# CSC 483 - Applied Biological Data Sciece W2D2

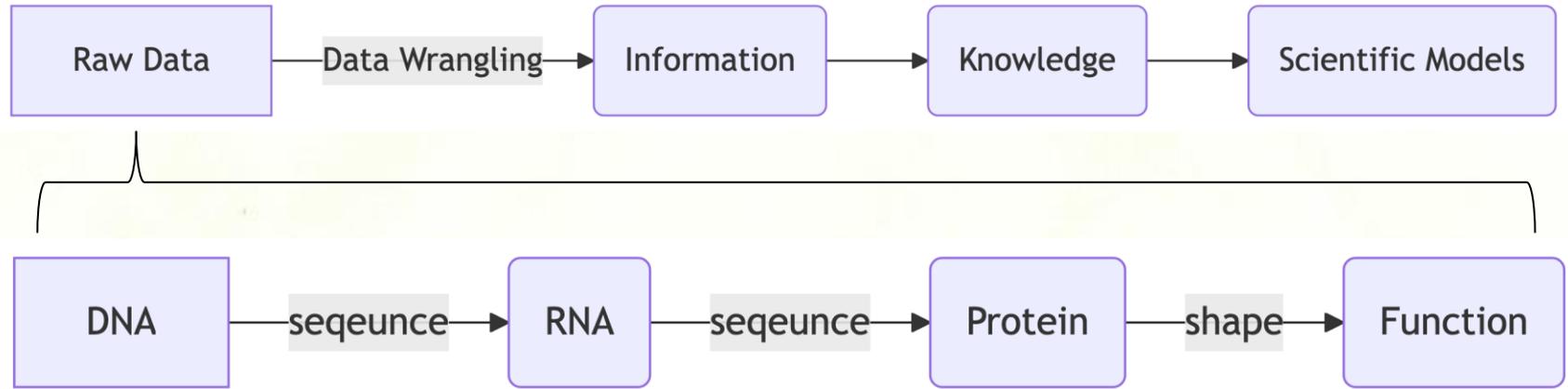# CSC 483: Applied Biological Data Science

- Georgia Doing, PhD (she/her/hers)
- email: doingg@union.edu
- office: Steinmetz 108B
  - office hours:
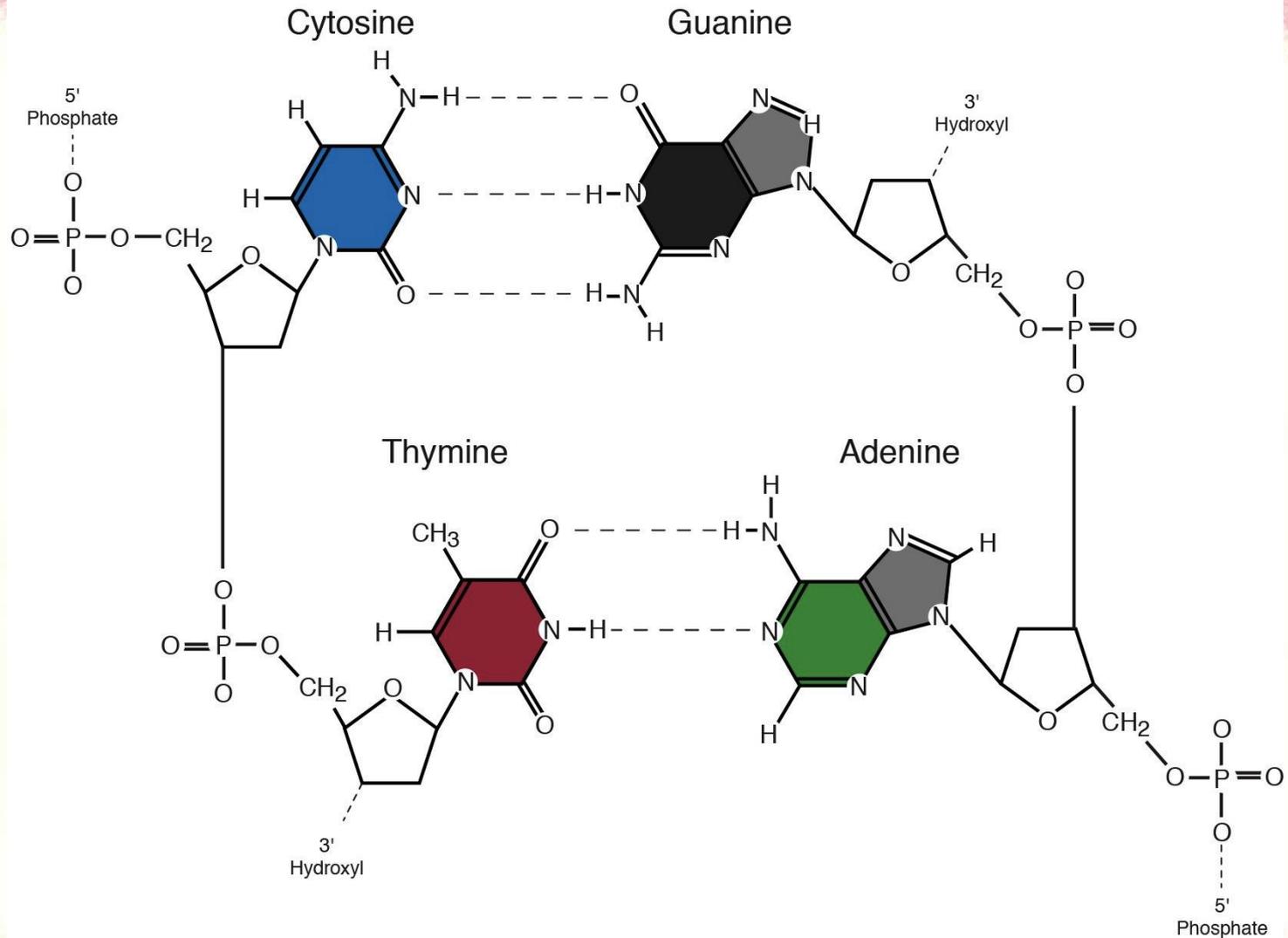    - Wednesday     2:00-3:30 pm
    - Thursday     4:00-5:30 pm

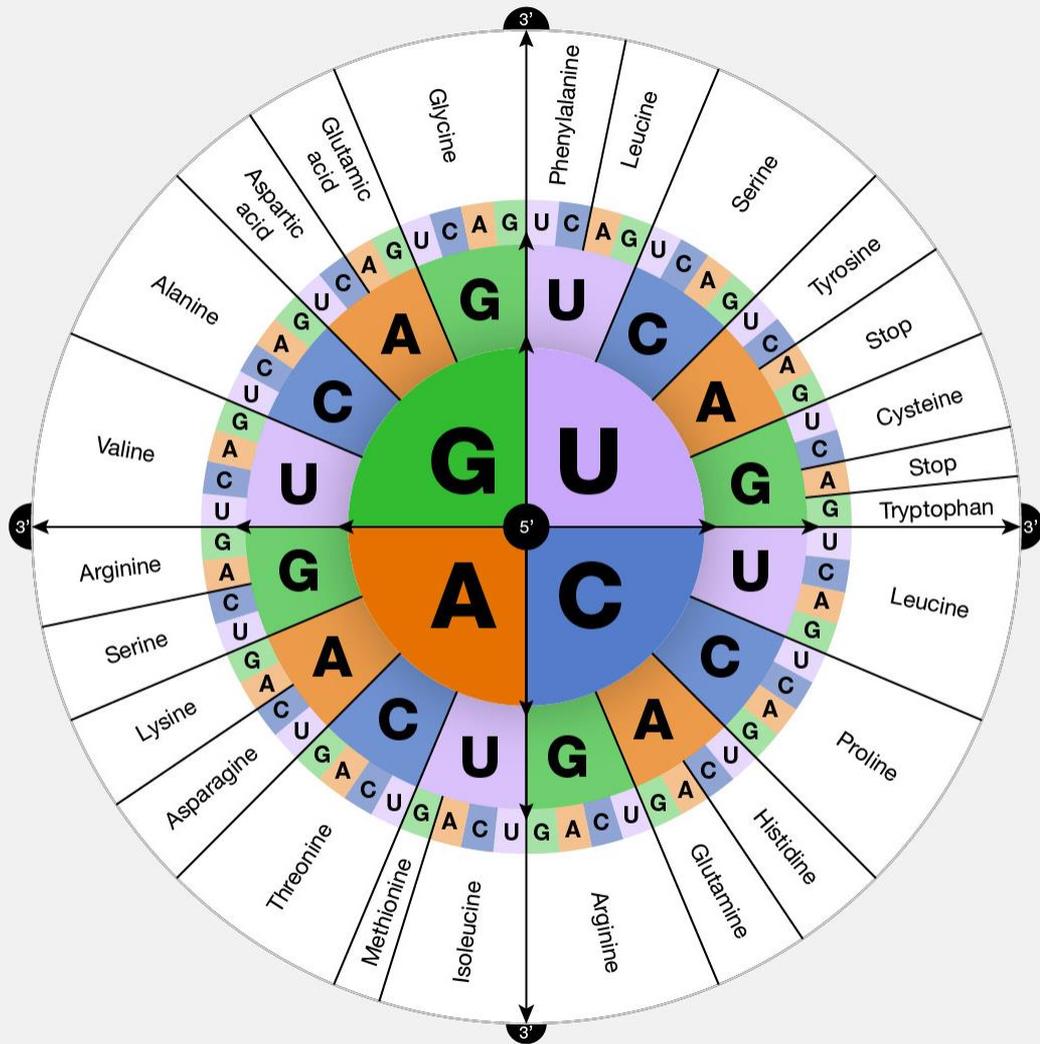Feel free to drop in whenever my office door is open or schedule an appointment.

