

**Example files**  
  
 Last update 25<sup>th</sup> January  
 Will contain problems and examples  
 not examinable

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0						
C 1							
A 2							
T 3							
G 4							
T 5							

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1					
C 1							
A 2							
T 3							
G 4							
T 5							

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1							
A 2							
T 3							
G 4							
T 5							

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1						
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

-----  
CATGT

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1					
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

A  
C

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1				
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

AC  
-C

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0			
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

ACG  
-C-

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0	-1		
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

ACGC  
---C

ACGC  
-C--

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0	-1	-2	-3
A 2	-2	1	0	0			
T 3	-3						
G 4	-4						
T 5	-5						

ACG  
-CA

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0	-1	-2	-3
A 2	-2	1	0	0	-1	-2	-3
T 3	-3	0	0	-1	-1	1	0
G 4	-4	-1	-1	2	1	0	3
T 5	-5	-2	-2	1	1	3	2

	0	A 1	C 2	G 3	C 4	T 5	G 6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0	-1	-2	-3
A 2	-2	1	0	0	-1	-2	-3
T 3	-3	0	0	-1	-1	1	0
G 4	-4	-1	-1	2	1	0	3
T 5	-5	-2	-2	1	1	3	2

	0	A	C	G	C	T	G
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

	0	A	C	G	C	T	G
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

ACGCTG-  
-C-ATGT

	0	A	C	G	C	T	G
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

ACGCTG-  
-CA-TGT

	0	A	C	G	C	T	G
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

-ACGCTG  
CATG-T-

$m = 1$   
 $s = 1$   
 $d = 2$

	x	A	G	C	
y		0	1	2	3
0		0	-2	-4	-6
A 1		-2			
A 2		-4			
A 3		-6			
C 4		-8			

We continue to fill the matrix using the recurrence rule

17

	x	A	G	C	
y		0	1	2	3
0		0	-2	-4	-6
A 1		-2	1		
A 2		-4			
A 3		-6			
C 4		-8			

F[0,0]	-1	F[0,1]
F[1,0]		F[1,1]

-2 -A  
A- versus

-2 (A- versus -A)

18

	x	A	G	C	
y		0	1	2	3
0		0	-2	-4	-6
A 1		-2	1	-1	-3
A 2		-4	-1	0	
A 3		-6	-3		
C 4		-8	-5		

19

	x	A	G	C	
y		0	1	2	3
0		0	-2	-4	-6
A 1		-2	1	-1	-3
A 2		-4	-1	0	-2
A 3		-6	-3	-2	-1
C 4		-8	-5	-4	-1

Conclusion:  $d(AAAC, AGC) = -1$

20

- To reconstruct the best alignment, we record which case(s) in the recursive rule maximized the score

	x	A	G	C
y	0	1	2	3
0	0	-2	-4	-6
A 1	-2	1	-1	-3
A 2	-4	-1	0	-2
A 3	-6	-3	-2	-1
C 4	-8	-5	-4	-1

21

- We now trace back a path that corresponds to the best alignment

	x	A	G	C
y	0	1	2	3
0	0	-2	-4	-6
A 1	-2	1	-1	-3
A 2	-4	-1	0	-2
A 3	-6	-3	-2	-1
C 4	-8	-5	-4	-1

AAAC  
AG-C

22

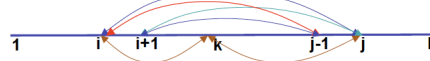
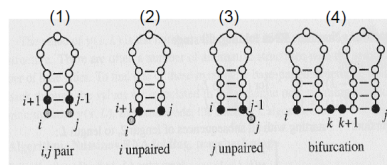
- Sometimes, more than one alignment has the best score

	x	A	G	C
y	0	1	2	3
0	0	-2	-4	-6
A 1	-2	1	-1	-3
A 2	-4	-1	0	-2
A 3	-6	-3	-2	-1
C 4	-8	-5	-4	-1

AAAC  
A-GC  
AAAC  
-AGC  
AAAC  
AG-C

23

### "Grow" from substructures



$$S(i, j) = \max \begin{cases} S(i+1, j-1) + w(i, j) & (1) \\ S(i+1, j) & (2) \\ S(i, j-1) & (3) \\ \max_{i < k < j} S(i, k) + S(k+1, j) & (4) \end{cases}$$

$w(i, j) = 1$  if  $i, j$  are complementary (i.e., GC, CG, AU or UA); 0 otherwise

### Initialization

	G	G	G	A	A	A	U	C	C
G	0								
G	0	0							
G	0	0	0						
A				0	0				
A				0	0	0			
A				0	0	0	0		
U							0	0	
C							0	0	0
C							0	0	0

