

Multiple Sequence Alignment (I)

Jianlin Cheng, PhD

School of Electrical Engineering and Computer Science
University of Central Florida



2006

Free for academic use. Copyright @ Jianlin Cheng & original sources for some materials

What's Multiple Sequence Alignment?

- A generalization of pairwise sequence alignment
- Align multiple sequences (e.g., a protein family) to maximize alignment score
- Biological goal is to identify evolutionary relationships, structures, and functions

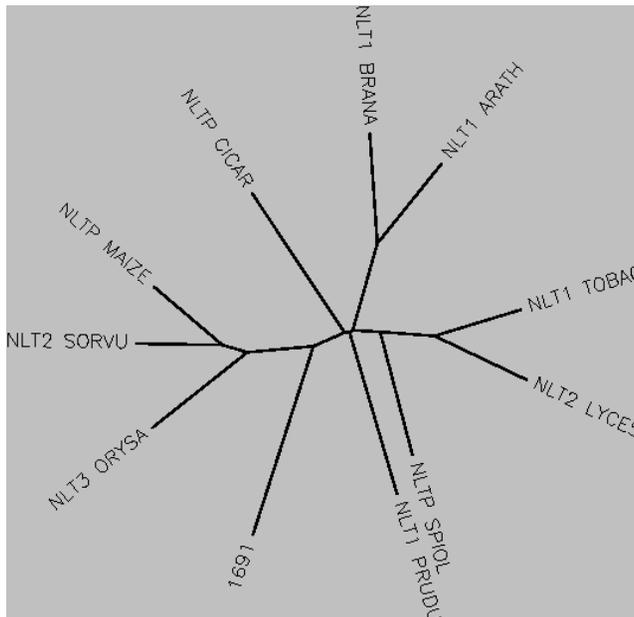
Why is MSA important?

- Construct phylogenetic tree
- Construct protein profile for database search and structure prediction
- Elucidate evolutionary conservation
- Identify DNA binding /regulatory sites
- Identify structural and functional important sites of proteins
- Improve alignment quality

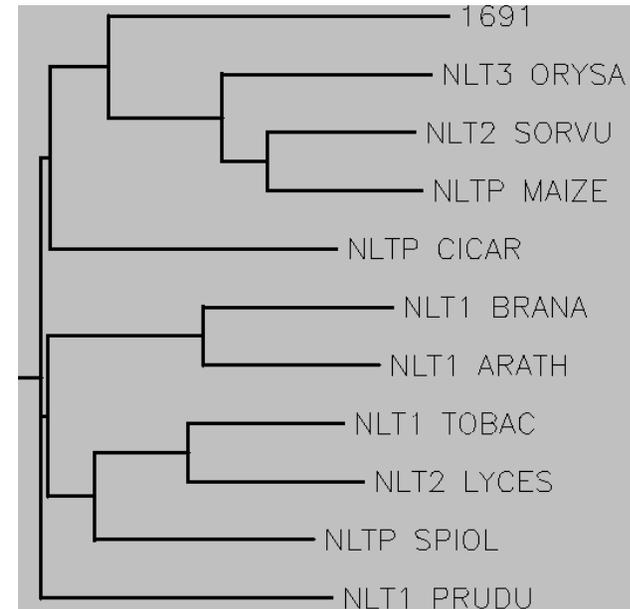
Example 1: phylogenetic tree

MSA

1691	1	----	MARAGHNKY	VARVMV	VALLLA	AARAPVT	CGQVVS	TWAPCI	MYADGE						
NLT1_BRANA	1	----	MAGLVKLS	SCLVLACM	IVAGPI	ATNAALS	CGTVSGN	LAACIG	YLTQ-						
NLT1_ARATH	1	----	MAGVMKLA	CLLLACM	IVAGPI	TSNAALS	CGSVNS	NLAACI	GYVLQ-						
NLTP_CICAR	1	----	MASKVVVC	VALIMC	IVIAPMA	ES-AITC	GRVDTA	LAPCLG	YLQG-						
NLT1_TOBAC	1	----	MEIAGKI	ACFVVL	CMVVAAP	CAEA--	ITCGQV	TSNLAP	CLAYLRN-						
NLTP_SPIOL	1	--	MASSAVI	KLACAV	LLCI	VVAAPY	AEAG-	ITCGM	VSSKLAP	CIGYLVKQ-					
NLT2_LYCES	1	----	MEMVSKI	ACFVLL	CMVVVAP	HAEA--	ITCGQV	TAGLAP	CLPYLQG-						
NLT1_PRUDU	1	--	MAYSAMT	KLALV	VALCMVV	VP	IAQA--	ITCGQV	SSNLAP	CIPYVRG-					
NLT3_ORYSA	1	----	MARAQLV	LVALV	AAALLL	AGPHT	TMAAIS	CGQVNS	AVSPCL	SYARG-					
NLT2_SORVU	1	MARSMK	LAVAI	AVVAAAA	VVLAAT	TSEAAV	T	CGQVSS	AI	GPCLSYARGQ					
NLTP_MAIZE	1	--	MARTQQ	LAVVAT	AVVAL	VLLAAAT	TSEAAIS	CGQVAS	AI	APCISYARGQ					
consensus			K	AC	V	CM	vv	AAP	A	AAit	Q	S	LaP	i	G



