# Patterns, Profiles and Multiple sequence alignments

## Arne Elofsson

To read:
http://en.wikipedia.org/wiki/Multiple_sequence_alignment
http://en.wikipedia.org/wiki/Hidden_Markov_model
Extra
http://en.wikipedia.org/wiki/HMMER
http://en.wikipedia.org/wiki/HHpred_/_HHsearch
http://www.ploscollections.org/article/info%3Adoi%2F10.1371%2Fjournal.pcbi.0030123
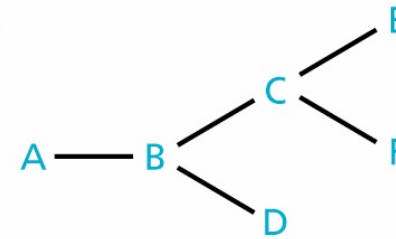http://www.ncbi.nlm.nih.gov/turorials/BLAST

# How to obtain MSAs

- Exact solution is impossible for a handful of sequences (2^N-1 alternatives)

- Popular methods include:

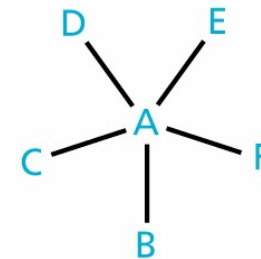  - ClustalW
  - T-coffee
  - kalign
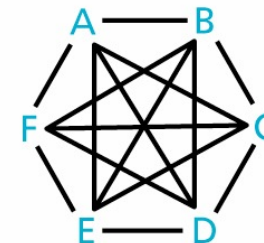  - PSIBLAST

$$2^N - 1$$

# Some scoring in MSAs

(A)



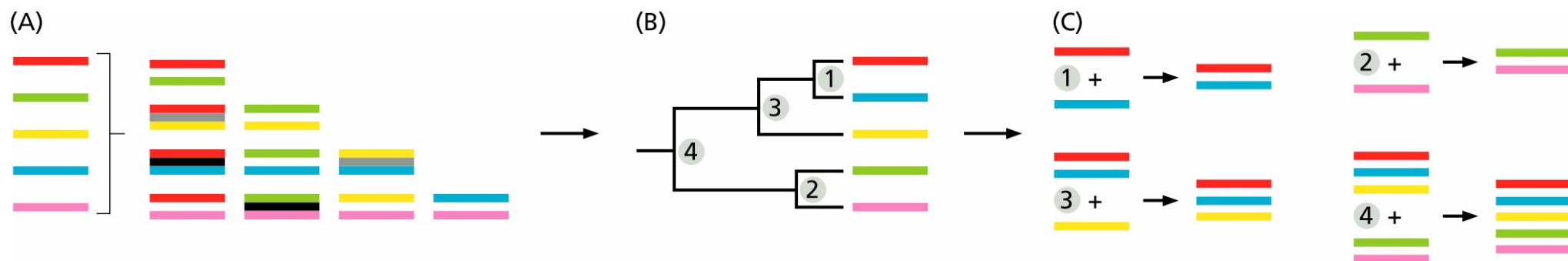$$score = S_{AB} + S_{BC} + S_{BD} + S_{CE} + S_{CF}$$

(B)



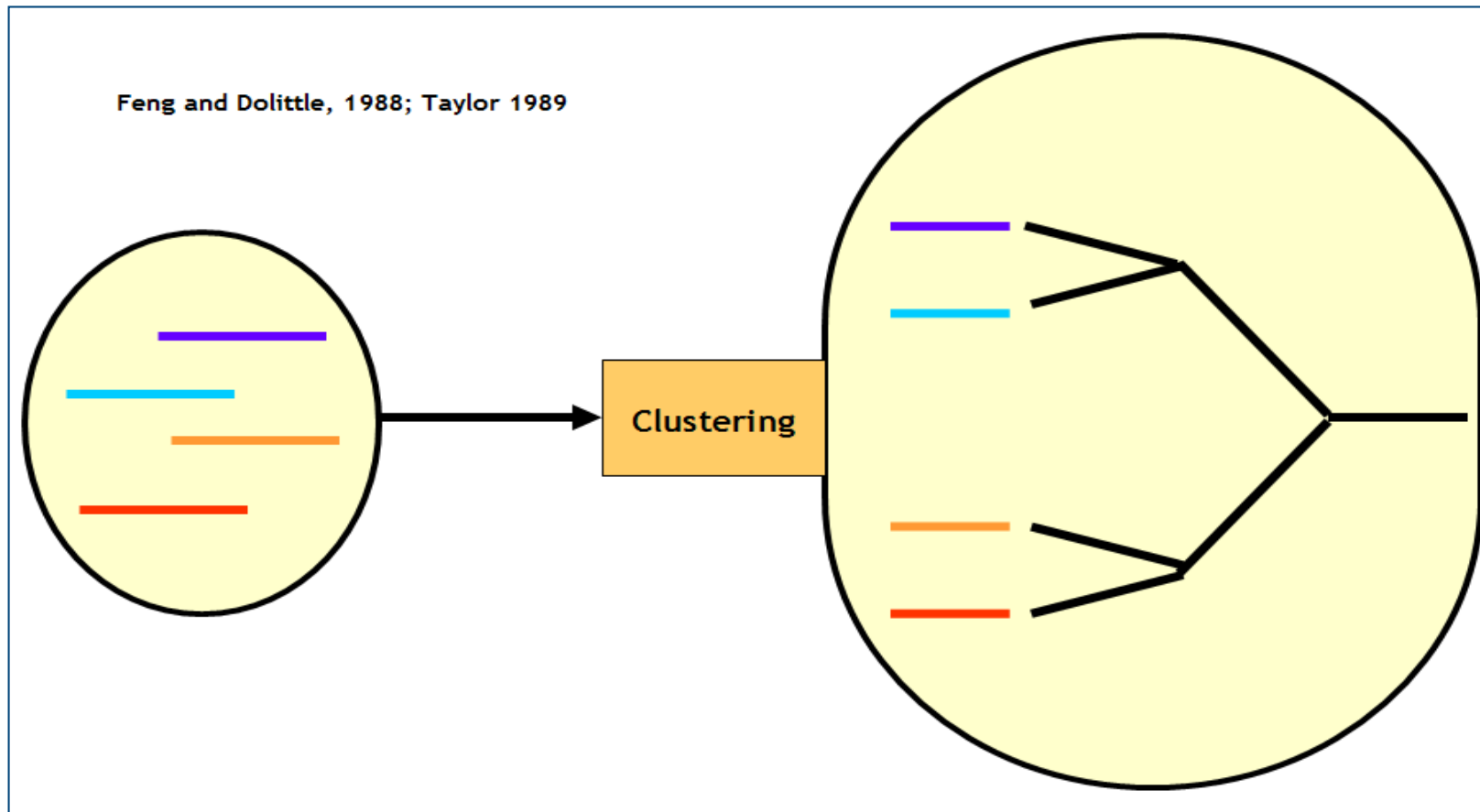$$score = S_{AB} + S_{AC} + S_{AD} + S_{AE} + S_{AF}$$

