

- ◆ Multiple sequence alignments
  - Definition
  - The need for MSA
  - The MSA problem
  - MSA methods
  
- ◆ Editing and formatting alignments
  - Software packages available

MSA definition

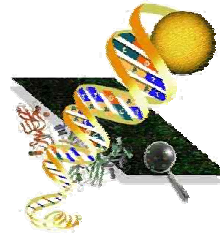
The need for MSA

The MSA problem

MSA methods

Sequence editors

# What is Multiple Sequence Alignment (MSA)?



- ◆ Multiple sequence alignment (MSA) can be seen as a generalization of Pairwise Sequence Alignment - instead of aligning two sequences,  $n$  sequences are aligned simultaneously, where  $n$  is  $> 2$

- ◆ **Definition:**

A multiple sequence alignment is an alignment of  $n > 2$  sequences obtained by inserting gaps ("-") into sequences such that the resulting sequences have all length  $L$  and can be arranged in a matrix of  $N$  rows and  $L$  columns where each column represents a homologous position

- ◆ **Note:**

MSA applies both to nucleotide and amino acid sequences

To construct a multiple alignment, one may have to introduce gaps in sequences at positions where there were no gaps in the corresponding pairwise alignment

→ multiple alignments typically contain more gaps than any given pair of aligned sequences

**MSA definition**

The need for MSA

The MSA problem

MSA methods

Sequence editors

XESEE.EXE

File Edit Search foRmat Command Options Page 1 Seq 1 Pos 1

Am\_li\_1060 N 1▶ ccatgcATGTCTAAGTTCACACTCTGGTACAGTGAAACCGCGAATGGCTCATTAAatCAG  
Me\_em\_988 N 2▶ ccatgcATGTCTAAGTTCACACCTCGTATGGTGAACCGCGAATGGCTCATTAAATCAG  
By\_au\_1002 N 3▶ ccatgcATGTCTAAGTTCACACTCTCGTACGGTGAACCGcgAATGGCTCATTAAATCAG  
Ma\_sc\_2599 N 4▶ ccatGCATGTCTAAGTTCACACTcTCGTACGGTGAACCGcgAATGGCTCATTAAATCAG  
Er\_sp\_643 N 5▶ ccatgcATGTCTAAGTACATACCTTTAAACGGTGAACCGcgAATGGCTCATTAAATCAG  
Mo\_ki\_1364 N 6▶ ccatgcATGTCTAAGTTCCTACTCTCGCACGGTGAACCGcgAATGGCTCATTAAATCAG  
\_\_\_\_\_10<sup>↓</sup>\_\_\_\_\_20<sup>↓</sup>\_\_\_\_\_30<sup>↓</sup>\_\_\_\_\_40<sup>↓</sup>\_\_\_\_\_50<sup>↓</sup>\_\_\_\_\_60<sup>↓</sup>

Am\_li\_1060 N 1▶ TCGAGGTTCTTTGATGATCCAAATCTACTTGGATAACTGTGGTAATTCTAGAGCTAATA  
Me\_em\_988 N 2▶ TCGAGGTTCTTTAGATGATCCAAATCTACTTGGATAACTGTGGTAATTCTAGAGCTAATA  
By\_au\_1002 N 3▶ TCGAGgtteCTTTAGATGATCCAAATCTACTTGGATAACTGTGGTAATTCTAGAGCTAATA  
Ma\_sc\_2599 N 4▶ TCGAGGTTCTTTAGATGATCCAAATCTACTTGGATAACTGTGGTAATTCTAGAGCTAATA  
Er\_sp\_643 N 5▶ CTATGGTTCTTTAGATCGTACCTACTACATGGATAACTGTAGTAATTCTAGAGCTAATAC  
Mo\_ki\_1364 N 6▶ TCGAGGTTCTTTAGATGATCCAAAGCTACTTGGATAACTGTGGTAATTCTAGAGCTAATA  
\_\_\_\_\_70<sup>↓</sup>\_\_\_\_\_80<sup>↓</sup>\_\_\_\_\_90<sup>↓</sup>\_\_\_\_\_100<sup>↓</sup>\_\_\_\_\_110<sup>↓</sup>\_\_\_\_\_120<sup>↓</sup>

Am\_li\_1060 N 1▶ CATGCCTACCAGCTCCGACCCGGTGGGCCTCGTTTCGGCTTTCCCTGTACAGGGGGGAG  
Me\_em\_988 N 2▶ CATGCCCAACCGCTCCGACCTGTAAGGAAAGAGCGCTTTTATCAGCTCAAACCAGTCT  
By\_au\_1002 N 3▶ CATGCCCAACCGCTCCGACCCCTTCGCAAGGAGGGGAAAGAGCGCTTTTATTAGTTCAA  
Ma\_sc\_2599 N 4▶ CATGCCCAACCGCTCCGACCCGCTTGGGGCCCTCCTCGCAAGGGGGCGGTGCCcGGCGG  
Er\_sp\_643 N 5▶ ATGCCACTATGCCCTGACCCGCAAGGGAACGGGTGGATTTATTAGAACAGAACCAATCGG  
Mo\_ki\_1364 N 6▶ CATGCCCGACAGCTCCGACCGTCGTGCGTAACAGCGGCGGGACGAGCGCTTTTATT  
\_\_\_\_\_130<sup>↓</sup>\_\_\_\_\_140<sup>↓</sup>\_\_\_\_\_150<sup>↓</sup>\_\_\_\_\_160<sup>↓</sup>\_\_\_\_\_170<sup>↓</sup>\_\_\_\_\_180<sup>↓</sup>

