
mousePB101(H-11C1)

Taeyoon kim

May 27, 2022

CONTENTS

1	Protein parameters analysis	3
1.1	Amino acid composition	4
1.2	Potential sites of chemical modification	5
1.3	Secondary structure fraction	6
1.4	Secondary structure prediction	7
1.5	Structural analysis	7
1.6	Protein Scales	8
2	Immunogenicity analysis	13
2.1	MHC class 1	13
2.2	MHC class 2	15
2.3	Appendix	17

Introduction

The data in the report is intended to be a resource to guide the development and lead-optimization of a therapeutic protein. These data can be used in a preemptive fashion - for example, in the decision to substitute an exposed residue on the surface that may be prone to the kind of chemical modification that might affect the stability of the protein. They can also be used to assist the troubleshooting of problems that can arise in the course of an clinical development - for example if an therapeutic protein displays stability issues in storage, or unacceptably high levels of immunogenicity in early clinical trials.

There is always a great deal of risk involved in the development of any therapeutic molecule but experience has shown that the kind of data presented in this report is an invaluable tool for mitigating that risk - either by helping to identify potential problems before they occur, or by guiding the troubleshooting of problems that can occur during the therapeutic protein development and lead-optimization.

Background

H-11C1 is an anti-VEGF antibody, which has previously been used as a surrogate for preclinical modelling of PB101 activity.

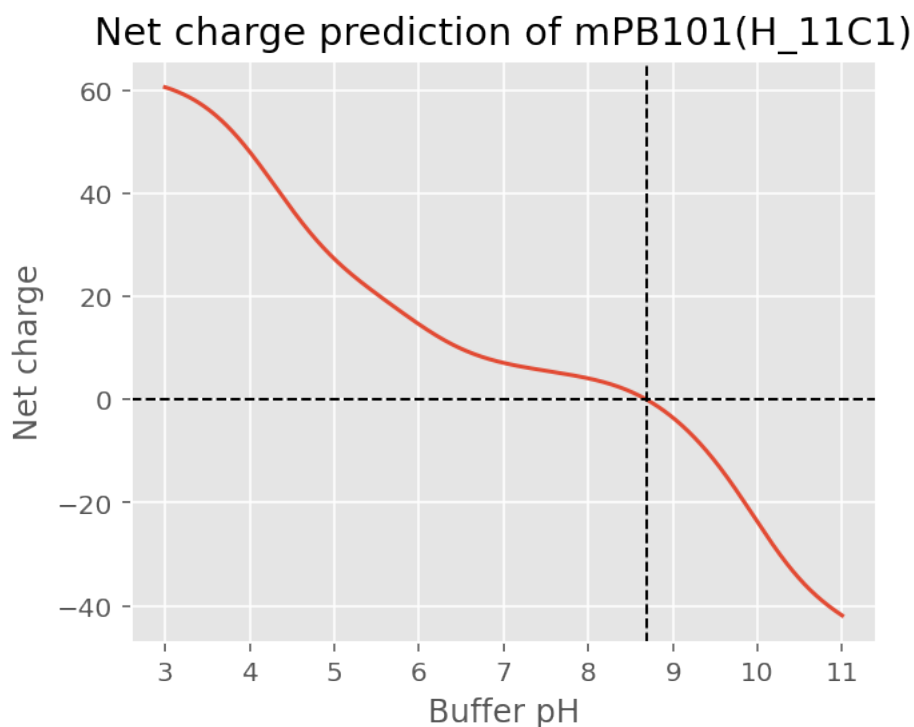
Protein sequence

```
>mPB101 (H_11C1)
GSPFIEMHTDIPKLVHMTGRQLIIPCRVTSPNVTITLKKFPFDLTLPDQGRI TWDSRRGFIIANATYKEIGLLNCEATV
NGHLYQTNYLTHRQNTITLDVQIRPPSPVTLLHGQTLVLNCTATTELNTRVQMSWNPYPGKATKNASIRQRIDRSHSHNNV
FHSV LKINNVESRDKGLYTCRVKSGSSFQSFNTSVHVYEGLEKPCICTVPEVSSVFIFPPKPKDVLTTITLTPKVT CVVV
DISKDDPEVQFSWFVDDVEVHTAQTQPREEQFNSTFRSVSELPIMHQDWLNGKEFKCRVNSAAFPAPIEKTISKTKGRPK
APQVYTIPPPKEQMAKDKVSLTCMITDFFPEDITVEWQWNGQPAENYKNTQPI MNNGSYFVYSKLVQKSNWEAGNTFT
CSVLHEGLHNHHTKSLSHSPGK
```


