

# Evolutionary and Structural Analysis of Pathogen Proteins.

Final year UG project 2024-25  
2024-10-16 (Week 4)

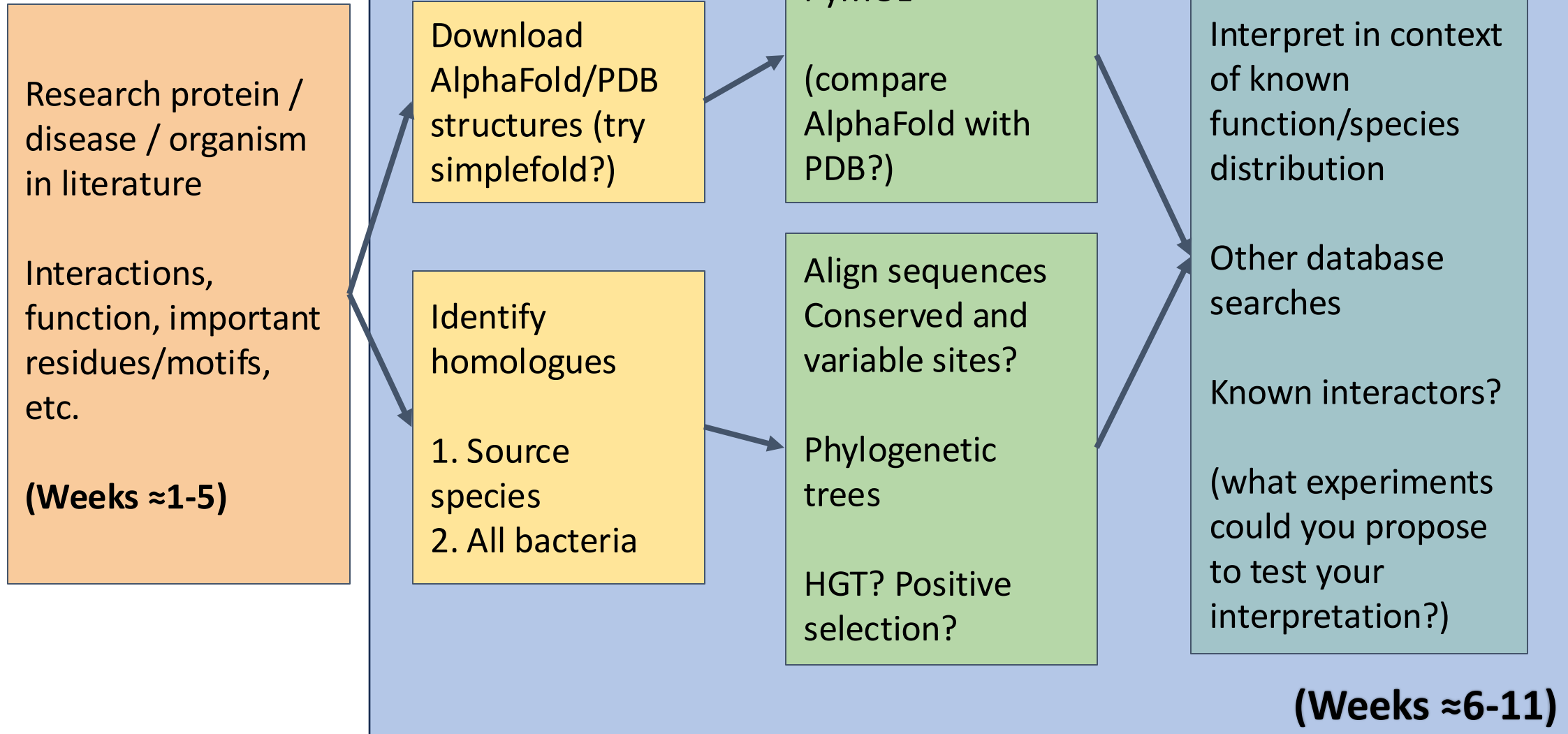
# Candidate proteins – start points

**Any changes needed?**

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

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

# Workflow



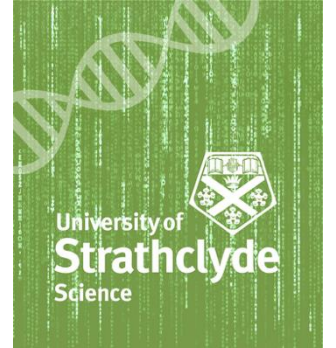
# Your questions/comments

(What would you like to talk about?)

# Thesis Introductions

# Examples

- Lsr2 (Ailsa)
- SpA (Yann)
- Note:
  - Use of figures
  - Density of references
  - Overall depth of description



# Sequence Alignment

# Pairwise sequence alignment

- We have two words we want to align:
  - COELACANTH
  - PELICAN
- We have a scoring scheme:
  - Matching identical letters scores **+1**
  - If we have different letters aligned, that scores **-1**
  - Inserting a gap in the alignment scores **-1**
- All sequence alignment is a mathematical operation: *the maximization of an alignment score*



# Initialise a matrix

		C	O	E	L	A	C	A	N	T	H
P E L I C A N	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7	← -8	← -9	← -10
	↑ -1										
	↑ -2										
	↑ -3										
	↑ -4										
	↑ -5										
	↑ -6										
	↑ -7										

