

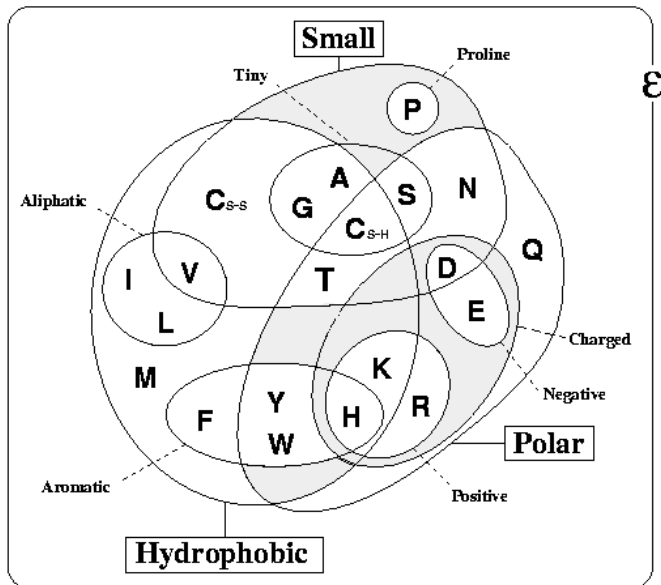
Sequence-based mutation analysis

SNP effects prediction from sequence

Sebastian Hollizeck, Robert Wagner

Practical Bioinformatics 'Protein Structure and Function Analysis'

