# CMSC423: Bioinformatic Algorithms, Databases and Tools
# Lecture 13

multiple alignment

motif finding

# Recap

- Multiple alignment is expensive – $O(n^k)$ for k sequences of length n (use same DP as for pairwise but on a k-dimensional matrix)

- Approximation algorithm (star alignment) can find a solution in $O(n^2 k^2)$ which is at most twice worse than the best alignment

# Consensus sequence

- For every column j in the alignment, pick the amino-acid AA that minimizes $\sum_i d(AA, S_i[j])$ (usually becomes majority rule)

- Intuitively – this is the sequence of the ancestor of all the sequences in the multiple alignment

- We can define the multiple alignment problem as:
  - find the multiple alignment that minimizes $\sum_i D(CO, S_i)$

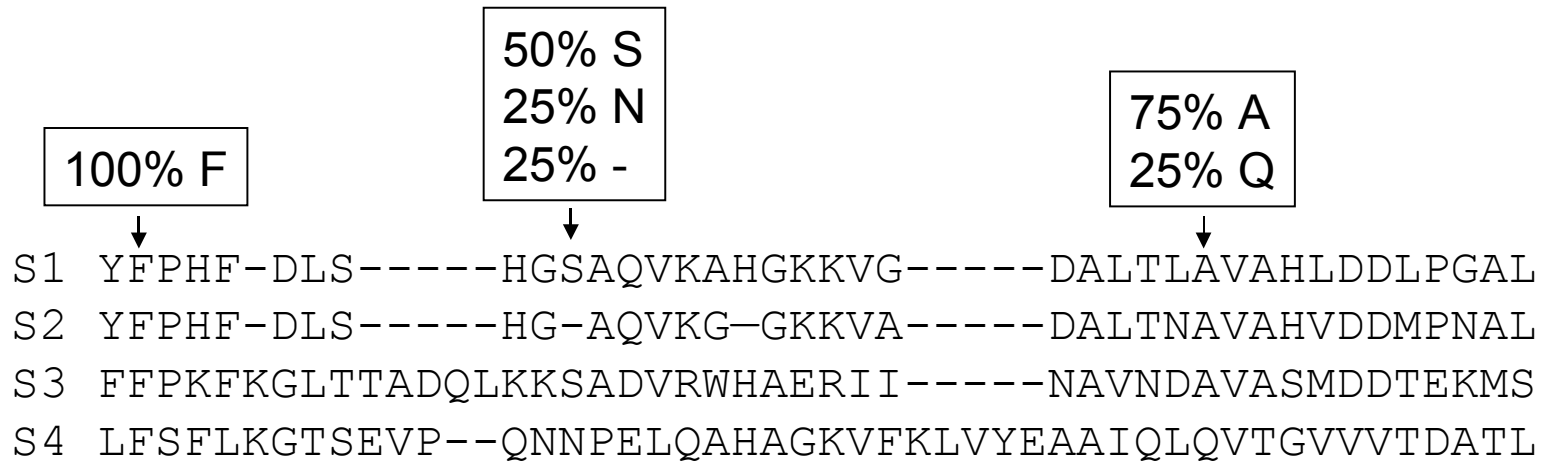- Still NP – hard, but consensus sequence useful on it's own.

```
CO  YFPHFKDLS-----HGSAQVKAHGKKVG-----DALTLAVAHVDDTPGAL
S1  YFPHF-DLS-----HGSAQVKAHGKKVG-----DALTLAVAHLDDLPGAL
S2  YFPHF-DLS-----HG-AQVKG—GKKVA-----DALTNAVAHVDDMPNAL
S3  FFPKFKGLTTADQLKKSADVRWHAERII-----NAVNDAVASMDDTEKMS
S4  LFSFLKGTSEVP--QNNPELQAHAGKVFKLVYEAAIQLQVTGVVVTDATL
```

# Iterative alignment revisited

- Pick a sequence (e.g. SC) as a starting point
- Align S1 to it & build consensus for the alignment
- Take S2 and align it to the consensus (instead of SC)
- repeat...
- Problem: consensus (or any single sequence) ignores the other sequences being aligned.
- Solution: keep track of % of each amino-acid aligned in each column
- score of alignment to profile – combination of scores to each AA.

