



Clustal 1988-2024

Des Higgins

UCD

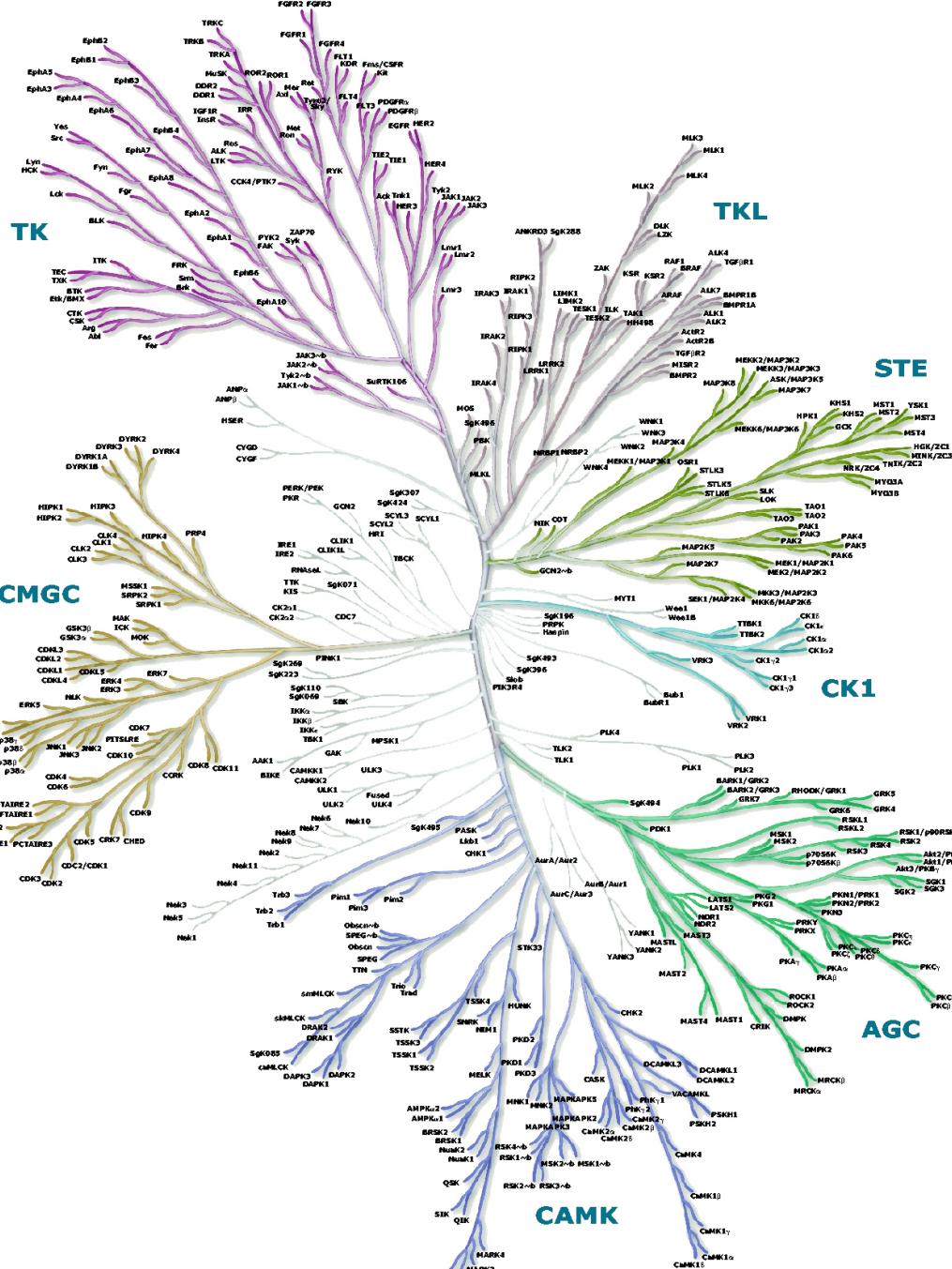
Multiple Alignment?

