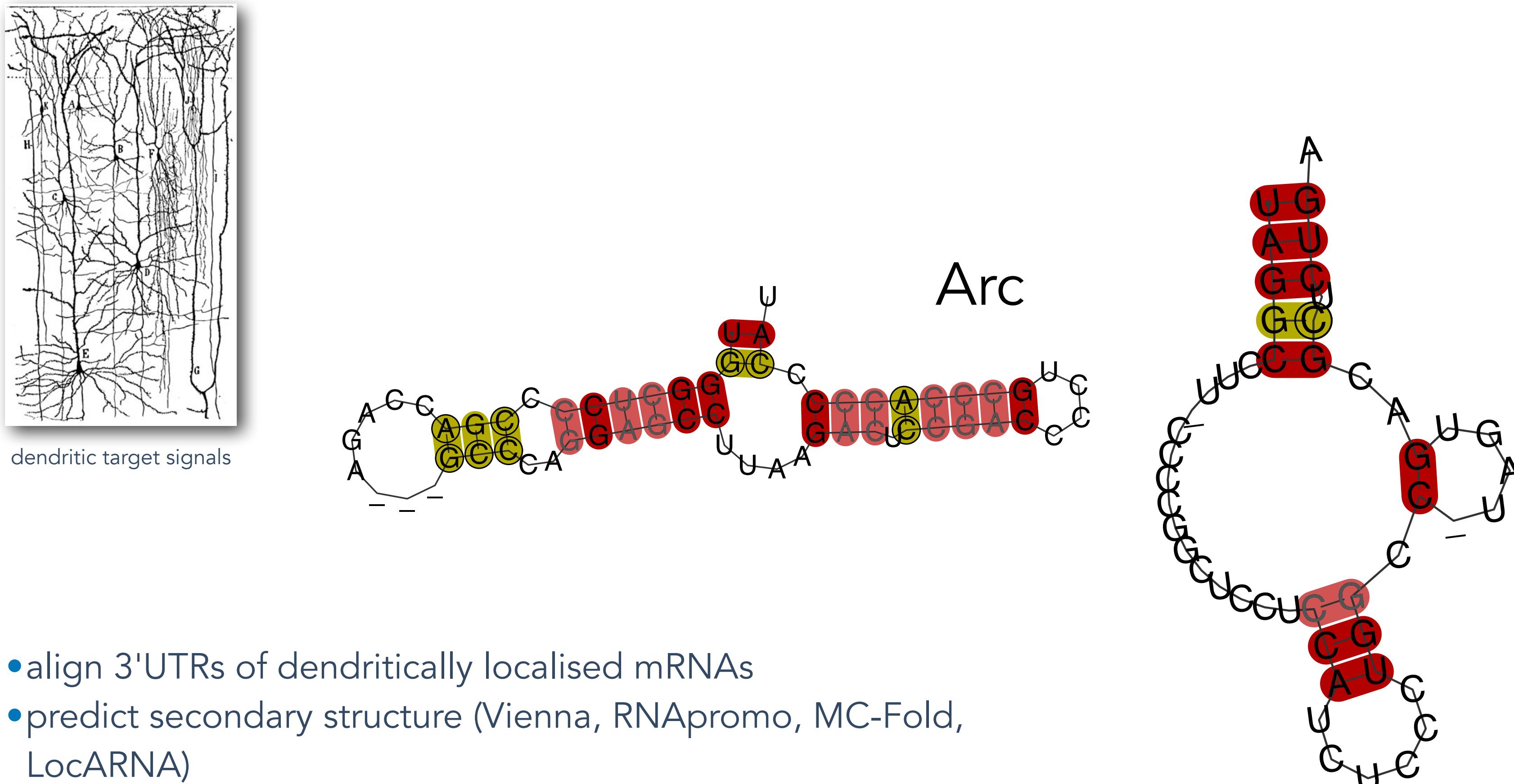
*Gria1*



expression coding in the entorhinal cortex



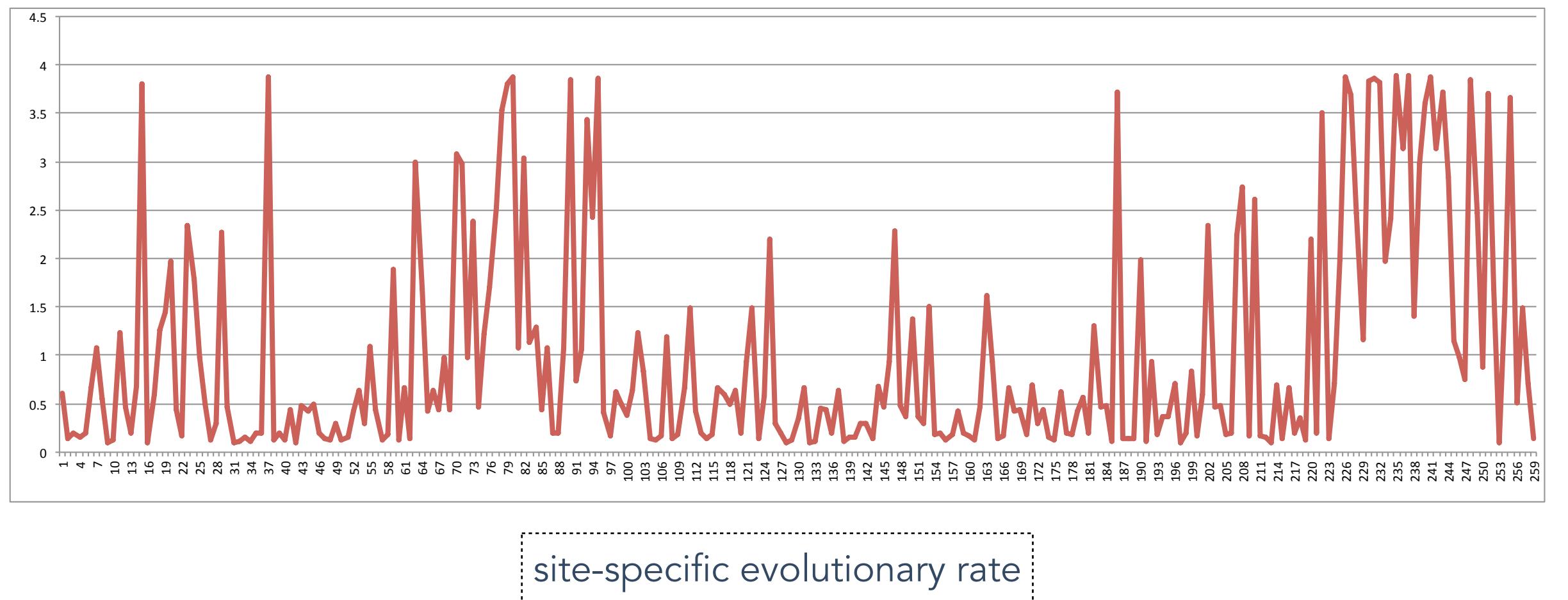
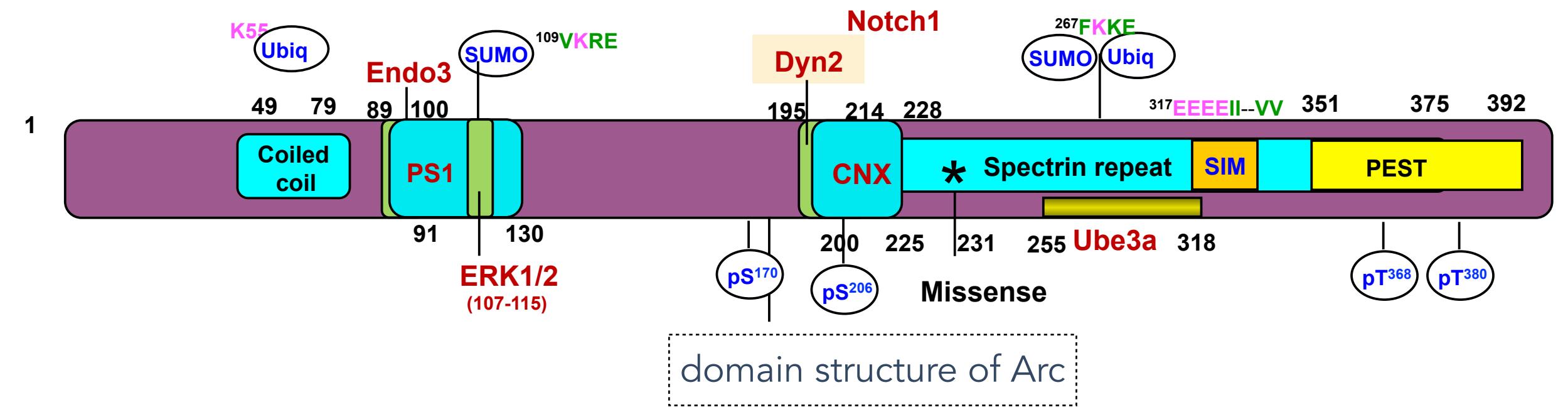
# Predicting mRNA localisation signals



- align 3'UTRs of dendritically localised mRNAs
- predict secondary structure (Vienna, RNAPromo, MC-Fold, LocARNA)
- combine sequence and structural (region) predictions (filter 60% id,  $\geq 35$ nt)
- rank candidates for experimental validation

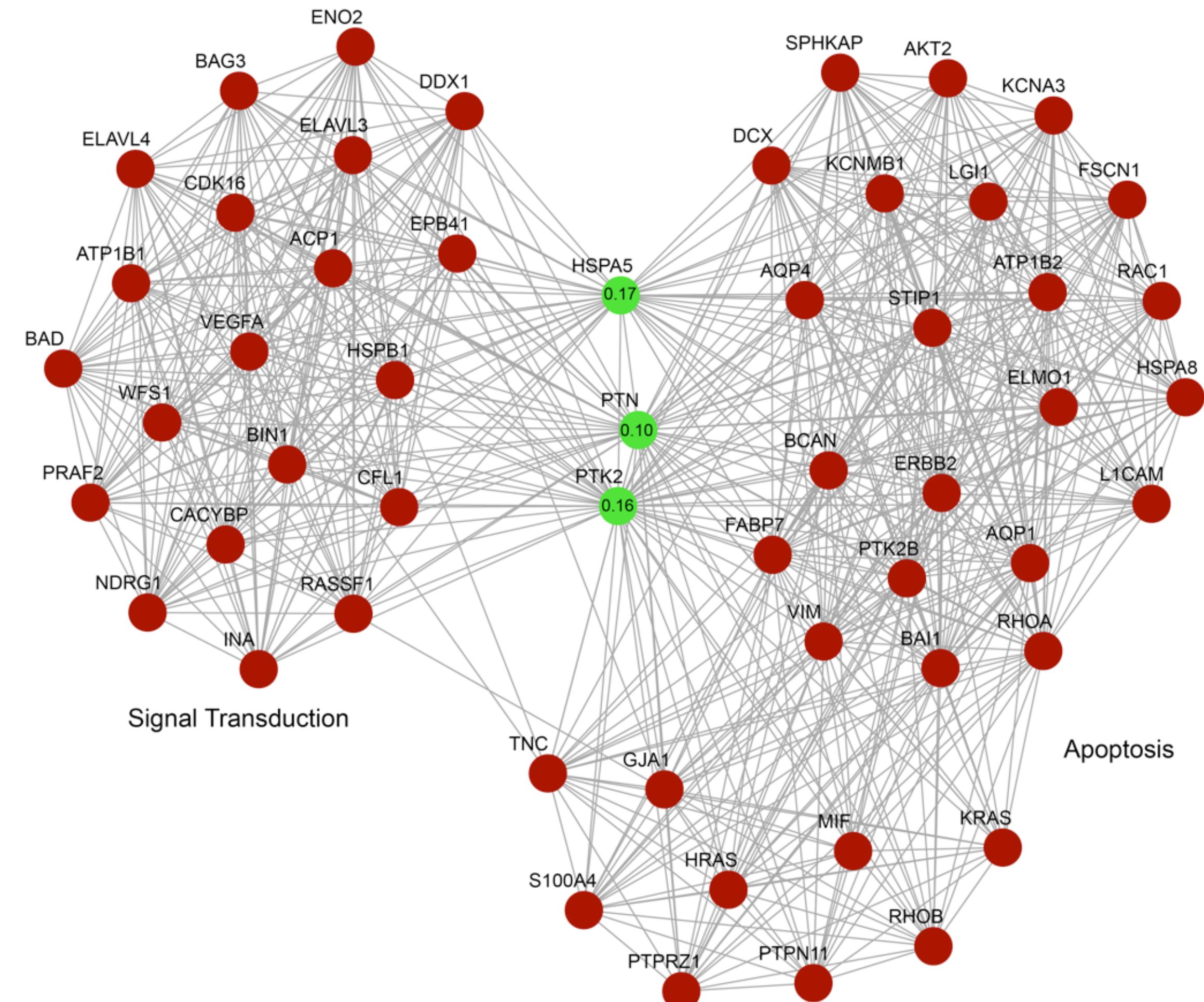
# Evolutionary Sequence Conservation & Neuroregulation

Which sites are most likely to be functional ?