(C)



$$score = S_{AB} + S_{AC} + S_{AD} + S_{AE} + S_{AF}$$
$$+ S_{BC} + S_{BD} + S_{BE} + S_{BF} + S_{CD}$$
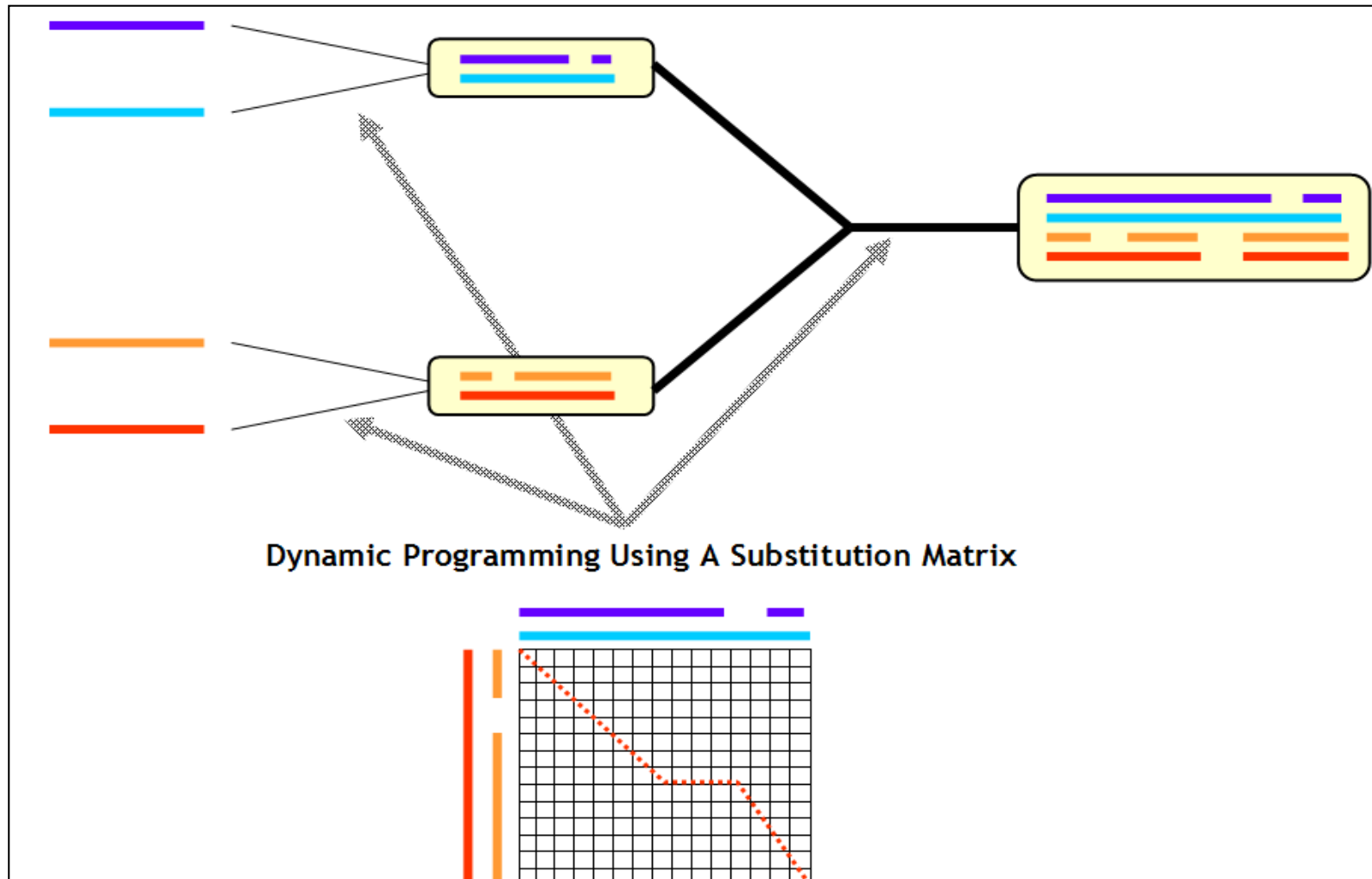$$+ S_{CE} + S_{CF} + S_{DE} + S_{DF} + S_{EF}$$

