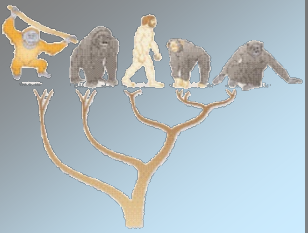


Sequence Analysis and Phylogenetics

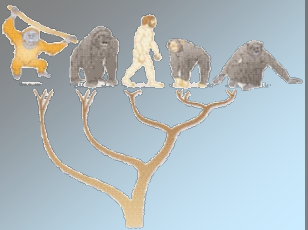
Part 3

Sepp Hochreiter



Contents

- 4 Multiple Alignment
 - 4.1 Motivation
 - 4.2 Multiple Sequence Similarities and Scoring
 - 4.2.1 Consensus and Entropy Score
 - 4.2.2 Tree and Star Score
 - 4.2.3 Weighted Sum of Pairs Score
 - 4.3 Multiple Alignment Algorithms
 - 4.3.1 Exact Methods
 - 4.3.2 Progressive Algorithms
 - 4.3.3 Other Multiple Alignment Algorithms
 - 4.4 Profiles and Position Specific Scoring Matrices



Motivation

4 Multiple Alignment

Compare more than two sequences: arranged sequences so that the amino acids for every the columns match as good as possible

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

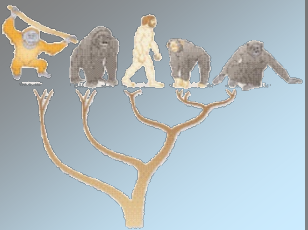
4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

		10		20		30																																		
<i>Human</i>	AP	S	R	K	F	F	V	G	G	N	W	K	M	N	G	R	K	Q	S	L	G	E	L	I	G	T	L	N	A	...	A	K	V	P	A	D	T			
<i>Chicken</i>	...	R	K	F	F	V	G	G	N	W	K	M	N	G	D	K	K	S	L	G	E	L	I	H	T	L	N	G	...	A	K	L	S	A	D	T				
<i>Yeast</i>	.	G	A	G	K	F	V	V	G	G	N	W	K	C	N	G	T	L	A	S	I	E	T	L	T	K	G	V	A	A	S	V	D	A	E	L	A	K	K	V
<i>E. coli</i>	..	A	R	T	F	F	V	G	G	N	F	K	L	N	G	S	K	Q	S	I	K	E	I	V	E	R	L	N	T	...	A	S	I	P	E	N	V			
<i>Amoeba</i>	..	M	R	H	P	L	V	M	G	N	W	K	L	N	G	S	R	H	M	V	H	E	L	V	S	N	L	R	K	..	E	L	A	G	V	A	G	C		
<i>Archaeon</i>	A	K	L	K	E	P	I	A	I	N	F	K	T	Y	I	E	A	T	G	K	R	A	L	E	I	A	K	A	A	...	E	K	V	Y	K	E	T			
<i>consensus</i>	...	r	.	f	.	v	g	g	n	w	k	l	n	g	.	k	.	s	i	.	e	l	v	.	l	.	a	...	a	.	v	...								

		40		50		60		70																													
<i>Human</i>	E	V	V	C	A	P	P	T	A	Y	I	D	F	A	R	Q	K	L	D	...	P	K	I	A	V	A	A	Q	N	C	Y	K	V	T	N	G	
<i>Chicken</i>	E	V	V	C	G	A	P	S	I	Y	L	D	F	A	R	Q	K	L	D	...	A	K	I	G	V	A	A	Q	N	C	Y	K	V	P	K	G	
<i>Yeast</i>	E	V	I	V	G	V	P	F	I	Y	I	P	K	V	Q	Q	I	L	A	G	E	A	N	I	L	V	S	A	E	N	A	W	T	K	S	.	G
<i>E. coli</i>	E	V	V	I	C	P	P	A	T	Y	L	D	Y	S	V	S	L	V	K	K	...	P	Q	V	T	V	G	A	Q	N	A	Y	L	K	A	S	G
<i>Amoeba</i>	A	V	A	I	A	P	P	E	M	Y	I	D	M	A	K	R	E	A	E	G	...	S	H	I	M	L	G	A	Q	N	V	N	L	N	L	S	G
<i>Archaeon</i>	G	V	T	I	V	V	A	P	Q	L	V	D	L	R	M	I	A	E	S	...	V	E	I	P	V	F	A	Q	H	I	D	P	I	K	P	G	
<i>consensus</i>	e	V	v	i	a	.p	.	y	i	d	...	l	i	.	v	g	A	q	n	.y	...	G							



Motivation

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

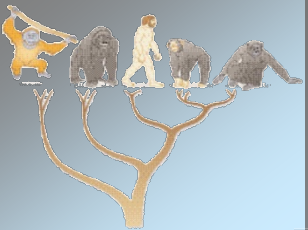
4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

		80		90		100		110																																	
<i>Human</i>		A	F	T	G	E	I	S	P	G	M	I	K	D	C	G	A	T	W	V	V	L	G	H	S	E	R	R	H	V	F	G	E	S	D	E	L	I	G	Q	Q
<i>Chicken</i>		A	F	T	G	E	I	S	P	A	M	I	K	D	I	G	A	A	W	V	I	L	G	H	S	E	R	R	H	V	F	G	E	S	D	E	L	I	G	Q	Q
<i>Yeast</i>		A	Y	T	G	E	V	H	V	G	M	L	V	D	C	Q	V	P	Y	V	I	L	G	H	S	E	R	R	Q	I	F	H	E	S	N	E	Q	V	A	E	K
<i>E. coli</i>		A	F	T	G	E	N	S	V	D	Q	I	K	D	V	G	A	K	Y	V	I	L	G	H	S	E	R	R	S	Y	F	H	E	D	D	K	F	I	A	D	K
<i>Amoeba</i>		A	F	T	G	E	T	S	A	A	M	L	K	D	I	G	A	Q	Y	I	I	I	G	H	S	E	R	R	T	Y	H	K	E	S	D	E	L	I	A	K	K
<i>Archaeon</i>		S	H	T	G	H	V	L	P	E	A	V	K	E	A	G	A	V	G	T	L	L	N	H	S	E	N	R	M	I	L	A	D	L	E	A	A	I	R	R	.
<i>consensus</i>		a	f	T	G	e	v	s	.	a	m	i	k	d	.	g	a	.	y	v	i	l	g	H	S	E	r	R	.	i	f	.	e	s	d	e	.	i	a	.	k