PROTEIN PARAMETERS ANALYSIS

The program performs most of the same functions as the Expasy ProtParam tool.

```
# Name of target protein: -----mPB101(H_11C1)
# Molecular weight(Dalton): -----47,931
# Total number of amino acid: -----423
# Chemical formula: -----C2133H3331N589O635S17
# Total number of atom is: -----6705
# Extinction coefficient(reduced): -----53400
# Reduced Abs 0.1%(=1 g/L): -----1.114
# Extinction coefficient(non-reduced): -----54025
# Non-reduced Abs 0.1%(=1 g/L): -----1.127
# Theoretical pI: -----8.687
# Aromaticity: -----8.51%
# The GRAVY value is -----0.483
mPB101(H_11C1) is more hydrophilic protein.
# The instability index: -----36.423
mPB101(H_11C1) is seems stable.
```

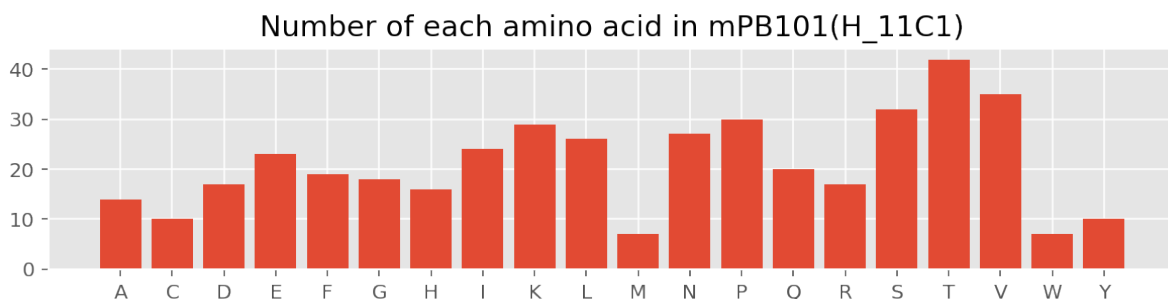


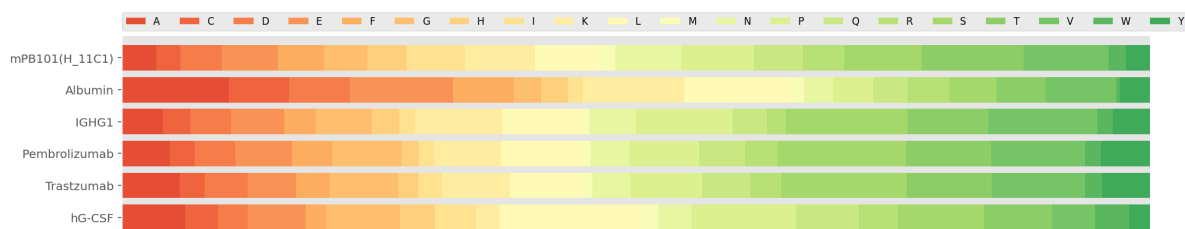
- Molecular weight
 - Amino acids are the building blocks that form polypeptides and ultimately proteins. Calculates the molecular weight of a protein.
- Chemical composition
 - A chemical formula is a way of presenting information about the chemical proportions of atoms that constitute a particular chemical compound or molecule, using chemical element symbols, numbers.
- Extinction coefficient
 - Extinction (or extinction coefficient) is defined as the ratio of maximum to minimum transmission of a beam of light that passes through a polarization optical train. extinction coefficient in units of $M^{-1}cm^{-1}$, at 280 nm measured in water.
- Theoretical pI
 - The isoelectric point (pI, pH(I), IEP), is the pH at which a molecule carries no net electrical charge or is electrically neutral in the statistical mean. The pI value can affect the solubility of a molecule at a given pH. Such molecules have minimum solubility in water or salt solutions at the pH that corresponds to their pI and often precipitate out of solution. Biological amphoteric molecules such as proteins contain both acidic and basic functional groups.
- Aromaticity
 - Calculate the aromaticity according to Lobry, 1994. Calculates the aromaticity value of a protein according to Lobry, 1994. It is simply the relative frequency of Phe+Trp+Tyr.
- GRAVY
 - The GRAVY value is calculated by adding the hydropathy value for each residue and dividing by the length of the sequence (Kyte and Doolittle; 1982). A higher value is more hydrophobic. A lower value is more hydrophilic.
- Instability_index
 - Implementation of the method of Guruprasad et al. (1990, Protein Engineering, 4, 155-161). This method tests a protein for stability. Any value above 40 means the protein is unstable (=has a short half life).