# A progressive MSA

# ClustalW

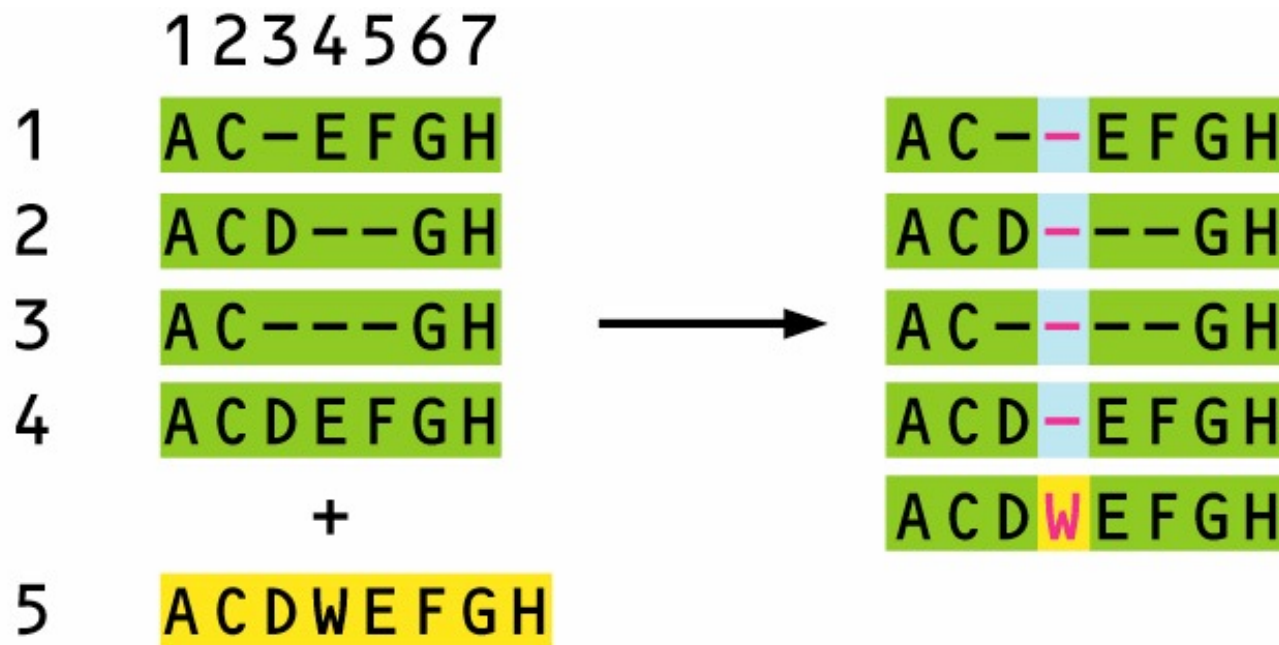

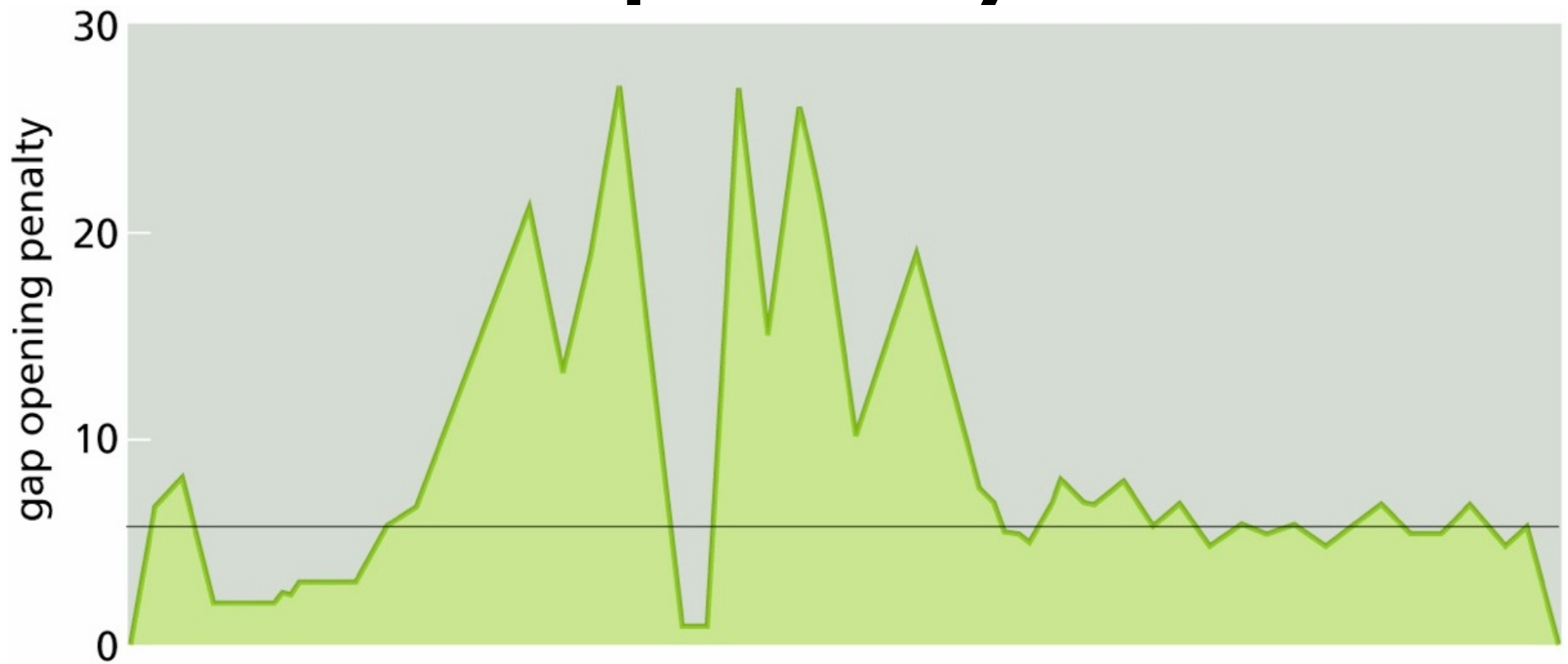Feng and Dolittle, 1988; Taylor 1989

Clustering

# ClustalW



Dynamic Programming Using A Substitution Matrix

# The gap scoring problem

# ClustalW gap-opening penalty

# Multiple sequence alignments
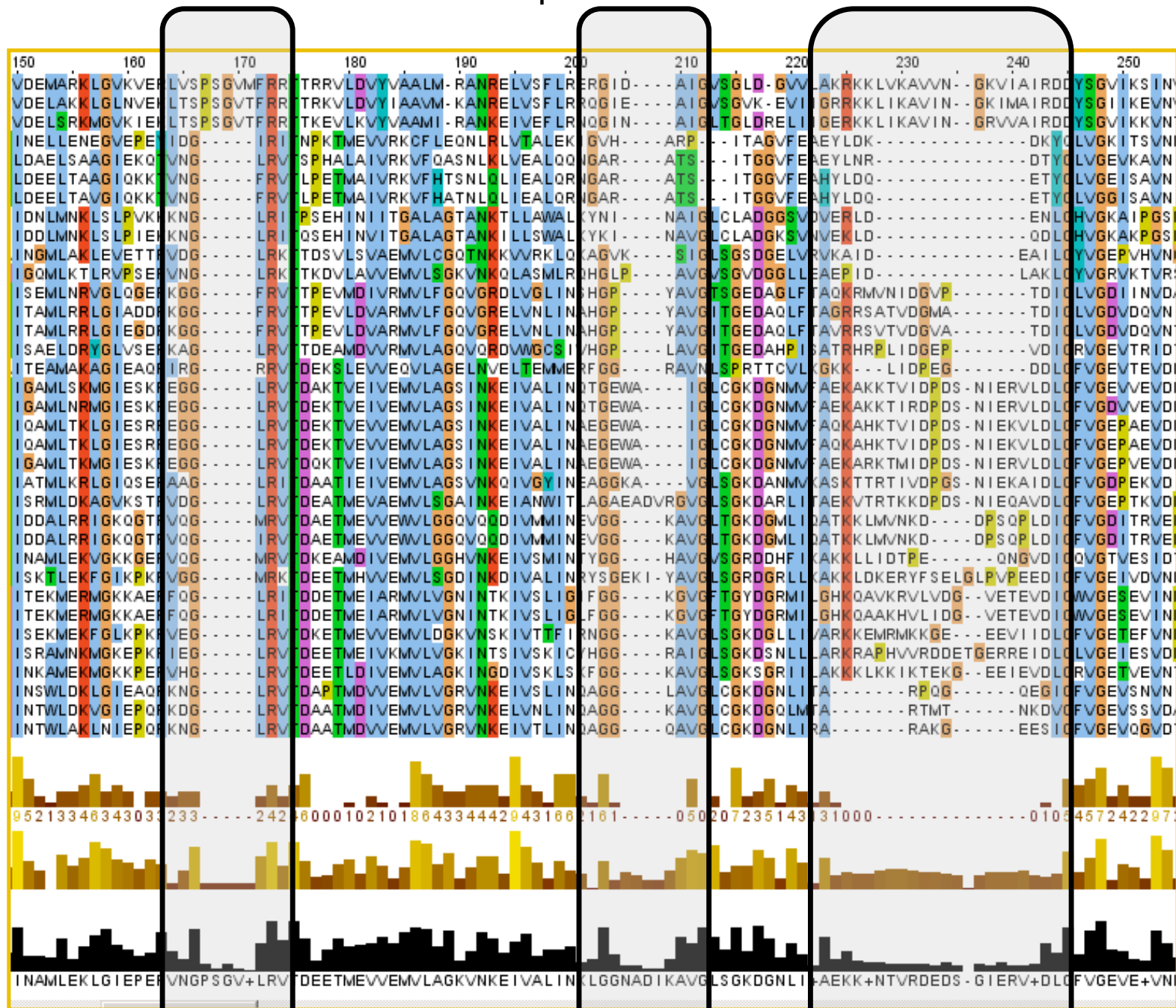
- Some information that can be obtained from a multiple sequence alignment

