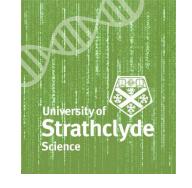


Evolutionary and Structural Analysis of Pathogen Proteins.

Final year UG project 2024-25 2024-10-27 (Week 6)





Any changes needed?

| Organism | Host | Gene/Protein | PHI accession | Student |
|-------------------|--------------|--------------|---------------|---------|
| Escherichia coli | Homo sapiens | espY | PHI:8647 | LB |
| Shigella flexneri | Homo sapiens | іраЈ | 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)

Visualise with **PyMOL**

Download

AlphaFold/PDB

structures (try

simplefold?)

homologues

2. All bacteria

Identify

1. Source

species

(compare AlphaFold with PDB?)

Align sequences Conserved and variable sites?

Phylogenetic trees

HGT? Positive selection?

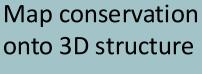
Interpret in context of known function/species distribution

Other database searches

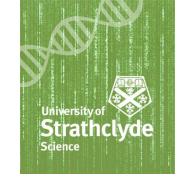
Known interactors?

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

(Weeks ≈6-11)

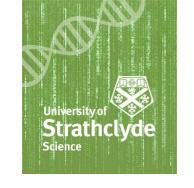




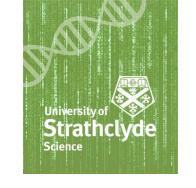


Your questions/comments

(What would you like to talk about?)

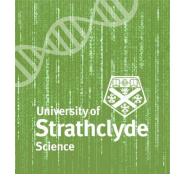


Tell everyone about your protein



Identifying and aligning homologues

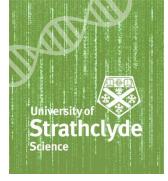
Useful steps - 1



- Identify UniProt record from PHI-Base
- UniProt has homologues: 100% (15), 90% (102), 50% (107) identity almost all *Mycobacterium* spp.
 - (what does this tell us?)
- (Similar Proteins -> View all entries -> Download)

Useful steps - 2

- BLASTP against
 ClusteredNR at NCBI gives
 over 100 clusters (1-71
 members per cluster)
 with ≈100% coverage,
 >60% identity
 - 97 are Mycobacteriales
 - 32 are *Mycobacterium*
- Also a rough phylogenetic tree
 - (what does this tell us?)
- (Alignments -> Download)







University of Strathclyde Science

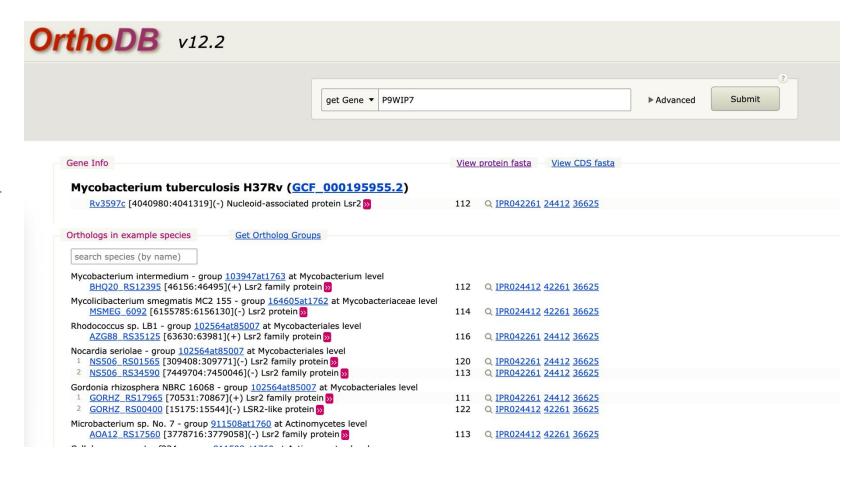
- Protein has a known domain homologous domains in Pfam (8k)
- Various domain architectures, many the same as my query (7576)
 - (What does this tell us?)
- (Click on link and Download)



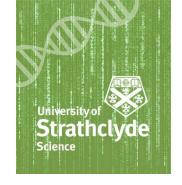


University of Strathclyde Science

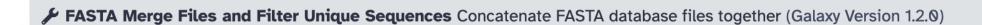
- Homologue databases
 - OrthoDB, InParanoid, EggNOG, etc.
- OrthoDB (5985)
- (Get ortholog groups -> select group View protein FASTA)







- Concatenate sequence files
 - Galaxy, cat, JalView
- Can also filter identical sequences











Tool Parameters

Run in batch mode?