- identified Arc orthologues
- modeled sequence evolution (RaxML)
- mapped selection rates to Arc sites
- candidate sites to be validated in the lab

# Network analysis reveals key disease genes



## Disease Gene Network

- Mine gene disease associations

## OMIM/HDO/GeneRIF

- Calculate pair-wise common disease counts between all gene pairs
- Set genes as nodes, disease counts as edges

## Cluster network by community

- Newman and Girvan 2002

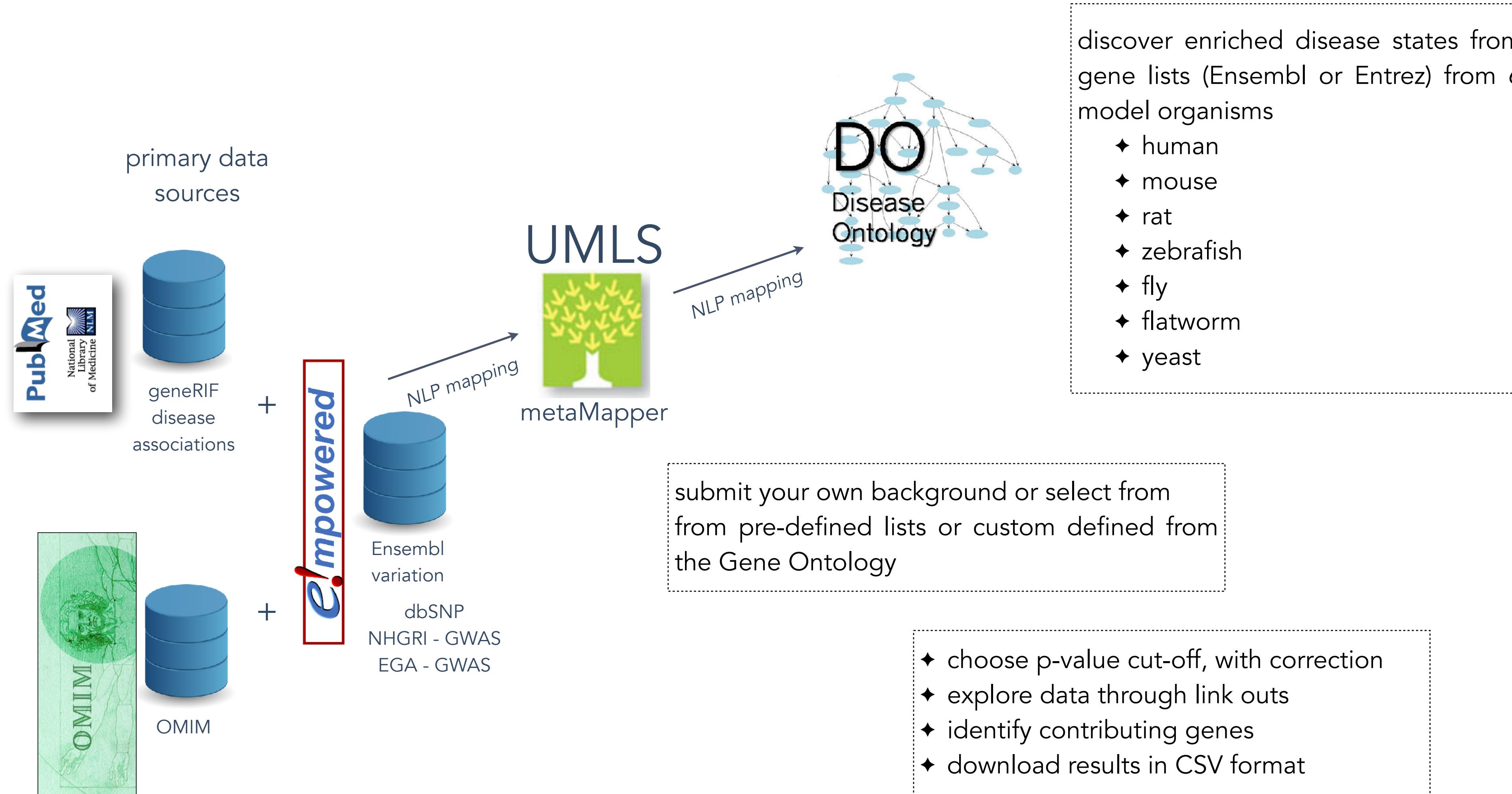
## Functional annotation

- semantic similarity, BINGO/GO, InterPro domain, KEGG, COG.

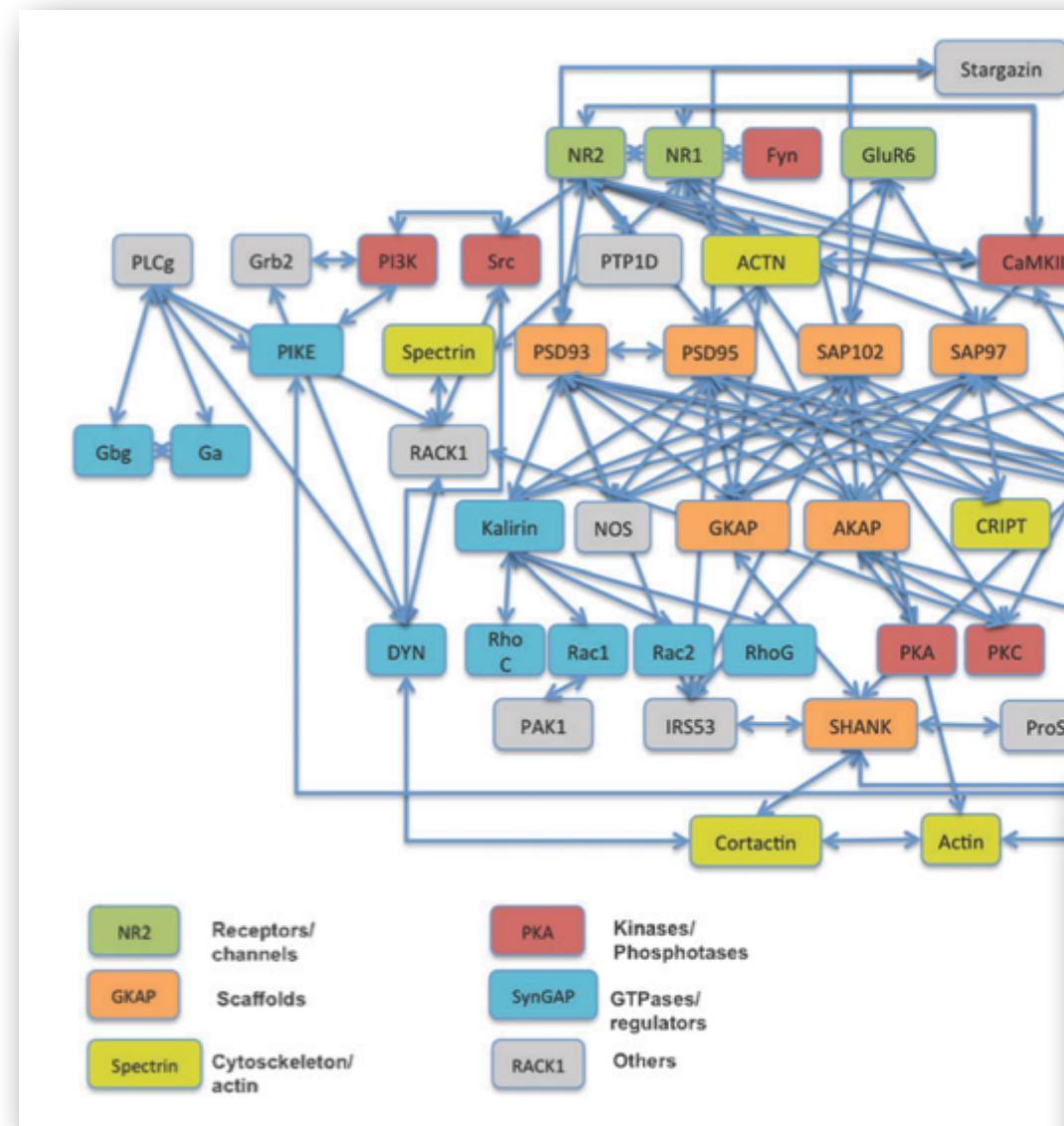
## Network analyses

- betweenness centrality
- modularity
- clustering coefficients
- power-law
- stress testing
- critical genes (targets)

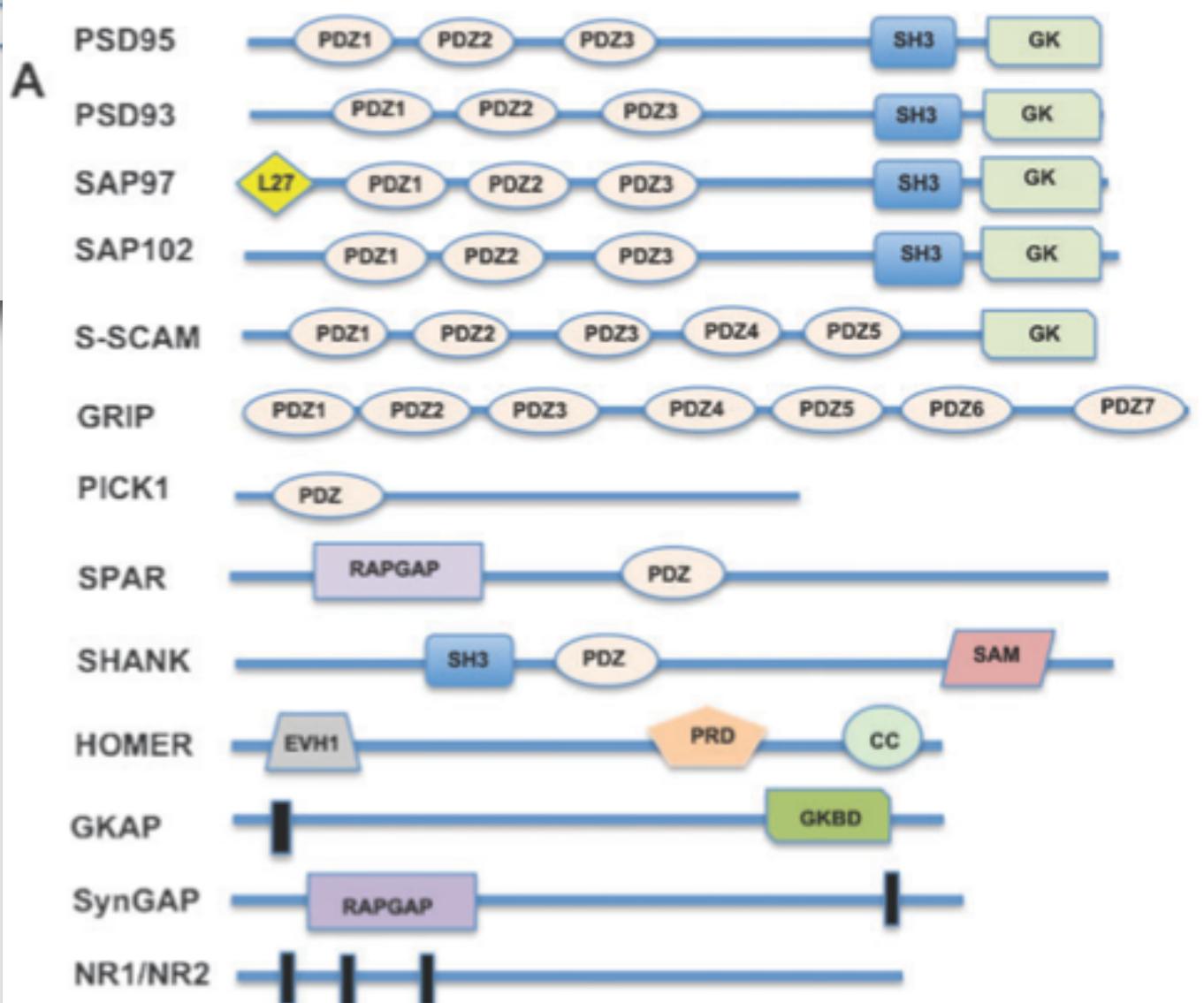
# Computational Pipelines - DisEnT - Disease Enrichment Root



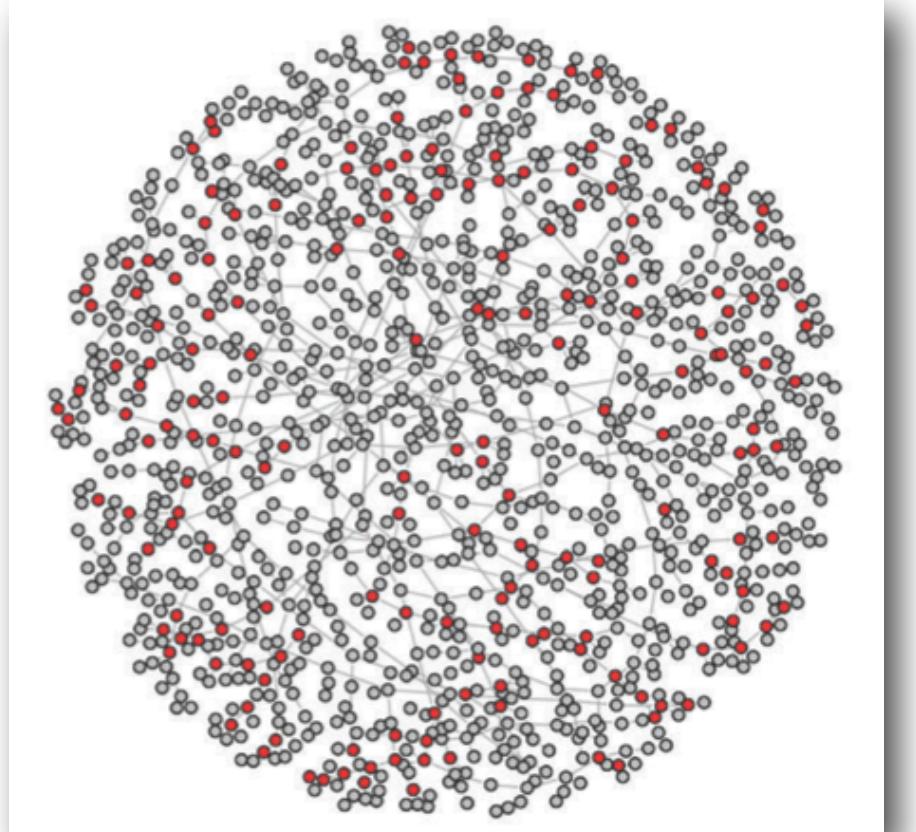
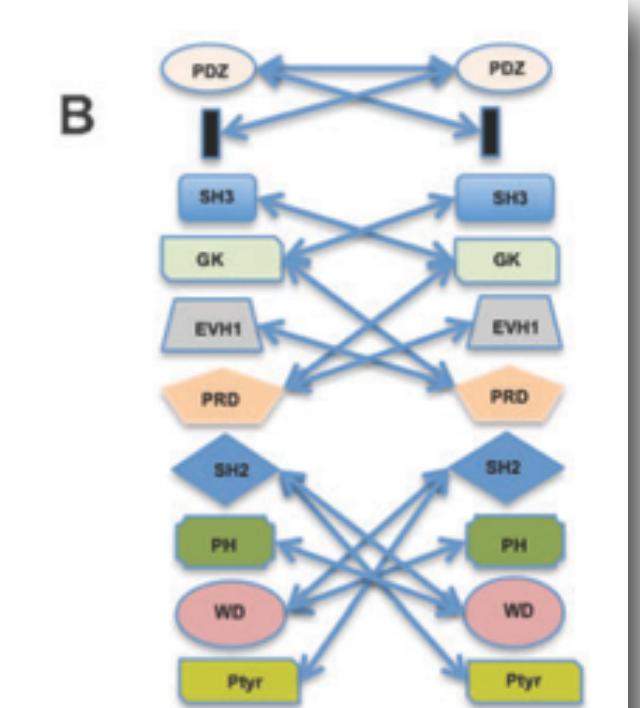
# Dynamical Modelling of Protein Complexes



