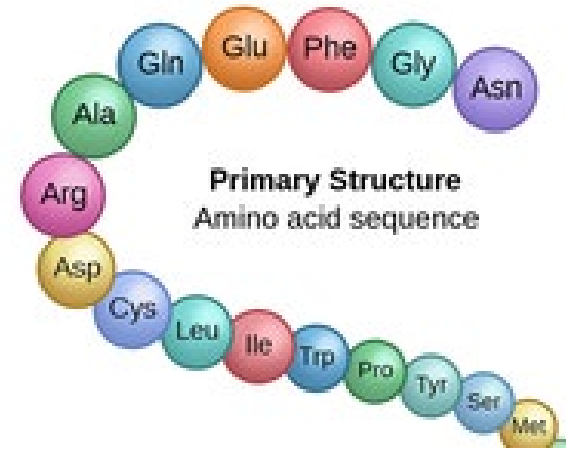# Bioinformatics for Protein Expression
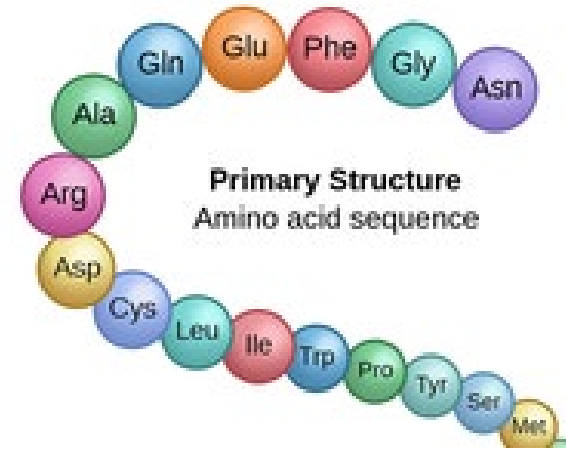
Day 2: Tuesday 21st March

# Today's talk

- How much can we learn and predict from sequence alone?

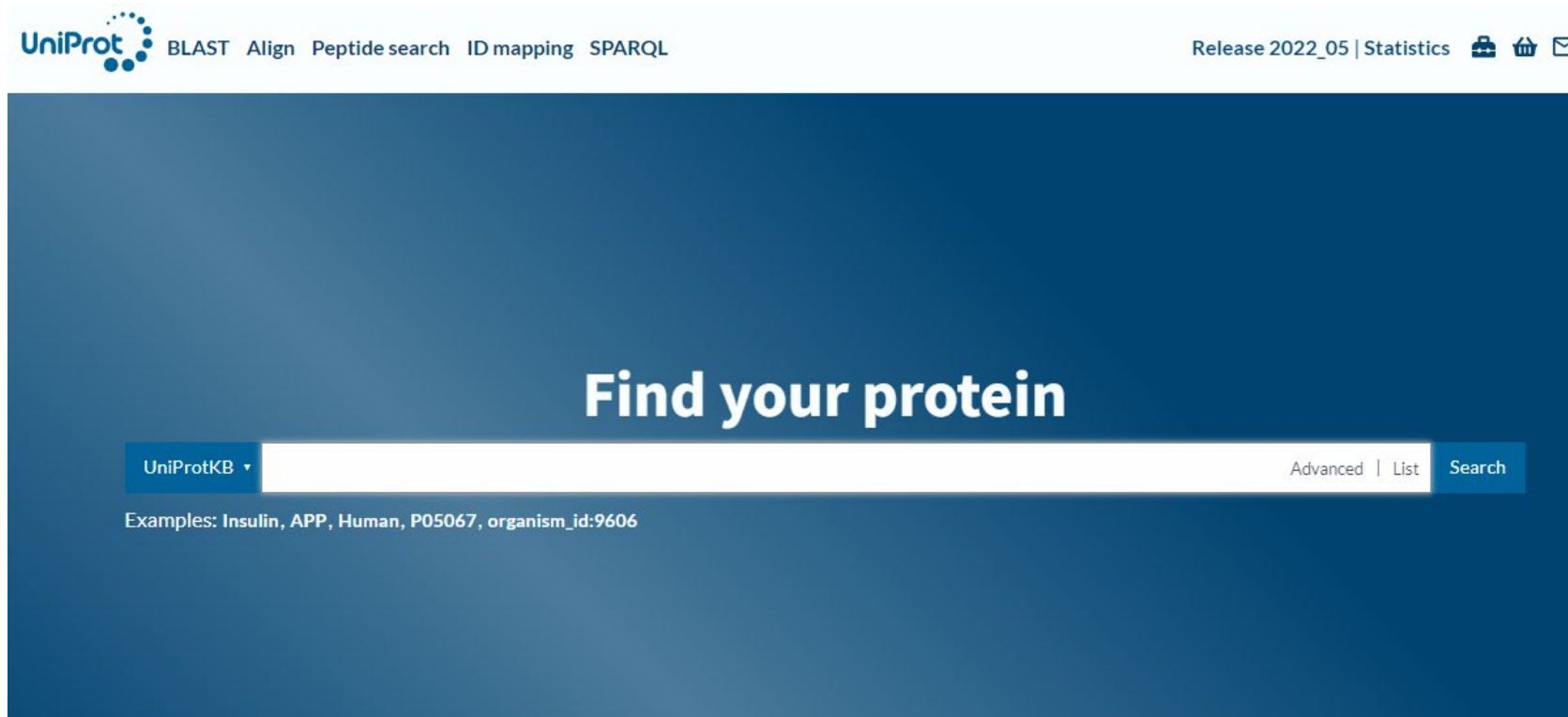- How does this help us with our purifications?