Example:  
GGGAAAUCC

$$S(i, i) = 0 \quad \forall 1 \leq i \leq L \quad \rightarrow \text{the main diagonal}$$

$$S(i, i-1) = 0 \quad \forall 2 \leq i \leq L \quad \rightarrow \text{the diagonal below}$$

$L$ : the length of input sequence

### Recursion

→  $j$

Fill up the table (DP matrix) -- diagonal by diagonal

	G	G	G	A	A	A	U	C	C
G	0	0	0	0					
G	0	0	0	0	0				
G		0	0	0	0	0			
A			0	0	0	0	?		
A				0	0	0	1		
A					0	0	1	1	
U							0	0	0
C								0	0
C									0

$$S(i, j) = \max \begin{cases} S(i+1, j-1) + w(i, j) & (1) \\ S(i+1, j) & (2) \\ S(i, j-1) & (3) \\ \max_{i < k < j} S(i, k) + S(k+1, j) & (4) \end{cases} \quad w(i, j) = \begin{cases} 1 & i, j \text{ are complementary} \\ 0 & \text{otherwise} \end{cases}$$

### Traceback

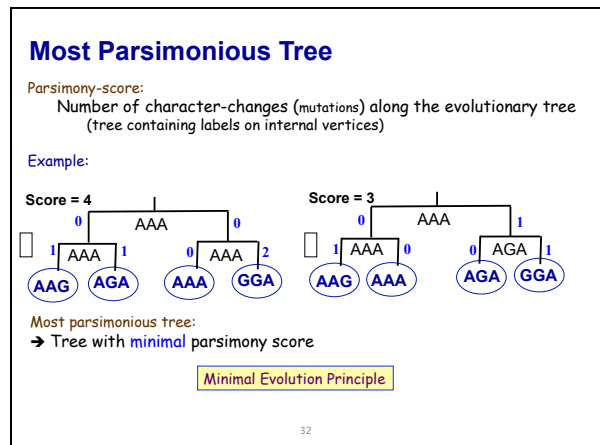
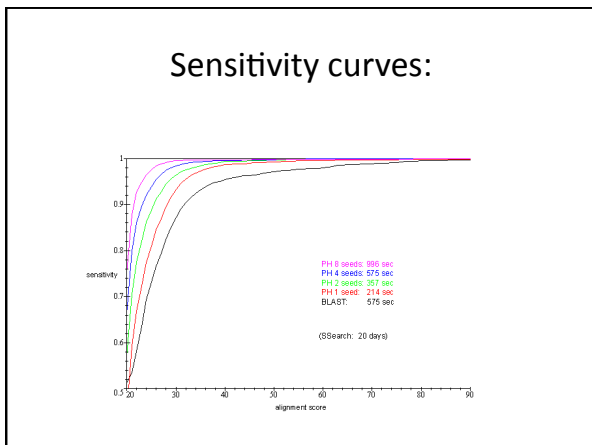
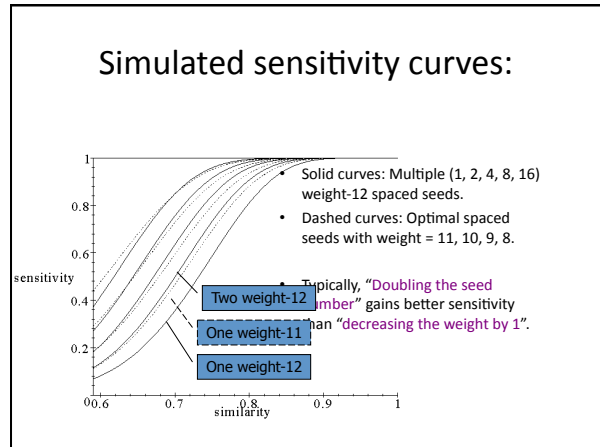
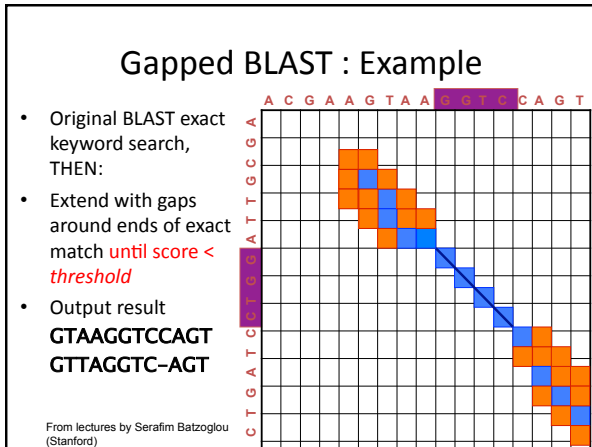
	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G		0	0	0	0	0	1	2	2
A				0	0	0	0	1	1
A					0	0	0	1	1
A						0	0	1	1
U							0	0	0
C								0	0
C									0