# Profile alignment

- Solution: keep track of % of each amino-acid aligned in each column

- score of alignment to profile – combination of scores to each AA.

```
       100% F         50% S                    75% A
                      25% N                    25% Q
                      25% -
        ↓              ↓                         ↓
S1  YFPHF-DLS-----HGSAQVKAHGKKVG-----DALTLAVAHLDDLPGAL
S2  YFPHF-DLS-----HG-AQVKG—GKKVA-----DALTNAVAHVDDMPNAL
S3  FFPKFKGLTTADQLKKSADVRWHAERII-----NAVNDAVASMDDTEKMS
S4  LFSFLKGTSEVP--QNNPELQAHAGKVFKLVYEAAIQLQVTGVVVTDATL
```
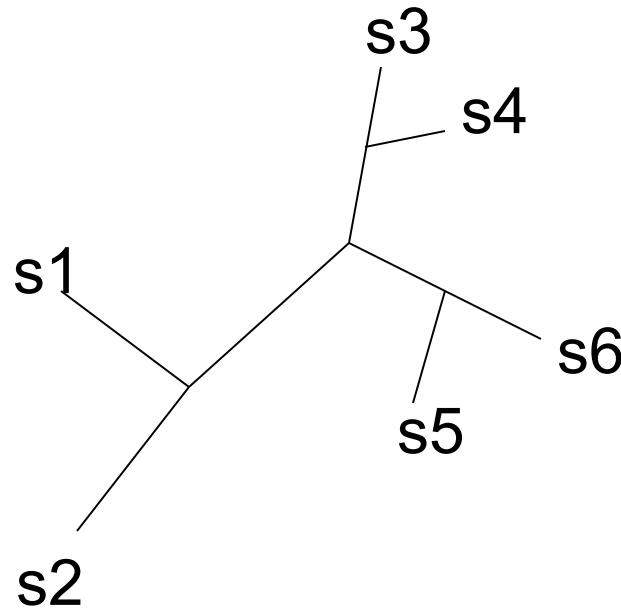
- Score(prof1, prof2) = weighted average of all pairs of amino-acids

# Practical algorithms

# Iterative alignment

```
SC YFPHFDLSHGSAQVKAHGKKVGDALTLAVGHLDDLPGAL
```

- ## Take sequences si in order:

  - align s1 with sc - results in gaps being inserted in both sequences

    ```
    SC YFPHFDLSHGSAQVKAHGKKVGDALTLAVGHLDDLPGAL
    S1 YFPHFDLSHG-AQVKG--KKVADALTNAVAHVDDMPNAL
    ```

  - align s2 with sc - if gaps must be inserted – insert in previously aligned sequences

    ```
    SC YFPHF-DLS-----HGSAQVKAHGKKVG-----DALTLAVAHLDDLPGAL
    S1 YFPHF-DLS-----HG-AQVKG—GKKVA-----DALTNAVAHVDDMPNAL
    S2 FFPKFKGLTTADQLKKSADVRWHAERII-----NAVNDAVASMDDTEKMS
    ```

  - and so on (note: if gaps coincide with previously introduced gaps no need to change previously aligned sequences)

    ```
    SC YFPHF-DLS-----HGSAQVKAHGKKVG-----DALTLAVAHLDDLPGAL
    S1 YFPHF-DLS-----HG-AQVKG—GKKVA-----DALTNAVAHVDDMPNAL
    S2 FFPKFKGLTTADQLKKSADVRWHAERII-----NAVNDAVASMDDTEKMS
    S3 LFSFLKGTSEVP--QNNPELQAHAGKVFKLVYEAAIQLQVTGVVVTDATL
    ```

# CLUSTALW

- Compute pairwise distances between strings
- Build phylogenetic tree
- Build iterative alignment by following tree edges

# MUSCLE

- Just like ClustalW but different
- Build pairwise distances – uses fast heuristic (just count # of k-mers in common)
- Build phylogenetic tree
- Build multiple alignment based on tree
- Re-estimate distances based on tree
- Re-build tree
- Re-build multiple alignment
- etc. etc. etc.