# How can bioinformatics help me?

- In this lecture you will learn how (using just your protein sequence) you can:

  - Predict domains
  - Identify post translational modifications
  - Calculate the molecular mass
  - Determine the isoelectric point
  - And most importantly, learn how to get your "extinction co-efficients" so you can calculate protein concentration accurately
  - Identify distant homologs



Primary Structure
Amino acid sequence

# First useful resource - Uniprot

- Uniprot.org

# Uniprot – brief functional description

# Uniprot – subcellular localisation

# Uniprot – post-translational modifications

# Uniprot – Domain composition

# Uniprot – sequence data

# Uniprot – sequence data

# Uniprot – sequence formats

- FASTA format



```
>sp|O94856|NFASC_HUMAN Neurofascin OS=Homo sapiens OX=9606 GN=NFASC PE=1 SV=4
MARQPPPPWVHAAFLLCLLSLGGAIEIPMDPSIQNELTQPPTITKQSAKDHIVDPRDNIL
IECEAKGNPAPSFHWTRNSRFFNIAKDPRVSMRRRSGTLVIDFRSGGRPEEYEGEYQCFA
RNKFGTALSNRIRLQVSKSPLWPKENLDPVVVQEGAPLTLQCNPPPGLPSPVIFWMSSSM
EPITQDKRVSQGHNGDLYFSNVMLQDMQTDYSCNARFHFTHTIQQKNPFTLKVLTTRGVA
ERTPSFMYPQGTASSQMVLRGMDLLLECIASGVPTPDIAWYKKGGDLPSDKAKFENFNKA
LRITNVSEEDSGEYFCLASNKMGSIRHTISVRVKAAPYWLDEPKNLILAPGEDGRLVCRA
NGNPKPTVQWMVNGEPLQSAPPNPNREVAGDTIIFRDTQISSRAVYQCNTSNEHGYLLAN
AFVSVLDVPPRMLSPRNQLIRVILYNRTRLDCPFFGSPIPTLRWFKNGQGSNLDGGNYHV
YENGSLEIKMIRKEDQGIYTCVATNILGKAENQVRLEVKDPTRIYRMPEDQVARRGTTVQ
LECRVKHDPSLKLTVSWLKDDEPLYIGNRMKKEDDSLTIFGVAERDQGSYTCVASTELDQ
DLAKAYLTVLADQATPTNRLAALPKGRPDRPRDLELTDLAERSVRLTWIPGDANNSPITD
YVVQFEEDQFQPGVWHDHSKYPGSVNSAVLRLSPYVNYQFRVIAINEVGSSHPSLPSERY
RTSGAPPESNPGDVKGEGTRKNNMEITWTPMNATSAFGPNLRYIVKWRRRETREAWNNVT
VWGSRYVVGQTPVYVPYEIRVQAENDFGKGPEPESVIGYSGEDYPRAAPTEVKVRVMNST
AISLQWNRVYSDTVQGQLREYRAYYWRESSLLKNLWVSQKRQQASFPGDRLRGVVSRLFP
YSNYKLEMVVVNGRGDGPRSETKEFTTPEGVPSAPRRFRVRQPNLETINLEWDHPEHPNG
IMIGYTLKYVAFNGTKVGKQIVENFSPNQTKFTVQRTDPVSRYRFTLSARTQVGSGEAVT
EESPAPPNEATPTAAPPTLPPTTVGATGAVSSTDATAIAATTEATTVPIIPTVAPTTIAT
TTTVATTTTTAAATTTTESPPTTTSGTKIHESAPDEQSIWNVTVLPNSKWANITWKHNF
GPGTDFVVEYIDSNHTKKTVPVKAQAQPIQLTDLYPGMTYTLRVYSRDNEGISSTVITFM
TSTAYTNNQADIATQGWFIGLMCAIALLVLILLIVCFIKRSRGGKYPVREKKDVPLGPED
PKEEDGSFDYSDEDNKPLQGSQTSLDGTIKQQESDDSLVDYGEGGEGQFNEDGSFIGQYT
VKKDKEETEGNESSEATSPVNAIYSLA
```

# Uniprot – editing the sequence

- But, what if I am only making one domain of this?

```
>sp|O94856|NFASC_HUMAN Neurofascin OS=Homo sapiens OX=9606 GN=NFASC PE=1 SV=4
MARQPPPPWVHAAFLLCLLSLGGAIEIPMDPSIQNELTQPPTITKQSAKDHIVDPRDNIL
IECEAKGNPAPSFHWTRNSRFFNIAKDPRVSMRRRSGTLVIDFRSGGRPEEYEGEYQCFA
RNKFGTALSNRIRLQVSKSPLWPKENLDPVVVQEGAPLTLQCNPPPGLPSPVIFWMSSSM
EPITQDKRVSQGHNGDLYFSNVMLQDMQTDYSCNARFHFTHTIQQKNPFTLKVLTTRGVA
ERTPSFMYPQGTASSQMVLRGMDLLLECIASGVPTPDIAWYKKGGDLPSDKAKFENFNKA
LRITNVSEEDSGEYFCLASNKMGSIRHTISVRVKAAPYWLDEPKNLILAPGEDGRLVCRA
NGNPKPTVQWMVNGEPLQSAPPNPNREVAGDTIIFRDTQISSRAVYQCNTSNEHGYLLAN
AFVSVLDVPPRMLSPRNQLIRVILYNRTRLDCPFFGSPIPTLRWFKNGQGSNLDGGNYHV
YENGSLEIKMIRKEDQGIYTCVATNILGKAENQVRLEVKDPTRIYRMPEDQVARRGTTVQ
LECRVKHDPSLKLTVSWLKDDEPLYIGNRMKKEDDSLTIFGVAERDQGSYTCVASTELDQ
DLAKAYLTVLADQATPTNRLAALPKGRPDRPRDLELTDLAERSVRLTWIPGDANNSPITD
YVVQFEEDQFQPGVWHDHSKYPGSVNSAVLRLSPYVNYQFRVIAINEVGSSHPSLPSERY
RTSGAPPESNPGDVKGEGTRKNNMEITWTPMNATSAFGPNLRYIVKWRRRETREAWNNVT
VWGSRYVVGQTPVYVPYEIRVQAENDFGKGPEPESVIGYSGEDYPRAAPTEVKVRVMNST
AISLQWNRVYSDTVQGQLREYRAYYWRESSLLKNLWVSQKRQQASFPGDRLRGVVSRLFP
YSNYKLEMVVVNGRGDGPRSETKEFTTPEGVPSAPRRFRVRQPNLETINLEWDHPEHPNG
IMIGYTLKYVAFNGTKVGKQIVENFSPNQTKFTVQRTDPVSRYRFTLSARTQVGSGEAVT
EESPAPPNEATPTAAPPTLPPTTVGATGAVSSTDATAIAATTEATTVPIIPTVAPTTIAT
TTTVATTTTTAAATTTTESPPTTTSGTKIHESAPDEQSIWNVTVLPNSKWANITWKHNF
GPGTDFVVEYIDSNHTKKTVPVKAQAQPIQLTDLYPGMTYTLRVYSRDNEGISSTVITFM
TSTAYTNNQADIATQGWFIGLMCAIALLVLILLIVCFIKRSRGGKYPVREKKDVPLGPED
PKEEDGSFDYSDEDNKPLQGSQTSLDGTIKQQESDDSLVDYGEGGEGQFNEDGSFIGQYT
VKKDKEETEGNESSEATSPVNAIYSLA
```