The structure is:

### Original BLAST: Example

- $w = 4$
- Exact keyword match of **GGTC**
- Extend diagonals with mismatches until score is under 50%
- Output result  
**GTAAGGTC**  
**GTTAGGTC**

From lectures by Serafim Batzoglou (Stanford)





### Fitch's Algorithm

Execute independently for each character:

- Bottom-up phase:** Determine set of possible states for each internal node
- Top-down phase:** Pick states for each internal node

Dynamic Programming framework

33

### Fitch's Algorithm

Bottom-up phase

Determine set of possible states for each internal node

- Initialization:  $R_i = \{s_i\}$
- Do a post-order (from leaves to root) traversal of tree
  - Determine  $R_i$  of internal node  $i$  with children  $j, k$ :

$$R_i = \begin{cases} R_j \cap R_k & \text{if } R_j \cap R_k \neq \emptyset \\ R_j \cup R_k & \text{otherwise} \end{cases}$$

34

### Fitch's Algorithm

Top-down phase

Pick states for each internal node

- Pick arbitrary state in  $R_{root}$  for the root
- Do pre-order (from root to leaves) traversal of tree
  - Determine  $s_j$  of internal node  $j$  with parent  $i$ :

$$s_j = \begin{cases} s_i & \text{if } s_i \in R_j \\ \text{arbitrary state} \in R_j & \text{otherwise} \end{cases}$$

35

### Weighted Parsimony

Sankoff's algorithm

- Each mutation  $a \rightarrow b$  costs differently -  $S(a,b)$ .

- Bottom-up phase:** Determine  $R_i(s)$  - cost of optimal state-assignment for subtree of  $i$ , when it is assigned state  $s$ .
- Top-down phase:** Pick optimal states for each internal node

Fitch's algorithm as special case:

- $R_i$  - set of states which yield minimal-cost subtree of  $i$

36

### Sankoff's Algorithm

Bottom-up phase

Determine  $R_i(s)$  for each internal node

- Initialization:  $R_i(s) = \begin{cases} 0 & \text{if } s_i = s \\ \infty & \text{otherwise} \end{cases}$
- Do a post-order (from leaves to root) traversal of tree
  - Determine  $R_i$  of internal node  $i$  with children  $j, k$ :
 
$$R_i(s) = \min_{s'} \{R_j(s') + S(s', s)\} + \min_{s''} \{R_k(s'') + S(s'', s)\}$$

Natural generalization For non-binary trees

Remember pointers  $s \rightarrow s'$

37

### Sankoff's Algorithm

Top-down phase

Pick states for each internal node

- Select minimal cost character for root ( $s$  minimizing  $R_{root}(s)$ )
- Do pre-order (from root to leaves) traversal of tree:
  - For internal node  $j$ , with parent  $i$ , select state that produced minimal cost at  $i$  (use pointers kept in 1<sup>st</sup> stage)

Complexity:  $O(mnk^2)$

#characters #states #taxa/nodes

38

### Large Parsimony Problem

- Input:** An  $n \times m$  matrix  $M$  describing  $n$  species, each represented by an  $m$ -character string
- Output:** A tree  $T$  with  $n$  leaves labeled by the  $n$  rows of matrix  $M$ , and a labeling of the internal vertices such that the parsimony score is minimized over all possible trees and all possible labelings of internal vertices
- Possible search space is huge, especially as  $n$  increases
- $(2n - 3)!!$  possible rooted trees
- $(2n - 5)!!$  possible unrooted trees
- Problem is NP-complete; Exhaustive search only possible w/ small  $n (< 10)$

39

### Solving NP-hard problems exactly is ... unlikely

- Number of (unrooted) binary trees on  $n$  leaves is  $(2n-5)!!$
- If each tree on **1000** taxa could be analyzed in **0.001** seconds, we would find the best tree in **2890 millennia**

#leaves	#trees
4	3
5	15
6	105
7	945
8	10395
9	135135
10	2027025
20	$2.2 \times 10^{20}$
100	$4.5 \times 10^{190}$
1000	$2.7 \times 10^{2900}$

### Genomes Evolve by Rearrangements

- Inversion (Reversal)
- Transposition
- Inverted Transposition

### The adjacency graph and the distance equation

$C$  = number of cycles  
 $I$  = number of odd paths  
 $G$  = number of "genes"

$D = G - (C + I/2)$

$D = 6 - (1 + 2/2) = 4$

Joint work with Julia Mixtacki and Jens Stoye

(A)