June 11, 2012



chemical properties of amionacids

- acidic
- aliphatic
- aromatic
- basic
- hydroxyl
- sulfur containing

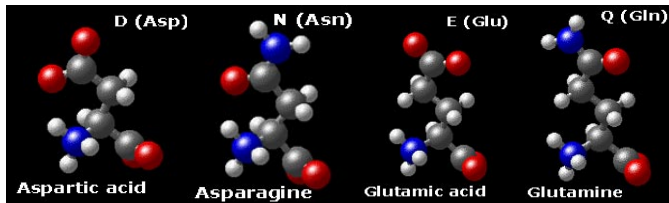


Figure: Acidic amino acids and their amides
all very hydrophilic

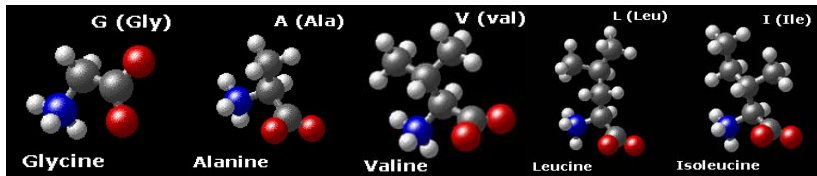


Figure: Aliphatic amino acids increase in hydrophobicity from left to right

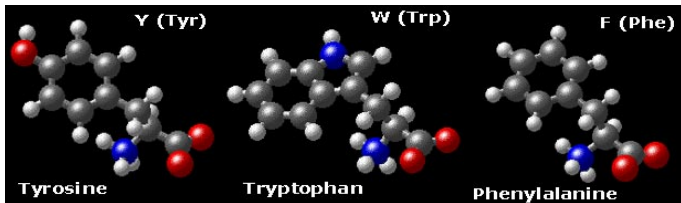


Figure: Aromatic amino acids increase in hydrophobicity from left to right

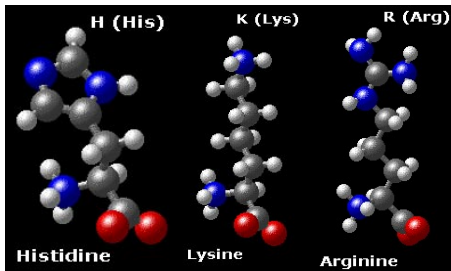


Figure: Basic amino acids

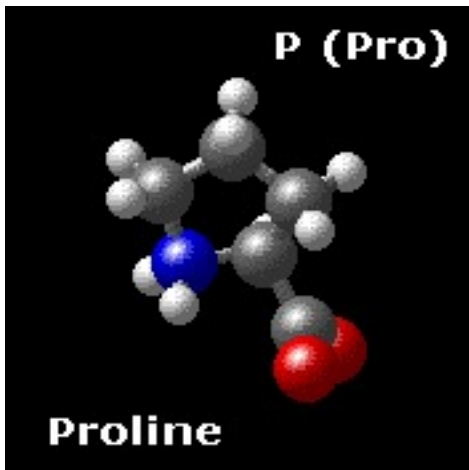


Figure: Proline - the helixbreaker

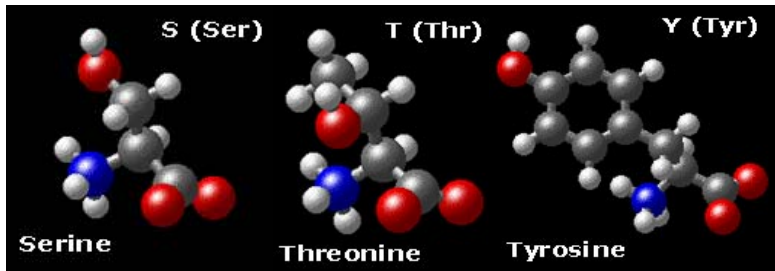


Figure: Hydroxyl amino acids

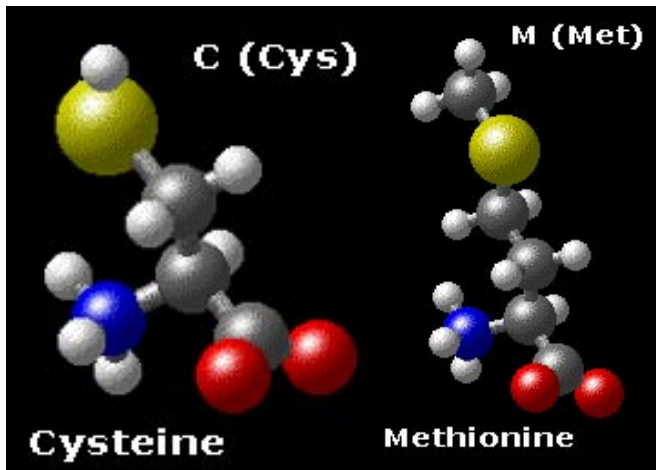


Figure: Sulphur containing amino acids

From Glycine to Leucine
From Alanine to Tryptophan
From Asparagine to Glutamine

common substitution matrices

- PAM
- BLOSUM

substitution matrices for special occasions

- secondary structure based
- context based

PAM means **P**oint **A**ccepted **M**utations

- assumes independence of mutations
- for closely related sequences
- for more divergent sequences power the original PAM-matrix

most common PAM1, PAM30, PAM70 and PAM250

BLOSUM means **BLO**ck **SU**bstitution **M**atrix
approach for more distant related sequences

- derived from distantly related sequences
- finding conserved regions in the multiple alignment
- clustering of sequences to reduce impact of high similarity

most common BLOSUM62 (used in BLAST)

| | | |
|------------------------|------------------|-------------------|
| aspects | PAM | BLOSUM |
| evolutionary model | explicit | implicit |
| based on | global alignment | conserved regions |
| weighting of mutations | equal | clustering |
| naming for similarity | lower number | higher number |

For the practical we will use the following tools to predict effects of SNPs:

- SIFT
- PolyPhen2
- SNAP

Sorting intolerant from tolerant substitutions

- Uses sequence homology
- Searches similar sequences
- Chooses closely related sequences
- Obtains multiple sequence alignment

- Calculates normalized possibilities for all possible substitutions for each position
- Every substitution below a given cutoff is assumed to be deleterious
- Based on BLOSSUM Matrix

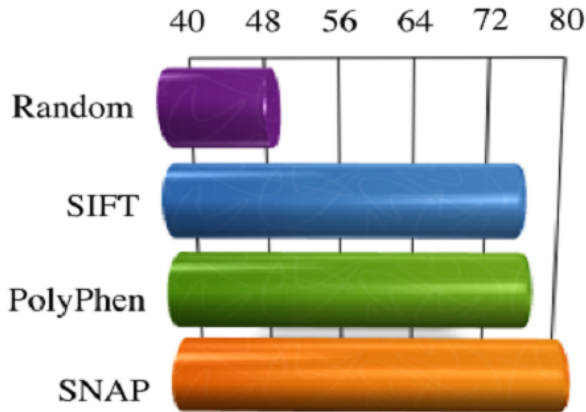
Polymorphism Phenotyping v2

- Basic idea very similar to SIFT - but using machine learning this time (Naive Bayes classifier)
- Searches homologous sequences
- Creates multiple sequence alignment
- Predicts how likely two alleles are to occupy a site, given the amino acid substitution pattern from the alignment

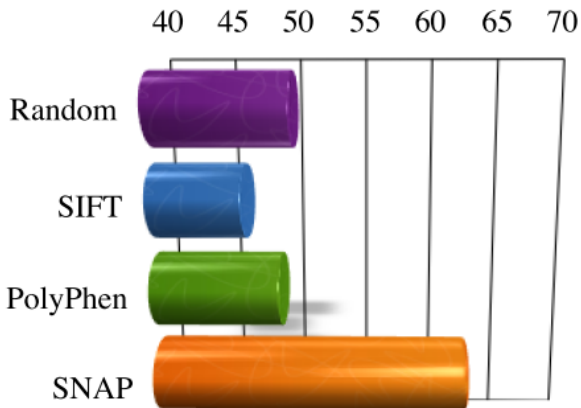
- Also uses other sources of information, like UniProt FT entries, whether or not the mutated position is annotated trans-membrane, signal, active site or binding site for example
- Amino acid substitution calculations based on PSIC score (psi blast)
- According to their paper they perform better than SIFT and SNAP

Screening for **non-acceptable polymorphisms**

- Uses a neural network to predict SNP effects
- Some of the features used:
 - psi-blast frequency profile
 - relative solvent accessibility prediction (PROFace)
 - secondary structure prediction (PROFsec)
 - Pfam information
 - PSIC scores
 - predicted residue flexibility (PROFbval)
 - window around mutation side (five amino acids)
- According to their paper they work better than SIFT and PolyPhen, especially on tough cases



SNAP evaluation on complete dataset.



SNAP evaluation on cases where methods disagree.

Usage example of SIFT, PolyPhen and SNAP.

Thank you for listening!

For further information about the methods (especially SNAP)
please see Protein Prediction 2 Lecture 15 - Slides or Video