Am\_li\_1060 N 1▶ TCGGGTGGGGACTCCGTTGGGGAAGAGCGCTTTTATTAGTTCAAaACCAGTCGGGCTTTC  
Me\_em\_988 N 2▶ GCCGGCTCAAACCAGTCCCTTGGTGAATCTGGATAACTTTTGGCGATCGCATGG  
By\_au\_1002 N 3▶ ACCAGTCGGGCCcTCACGGGTCCGTCCTTGGTGAATCTGGATAACTTTGTGCCGATCG  
Ma\_sc\_2599 N 4▶ GGAAGAGCGCTTTTATTAGTTCAAACCAGTCGGGgTCCCAGCCCcGTCTCTTTGGTG  
Er\_sp\_643 N 5▶ TGGTGGCTTCGGCTGCTGCTGTTGCAATCTGGATGACTCTGGATAACTTCACTGATCGCG  
Mo\_ki\_1364 N 6▶ AGTTGAAAACCAGtcggCCTCGCGGCCGTCCCCTTGGTGAATCTGGATAACTTTGAGCCG  
\_\_\_\_\_190<sup>↓</sup>\_\_\_\_\_200<sup>↓</sup>\_\_\_\_\_210<sup>↓</sup>\_\_\_\_\_220<sup>↓</sup>\_\_\_\_\_230<sup>↓</sup>\_\_\_\_\_240<sup>↓</sup>

XESEE.EXE

File Edit Search foRmat Command Options Page 1 Seq 7 Pos 4

```

Am_li_1060 N 1▶ ccatgcATGTCTAAGTTCACACTCTGGTACAGTGAAACCGCGAATGGCTCATTAAatCAG
Me_em_988 N 2▶ ccatgcATGTCTAAGTTCACACCTCGTATGGTGAACCGCGAATGGCTCATTAAATCAG
By_au_1002 N 3▶ ccatgcATGTCTAAGTTCACACTCTCGTACGGTGAACCGcgAATGGCTCATTAAATCAG
Ma_sc_2599 N 4▶ ccatGCATGTCTAAGTTCACACTcTCGTACGGTGAACCGcgAATGGCTCATTAAATCAG
Er_sp_643 N 5▶ ccatgcATGTCTAAGTACATACCTTTAAACGGTGAACCGcgAATGGCTCATTAAATCAG
Mo_ki_1364 N 6▶ ccatgcATGTCTAAGTTCCTACTCTCGCACGGTGAACCGcgAATGGCTCATTAAATCAG
      10┘ 20┘ 30┘ 40┘ 50┘ 60┘

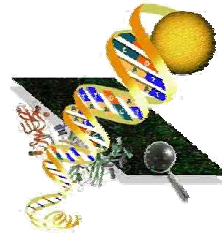
Am_li_1060 N 1▶ TCGAGGTTCTTTGATGATCCAAT--CTACTTGGATAACTGTGGTAATTCTAGAGCTAA
Me_em_988 N 2▶ TCGAGGTTCTTAGATGATCCAAT--CTACTTGGATAACTGTGGTAATTCTAGAGCTAA
By_au_1002 N 3▶ TCGAGgtteCTTAGATGATCCAAT--CTACTTGGATAACTGTGGTAATTCTAGAGCTAA
Ma_sc_2599 N 4▶ TCGAGGTTCTTAGATGATCCAAT--CTACTTGGATAACTGTGGTAATTCTAGAGCTAA
Er_sp_643 N 5▶ CTATGGTTCTTAGATCGTA-CCTA--CTACATGGATAACTGTAGTAATTCTAGAGCTAA
Mo_ki_1364 N 6▶ TCGAGGTTCTTAGATGATCCAAG--CTACTTGGATAACTGTGGTAATTCTAGAGCTAA
      70┘ 80┘ 90┘ 100┘ 110┘ 120┘

Am_li_1060 N 1▶ TACATGCCTACCAGCTCCGACCCG-GGG-----AAGAGCGCTTTTATTAGTTCAA-AaC
Me_em_988 N 2▶ TACATGCCCAACCGCTCCGACCTA-GGA-----AAGAGCGCTTTTATCAGCTCAA-AAC
By_au_1002 N 3▶ TACATGCCCACCAGCTCCGACCCG-GGA-----AAGAGCGCTTTTATTAGTTCAA-AAC
Ma_sc_2599 N 4▶ TACATGCCCAACAGCTCCGACCCGCCcGGCGGGGAAGAGCGCTTTTATTAGTTCAA-AAC
Er_sp_643 N 5▶ TACATGCCACTATGCCCTGACCCG-GGA-----ACGGGTGGATTTATTAGAACAG-AAC
Mo_ki_1364 N 6▶ TACATGCCCGACAGCTCCGACCCGGCGGCGGGACGAGCGCTTTTATTAGTTGAA-AAC
      130┘ 140┘ 150┘ 160┘ 170┘ 180┘

Am_li_1060 N 1▶ CAGTCGG-TCCTTTTG-GT---GACTCTG--G-ATAACTTTGTGCCGATCGCATCGGTC
Me_em_988 N 2▶ CAGTCTG-TCCCTT--G-GT---GAATCTG--G-ATAACTTTTTGCCGATCGCA-TGGC-
By_au_1002 N 3▶ CAGTCGG-TCCCCTT-G-GT---GACTCTG--G-ATAACTTTGTGCCGATCGCA-CGGC-
Ma_sc_2599 N 4▶ CAGTCGG-TCCTCTTTG-GT---GACTCTG--G-ATAACTTTGTGCCGATCGCA-CGGC-
Er_sp_643 N 5▶ CAATCGG-GCAATCT-GGAT---GACTCTG--G-ATAACTTCA--CTGATCGCGTCGGC-
Mo_ki_1364 N 6▶ CAGtegg-TCCCCTT-G-GT---GACTCTG--G-ATAACTTTGAGCCGATCGCA-CGGC-
      190┘ 200┘ 210┘ 220┘ 230┘ 240┘

```

# Why do we need MSA?



- ◆ Multiple sequence alignment can help to develop a sequence “finger print” which allows the identification of members of distantly related protein family (motifs)
- ◆ Formulate & test hypotheses about protein 3-D structure
- ◆ MSA can help us to reveal biological facts about proteins, e.g.:  
(e.g. how protein function has changed or evolutionary pressure acting on a gene)
- ◆ Crucial for genome sequencing:
  - Random fragments of a large molecule are sequenced and those that overlap are found by a multiple sequence alignment program.
  - There should be one correct alignment that corresponds to the genomic sequence rather than a range of possibilities
  - Sequence may be from one strand of DNA or the other, so complements of each sequence must also be compared
  - Sequence fragments will usually overlap, but by an unknown amount and in some cases, one sequence may be included within another
  - All of the overlapping pairs of sequence fragments must be assembled into large composite genome sequence
- ◆ To establish homology for phylogenetic analyses
- ◆ Identify primers and probes to search for homologous sequences in other organisms

MSA definition

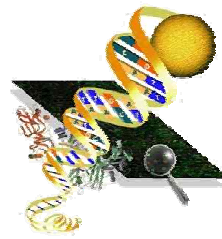
The need for MSA

The MSA problem

MSA methods

Sequence editors

# The alignment problem



- ◆ How do we generate a multiple alignment? Given a pairwise alignment, just add the third, then the fourth, and so on, until all have been aligned. Does it work?