		120		130		140		150																																			
<i>Human</i>		V	A	H	A	L	A	E	G	L	G	V	I	A	C	I	G	E	K	L	D	E	R	E	A	G	I	T	E	K	V	V	F	E	Q	T	K	V	I	A	D		
<i>Chicken</i>		V	A	H	A	L	A	E	G	L	G	V	I	A	C	I	G	E	K	L	D	E	R	E	A	G	I	T	E	K	V	V	F	E	Q	T	K	A	I	A	D		
<i>Yeast</i>		V	K	V	A	I	D	A	G	L	K	V	I	A	C	I	G	E	T	E	A	Q	R	I	A	N	Q	T	E	E	V	V	A	A	Q	L	K	A	I	N	N		
<i>E. coli</i>		T	K	F	A	L	G	Q	G	V	G	V	I	L	C	I	G	E	T	L	E	E	K	K	A	G	K	T	L	D	V	V	E	R	Q	L	N	A	V	L	E		
<i>Amoeba</i>		F	A	V	L	K	E	Q	G	L	T	P	V	L	C	I	G	E	T	E	A	E	N	E	A	G	K	T	E	E	V	C	A	R	Q	I	D	A	V	L	K		
<i>Archaeon</i>		.	.	.	A	E	E	V	G	L	M	T	M	V	C	S
<i>consensus</i>		.	.	.	a	l	.	.	G	l	.	v	i	.	C	i	g	e	.	.	e	r	.	a	g	.	t	e	.	v	v	.	q	l	.	a	i	.	.				



Motivation

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

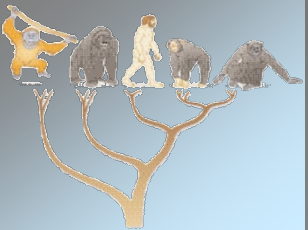
4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

	160	170	180	190
<i>Human</i>	NV . . KD WSK V V L A Y E P V W A I G T G K T A T P Q Q A Q E V H E K L R G			
<i>Chicken</i>	NV . . KD WSK V V L A Y E P V W A I G T G K T A T P Q Q A Q E V H E K L R G			
<i>Yeast</i>	A I S K E A W K N I I L A Y E P V W A I G T G K T A T P D Q A Q E V H Q Y I R K			
<i>E. coli</i>	E V . . K D F T N V V V A Y E P V . A I G T G L A A T P E D A Q D I H A S I R K			
<i>Amoeba</i>	T Q G A A A F E G A V I A Y E P V W A I G T G K S A T P A Q A Q A V H K F I R D			
<i>Archaeon</i>	D Y V A V E P P E L I G T G I P V S K A K P E V I T N			
<i>consensus</i>	. v w . . v v l A y E P v w a I G T G k t a t p . q a q e v h . . i r .			
	200	210	220	230
<i>Human</i>	W L K S N V S D A V A Q S T R I I Y G G S V T G A T C K E L A S Q P D V D G F L			
<i>Chicken</i>	W L K T H V S D A V A Q S T R I I Y G G S V T G G N C K E L A S Q H D V D G F L			
<i>Yeast</i>	W M T E N I S K E V A E A T R I Q Y G G S V N P A N C N E L A K K A D I D G F L			
<i>E. coli</i>	F L A S K L G D K A A S E L R I L Y G G S A N G S N A V T F K D K A D V D G F L			
<i>Amoeba</i>	H I A K . V D A N I A E Q V I I Q Y G G S V N A S N A A E L F A Q P D I D G A L			
<i>Archaeon</i>	. . T V E L V K K V N P E V K V L C G A G I S T G E D V K K A I E L G T V G V L			
<i>consensus</i>	w l . . . v . . . v a . . . r i l y G g s v . g g n . . e l a . . . d v d G f L			



Motivation

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

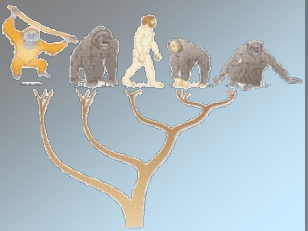
4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

	240
<i>Human</i>	VGGASLKP . EFVVDIINAKQ
<i>Chicken</i>	VGGASLKP . EFVVDIINAKH
<i>Yeast</i>	VGGASLDAAKFKTIINSVSEKL . .
<i>E. coli</i>	VGGASLKP . EFVVDIINSRN
<i>Amoeba</i>	VGGASLKADAFAVIVKAAEAAKQA
<i>Archaeon</i>	LASGVTKAKDPEKAIWDLVSGI . .
<i>consensus</i>	v g g a s l k . . e f . . i i n



Motivation

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

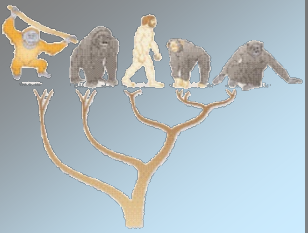
4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Multiple sequence alignment is used to

- ↳ detect remote homologous regions
- ↳ detect motifs (regular patterns) in protein families
- ↳ detect conserved regions or positions (disulfide bonds)
- ↳ detect structural blocks like helices or sheets
- ↳ construct phylogenetic trees
- ↳ construct a profiles (search or averages)
- ↳ sequence genomes by superimposing fragments (nucleotides)
- ↳ cluster proteins according to similar regions