Phylogenetic unrooted tree



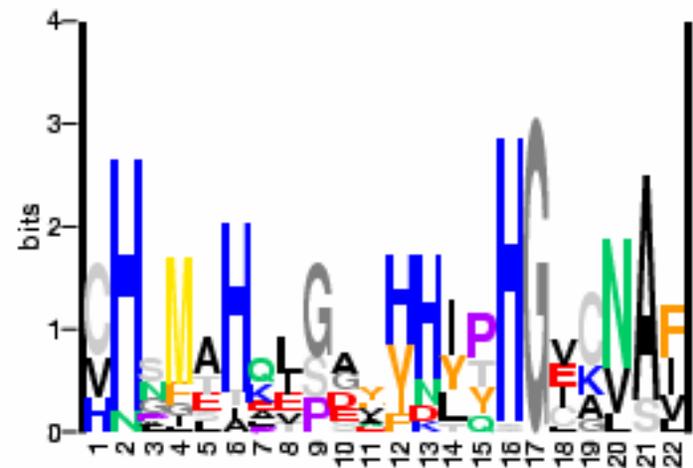
Phylogenetic rooted tree

Example 2: Functional Important Residues / Site

```

FUCO_ECOLI  VHGM1AH2PLGAFYNTPHGVANAI
GLDA_BACST  HNGFTALEGEIHHLTHGEK3VA4F
GLDA_ECOLI  VHNGLTAI5PDAH6HYYHGEK7VA8F
MEDH_BACMT  VHSISHQVGGVYKLQHGICNSV
ADH1_CLOAB  CHSMAHKTGAVFHI9PHGCANAI
ADHE_ECOLI  CHSMAHKLGSQFHI10PHGLANAL
ADH2_ZYMMO  VHAMA11HQLGGY12NLPHGVCNAV
ADH4_YEAST  VHALA13HQLGGFY14HLPHGVCNAV
ADHA_CLOAB  CHPMEHEL15SAYYDITHGVGLAI
ADHB_CLOAB  VHLMEHEL16SAYYDITHGVGLAI
  
```

iron containing alcohol dehydrogenases



PSSM of BL00913C (ADH_IRON_1); 11 sequences.

conserved histidines probably bind the ferrous ion(s) required for these enzymes activity

Source: http://bioinformatics.weizmann.ac.il/~pietro/Making_and_using_protein_MA/

Scoring Function of MSA

- Amino acid level: PAM and BLOSUM
- Score of a column: sum of the scores of all residue pairs in the column.
- Score of an MSA: sum of the scores of all columns.

DP for Multiple Sequence Alignment

To align prefixes:
 $P[1..i]$, $Q[1..j]$, $R[1..k]$

Consider $2^3 - 1 = 7$ situations

(1) Given $P[1..i-1]$, $Q[1..j-1]$, $R[1..k-1]$, match $P[i]$, $Q[j]$, and $R[k]$