Example:

It is not self-evident how these sequences are to be aligned together. Here are some possibilities:

Taxon A AGAC  
Taxon B --AC

Taxon A AGAC  
Taxon C AG--

Taxon B AC  
Taxon C AG

Taxon A AGAC  
Taxon B --AC  
Taxon C AG--

Taxon A AGAC  
Taxon C AG--  
Taxon B --AC

Taxon B AC--  
Taxon C AG--  
Taxon A AGAC

Taxon B --AC  
Taxon C --AG  
Taxon A AGAC

- ◆ It depends not only on the various alignment parameters but also on the order in which sequences are added to the multiple alignment

MSA definition

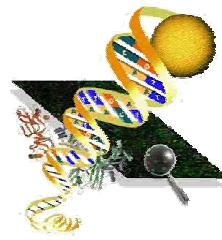
The need for MSA

The MSA problem

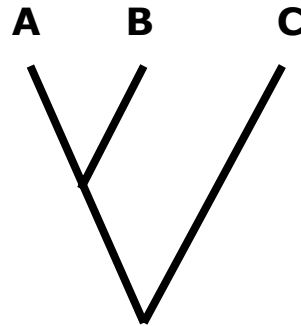
MSA methods

Sequence editors

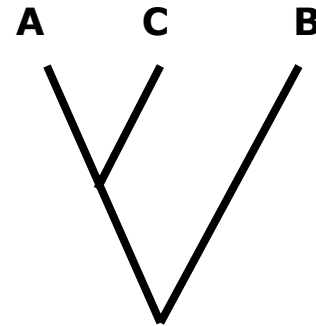
# The alignment problem



- ◆ What happens when a sequence alignment is wrong?

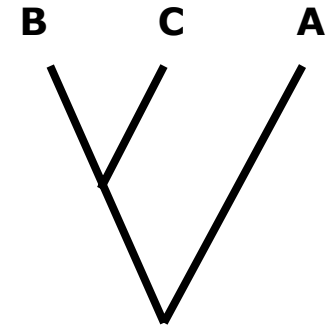


A: AGT  
B: AT  
C: ATC



A: AGT  
B: A -T  
C: ATC

A: AGT  
B: AT -  
C: ATC



A: AGT -  
B: A -T -  
C: A -TC

MSA definition

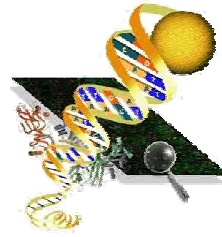
The need for MSA

**The MSA problem**

MSA methods

Sequence editors

# From pairwise to multiple alignments



- ◆ In pairwise alignments, one has a two-dimensional matrix with the sequences on each axis. The number of operations required to locate the best “path” through the matrix is approximately proportional to the product of the lengths of the two sequences
- ◆ A possible general method would be to extend the pairwise alignment method into a simultaneous N-wise alignment, using a complete dynamical-programming algorithm in N dimensions. Algorithmically, this is not difficult to do

**But what about execution time?**

MSA definition

The need for MSA

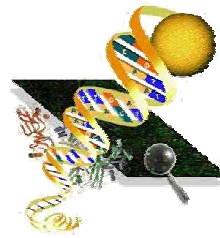
**The MSA problem**

MSA methods

Sequence editors



# The big-O notation



- ◆ One of the most important properties of an algorithm is how its execution time increases as the problem is made larger (e.g. more sequences to align). This is the so-called **algorithmic** (or computational) **complexity** of the algorithm
- ◆ There is a notation to describe the algorithmic complexity, called **the big-O notation**. If we have a problem size (number of input data points)  $n$ , then an algorithm takes  **$O(n)$**  time if the time increases linearly with  $n$ . If the algorithm needs time proportional to the square of  $n$ , then it is  **$O(n^2)$**
- ◆ It is important to realize that an algorithm that is quick on small problems may be totally useless on large problems if it has a bad  $O()$  behavior. As a rule of thumb one can use the following characterizations, where  $n$  is the size of the problem, and  $c$  is a constant:

$O(c)$	utopian
$O(\log n)$	excellent
$O(n)$	very good
$O(n^2)$	not so good
$O(n^3)$	pretty bad
$O(c^n)$	disaster

MSA definition

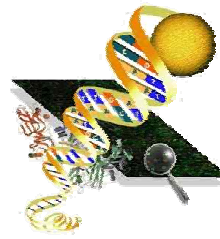
The need for MSA

**The MSA problem**

MSA methods

Sequence editors

# The big-O notation



- ◆ To compute a N-wise alignment, the algorithmic complexity is something like  $O(c^{2n})$ , where  $c$  is a constant, and  $n$  is the number of sequences
- ◆ Example:  
A pairwise alignment of two sequences [ $O(c^{2 \times 2})$ ], takes 1 second, then four sequences [ $O(c^{2 \times 4})$ ], would take  $10^4$  seconds (2.8 hours), five sequences [ $O(c^{2 \times 5})$ ],  $10^6$  seconds (11.6 days), six sequences [ $O(c^{2 \times 6})$ ],  $10^8$  seconds (3.2 years), seven sequences [ $O(c^{2 \times 7})$ ],  $10^{10}$  seconds (317 years), and so on