$C_{ij}$	a	g	c	t
a	0	1	3	3
g	1	0	3	3
c	3	3	0	1
t	3	3	1	0

(B)

Original Sankoff parsimony: up phase. (A) A matrix defining transition costs between states. (B) A sample phylogenetic tree with cost vectors calculated for all nodes using Algorithm 1. Semerari et al. BMC Bioinformatics 2009 10:51 doi:10.1186/1471-2105-10-51

Original Sankoff parsimony: down phase. Ancestral states for the phylogeny, cost matrix and observed states of Figure 1, as obtained by Algorithm 2. Without solving ties, unrooted nodes are written next to each inner node. Red lines indicate the most parsimonious reconstructions after solving ties. Semerari et al. BMC Bioinformatics 2009 10:51 doi:10.1186/1471-2105-10-51

```

Algorithm 1 (Original Sankoff algorithm: Up phase). A procedure call Sankoff_Up(T, C, S) calculates the cost vector  $S^{(p)}$  of all nodes  $p$  of the phylogeny  $T$ , given a cost matrix  $C = (c_{ij})$ .

procedure Sankoff_Up( $T, C, S$ )
  for all nodes  $p$  of  $T$  in postorder do
    if  $p$  is a leaf then
      for all  $i$  in  $1, \dots, n$  do
        if state  $i$  observed at leaf  $p$  then
           $S^{(p)}_i \leftarrow 0$ 
        else
           $S^{(p)}_i \leftarrow \infty$ 
      else
         $(q, r) \leftarrow$  children of  $p$ 
        for all  $i$  in  $1, \dots, n$  do
           $S^{(p)}_i \leftarrow \min_j \{ \text{cost}(q, i) + S^{(q)}_j, \text{cost}(r, i) + S^{(r)}_j \}$ 
function cost( $x, i$ )
   $\min \leftarrow \infty$ 
  for all  $j$  in  $1, \dots, n$  do
    if  $c_{ij} + S^{(x)}_j < \min$  then
       $\min \leftarrow c_{ij} + S^{(x)}_j$ 
  return  $\min$ 
    
```

```

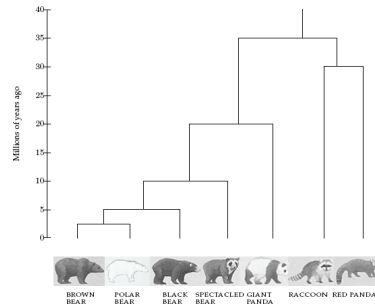
Algorithm 2 (Original Sankoff algorithm: Down phase). A procedure call Sankoff_Down(x, T, C, S, Sanc) calculates the ancestral states  $S^{(x)}$  of all nodes  $p$  of the phylogeny  $T$ , given the root  $x$  of  $T$ , a cost matrix  $C = (c_{ij})$  of transition costs between states, and the cost vectors  $S^{(p)}$  for all nodes  $p$  of  $T$  as calculated by Sankoff_Up(T, C, S).

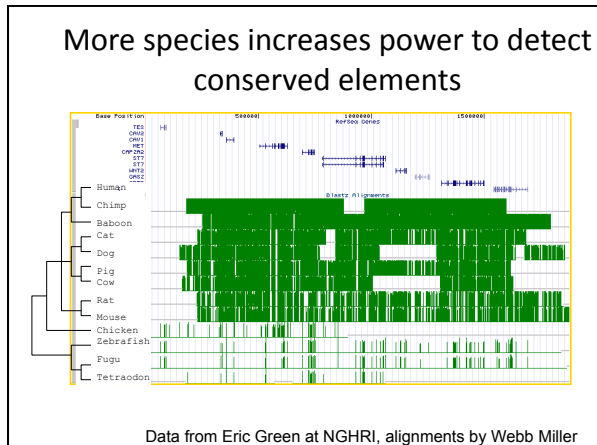
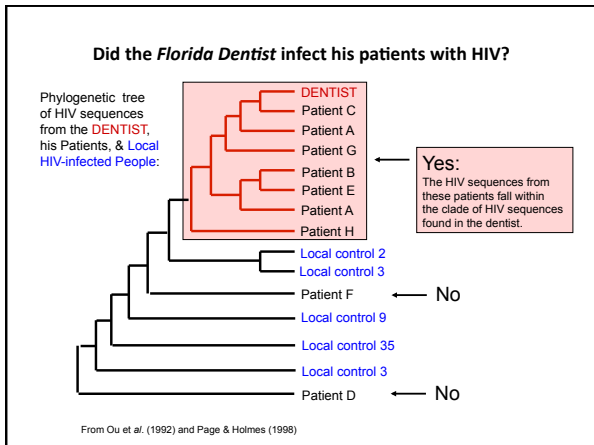
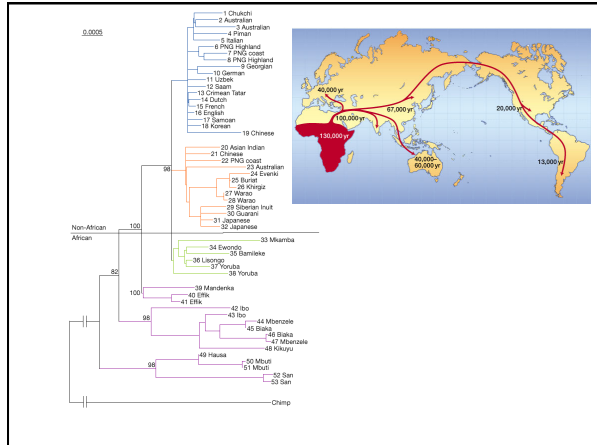
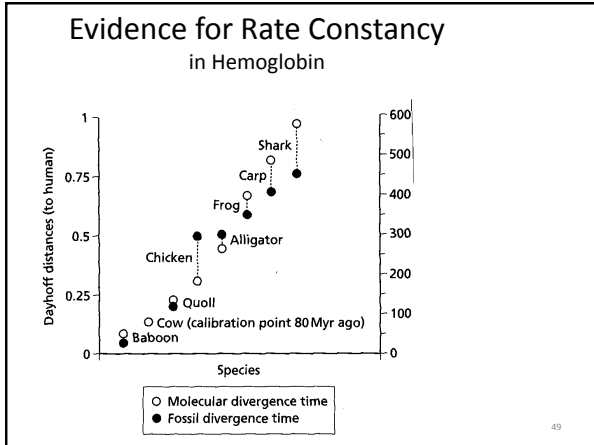
procedure Sankoff_Down( $x, T, C, S, S_{anc}$ )
   $S^{(x)}_i \leftarrow \arg \min_j S^{(x)}_j$ 
  for all  $j$  in  $S^{(x)}$  do
    for all child  $y$  of  $x$  do
      Sankoff_Down( $y, T, C, S, S_{anc}$ )
procedure Sankoff_Down( $x, T, C, S, S_{anc}$ )
   $\min\_state(x, x, C, S, S_{anc})$ 
  for all  $j$  in  $S^{(x)}$  do
    for all child  $y$  of  $x$  do
      Sankoff_Down( $y, T, C, S, S_{anc}$ )
procedure min_state( $x, x, C, S, S_{anc}$ )
   $\min \leftarrow \infty$ 
  for all  $j$  in  $1, \dots, n$  do
    if  $x = \text{root}(T)$  then
       $\text{trans\_cost} \leftarrow S^{(x)}_j$ 
    else
       $\text{trans\_cost} \leftarrow c_{ij} + S^{(y)}_j$ 
    if  $\text{trans\_cost} < \min$  then
       $\min \leftarrow \text{trans\_cost}$ 
       $S^{(x)}_{\text{anc}} \leftarrow j$ 
    else if  $\text{trans\_cost} = \min$  then
       $S^{(x)}_{\text{anc}} \leftarrow S^{(x)}_{\text{anc}} \cup \{j\}$ 
    
```