1.1 Amino acid composition

We can easily count the number of each type of amino acid.

```
Total number of positively charged residues(Arg + Lys) :-----46
Total number of negatively charged residues(Asp + Glu) :-----40
```





1.2 Potential sites of chemical modification

An initial scan of the protein sequences is presented based purely upon sequence. If a structural analysis was also requested, this section should be used in conjunction with the molecular surface analysis described in a subsequent section. Any of the sites listed below could be candidates for further consideration if the molecular surface analysis shows that they are significantly exposed on the surface of the protein, increasing their propensity for chemical modification. The canonical sequence analysis is also helpful here, since each of these sites can also be considered in the context of their frequency of occurrence within the canonical library of homologous sequences.

1.2.1 Potential deamidation positions

Asparagine (N) and glutamine (Q) residues are particularly prone to deamidation when they are followed in the sequence by amino acids with smaller side chains, that leave the intervening peptide group more exposed. Deamidation proceeds much more quickly if the susceptible amino acid is followed by a small, flexible residue such as glycine whose low steric hindrance leaves the peptide group open for attack.

- Search patterns: ASN/GLN-ALA/GLY/SER/THR

```

65-NA-66
81-NG-82
86-QT-87
94-QT-95
96-NT-97
115-QT-116
128-NT-129
144-NA-145
189-QS-190
192-NT-193
264-QT-265
273-NS-274
291-NG-292
300-NS-301
360-NG-361
369-NT-370
375-NT-376
377-NG-378
397-NT-398

```

1.2.2 Potential o-linked glycosylation sites

The O-linked glycosylation of serine and threonine residues seems to be particularly sensitive to the presence of one or more proline residues in their vicinity in the sequence, particularly in the -2 and +3 positions.