Scoring and Similarity

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

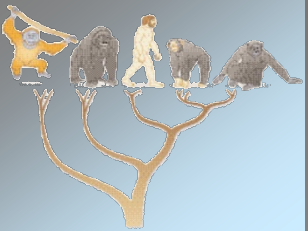
4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Similarity measures can be based on:

- ↳ the similarity of all sequences to a reference sequence
- ↳ the similarities between evolutionary adjacent sequences
- ↳ all pairwise similarities



Consensus and Entropy

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

consensus sequence: obtained if for each column in the alignment

1. the **most frequent** amino acid or
2. the amino acid which has the **highest score to all other** amino acids is chosen

4.2.2 Tree and Star

consensus score: sum of the pairwise score between sequences and the consensus sequence

4.2.3 Sum of Pairs

4.3 Algorithms

generalized by profiles instead of sequences

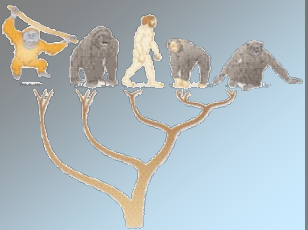
4.3.1 Exact Methods

4.3.2 Progressive

profile: relative frequency instead of most frequent

4.3.3 Other

4.4 Profiles / PSSMs



Consensus and Entropy

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

high entropy of the letter distribution: all letter are equally probable
zero entropy: one letter in the column

good alignment correlates with a low accumulative entropy

entropy score:
$$- \sum_i \sum_a f_{i,a} \log f_{i,a}$$

$f_{i,a}$: relative frequency of letter a in column i



Tree and Star Score

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

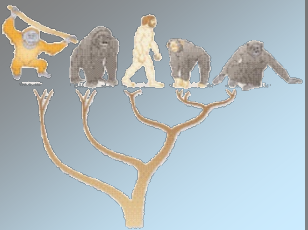
4.3.3 Other

4.4 Profiles / PSSMs

To count the number of mutations only those pairs should be compared which are evolutionary adjacent

E
E
E
E
D
D
D
D

evolutionary adjacent sequences are represented through a phylogenetic tree, which must be constructed



Tree and Star Score

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

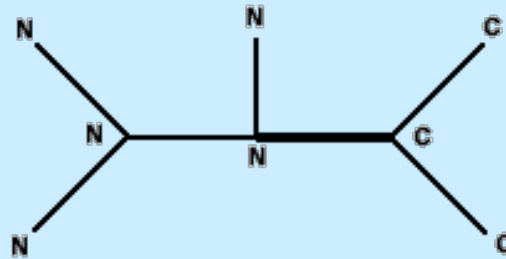
NNN

NNN

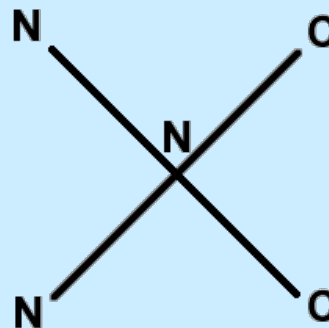
NNN

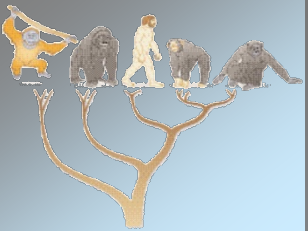
NNC

NCC



phylogenetic star: one sequence is considered as ancestor





Weighted Sum of Pairs

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

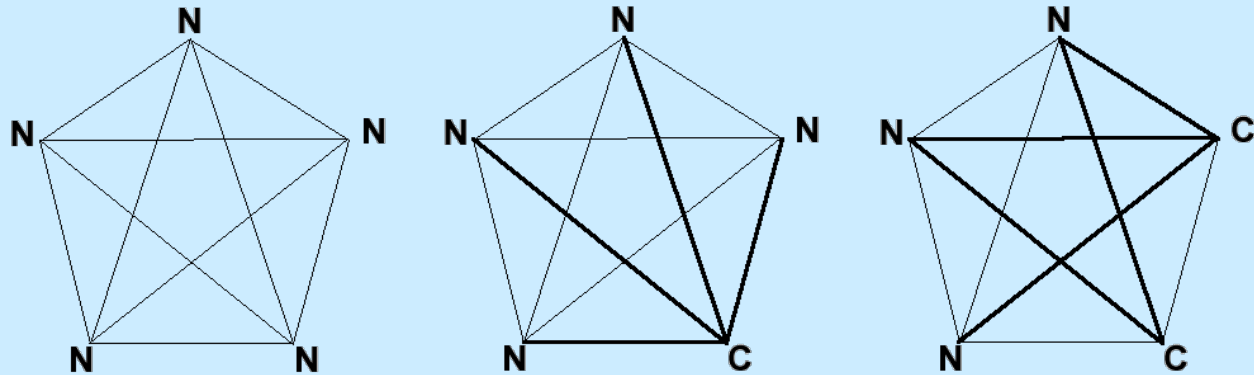
4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

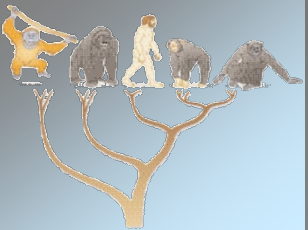
weighted sum of pairs: all pairwise comparisons



alignment length: L
number sequences: N

$$\sum_{i=1}^L \sum_{l=1}^{N-1} \sum_{j=l+1}^N w_{l,j} s(x_{i,l}, x_{i,j})$$

weights: reduce the influence of closely related sequences



Weighted Sum of Pairs

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Disadvantage: score relatively decreases with increasing N for conservative regions; but larger N means more conservative

C	C
C	C
C	C
C	C
C	D

$$S_{\text{old}} = \frac{N(N-1)}{2} s(C, C) \quad N \text{ Cs vs. } (N-1) \text{ Cs and one D}$$