### The Giant Panda Riddle

- Giant pandas look like bears but have features that are unusual for bears and typical for raccoons, e.g., they do not hibernate
- Is the Giant panda closer to a bear or to a raccoon?**

### Evolutionary Tree of Bears and Raccoons





### Approximate methods

- For larger data sets computing time becomes prohibitive and we only explore some subset of all possible trees (hoping that the optimal trees will be found in the subset explored)
- *Heuristic approaches* sacrifice the guarantee of optimality in favor of reduced computer time
- Use “hill climbing” methods. Initial tree starts the process, then we seek to improve its score
- When we can find no way to further improve the score, we stop. We don’t know if we reached a *local* or a *global optimum*

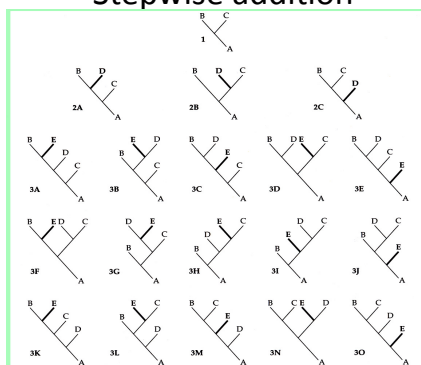
53

### Initial trees

- May be obtained by *stepwise addition*, the most commonly used method
- Similar to exhaustive search but evaluate trees at every step, each time you add a new taxon and only follow the path derived from the optimal tree
- Which taxa do you choose first? Which do you connect next?
- These are “greedy algorithms”