PSD core-complex



Protein domain binding preferences

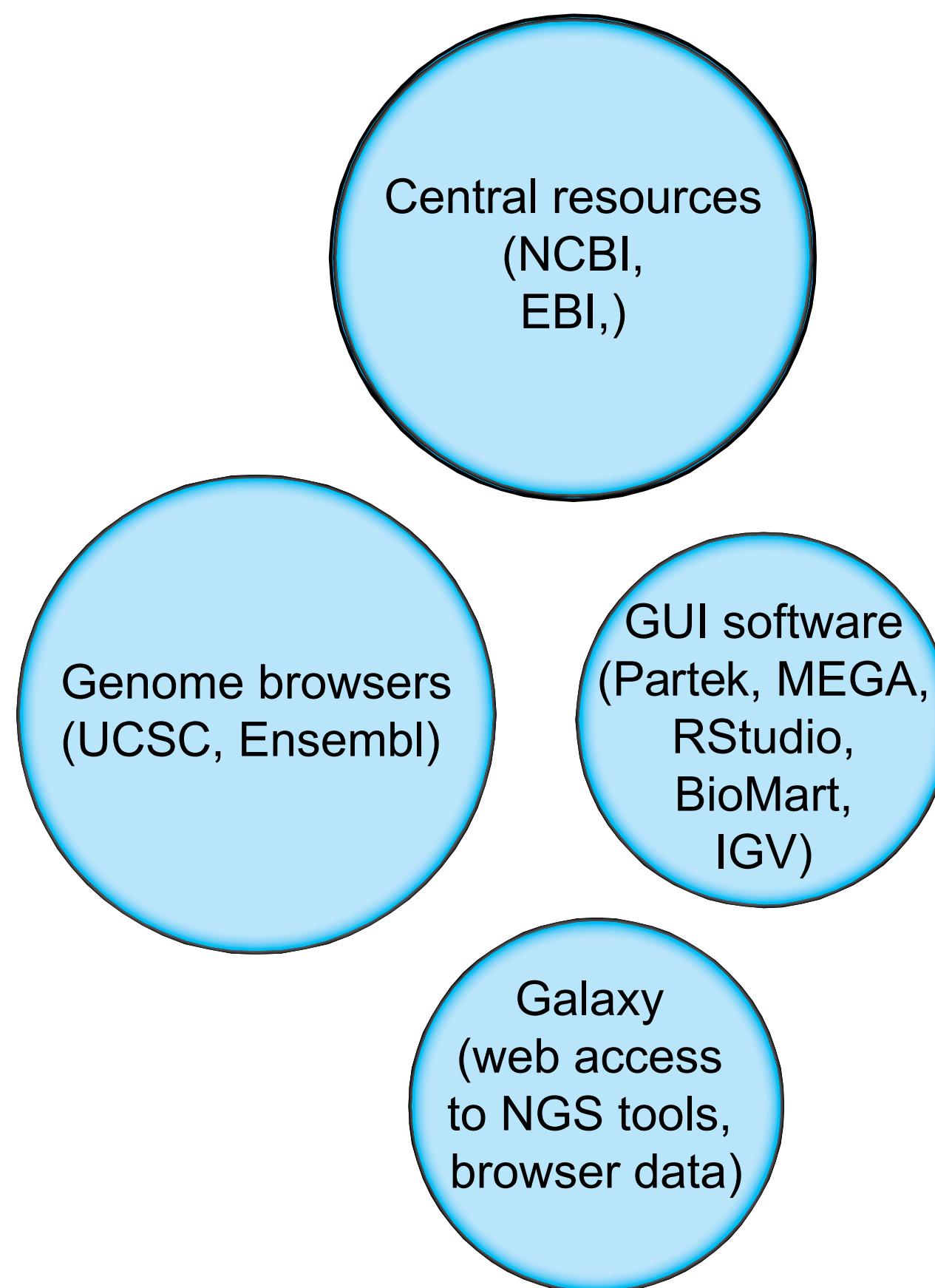


simulated complex

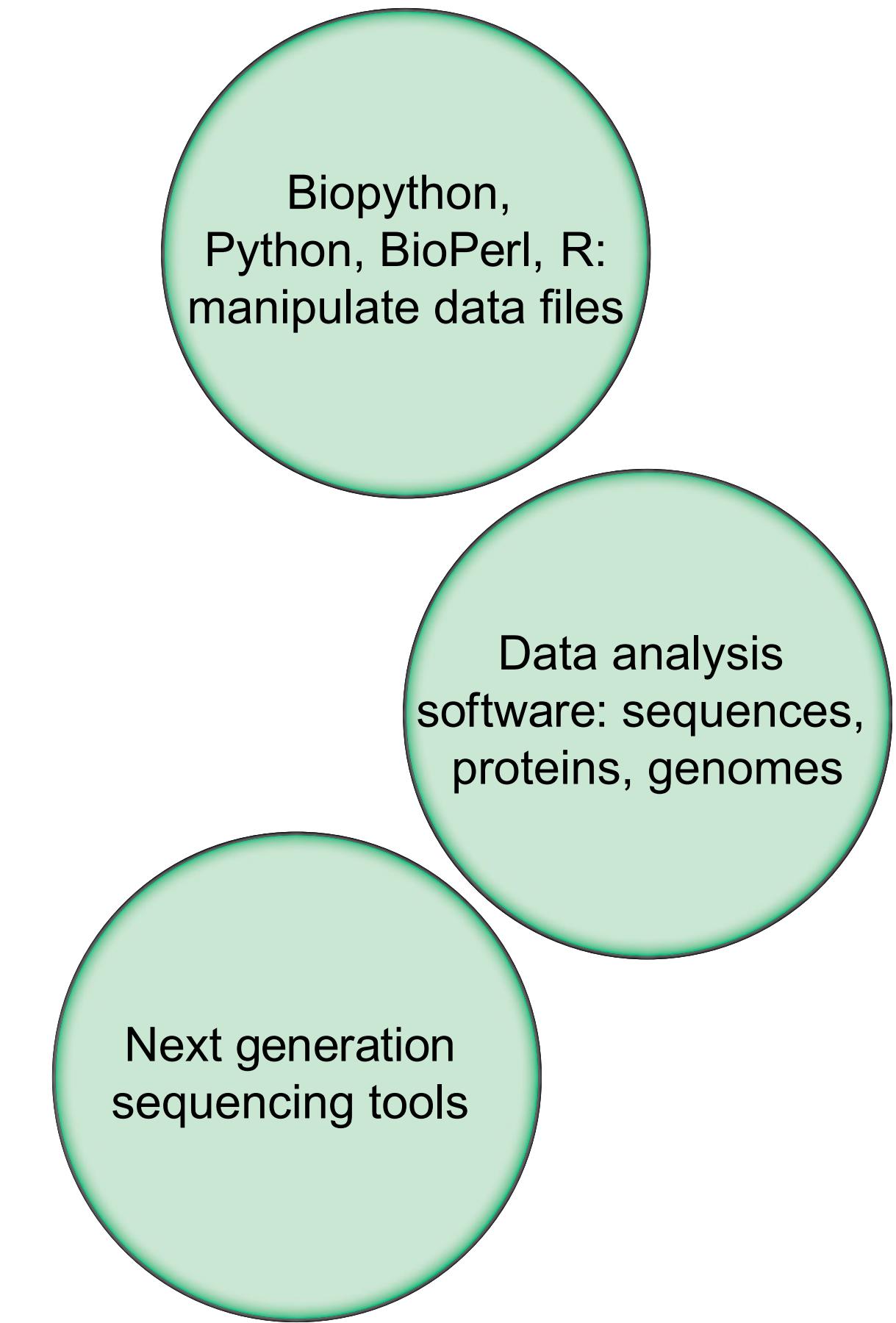
Sorokina et al. 2011

# Web Tools or Programming?

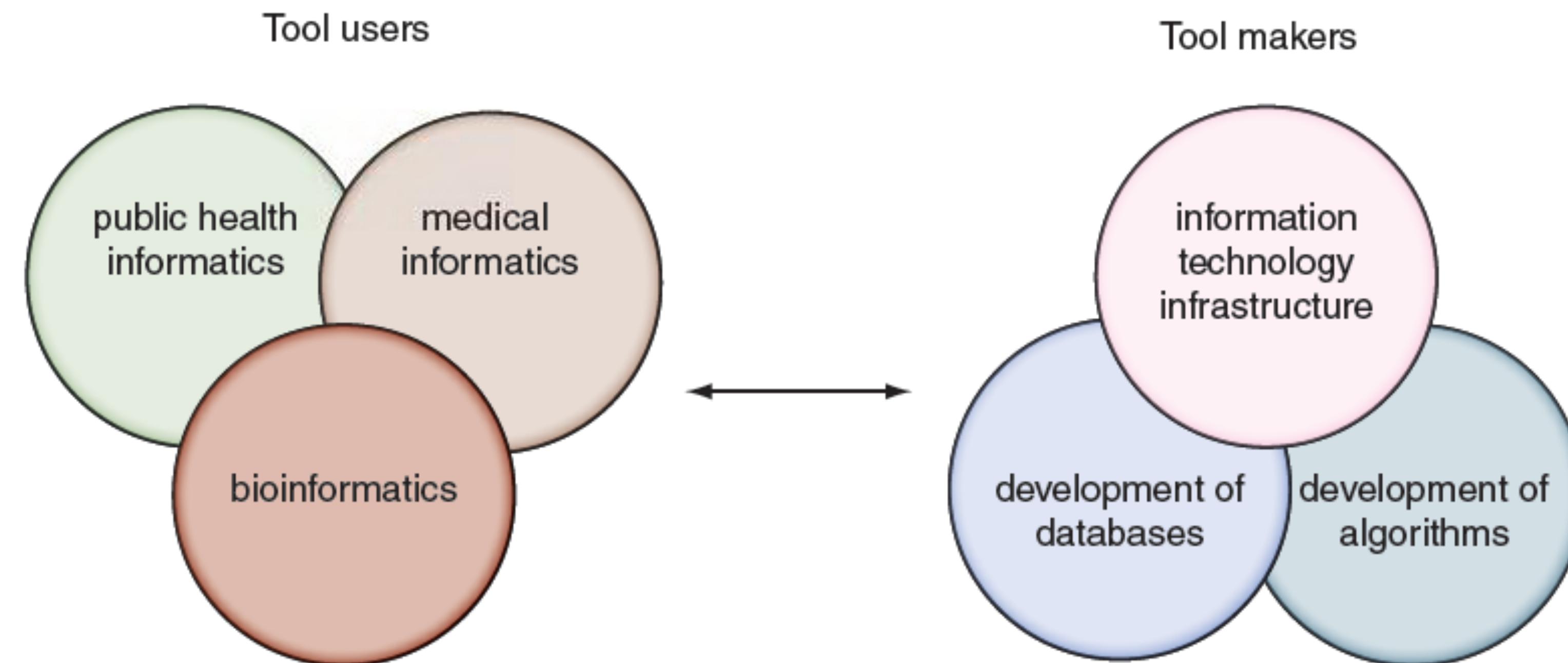
Web-based or  
graphical user interface (GUI)



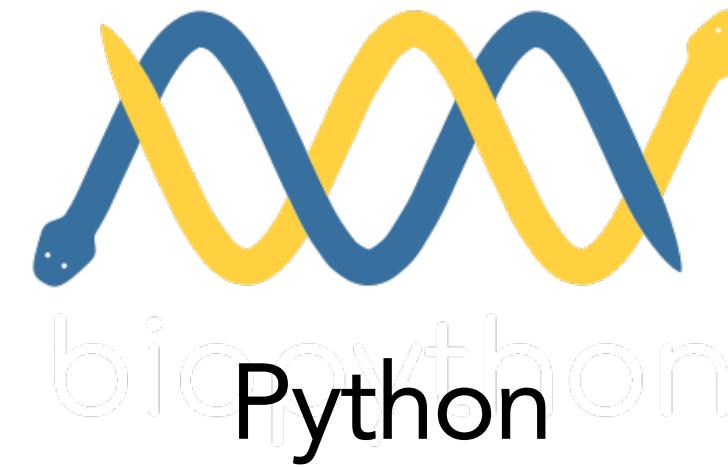
Command line (often Linux)



# Tool Users & Tool Makers



# Coding in Bioinformatics



R/Bioconductor



BioPerl

- ◆ Bioinformatics requires extensive programming skills in several languages
  - Dominated by
    - R (statistical programming) - <https://cran.r-project.org>, <https://www.bioconductor.org>
    - Python - <https://www.python.org>, <https://biopython.org>, Pandas, Numpy, Scipy etc.
  - Increasing use of Machine Learning and “Data Science” methods
    - Supervised/Unsupervised Classification
    - Training of Statistical Models
    - Modelling of Data in Graph (Network) Structures
    - Bayesian Inference
    - Computational Statistics
    - Simulation
  - Ability to integrate multiple packages/tools/data sources is key

