# Evolutionary and Structural Analysis of Pathogen Proteins.

Final year UG project 2024-25
2024-10-21 (Week 5)
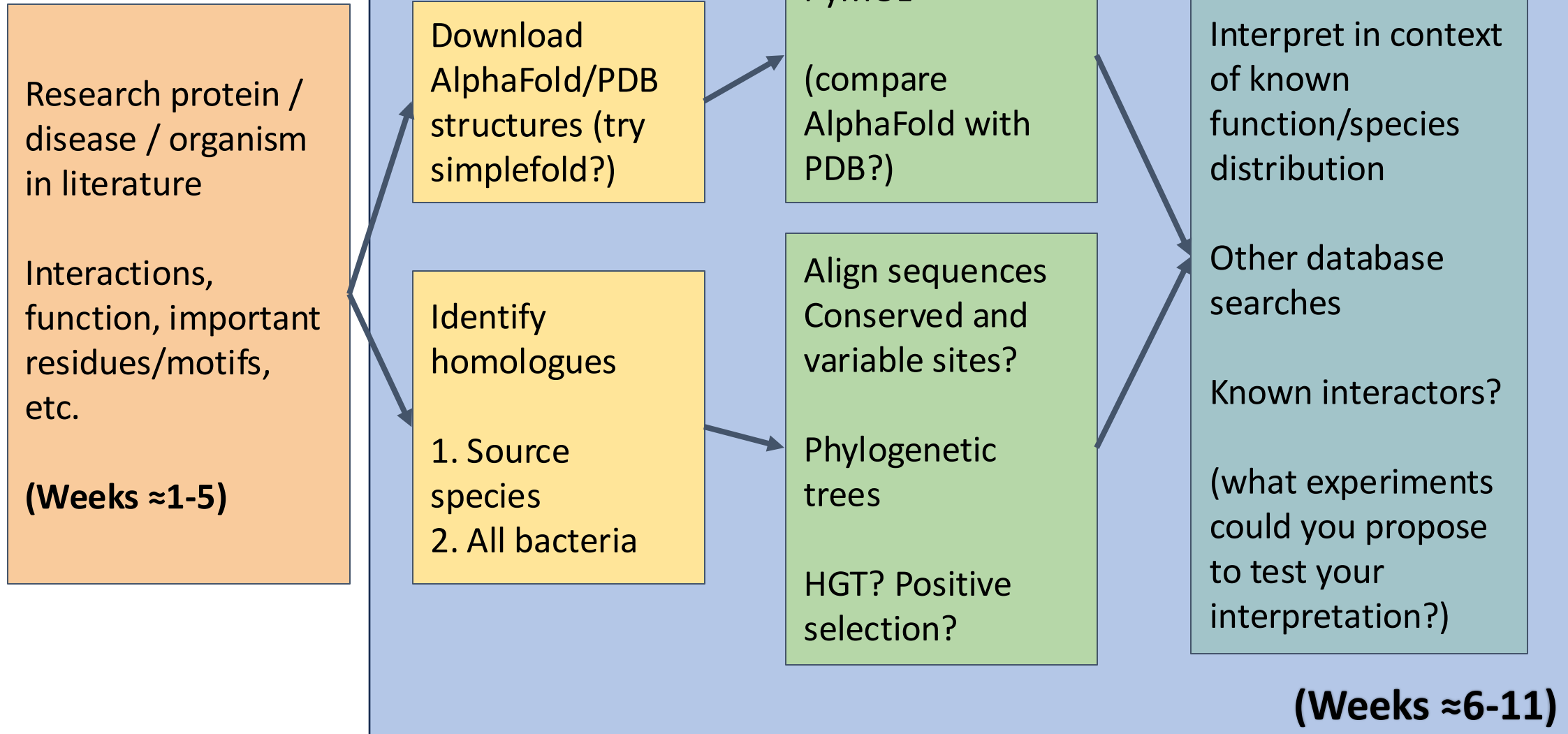
# Candidate proteins – start points

**Any changes needed?**

| Organism | Host | Gene/Protein | PHI accession | Student |
|----------|------|--------------|---------------|---------|
| *Escherichia coli* | *Homo sapiens* | *espY* | PHI:8647 | LB |
| *Shigella flexneri* | *Homo sapiens* | *ipaJ* | PHI:9253 | LT |
| *Candida albicans* | *Mus musculus* | *sap6* | PHI:10193 | IM |
| *Pseudomonas aeruginosa* | *Homo sapiens* | *tplE* | PHI:6646 | AE |
| *Vibrio vulnificus* | *Mus musculus* | *vvhA* | PHI:6877 | JT |

http://www.phi-base.org/

# Workflow

**Research protein / disease / organism in literature**

Interactions, function, important residues/motifs, etc.