- Search patterns: PRO-SER/THR

106-PS-107

- Search patterns: SER/THR-X-X-PRO

9-TDIP-12
45-TLTP-48
230-TLTP-233
280-SELP-283
326-TIPP-329
418-SHSP-421

1.2.3 Potential n-linked glycosylation sites

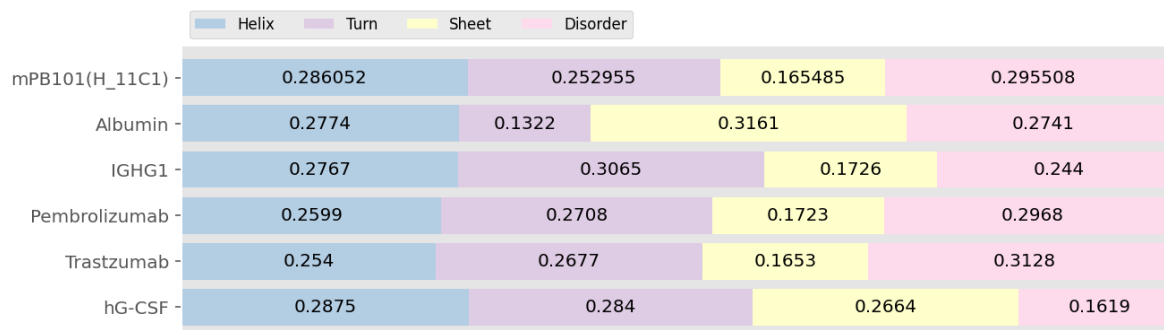
- Search patterns: ASN-X-SER/THR

33-NVT-35
65-NAT-67
120-NCT-122
144-NAS-146
192-NTS-194
273-NST-275
377-NGS-379

1.3 Secondary structure fraction

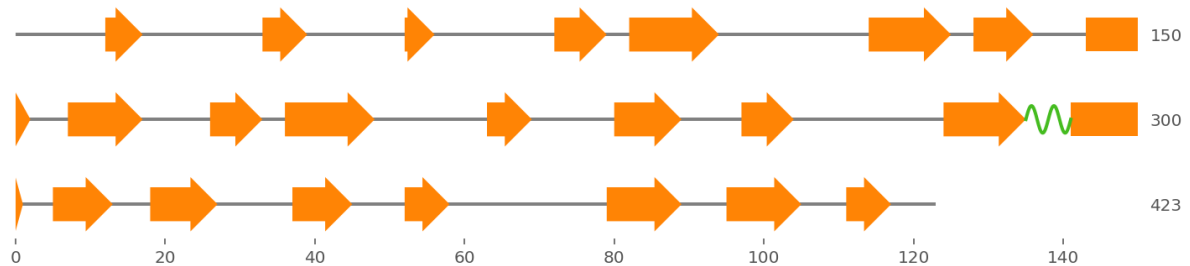
The fraction of amino acids that tend to be found in the three classical secondary structures. These are beta sheets, alpha helices, and turns (where the residues change direction).

- Amino acids in helix: V, I, Y, F, W, L.
- Amino acids in turn: N, P, G, S.
- Amino acids in sheet: E, M, A, L.



1.4 Secondary structure prediction

Protein secondary structure prediction is one of the most important and challenging problems in bioinformatics. Here in, the P-SEA algorithm that to predict the secondary structures of proteins sequences based only on knowledge of their primary structure.



1.5 Structural analysis

1.5.1 Detection of disulfide bonds

This function detects disulfide bridges in protein structures. Then the detected disulfide bonds are visualized and added to the bonds attribute of the AtomArray. The employed criteria for disulfide bonds are quite simple in this case: the atoms of two cysteine residues must be in a vicinity of Å and the dihedral angle of must be .

C	121	CYS	SG	S	-18.443	-23.604	-18.026
C	180	CYS	SG	S	-19.369	-25.116	-19.127



The found disulfide bonds are visualized with the help of Matplotlib: The amino acid sequence is written on the X-axis and the disulfide bonds are depicted by yellow semi-ellipses.

1.5.2 Calculation of protein diameter

This calculates the diameter of a protein defined as the maximum pairwise atom distance.

```
# Diameter of mPB101(H_11C1) is: -----136.385 Angstrong.
```