# Today (145 min)

- Paper on methods/background

- Metadata summaries

- Explore data ,brainstorm ideas

- HW:
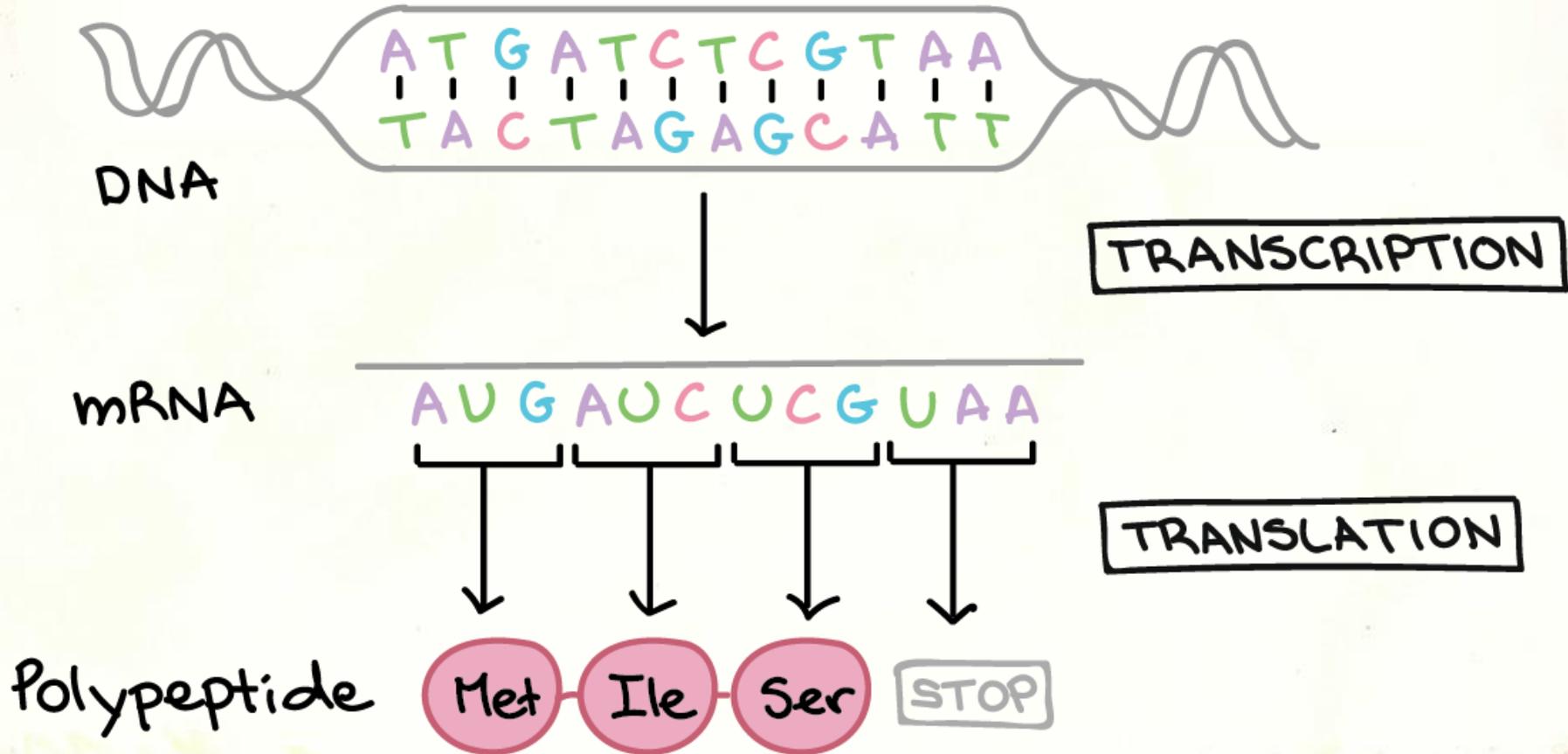  - rosalind problems on sequence analysis (w/ colab notebooks)
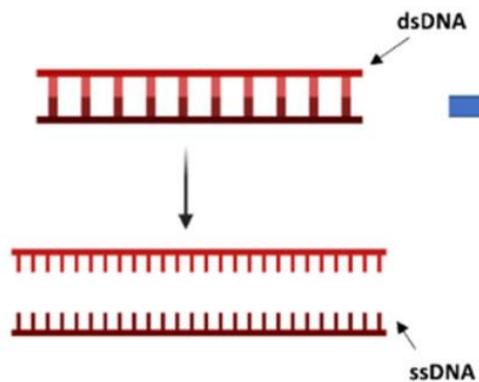  - due: tues
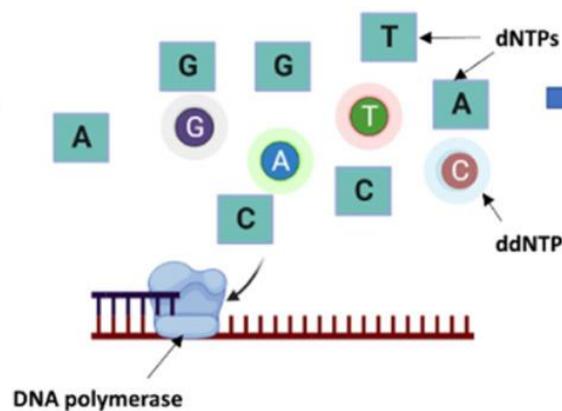
# What is Biological Data Science?

# THE CENTRAL DOGMA

DNA

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

TRANSCRIPTION

mRNA

A U G A U C U C G U A A
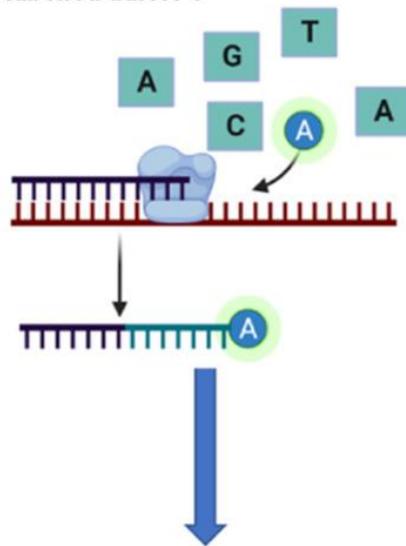
TRANSLATION

Polypeptide

Met — Ile — Ser — STOP

① Denaturation of dsDNA into ssDNA template through heat treatment

dsDNA

ssDNA

② Primer annealing and extension by the DNA polymerase by addition of contemporary dNTPs
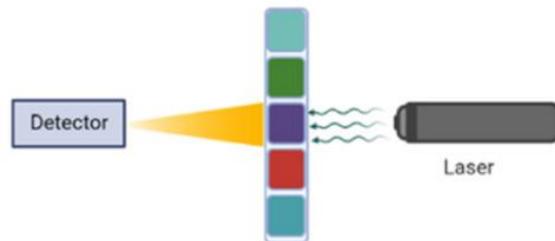
dNTPs

ddNTP

DNA polymerase

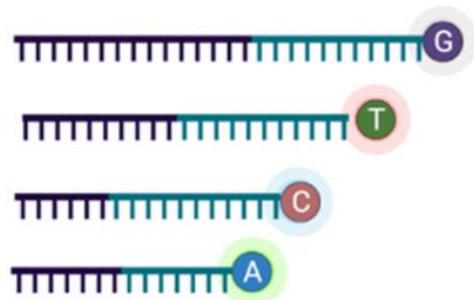③ Termination of each round of primer extension by the fluorescently labelled ddNTPs