**This is disastrous!**

MSA definition

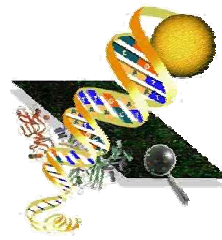
The need for MSA

**The MSA problem**

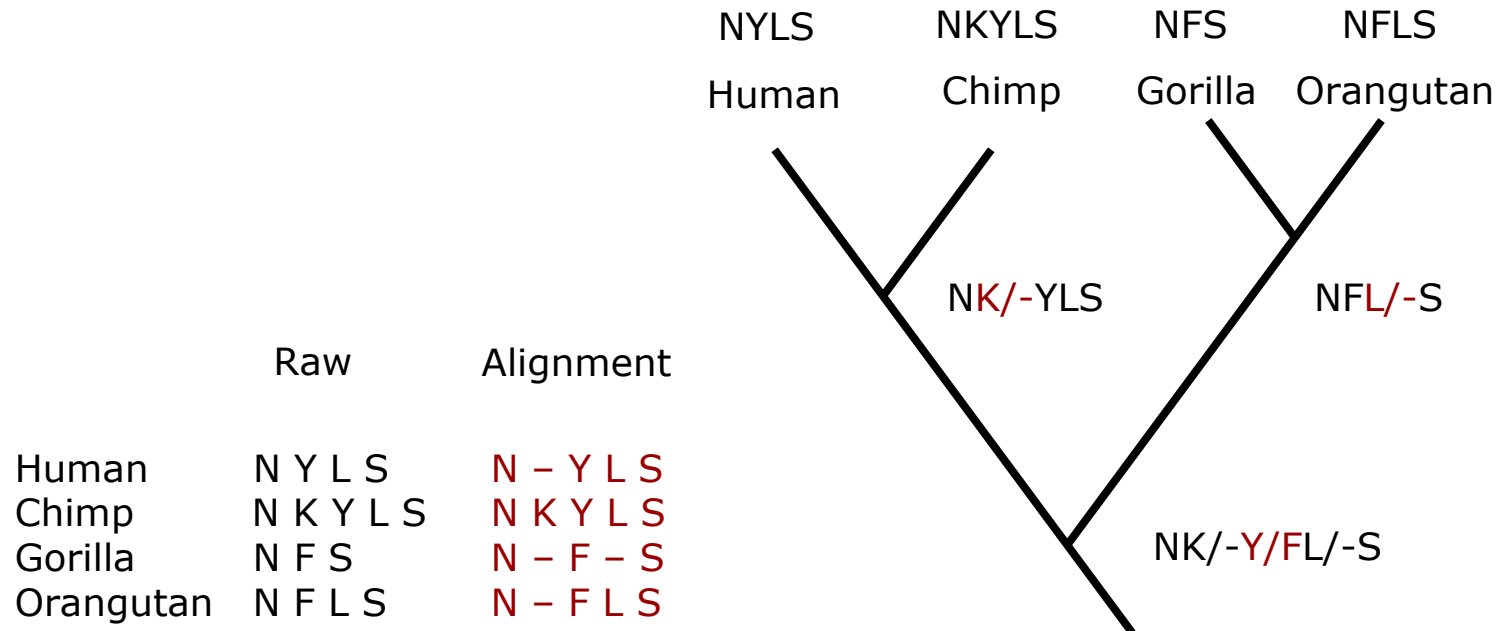
MSA methods

Sequence editors

# How to optimize alignment algorithms?



- ◆ Use structural information:
  - reading frame
  - protein structure
- ◆ Sequence elements are not truly independent but related by phylogeny



- ◆ Sequences often contain highly conserved regions

MSA definition

The need for MSA

**The MSA problem**

MSA methods

Sequence editors

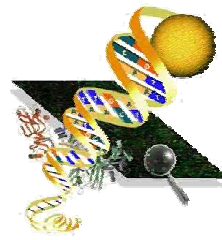
60 total sequences

Mode: Select / Slide Selection: 0 Position: 352 Sequence Mask: None Numbering Mask: None Start ruler at: 1

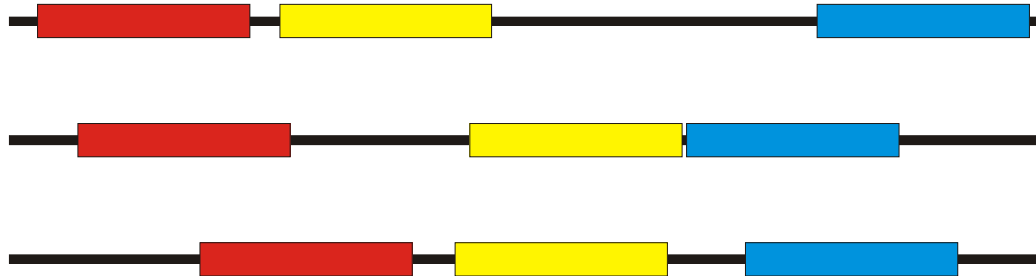
Rich toolbar with icons for alignment, editing, and visualization.