# Uniprot

- You can manually edit this to contain your sequence of interest
- This can be done in any editing software, eg. NotePad

NFASC_truncation.fasta - Notepad

File   Edit   Format   View   Help

>NFASC truncation
IECEAKGNPAPSFHWTRNSRFFNIAKDPRVSMRRRSGTLVIDFRSGGRPEEYEGEYQCFA
RNKFGTALSNRIRLQVSKSPLWPKENLDPVVVQEGAPLTLQCNPPPGLPSPVIFWMSSSM
EPITQDKRVSQGHNGDLYFSNVMLQDMQTDYSCNARFHFTHTIQQKNPFTLKVLTTRGVA
ERTPSFMYPQGTASSQMVLRGMDLLLECIASGVPTPDIAWYKKGGDLPSDKAKFENFNKA
LRITNVSEEDSGEYFCLASNKMGSIRHTISVRVKAAPYWLDEPKNLILAPGEDGRLVCRA
NGNPKPTVQWMVNGEPLQSAPPNPNREVAGDTIIFRDTQISSRAVYQCNTSNEHGYLLAN
AFVSVLDVPPRMLSPRNQLIRVILYNRTRLDCPFFGSPIPTLRWFKNGQGSNLDGGNYHV
YENGSLEIKMIRKEDQGIYTCVATNILGKAENQVRLEVKDPTRIYRMPEDQVARRGTTVQ
LECRVKHDPSLKLTVSWLKDDEPLYIGNRMKKEDDSLTIFGVAERDQGSYTCVASTELDQ
DLAKAYLTVLADQATPTNRLAALPKGRPDRPRDLELTDLAERSVRLTWIPGDANNSPITD
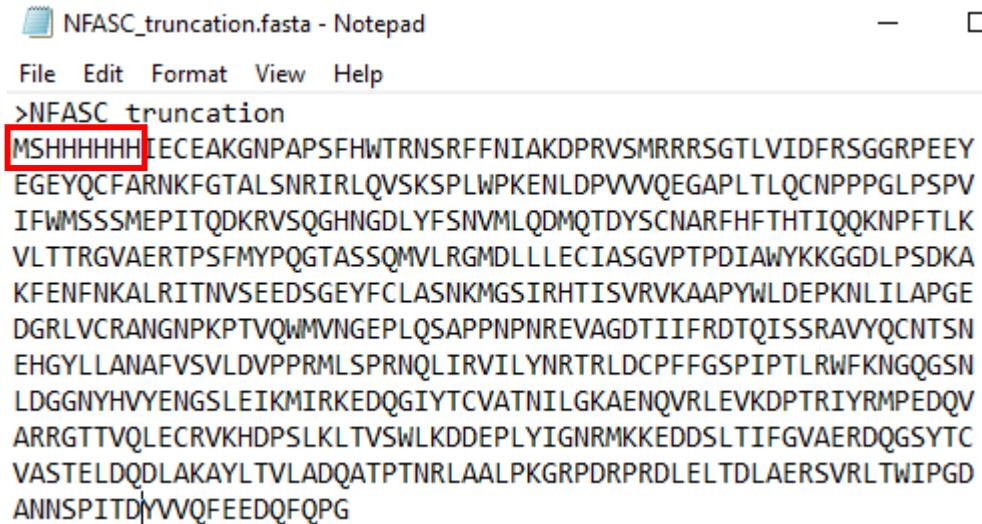YVVQFEEDQFQPG

UNIVERSITY OF CAMBRIDGE

# Uniprot

- You can manually edit this to contain your sequence of interest
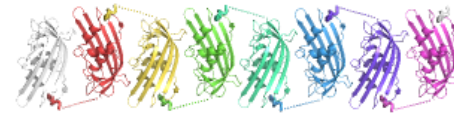- This can be done in any editing software, eg. NotePad



- You can also add the sequence of tags such as His or GST-tags

# Or you can get the sequence from a plasmid map

- For this course, there is a link to the software and plasmids files we're using: http://www.atomicvirology.path.cam.ac.uk/brazil



**Atomic Virology Lab**
**University of Cambridge**

**Theoretical and practical course in protein biochemistry, biophysics and structural biology**

*Departamento de Biologia Celular e Molecular e Bioagentes Patogênicos*
*Faculdade de Medicina de Ribeirão Preto, 20 to 31 March 2023*

**Useful software links**

ProtParam: https://web.expasy.org/protparam/
ApE plasmid editor: https://jorgensen.biology.utah.edu/wayned/ape/
UniProt knowledgebase: https://www.uniprot.org/
NCBI BLAST: https://blast.ncbi.nlm.nih.gov/Blast.cgi
AlphaFold: https://github.com/deepmind/alphafold
ColabFold: https://github.com/sokrypton/ColabFold

**Vector maps**

$His_6$-EGFP   $His_6$-mTurquoise2   $His_6$-mVenus   $His_6$-mCherry   antiGFPnanobody-GST   GST-3Cprotease

**UNIVERSITY OF CAMBRIDGE**

# Or you can get the sequence from a plasmid map

- Use ApE to open the plasmid file and extract the sequence data

**Atomic Virology Lab**
University of Cambridge

**Theoretical and practical course in protein biochemistry, biophysics and structural biology**

*Departamento de Biologia Celular e Molecular e Bioagentes Patogênicos*
*Faculdade de Medicina de Ribeirão Preto, 20 to 31 March 2023*

**Useful software links**

ProtParam: https://web.expasy.org/protparam/
ApE plasmid editor: https://jorgensen.biology.utah.edu/wayned/ape/
UniProt knowledgebase: https://www.uniprot.org/
NCBI BLAST: https://blast.ncbi.nlm.nih.gov/Blast.cgi
AlphaFold: https://github.com/deepmind/alphafold
ColabFold: https://github.com/sokrypton/ColabFold