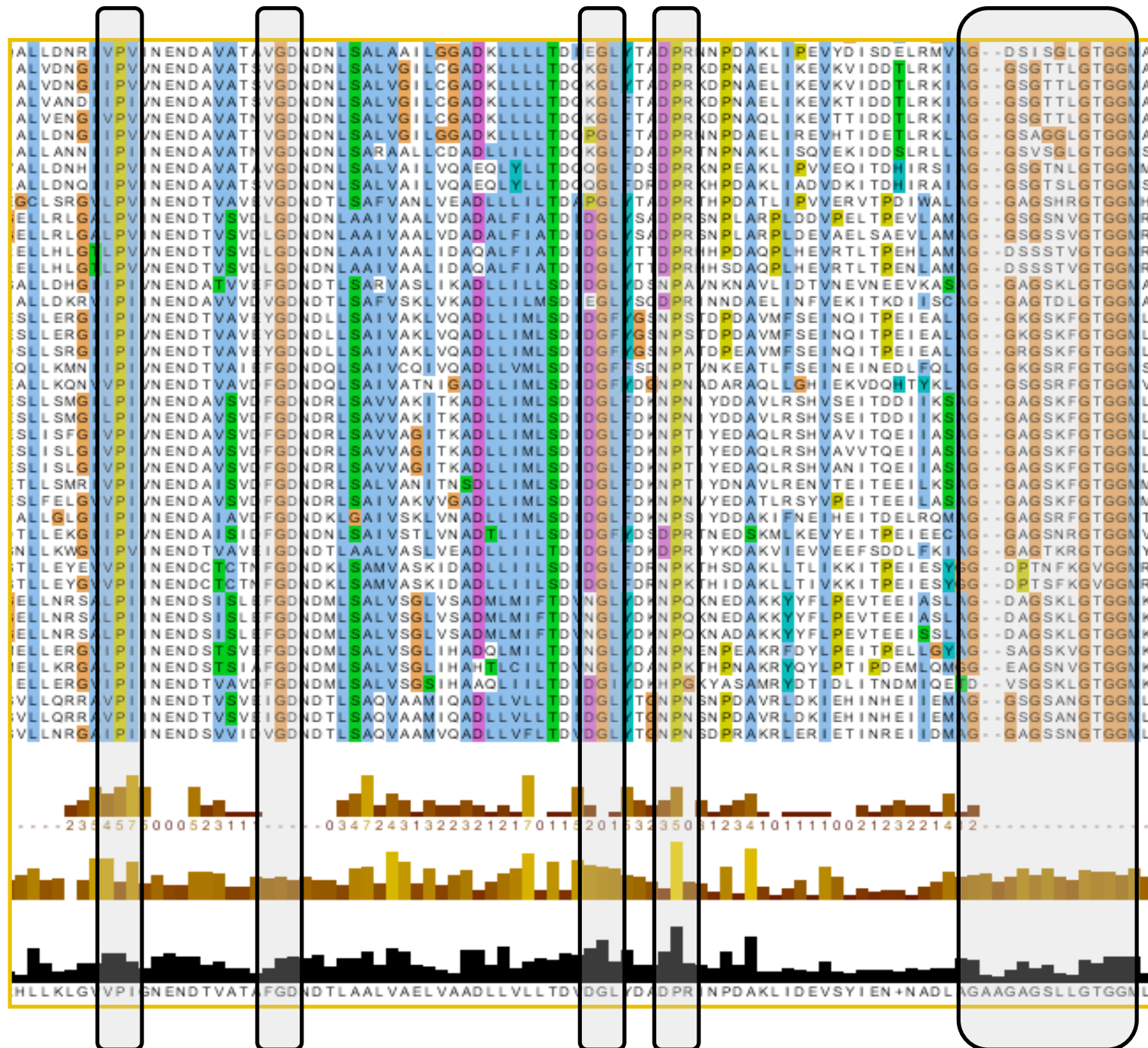# Multiple sequence alignments
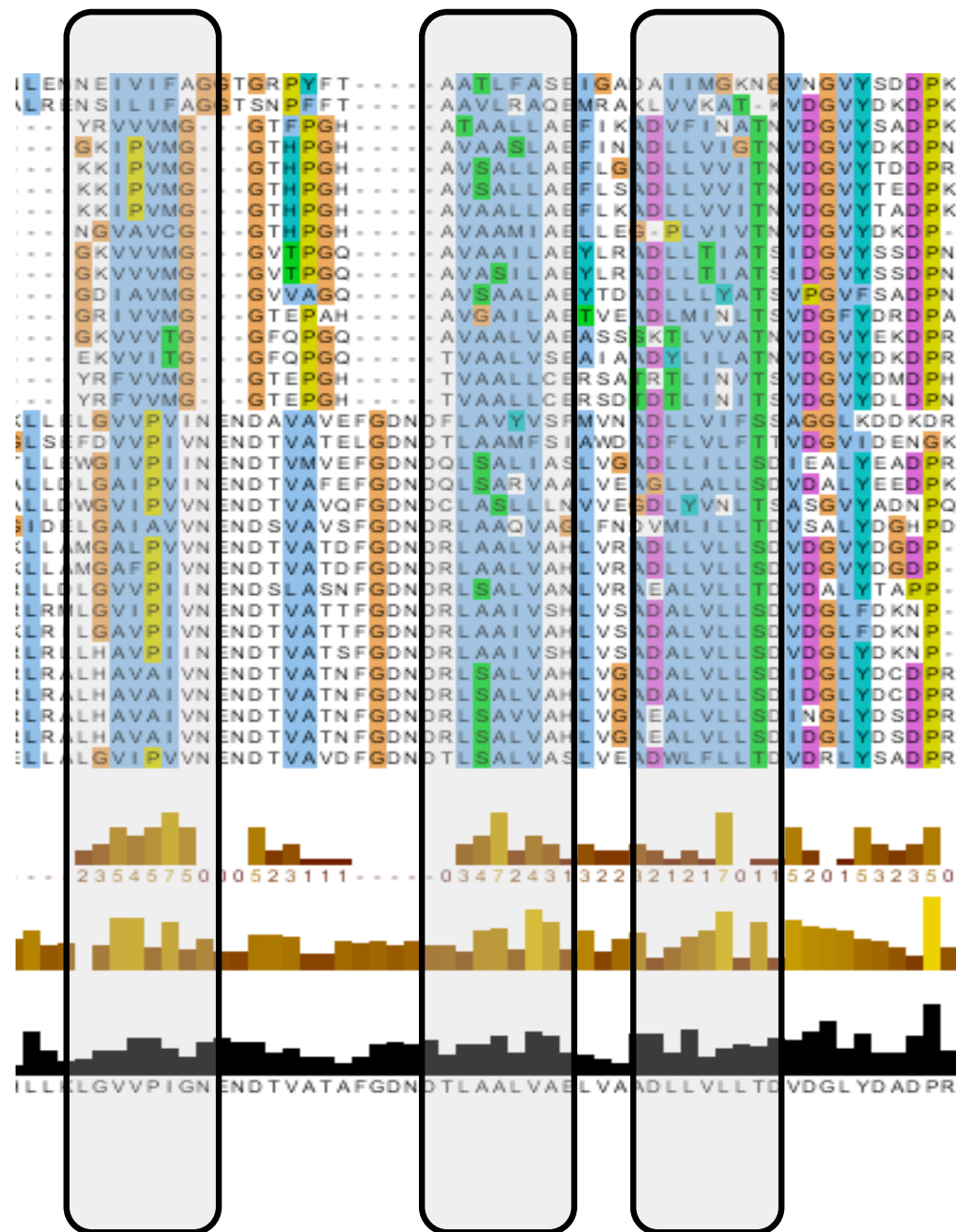## (good ones look pretty !!)