- Align 3 or more sequences together

TPIS_THEPX	TAKDANEVIKAIRNTIASLYGKEKADLVRIQYCGSVKPENISELIAESDIDGALVGGASL
TPIS_GEOTN	TAEDANNVCGHIRSVARLFGAEEAAEAIRIQQYCGSVKPENISFLAQEHIDGALVGGASL
TPIS_SYN P2	-ATEANR VIGLIR---EKLT NKN---VTI QYCGSVKPNNVDEI IAQPEIDGALVGGASL
	* :***.* ** * . : : *****:*****:*****: .*****

Homologous residues lined up in columns

Needed because of
Orthologues from different species
But mainly:
Paralogues from Gene duplications
Multi-gene families
e.g. humans have approx. 500 protein kinases



Human Protein Kinases

The human kinome comprises 40 atypical PKs and 478 classical PKs. The latter consist of 388 serine/threonine kinases, 90 tyrosine kinases and 50 sequences which lack a functional catalytic site.

(Manning et al., Science, 2002)

Globin Multiple Alignment

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

-----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYPWTQRFFESFGDLST
-----VQLSGEEKAAVLALWDKVN--EEEVGGEALGRLLVVYPWTQRFFDSFGDLSN
-----VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPFH-DLS-
-----VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFILGFPTTKTYFPFH-DLS-
-----VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKT
PIVDTGSVAPLSAAEKTIRSAWAPVYSTYETSGVDILVKFFTSTPAAQEFFFPKFKGLTT
-----GALTESQAALVKSSWEFNANIPKHTHRFFILVLEIAPAAKDLFSFLKGTE

* : : : * . : . : * : * : * : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

PDAVMGNPKVKAHGKKVLGAFSDGLAHLDN-----LKGTFATLSELHCDKLHVD PENFRL
PGAVMGNPKVKAHGKKVLHSFGEGVHLDN-----LKGTFAALSELHCDKLHVD PENFRL
----HGSAQVKGHGKKVADALTNAAHVDD-----MPNALSALSSDLHAHKLRDPVNFKL
----HGSAQVKAHGKKVGDALTLAVGHLD-----LPGALSNLSDLHAHKLRDPVNFKL
EAEMKASEDLKKHGVTVL TALGAILKKGH-----HEAEELKPLAQSHATKH KIPIKYLEF
ADQLKKSADVRWHAERIINAVNDAVASMDT--EKMSMKLRLDLSGKHAKS FQVDPQYFKV
VP--QNNPELQAHAGKVF KL VYEAAIQLQVTGVVVTDATLKNLGSVHVS KGVA D-AHF PV

. . : * . : . : * . * . : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

LGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH-----
LGNVLVVVLARHFGKDFTPELQASYQKVVAGVANALAHKYH-----
LSHCLLVTLAAHLPAEFTP AVHASLDKFLASVSTVLT SKYR-----
LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLT SKYR-----
ISEAI IHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
LAAVIADTV AAG---D-----AGFEKLM SMICILLRSAY-----
VKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA---

Globin Multiple Alignment

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

-----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYPWTQRFFESFGDLST
-----VQLSGEEKAAVLALWDKVN--EEEVGGEALGRLLVVYPWTQRFFDSFGDLSN
-----VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPFH-DLS-
-----VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFILGFPTTKTYFPFH-DLS-
-----VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKT
PIVDTGSVAPLSAAEKTIRSAWAPVYSTYETSGVDILVKFFTSTPAAQEFFFPKFKGLTT
-----GALTESQAALVKSSWEFNANIPKHTHRFFILVLEIAPAAKDLFSFLKGTE

* : : : * . : . : * : * : * : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

PDAVMGNPKVKAHGKKVLGAFSDGLAHLDN-----LKGTFATLSELHCDKLHVD PENFRL
PGAVMGNPKVKAHGKKVLHSFGEGVHLDN-----LKGTFAALSELHCDKLHVD PENFRL
----HGSAQVKGHGKKVADALTNAAHVDD-----MPNALSALSSDLHAHKLRDPVNFKL
----HGSAQVKAHGKKVGDALTIAVGHLD-----LPGALSNLSDLHAHKLRDPVNFKL
EAEMKASEDLKKHGVTVL TALGAILKKGH-----HEAEELKPLAQSHATKH KIPIKYLEF
ADQLKKSADVRWHAERI I NAVNDAVASMDT--EKMSM KLRDLSGKHA KSFQVDPQYFKV
VP--QNNPELQAHAGKVF KL VYEAAIQLQVTGVVVT DATLKNLGSVHVS KGVA D-AHF PV

. . : * . : . : * . * . : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

LGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH-----
LGNVLVVVLARHFGKDFTPELQASYQKVVAGVANALAHKYH-----
LSHCLLVTLAAHLPAEFTP AVHASLDKFLASVSTVLT SKYR-----
LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLT SKYR-----
ISEAI IHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
LAAVIADTV AAG---D-----AGFEKLM SMICILLRSAY-----
VKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA---

Globin Multiple Alignment

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

-----VHLT**PEEKSAVTALWGKV**N--**VDEVGGEALGRLLVVYPWTQRFFESFGDL**ST
-----VQLS**GEEKAAVLALWDKV**N--**EEEVGGEALGRLLVVYPWTQRFFDSFGDLSN**
-----VLS**PADKTNVKAAWGKVGAH**AGEYGAEALERMFLSFPTTKTYFPFH-DLS-
-----VLS**AADKTNVKAAWSKVGGH**AGEYGAEALERMFILGFPTTKTYFPFH-DLS-
-----VLS**EGEWQLVLHWAKVEAD**VAGHGQDILIRLFKSHPETLEKFDRFKHLKT
PIVDTGSVAPLS**AAEKT**KIRSAWAPVYSTYETSGVDIILVKFFTSTPAAQE~~FFPKFKGLTT~~
-----GALT**ESQAALVKSSWEENANIPKHTHRFF**FILVLEIAPA~~AKD~~LFSFLKGTE

* : : : * . : . : * : * : * : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

PDAVMGN**PKVKAHGKKVLGAFSDGLAHLDN**--**LKGTFATLSELHCDKLHVD**PENFRL
PGAVMGN**PKVKAHGKKVLHSFGEGVHHLDN**--**LKGTF**AALSELHCDKLHVD**PENFRL**
----HGS**AQVKGHGKKVADALTNAV**AHVDD----**MPNALSALS**DLHAHKLRVD**PVNFKL**
----HGS**AQVKAHGKKVGDALT**LAVGHLDD----**LPGALSNLSDLH**AHKLRVD**PVNFKL**
EAEMKAS**EDLKKHGVTVL**TALGAILKKGH----**HEAELKPLAQSHAT**KHKIP**I**KYLEF
ADQLKKS**ADVRWHAERI**I NAVNDAVASMDDT--EKM**SMKLRDLSGKHAKSFQVD**PQYFKV
VP--QNN**PELQAHAGKVFKLVYEAAI**QLQVTGVVVT**DATLKNLGSVHVS**KGVAD-AHFV

. . : * . : . : * . * . : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

LGNVLVCVLAHH**FGKEFTPPVQA**AYQKVVAGVANALA**HKYH**-----
LGNVLVVVLARH**FGKDFTP**ELQA**SYQKVVAGVANALA****HKYH**-----
LSHCLLVTLAAHLP**AEFTP**AVHASLDKFLASVSTVLT**SKYR**-----
LSHCLLSTLAVHL**PNDFTPAVHA**SLDKFLSSVSTVLT**SKYR**-----
ISEAIIHVLHSR**HPGDFGADAQG**AMNKALELFRKDIA**AKYKELGYQG**
LAAVIADTVAAAG---D-----AGFEKLM~~S~~MICILLR**SAY**-----
VKEAILKTIKEVVGAKWSEELNS**AWTIAYDELAIVIK**KEMNDAA---

Alpha helices

Globin Multiple Alignment

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

-----VHLT**PEEKSAVTALWGKVN**--**VDEVGGEALGRLLVVYPWTQRFFESFGDLST**
-----VQLS**GEEKAAVLALWDKVN**--**EEEVGGEALGRLLVVYPWTQRFFDSFGDLSN**
-----VLS**PADKTNVKAAWGKVGAH****AGEYGAEALERMFLSFPTTKTYFPHF-DLS-**
-----VLS**AADKTNVKAAWSKVGGH****AGEYGAEALERMFILGFPTTKTYFPHF-DLS-**
-----VLS**EGEWQLVLHVWAKVEAD****VAGHGQDILIRLFKSHPETLEKFDRFKHLKT**
PIVDTGSVAPLS**AAEKT**KIRSAWAPVYST**YETSGVDIILVKFFTSTPAAQE**FFPKFKGLTT
-----GALT**ESQAALVKSSWEENANIPKHTHRFFILVLEIAPAAKD**LFSFLKGTE

* : : : * . : . : * : * : * : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