⑥ Sequence analysis and reconstruction

⑤ Separation of chain-terminated oligonucleotides using gel electrophoresis, preferably single capillary gel electrophoresis
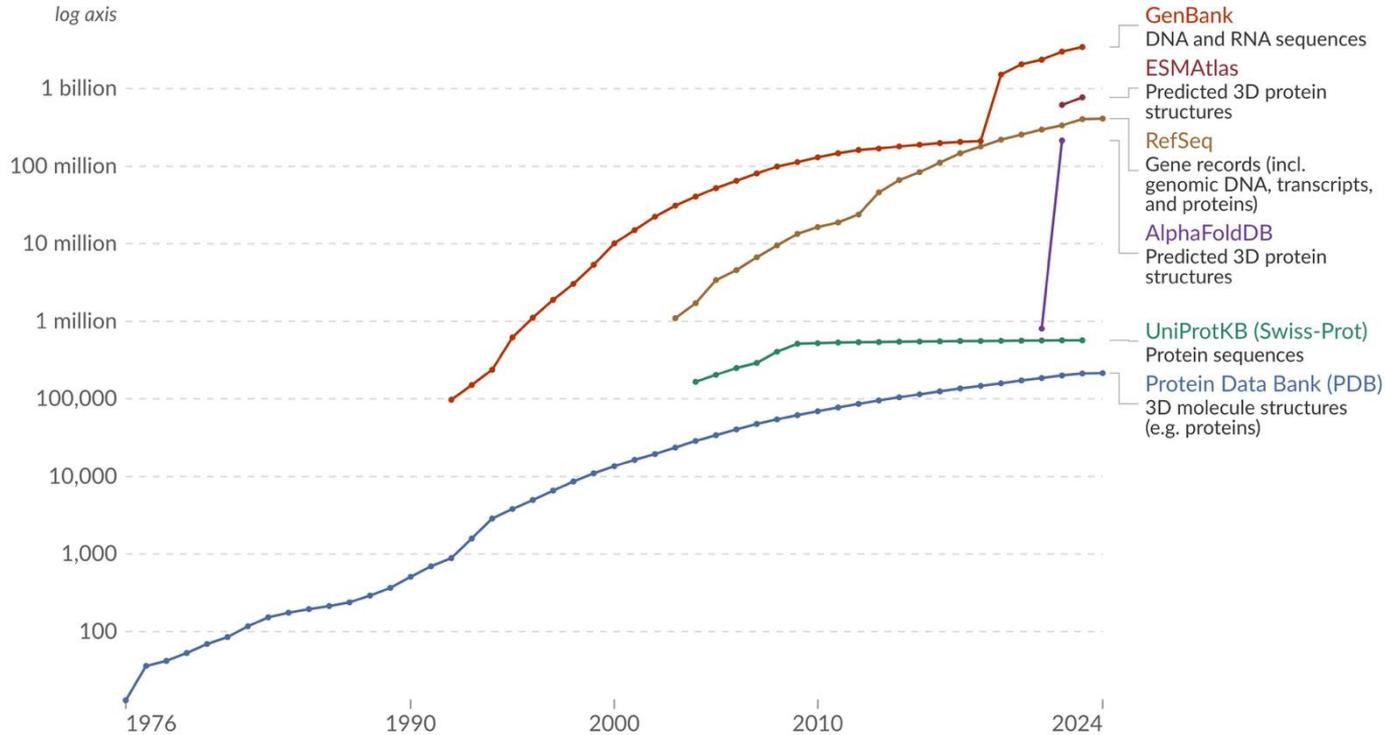
Detector

Laser

④ Fluorescently labelled DNA sample

# Number of entries in biological sequence databases

Biological sequence databases store data such as DNA, RNA, and amino acid sequences and 3D protein structures. This data includes entries from GenBank, RefSeq, PDB, UniProtKB/SwissProt, as well as predicted protein structures in AlphaFoldDB and ESMAtlas.

**GenBank** — DNA and RNA sequences
**ESMAtlas** — Predicted 3D protein structures
**RefSeq** — Gene records (incl. genomic DNA, transcripts, and proteins)
**AlphaFoldDB** — Predicted 3D protein structures
**UniProtKB (Swiss-Prot)** — Protein sequences
**Protein Data Bank (PDB)** — 3D molecule structures (e.g. proteins)

log axis

1 billion
100 million
10 million
1 million
100,000
10,000
1,000
100

1976   1990   2000   2010   2024

**Data source:** Epoch AI (2024)

CC BY

https://ourworldindata.org/grapher/number-of-entries-in-biological-sequence-databases

1000s Genomes Project: https://www.internationalgenome.org/

# Experimental Workflow

Oral 'animacules' (1676)

Oral microbiome (2011)



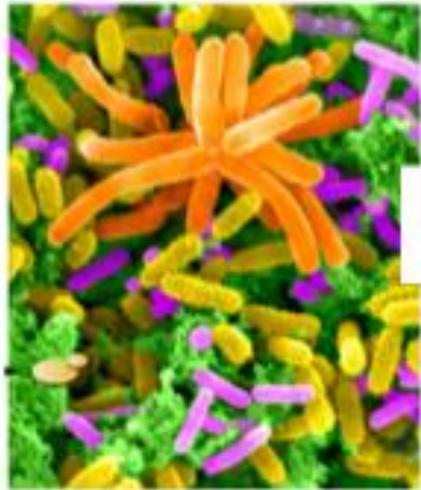A. van Leeuwenhoek

Valm et al., PNAS
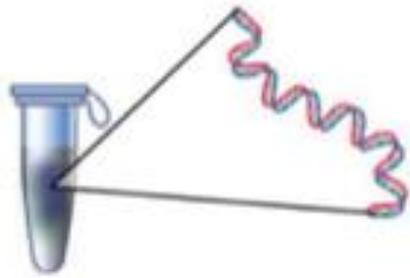
**Isolate**

Genomics / Transcriptomics
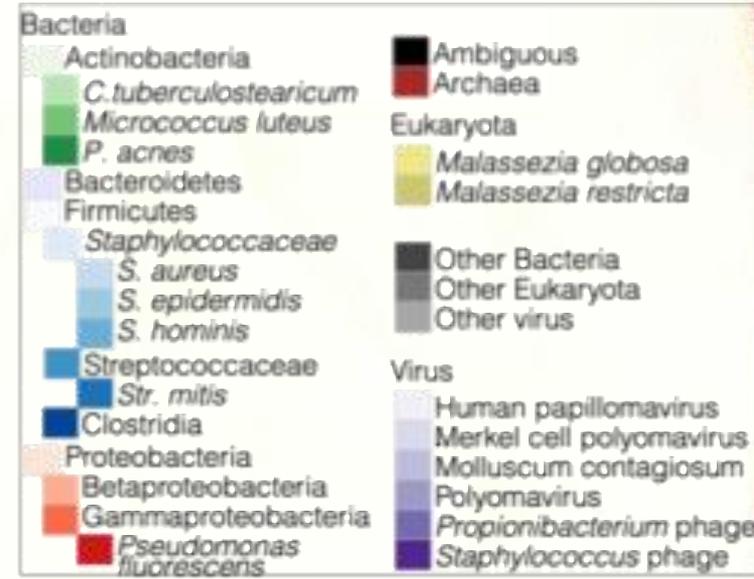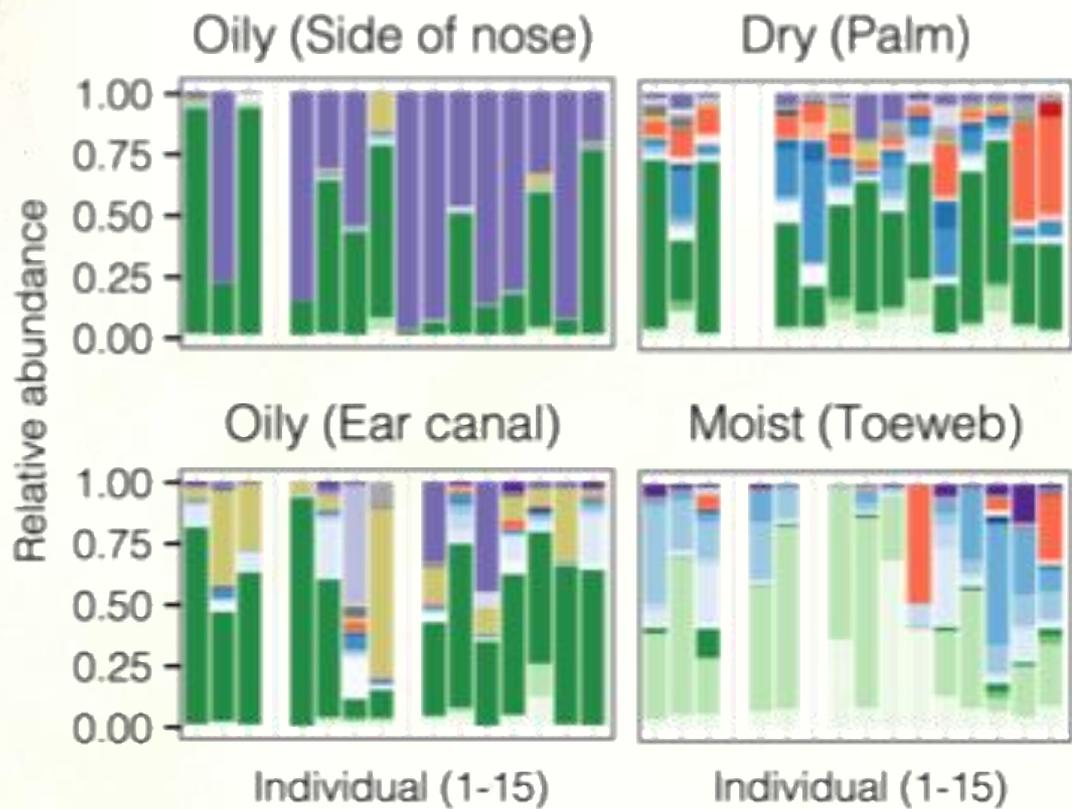
**Community**