**Vector maps**

His$_6$-EGFP   His$_6$-mTurquoise2   His$_6$-mVenus   His$_6$-mCherry   antiGFPnanobody-GST   GST-3Cprotease

UNIVERSITY OF CAMBRIDGE

http://www.atomicvirology.path.cam.ac.uk/brazil

# ApE – A plasmid Editor

- Open your vector map – different features are highlighted

# ApE – A plasmid Editor

- You can also view this as a graphic map

# ApE – A plasmid Editor

- Identify the relevant open reading frame (ORF)

# ApE – A plasmid Editor

- Translate the ORF

# ApE – A plasmid Editor

- Translate as 1-letter code

# ApE – A plasmid Editor

- Now you have the exact sequence of your construct including tags etc

# Second useful resource - ProtParam

- Protparam uses a sequence to calculate:

  - molecular mass
  - isoelectric point
  - extinction coefficients

    https://web.expasy.org/protparam/

# Second useful resource  - ProtParam

# Second useful resource - ProtParam

```
Number of amino acids: 510

Molecular weight: 57364.86

Theoretical pI: 6.37

Amino acid composition:
Ala (A)   31      6.1%
Arg (R)   34      6.7%
Asn (N)   30      5.9%
Asp (D)   35      6.9%
Cys (C)    9      1.8%
Gln (Q)   24      4.7%
Glu (E)   27      5.3%
Gly (G)   32      6.3%
His (H)    9      1.8%
Ile (I)   23      4.5%
Leu (L)   45      8.8%
Lys (K)   25      4.9%
Met (M)   11      2.2%
Phe (F)   16      3.1%
Pro (P)   34      6.7%
Ser (S)   32      6.3%
Thr (T)   33      6.5%
Trp (W)    7      1.4%
Tyr (Y)   18      3.5%
Val (V)   35      6.9%
```

# Second useful resource - ProtParam

```
Number of amino acids: 510

Molecular weight: 57364.86          ←  Mass in daltons, 57.4 kDa

Theoretical pI: 6.37

Amino acid composition:
Ala (A)   31      6.1%
Arg (R)   34      6.7%
Asn (N)   30      5.9%
Asp (D)   35      6.9%
Cys (C)    9      1.8%
Gln (Q)   24      4.7%
Glu (E)   27      5.3%
Gly (G)   32      6.3%
His (H)    9      1.8%
Ile (I)   23      4.5%
Leu (L)   45      8.8%
Lys (K)   25      4.9%
Met (M)   11      2.2%
Phe (F)   16      3.1%
Pro (P)   34      6.7%
Ser (S)   32      6.3%
Thr (T)   33      6.5%
Trp (W)    7      1.4%
Tyr (Y)   18      3.5%
Val (V)   35      6.9%
```

UNIVERSITY OF
CAMBRIDGE

# Second useful resource  - ProtParam

```
Number of amino acids: 510

Molecular weight: 57364.86

Theoretical pI: 6.37    ⬅

Amino acid composition:
Ala (A)   31        6.1%
Arg (R)   34        6.7%
Asn (N)   30        5.9%
Asp (D)   35        6.9%
Cys (C)    9        1.8%
Gln (Q)   24        4.7%
Glu (E)   27        5.3%
Gly (G)   32        6.3%
His (H)    9        1.8%
Ile (I)   23        4.5%
Leu (L)   45        8.8%
Lys (K)   25        4.9%
Met (M)   11        2.2%
Phe (F)   16        3.1%
Pro (P)   34        6.7%
Ser (S)   32        6.3%
Thr (T)   33        6.5%
Trp (W)    7        1.4%
Tyr (Y)   18        3.5%
Val (V)   35        6.9%
```

**pI is the pH where your protein has no charge**

**Your protein can be unstable at its pI**

**Keep your purification buffers at least 1 pH unit from the pI**

# Extinction coefficients for protein concentration

```
Extinction coefficients:

Extinction coefficients are in units of  M⁻¹ cm⁻¹, at 280 nm measured in water.

Ext. coefficient     65820
Abs 0.1% (=1 g/l)    1.147, assuming all pairs of Cys residues form cystines


Ext. coefficient     65320
Abs 0.1% (=1 g/l)    1.139, assuming all Cys residues are reduced
```

# Extinction coefficients for protein concentration

Extinction coefficients:

Extinction coefficients are in units of $M^{-1} cm^{-1}$, at 280 nm measured in water.

Ext. coefficient     65820
Abs 0.1% (=1 g/l)   1.147, assuming all pairs of Cys residues form cystines

Ext. coefficient     65320
Abs 0.1% (=1 g/l)   1.139, assuming all Cys residues are reduced

$E_{0.1\%}$ in units of $(mg/mL)^{-1}$

UNIVERSITY OF CAMBRIDGE

# What determines the E value?

Ext. coefficient         65320
Abs 0.1% (=1 g/l)        1.139

Amino acid composition:
Ala (A)   31        6.1%
Arg (R)   34        6.7%
Asn (N)   30        5.9%
Asp (D)   35        6.9%
Cys (C)    9        1.8%
Gln (Q)   24        4.7%
Glu (E)   27        5.3%
Gly (G)   32        6.3%
His (H)    9        1.8%
Ile (I)   23        4.5%
Leu (L)   45        8.8%
Lys (K)   25        4.9%
Met (M)   11        2.2%
Phe (F)   16        3.1%
Pro (P)   34        6.7%
Ser (S)   32        6.3%
Thr (T)   33        6.5%
Trp (W)    7        1.4%
Tyr (Y)   18        3.5%
Val (V)   35        6.9%

Phe

Trp

Tyr

# What determines the E value?

Ext. coefficient      65320
Abs 0.1% (=1 g/l)     1.139

```
Amino acid composition:
Ala (A)  31      6.1%
Arg (R)  34      6.7%
Asn (N)  30      5.9%
Asp (D)  35      6.9%
Cys (C)   9      1.8%
Gln (Q)  24      4.7%
Glu (E)  27      5.3%
Gly (G)  32      6.3%
His (H)   9      1.8%
Ile (I)  23      4.5%
Leu (L)  45      8.8%
Lys (K)  25      4.9%
Met (M)  11      2.2%
Phe (F)  16      3.1%
Pro (P)  34      6.7%
Ser (S)  32      6.3%
Thr (T)  33      6.5%
Trp (W)   7      1.4%
Tyr (Y)  18      3.5%
Val (V)  35      6.9%
```