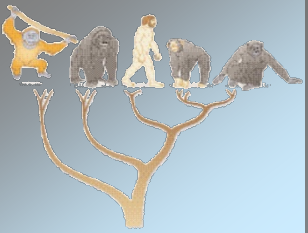
$$S_{\text{new}} = \frac{N(N-1)}{2} s(C, C) - (N-1)s(C, C) + (N-1)s(C, D)$$

$$\frac{S_{\text{old}} - S_{\text{new}}}{S_{\text{old}}} = \frac{2(N-1)s(C, C) - 2(N-1)s(C, D)}{N(N-1)s(C, C)} =$$

$$\frac{2}{N} \left(1 - \frac{s(C, D)}{s(C, C)} \right) \quad \text{the larger } N, \text{ the smaller the difference (paradox!)}$$

reasonable scoring matrices: $s(C, D) < s(C, C)$

$$\left(1 - \frac{s(C, D)}{s(C, C)} \right) > 0$$



Weighted Sum of Pairs

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

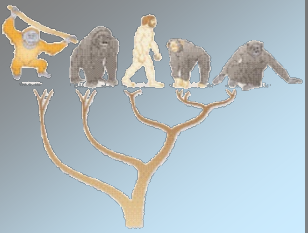
4.3.3 Other

4.4 Profiles / PSSMs

contra-intuitive: a new letter in a column of 100 equal letters is more surprising as a new letter in a column of 3 equal letters

Information gain: $-\log f_{i,a} = \log(N)$

Gaps: as for pairwise algorithms, linear gaps more efficient



Multiple Alignment Algorithms

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

multiple alignment optimization problem: NP-hard

Exact solution: only 10 to 15 sequences

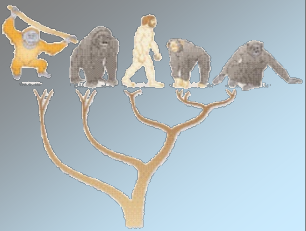
algorithm classes:

↳ global and progressive methods: MSA, COSA, GSA, clustalW, TCoffee

↳ iterative and search algorithms: DIALIGN, MultAlin, SAGA, PRRP, Realigner

↳ local methods (motif/profile): eMotif, Blocks, Dialign, Prosite, HMM, Gibbs sampling

↳ divide-and-conquer algorithms: DCA, OMA



Multiple Alignment Algorithms

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

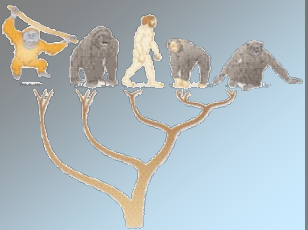
4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Global progressive alignments methods		
CLUSTALW	ftp://ftp.ebi.ac.uk/pub/software	Thompson et al. (1994/97) Higgins et al. (1996)
MSA	http://www.psc.edu/ http://www.ibr.wustl.edu/ibr/msa.html ftp://fastlink.nih.gov/pub/msa	Lipman et al. (1989) Gupta et al. (1995)
PRALINE	http://mathbio.nimr.mrc.ac.uk/~jhering/praline	Heringa (1999)
Iterative and search algorithms		
DIALIGN segment alignment	http://www.gsf.de/biodv/dialign.html	Morgenstern et al. (1996)
MultAlin	http://protein.toulouse.inra.fr/multalin.html	Corpet (1988)
PRRP progressive global alignment	ftp://ftp.genome.ad.jp/pub/genome/saitamacc	Gotoh (1996)
SAGA genetic algorithm	http://igs-server.cnrs-mrs.fr/~cnotred/Projects_home_page/saga_home_page.html	Notredame and Higgins (1996)
Local alignments / motif / profile		
Aligned Segment Statistical Eval. Tool (Asset)	ftp://ncbi.nlm.nih.gov/pub/neuwald/asset	Neuwald and Green (1994)
BLOCKS	http://blocks.fhcrc.org/blocks/	Henikoff and Henikoff (1991, 1992)
eMOTIF	http://dna.Stanford.EDU/emotif/	Nevill-Manning et al. (1998)
GIBBS (Gibbs sampler)	ftp://ncbi.nlm.nih.gov/pub/neuwald/gibbs9_95/	Lawrence et al. (1993) Liu et al. (1995) Neuwald et al. (1995)
HMMER hidden Markov model	http://hmmcr.wustl.edu/	Eddy (1998)
MACAW	ftp://ncbi.nlm.nih.gov/pub/macaw	Schuler et al. (1991)
MEME (EM method)	http://meme.sdsc.edu/meme/website/	Bailey and Elkan (1995) Grundy et al. (1996, 1997) Bailey and Gribskov (1998)
Profile (UCSD)	http://www.sdsc.edu/projects/profile/	Gribskov and Veretnik (1996)
SAM hidden Markov model	http://www.cse.ucsc.edu/research/comp/bio/sam.html	Krogh et al. (1994) Hughey and Krogh (1996)