	320	330	340	350	360	370																																																														
Adriohydrobi	C	T	A	C	G	T	T	A	C	G	T	T	T	G	A	T	C	C	A	A	A	C	A	T	T	T	G	A	T	T	A	A	C	A	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A					
Alvania	A	T	C	C	A	A	T	A	A	T	A	T	T	G	A	T	T	A	A	A	A	A	A	T	A	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	G	T	A	A	T	T	T	T	C	T	T	T	A	A	G	A	G	T	T			
Amnicola 106	T	A	A	A	T	A	T	A	A	A	T	A	A	T	G	A	T	C	C	A	A	A	T	A	T	T	T	G	A	T	T	A	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T						
Antroselates	C	A	T	G	C	A	T	A	A	T	G	A	T	C	C	A	A	A	T	A	T	T	T	G	A	T	T	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A	A	T	C	T	T					
Ascorhis	T	T	A	T	A	A	A	A	G	T	A	G	A	T	T	G	A	C	T	T	T	G	A	T	C	C	A	T	A	A	G	T	T	T	T	G	A	T	T	A	A	C	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C		
Baikalia	T	A	A	A	A	A	C	T	T	A	T	A	T	G	T	A	T	T	G	A	T	C	C	A	A	A	T	A	T	T	T	G	A	T	T	A	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G					
Beddomeia	A	A	G	T	T	T	G	T	G	C	A	T	T	T	G	T	G	C	T	A	T	G	A	T	C	C	A	A	A	G	T	T	T	T	G	A	T	T	A	A	A	G	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C			
Bithynia	T	A	A	A	C	C	T	T	A	T	T	A	G	C	A	A	T	G	A	T	C	C	G	A	A	T	A	T	T	C	G	A	T	T	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C							
Bythinella 1	A	A	T	A	C	C	A	A	T	A	A	G	T	A	A	T	G	A	T	C	C	A	A	A	A	C	T	G	A	T	T	A	A	A	A	A	A	T	T	A	G	C	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A					
Cincinnatia	A	C	T	T	T	T	A	A	G	T	T	T	G	A	T	C	C	A	A	A	A	T	T	T	G	A	T	T	A	A	C	A	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A							
Coxiella	T	A	A	A	A	A	T	A	C	T	A	A	T	A	A	G	T	T	T	G	A	T	C	C	A	G	A	A	A	T	C	T	G	A	T	T	A	A	A	A	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A				
Eatoniella	T	A	A	T	A	A	G	A	T	C	C	C	T	T	A	A	T	T	A	A	A	G	G	A	C	T	A	A	C	A	G	A	A	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A	A	T	T	T	T					
Emmericia	A	T	A	C	T	A	A	T	A	G	G	T	A	A	T	G	A	T	C	C	A	A	C	A	A	G	G	T	T	T	G	A	T	T	A	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T			
Emmericia 30	A	T	A	C	T	A	A	T	A	G	G	T	A	A	T	G	A	T	C	C	A	A	C	A	A	G	G	T	T	T	G	A	T	T	A	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T			
Erhaia 652	T	T	A	A	T	A	T	A	A	A	C	A	C	C	T	A	T	A	T	G	T	A	T	T	G	A	T	C	C	A	A	A	T	T	A	T	T	T	G	A	T	T	A	A	A	A	A	A	C	T	A	G	T	T	A	C	C	G	T	A	G	G						
Fairbankia	T	T	A	T	T	T	A	T	A	T	T	T	G	A	A	T	T	A	G	G	G	T	A	G	C	T	A	C	C	T	A	G	G	G	A	T	A	A	C	A	A	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T	T
Fluvidona	A	A	T	A	T	T	A	T	A	G	C	C	T	A	T	T	A	G	C	T	A	T	G	A	T	C	C	A	A	A	T	T	T	T	T	T	G	A	T	C	A	A	G	A	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A			
Fluvipupa	T	T	T	A	T	A	A	A	C	C	T	T	G	A	T	C	C	A	A	A	A	A	T	T	T	G	A	T	T	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A	A	T	C					
Gammatrixula	A	A	A	A	T	A	T	G	C	T	C	A	T	A	A	G	C	T	A	T	G	A	T	C	C	A	A	A	A	A	T	T	T	G	A	T	T	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A					
Geomelania 8	T	A	A	A	T	A	A	A	C	T	A	A	T	A	A	G	T	T	T	A	G	A	T	C	C	A	G	A	A	G	A	T	T	C	T	G	A	T	T	A	A	A	A	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C		
Geomelania 8	T	A	A	A	T	A	A	A	C	T	A	A	T	A	A	G	T	T	T	A	G	A	T	C	C	A	G	A	A	T	T	T	C	T	G	A	T	T	A	A	A	A	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C			
Graziana 256	T	A	A	A	C	T	C	A	T	A	A	G	T	T	T	T	G	A	T	C	C	A	A	A	A	T	T	T	G	A	T	T	A	A	C	A	G	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A			
Hauffenia 25	A	T	A	A	A	C	C	T	A	T	A	A	G	T	T	T	T	G	A	T	C	C	A	A	A	G	T	T	T	T	G	A	T	T	A	A	C	A	G	A	A	t	t	a	a	c	c	g	t	a	g	g	g	a	t	a	a	c	a	g	c	a						
Heleobops	C	T	A	A	T	A	A	G	T	T	T	T	G	A	T	C	C	A	A	A	T	A	T	T	T	G	A	T	T	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A	C	A	G	C	A	T	A	A	T	T					
Hemistomia	A	T	A	T	T	A	A	T	A	A	C	A	T	A	A	G	T	T	T	T	G	A	T	C	C	A	A	A	T	T	T	T	T	G	A	T	T	A	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T	A	A						
Heterocyclus	C	A	T	T	T	T	T	A	A	T	A	C	T	T	A	C	A	A	G	T	C	A	T	G	A	T	C	C	A	A	A	T	T	T	T	T	G	A	T	T	A	A	A	A	A	T	T	A	G	T	T	A	C	C	G	T	A	G	G	G	A	T						

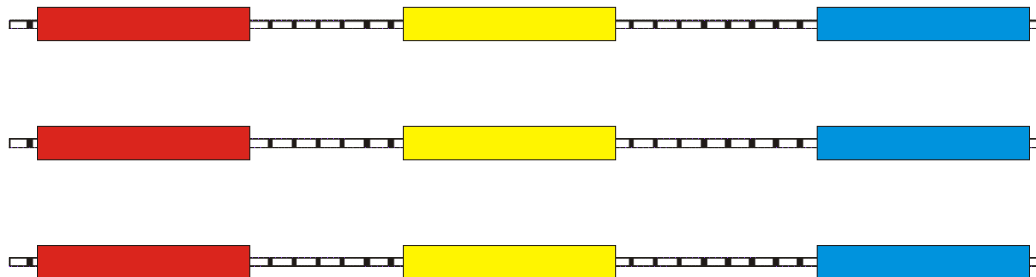
# How to optimize alignment algorithms?



- ◆ Sequences often contain highly conserved regions



**These regions can be used for an initial alignment**



By analyzing a number of small, independent fragments,  
the algorithmic complexity can be drastically reduced!