# Features found in MSA:

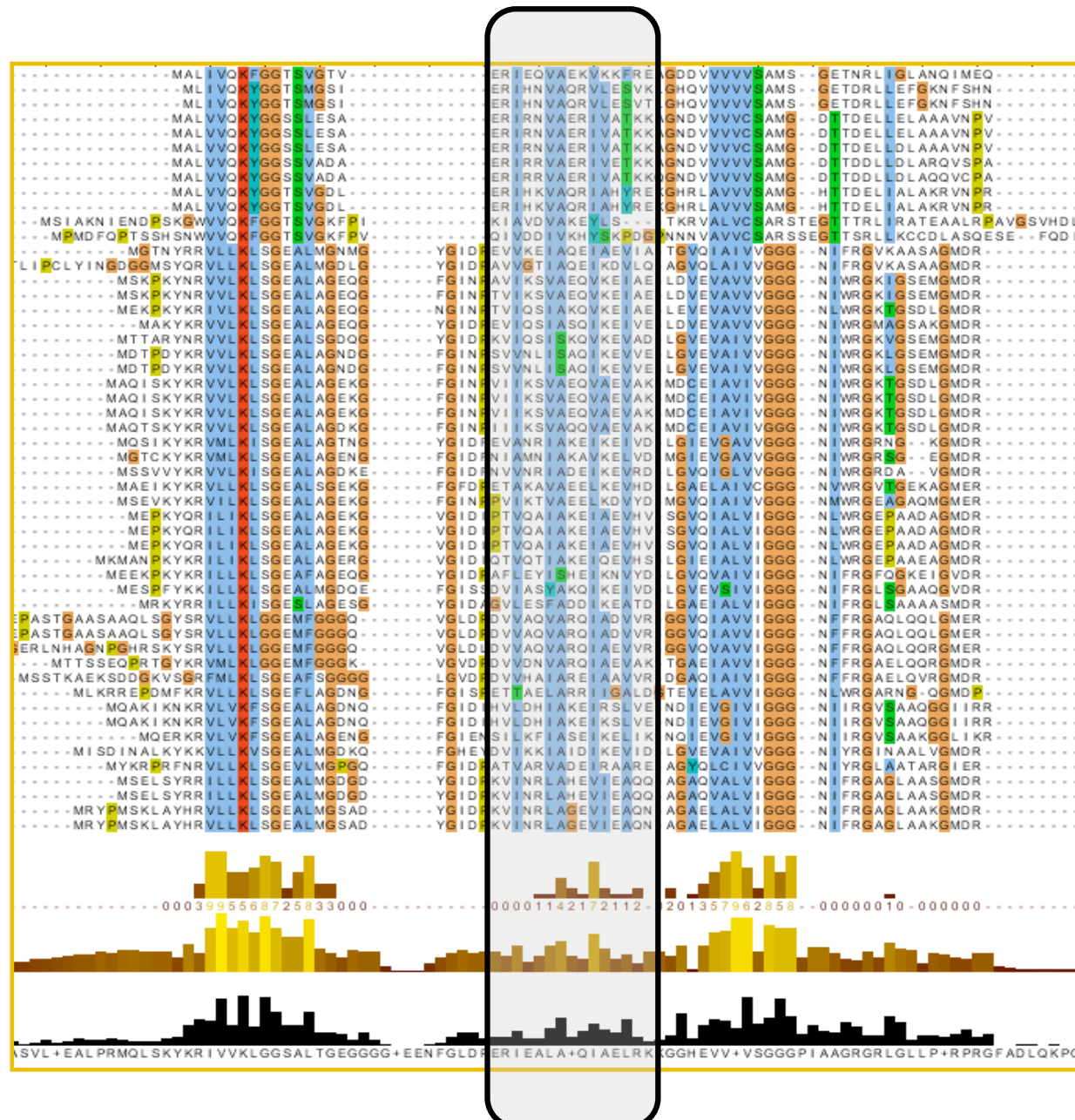The position of insertions and deletions suggests regions where surface loops exist…

Conserved glycine or proline suggests a β-turn.

Residues with hydrophobic properties conserved at i, i+2, i+4 (etc) separated by unconserved or hydrophilic residues suggests a surface β-strand…

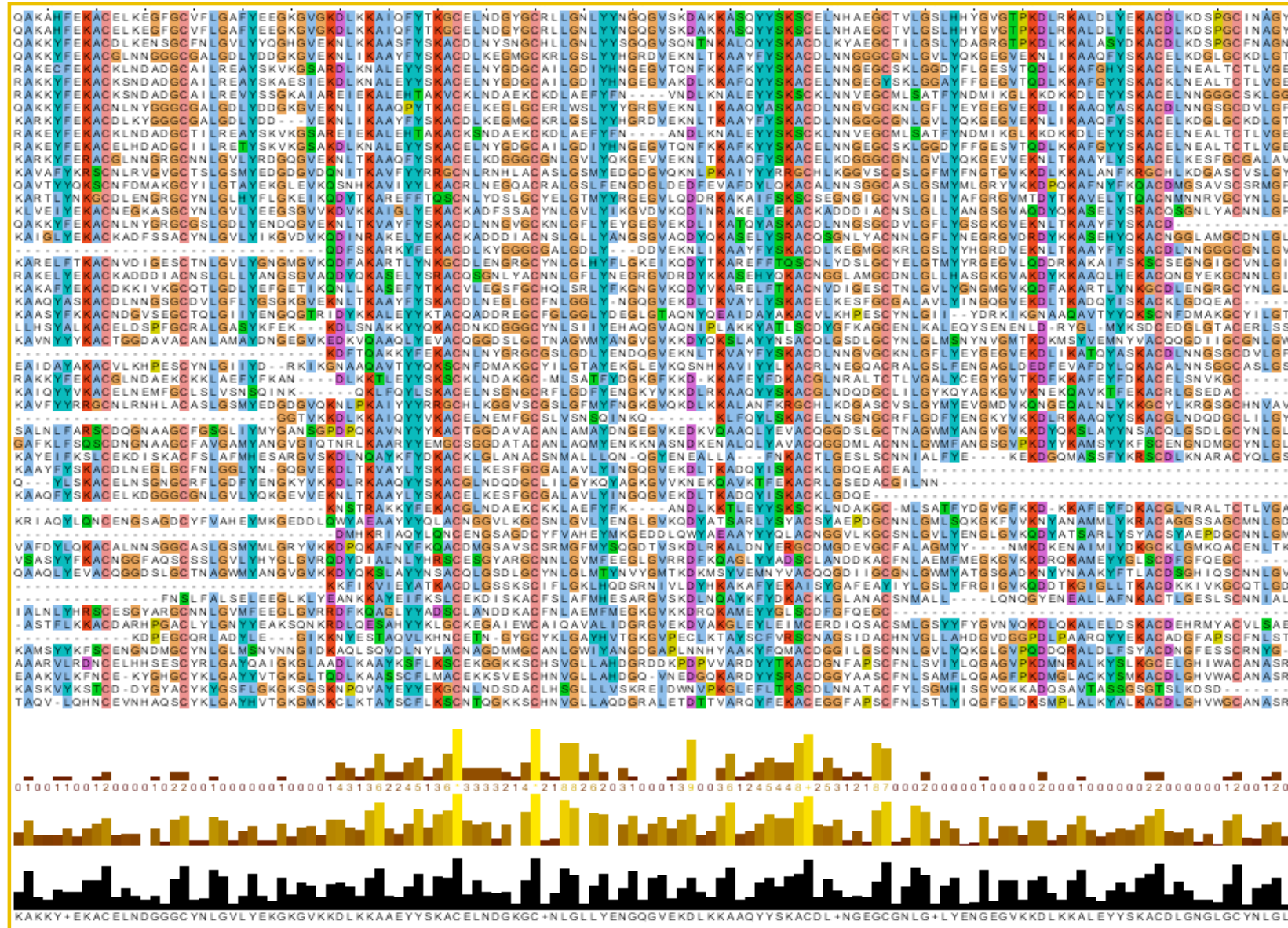A short run of hydrophobic amino acids (4 or 5 residues) suggests a buried β-strand…

Pairs of conserved hydrophobic amino acids separated by pairs of unconserved or hydrophilic residues suggests an α-helix with one face packed in the protein core. Similarly, an i, i+3, i+4, i+7 pattern of conserved residues."