# Bioinformatics - crossing disciplines

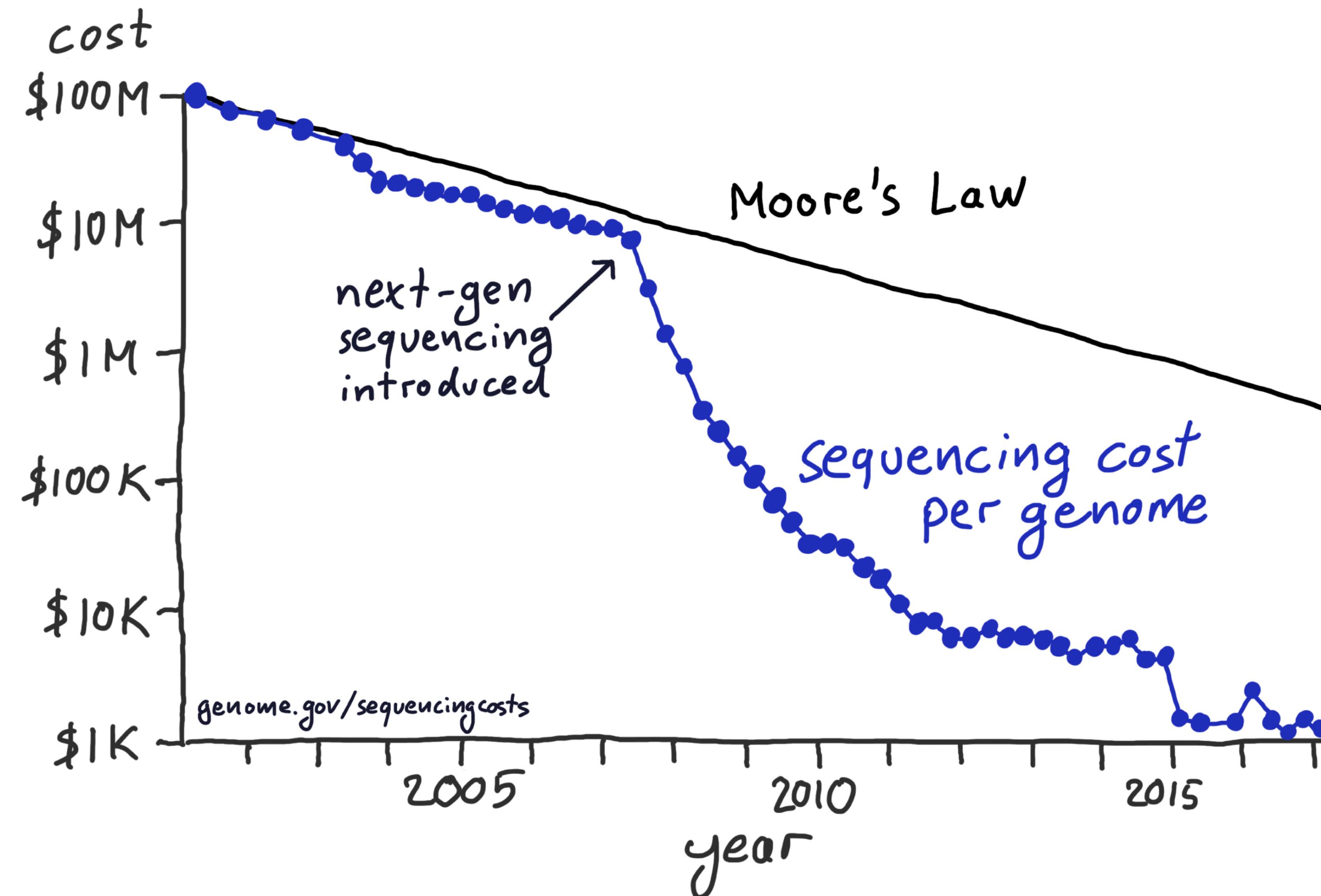
## Computer Science

- Biology provides CS with new challenges with clear medical significance.
- Complex and large datasets, sometimes very noisy with hidden structures.
- Can biological solutions be used to inspire new computational tools and methods?

## Biology

- Bioinformatics enables analysis of increasingly large and complex data sets generated in Biological experiments. (100k Genomes, UKB-Imaging, Microbiomes, Single-cell sequencing)
- Bioinformatics facilitates the structuring of often unstructured and heterogeneous data to allow data integration for insight
- Bridges CS to Biology especially when realised in Web services & Tools

# Next-generation Sequencing



# Genomic (Open) Data

**Genome Information by organism**

Search by organism | Clear

2015 Overview [13572] Eukaryotes [2398] Prokaryotes [48665] Viruses [4905] Plasmids [6121] Organelles [7134]

2016 Overview [17110] Eukaryotes [3513] Prokaryotes [74033] Viruses [5699] Plasmids [7759] Organelles [8566]

2017 Overview [25795] Eukaryotes [4739] Prokaryotes [107512] Viruses [7475] Plasmids [10020] Organelles [10722]

Genome > Genome Information by Organism

Organism name (common or scientific) or Accession (Assembly, BioProject or replicon) ... Q Search

Overview (47211); Eukaryotes (9210); Prokaryotes (210400); Viruses (32535); Plasmids (18282); Organelles (14660)

Genome > Genome Information by Organism

Organism name (common or scientific) or Accession (Assembly, BioProject)

Overview (72067); Eukaryotes (24955); Prokaryotes (437192); Viruses (51007)

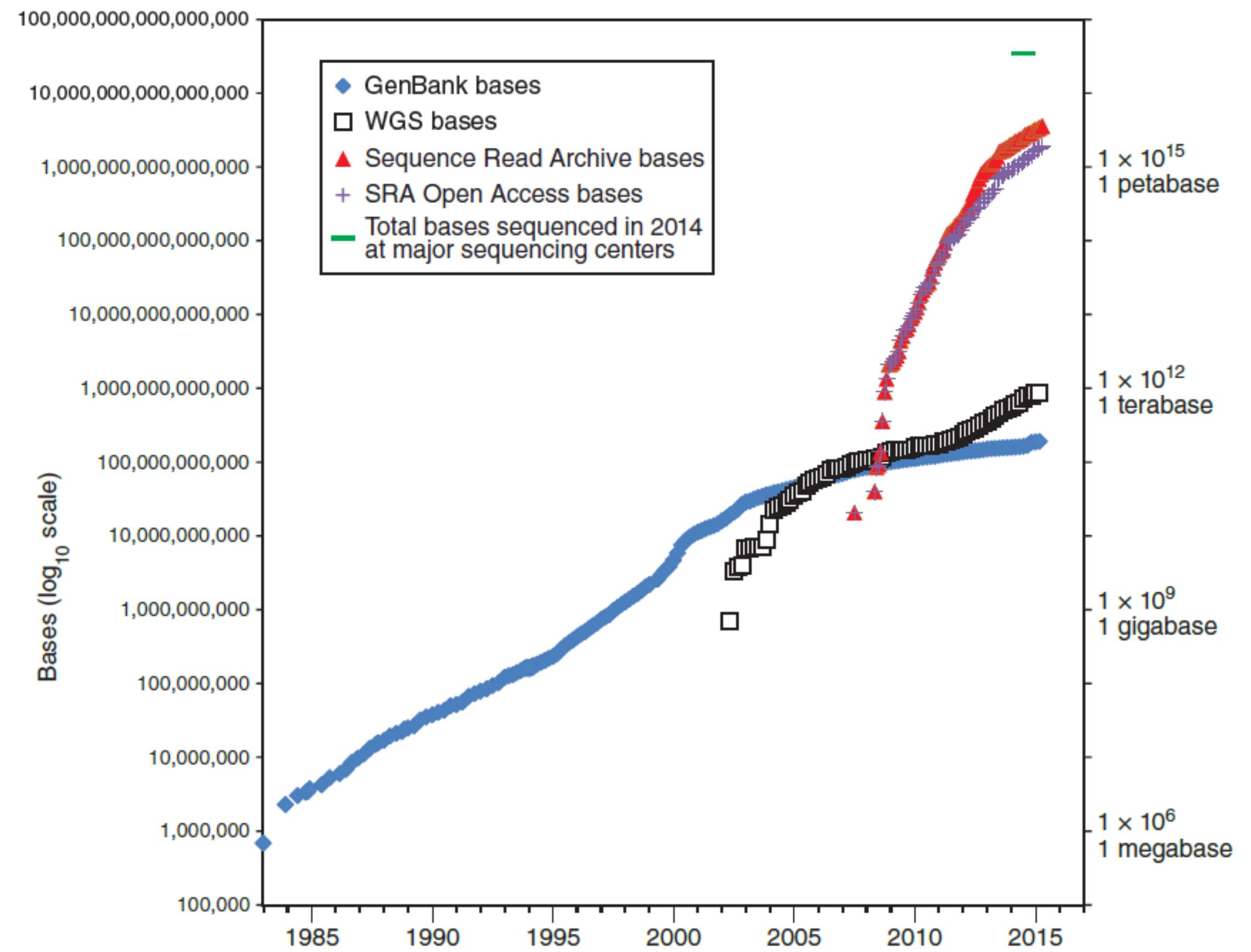
Choose Columns

#	Organism Name	Organism ID
1	'Brassica napus' phytoplasma	Bacteria;Terrabacteria group;T
2	'Candidatus Kapabacteria' thiocyanatum	Bacteria;FCB group;Bacteroid