**(Weeks ≈1-5)**

Download AlphaFold/PDB structures (try simplefold?)

Identify homologues

1. Source species
2. All bacteria

Visualise with PyMOL

(compare AlphaFold with PDB?)

Align sequences Conserved and variable sites?

Phylogenetic trees

HGT? Positive selection?

Map conservation onto 3D structure

Interpret in context of known function/species distribution

Other database searches

Known interactors?

(what experiments could you propose to test your interpretation?)

**(Weeks ≈6-11)**

# Your questions/comments

(What would you like to talk about?)

# Thesis Introductions

# A very quick introduction to building a phylogenetic tree

# A Family Tree

- Family trees are not a good model of how bacteria (or species in general) evolve

# A Brief Introduction to Phylogenetics

- Online introductory course:
https://www.ebi.ac.uk/training/online/courses/introduction-to-phylogenetics/
- Conor Meehan's introductory course:
https://conmeehan.github.io/PathogenDataCourse/IntroToPhylogenetics.html

- Phylogenetics is the reconstruction of evolutionary history from genetic/genomic data
  - **Input**: Aligned protein sequence data
  - **Output**: A tree estimating evolutionary relationships
- **Phylogenetic reconstruction is a mathematical activity**
- The biology in phylogenetics comes from three places:
  - Aligning the input sequence set correctly (evolutionary equivalence)
  - The model of substitution used (e.g. how likely is residue/base A to be replaced/substituted by residue/base B?)
  - The assumption of a bifurcating tree (this doesn't apply to some methods, e.g. splitstree, but other assumptions do apply there)

# Phylogenetic Trees (Topology)

- We assume that species evolve by a series of branching events
  - e.g. assume that species cannot interbreed so, when one species splits into two, it is an **irrevocable** branching event
- https://sipbs-compbiol.github.io/BM211-Workshop-5/



(A, B) are more closely related to each other than to C

Three different relationships:

(A, B) more closely related to each other than to C
(A, C) more closely related to each other than to B
(B, C) more closely related to each other than to A

# Input sequence alignment

- The goal is that each column in the alignment represents a single *evolutionarily equivalent* position – subject to similar selection pressures
  - We can then make inferences based on what changes are permitted at that position
  - Structural equivalence can (but does not always) imply functional equivalence
- Many gaps in a column mean information is missing and inference is less robust (they bias the tree)
  - Remove "gappy columns", e.g. https://vicfero.github.io/trimal/
- The larger amino acid alphabet means that alignments are more robust than nucleotide sequence alignments
- Codon degeneracy means that amino acid alignments can mask relevant evolutionary change, or mask "saturation"
  - Best practice for low/moderate divergence: align amino acid sequences and backtrace the coding sequence to nucleotide to make the tree
  - For highly divergent sequences, amino acid-based trees may be more robust

# How To Make a (Simple) Tree in Galaxy

- Start with a FASTA (protein) sequence alignment

# How To Make a (Simple) Tree in Galaxy

- Use IQ-Tree to generate a phylogenetic tree: **SPECIFY AA sequence type!**

# How To Make a (Simple) Tree in Galaxy

- Produces more than one tree
  - BioNJ: Neighbour-Joining (tree-building algorithm)
  - Maximum Likelihood (fitting a tree to the data)
- **Trees produced by different approaches (or with different parameters/inputs) are often different**
  - **This is not bad!** It's something to note in the discussion – is the tree *robust*?
  - **Justify choices** (as much as possible) in the thesis
  - There are many parameters/options to choose – it's fine to use defaults, but state clearly that you did so in your Methods.