MSA definition

The need for MSA

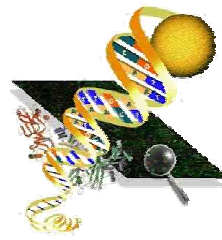
**The MSA problem**

MSA methods

Sequence editors

	140	150	160	170	180	190	200	210	220	230																			
Adrio hydrobi	C	C	A	A	C	C	G	U	U	C	G	A	C	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	
Adrio insulan	C	C	A	A	C	C	G	U	U	C	A	C	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Alvania	C	A	A	A	C	A	G	C	U	U	C	G	A	C	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A
Alzoniella 2	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Ammnicola 106	C	U	A	C	C	A	G	C	U	U	C	G	A	C	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A
Amphithalamu	C	U	A	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Antroselates	C	U	A	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Ascorhis	A	U	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Assimineae	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Assimineae 16	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Assimineae 22	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Avenionia 22	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Baicalia	C	U	A	A	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	A	U	C	A	G	C	U	C	A	G
Barleeia	C	C	C	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Beddomeia	C	U	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Belgrandia	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Bithynia	C	C	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Bythinella 1	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Bythiospeum	C	U	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Calopia	C	U	A	G	U	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Cecina 2522	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Clenchiella	C	U	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Coxiella	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Eatoniella	G	C	A	A	A	C	C	U	A	G	U	C	U	C	G	U	C	G	U	A	U	C	A	G	C	U	C	A	G
Emmericia	C	U	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Emmericia 30	C	U	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Erhaia 652	C	C	C	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Fairbankia	C	C	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Fissuria 243	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Fluvidona	C	C	C	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Fluvipupa	C	C	C	C	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	A	U	C	A	G	C	U	C	A	G
Fontigens	A	C	C	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Gammatricula	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Geomelania 8	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Geomelania 8	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Graziana 256	C	C	A	A	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Hauffenia 25	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Heleobops	C	U	A	A	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	A	U	C	A	G	C	U	C	A	G
Hemistomia	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Heterocyclus	C	C	G	C	C	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Horatia 2598	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Hydrobia 653	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C
Hydrococcus	C	C	A	C	A	A	G	C	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G
Islamia 2327	C	C	A	A	C	C	G	U	U	C	G	A	C	C	U	U	U	U	A	U	C	A	G	C	U	C	A	G	C

# MSA methods



- ◆ Progressive global alignment of the sequences starting with an alignment of the most alike sequences and then building an alignment by adding more sequences
- ◆ Iterative methods that make an initial alignment of groups of sequences and then revise the alignment to achieve a better result
- ◆ Alignments based on locally conserved patterns found in the same order in the sequences

MSA definition

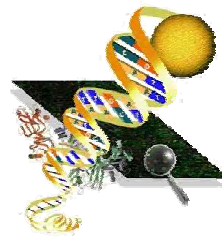
The need for MSA

The MSA problem

**MSA methods**

Sequence editors

# "Optimal" vs. "correct" alignment



- ◆ For a given group of sequences, there is no single "correct" alignment, only an alignment that is "optimal" according to some set of calculations  
This is partly due to:
  - the complexity of the problem,
  - limitations of the scoring systems used,
  - our limited understanding of life and evolution
- ◆ Determining what alignment is best for a given set of sequences is really up to the judgment of the investigator
- ◆ Success of the alignment will depend on the similarity of the sequences. If sequence variation is great it will be very difficult to find an optimal alignment

MSA definition

The need for MSA

The MSA problem

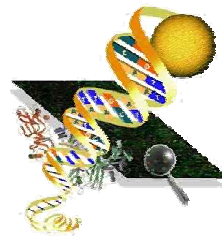
**MSA methods**

Sequence editors





# MSA and gaps



## ◆ Gaps can occur:

Before the first character of a string

```
CTGCGGG---GGTAAT
  ||||      ||  ||
--GCGG-AGAGG-AA-
```

Inside a string

```
CTGCGGG---GGTAAT
  ||||      ||  ||
--GCGG-AGAGG-AA-
```

After the last character of a string

```
CTGCGGG---GGTAAT
  ||||      ||  ||
--GCGG-AGAGG-AA-
```

- ◆ **Note:** In protein-coding nucleotide sequences most gaps have a length of 3N

MSA definition

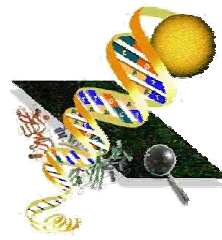
The need for MSA

The MSA problem

**MSA methods**

Sequence editors

# MSA and gaps



## Gap Penalties

- ◆ In the MSA scoring scheme, a penalty is subtracted for each gap introduced into an alignment because the gap increases uncertainty into an alignment
- ◆ The gap penalty is used to help decide whether or not to accept a gap or insertion in an alignment
- ◆ Biologically, it should in general be easier for a sequence to accept a different residue in a position, rather than having parts of the sequence chopped away or inserted. Gaps/insertions should therefore be more rare than point mutations (substitutions)
- ◆ In general, the lower the gapping penalties, the more gaps and more identities are detected but this should be considered in relation to biological significance
- ◆ Most MSA programs allow for an adjustment of gap penalties

MSA definition

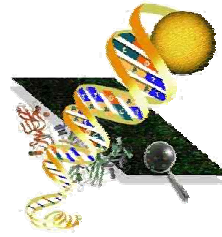
The need for MSA

The MSA problem

**MSA methods**

Sequence editors

# MSA with ClustalW



- ◆ Works by progressive alignment: it aligns a pair of sequences then aligns the next one onto the first pair
- ◆ Most closely related sequences are aligned first, and then additional sequences and groups of sequences are added, guided by the initial alignments  
Uses alignment scores to produce a phylogenetic tree
- ◆ Aligns the sequences sequentially, guided by the phylogenetic relationships indicated by the tree
- ◆ Gap penalties can be adjusted based on specific amino acid residues, regions of hydrophobicity, proximity to other gaps, or secondary structure
- ◆ Is available with a great web interface: <http://www.ebi.ac.uk/clustalw/>
- ◆ Also available as **ClustalX** (stand-alone MS-Windows software)

MSA definition

The need for MSA

The MSA problem

**MSA methods**

Sequence editors

### ClustalW Submission Form

Clustal W is a general purpose multiple sequence alignment program for DNA or proteins. It produces biologically meaningful multiple sequence alignments of divergent sequences. It calculates the best match for the selected sequences, and lines them up so that the identities, similarities and differences can be seen. Evolutionary relationships can be seen via viewing Cladograms or Phylograms.