i
YNT**PHG**VANAI
 j
HNG**F**TALEGEI
 k
VHNG**LTA**IPDA

i
YNT**P** - **H****G**VANAI
 j
HN - **G****F** - **T**ALEGEI
 k
VHNG**LTA**IPDA

(2) Given $P[1..i], Q[1..j-1], R[1..k-1]$,
 match gap, $Q[j]$, and $R[k]$

YNTPHG	-	VANAI
HN - GF	-	T ALEGEI
VHNGLT	-	A IPDA

(3) Given $P[1..i], Q[1..j], R[1..k-1]$,
 match gap, gap, and $R[k]$

YNTPHG	-	VANAI
HN - GF	-	T ALEGEI
VHNGLT	-	A IPDA

(4) Given $P[1..i], Q[1..j-1], R[1..k]$,
 match gap, $Q[j]$, and gap

YNTPHG	-	-VANAI
HN - GF	-	T ALEGEI
VHNGLT	-	-IPDA

(5) Given $P[1..i], Q[1..j-1], R[1..k-1]$,
 match gap, $Q[j]$, and $R[k]$

YNTPHG	-	VANAI
HN - GF	-	T ALEGEI
VHNGLT	-	A IPDA

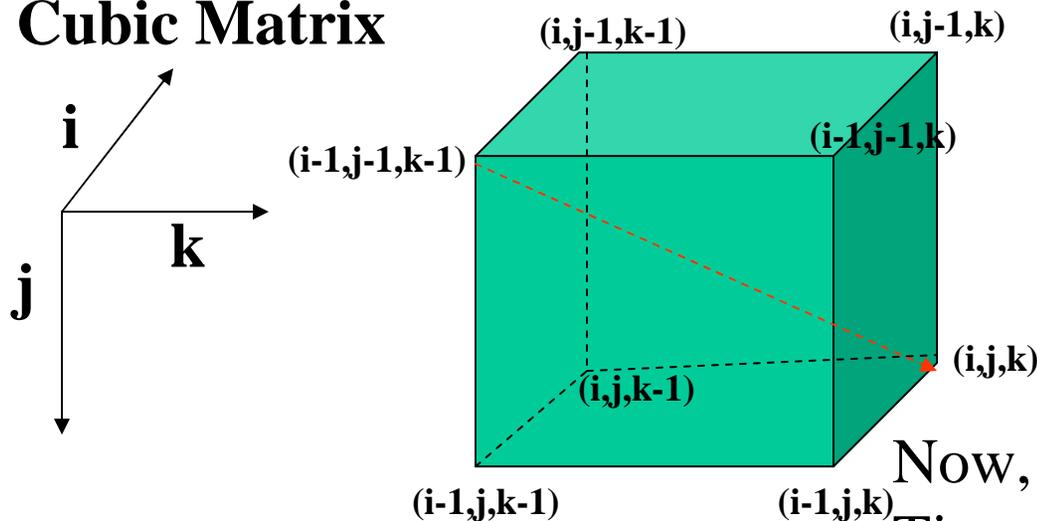
(6) Given $P[1..i-1], Q[1..j], R[1..k-1]$,
match $P[i]$, gap, and $R[k]$

		i	
Y	N	T	P - H G V A N A I
		j	
H	N - G	F	T - A L E G E I
		k	
V	H	N	G L T A I P D A

(7) Given $P[1..i-1], Q[1..j-1], R[1..k]$,
match $P[i]$, $Q[j]$, and gap

		i	
Y	N	T	P - H - G V A N A I
		j	
H	N - G	F - -	T - A L E G E I
		k	
V	H	N	G L T A - I P D A

Cubic Matrix



Now, fill a 3D matrix.
Time = n^3 . Each
Step has $2^3 - 1$ paths.

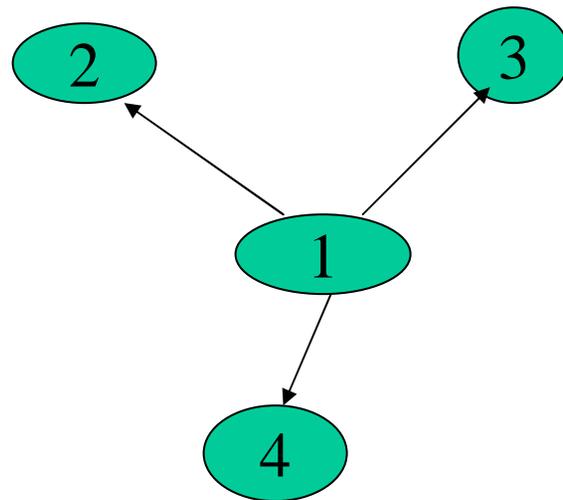
Now, generalized to m sequences
Time Complexity = $n^m * 2^m * m^2$

Heuristic MSA Approaches

- DP can't be used for more than three sequences
- Star approach
- Progressive approach
- Iterative approach
- Hidden Markov Model

Star Approach

Use one sequence of your interest or closest to all other sequences as seed. Align the seed to other sequences and generate a MSA from pairwise alignments.