Cysteine is a rare amino acid, and is often used in disulphide bonds ( pairs of conserved cysteines )

Charged residues ( histidine, aspartate, glutamate, lysine, arginine ) and other polar residues embedded in a conserved region indicate functional importance

# Coloring your alignments

# Coloring your alignment

Alignment of 27 avian influenza hemagglutinin protein sequences colored by residue conservation (top) and residue properties (bottom)

# PSSMs

- How to use evolution to detect (more) homologs ?

  - Iterated sequence search

  - Patterns

  - Profiles (PSSMs, PSIBLAST)

  - HMMs

# Searching with PSSMs

- PSSMs are profiles

  - without gap information

  - Without substitution table information

- Dynamic programming can be used

  - Identical algorithm as for Smith-Waterman

- Profile alignment vs. single sequence

  - Better alignments

  - Better detection

# Profile-sequence alignment



sequence

ACD......VWY

# Average profiles

- Gribskov, McLachlan and Eisenberg 1987

  - No underlying probabilistic model, but rather assigned position specific scores for each match state and gap penalty

  - The score for each consensus position is set to the average of the standard substitution scores from all the residues in the corresponding multiple sequence alignment column

  - Gap costs

# The "average" profile method

- Score for each residue is average score for that residue with all sequence in MSA

- Average score over all replacements:

# Non-probabilistic or

```
HBA_HUMAN     ...VGA--HAGEY...
HBB_HUMAN     ...V----NVDEV...
MYG_PHYCA     ...VEA--DVAGH...
GLB3_CHITP    ...VKG-------D...
GLB5_PETMA    ...VYS--TYETS...
LGB2_LUPLU    ...FNA--NIPKH...
GLB1_GLYDI    ...IAGADNGAGV...
              ***   *****
```

The score for residue 'a' in column 1

$$\frac{5}{7}s(V,a) + \frac{1}{7}s(F,a) + \frac{1}{7}s(I,a)$$

**s(a,b) : standard substitution matrix**

# Average profiles – example

- One position contains
  - 50% ILE ; 30% THR ; 20% VAL
- Calculate the score for ILE in this position
- Use the PAM250 Matrix
  - I-I=5
  - I-T=0
  - I-V=4
- Calculate
  - 0.5*5+0.3*0+0.2*4=3.3
- Integer

# Average Profiles

- They also set gap penalties for each column using a heuristic equation that decrease the cost of a gap according to the length of the longest gap observed in the multiple alignment spanning the column

# Problem With Average profiles

- If we had an alignment with 100 sequences, all with a cysteine (C), at some position, the probability distribution for that column for an "average" profile would be exactly the same as would be derived from a single sequence

- Doesn't correspond to our expectation that the likelihood of a cysteine should go up as we see more confirming examples

# Similar Problem With Gaps

```
HBA_HUMAN        ...VGA--HAGEY...
HBB_HUMAN        ...V----NVDEV...
MYG_PHYCA        ...VEA--DVAGH...
GLB3_CHITP       ...VKG------D...
GLB5_PETMA       ...VYS--TYETS...
LGB2_LUPLU       ...FNA--NIPKH...
GLB1_GLYDI       ...IAGADNGAGV...
                    ***   *****
```

Scores for a deletion in columns 2 and 4 would be set to the same value

More reasonable to set the probability of a new gap opening to be higher in column 4

# The amino acid frequency can be used for scoring

$$f_{u,b} = \frac{n_{u,b}}{N_{seq}} \, (EQ6.1)$$

$$m_{u,a} = \sum_b f_{u,b} s_{a,b} \, (EQ6.2)$$

# Higher scoring for more conserved positions

$$m_{u,a} = \sum_b \frac{\ln(1 - f'_{u,b})}{\ln(1/(N_{seq} + 1))} s_{a,b} \, (EQ6.3)$$

$$m_{u,a} = \log \frac{q_{u,a}}{p_a} \, (EQ6.4)$$

# PSI-BLAST algorithm

- Input a single protein sequence and compares it to a protein database, using BLAST

- The program constructs a multiple alignment, and then a profile,

- The profile is compared to the protein database, again seeking local alignments.

- PSI-BLAST estimates the statistical significance of the local alignments found.

- Finally, PSI-BLAST iterates, by returning to step (2), an arbitrary number of times or until convergence.

# Psiblast

# PSI-BLAST

- Advantages
  - Fast (40 times faster than DP)
  - Significant better than DP
  - Good E-value estimates
- Disadvantages
  - Not optimal alignments

# PSI-BLAST

- Important parameters
  - E-value cutoff
  - Number of iterations
  - Low complexity sequence filtering

# PSI-BLAST in a nutshell

- With a protein sequence as query, use BLAST to search a protein sequence database.

- Collapse significant local alignments (those with $E$-value less than or equal to a set threshold $h$) into a multiple alignment, using the residue of the query sequence as alignment-column placeholders.

- Abstract a position-specific score matrix from the multiple alignment.

- Search the database with the score matrix as query.

- Iterate a fixed number of times, or until convergence.

# PSI-BLAST live example

Sequence;

SISSRVKSVLLLGLQNAELAQKVGTTQQSIEQLENGKTRPRFLPELASAILGVSVDWLLNGT

Server:

http://www.ncbi.nlm.nih.gov/blast/

Run against Swissprot (faster)

# Markov Chains

Rain



Sunny                    Cloudy

**States** : Three states - sunny, cloudy, rainy.

weather today

|  |  | Sun | Cloud | Rain |
|---|---|---|---|---|
| weather yesterday | Sun | 0.5 | 0.25 | 0.25 |
|  | Cloud | 0.375 | 0.125 | 0.375 |
|  | Rain | 0.125 | 0.625 | 0.375 |

**State transition matrix** : The probability of the weather given the previous day's weather.

| Sun | Cloud | Rain |
|---|---|---|
| 1.0 | 0.0 | 0.0 |

**Initial Distribution** : Defining the probability of the system being in each of the states at time 0.

# Hidden Markov Models



**Hidden states** : the (TRUE) states of a system that may be described by a Markov process (e.g., the weather).

**Observable states** : the states of the process that are `visible' (e.g., seaweed dampness).

# Components of HMM

Seaweed

|  | Dry | Dryish | Damp | Soggy |
|---|---|---|---|---|
| Sun | 0.60 | 0.20 | 0.15 | 0.05 |
| Cloud | 0.25 | 0.25 | 0.25 | 0.25 |
| Rain | 0.05 | 0.10 | 0.35 | 0.50 |

weather

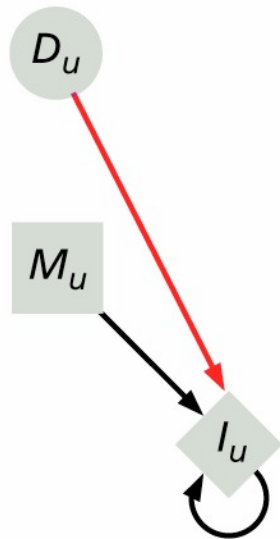**Output matrix** : containing the probability of observing a particular observable state given that the hidden model is in a particular hidden state.