Exact Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

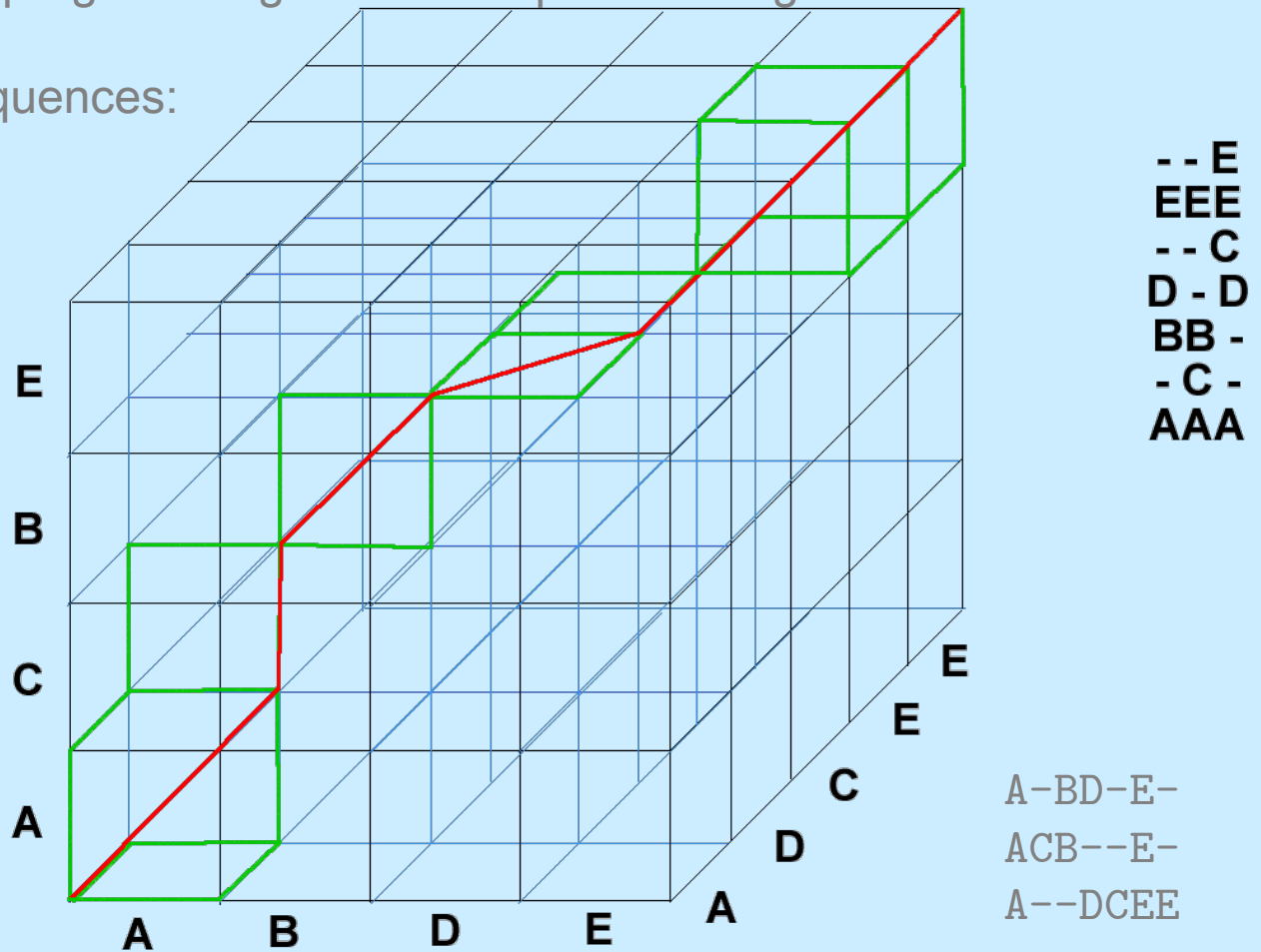
4.3.2 Progressive

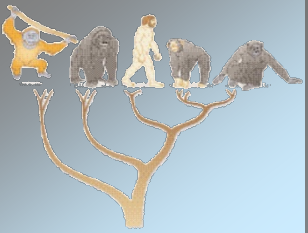
4.3.3 Other

4.4 Profiles / PSSMs

MSA (Lippman et al., 1989, Gupa et al., 1995): generalizes the dynamic programming ideas from pairwise alignment

three sequences:





Exact Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

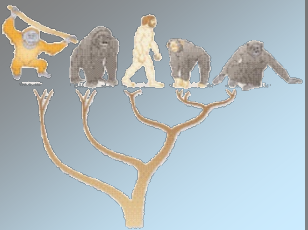
4.4 Profiles / PSSMs

memory and computational complexity: exponentially with N

Gupa et al., 1995: pairwise alignments constrain the path and not the whole hypercube must be filled

MSA (Gupa):

1. compute all pairwise alignment scores $S_{k,l}$
2. predict a phylogenetic tree based on the pairwise scores
3. compute pairwise weights based on the tree
4. construct a temporary multiple alignment with score S_t
5. Compute $B_{k,l}$, a lower bound on $S[k,l]$, the score of the projection of the optimal multiple alignment to k and l
6. Compute space constraints similar to the Baum-Welch
7. compute the optimal alignment on the constraint cube; Dijkstra's shortest path algorithm for nonnegative edges; priority queue; non-negativity guarantees monotone increasing **costs**
8. compare the weight in the alignment with the maximal weight



Exact Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

last step compares actual and maximal weight, if actual is larger then a better alignment may be possible, larger maximal weight means more computational costs

Carillo-Lipman bound:

$$B_{k,l} = S_t + S_{k,l} - \sum_{i,j} S_{i,j}$$

$$S \geq S_t$$

$$\Leftrightarrow \sum_{i,j} S[i,j] \geq S_t$$

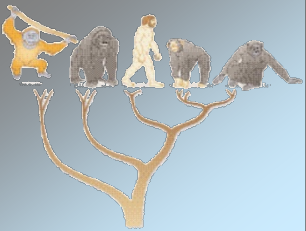
$$\Rightarrow \sum_{(i,j) \neq (k,l)} S_{i,j} + S[k,l] \geq S_t$$

$$\Leftrightarrow S[k,l] \geq S_t - \sum_{(i,j) \neq (k,l)} S_{i,j}$$

$$\Leftrightarrow S[k,l] \geq S_t + S_{k,l} - \sum_{i,j} S_{i,j}$$

$$\Leftrightarrow S[k,l] \geq B_{k,l}$$

$$S[k,l] \leq S_{k,l}$$
$$S_t \leq S$$



Exact Methods

4 Multiple Alignment

MSA improved by the A^* algorithm (Lermen and Reinert, 1997)

4.1 Motivation

A^* -algorithm

Lower bound

4.2 Scoring

Input: graph (the graph), start (start node), goal (goal node), $h(s)$ approximation of the distance of node s to the goal, S (priority queue), N (list of visited nodes)