(Used by PSI-BLAST)

Select a seed

1. **VHGMAHPLAFYNT**HGVANAI
2. HNGFTEGEIHHHLTHGEKVF
3. VHNGLAIPDAH**YHGEK**VAF
4. VSHQVGGVYKLQHGICNSV

Pairwise

1. VHGMAHPLAFYNT-HGVANAI
2. HNGFT-E-GEIHHHLTHGEKVF

1. VHGMAHPL-AFYNT**H**GVANAI
3. VHNGL-AIPDA-H**Y**HGEKVAF

1. VHGMAHPL-AFYNT-HGVANAI
4. VHS--HQVGGVYKLQHGICNSV

Do combination step by step using seed as an anchor

1. VHGMAHPLAFYNT-HGVANAI
2. HNGFT-E-GEIHHHLTHGEKVF

1. VHGMAHPL-AFYNT-HGVANAI
2. HNGFT-E--GEIHHHLTHGEKVF
3. VHNGL-AIPDA-HY-HGEKVAF

1. VHGMAHPL-AFYNT**H**GVANAI
3. VHNGL-AIPDA-H**Y**HGEKVAF

1. VHGMAHPL-AFYNT-HGVANAI
4. VHS--HQVGGVYKLQHGICNSV

1. VHGMAHPL-AFYNT-HGVANAI
2. HNGFT-E--GEIHHHLTHGEKVF
3. VHNGL-AIPDA-HY-HGEKVAF
4. VHS--HQVGGVYKLQHGICNSV

Progressive Approach

- Basic idea: align similar sequences first
 - (1) construct a coarse phylogenetic tree
 - (2) use the tree as guide tree to align sequences progressively.
- Representative methods: CLUSTALW and T-COFFEE

Simplified Example of Progressive MSA

1. VHGMALPLAFYNTHGVAIAI
2. HNGFTEGEIHHLTHGEKVF
3. VHNGLAIPDAHGHGEKVAF
4. VHSHQVGGVYKLGHGICNSV

Align sequences pairwise
to construct a similarity matrix

	1	2	3	4
1	10	5	6	8
2	5	9	7	4
3	6	7	9	4
4	8	4	4	10

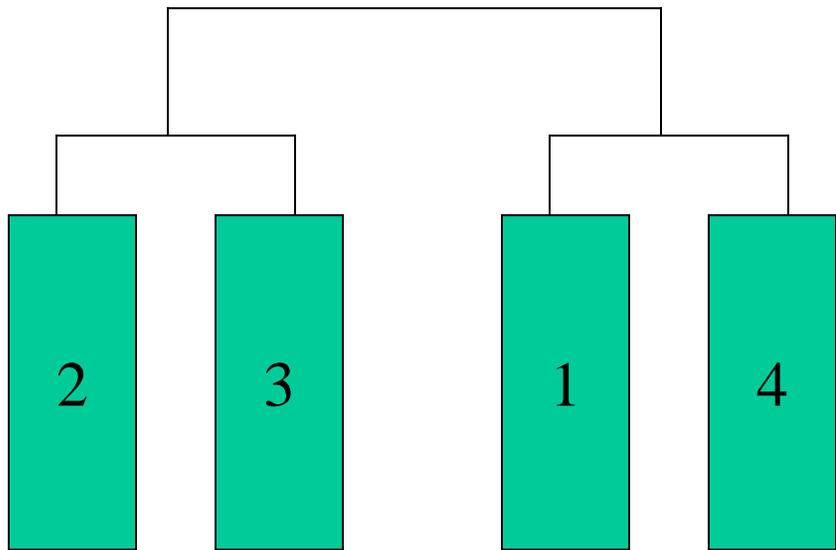
Convert similarity matrix to distance matrix if necessary

$$(S_{aa} + S_{bb})/2 - S_{ab}$$

	1	2	3	4
1	0	4.5	3.5	2
2	4.5	0	2	5.5
3	3.5	2	0	5.5
4	2	5.5	5.5	0

Step 1

Step 2



Construct a guide tree from distance/similarity matrix

Step 3

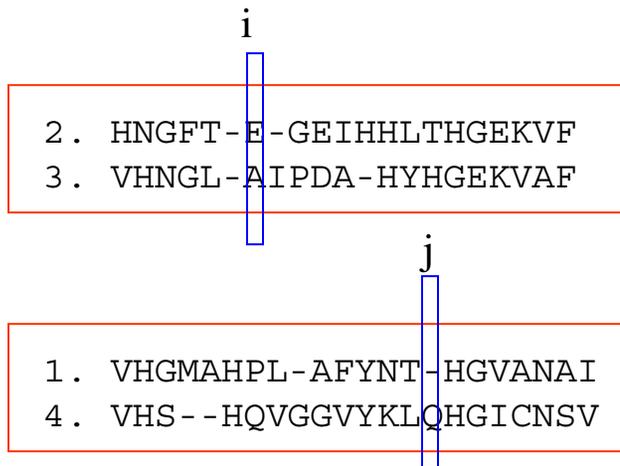
Step 4

Align sequence progressively according to guide tree.