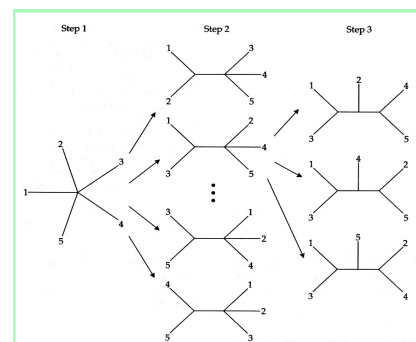
54

### Stepwise addition



55

- Initial trees also may be obtained by *star decomposition*, another greedy algorithm



56

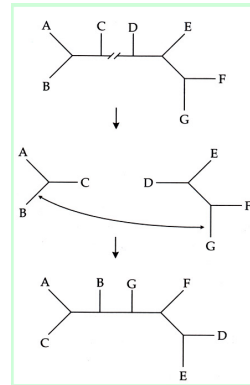
### Branch swapping

- To improve the initial estimate we can perform sets of predefined rearrangements on the tree
- Any of these rearrangements amounts to a 'stab in the dark'
- Globally optimal trees may be several rearrangements away from the starting tree
- If a better tree is found, a new round of rearrangements is then performed in the new tree
- Several branch-swapping algorithms are available

57

### Branch swapping by tree bisection and reconnection (TBR)

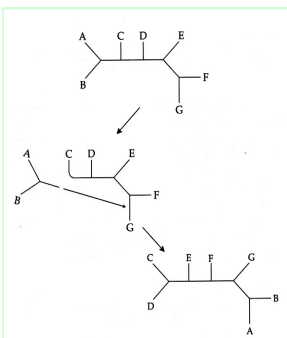
1. Tree is bisected along a branch, yielding two disjunct subtrees
2. The subtrees are reconnected by joining a pair of branches, one from each subtree
3. All possible bisections and pairwise reconnections are evaluated



58

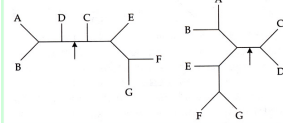
### Branch swapping by subtree pruning and regrafting

1. A subtree is pruned from the tree (e.g. A,B)
2. The subtree is then regrafted to a different location on the tree
3. All possible subtree removals and reattachment points are evaluated



59

### Branch swapping by nearest-neighbor interchanges (NNI)



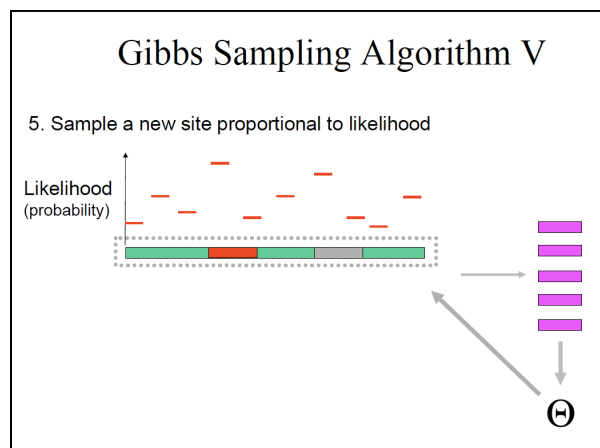
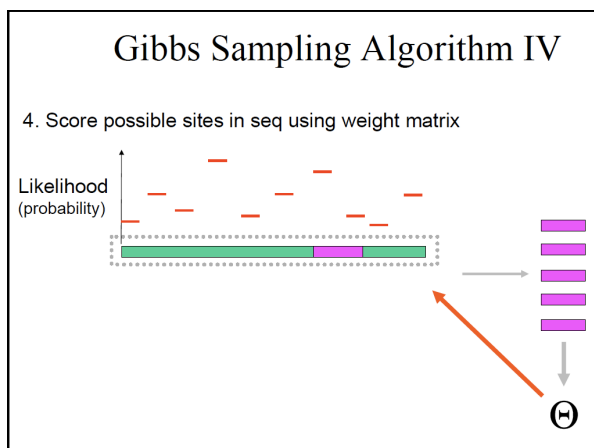
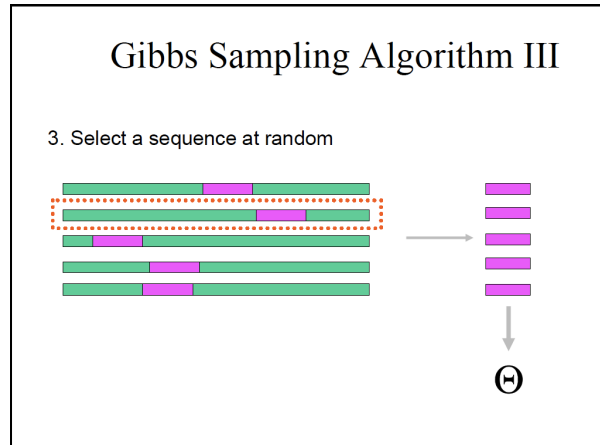
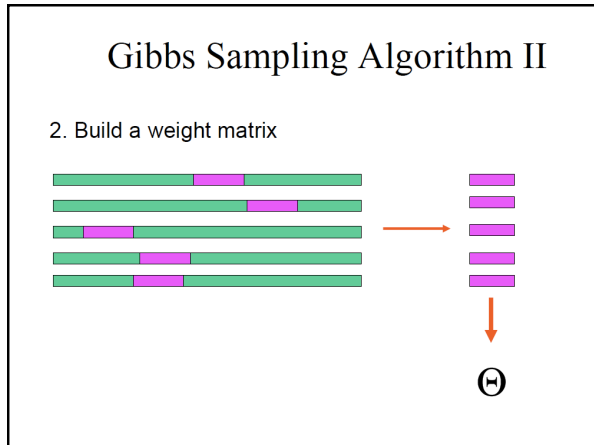
1. Each interior branch of the tree defines a local region of four subtrees