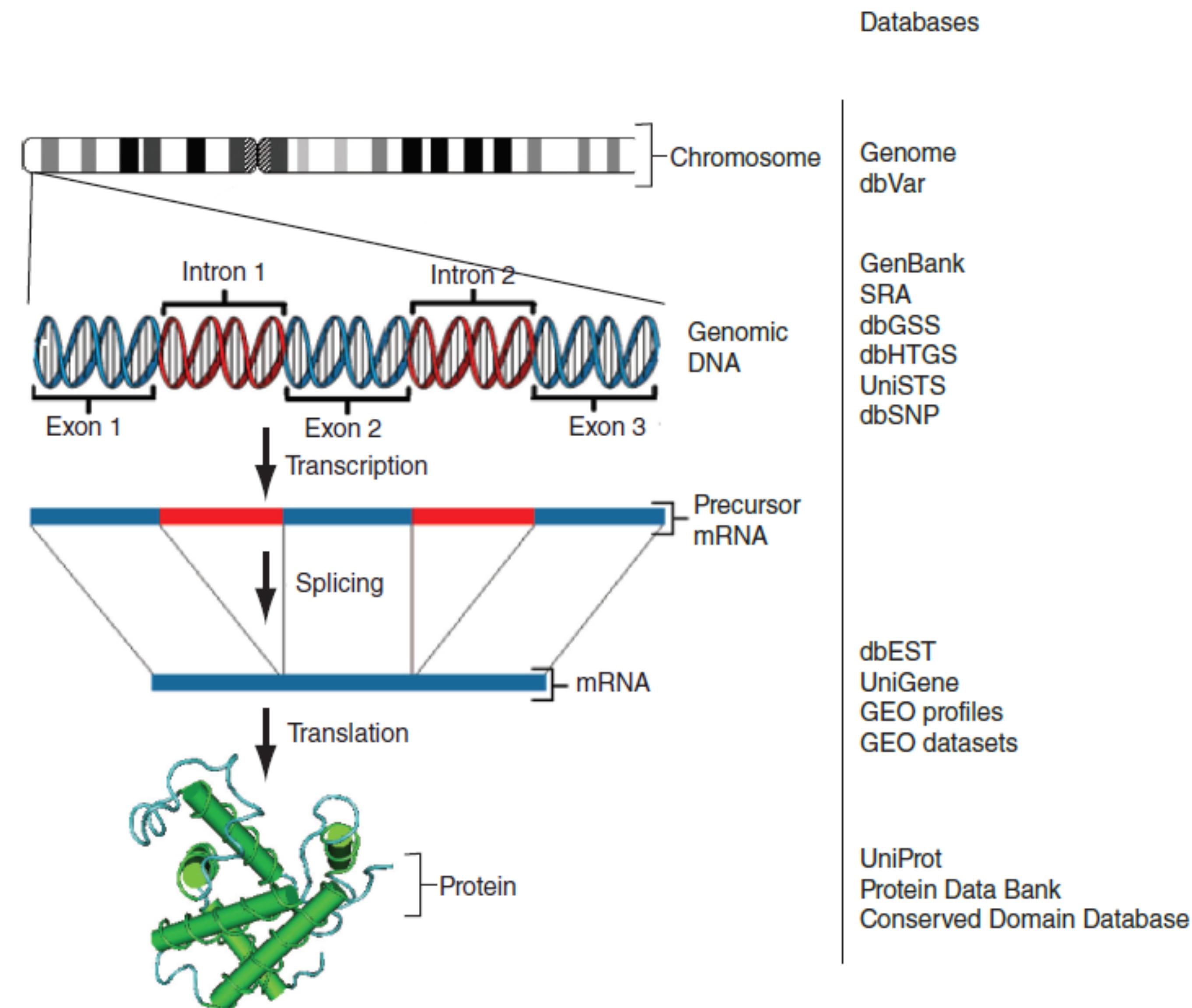
2022

Year	#Genomes
2015	13572
2016	17110
2017	25795
2019	47211
2022	72067

## Growth of Sequence Databases

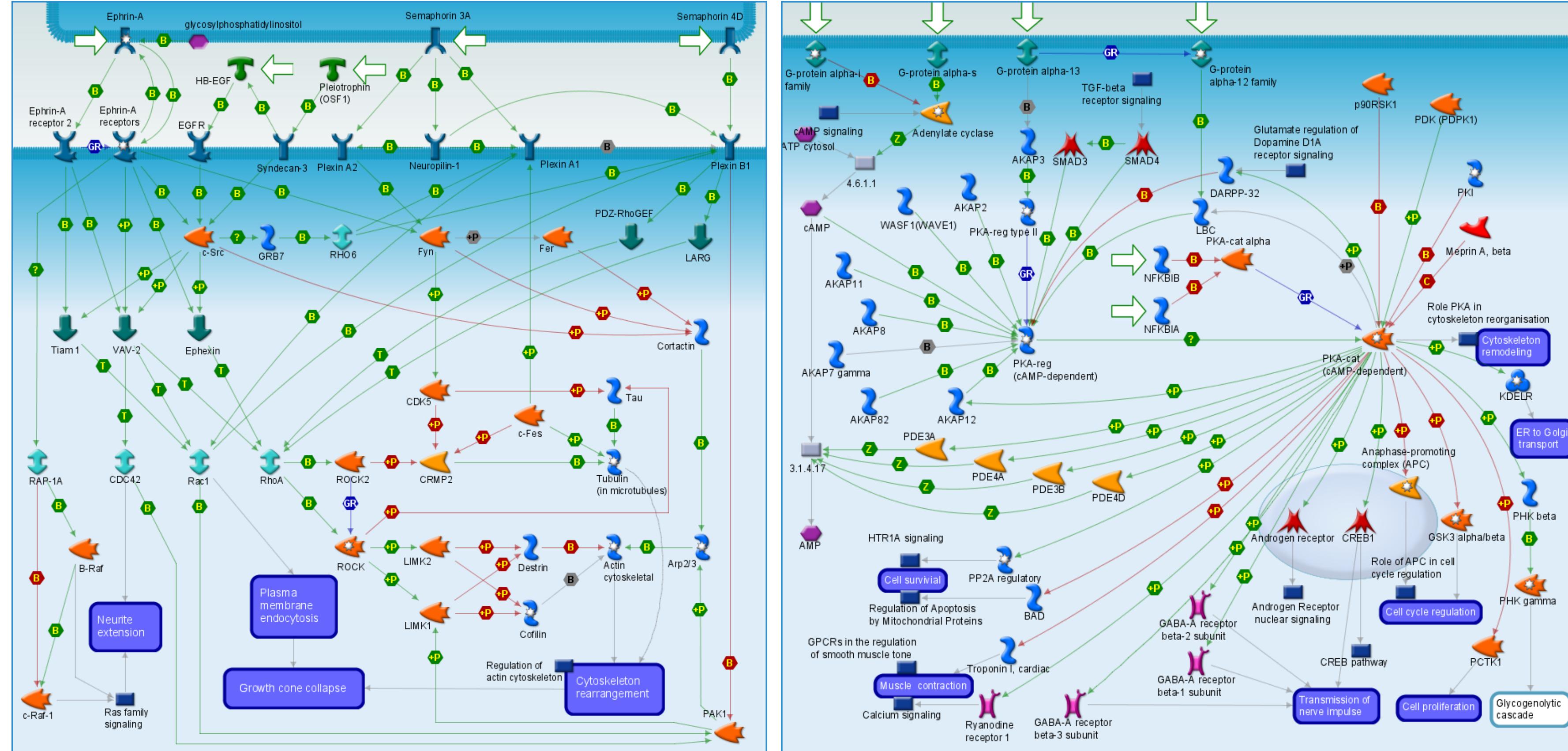


# Data Types & Databases



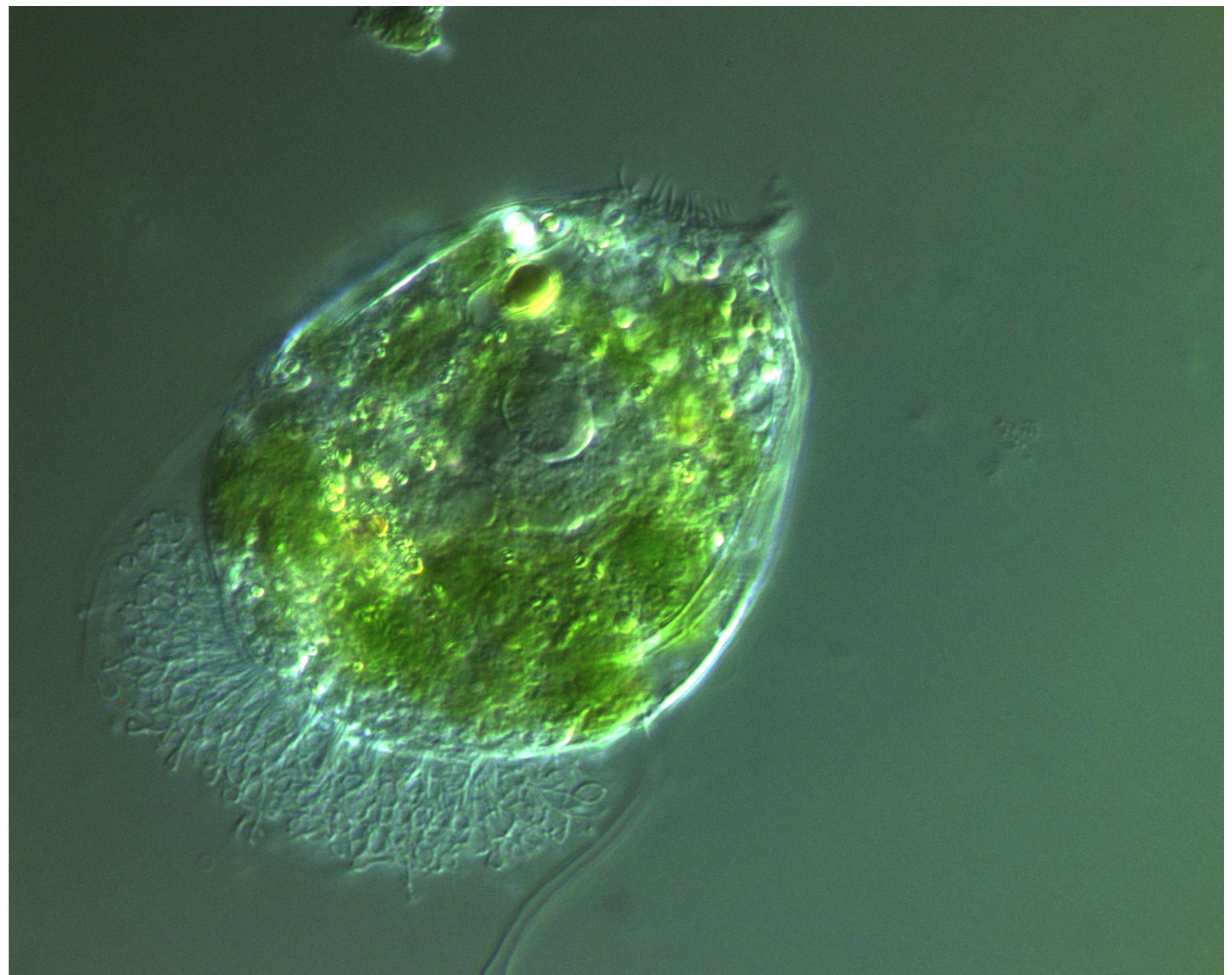
B&FG3e

# How much biology should the bioinformatician know?

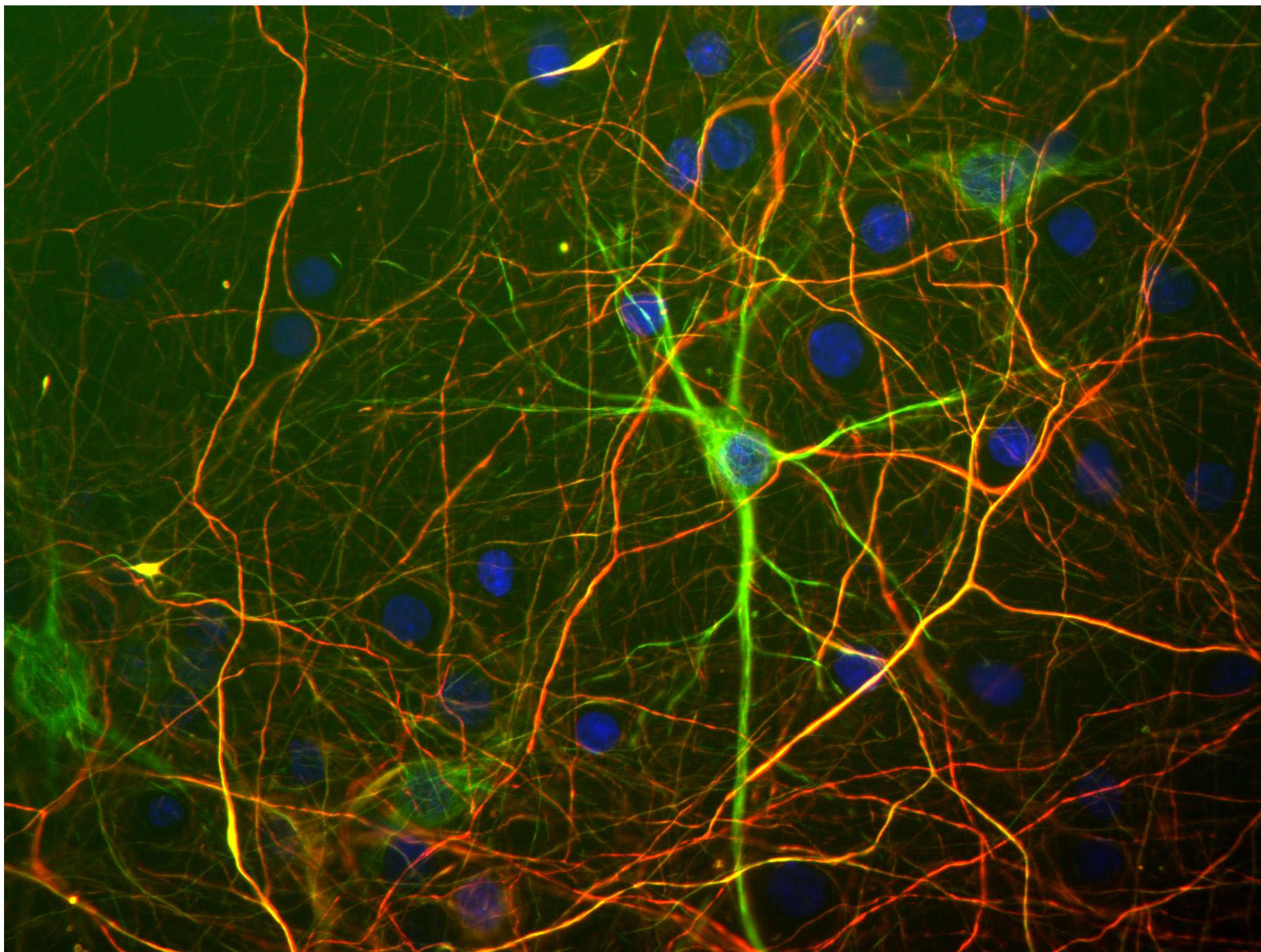


Bioinformaticians and Computational Biologists **develop** and **test** new biological hypotheses, not just analyse data for others.

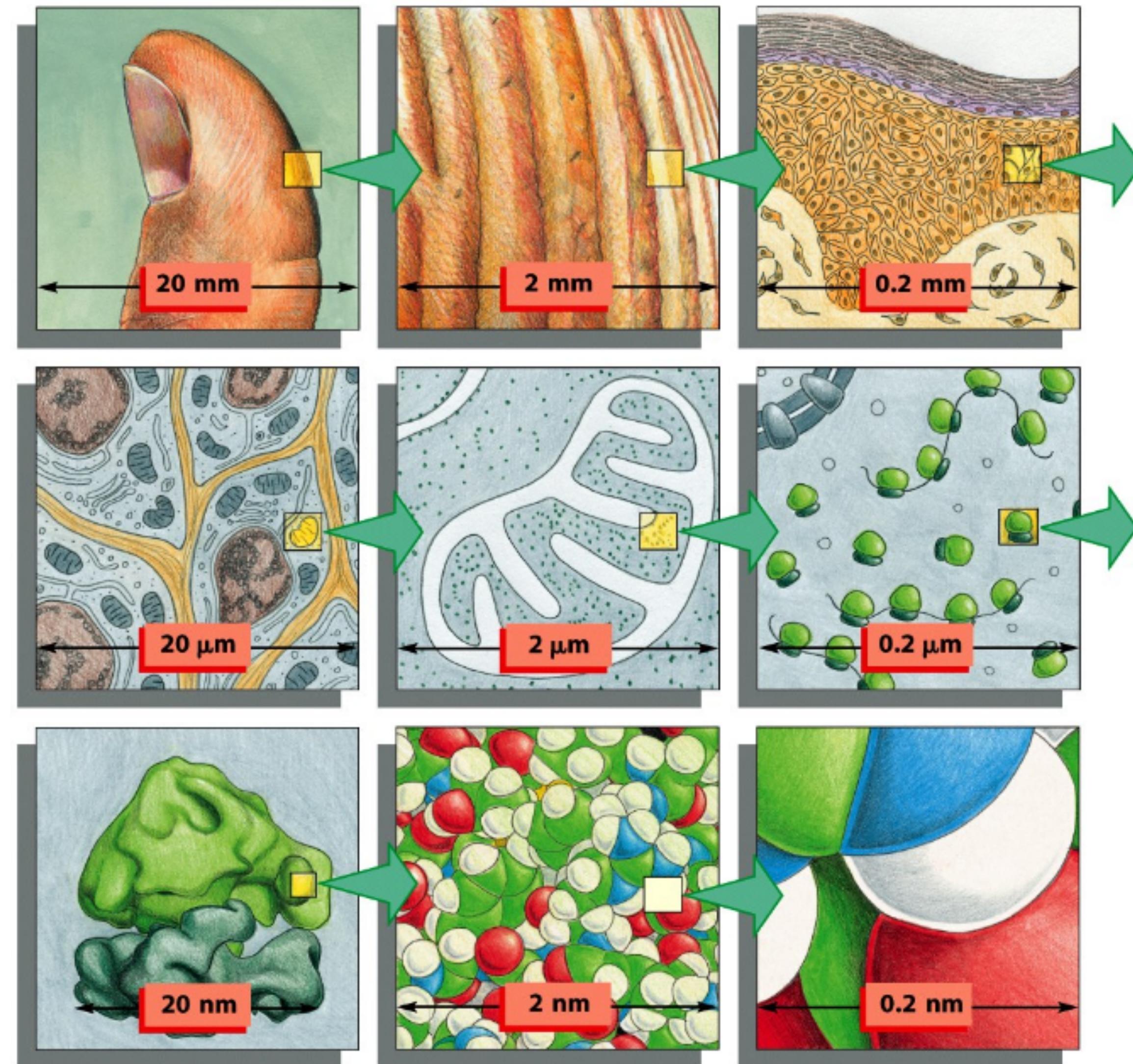
# A Unicellular Organism - *Strombidium stylifer*