Phe

Trp

Tyr

# Calculating protein concentration

- Spectrophotometry, measure absorbance at 280 nm

$$A = log_{10}\left(\frac{I_0}{I}\right)$$

# Calculating protein concentration

- Beer-Lambert Law

$$A = log_{10}\left(\frac{I_0}{I}\right) = \varepsilon cl$$

$\varepsilon$ = *extinction co-efficient*

$c$ = *concentration*

$l$ = *path length*

# Calculating protein concnetration

- Beer-Lambert Law and molecular extinction coefficients

$A = \varepsilon c l$

Measure the absorbance

Use the theoretical extinction co-efficient

$c = \dfrac{A}{\varepsilon l}$

Calculate concentration

# Calculating protein concentration

- Beer-Lambert Law and molecular extinction coefficents

Calculate concentration in M or mg/mL

$A = 0.5$

$$c = \frac{A}{\varepsilon l}$$

```
Ext. coefficient      65320
Abs 0.1% (=1 g/l)     1.139
```

# Calculating protein concentration

- Beer-Lambert Law and molecular extinction coefficients

A = 0.5

$$c = \frac{A}{\varepsilon l}$$

Ext. coefficient    65320
Abs 0.1% (=1 g/l)   1.139

Calculate concentration in M or mg/mL

Molar conc,    65320 (M)
               = 0.5/65320 = 7.65 x 10$^{-6}$
                            = 7.65 µM

In mg/mL,      1.139
               = 0.5/1.139 = 0.44 mg/mL

UNIVERSITY OF
CAMBRIDGE

# This correction is super important

- Always correct your absorbance measurements using your extinction coefficient

- Some instruments will quote concentration from absorbance but if you haven't entered an extinction coefficient it will be wrong!

# Sequence alignments

- Pairwise sequence alignments
  - mouse vs human
  - Proteins from different viruses

https://www.ebi.ac.uk/Tools/psa/emboss_needle/



UNIVERSITY OF CAMBRIDGE

# Pairwise alignment

- Identity calculation

- Use similarity/homology carefully, calculated differently by programs

- Useful for comparing animal models of human diseases

```
 ‾
# Length: 685
# Identity:      566/685 (82.6%)   ←
# Similarity:    615/685 (89.8%)
# Gaps:          1/685  ( 0.1%)
# Score: 3091.0
#
#
#=======================================

Human       1 MAEWLLSASWQRRAKAMTAAAGSAGRAAVPLLLCALLAPGGAYVLDDSDG     50
              ||.....||.||:||-||||||||.|-||||||||||-||||||||||||
Mouse       1 MANSQPKASQQRQAKVMTAAAGSASRVAVPLLLCALLVPGGAYVLDDSDG     50

Human      51 LGREFDGIGAVSGGGATSRLLVNYPEPYRSQILDYLFKPNFGASLHILKV    100
              |||||||||||||||||||||||||||||:||||||||||||||||||||
Mouse      51 LGREFDGIGAVSGGGATSRLLVNYPEPYRSEILDYLFKPNFGASLHILKV    100

Human     101 EIGGDGQTTDGTEPSHMHYALDENYFRGYEWWLMKEAKKRNPNITLIGLP    150
              |||||||||||||||||||-||||||||||||||||||||||:|.|:|||
Mouse     101 EIGGDGQTTDGTEPSHMHYELDENYFRGYEWWLMKEAKKRNPDIILMGLP    150

Human     151 WSFPGWLGKGFDWPYVNLQLTAYYVVTWIVGAKRYHDLDIDYIGIWNERS    200
              |||||||||-||||||||||||||||-||:|||-|||||||||||||||-
Mouse     151 WSFPGWLGKGFSWPYVNLQLTAYYVVRWILGAKHYHDLDIDYIGIWNERP    200

Human     201 YNANYIKILRKMLNYQGLQRVKIIASDNLWESISASMLLDAELFKVVDVI    250
              ::||||||-||||:|||||||:||||||||-||:|:|||-||:||||||
Mouse     201 FDANYIKELRKMLDYQGLQRVRIIASDNLWEPISSSLLLDQELWKVVDVI    250

Human     251 GAHYPGTHSAKDAKLTGKKLWSSEDFSTLNSDMGAGCWGRILNQNYINGY    300
              |||||||::..:||::||||||||||||::||:|||||-|||||||||-
Mouse     251 GAHYPGTYTVWNAKMSGKKLWSSEDFSTINSNVGAGCWSRILNQNYINGN    300

Human     301 MTSTIAWNLVASYYEQLPYGRCGLMTAQEPWSGHYVVESPVWVSAHTTQF    350
              ||||||||||||||:||||||-|||||||||||||||-||:|||||||
Mouse     301 MTSTIAWNLVASYYEELPYGRSGLMTAQEPWSGHYVVASPIWVSAHTTQF    350

Human     351 TQPGWYYLKTVGHLEKGGSYVALTDGLGNLTIIIETMSHKHSKCIRPFLP    400
              ||||||||||||||||||||||||||||||||||||||:||-||||-|
Mouse     351 TQPGWYYLKTVGHLEKGGSYVALTDGLGNLTIIIETMSHQHSMCIRPYLP    400

Human     401 YFNVSQQFATFVLKGSFSEIPELQVWYTKLGKTSERFLFKQLDSLWLLDS    450
              |:|||-|.|||-||||..||-|||||||||||...:|..||||:|||||-
Mouse     401 YYNVSHQLATFTLKGSLREIQELQVWYTKLGTPQQRLHFKQLDTLWLLDG    450

Human     451 DGSFTLSLHEDELFTLTTLTTGRKGSYPLPPKSQPFPSTYKDDFNVDYPF    500
              -||||||-|-||||:||||||||||||||||||-||.|:|||:-||||||||:||-
Mouse     451 SGSFTLELEEDEIFTLTTLTTGRKGSYPPPPSSKPFPTNYKDDFNVEYPL    500

Human     501 FSEAPNFADQTGVFEYFTNIEDPGEHHFTLRQVLNQRPITWAADASNTIS    550
              |||||||||||||||||:.|.|| .||-|||||||||||||||||:|||
Mouse     501 FSEAPNFADQTGVFEYYMNNED-REHRFTLRQVLNQRPITWAADASSTIS    549
```