**Initial Distribution** : contains the probability of the (hidden) model being in a particular hidden state at time t = 1.
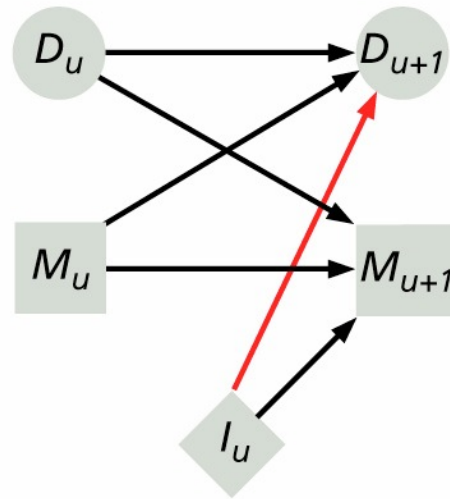
**State transition matrix** : holding the probability of a hidden state given the previous hidden state.
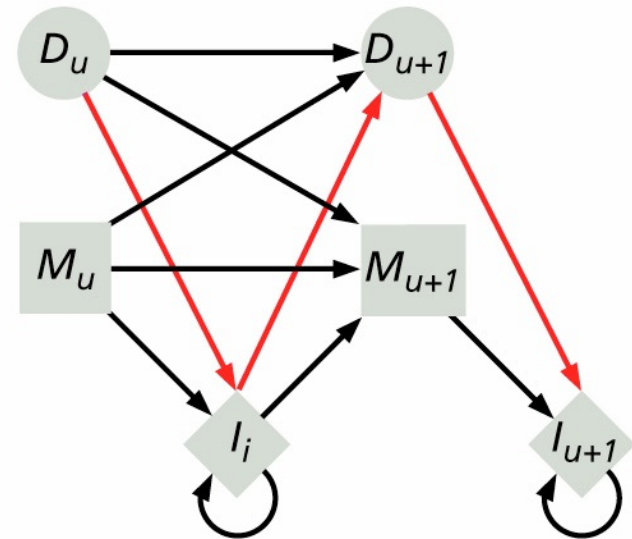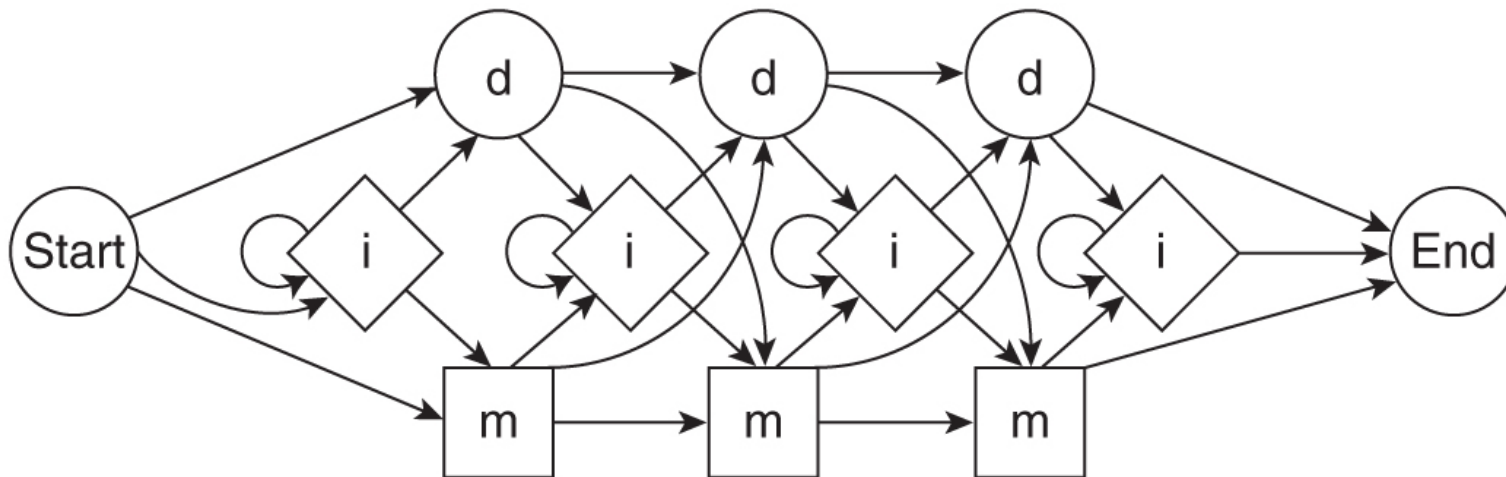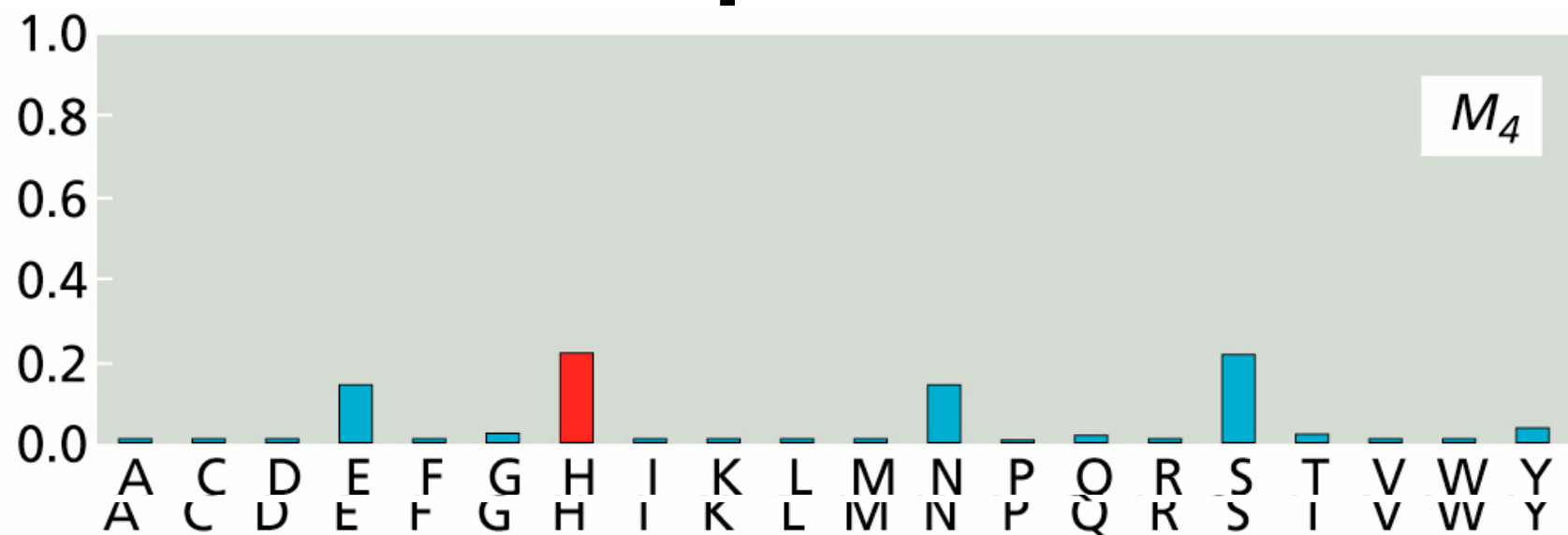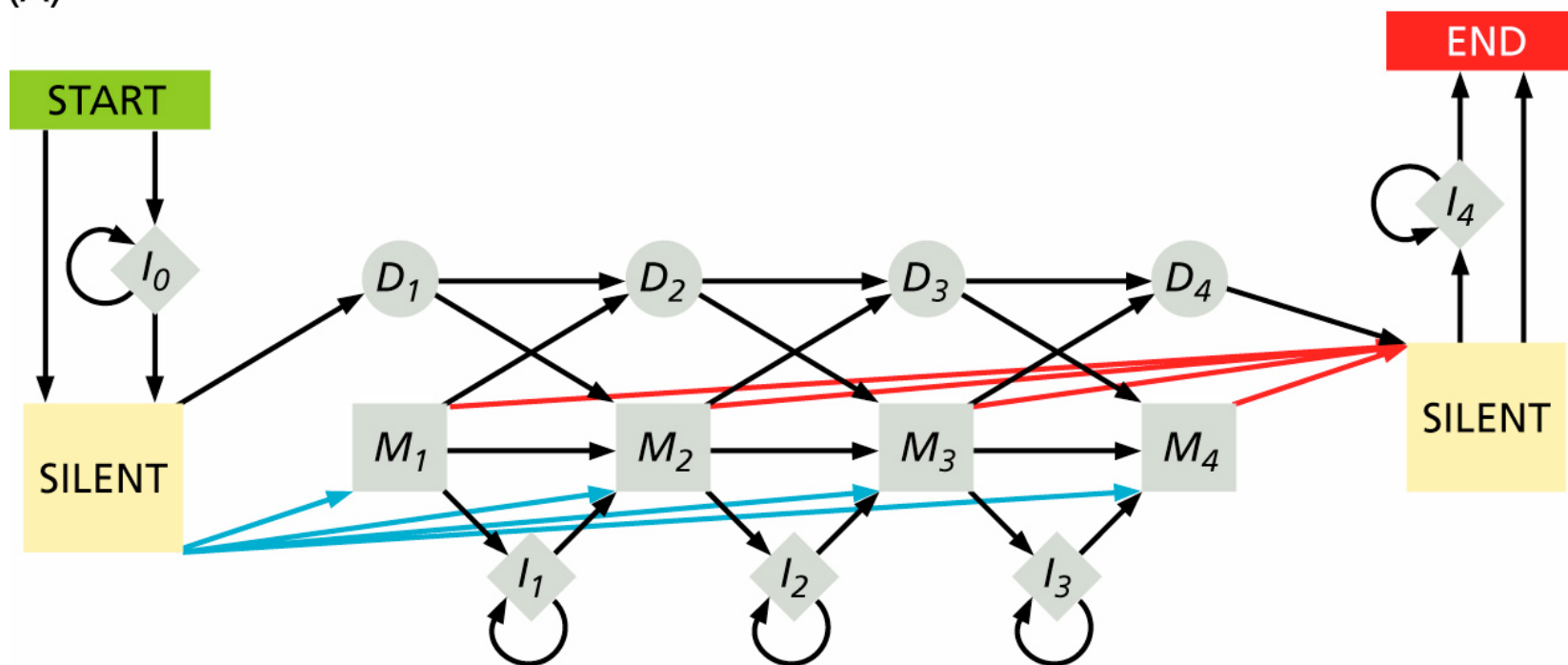
# Profile HMMs