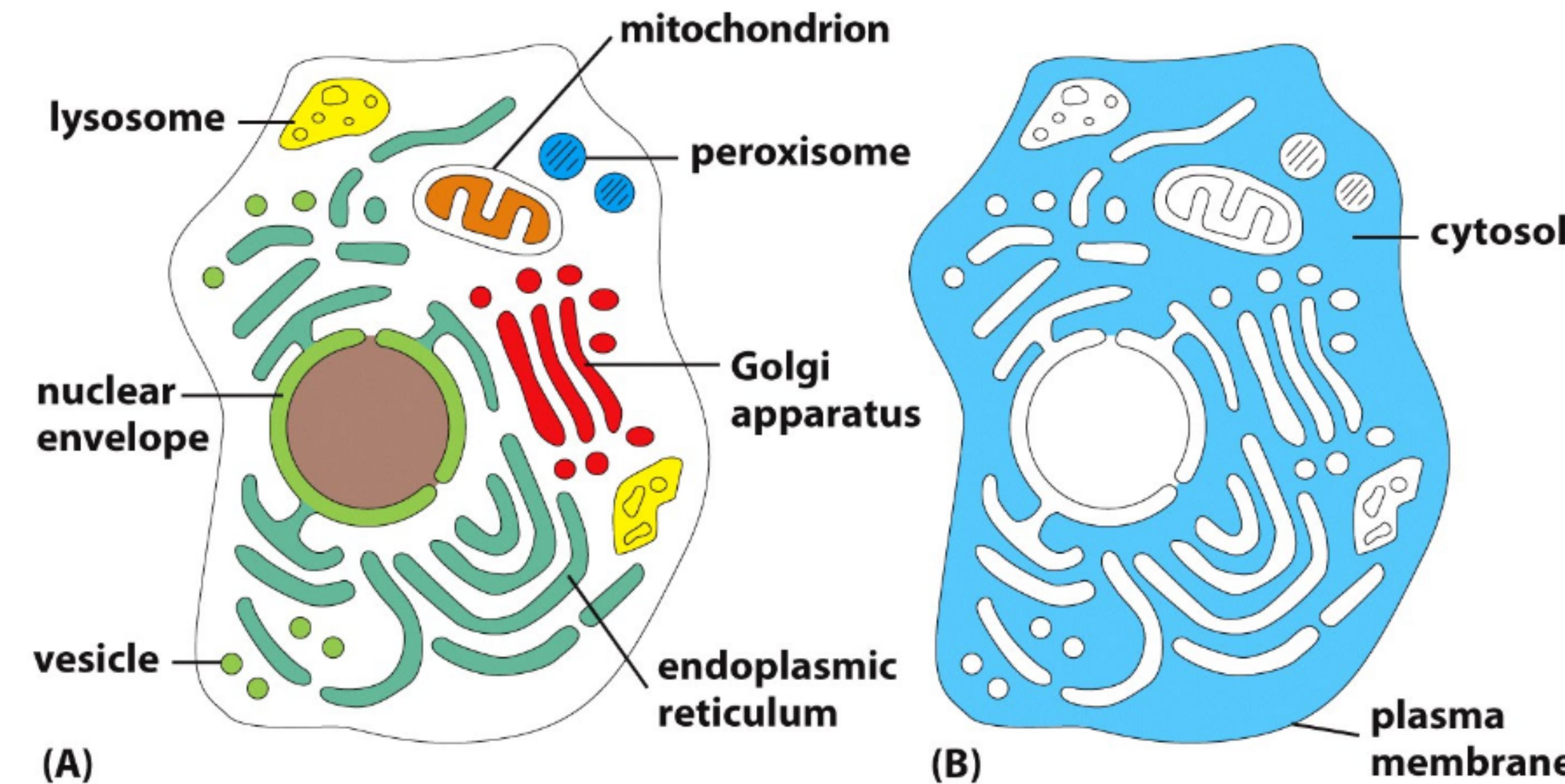
# Cultured Neurons



# Studying Across Scales

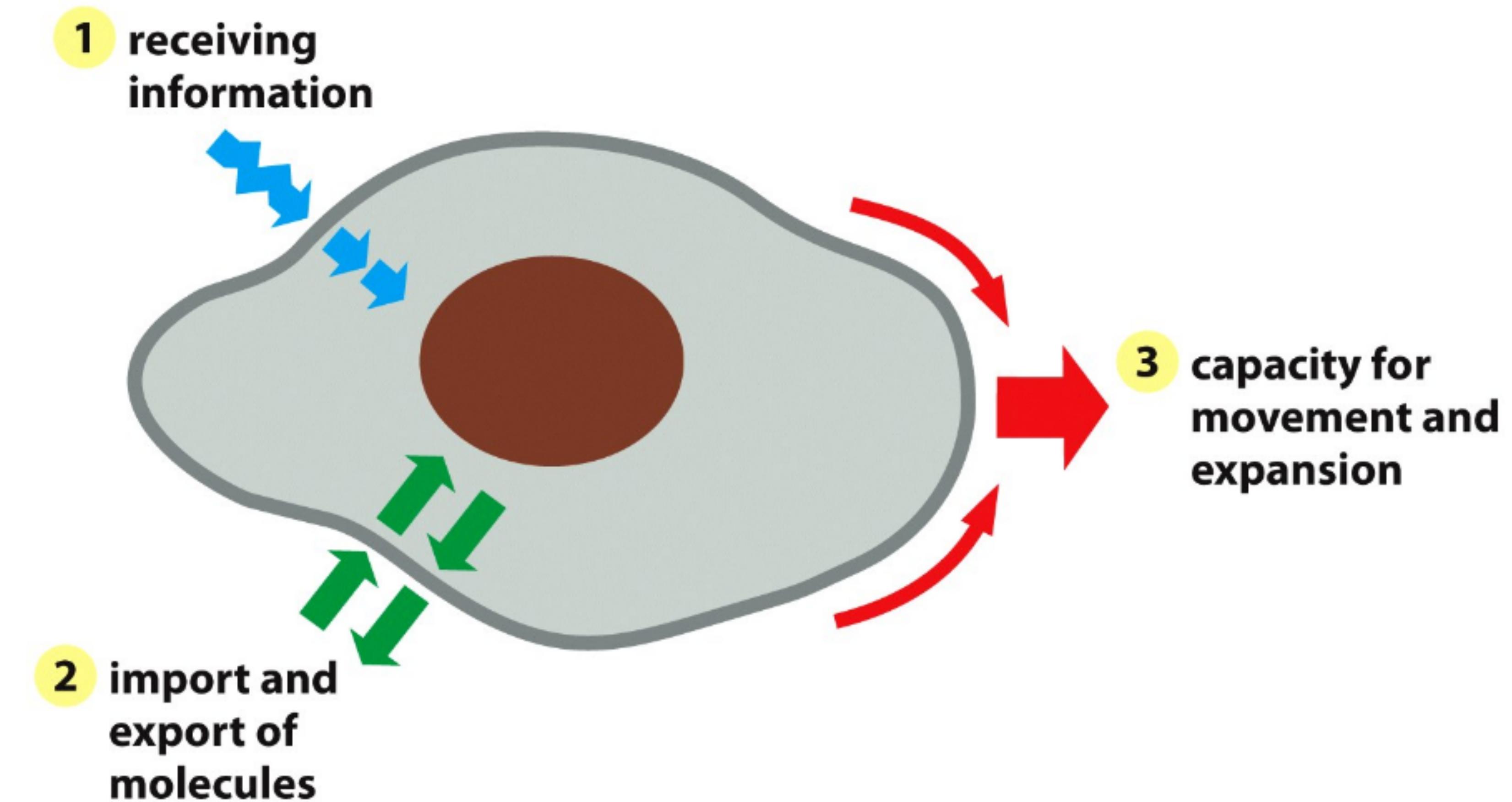


# Cell Structure



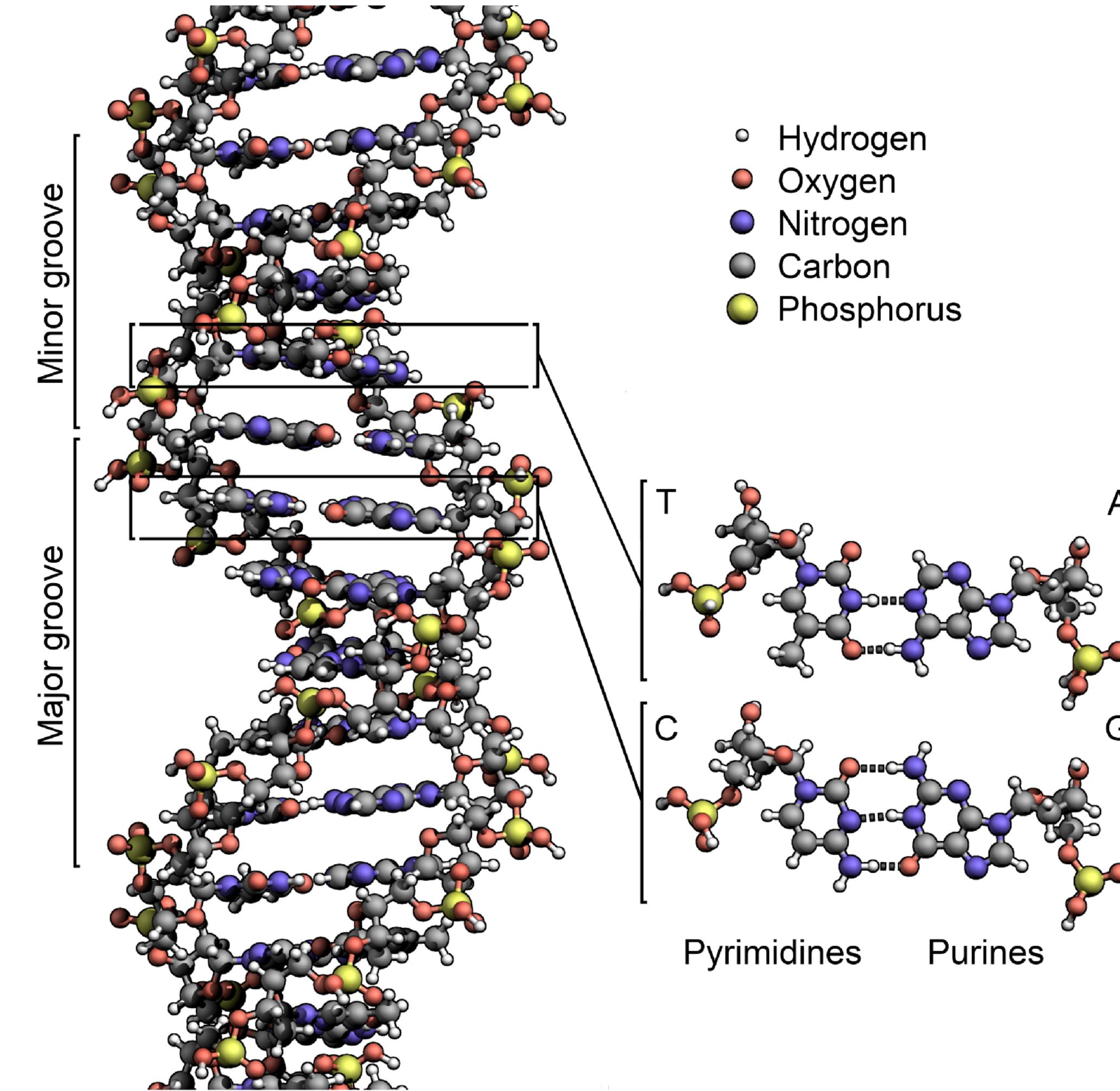
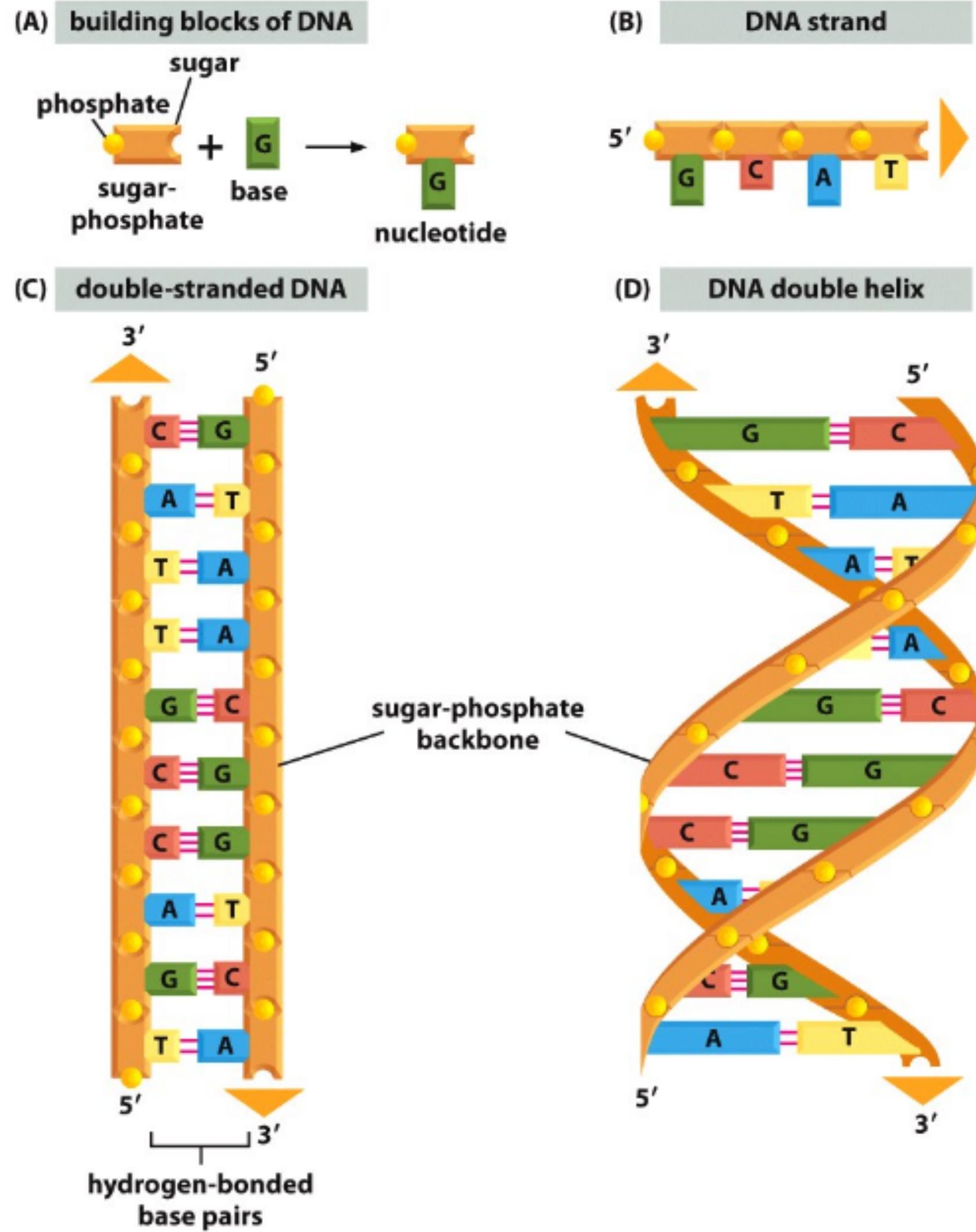
Alberts, B. (2014). *Essential cell biology*. New York: Garland Science.