- Align 2 and 3
- Align 1 and 4
- Align group of 2-3 and group of 1-4

Question:

How to align a group of sequences against another group of sequences?



$$\text{Score}(i,j) = S(E,-)+S(E,Q)+S(A,-)+S(A,Q)$$

Align one group with another group:

- Treat each group as a single generalized sequence consisting of columns
- Use the same dynamic programming algorithm
- Once add one gap into a position in one group, all sequences in the group is added the gap
- $\text{Score}(\text{Col}_i, \text{Col}_j)$ is the sum score of all residue pairs across column i and column j)

Some Progressive MSA Tools

- **CLUSTALW(X)** (neighbor-joining algorithm, similarity matrix)
- **PILEUP** (UPGMA algorithm for guide tree)
- **MULTALIGN**(UPGMA algorithm)
- **T-COFFEE** (improved over CLUSTALW, but slower)
- **3D-COFFEE** (use structure information)

Iterative Approach

- Progressive approach is a heuristic approach. Optimal MSA is not guaranteed.
- Alignments in previous steps are fixed. Error propagation.
- Iterative approach tries to adjust alignments according to an objective function to improve alignments.

Iterative Methods

- PRRP: optimize a progressive alignment by iteratively dividing the sequences into two groups and realigning them using a group-to-group alignment algorithm
- SAGA: use a genetic algorithm
- Iterative approach often improves alignment accuracy at the expense of computation time

Hidden Markov Model Approach

- MUSCLE (<http://www.drive5.com/muscle/>, Robert Edgar)
- State of Art Multiple Sequence Alignment
- Online service:
http://phylogenomics.berkeley.edu/cgi-bin/muscle/input_muscle.py
- Not discussed here, but in Advanced Bioinformatics (Spring, 2007)

A Case Study of Using CLUSTALW

HIV-1 Capsid Protein

VHQAISPRTLNAWVKVVEEKAFSPEVIPMFSALSEGATPQDLNTMLNTVGGHQA
AMQMLKETINEEAAEWDRVHPVHAGPIAPGQMREPRGSDIAGTTSTLQEQIGW
MTNNPIPVGEIYKRWIILGLNKIVRMYSPTSILDIRQGPKEPFRDYVDRFYKTLR
AEQASQEVKNWMTETLLVQNANPDCKTILKALGPAATLEEMMTACQ

Search Database

POL-HV1BR

VHQAI SPRTLNAWVKVVEEKAFSPEVIPMFSALSEGATPQDLNTMLNTVGGHQAAMQMLK
ETINEEAAEWDRVHPVHAGPIAPGQMREPRGSDIAGTTSTLQEQIGWMTNPPPIPVGEIY
KRWILGLNKIVRMYSPTSILDIRQGPKEPFRDYVDRFYKTLRAEQASQEVKNWMTETLL
VQANPDCKTILKALGPAATLEEMMTACQ

AAN23474

VHQAI SPRTLNAWVKVIEEKAFSPEVIPMFSALSEGATPQDLNTMLNTVGGHQAAMQMLK
ETINEEAAEWDRIHPVHAGPIAPGQMREPRGSDIAGTTSTLQEQIGWMTSNPPPIPVGEIY
KRWILGLNKIVRMYS PVSILDIRQGPKEPFRDYVDRFYKTLRAEQASQDVKNWMTETLL
VQANPDCKTILKALGPAATLEEMMTACQ

AAN73808

VHQAI SPRTLNAWVKVIEEKAFSPEVIPMFSALSEGATPQDLNTMLNTVGGHQAAMQMLK
ETINEEAAEWDRLHPVHAGPVAPGQMRDPRGSDIAGTTSTLQEQIGWMTSNPPPIPVGEIY
KRWILGLNKIVRMYS PVSILDIRQGPKEPFRDYVDRFYKTLRAEQATQEVKNWMTETLL
IQANPDCKTILKALGPAATLEEMMTACQ