# Start filling the matrix

	C	O	E	L	A	C	A	N	T	H	
P	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7	← -8	← -9	← -10
	↑ -1	↖ -1	↖ -2								

CO  
-P

CO  
P-

# A full matrix

		C	O	E	L	A	C	A	N	T	H
	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7	← -8	← -9	← -10
P	↑ -1	↖ -1	↖ -2	↖ -3	↖ -4	↖ -5	↖ -6	↖ -7	↖ -8	↖ -9	↖ -10
E	↑ -2	↖ -2	↖ -2	↖ -1	← -0	← -3	← -4	← -5	← -6	← -7	← -8
L	↑ -3	↖ -3	↖ -3	← -2	↖ -2	← -1	← -2	← -3	← -4	← -5	← -6
I	↑ -4	↖ -4	↑ -4	↑ -3	↑ -1	↖ -1	↖ -2	↖ -1	↖ -4	↖ -5	↖ -6
C	↑ -5	↖ -3	← -4	↑ -4	↑ -2	↖ -2	↖ -0	← -1	← -2	← -3	← -4
A	↑ -6	↑ -4	↖ -4	↖ -5	↑ -3	↖ -1	↑ -1	↖ -1	← -0	← -1	← -2
N	↑ -7	↑ -5	↖ -5	↖ -5	↑ -4	↑ -2	↖ -2	← -0	↖ -2	← -1	← -0

# Traceback

		C	O	E	L	A	C	A	N	T	H	
		0	← -1	← -2	← -3	← -4	← -5	← -6	← -7	← -8	← -9	← -10
P	↑ -1	↖ -1	↖ -2	↖ -3	↖ -4	↖ -5	↖ -6	↖ -7	↖ -8	↖ -9	↖ -10	
E	↑ -2	↖ -2	↖ -2	↖ -1	← -0	← -3	← -4	← -5	← -6	← -7	← -8	
L	↑ -3	↖ -3	↖ -3	← -2	↖ -2	← -1	← -2	← -3	← -4	← -5	← -6	
I	↑ -4	↖ -4	↑ -4	↑ -3	↑ -1	↖ -1	↖ -2	↖ -1	↖ -4	↖ -5	↖ -6	
C	↑ -5	↖ -3	← -4	↑ -4	↑ -2	↖ -2	↖ -0	← -1	← -2	← -3	← -4	
A	↑ -6	↑ -4	↖ -4	↖ -5	↑ -3	↖ -1	↑ -1	↖ -1	← -0	← -1	← -2	
N	↑ -7	↑ -5	↖ -5	↖ -5	↑ -4	↑ -2	↖ -2	← -0	↖ -2	← -1	← -0	

COELACANTH  
- PELICAN -

# Where was the biology in that?

Honest question. Where do you think it was?

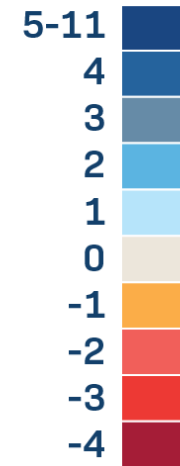
# Pairwise sequence alignment

- We have two words we want to align:
  - COELACANTH
  - PELICAN
- **We have a scoring scheme:**
  - Matching identical letters scores **+1**
  - If we have different letters aligned, that scores **-1**
  - Inserting a gap in the alignment scores **-1**
- All sequence alignment is a mathematical operation: *the maximization of an alignment score*

**This was all the biology:  
What score do we give the same/  
similar symbols in the alignment?  
Represents evolutionary pressure  
on the sequence**

# BLOSUM62 Substitution Matrix

j																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																					</
---	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	----



polar

nonpolar

aromatic

BLOSUM62: a common default scoring scheme for protein sequences

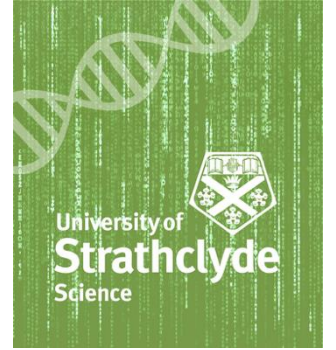
## Next Week's Group Meetings

Tuesday 20<sup>th</sup> October 13:30 HW324

Thursday 23<sup>rd</sup> October 10:30 HW324



# Topics to Discuss at Next Meeting



- Tell me about your protein
  - What scientific questions do you think you could address in your research? (This could be inspired by something you've read in a paper describing your protein – or something more general)
- Phylogenetic Trees?
- 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

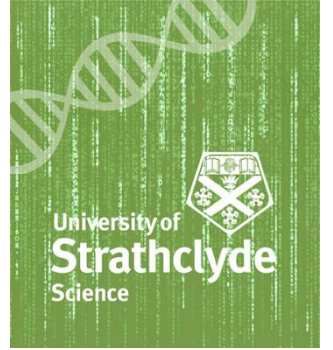
UniProt: <https://www.uniprot.org/>

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

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

- Repository of record for protein structures

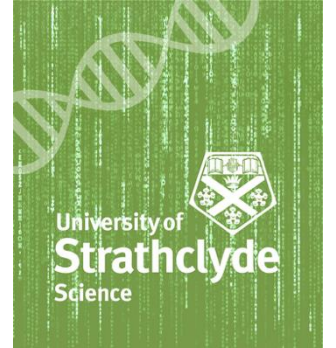
# SIPBS CompBiol Sites



- 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](#)