PDAVMGN**PKVKAH**GKKVLGAFSDGLAHLDN-----L**KGT**FATLSELHCDKLHVD**PENFRL**
PGAVMGN**PKVKAH**GKKVLHSFGEGVHHLDN-----L**KGT**FAALSELHCDKLHVD**PENFRL**
----HGS**AQVKGH**GKKVADALTNAVAVDD-----MPNALSAISDL**HAHKLRVD**PVNFKL
----HGS**AQVKAH**GKKVGDALTIAVGHLDD-----LPGALSNLSDL**HAHKLRVD**PVNFKL
EAEMKAS**EDLKKH**GVTVLTALGAILKKGH-----HEAELKPLAQSHATKHKIP**IKYLEF**
ADQLKKS**ADVRWHAERI**I NAVNDAVASMDTT--EKM**SMKLRDLSGKHAKSFQVD**PQYFKV
VP--QNN**PELQAHAGKVFKLVYEAAI**QLQVTGVVVT**DATLKNLGSVHVSKGVAD-AHF**PV

. . : * . : . : * . * . : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

LGNVLVCVLAHH**FGKEFTPPVQA**AYQKVVAGVANALA**HKYH**-----
LGNVLVVVLARH**FGKDFTPELQA**SYQKVVAGVANALA**HKYH**-----
LSHCLLVTLAAH**LPAEFTP**AVHASLDKFLASVSTVLT**SKYR**-----
LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLT**SKYR**-----
ISEAIIHVLHSR**HPGDFGADAQG**AMNKALELFRKDIA**AKYKELGYQG**
LAAVIADTVAG---D-----AGFEKLMMSMICILLRSAY-----
VKEAILKTIKEV**VGAKWSEELNS**AWTIAYDELAIVIK**KEMNDAA**---

Haem binding Histidines

Globin Multiple Alignment

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

-----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYPWTQRFFESFGDLST
-----VQLSGEEKAAVLALWDKVN--EEEVGGEALGRLLVVYPWTQRFFDSFGDLSN
-----VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPFH-DLS-
-----VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFILGFPTTKTYFPFH-DLS-
-----VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKT
PIVDTGSVAPLSAAEKTIRSAWAPVYSTYETSGVDILVKFFTSTPAAQEFFFPKFKGLTT
-----GALTESQAALVKSSWEFNANIPKHTHRFFILVLEIAPAAKDLFSFLKGTE

* : : : * . : . : * : * : * : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

PDAVMGNPKVKAHGKKVLGAFSDGLAHLDN-----LKGTFATLSELHCDKLHVD PENFRL
PGAVMGNPKVKAHGKKVLHSFGEGVHLDN-----LKGTFAALSELHCDKLHVD PENFRL
----HGSAQVKGHGKKVADALTNAAHVDD-----MPNALSALSSDLHAHKLRDPVNFKL
----HGSAQVKAHGKKVGDALTIAVGHLD-----LPGALSNLSDLHAHKLRDPVNFKL
EAEMKASEDLKKHGVTVL TALGAILKKGH-----HEAEELKPLAQSHATKH KIPIKYLEF
ADQLKKSADVRWHAERI I NAVNDAVASMDT--EKMSM KLRDLSGKHA KSFQVDPQYFKV
VP--QNNPELQAHAGKVF KL VYEAAIQLQVTGVVVT DATLKNLGSVHVS KGVA D-AHF PV

. . : * . : . : * . * . : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

LGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH-----
LGNVLVVVLARHFGKDFTPELQASYQKVVAGVANALAHKYH-----
LSHCLLVTLAAHLPAEFTP AVHASLDKFLASVSTVLT SKYR-----
LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLT SKYR-----
ISEAI IHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
LAAVIADTV AAG---D-----AGFEKLM SMICILLRSAY-----
VKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDA---

The first MSA?