AAD28894

VHQAL SPRTLNAWVKVVEEKAFSPEVIPMFSALSEGATPQDLNTMLNTVGGHQAAMQMLK
ETINEEAAEWDRLHPVHAGPIAPGQMREPRGSDIAGTTSTLQDQIGWMTNPPPIPVGEIY
KRWILGLNKIVRMYSPTSILDIKQGPKEPFRDYVDRFYKTLRAEQASQEVKNWMTETLL
VQANPDCKTILKALGPAATLEEMMTACQ

ClustalW

- Web service:

<http://www.ebi.ac.uk/clustalw/>

ClustalW Submission Form

ClustalW 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. [New users, please read the FAQ.](#)

>> [Download Software](#)   

YOUR EMAIL	ALIGNMENT TITLE	RESULTS	ALIGNMENT	CPU MODE
<input type="text" value="@cs.ucf.edu"/>	<input type="text" value="Sequence"/>	<input type="text" value="interactive"/>	<input type="text" value="full"/>	<input type="text" value="single"/>
KTUP (WORD SIZE)	WINDOW LENGTH	SCORE TYPE	TOPDIAG	PAIRGAP
<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="percent"/>	<input type="text" value="def"/>	<input type="text" value="def"/>
MATRIX	GAP OPEN	END GAPS	GAP EXTENSION	GAP DISTANCES
<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="def"/>	<input type="text" value="def"/>

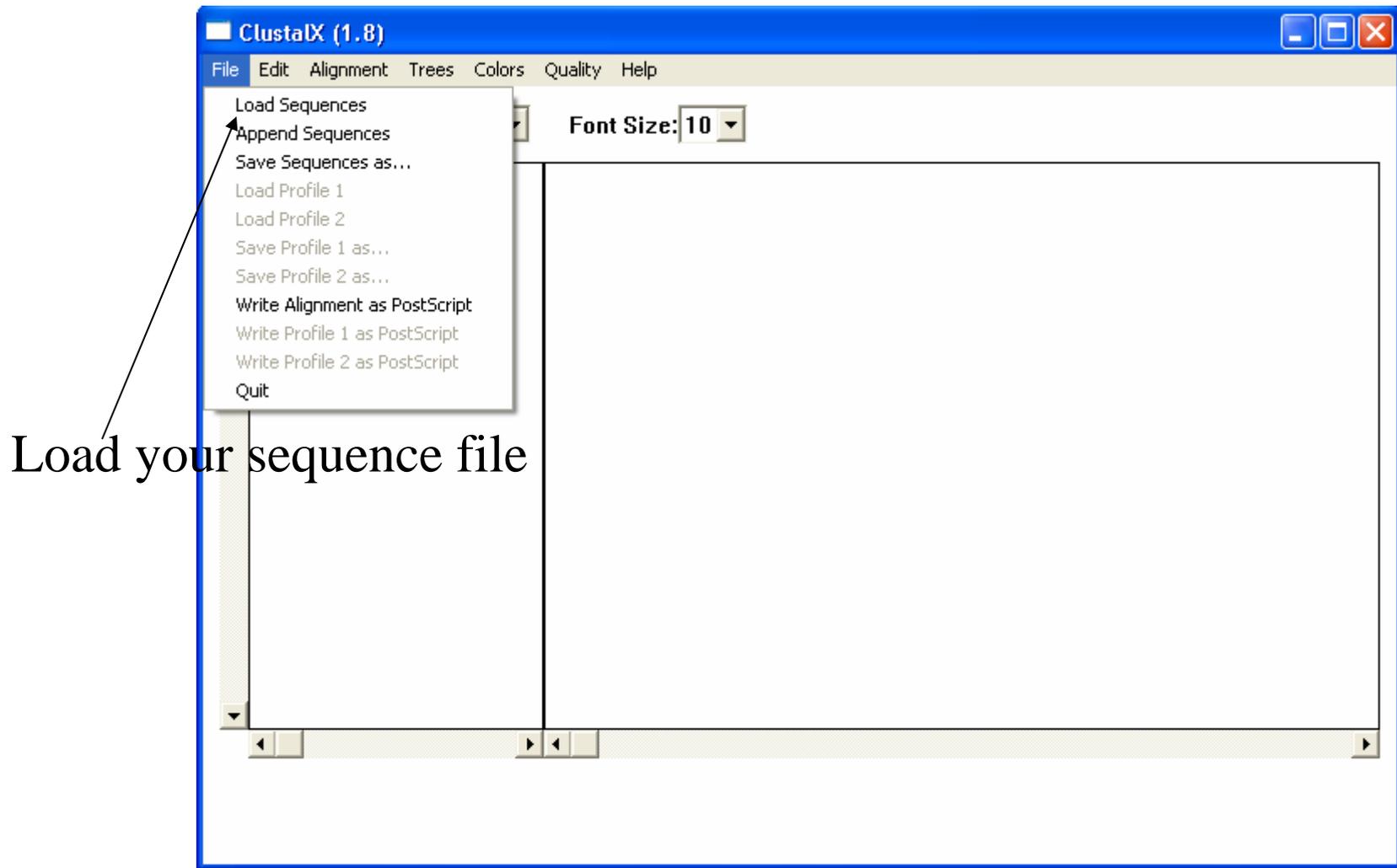
OUTPUT		PHYLOGENETIC TREE		
OUTPUT FORMAT	OUTPUT ORDER	TREE TYPE	CORRECT DIST.	IGNORE GAPS
<input type="text" value="aln w/numbers"/>	<input type="text" value="aligned"/>	<input type="text" value="nj"/>	<input type="text" value="on"/>	<input type="text" value="off"/>