# Cells Respond to their Environment

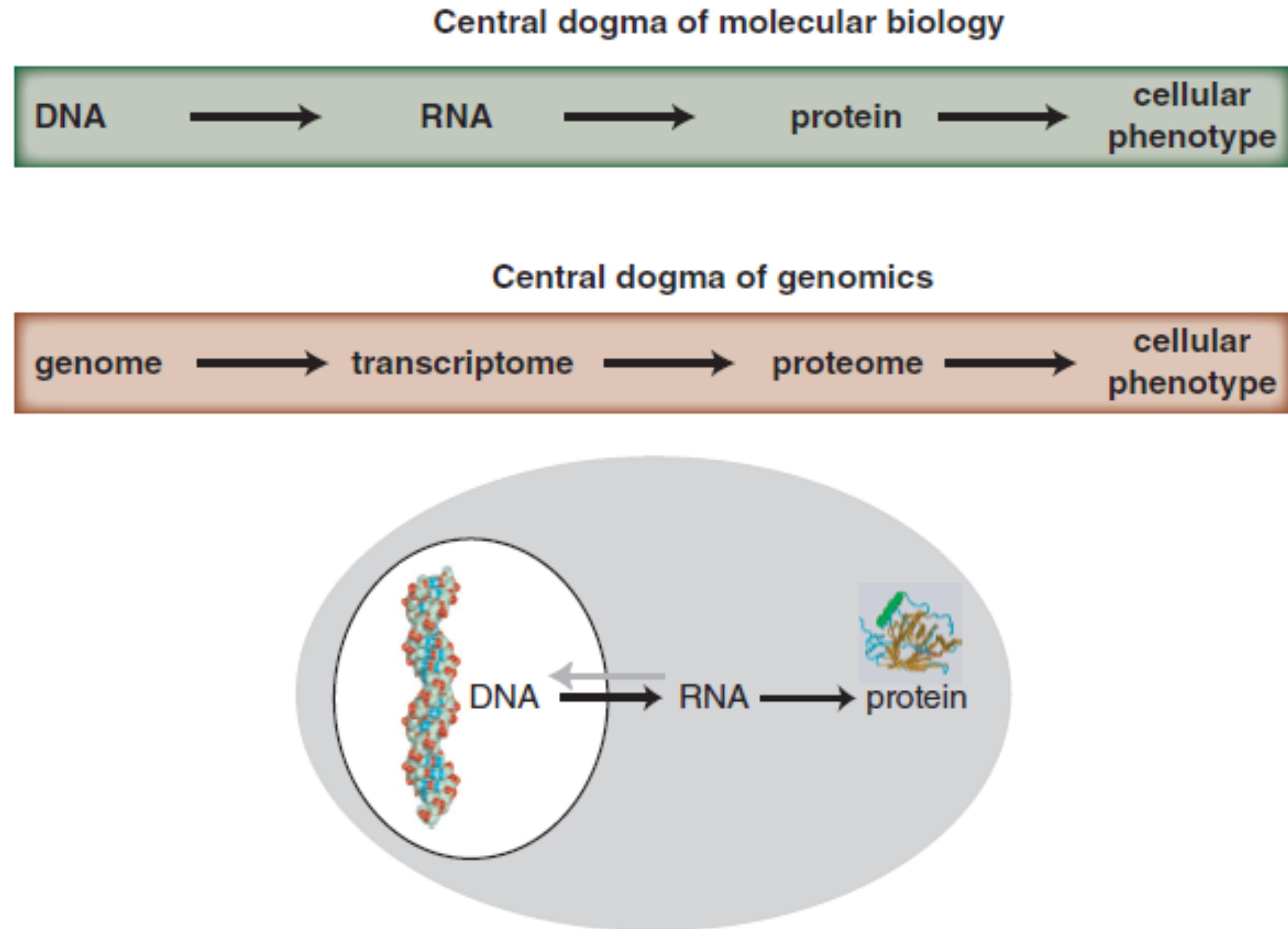


Alberts, B. (2014). *Essential cell biology*. New York: Garland Science.

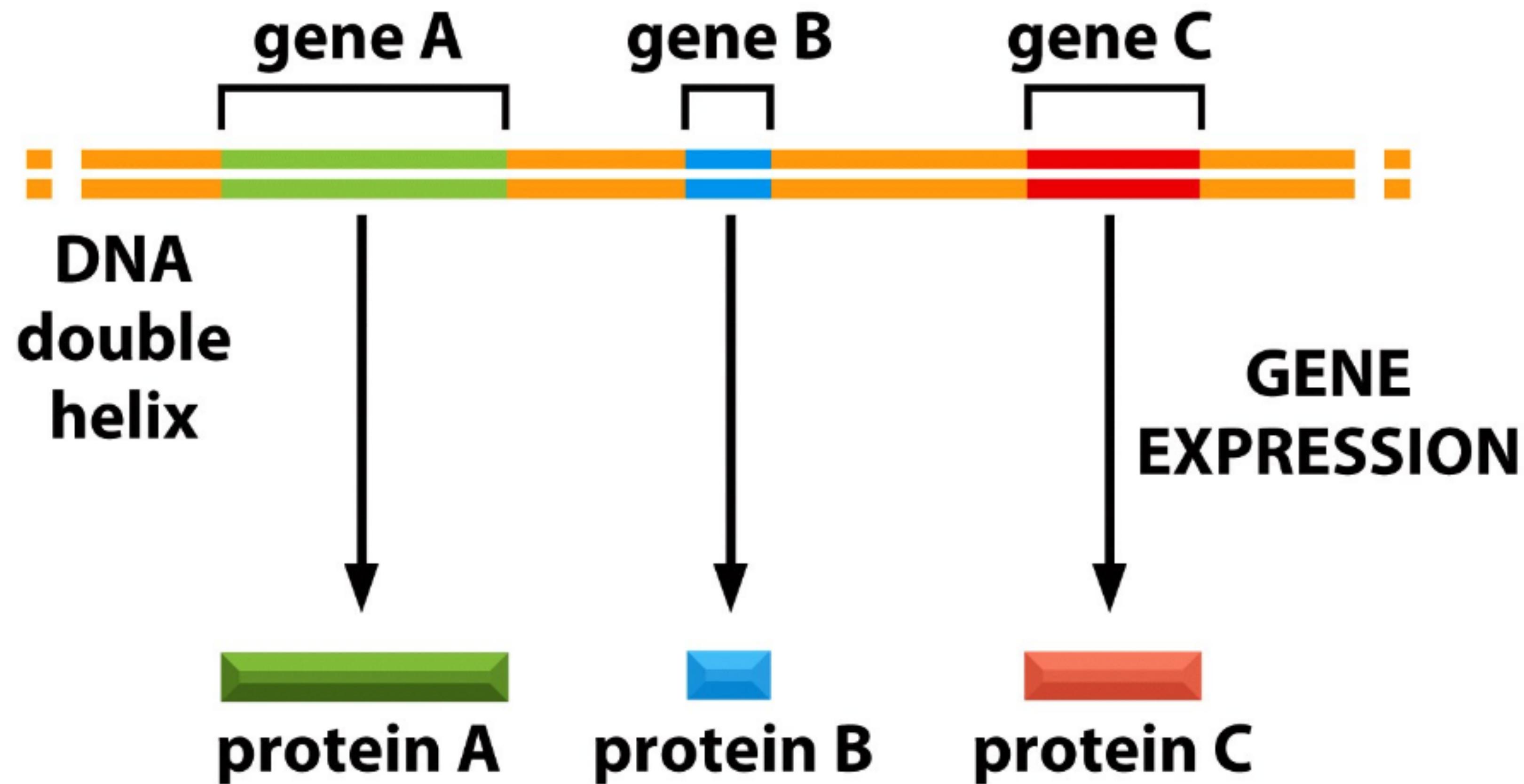
# Deoxy-riboNucleicAcid (DNA)



# The “Central Dogma”



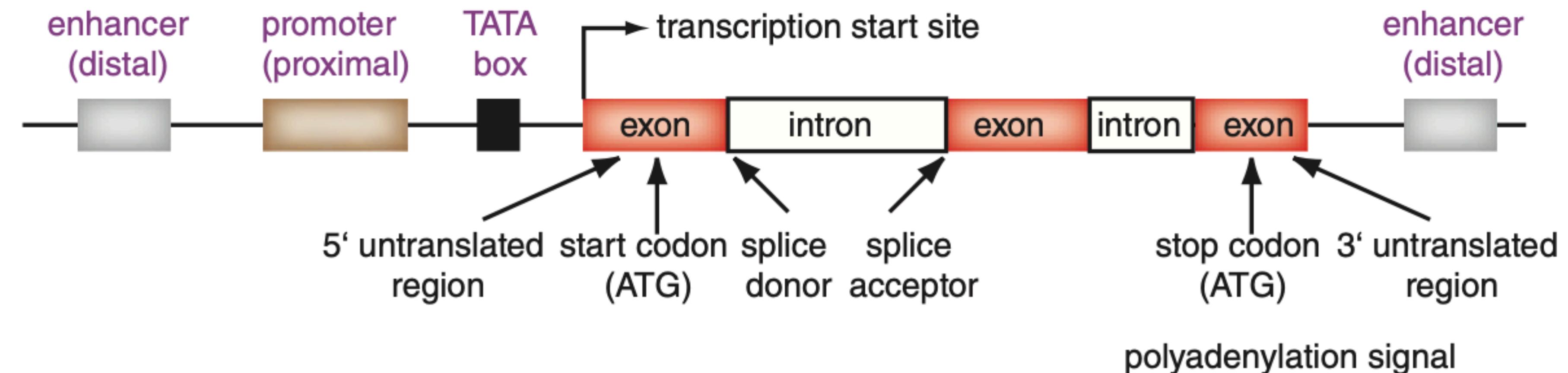
# Protein Coding Genes



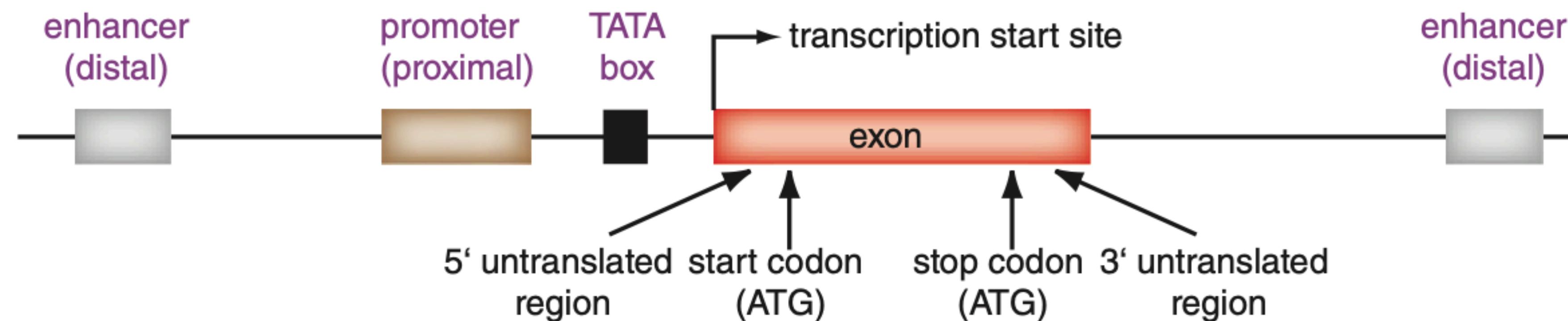
Alberts, B. (2014). *Essential cell biology*. New York: Garland Science.