15: IQ-TREE on data 8: MaxLikeli 👁 ✏ 🗑
hood Tree

14: IQ-TREE on data 8: BIONJ Tr 👁 ✏ 🗑
ee

# How To Make a (Simple) Tree in Galaxy

- Raw tree data (Newick format) looks cryptic – you don't have to read this data yourself: it's for computers

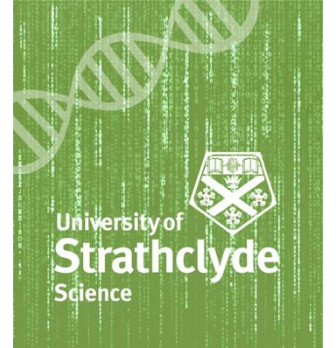# How To Make a (Simple) Tree in Galaxy

- You can visualize your tree(s) directly in Galaxy (e.g. with Newick Display)

# How To Make a (Simple) Tree in Galaxy

- You can visualize your tree(s) directly in Galaxy (e.g. with Newick Display)
    - TBH it doesn't do a great job

# How To Make a (Simple) Tree in Galaxy

- Download your (Newick) tree file
  - This will let you use better visualization tools
- FigTree
  - http://tree.bio.ed.ac.uk/software/Figtree/
- DendroScope
  - https://uni-tuebingen.de/en/fakultaeten/mathematisch-naturwissenschaftliche-fakultaet/fachbereiche/informatik/lehrstuehle/algorithms-in-bioinformatics/software/dendroscope/
- iTOL (Interactive Tree of Life)
  - https://itol.embl.de/

# FigTree

- Open the downloaded Newick (.nhx) file

# FigTree

- Use the sliders and tree rooting options to make the tree more legible

# FigTree

- Use tree layout options to make the whole tree easier to see/interpret

# FigTree

- Use selection and colour options to highlight groups of sequences for the reader

# FigTree

- Export .png, .jpg, .pdf, .svg files for inserting into your thesis (or presentation)