4.2.1 Consensus

Output: list P of the shortest path

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

BEGIN FUNCTION

insert (start,S)

while not isEmpty(S) **do**

current_node = pop(S)

if current_node in N **then** {no path from start to goal}
return "no path"

end if

insert (current_node, N)

if current_node = goal **then**

reconstruct_shortest_path(start,goal, graph)

else {find all nodes accessible from current node}

successors = expand(current_node, graph)

save_predecessor_in_graph(current_node, graph)

for all s in successors **do** {save node which lead to s}

predecessor(s) = current_node {compute and store costs}

cost(s) = cost(current_node) + edge(graph,current_node,s)

all_cost(s) = cost(s) + h(s)

insert(s,S) {according to all_cost(s)}

end for

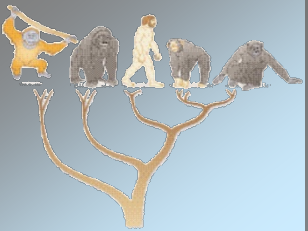
end if

end while

return "no path found"

END FUNCTION

BEGIN SUBFUNCTION {shortest path P as list}
reconstruct_shortest_path (start, node, graph)
if node not= start **then**
push(node, P) {get predecessor}
predecessor = getPredecessor(node, graph)
reconstruct_shortest_path (start, predecessor, graph)
else
return P
end if
END SUBFUNCTION



Exact Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

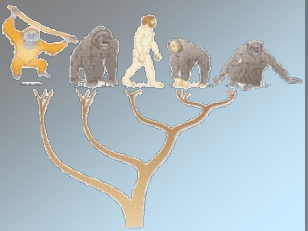
MSA: weighted sum of pairs and a linear gap penalty

Weight: difference pairwise and projected multiple alignment (larger difference means higher weight)

similar sequences: pull the multiple alignment towards them which down-weights them

weights through the phylogenetic tree remove weights between distant sequences

Summing up all the weights: overall divergence of the sequences



Progressive Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

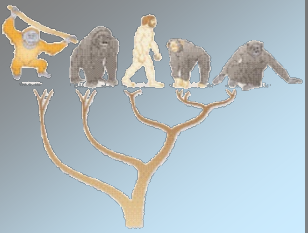
Progressive methods are the most popular methods for multiple alignment: ClustalW (Thomson, Higgins, Gibson, 1994) and TCOFFEE (Notredame, Higgins, Heringa, 2000)

ClustalW and TCOFFEE:

- ↳ perform pairwise alignment for each pair
- ↳ weight matrix: one minus the ratio of perfect matches
- ↳ construct a phylogenetic tree (Neighbor-Joining method)
- ↳ alignments between pairs sequences/alignments (start with closest distance); alignments are propagated through the tree

Initial alignments may be found through local alignment

phylogenetic tree supplies the weighting factors



Progressive Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

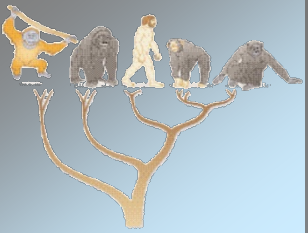
Disadvantage progressive methods:

- ↳ local minima
- ↳ same scoring matrix for close and remote related sequences and same gap parameters

ClustalW

gap penalties context dependent:

- ↳ gaps in hydrophobic regions are more penalized
- ↳ gaps which are within eight amino acids to other gaps are more penalized
- ↳ gaps in regions of other gaps have lower gap opening penalty
- ↳ gap penalties are amino acid dependent



Progressive Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

scoring matrices are adapted:

↳ scoring matrix from the PAM or the BLOSUM families

sequences are weighted through a phylogenetic tree:

↳ similar sequences lower weights (unbalanced data sets)

↳ phylogenetic tree weights with w_i as the weight of sequence i

$$\sum_{i=1}^{N-1} \sum_{j=i+1}^N w_i w_j s(i, j)$$

adaptive phylogenetic tree:

↳ insufficient scores change the tree

initial gap penalty parameters:

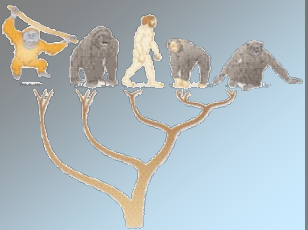
↳ according to scoring matrix

↳ similarity of the sequences (% identity)

↳ length of the sequences (log of the shorter sequences is added)

↳ difference of the length to avoid gaps in the shorter sequence

$$\cdot (1 + |\log(n/m)|)$$



Progressive Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

TCoffee (Tree based Consistency Objective Function For alignmEnt Evaluation) often better alignment than clustalW

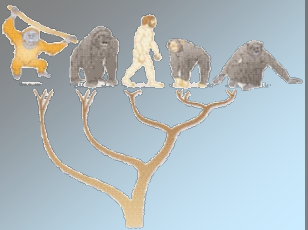
TCoffee work as follows:

↳ libraries of pairwise alignments based on both global (clustalW) and local (FASTA) alignments (combination is more reliable)

↳ library weights are computed according to % identity

↳ libraries are combined and extended; arithmetic mean of weights; extension by aligning two sequences through a third sequence

↳ progressive alignment with a distance based on extended library



Other Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Center Star Alignment

center sequence \bar{i} : $\bar{i} = \arg \min_i \sum_j C_{i,j}$

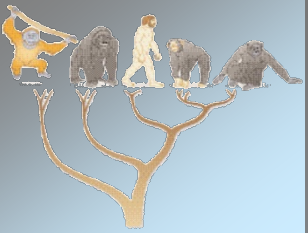
pairwise alignment costs $C_{i,j}$

$$\bar{i} = 1$$

new sequence is added to the set of aligned sequences by a pairwise alignment to the center sequence introducing new gaps

