

## Substitution matrices

Need mainly two things:

- For every pair *a*, *b*: *p<sub>ab</sub>*, the probability of observing *a* aligned with *b*. *p<sub>ab</sub>* = *p<sub>ba</sub>*
- For every *a*: *p<sub>a</sub>*, the probability of observing an
  *a*.



#### Aligned sequences Related / Unrelated

- Let *M* be the model in which *x* and *y* are related and obtained according to the joint probabilities  $p_{ab}$ .
- Let *R* be the model in which *x* and *y* are unrelated and obtained independently at random according to the individual probabilities  $p_a$ .











- Stands for Point Accepted Mutations.
- An accepted mutation is a mutation that was positively selected by the environment and did not cause the death of the organism.
- Given a PAM matrix M,  $M_{ab} = p[a \rightarrow b]$  in a certain *evolutionary time period*.

# Unit of Evolution

- It is difficult to capture from statistical data the relation of proteins that are evolutionary very far apart. If  $a \rightarrow b$ , we don't capture the intermediate mutations.
- Define 1 unit of evolution as the amount of evolution that will change 1 in 100 amino acids on average.
- Compute the 1-PAM matrix corresponding to 1 unit of evolution from short time interval statistical data.
- Obtain other *k*-PAM matrices from the first one.







# 2-PAM matrix

- *p*<sub>2</sub>[*a*→*b*] in two units of evolution will be the probability of *a* mutating to some character *c* in one unit of evolution and *c* mutating to *b* in another unit of evolution.
- $p_2[a \rightarrow b] = \Sigma_c p[a \rightarrow c] \cdot p[c \rightarrow b] = \Sigma_c M_{ac} \cdot M_{cb}$
- 2-PAM matrix = M<sup>2</sup>



## k-PAM matrix

- k-PAM =  $M^k$
- $S_k(a,b) = 10 \log_{10} (M_{ab}^k/p_b)$



### BLOSUM matrices (BLOCKS substitution matrices)

- BLOSUM matrices are derived from a database of BLOCKS (the BLOCKS database) where each block is a multiple ungapped alignment of related protein sequences.
- The goal is to obtain a scoring for protein sequences that are evolutionary far apart. How far?
- The sequences from each block are clustered, putting two sequences in the same cluster if they have more than L% similarity (percentage of aligned matching residues).
- Distant sequences → occur in different clusters



![](_page_5_Figure_0.jpeg)

![](_page_5_Figure_1.jpeg)