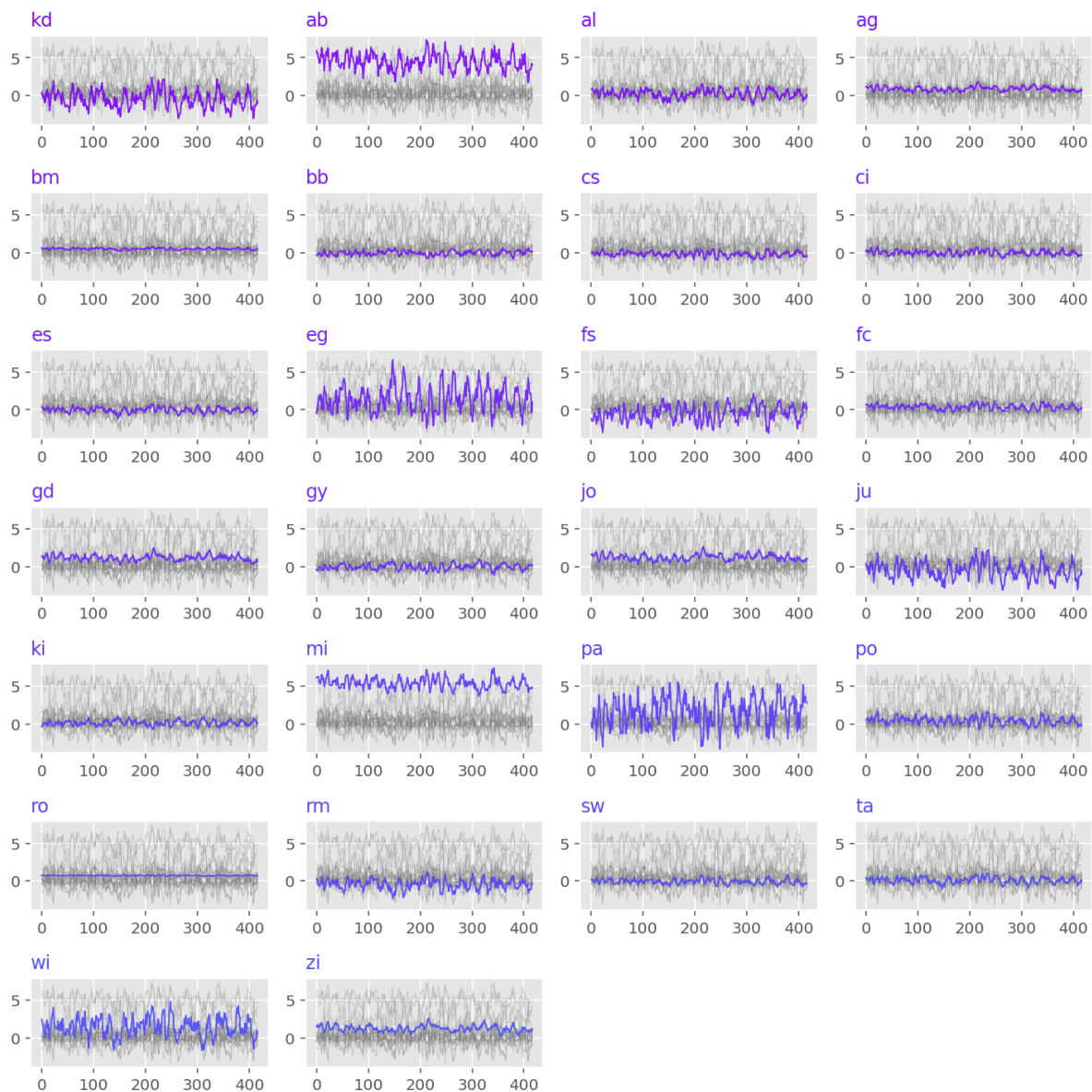
1.6 Protein Scales

Protein scales are a way of measuring certain attributes of residues over the length of the peptide sequence using a sliding window. Scales are comprised of values for each amino acid based on different physical and chemical properties, such as hydrophobicity, secondary structure tendencies, and surface accessibility. As opposed to some chain-level measures like overall molecule behavior, scales allow a more granular understanding of how smaller sections of the sequence will behave.

- kd → Kyte & Doolittle Index of Hydrophobicity
- hw → Hopp & Wood Index of Hydrophilicity
- em → Emini Surface fractional probability (Surface Accessibility)
-

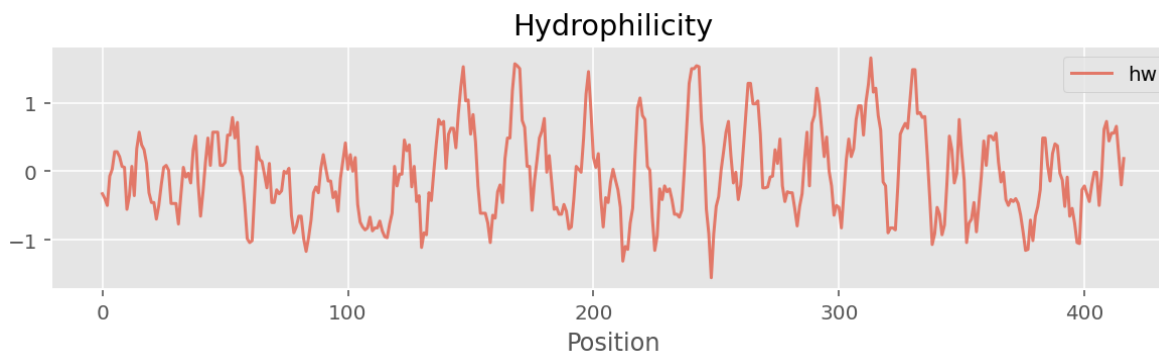
1.6.1 Hydrophobicity index

hydrophobicity is the physical property of a molecule that is seemingly repelled from a mass of water (known as a hydrophobe).