# You can also use UniProt to do alignments

- Search for your gene/protein of interest

- Select the ones you want to compare

- Click Align

# UniProt Alignment

# You can also make phylogenetic trees

- Select all genes

- Do Alignment

# You can also make phylogenetic trees

- Select all genes

- Do Alignment

- Select Trees

# Sequence homology

- Identify conservation across species

- Highlight important functional regions of high conservation

# Sequence homology



| | Description | Scientific Name | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession |
|---|---|---|---|---|---|---|---|---|---|
| ☑ | galactocerebrosidase isoform a precursor [Homo sapiens] | Homo sapiens | 1414 | 1414 | 100% | 0.0 | 100.00% | 685 | NP_000144.2 |
| ☑ | galactocerebrosidase isoform c precursor [Homo sapiens] | Homo sapiens | 1355 | 1355 | 100% | 0.0 | 96.64% | 662 | NP_001188330.1 |
| ☑ | galactocerebrosidase isoform X1 [Homo sapiens] | Homo sapiens | 1285 | 1285 | 90% | 0.0 | 100.00% | 629 | XP_011534920.1 |
| ☑ | galactocerebrosidase isoform d [Homo sapiens] | Homo sapiens | 1285 | 1285 | 90% | 0.0 | 99.84% | 659 | NP_001188331.1 |
| ☑ | galactocerebrosidase isoform X2 [Homo sapiens] | Homo sapiens | 1181 | 1181 | 83% | 0.0 | 100.00% | 569 | XP_047287154.1 |
| ☑ | galactocerebrosidase precursor [Mus musculus] | Mus musculus | 1157 | 1157 | 100% | 0.0 | 82.63% | 684 | NP_032105.2 |
| ☑ | galactocerebrosidase isoform X1 [Mus musculus] | Mus musculus | 951 | 951 | 83% | 0.0 | 81.20% | 568 | XP_017170449.1 |
| ☑ | galactocerebrosidase precursor [Danio rerio] | Danio rerio | 856 | 856 | 95% | 0.0 | 64.08% | 660 | NP_001005921.1 |
| ☑ | galactocerebrosidase precursor [Danio rerio] | Danio rerio | 842 | 842 | 94% | 0.0 | 64.96% | 664 | NP_998276.1 |
| ☑ | galactocerebrosidase isoform X2 [Mus musculus] | Mus musculus | 759 | 759 | 60% | 0.0 | 86.81% | 469 | XP_006515535.1 |
| ☑ | Putative galactocerebrosidase [Caenorhabditis elegans] | Caenorhabditis elegans | 285 | 285 | 85% | 2e-85 | 34.29% | 645 | NP_498726.3 |

UNIVERSITY OF CAMBRIDGE

# Sequence homology



| | Description | Scientific Name | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession |
|---|---|---|---|---|---|---|---|---|---|
| ☑ | galactocerebrosidase isoform a precursor [Homo sapiens] | Homo sapiens | 1414 | 1414 | 100% | 0.0 | 100.00% | 685 | NP_000144.2 |
| ☑ | galactocerebrosidase isoform c precursor [Homo sapiens] | Homo sapiens | 1355 | 1355 | 100% | 0.0 | 96.64% | 662 | NP_001188330.1 |
| ☑ | galactocerebrosidase isoform X1 [Homo sapiens] | Homo sapiens | 1285 | 1285 | 90% | 0.0 | 100.00% | 629 | XP_011534920.1 |
| ☑ | galactocerebrosidase isoform d [Homo sapiens] | Homo sapiens | 1285 | 1285 | 90% | 0.0 | 99.84% | 659 | NP_001188331.1 |
| ☑ | galactocerebrosidase isoform X2 [Homo sapiens] | Homo sapiens | 1181 | 1181 | 83% | 0.0 | 100.00% | 569 | XP_047287154.1 |
| ☑ | galactocerebrosidase precursor [Mus musculus] | Mus musculus | 1157 | 1157 | 100% | 0.0 | 82.63% | 684 | NP_032105.2 |
| ☑ | galactocerebrosidase isoform X1 [Mus musculus] | Mus musculus | 951 | 951 | 83% | 0.0 | 81.20% | 568 | XP_017170449.1 |
| ☑ | galactocerebrosidase precursor [Danio rerio] | Danio rerio | 856 | 856 | 95% | 0.0 | 64.08% | 660 | NP_001005921.1 |
| ☑ | galactocerebrosidase precursor [Danio rerio] | Danio rerio | 842 | 842 | 94% | 0.0 | 64.96% | 664 | NP_998276.1 |
| ☑ | galactocerebrosidase isoform X2 [Mus musculus] | Mus musculus | 759 | 759 | 60% | 0.0 | 86.81% | 469 | XP_006515535.1 |
| ☑ | Putative galactocerebrosidase [Caenorhabditis elegans] | Caenorhabditis elegans | 285 | 285 | 85% | 2e-85 | 34.29% | | 498726.3 |

# Comparing Outputs – trees and alignments

# Including conservation improves search sensitivity

- Standard sequence alignment maximises the correspondence of residues across both proteins

# Including conservation improves search sensitivity

- Standard sequence alignment maximises the correspondence of residues across both proteins

- But some residues will be *evolutionarily conserved* while others won't
  - Conserved residues more likely to be important for function

- Weighting the alignment by conservation makes it more sensitive and accurate
  - Profiles and Hidden Markov Models (HMMs)



Logo Plot

UNIVERSITY OF
CAMBRIDGE

# Profile/HMM searching for identifying distant homologs

- Use the *query sequence* to perform an initial search of the *sequence database*

- Take all hits from initial search and build a profile or HMM

- Use this profile/HMM to perform a subsequent search of the *same sequence database*

  - Will be more sensitive and accurate

- Profile searching: PSI-BLAST

- HMM searching: HMMER (HMMsearch)

UNIVERSITY OF CAMBRIDGE

# HMM:HMM alignments

- Can also use HMM/HMM comparisons to improve specific multiple sequence alignments
  - Upweight the alignment of regions that are evolutionarily conserved
- Clustal Omega