Therefore for center star cost C with projection $C(i,j)$:

$$C(1,j) = C_{1,j}$$



Other Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Gusfield, 1993: cost is less than twice the optimal cost, if

$$C(i, i) = 0 \quad \text{and} \quad C(i, j) \leq C(i, k) + C(k, j)$$

scoring matrix s with

$$s(-, -) = 0$$

$$s(-, i) < 0$$

$$s(k, k) \geq s(i, k) + s(k, j) - s(i, j)$$

A	B		A	B
		>		
A	C		C	A

Then $C(i, j) = S_{i,i} - 2 S_{i,j} + S_{j,j}$ fulfills above conditions

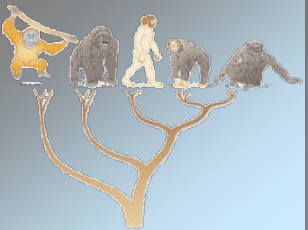
Then the condition $C(i, j) \leq C(i, k) + C(k, j)$ is equivalent to

$$S_{i,i} - 2 S_{i,j} + S_{j,j} \leq S_{i,i} - 2 S_{i,k} + S_{k,k} +$$

$$S_{k,k} - 2 S_{k,j} + S_{j,j}$$

$$\Leftrightarrow S_{i,j} \geq S_{i,k} + S_{k,j} - S_{k,k}$$

$$\Leftrightarrow S_{k,k} \geq S_{i,k} + S_{k,j} - S_{i,j}$$



Other Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

align i to k and j to k then align i , j , and k based on the pairwise alignments, the alignment has a gap if a gap was in one alignment

S is score of the multiple alignment

Per construction: $S[i, k] = S_{i,k}$, $S[k, j] = S_{k,j}$ and $S[k, k] = S_{k,k}$

Componentwise holds: $s(i, j) \geq s(i, k) + s(k, j) - s(k, k)$

Therefore $S[i, j] \geq S[i, k] + S[k, j] - S[k, k]$ and

$$S[i, j] \geq S_{i,k} + S_{k,j} - S_{k,k}$$

inequality to show follows from $S_{i,j} \geq S[i, j]$

→triangle inequality for costs shown by previous equivalences



Other Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

idea of the proof of Gusfield center sequence alignment with cost C and the optimal cost C^*

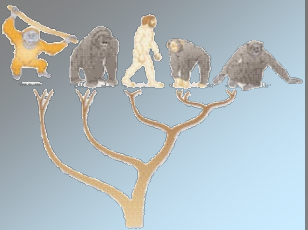
$$C = \sum_{i=1}^N \sum_{j=1, j \neq i}^N C(i, j) \leq$$

$$\sum_{i=1}^N \sum_{j=1, j \neq i}^N C(i, 1) + C(1, j) = 2(N-1) \sum_{i=2}^N C_{i,1}$$

$$C^* = \sum_{i=1}^N \sum_{j=1, j \neq i}^N C^*(i, j) \geq \sum_{i=1}^N \sum_{j=1, j \neq i}^N C_{i,j} \geq$$