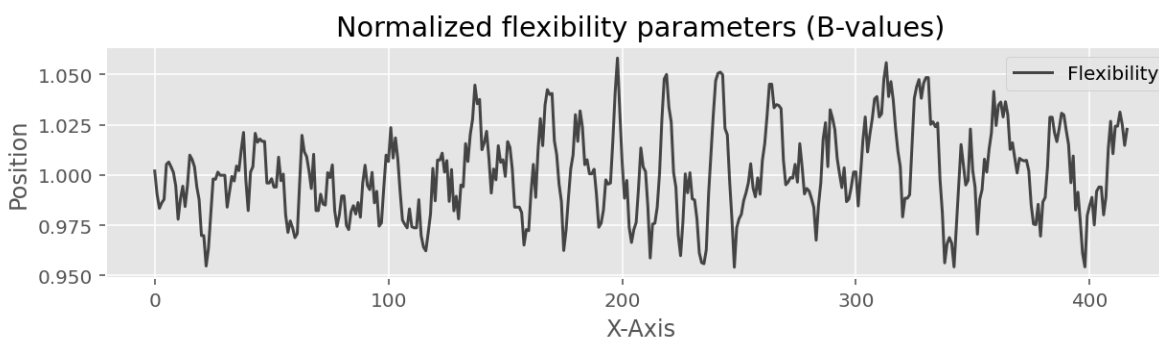
1.6.2 Hydrophilicity index

Hydrophilicity is the tendency of a molecule to be solvated by water.



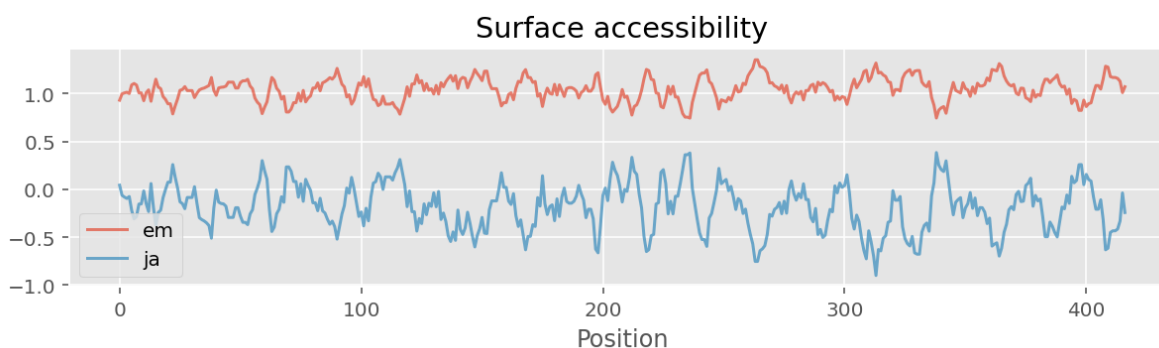
1.6.3 Flexibility index

Proteins are dynamic entities, and they possess an inherent flexibility that allows them to function through molecular interactions within the cell.