Metagenomics / Metatranscriptomics

Microbial Community       Molecular Extraction       Sequencing

Oily (Side of nose)　　Dry (Palm)

Oily (Ear canal)　　Moist (Toeweb)

Individual (1-15)　　Individual (1-15)

Relative abundance

Bacteria
Actinobacteria
  *C. tuberculostearicum*
  *Micrococcus luteus*
  *P. acnes*
Bacteroidetes
Firmicutes
  *Staphylococcaceae*
    *S. aureus*
    *S. epidermidis*
    *S. hominis*
  Streptococcaceae
    *Str. mitis*
  Clostridia
Proteobacteria
  Betaproteobacteria
  Gammaproteobacteria
    *Pseudomonas fluorescens*

Ambiguous
Archaea

Eukaryota
  *Malassezia globosa*
  *Malassezia restricta*

Other Bacteria
Other Eukaryota
Other virus

Virus
  Human papillomavirus
  Merkel cell polyomavirus
  Molluscum contagiosum
  Polyomavirus
  *Propionibacterium* phage
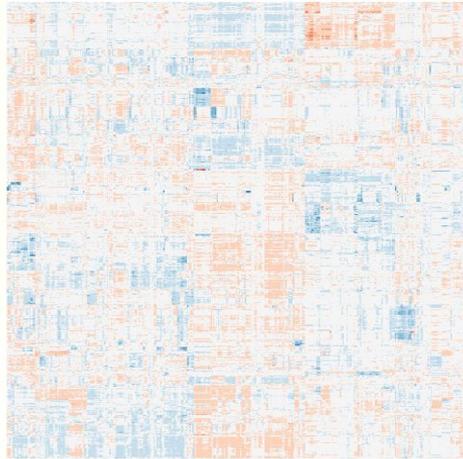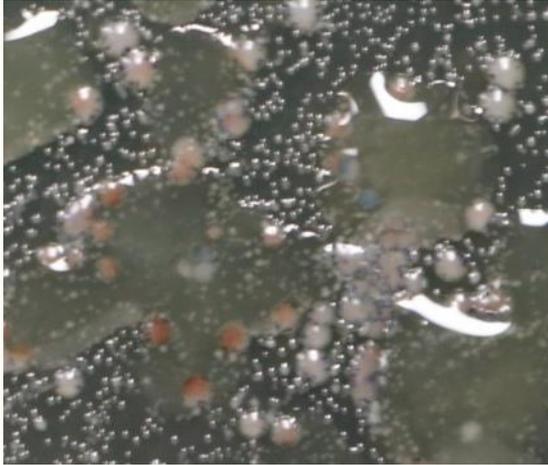  *Staphylococcus* phage
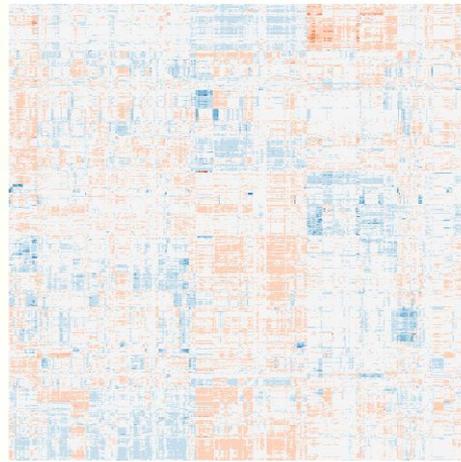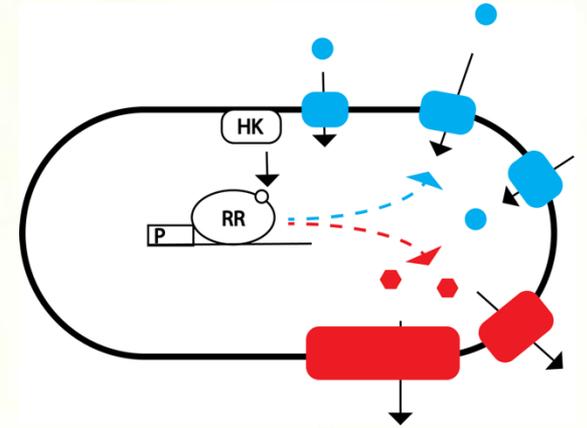
Oh et al., *Nature*

Microbial Culture

Microbial Culture

RNA-seq Data

Microbial Culture

RNA-seq Data

Biological Model
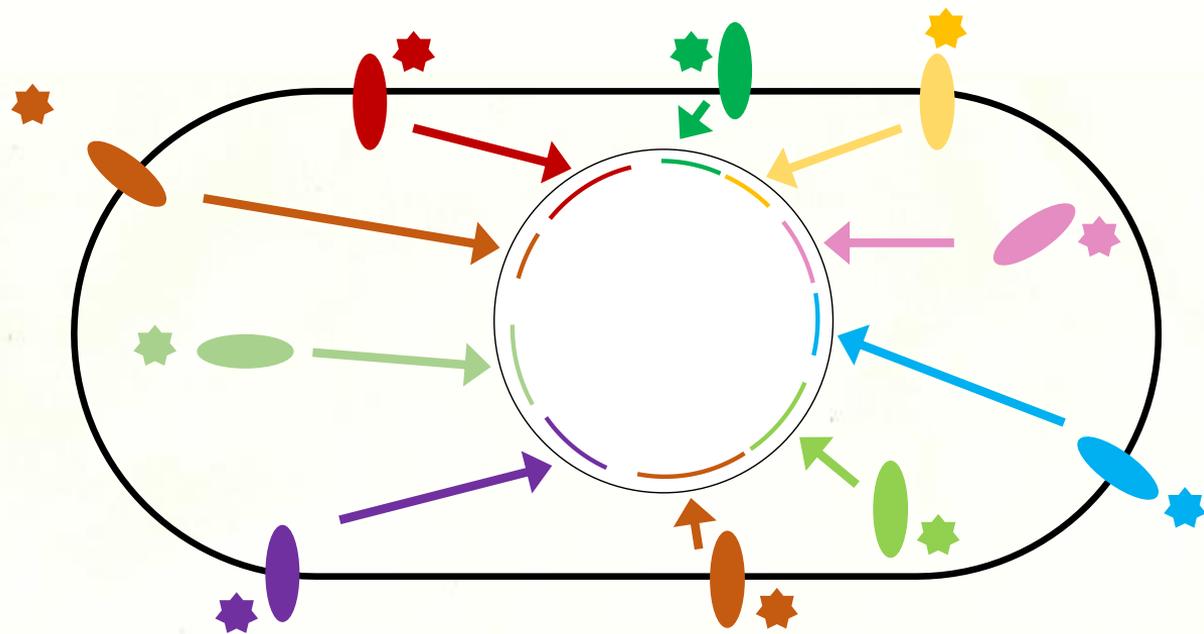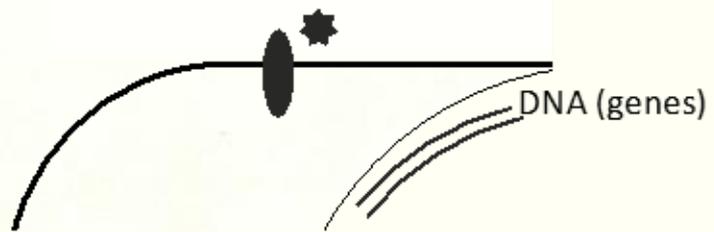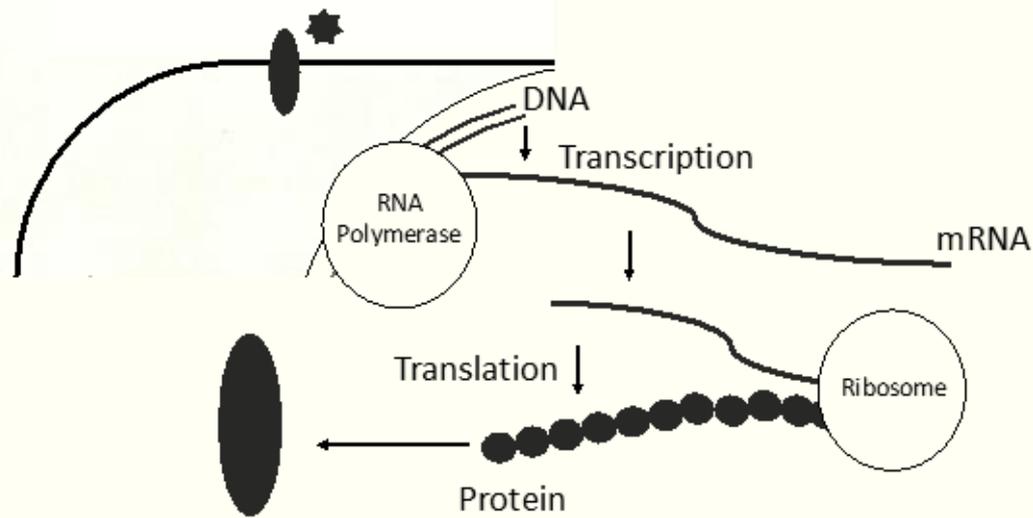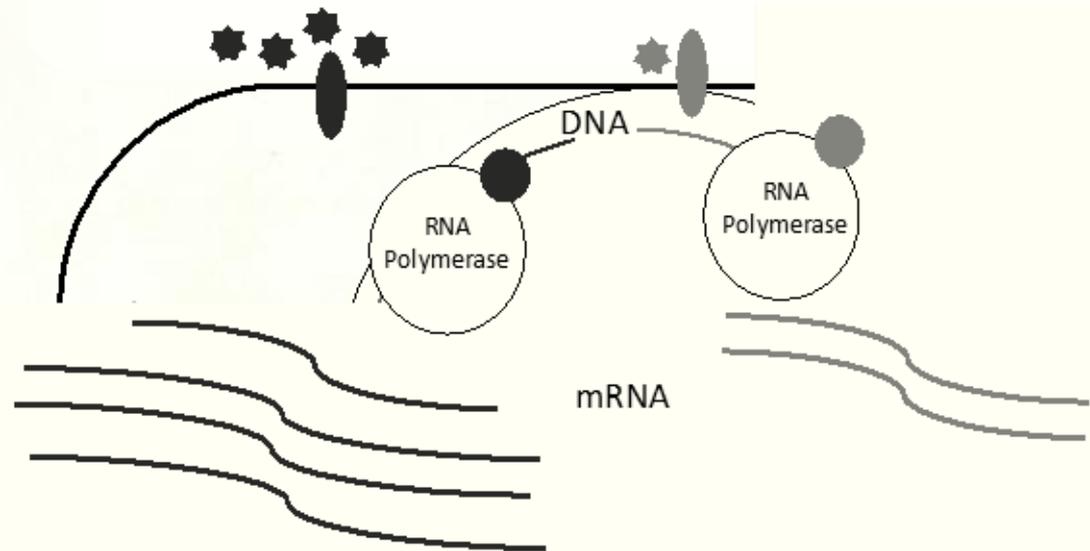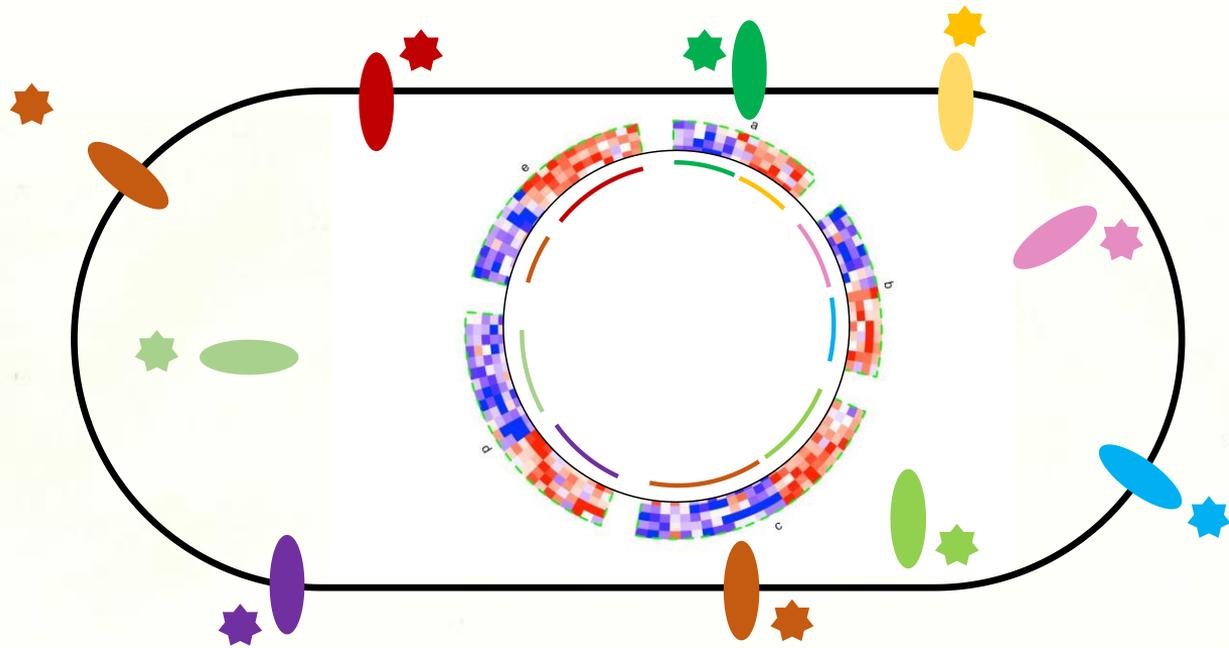
HK