$$\sum_{i=1}^N \sum_{j=2}^N C_{i,1} = N \sum_{i=2}^N C_{i,1}$$

$$\Rightarrow \frac{C}{C^*} \leq \frac{2(N-1)}{N} \leq 2$$



Other Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

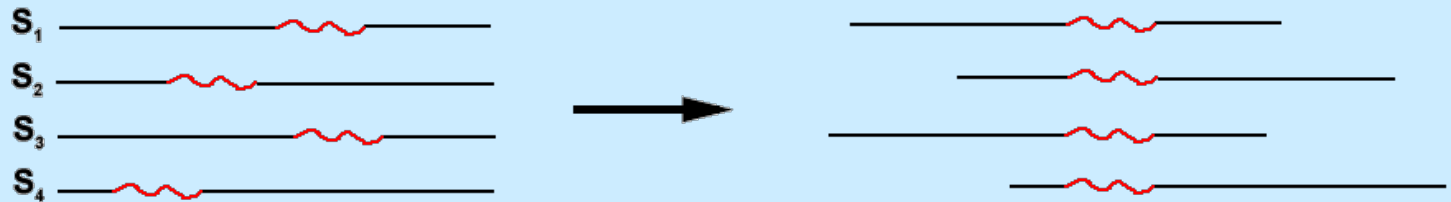
4.3.1 Exact Methods

4.3.2 Progressive

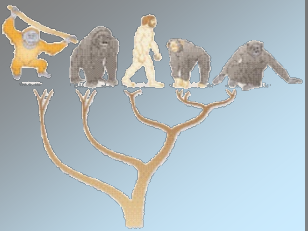
4.3.3 Other

4.4 Profiles / PSSMs

Motifs or pattern can be superimposed for alignment landmarks



Profiles and blocks can be derived from multiple alignments



Other Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

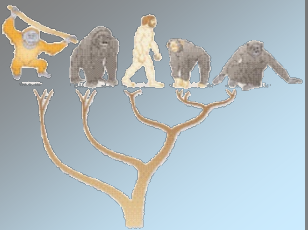
4.4 Profiles / PSSMs

SAGA (Sequence Alignment by Genetic Algorithm): genetic algorithm

MSASA (Multiple Sequence Alignment by Simulated Annealing): simulated annealing

Gibbs sampling

HMMs (hidden Markov models) can be used to find motifs



Other Methods

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

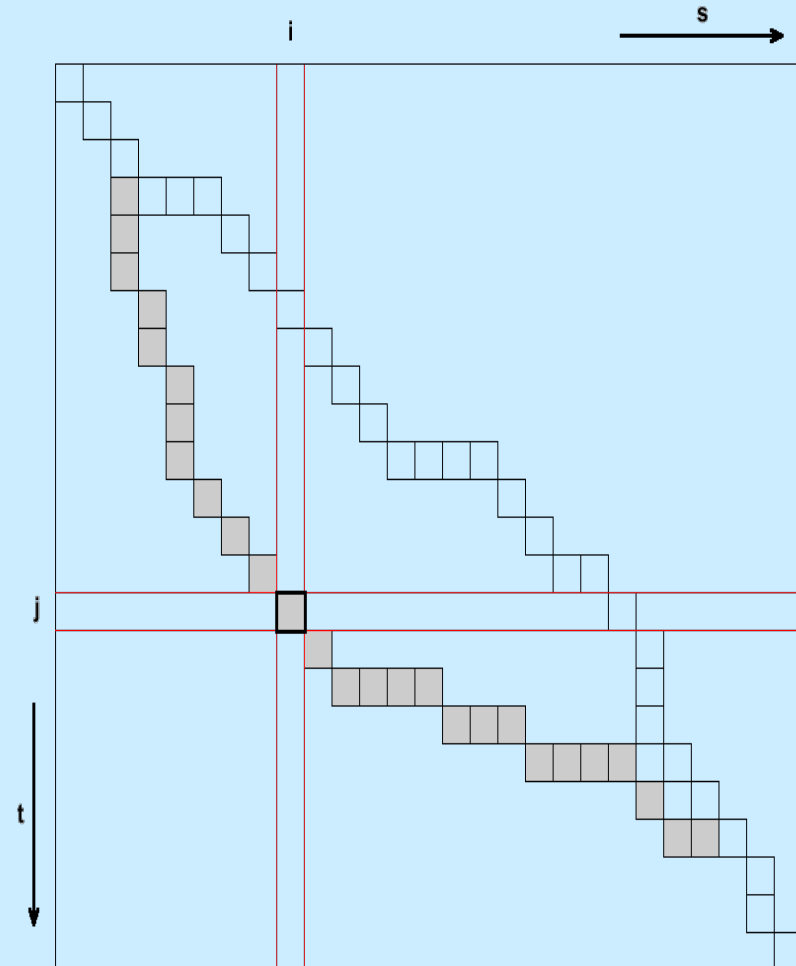
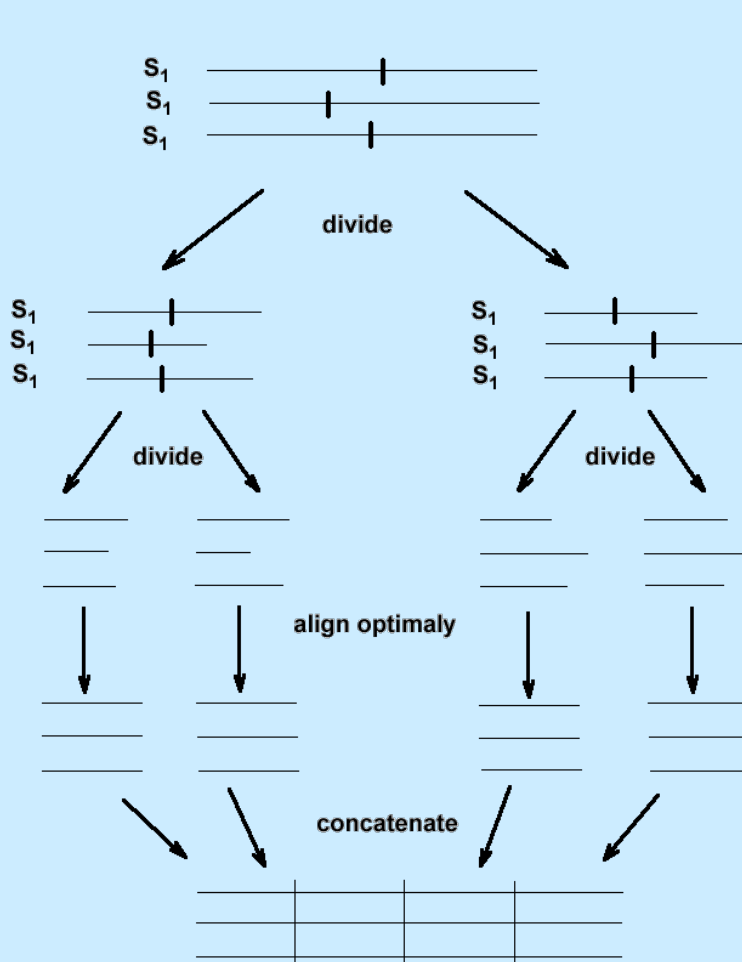
4.3.1 Exact Methods

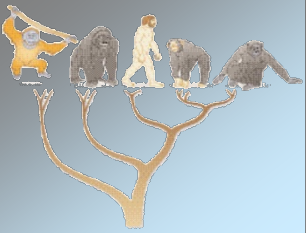
4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Divide-and-conquer Algorithms





Profiles and PSSMs

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

Profiles and Position Specific Scoring Matrices

Assumptions:

- ↪ \mathbf{x} is i.i.d. in its elements according to p_x
- ↪ n the length of \mathbf{x} is large
- ↪ expected letter score for random sequences $\sum_i p_x(i) s(i) < 0$
- ↪ exist i for which $s(i) > 0$

$$S_n = \sum_{i=1}^n s(i) \quad \text{centered value: } \tilde{S}_n = S_n - \frac{\ln n}{\lambda}$$

$$P(\tilde{S}_n > S) \approx 1 - \exp(-K e^{-\lambda S}) \approx K e^{-\lambda S}$$

$$\sum_i p_x(i) \exp(\lambda s(i)) = 1$$



Profiles and PSSMs

4 Multiple Alignment

4.1 Motivation

4.2 Scoring

4.2.1 Consensus

4.2.2 Tree and Star

4.2.3 Sum of Pairs

4.3 Algorithms

4.3.1 Exact Methods

4.3.2 Progressive

4.3.3 Other

4.4 Profiles / PSSMs

q_i : frequency of a letter a_i in a column of a multiple alignment

for sufficient high scoring segments

$$\lim_{n \rightarrow \infty} q_i = p_x(i) \exp(\lambda s(i))$$

$$\Rightarrow s(i) = \ln \left(\frac{q_i}{p_x(i)} \right) / \lambda$$

“Position Specific Scoring Matrices” (PSSMs) or profiles

new sequence: high scores mean similar alignment sequences