# Gene Structures (Eukaryotic)

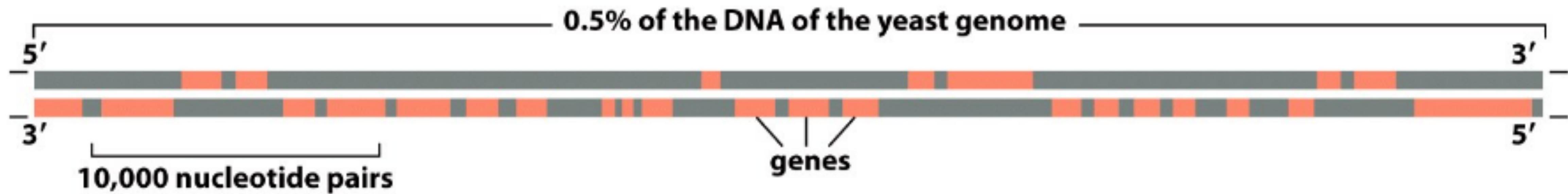
(a) Gene with multiple exons



(b) Single exon gene

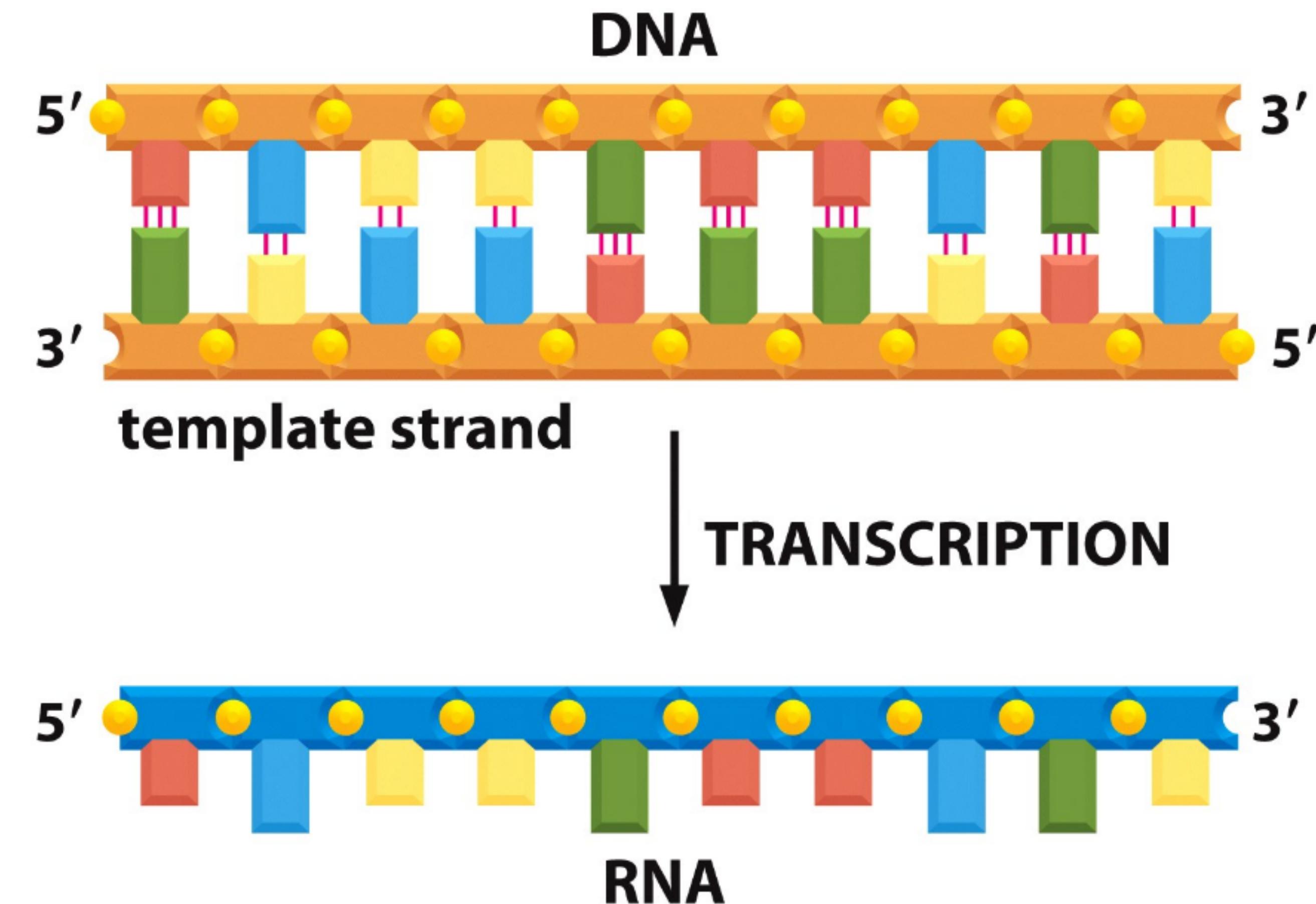


# DNA Strandedness



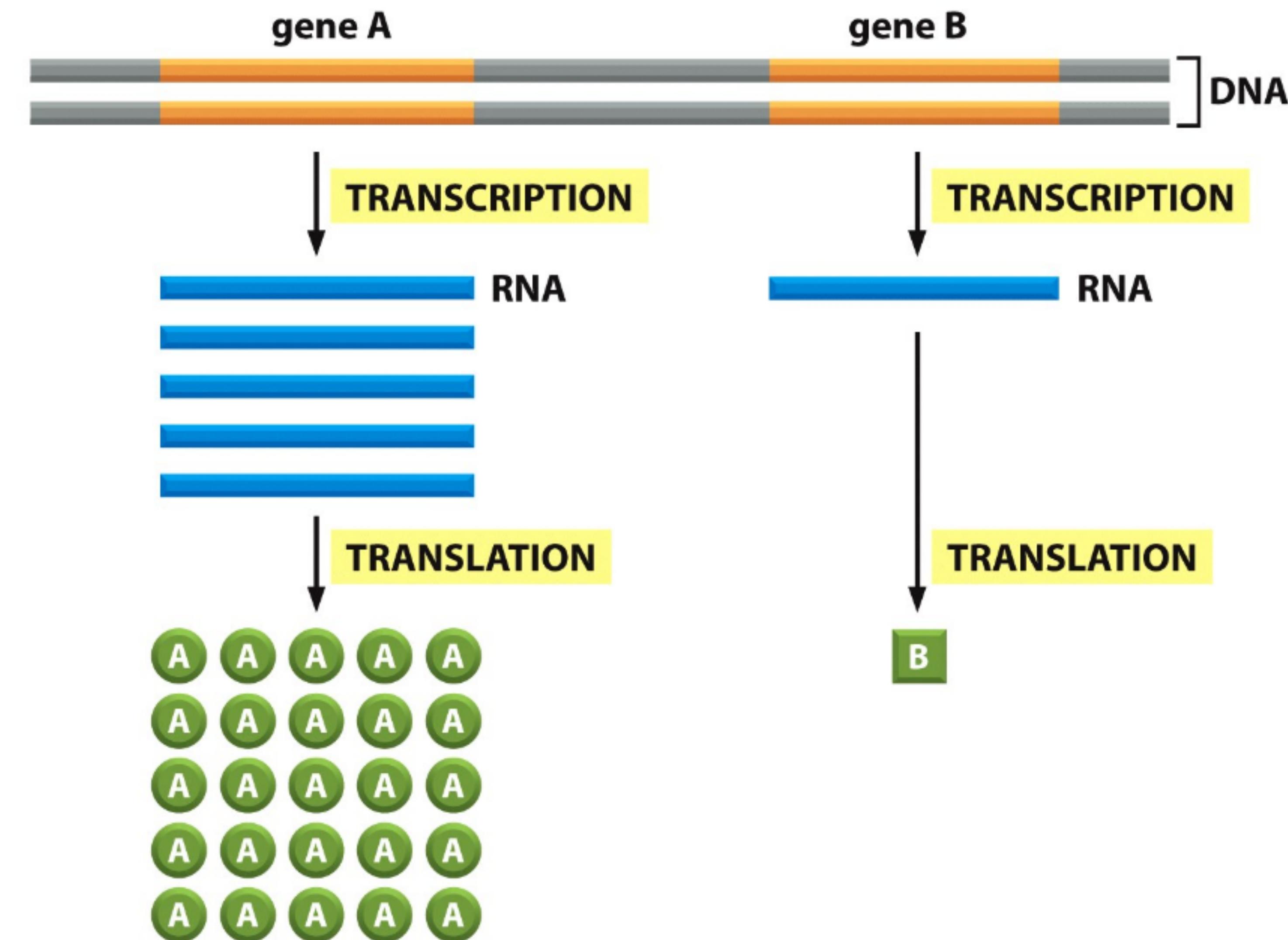
Alberts, B. (2014). *Essential cell biology*. New York: Garland Science.

# Transcription



Alberts, B. (2014). *Essential cell biology*. New York: Garland Science.

# Translation



Alberts, B. (2014). *Essential cell biology*. New York: Garland Science.

# Noteable & Jupyter Notebooks

The University of Edinburgh

Please select a personal notebook server

Standard Notebook (Python 3) ▾

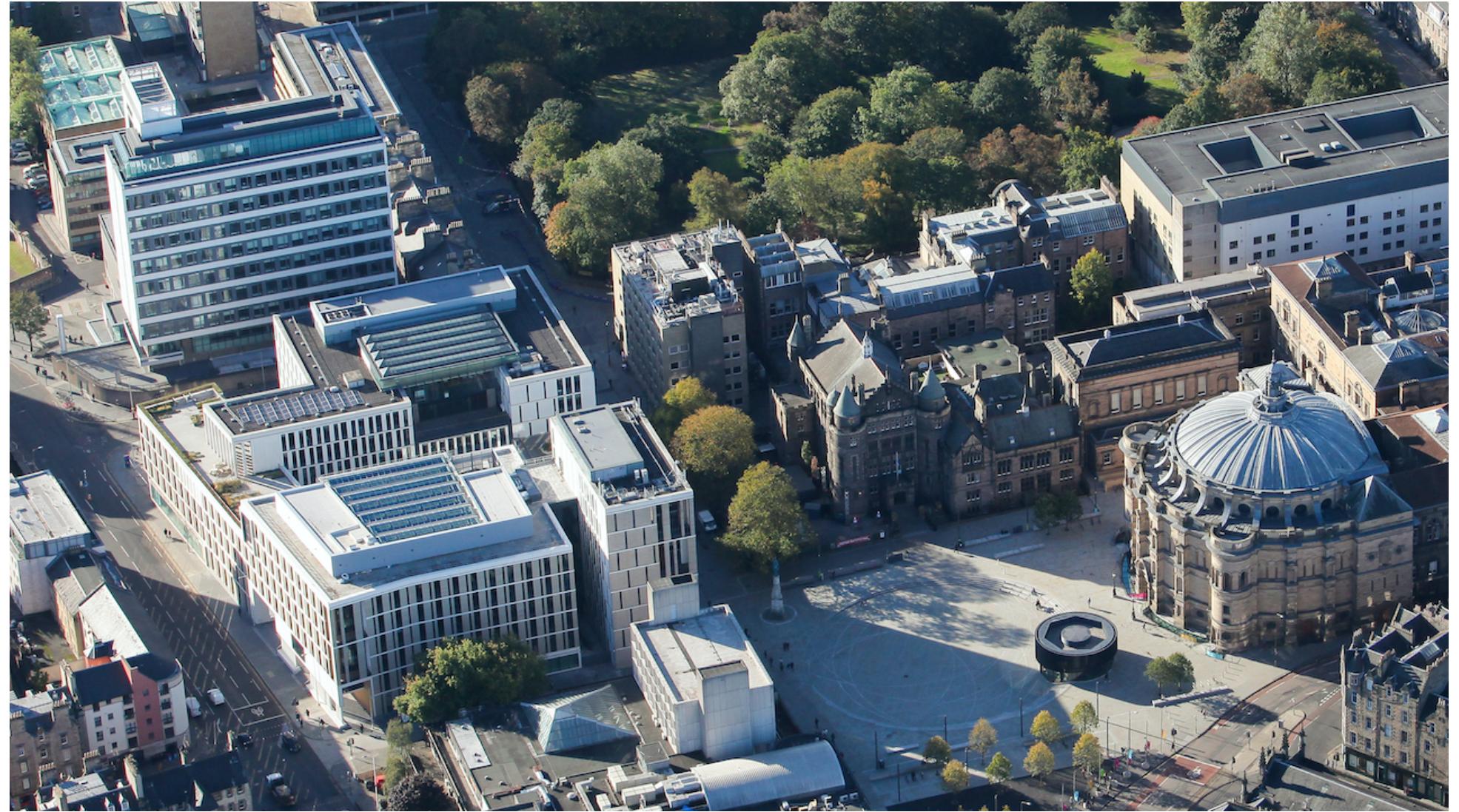
Start

What does a *Collaborative Session* notebook do?

"Collaborative Session" notebooks allow for two or more people to edit the same document by sharing a link. Collaborative Sessions launch in the newer JupyterLab interface. See our documentation for more information.

<https://noteable.edina.ac.uk/>

# Bioinformatics at Edinburgh



## Selected Research & Facilities

<http://web.inf.ed.ac.uk/anc/research/bioinformatics>

<https://www.ed.ac.uk/biology/research/themes/research-theme-systems>

<https://www.ed.ac.uk/mrc-human-genetics-unit/research/semple-group>

[http://bifx-core.bio.ed.ac.uk/wiki/index.php/Main\\_Page](http://bifx-core.bio.ed.ac.uk/wiki/index.php/Main_Page)

<https://genomics.ed.ac.uk>

## Seminar series:

<https://www.ed.ac.uk/biology/synthsys/news-and-events/events>

<http://www.ed.ac.uk/biology/news-events/events-and-seminars>

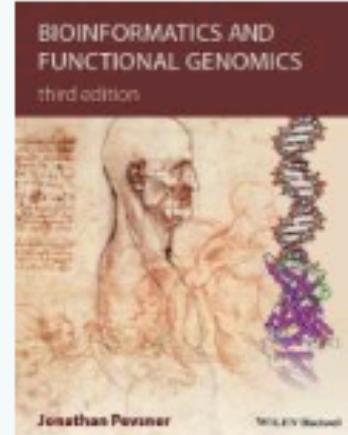
<https://www.ed.ac.uk/usher/news-events/seminars>

## Selected Topic Areas

- Database integration
- Data ontologies and provenance
- Evolutionary and genetic computation
- Gene expression databases
- High performance data structures for semi-structured data
- Microarray data analysis
- Epigenomics
- Natural language and bio-text mining
- Neural computation, visualisation and simulation
- Protein complex modelling
- Systems Biology, modelling
- Synthetic Biology
- Machine learning
- Medical Image Analysis

## Week 2 - Pairwise Sequence Alignment

If you would like to read ahead then please look at the following chapter, if not we will cover this in next week's lecture  
This is available from the Bio1 course "Resource List"



### BOOK Bioinformatics and functional genomics

Pevsner, Jonathan, 1961-, 3rd ed., Chichester, West Sussex, UK ; Hoboken, NJ, USA,  
John Wiley and Sons, Incorporated, 2015

*Note: Read Chapter 3, "Pairwise Sequence Alignment"*

 Add tags to item

**Complete**