RR

P

Molecular cues

Proteins sensors

DNA (genes)

DNA

RNA
Polymerase

mRNA

Translation

Ribosome

Protein

DNA

RNA
Polymerase

RNA
Polymerase

mRNA

Control
PAO1
PA14
0 µg ml⁻¹RA
0 µg ml⁻¹RA
0 µg ml⁻¹RA

Cystic Fibrosis

MULTIDRUG-RESISTANT
PSEUDOMONAS AERUGINOSA

# Computationally Efficient Assembly of *Pseudomonas aeruginosa* Gene Expression Compendia

Georgia Doing,[a] Alexandra J. Lee,[b] Samuel L. Neff,[a] Taylor Reiter,[c] Jacob D. Holt,[a] Bruce A. Stanton,[a] Casey S. Greene,[c,d] Deborah A. Hogan[a]

[a]Department of Microbiology and Immunology, Geisel School of Medicine at Dartmouth, Hanover, New Hampshire, USA
[b]Genomics and Computational Biology Graduate Program, University of Pennsylvania, Philadelphia, Pennsylvania, USA
[c]Department of Biochemistry and Molecular Genetics, University of Colorado School of Medicine, Denver, Colorado, USA
[d]Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, USA

https://pubmed.ncbi.nlm.nih.gov/36541761/

# ANTIMICROBIAL RESISTANCE THREATS
## in the United States, 2021-2022

CDC used new data[1] to analyze the U.S. burden of the following antimicrobial-resistant pathogens typically found in healthcare settings:


**Carbapenem-resistant Enterobacterales (CRE)**


**Carbapenem-resistant *Acinetobacter***


***Candida auris* (*C. auris*)**


**Methicillin-resistant *Staphylococcus aureus* (MRSA)**


**Vancomycin-resistant Enterococcus (VRE)**


**Extended-spectrum beta-lactamase (ESBL)-producing Enterobacterales**
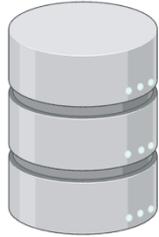

**Multidrug-resistant (MDR) *Pseudomonas aeruginosa***

CDC previously reported that the burden of these pathogens increased in the United States in 2020 in the COVID-19 Impact Report. The information below describes the burden in the two following years, 2021 and 2022, and compares against 2019 data.
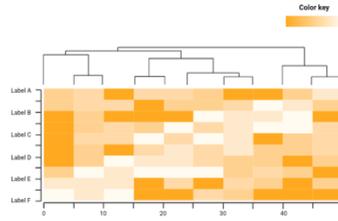
CDC https://www.cdc.gov/ecoli/about/index.html

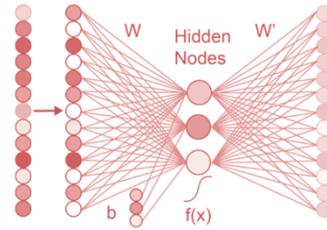| SRA Search | Compendia | DAE Models | Bi-partite Graphs | Comparative Genetics |
|---|---|---|---|---|
|  |  |  |  |  |

- organism: *E. coli*
- molecule: mRNA
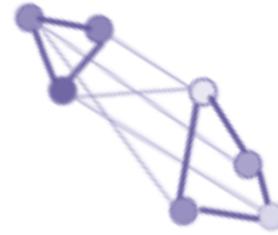- pull data from Short Read Archive (SRA) and Gene Expression Omnibus (GEO)

- metadata
- *E. coli* (20 strains) reference pan-genomes
- salmon k-mer mapping (k=15)
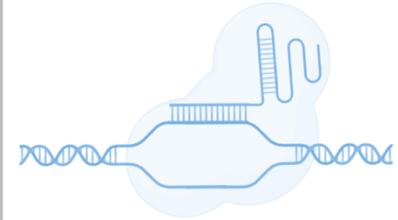
- denoising autoencoder (tied)
- 1 layer of 50 hidden nodes
- fully connected
- objective): reconstruction error

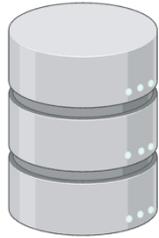- $G(E,V)$ where $E \in Corr_{gene, gene}$ and $V \in Genes_{S.e.} \cup Genes_{S.a.}$
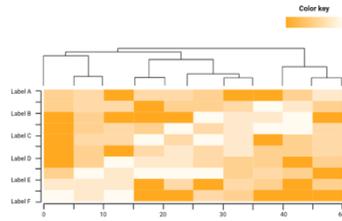- Weight edges by correlation strength
- Color genes by homology

- knock-down genes of interest under function-relevant conditions
- look for fitness effects and/or phenotypes