# Going further

- Using nucleotide alignments can be more informative and accurate, but backtracing can be tricky – not all proteins have a known coding sequence
  - https://ncfp.readthedocs.io/en/stable/
- Identify the most appropriate substitution model before building the tree
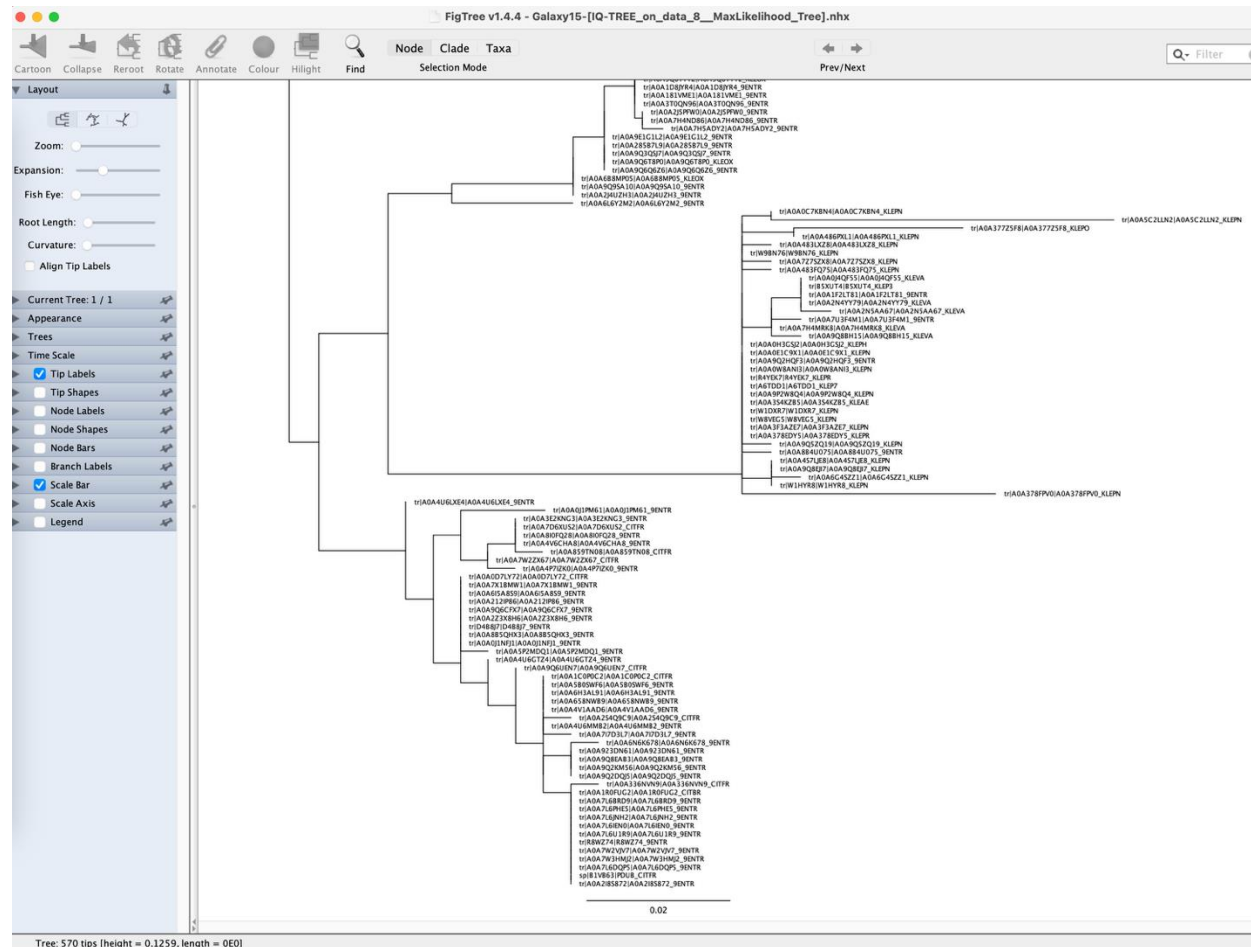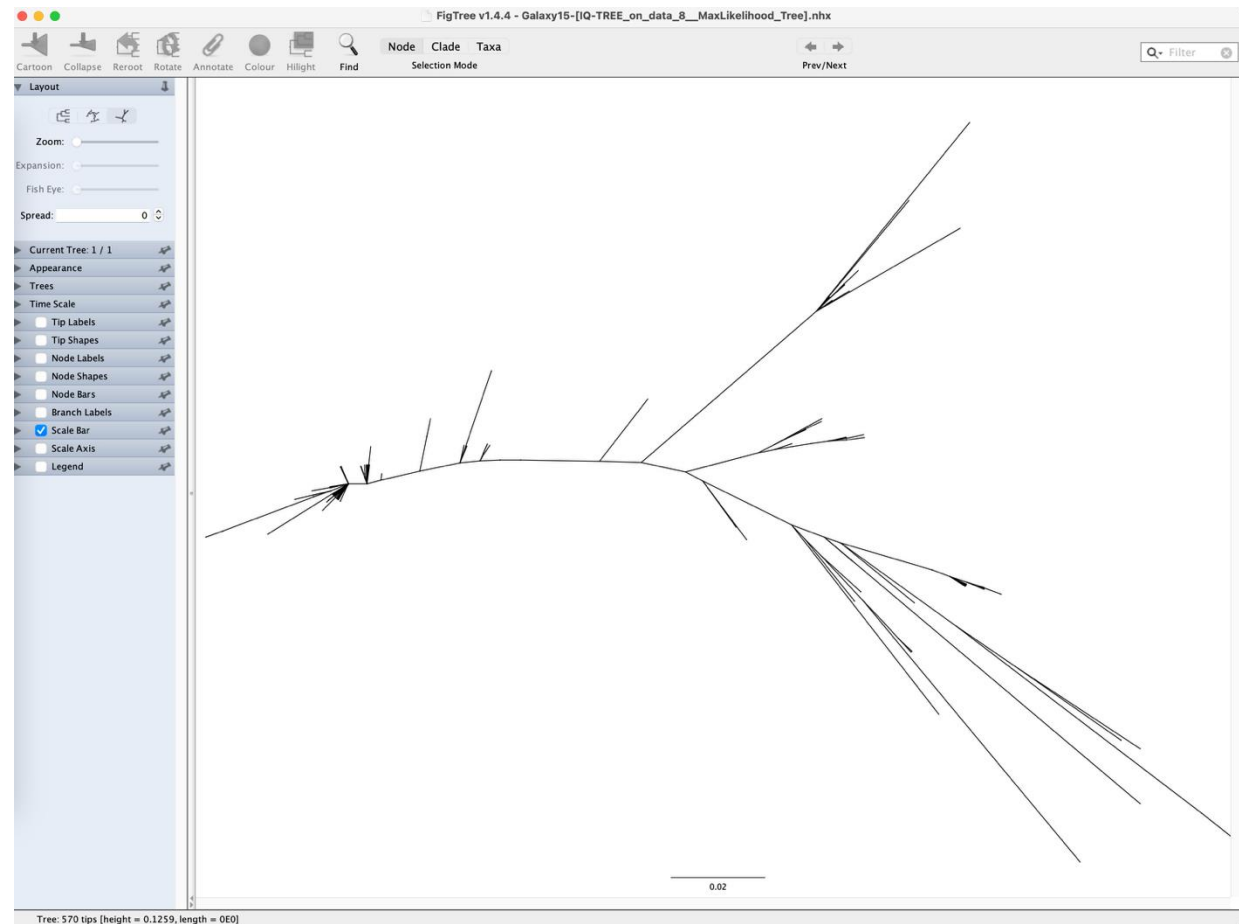  - (built-in for IQ-tree and RAxML)
- Maximum likelihood methods are the baseline standard
  - IQ-tree, RAxML, etc.
- Bayesian methods are statistically more robust, but are computationally very intensive
  - RevBayes, BEAST
- Bootstrapping gives an estimate of the robustness of your tree to changes in the input data