# Biological relevance of multiple alignments

# Motif finding

# Motif finding

- Special case of multiple alignment – find short "motif" that occurs almost identically in multiple DNA sequences

- Local multiple alignment (the definition of multiple alignment sofar was global)

- Motif finding – special requirements
  - inexact alignment sought
  - but no gaps allowed

- Biological significance
  - gene promoters
  - transcription factor binding sites
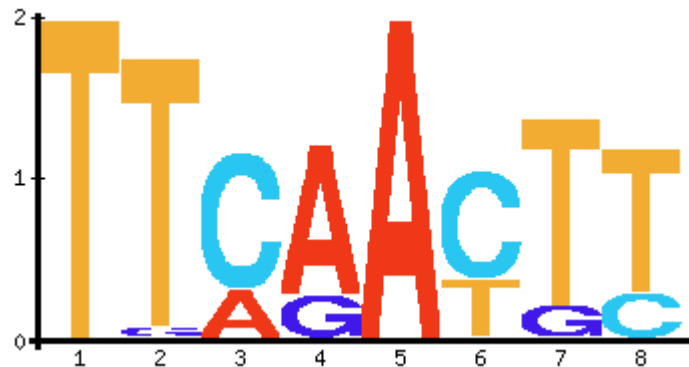  - other elements involved in gene regulation

# Motif finding...example

TTAGAGGTTGACTA**TTCAACTT**TTGAGGAGGCCTAG*TAGAGC*
AGCCGACT**TGCAACTT**AGGCGTGGTCAGTGCCCTAA*TAGAGC*
GGCCTATTTGGGCCACTTAGACC**TTCAACTT**TTGCA*TAGAGC*
CCACAG**TTAGATGT**CCAAAAGACAAATATAGAGGGC*TAGAGC*
ACACGGACTGCG**TTCAATGC**TTACAGCAGATTGAGT*TAGAGC*
TTCAAAGACTTGACTATTG**TTCAACTT**TGAAGACTA*TAGAGC*

Promoter region                                    Gene

Motif "sequence logo"

From genetics.mgh.harvard.edu/sheenlab/

# Finding motifs – Gibbs sampling

- Observations:
  - since no gaps – all motifs have equal length (assume known value - m)
  - exhaustive search of promoter region is impractical: all combinations of substrings of length m among k sequences of length L = $(L - m + 1)^k$