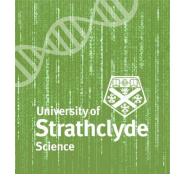
- Merge individual FASTAs (output collection if input is collection)
- Merge all FASTAs (always output a single FASTA)

The 'merge all' mode produces one output FASTA for all input FASTA files. The individual mode generates one FASTA file for each set of input FASTAs. For example, if the tool is given 2 collections of 10 FASTAs, it will merge the collections pairwise to create an output collection of 10 FASTAs.

Input FASTA File(s)

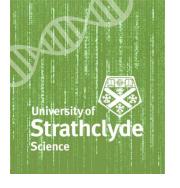
+ Insert Input FASTA File(s)

Useful steps - 5

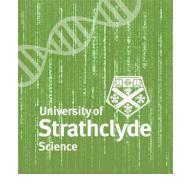


- Align sequences
 - MAFFT, Clustal Omega, MUSCLE, MMSeqs2
- Use standard scoring matrix or incorporate sequence composition?
- Does the alignment look reasonable?
 - If not, what to do?
 - Exclude distant sequences?
 - Align clusters?
 - Can it be "rescued" with trimAl?
 - Would using a structure or other pre-alignment to "anchor" the alignment help?

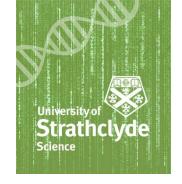
Topics/Progress to Discuss at Next Meeting



- What would you like to cover?
- Keep investigating your protein
- Make a tree for your sequences
 - o How does the tree look?
 - Are organisms/species grouped together how you'd expect? Any surprises?
 - \circ Is there evidence supporting gene duplication (two clades, with the same organisms repeated)?
 - Is there evidence supporting horizontal gene transfer (organism in the "wrong" clade)?
- Compare to sequence data
 - Do patterns of sequence variation map onto your tree?
- Galaxy workflow or scripting?
 - Have you created a workflow for your analysis: alignment? making a tree?
 - Have you noted/saved the tool versions/parameter settings for writing your thesis?

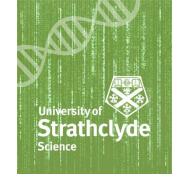


Next Week's Group Meeting Monday 3rd November 13:30 HW324



Useful Links





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

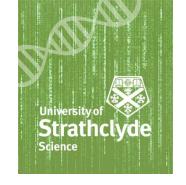
PyMOL: https://pymol.org/2/ and/or ChimeraX: https://www.cgl.ucsf.edu/chimerax/

Protein structure visualisation/annotation

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

- Visualisation of multiple sequence alignments





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

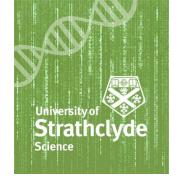
UniProt: https://www.uniprot.org/

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

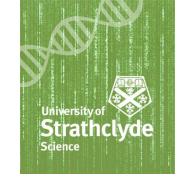
RCSB/PDB: https://www.rcsb.org/

- Repository of record for protein structures

SIPBS CompBiol Sites

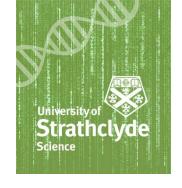


- BM432 Project Pages
 - https://sipbs-compbiol.github.io/bm432-project/
- An incomplete little book of bioinformatics
 - https://sipbs-compbiol.github.io/little-bioinformatics-book/



Project Management Tools

You may want tools to...



- Manage your time
 - E.g. Pomodoro technique (e.g. BeFocused, <u>Pomofocus</u>, <u>Forest</u>)
- 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