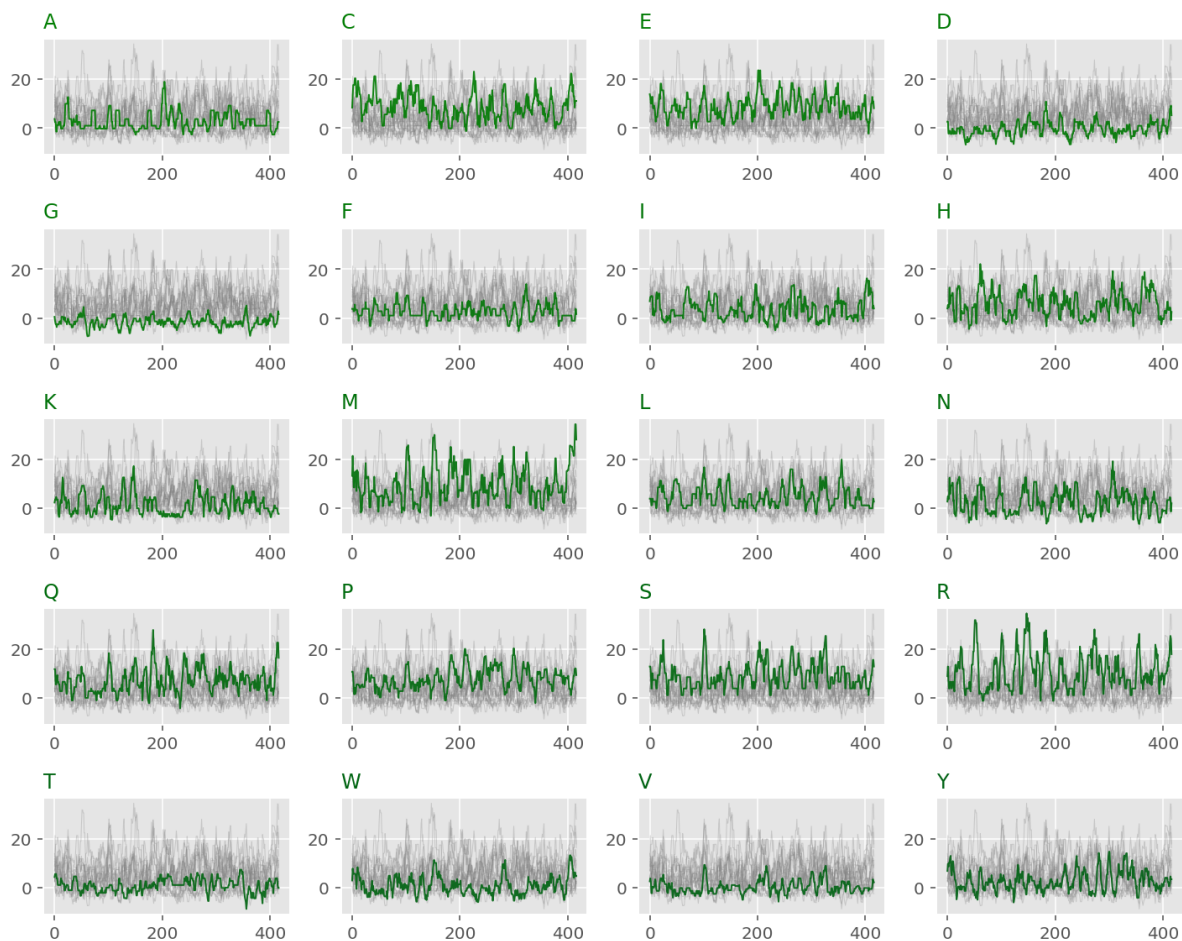
1.6.4 Surface accessibility

Data describing the solvent-accessible surface of a molecule is of great utility in the development of that molecule as a therapeutic, particularly in the case of antibodies. In the context of this report, the most obvious application of molecular surface data is in combination with the potential sites of chemical modification, described in the previous section. Proteins are known to undergo many different chemical modifications as a result of interactions with their aqueous environment. The probability and kinetic rate of such a modification is greatly enhanced by the degree of exposure of the potential modification site to the solvent environment. The solvent-accessible surface for each residue depends upon the degree of exposure of the residue on the surface, but also on the size of the residue side chain.



1.6.5 Instability index

The instability index provides an estimate of the stability of your protein in a test tube. Statistical analysis of 12 unstable and 32 stable proteins has revealed that there are certain dipeptides, the occurrence of which is significantly different in the unstable proteins compared with those in the stable ones.



IMMUNOGENICITY ANALYSIS

We use the method of removing and/or reducing potential T-cell epitopes, as an approach to the management of the immunogenicity of biologics. The protein sequence is scanned *in silico*, for sequences that have a strong binding signature for a family of 50 MHC Class II receptors, whose alleles cover 96 – 98% of the human population. The presented histograms for each variable region sequence, show the average (for the *n* positively-testing MHC II alleles) of epitope strength at each position as a percentage for all epitopes above a threshold of 20%. At each position in the sequence, the number of alleles scoring above the threshold is shown above the histogram at that position. The epitopes of most concern for the antibody's immunogenicity are therefore those that have not just the highest average score per allele (as shown by the histogram), but which also score above the threshold across more alleles, since these epitopes are more likely to engender an immune response in a larger fraction of the patient population.

Experience using *in silico* algorithms of this kind in conjunction with laboratory immunogenicity assays has shown that epitopes below this threshold do not generally contribute significantly to the protein's immunogenicity. The number of alleles, the affected alleles and their individual scores are also listed in the detailed analyses below each histogram figure.