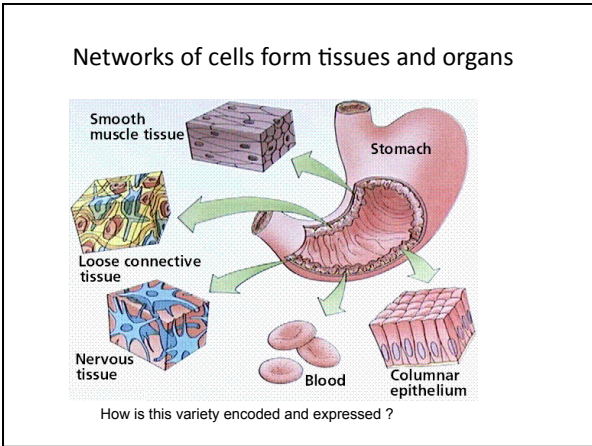
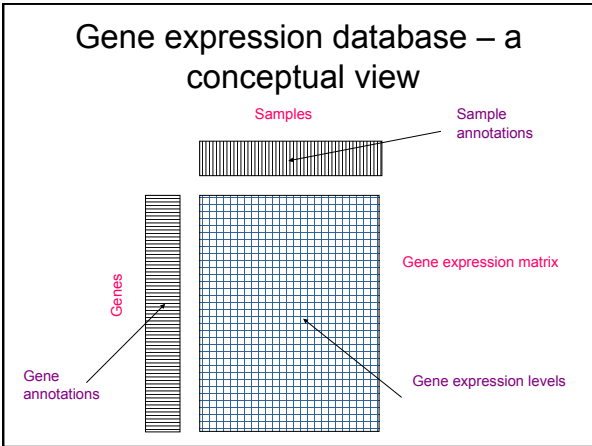
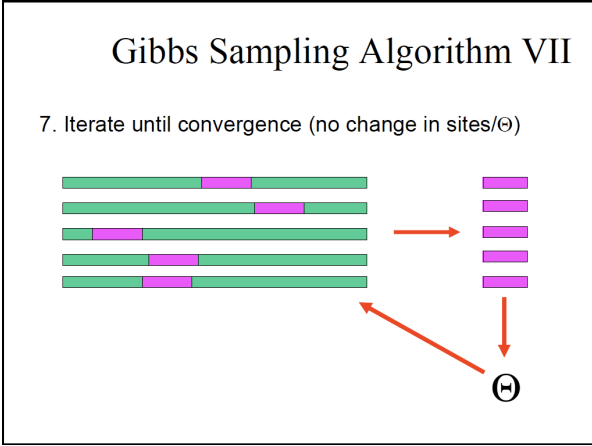
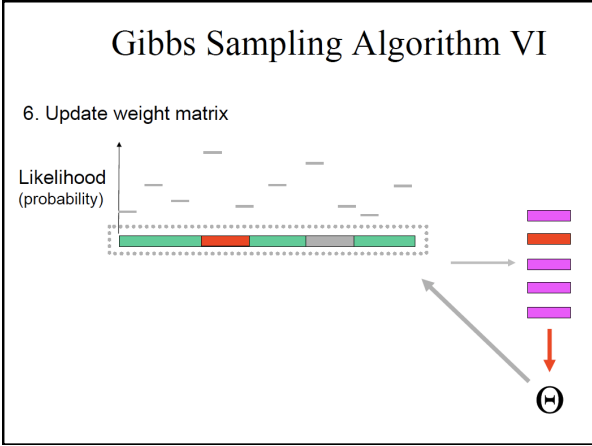
2. Interchanging a subtree on one side of the branch with one from the other constitutes an NNI
3. Two such rearrangements are possible for each interior branch (all interior branches are swapped)