YOUR EMAIL	ALIGNMENT TITLE	CPU MODE	ALIGNMENT	OUTPUT FORMAT
<input type="text"/>	-NONE-	clustalw_mp	full	aln w/numbers
OUTPUT ORDER	COLOR ALIGNMENT	KTUP (WORD SIZE)	WINDOW LENGTH	SCORE TYPE
aligned	no	def	def	percent
TOPDIAG	PAIRGAP	PHYLOGENETIC TREE	MATRIX	GAP OPEN
def	def	TREE TYPE: none CORRECT DIST.: off IGNORE GAPS: off	def	def
END GAPS	GAP EXTENSION	GAP DISTANCES	TREE TYPE:	TREE GRAPH DISTANCES:
def	def	def	cladogram	hide

} Operational options

} Output options

} Input options, matrix choice, gap opening penalty

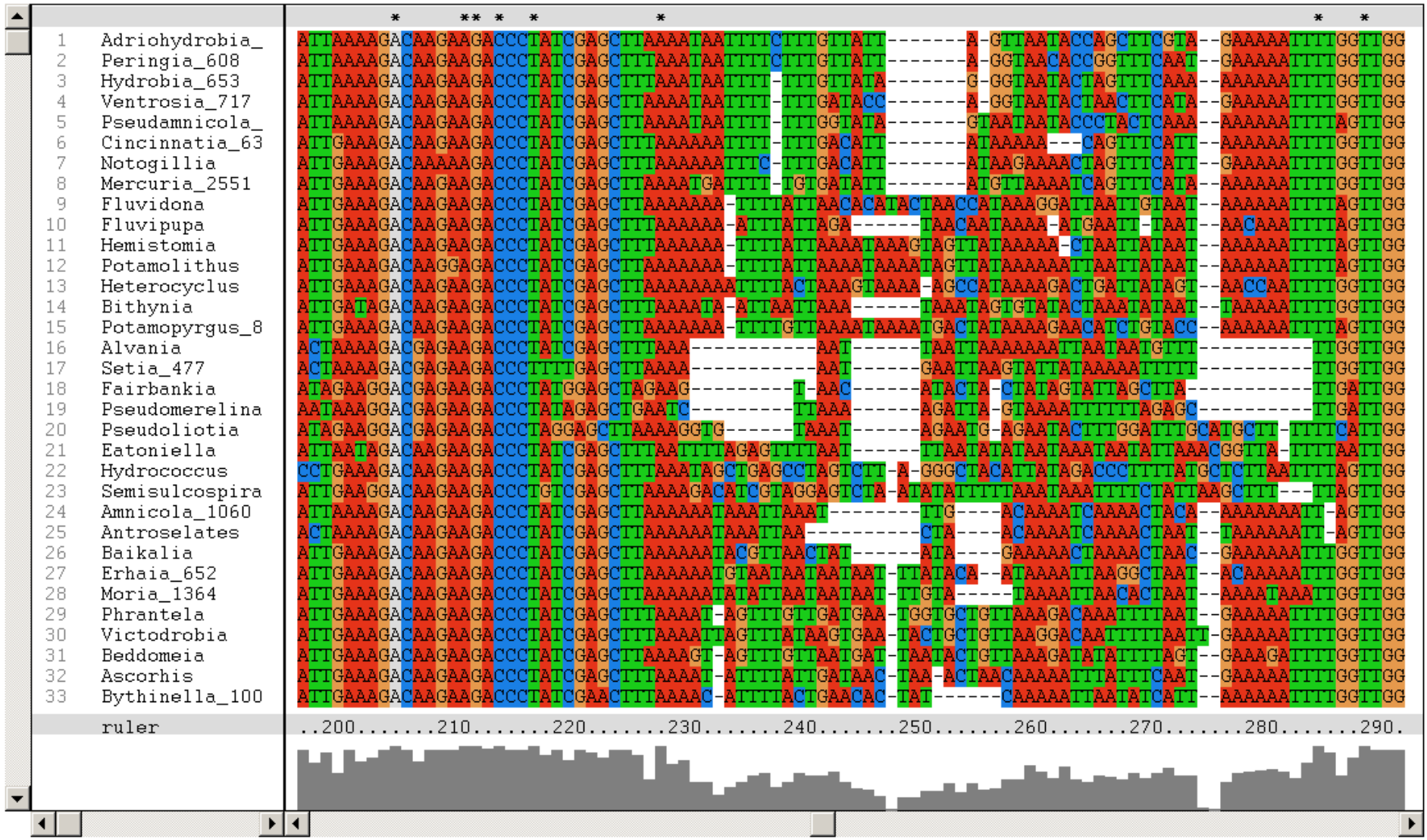
} Gap information, output tree type

Enter or Paste a set of Sequences in any supported format :

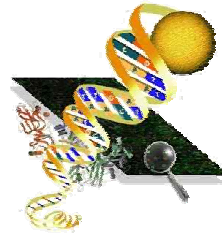
} File input in GCG, FASTA, EMBL, GenBank, Phylip, or several other formats

Multiple Alignment Mode

Font Size: 10



# MSA with PILEUP



- ◆ PILEUP is the MSA program that is part of the Genetics Computer Group (GCG) sequence analysis package
- ◆ Sequences are aligned pairwise using dynamic programming algorithm
- ◆ The scores are used to produce a phylogenetic tree, which is then used to guide the alignment of the most closely related sequences and groups of sequences
- ◆ Resulting alignment is a global alignment produced by the Needleman-Wunsch algorithm