# What else can we predict from sequence?

- Secondary structure prediction using NetSurfP

  https://dtu.biolib.com/NetSurfP-3/

# Secondary structure and disorder

# Secondary structure and disorder

| | | | | | | | Probability of | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | | Alpha | Beta | Coil |

# Column 1: Class assignment - B for buried or E for Exposed - Threshold: 25% exposure, but not based on RSA
# Column 2: Amino acid
# Column 3: Sequence name
# Column 4: Amino acid number
# Column 5: Relative Surface Accessibility - RSA
# Column 6: Absolute Surface Accessibility
# Column 7: Not used
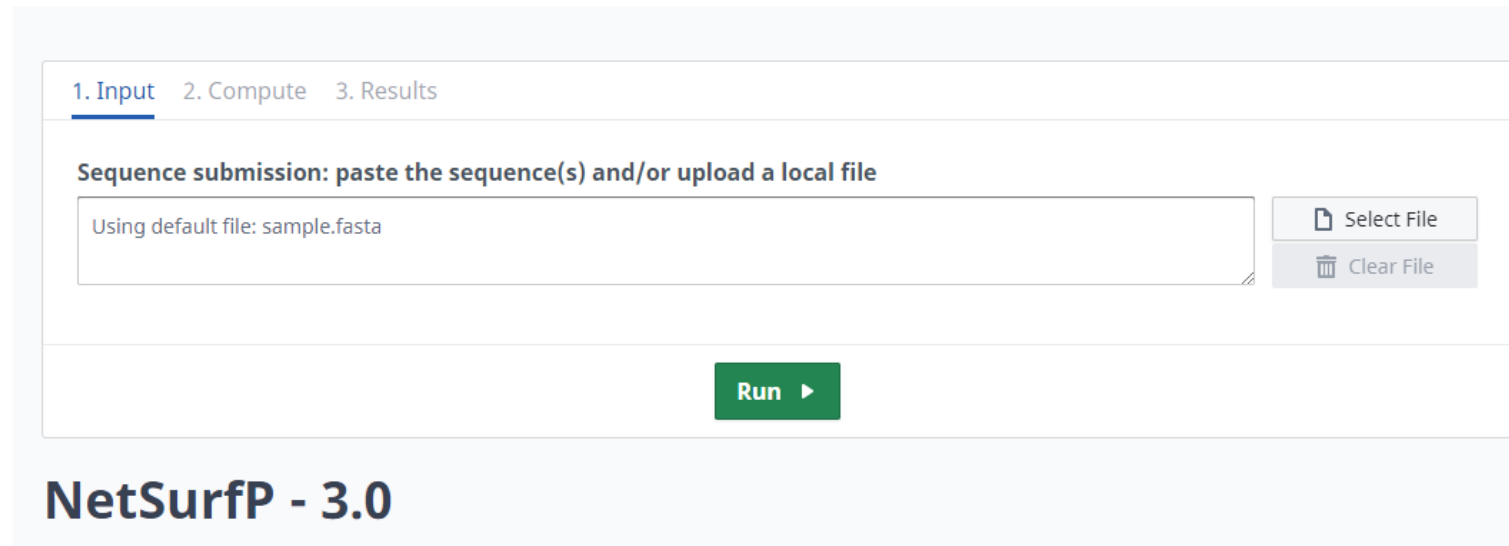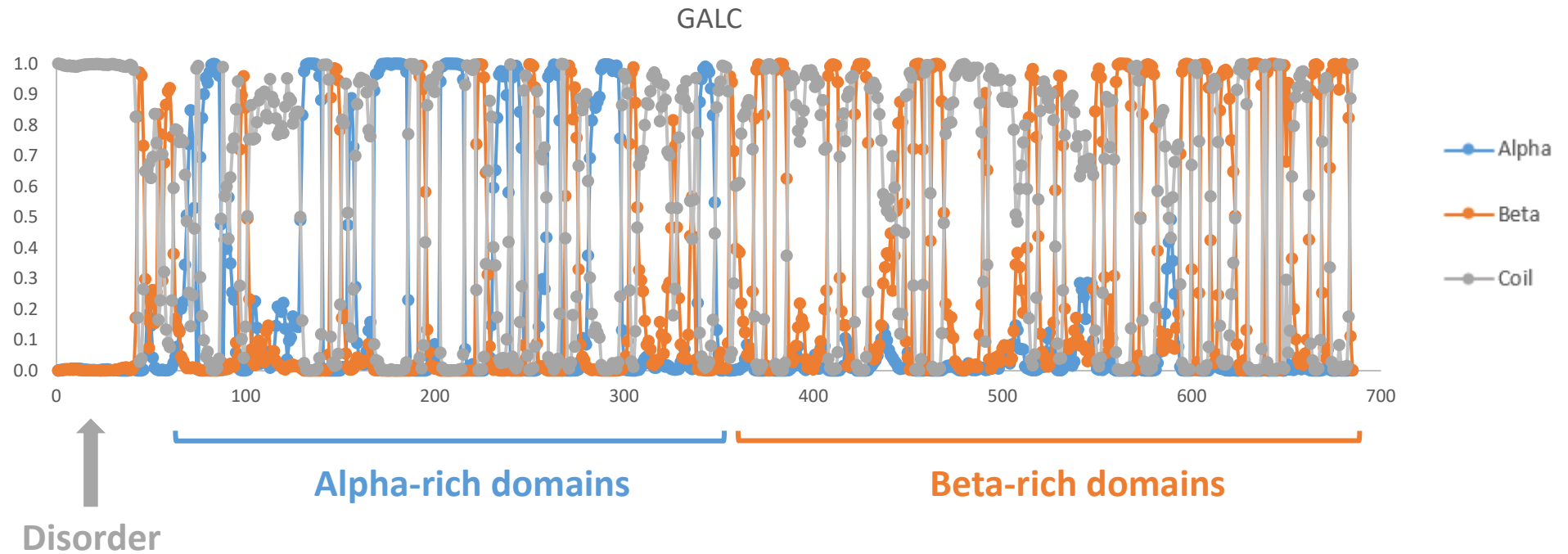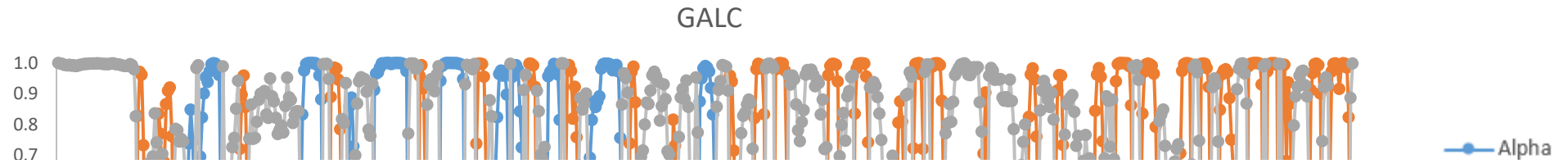# Column 8: Probability for Alpha-Helix
# Column 9: Probability for Beta-strand
# Column 10: Probability for Coil

| | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| E | M | NFASC_re | 1 | 0.890999 | 199.5837 | 0 | 2.52E-05 | 8.46E-05 | 0.99989 |
| E | A | NFASC_re | 2 | 0.822632 | 106.1195 | 0 | 0.000178 | 0.000218 | 0.999604 |
| E | R | NFASC_re | 3 | 0.789607 | 216.3524 | 0 | 0.000201 | 0.000312 | 0.999488 |
| E | Q | NFASC_re | 4 | 0.780665 | 175.6495 | 0 | 5.27E-05 | 0.000324 | 0.999623 |
| E | P | NFASC_re | 5 | 0.716381 | 113.9045 | 0 | 2.86E-05 | 0.00016 | 0.999811 |
| E | P | NFASC_re | 6 | 0.703382 | 111.8378 | 0 | 3.21E-05 | 0.000161 | 0.999807 |
| E | P | NFASC_re | 7 | 0.6822 | 108.4698 | 0 | 7.56E-05 | 0.000202 | 0.999722 |
| E | P | NFASC_re | 8 | 0.68275 | 108.5572 | 0 | 0.000286 | 0.000226 | 0.999488 |
| E | W | NFASC_re | 9 | 0.639935 | 182.3814 | 0 | 0.000634 | 0.000545 | 0.998821 |
| E | V | NFASC_re | 10 | 0.579066 | 100.7574 | 0 | 0.000555 | 0.00063 | 0.998815 |
| E | H | NFASC_re | 11 | 0.641862 | 143.7771 | 0 | 0.000518 | 0.000702 | 0.99878 |
| E | A | NFASC_re | 12 | 0.601697 | 77.61896 | 0 | 0.00058 | 0.000626 | 0.998793 |
| E | A | NFASC_re | 13 | 0.548697 | 70.7819 | 0 | 0.000435 | 0.001032 | 0.998533 |
| E | F | NFASC_re | 14 | 0.536592 | 128.7821 | 0 | 0.000239 | 0.002572 | 0.99719 |
| E | L | NFASC_re | 15 | 0.566734 | 113.9135 | 0 | 0.000126 | 0.003813 | 0.996061 |
| E | L | NFASC_re | 16 | 0.573608 | 115.2952 | 0 | 7.97E-05 | 0.003012 | 0.996908 |
| E | C | NFASC_re | 17 | 0.54777 | 91.47765 | 0 | 0.000106 | 0.003483 | 0.996411 |
| E | L | NFASC_re | 18 | 0.644464 | 129.5373 | 0 | 9.82E-05 | 0.003267 | 0.996635 |
| E | L | NFASC_re | 19 | 0.66675 | 134.0167 | 0 | 8.86E-05 | 0.002037 | 0.997875 |
| E | S | NFASC_re | 20 | 0.747813 | 115.911 | 0 | 0.000147 | 0.001159 | 0.998693 |
| E | L | NFASC_re | 21 | 0.736005 | 147.937 | 0 | 0.000195 | 0.000498 | 0.999306 |
| E | G | NFASC_re | 22 | 0.797154 | 82.904 | 0 | 0.000467 | 0.000326 | 0.999207 |
| E | G | NFASC_re | 23 | 0.817293 | 84.99846 | 0 | 0.001605 | 0.000448 | 0.997946 |
| E | A | NFASC_re | 24 | 0.804079 | 103.7262 | 0 | 0.002641 | 0.0008 | 0.996559 |

# Secondary structure and disorder



GALC

Alpha-rich domains

Beta-rich domains
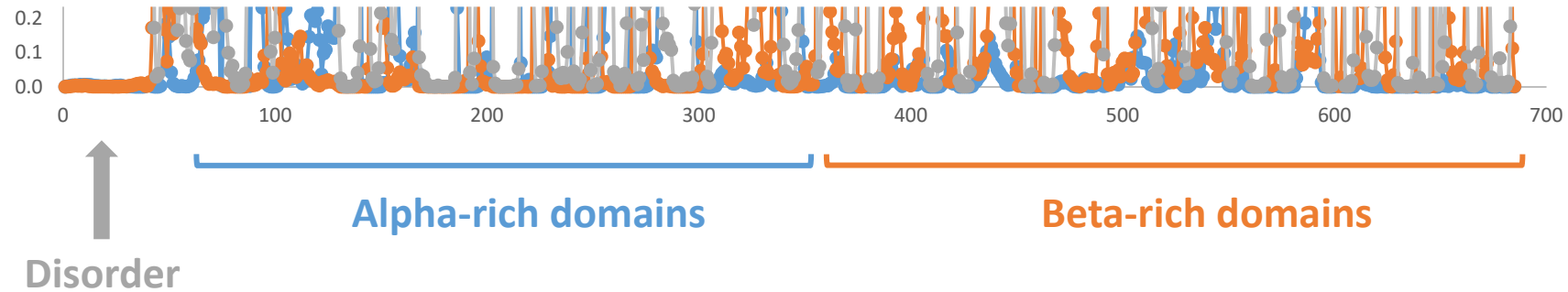
Disorder

# Secondary structure and disorder



GALC

We can (often) do even better than this with AlphaFold, but that's for next week!

Alpha-rich domains

Beta-rich domains

Disorder

UNIVERSITY OF CAMBRIDGE

# Summary

- Today you learnt how to use online resources to:
  - Predict domains
  - Identify post translational modifications
  - Calculate the molecular mass
  - Determine the isoelectric point
  - How to calculate extinction co-efficients
  - Identify distant homologs

- Try this out with your favourite protein!

UniProt

Expasy

ProtParam tool

ApE
A plasmid Editor

UNIVERSITY OF CAMBRIDGE