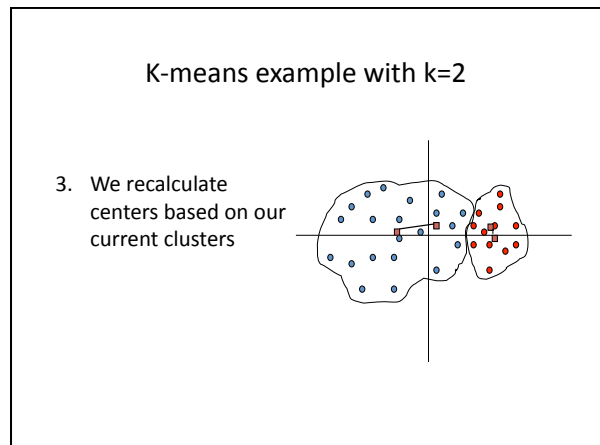
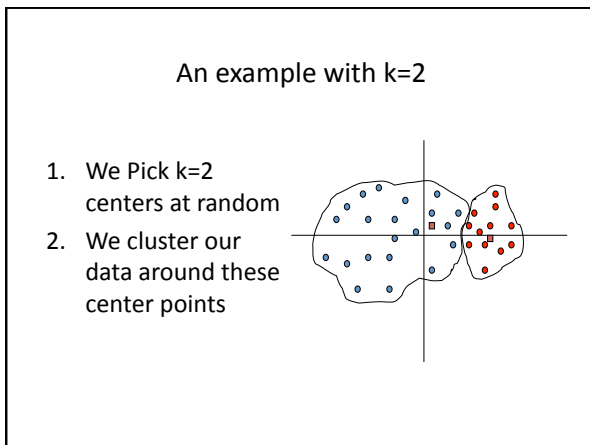
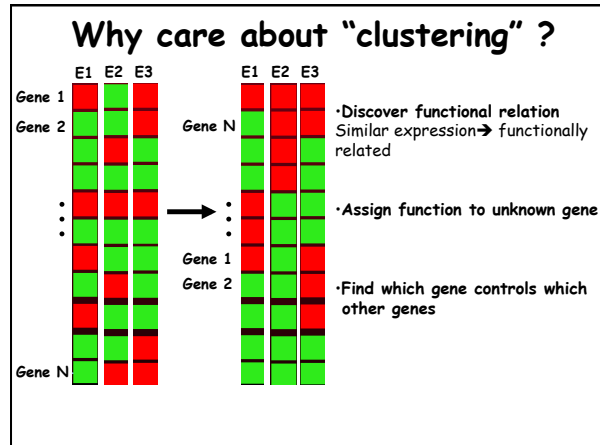
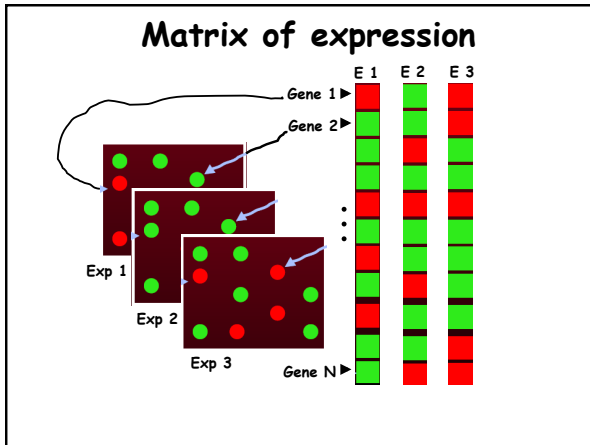
60











### K-means example with k=2

4. We re-cluster our data around our new center points

### K-means example with k=2

5. We repeat the last two steps until no more data points are moved into a different cluster

### Cluster Quality

Quality of cluster assessed by ratio of distance to nearest cluster and cluster diameter

### PCA algorithm

- "Column-centered" matrix:  $A$
- Covariance matrix:  $A^T A$
- Eigenvalue Decomposition
  - $A^T A = U \Lambda U^T$
  - $U$ : Eigenvectors (principle components)
  - $\Lambda$ : Eigenvalues