MSA definition

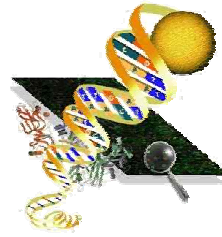
The need for MSA

The MSA problem

**MSA methods**

Sequence editors

# Iterative MSA methods



- ◆ Attempt to correct initial alignment problems by repeatedly aligning subgroups of the sequences and then by aligning these subgroups into a global alignment of all the sequences
- ◆ MultAlin – recalculates pair-wise scores during the production of the progressive alignment and uses these scores to recalculate the tree
- ◆ PRRP – initial alignment is made to predict a tree, the tree is used to produce weights where the sequences are analyzed for the presence of aligned regions that include gaps
- ◆ SAGA – based on genetic algorithm that is a machine-learning algorithm that attempts to produce alignments by the simulations of evolutionary changes in sequences

MSA definition

The need for MSA

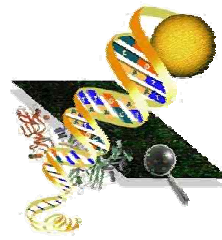
The MSA problem

**MSA methods**

Sequence editors



# Editing and formatting alignments



## ◆ Sequence editors are used for:

- manual alignment/editing of sequences
- visualization of data
- data management
- import/export of data
- graphical enhancement of data for presentations

## ◆ Examples:

- **CINEMA** (Color Interactive Editor for Multiple Alignments) web applet  
<http://www.biochem.ucl.ac.uk/bsm/dbbrowser/CINEMA2.02/kit.html>
- **GDE** (Genetic Data Environment) - UNIX based  
[http://bimas.dcrtnih.gov/gde\\_sw.html](http://bimas.dcrtnih.gov/gde_sw.html)
- **GeneDoc** - MS Windows <http://www.psc.edu/biomed/genedoc/>
- **MACAW** - local multiple sequence alignment program and sequence editing tool available by anonymous FTP from [ncbi.nih.gov/pub/schuler/macaw](http://ncbi.nih.gov/pub/schuler/macaw)
- **BioEdit** - sequence alignment editor for MS Windows with web access and accessory applications (BLAST, local BLAST, ClustalW, Phylip and more)

MSA definition

The need for MSA

The MSA problem

MSA methods

Sequence editors

BioEdit Sequence Alignment Editor - [C:\Thommy\Academy\Hydrobiidae project\1BS\_aligned.bio]

File Edit Sequence Alignment View World Wide Web Accessory Application RNA Options Window Help

Add / Modify / Remove an Accessory Application

ClustalW Multiple alignment  
BLAST

CAP contig assembly program  
DNADist ---> Neighbor phylogenetic tree  
DNADist DNA distance matrix  
DNAMik DNA Maximum Likelihood program with molecular clock  
FastDNAMl DNA maximum likelihood  
Fitch -- Fitch-Margolash and Least-Squares Distance Methods  
IdPlot identity plotter  
Kitsch -- Fitch-Margolash and Least Squares Methods with Evolutionary Clock  
NEIGHBOR -- Neighbor-Joining and UPGMA methods  
Protdist ---> Fitch phylogenetic tree  
Protdist protein distance matrix  
Protpars protein parsimony method

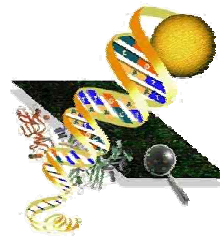
Courier New 11 B 84 t

Mode: Select / Slide Selection:0 Position: 19

10 20 80 90 100

Adrio hydrobi	CCAUGCAUGUCUAAGUUCACACUCU	GCACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Adrio insulan	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Alvania	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Alzoniella 2	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Amnicola 106	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Amphithalamu	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Antroselates	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Ascorhis	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Assimineea	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Assimineea 16	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Assimineea 22	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Avenionia 22	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Baicalia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Barleesia	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Beddomeia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Belgrandia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Bithynia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Bythinella 1	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Bythiospeum	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Calopia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Cecina 2522	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Clenchiella	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Coxiella	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Eatoniella	CCAUGCAUGUACAAGUUCACACCCU	GCACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	UGCGAAUCUUGGA
Emmericia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Emmericia 30	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Erhaia 652	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Fairbankia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Fissuria 243	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Fluvidona	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Fluvipupa	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Fontigens	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCGA	---AGCUACUUGGA
Gammatricula	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Geomelania 8	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Geomelania 8	CCAUGCAUGUCUAAGUUCACACCCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Graziana 256	CCAUGCAUGUACAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Hauffenia 25	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Heleobops	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Hemistomia	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Heterocyclus	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Horatia 2598	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Hydrobia 653	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Hydrococcus	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AUCUACUUGGA
Islamia 2327	CCAUGCAUGUCUAAGUUCACACUCU	GUACGG	UGAAACCGCGAAUGGCUCAUUAUUCAGUCGAGG	UUCCUUAGAUGAUCCAA	---AGCUACUUGGA

# Summary MSA



## ◆ Definition:

A multiple sequence alignment is an alignment of  $n > 2$  sequences obtained by inserting gaps (" - ") into sequences such that the resulting sequences have all length  $L$  and can be arranged in a matrix of  $N$  rows and  $L$  columns where each column represents a homologous position

## ◆ Why do we need MSA?

- Formulate & test hypotheses about protein 3-D structure
- MSA can help us to reveal biological facts about proteins
- Crucial for genome sequencing
- To establish homology for phylogenetic analyses
- Identify primers and probes to search for homologous sequences in other organisms

## ◆ The MSA problem

- Most pairwise alignment algorithms are too complex to be used for  $n$ -wise alignments
- Alignment algorithms need to be optimized
  - \* use structural information
  - \* use phylogenetic information
  - \* use conserved regions

## ◆ MSA methods

- Progressive global alignment (starts with the most alike sequences)
  - \* e.g., ClustalW, ClustalX, Pileup
- Iterative methods (initial alignment of groups of sequences that are revised)
  - \* MultAlin, PRRP, SAGA
- Alignments based on locally conserved patterns

## ◆ Sequence editors

- CINEMA GDE, GeneDoc, MACAW, BioEdit

MSA definition

The need for MSA

The MSA problem

MSA methods

Sequence editors