# Next Week's Group Meetings
Tuesday 20$^{th}$ October 13:30 HW324
Thursday 23$^{rd}$ October 10:30 HW324

# Topics to Discuss at Next Meeting

- What would you like to cover?

# Useful Links

# Useful tools (many others are available)

GalaxyEU: https://usegalaxy.eu/

- Sequence alignment (e.g. MAFFT), phylogenetics (e.g. RaxML), positive selection (e.g. codeML)

iTOL: https://itol.embl.de/

- Visualisation/annotation of phylogenetic trees

PyMOL: https://pymol.org/2/ and/or ChimeraX: https://www.cgl.ucsf.edu/chimerax/

- Protein structure visualisation/annotation

Jalview: http://www.jalview.org/

- Visualisation of multiple sequence alignments

# Useful sites/databases

PHI-base: http://www.phi-base.org/

- Proteins involved in host-pathogen interactions, with linked evidence

EMBL AlphaFold: https://www.alphafold.ebi.ac.uk/

- AlphaFold predictions for proteins from model organisms
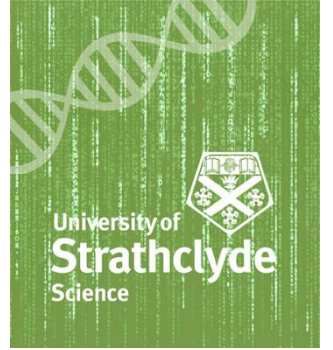
UniProt: https://www.uniprot.org/

- Protein sequence (including homologous sequences) and functional information with evidence

RCSB/PDB: https://www.rcsb.org/

- Repository of record for protein structures
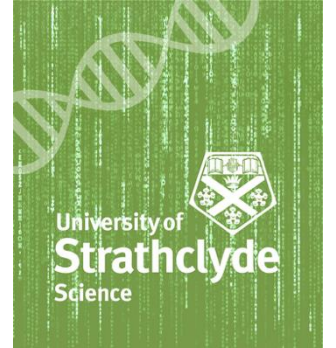
# SIPBS CompBiol Sites

- BM432 Project Pages
  - https://sipbs-compbiol.github.io/bm432-project/

- An incomplete little book of bioinformatics
  - https://sipbs-compbiol.github.io/little-bioinformatics-book/

# Project Management Tools

# You may want tools to…

- Manage your time
  - E.g. Pomodoro technique (e.g. BeFocused, Pomofocus, Forest)
- Schedule work
  - Reminders (macOS, MS Office)
  - Calendar (macOS, MS Office), with email alerts
  - Trello, Asana, etc.
- Manage your project data and information effectively
  - How to name files
  - Project management guidelines (BM432, 2022-23 session; me and Dr Feeney)
  - How to keep a lab notebook
  - Keeping a computational biology lab notebook: https://doi.org/10.1371/journal.pcbi.1004385
  - Organising a lab book