The raw immunogenicity score quoted is the total over all epitopes above the threshold for all affected alleles. The normalized immunogenicity score is this raw score divided by the sequence length, and represents epitope strength per unit sequence to enable comparisons of protein sequences of different lengths.

2.1 MHC class 1

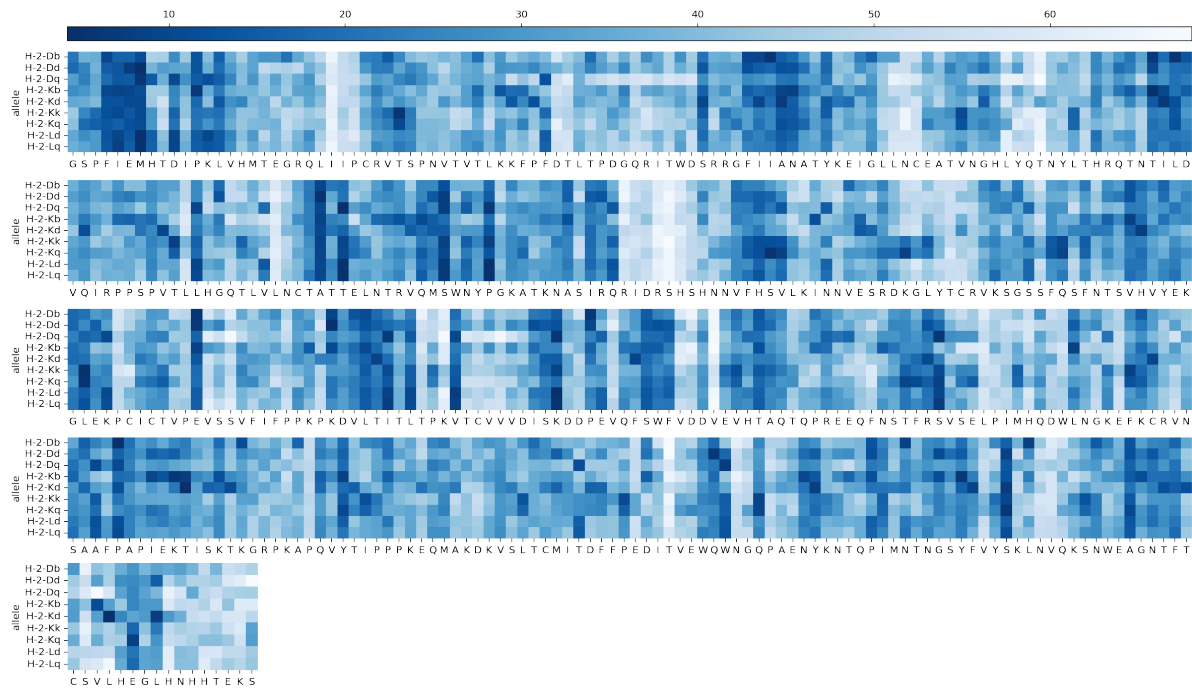
Class I major histocompatibility complex (MHC) molecules bind, and present to T cells, short peptides derived from intracellular processing of proteins. The peptide repertoire of a specific molecule is to a large extent determined by the molecular structure accommodating so-called main anchor positions of the presented peptide.

Their function is to display peptide fragments of proteins from within the cell to cytotoxic T cells; this will trigger an immediate response from the immune system against a particular non-self antigen displayed with the help of an MHC class I protein. Because MHC class I molecules present peptides derived from cytosolic proteins, the pathway of MHC class I presentation is often called cytosolic or endogenous pathway.¹

- MHC class 1 superset
 - HLA-A01:01, HLA-A02:01, HLA-A03:01, HLA-A24:02, HLA-B07:02, HLA-B40:01

¹ Kimball's Biology Pages, Histocompatibility Molecules

2.1.1 Predicts binding of peptides to MHC class1



2.1.2 Top10 strong binding peptide

	allele
Core	
FPEDITVEW	5
FPPKPKDVL	4
RPPSPVTLL	4
FPAPIEKT	4
IRPPSPVT	3
KPDVLTITL	3
TSPNVTVTL	3
STFRSVSEL	3
SPVTLLHTL	3
SPNVTVTLF	3

