# CMSC423: Bioinformatic Algorithms, Databases and Tools
## Lecture 10

Sequence alignment: inexact
alignment, multiple sequence
alignment

---

# 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
      ```
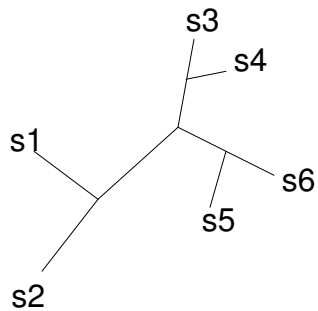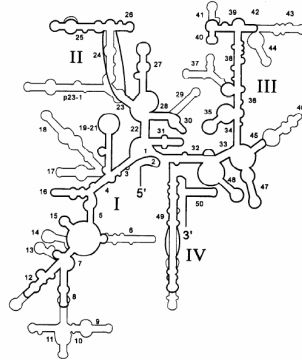
# CLUSTALW

- 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
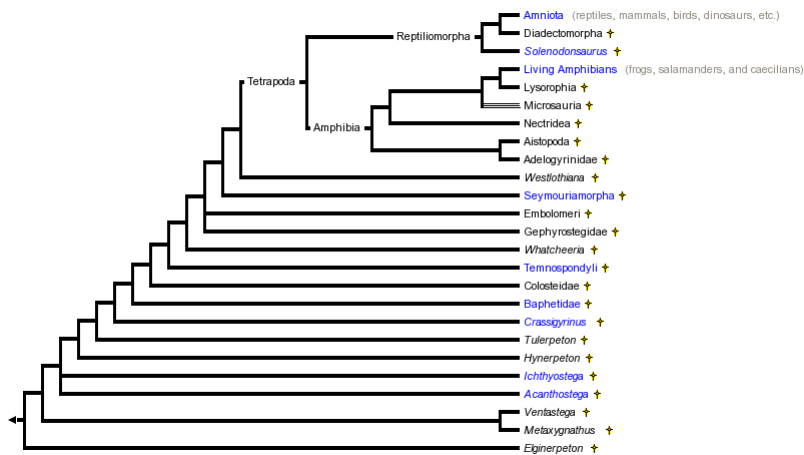- Re-build multiple alignment
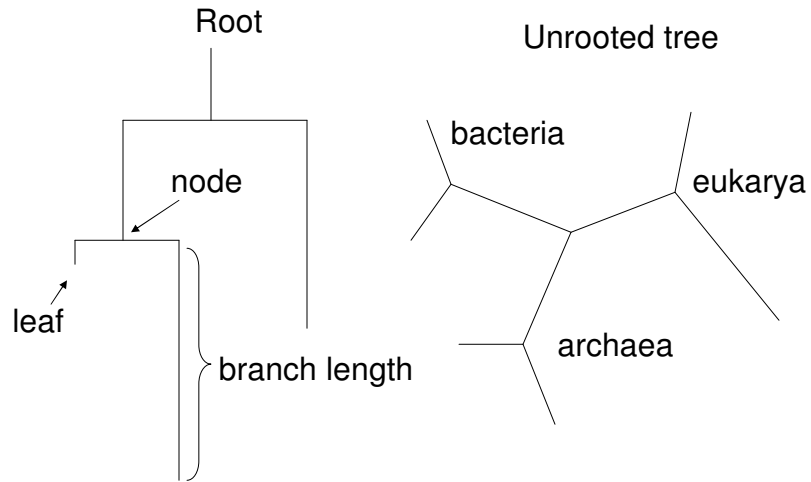- etc. etc. etc.

# Biological relevance of multiple alignments



# 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

branch length

archaea

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

# Phylogeny questions

- A. Easy-ish – can be done with dynamic programming
- B. Hard – Many possible trees

$$\frac{(2n-3)!}{2^{n-2}(n-2)!}$$    rooted trees with n leaves

# Scoring a tree – Sankoff's algorithm

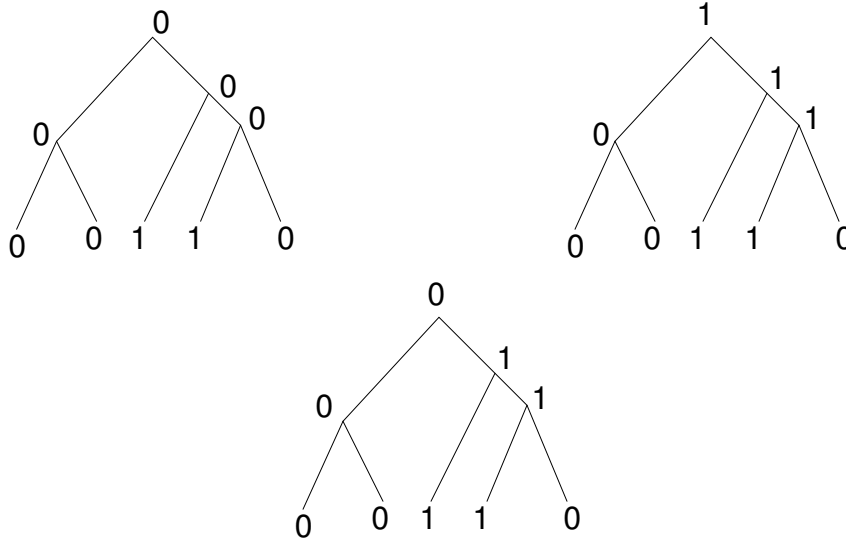- Assumption – we try to minimize # of state changes from root to leaves – Parsimony approach
- Small parsimony
  - given a tree where leaves are labeled with m-character strings
  - find labels at internal nodes s.t. # of state transitions is minimzed
- Weighted small parsimony
  - same as parsimony except that state transitions are assigned weights
  - minimize the overall weight of the tree

# Example

0
0
0
0
0 0 1 1 0

1
1
1
0
0 0 1 1 0

0
1
1
0
0 0 1 1 0

# Sankoff's algorithm

- At each node v in the tree store $s(v,t)$ – best parsimony score for subtree rooted at v if character stored at v is t
- Traverse the tree in post-order and update $s(v,t)$ as follows
  - assume node v has children u and w
  - $s(v,t) = \min_i \{s(u,i) + score(i,t)\} + \min_j \{s(w,j) + score(j,t)\}$
- Character at root will be the one that maximizes $s(root, t)$

- Note – this solves the weighted version. For unweighted set score $(i,i) = 0$, $score(i,j) = 1$ for any $i,j$