Enter or Paste a set of Sequences in any supported format: [Help](#)

```
DIAGTTSTLQEQIGWMTSNPPIPVGEIYKRWII LGLNKIVRMYS PV
SILD IRQGPKEPFRDYVDRFYKTLRAEQATQEVKNWMTETLLIQNA
NPDCKTILKALGPAATLEEMMTACQ
>AAD28894
VHQALSPRTLNAWVKVVEEKAFSP EVIPMFSALSEGATPODLN TML
NTVGGHQAA MQMLKETINEEAAEWDR LHPVHAGPIAPGOMREPRGS
DIAGTTSTLQDQIGWMTNPPIPVGEIYKRWII LGLNKIVRMYSPT
SILD IKQGPKEPFRDYVDRFYKTLRAEQASQEVKNWMTETLLVQNA
NPDCKTILKALGPAATLEEMMTACQ
```

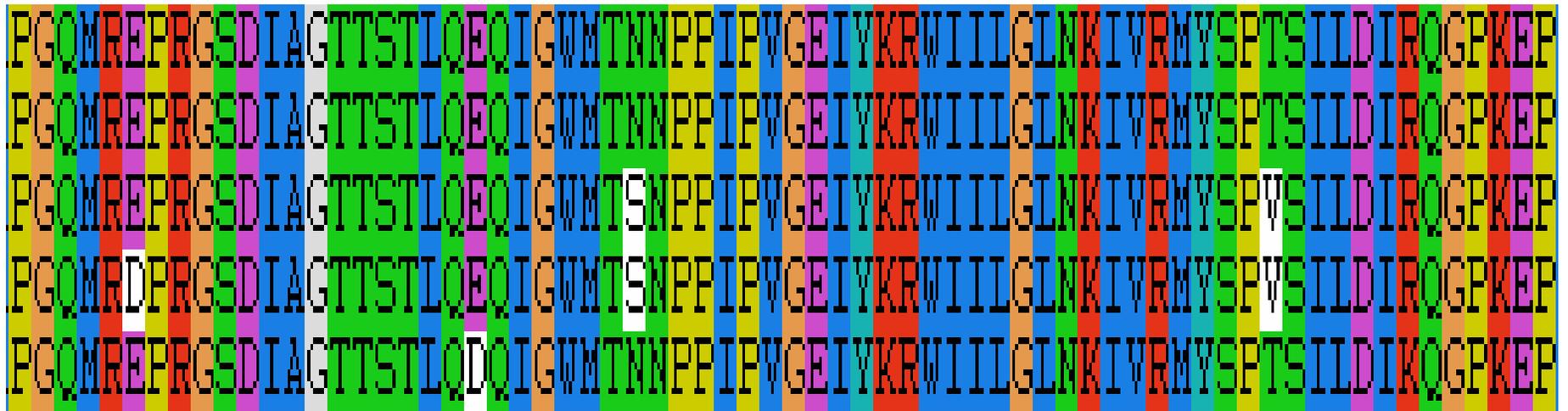
Upload a file:

ClustalX with GUI

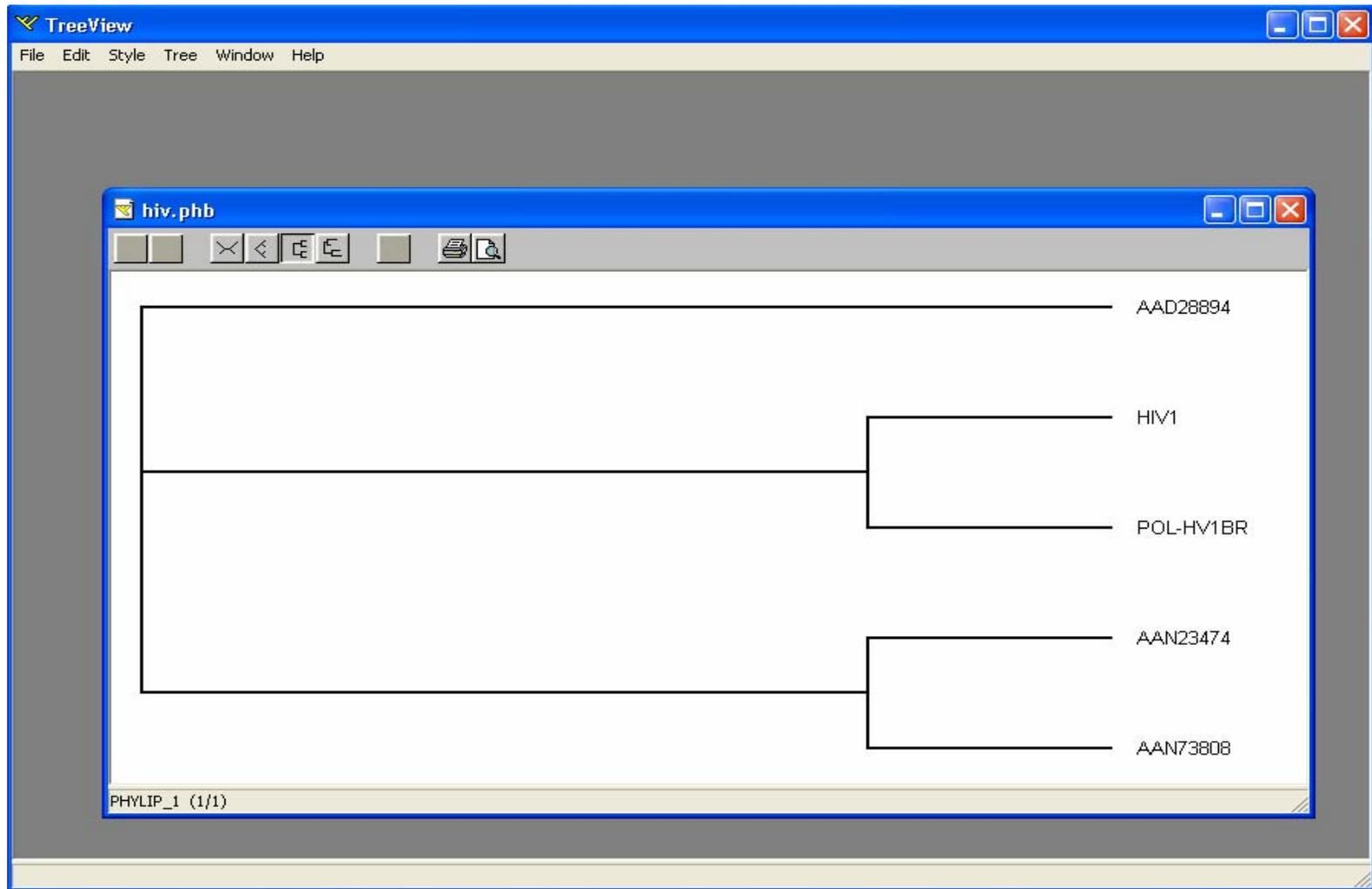


Load your sequence file

<http://bioinformatics.unc.edu/software/opensource/index.htm>



Visualized MSA (conservation and variation)



Distance Based Phylogenetic Tree (visualized by treeview:

<http://taxonomy.zoology.gla.ac.uk/rod/treeview.html>)

T-Coffee

 ch.EMBLnet.org

Home

Services

Courses

Links

Contacts

T-COFFEE

This program is more accurate than ClustalW for sequences with less than 30% identity,
but it is slower...

For any question please contact the author: [Cedric Notredame](#)

Please cite:

"T-Coffee: A novel method for multiple sequence alignments"
C. Notredame, D. Higgins, J. Heringa
Journal of Molecular Biology, **302**, 205-217, (2000)

Valid format for input is: FASTA(Pearson)

max number of sequences = 30

max total length of sequences = 10000

[Help page](#)

Output order:

Sequence type: Protein DNA

Input
sequences:
(see above
for valid
formats)

Run T-COFFEE

Clear Input

Address: <http://www.ch.embnet.org/software/TCoffee.html>