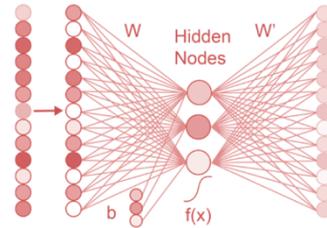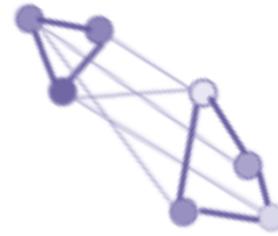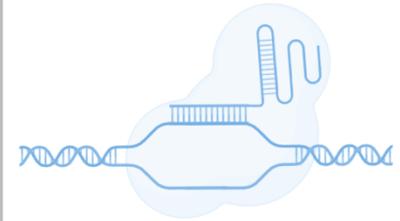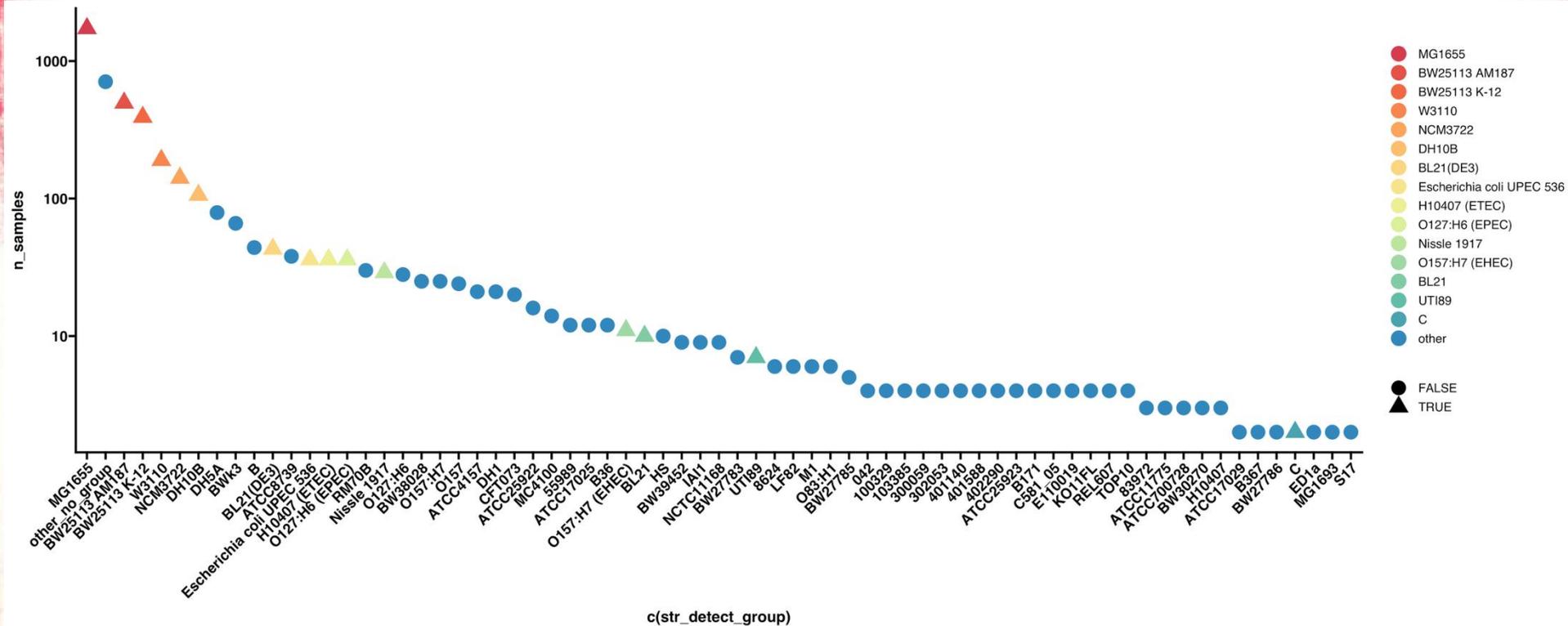| SRA Search | Compendia | DAE Models | Bi-partite Graphs | Comparative Genetics |
|---|---|---|---|---|
| • organism: *E. coli*<br>• molecule: DNA, RNA or text<br>• pull data from Short Read Archive (SRA) and Gene Expression Omnibus (GEO) | • metadata<br>• subset of strains and samples reference pan-genomes<br>• salmon k-mer mapping (k=15) | Dimensionality reduction:<br><br>PCA<br>ICA<br>NMD<br>DAE | Pattern detection:<br><br>correlation regression clustering classification | Hypothesis testing:<br><br>CRISPR data phenotypes (future directions) |

# Ecoli Pangenome Compendium

samples: 20,292

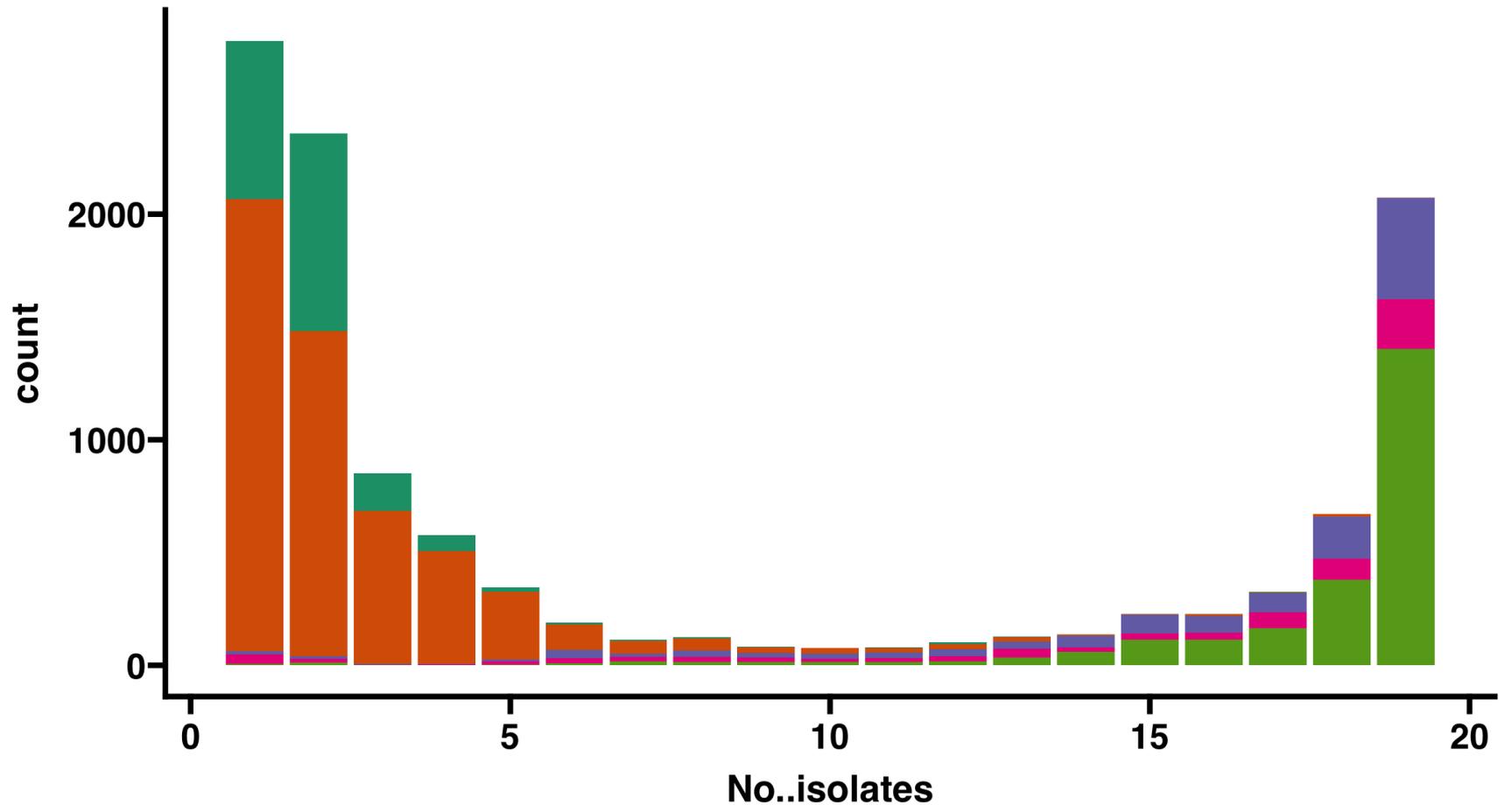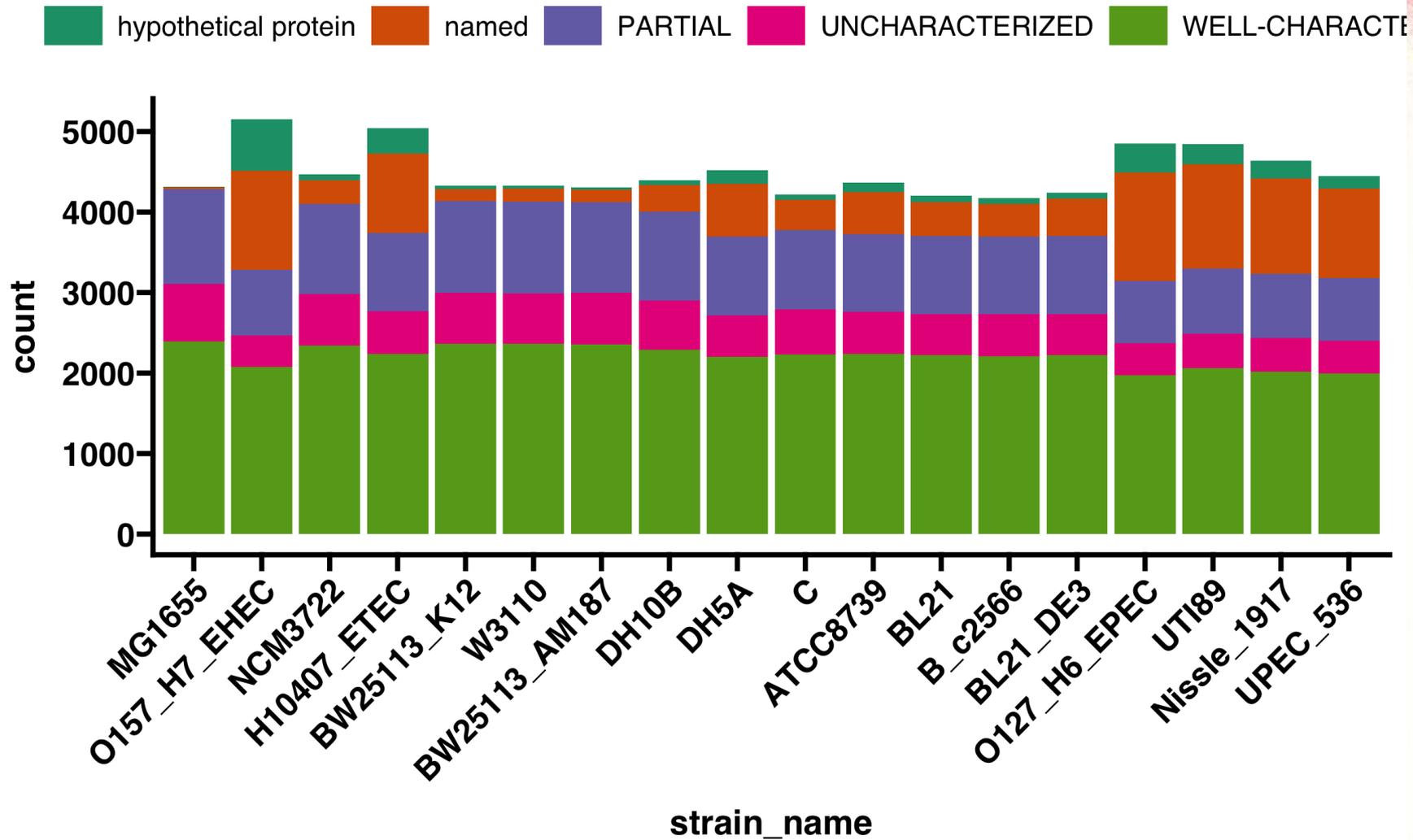unique studies: 1,365