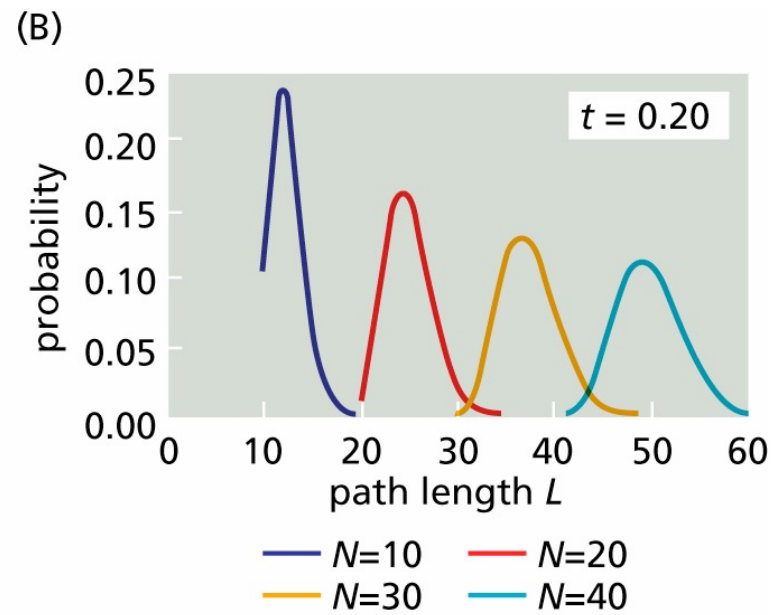# A profile HMM
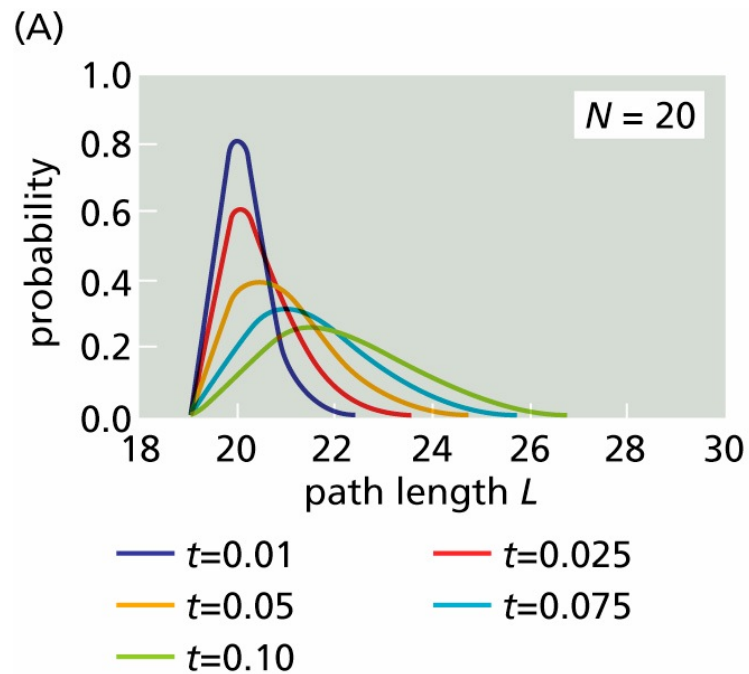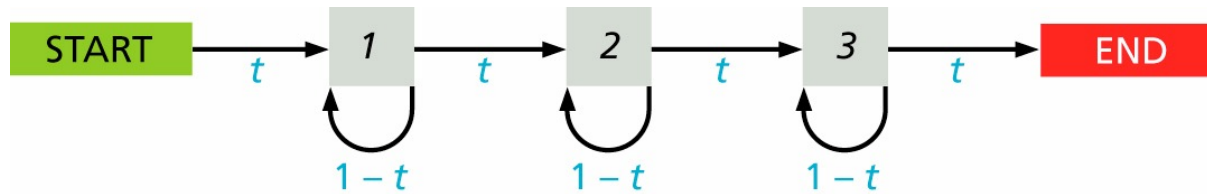
# Emission probabilities

# A local HMM



(A)

# Length dependency of HMMs

# Scoring an HMM

- The most probable path (Viterbi)

- Scoring for all paths (forward/backward)

# Training an HMM

- Using unaligned sequences

- Baum-Welch expectations maximization

  - Estimating the number each emission and transition is used

- Starting with rough estimates

# Profile HMMs: Effectiveness

- Advantages:

  - Expressive profiling method

  - Transparent method: You can view and interpret the model produced

  - Very effective at detecting remote homologs

- Disadvantages:

  - Slow – full search on a database of 400,000 sequences can take 15 hours

  - Have to avoid over-fitting and locally optimal models