- Thanks to Lior Pachter
<https://liorpachter.wordpress.com/2013/12/08/molecular-restoration-studies-of-extinct-forms-of-life/>
- L. Pauling and E. Zuckerkandl (1963)
Chemical Paleogenetics: Molecular ‘restoration studies’ of extinct forms of life. Acta Chem. Scand. 17, S9–S16.

```
alpha      -VLSPADKTNVKAAGKVGAGAEGEYGAELERMFLSFPTTKTYFPHFDL-----SHGSA
gamma     GHFTEEDKATITSLWGKVNV--EDAGGETLGRLLLVVYPWTQRFFDSFGNLSSASAIMGNP
beta      VHLTPEEKSAVTALWGKVNV--DEVGGEALGRLLLVVYPWTQRFFESFGDLSTPDAVMGNP
delta     VHLTPEEKTAVNALWGKVNV--DAVGGEALGRLLLVVYPWTQRFFESFGDLSSPDAVMGNP
                   : : *: :.. **** .    *.*: *::: :* *: :*   *       *.

alpha      QVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRDPVNFKLLSHCLLVTLAAHL
gamma     KVKAHGKKVLTSLGDAKHLDDLKGTFAQSELHCDKLHVDOPENFKLLGNVLVTVLAIHF
beta      KVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDOPENFRLLGNVLVCVLAHHF
delta     KVKAHGKKVLGAFSDGLAHLDNLKGTFSQLSELHCDKLHVDOPENFRLLGNVLVCVLAERNF
                   :**.***** : : .. *:***: : : ***:***:***:***:***:***:***: .: .** ..

alpha      PAEFTPAAVHASLDKFLASVSTVLTSKYR
gamma     GKEFTPEVQASWQKMVTAVASALSSRYH
beta      GKEFTPQVQAAYQKVVAGVANALAHKYH
delta     GKEFTPQMQAAAYQKVVAGVANALAHKYH
                   ***** :*: : *.: .*: *: .*: :
```



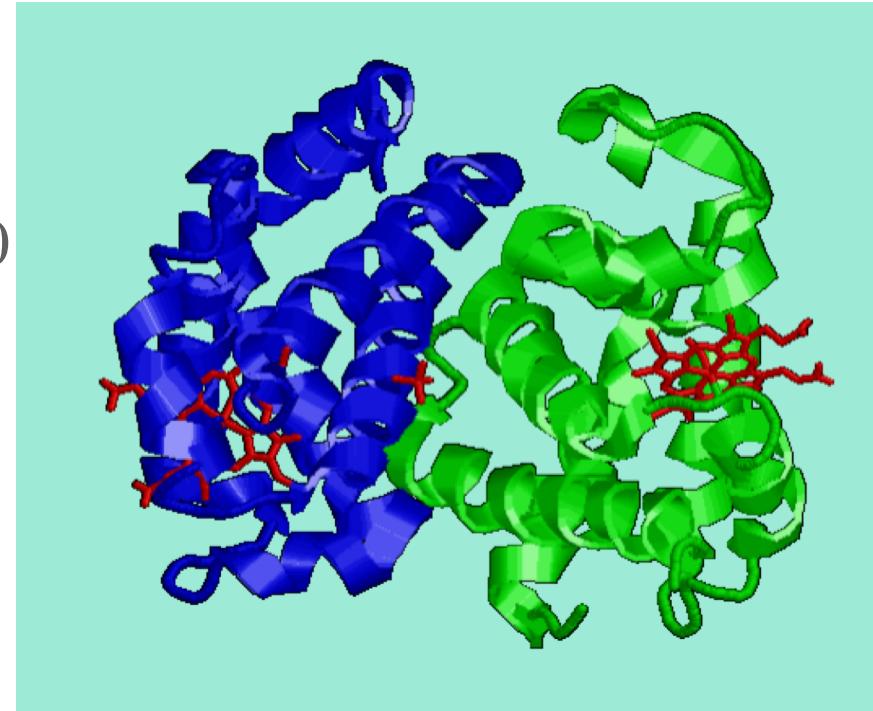
Pauling and Zuckerkandl, Japan, 1986

VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP
-VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLS-----HGSA
* *

KVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHF
QVKHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVPVNFKLLSHCLLVTLA AHL
** *

GKEFTPQPVQAAYQKV VAGVANALAHKYH
PAEFTPQAVHASLDKFLASVSTVLT SKYR
***** * * * * * * * * *

- Dynamic Programming
 - Needleman and Wunsch, 1970
 - $O(L^2)$ algorithm
- Maximise score
(or minimise distance)
 - Gap penalties
 - Amino acid weight matrix