- Solution: random search

1. Pick random substring of length m from each of the strings

2. Construct multiple alignment (easy since no gaps) and compute profile

3. Pick random sequence s and remove from multiple alignment. Recompute profile.

4. Within removed sequence, search for best fit to profile and insert into alignment

5. Repeat until profile does not improve
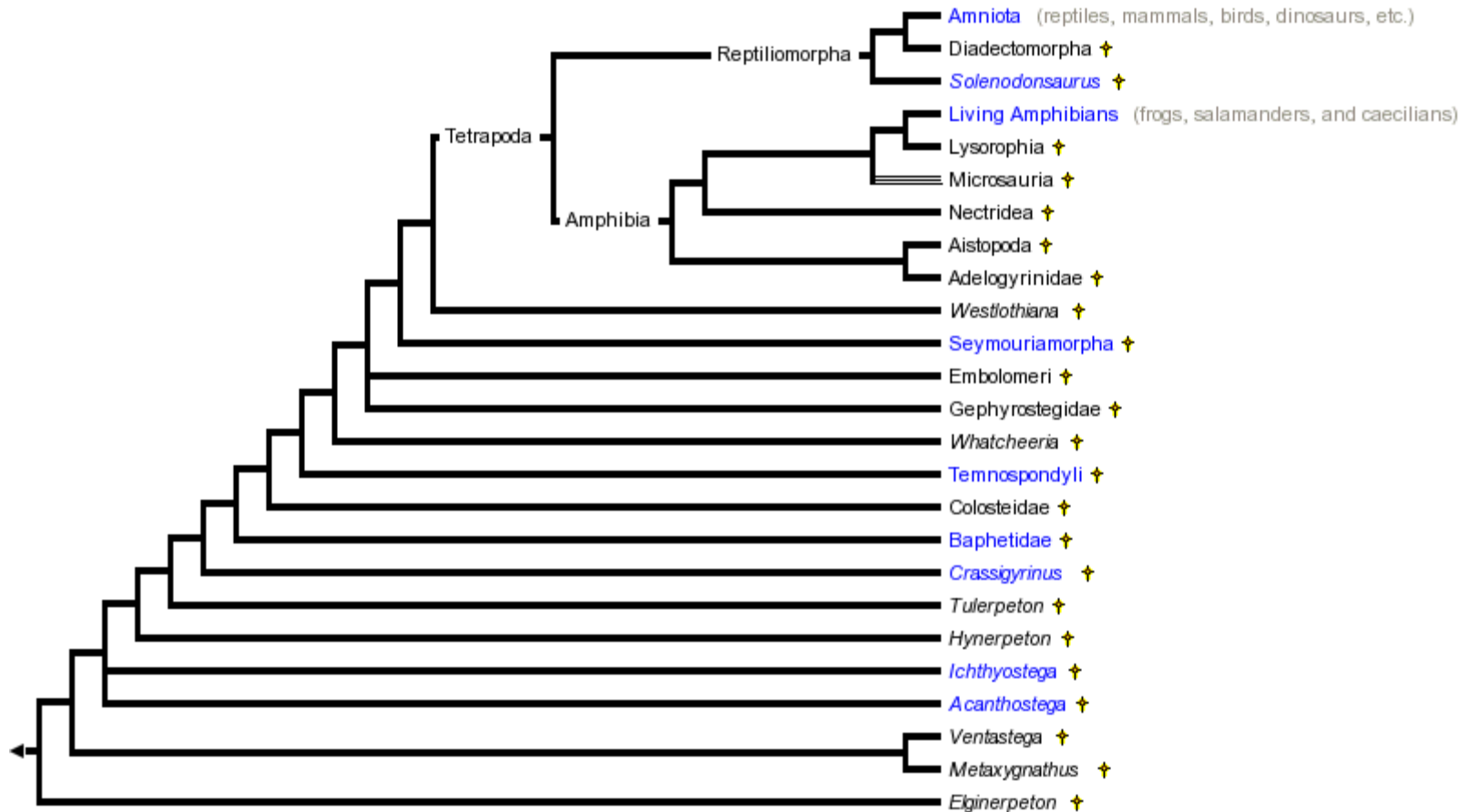
# Gibbs sampling...cont

- How do you find best match to profile?
- What is overall running time of algorithm?
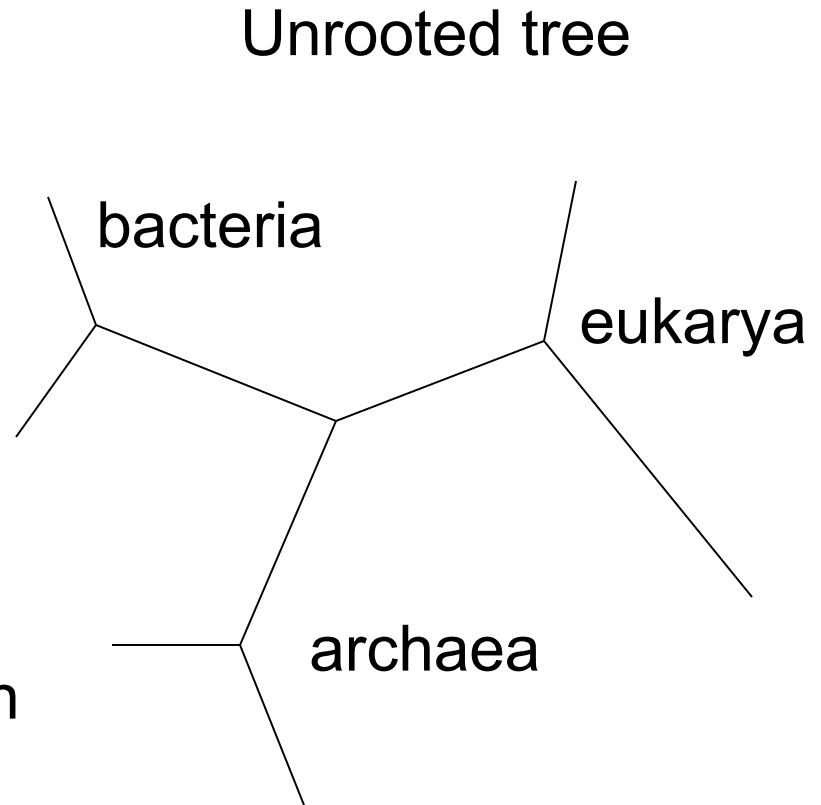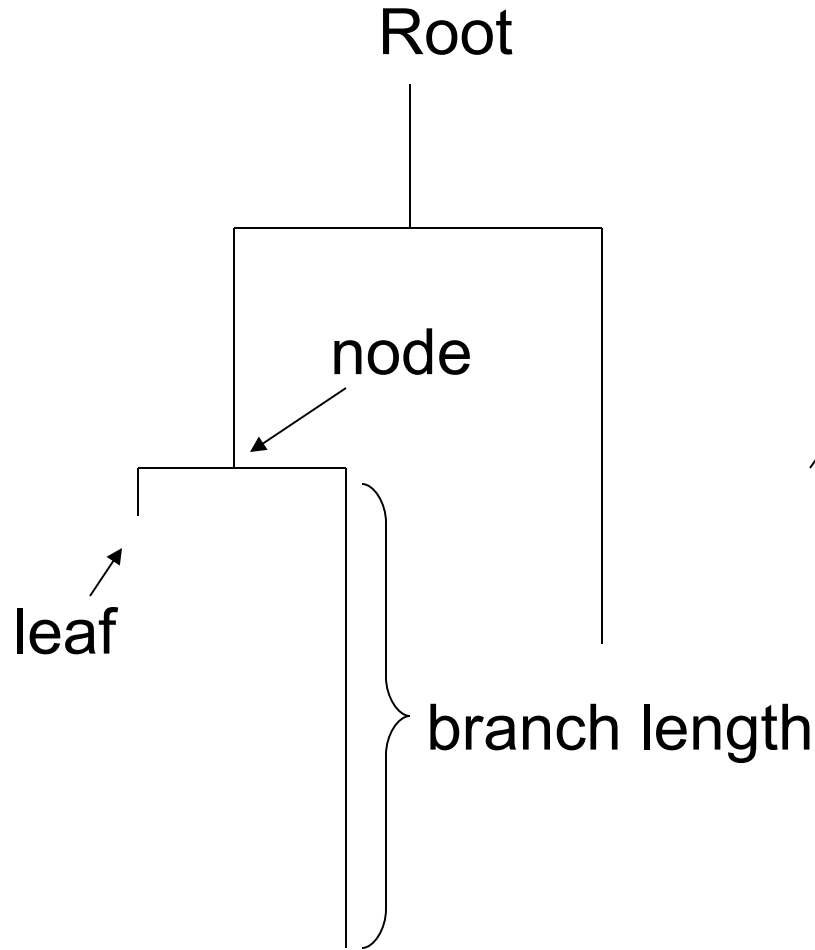
# Phylogenetic trees

# Phylogenetic trees – how evolution works

- http://www.tolweb.org/tree/ - the tree of life

# Anatomy of a tree

Root

Unrooted tree

node

bacteria

eukarya

leaf

archaea

branch length

Phylogenetic trees are usually binary (though they don't have to)

# Phylogeny questions

- Given several organisms & a set of features (usually sequence, but also morphological: wing shape/color...)

- A. Given a phylogenetic tree – figure out what the ancestors looked like (what are the features of internal nodes)

?  wings, feathers, teeth
claws, no wings, fur

- B. Find the phylogenetic tree that best describes the common evolutionary heritage of the organisms

A
B
C

C
A
B

B
A
C