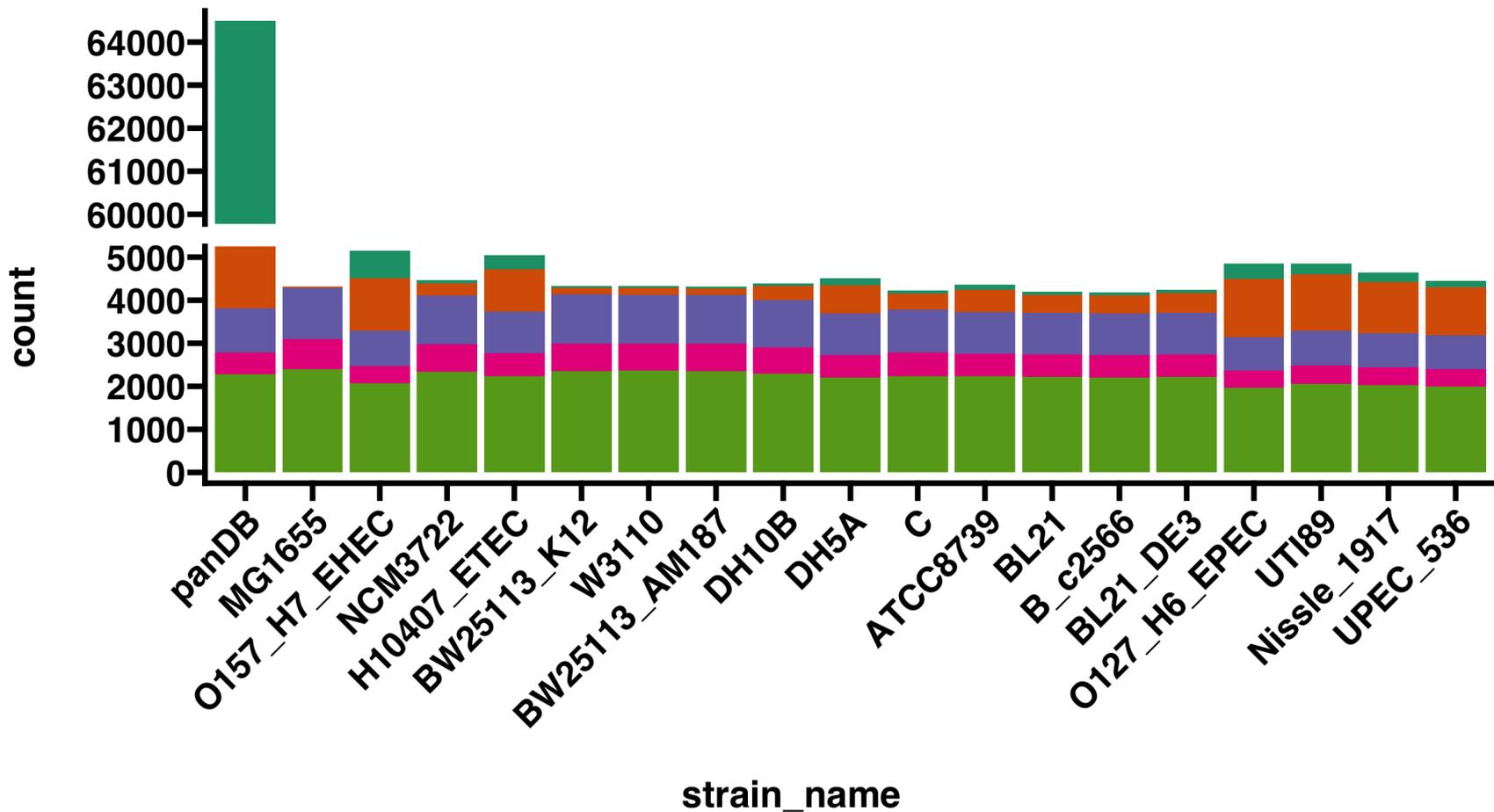samples with high confidence strain: 4,260

samples with putative strain: 11,859

# Head of "count" table

| X | Name | ERX1805690 | ERX1805691 | ERX1805692 | ERX2279850 | ERX2279851 | ERX2279852 | ERX2279853 | ERX2279854 |
|---|------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 0ed86576d343ce011a61773e0620e335_1045 | character(0) | 1.000 | 0.000 | 0.00 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 0ed86576d343ce011a61773e0620e335_1047 | character(0) | 0.000 | 0.000 | 0.00 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 0ed86576d343ce011a61773e0620e335_1120 | character(0) | 0.000 | 0.000 | 0.00 | 6.476 | 6.236 | 0.000 | 0.000 | 7.554 |
| 0ed86576d343ce011a61773e0620e335_1175 | character(0) | 0.000 | 0.000 | 0.00 | 2.000 | 2.000 | 3.000 | 2.000 | 0.000 |
| 0ed86576d343ce011a61773e0620e335_1179 | character(0) | 0.000 | 0.000 | 0.00 | 2.000 | 4.000 | 0.000 | 4.000 | 2.000 |
| 0ed86576d343ce011a61773e0620e335_1311 | character(0) | 0.000 | 0.000 | 0.00 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 0ed86576d343ce011a61773e0620e335_1450 | character(0) | 0.000 | 0.000 | 0.00 | 3.000 | 0.000 | 0.000 | 0.000 | 2.000 |
| 0ed86576d343ce011a61773e0620e335_1476 | character(0) | 5.403 | 2.351 | 0.00 | 15.370 | 11.212 | 23.317 | 5.232 | 12.820 |
| 0ed86576d343ce011a61773e0620e335_1484 | character(0) | 0.000 | 9.229 | 10.55 | 2.000 | 2.000 | 0.000 | 0.000 | 3.000 |
| 0ed86576d343ce011a61773e0620e335_2071 | character(0) | 0.000 | 0.000 | 0.00 | 0.000 | 0.000 | 36.005 | 0.000 | 0.000 |

From DEseq2, median of ratios normalization:

$(Y_{ij})$: The raw count for gene $(i)$ in sample $(j)$.
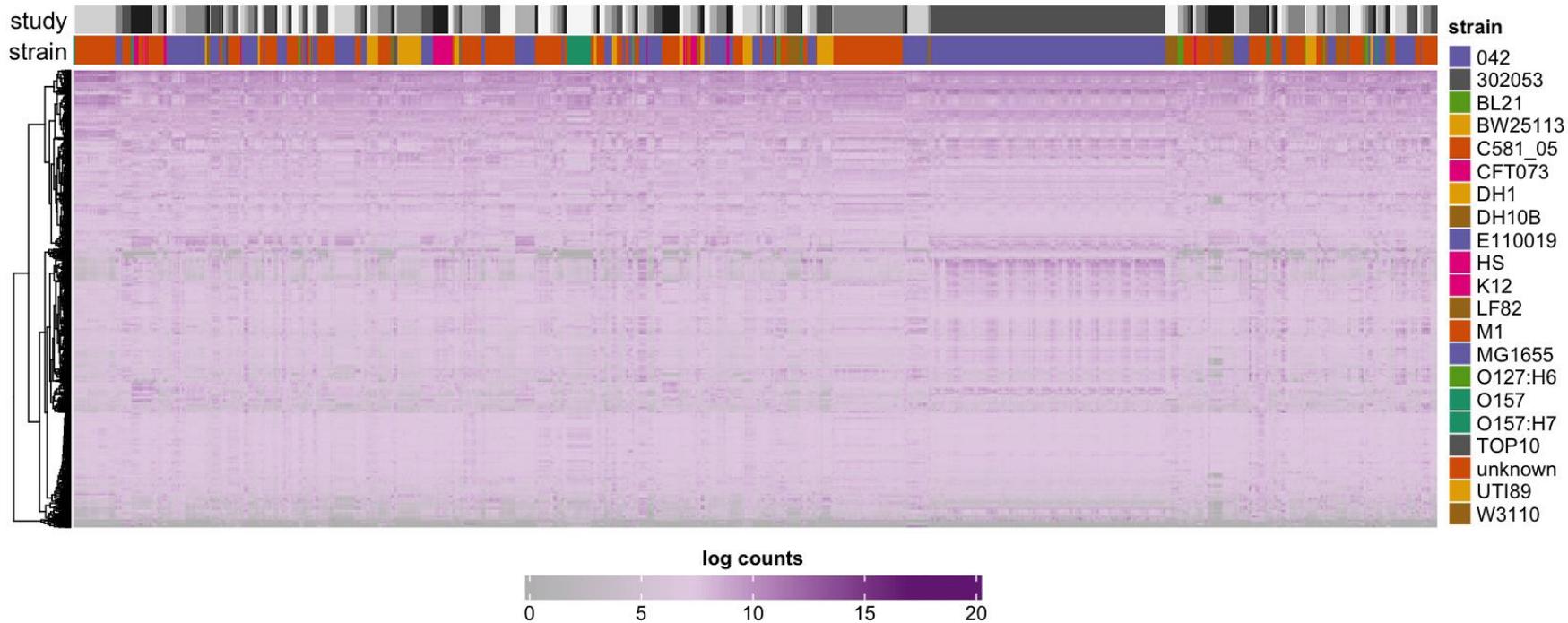$(g_i)$ :The geometric mean of counts for gene $(i)$ across all samples.
The size factor for sample $(j)$ ,denoted as $(S_j)$ ,is calculated as:

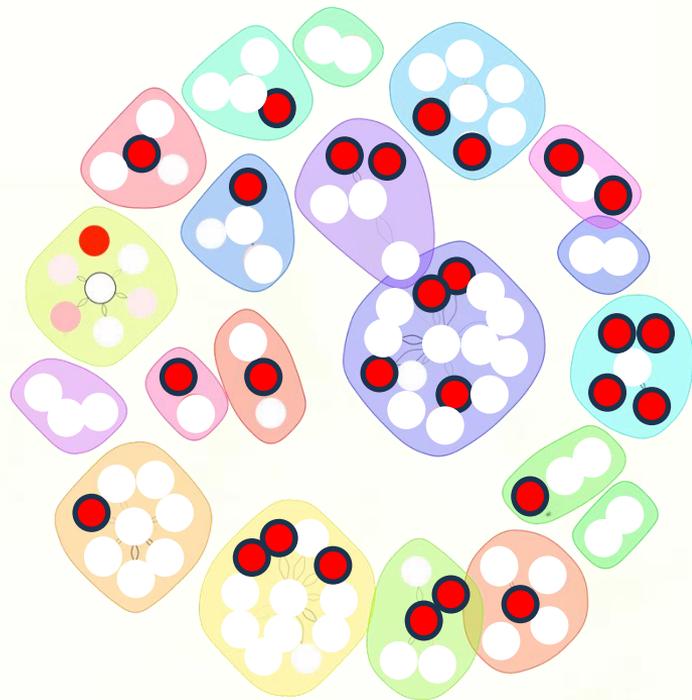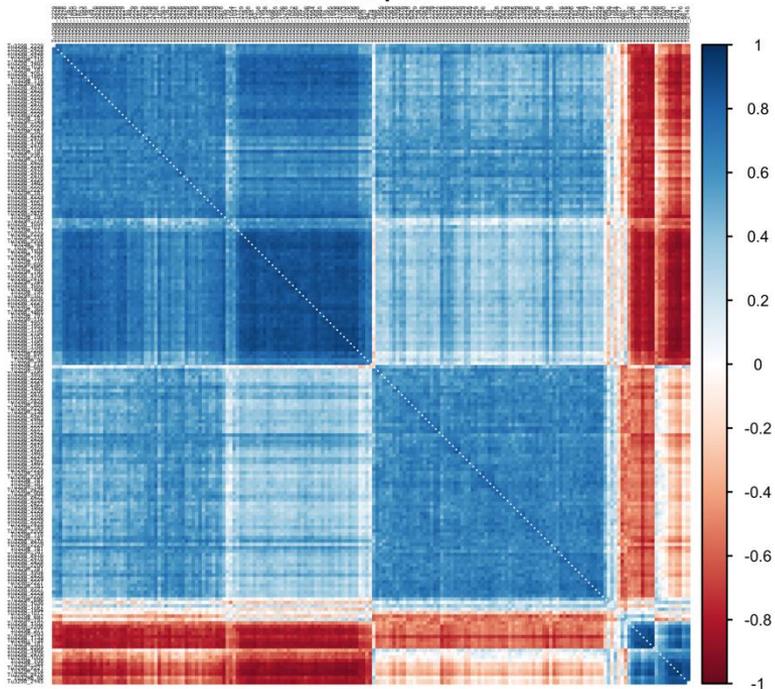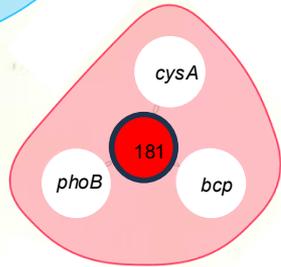$$\left[ S_j = \text{median}\left(\frac{Y_{ij}}{g_i}\right)_{i \in 1,...,m} \right]$$

Where:
$\left( g_i = \left(\prod_{j=1}^{n} Y\,ij\right)^{1/n} \right)$ ,the geometric mean of counts for gene $(i)$ across $(n)$ samples.

The ratio $\left(\frac{Y_{ij}}{g_i}\right)$ is computed for each gene $(i)$ in sample $(j)$ .The median of these ratios is taken as the size factor $(S_j)$ for sample $(j)$.
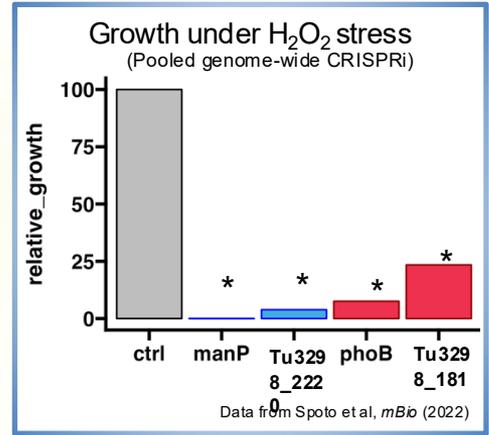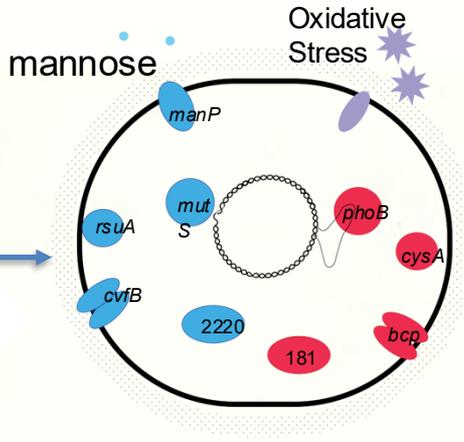
S.e. model post fine-tune

# Lab machines: getting started

- Log on:
    - Your Union usernames
    - Default password: the word union followed by your ID number; e.g., union12345

**Always log off before leaving the lab**