	A	C	D	E	F	G	H	I	K	L	M	N	P	O	R	S	T	V	W	Y
A	4																			
C	0	9																		
D	-2	-3	6																	
E	-1	-4	2	5																
F	-2	-2	-3	-3	6															
G	0	-3	-1	-2	-3	6														
H	-2	-3	-1	0	-1	-2	8													
I	-1	-1	-3	-3	0	-4	-3	4												
K	-1	-3	-1	1	-3	-2	-1	-3	5											
L	-1	-1	-4	-3	0	-4	-3	2	-2	4										
M	-1	-1	-3	-2	0	-3	-2	1	-1	2	5									
N	-2	-3	1	0	-3	0	1	-3	0	-3	-2	6								
P	-1	-3	-1	-1	-4	-2	-2	-3	-1	-3	-2	-2	7							
Q	-1	-3	0	2	-3	-2	0	-3	1	-2	0	0	-1	5						
R	-1	-3	-2	0	-3	-2	0	-3	2	-2	-1	0	-2	1	5					
S	1	-1	0	0	-2	0	-1	-2	0	-2	-1	1	-1	0	-1	4				
T	0	-1	-1	-1	-2	-2	-2	-1	-1	-1	-1	0	-1	-1	-1	1	5			
V	0	-1	-3	-2	-1	-3	-3	3	-2	1	1	-3	-2	-2	-3	-2	0	4		
W	-3	-2	-4	-3	1	-2	-2	-3	-3	-2	-1	-4	-4	-2	-3	-3	-2	-3	11	
Y	-2	-2	-3	-2	3	-3	2	-1	-2	-1	-1	-2	-3	-1	-2	-2	-2	-1	2 7	

Blosum 62 Matrix

- +ve: likely substitutions -ve: unlikely subs.
 - from counting matches in blocks of alignment
 - Henikoff and Henikoff, 1992

	A	T	T	A	C	C
G	0	$0.1 \rightarrow 0.3 \rightarrow 0.5 \rightarrow 0.6 \rightarrow 0.7 \rightarrow 0.8$				
	\downarrow	\searrow	\downarrow	\downarrow	\searrow	\downarrow
	0.1	$0.2 \rightarrow 0.4$		$0.6 \rightarrow 0.7 \rightarrow 0.8 \rightarrow 0.9$		
A	\downarrow	\searrow		\searrow		
	0.3	0.1	$0.3 \rightarrow 0.5 \rightarrow 0.6 \rightarrow 0.7 \rightarrow 0.8$			
T	\downarrow	\downarrow	\searrow	\searrow		
	0.4	0.2	0.1	$0.3 \rightarrow 0.4 \rightarrow 0.5 \rightarrow 0.6$		
	\downarrow	\downarrow	\downarrow	\downarrow	\searrow	\downarrow
G	0.5	0.3	0.2	$0.4 \rightarrow 0.5 \rightarrow 0.6 \rightarrow 0.7$		
	\downarrow	\downarrow	\searrow	\downarrow	\searrow	\downarrow
	0.6	0.4	0.3	0.2	$0.3 \rightarrow 0.4 \rightarrow 0.5$	
T	\downarrow	\searrow		\downarrow	\searrow	
	0.8	0.6	0.5	0.4	0.2	$0.3 \rightarrow 0.4$
A	\downarrow	\searrow	\downarrow	\downarrow	\searrow	\searrow
	1	0.8	0.7	0.6	0.4	0.4
A	\downarrow	\downarrow	\downarrow	\searrow	\downarrow	\searrow
	1.2	1	0.9	0.8	0.6	0.4
C	\downarrow	\searrow	\downarrow	\downarrow	\searrow	\downarrow
	1.4	1.2	1.1	1	0.8	0.6
A	\downarrow	\searrow	\downarrow	\searrow	\downarrow	\searrow
						0.6

Human beta	-----VHLT PEEKSAVTALWGKV N-- VDEVGGEALGRLLVVYPWTQR FFESFGDLST
Horse beta	-----VQLS GEEKAAVLALWDKV N-- EEEVGGEALGRLLVVYPWTQR FFDSFGDLSN
Human alpha	-----VLS PADKTNVKAAWGKV GAH AGEYGAEALERMFLSFPTTKT YFPHF-DLS-
Horse alpha	-----VLS AADKTNVKAAWSKV GGH AGEYGAEALERMFLGFPTTKT YFPHF-DLS-
Whale myoglobin	-----VLS EGEWQLVLHVWAKV EAD VAGHGQDILIRLFKSHPETILE KFDRFKHLKT
Lamprey globin	PIVDTGSVAPLS AAEKTKIRSAWAPV YST YETSGVDILVKFFTTPAAQE FFPKFKGLTT
Lupin globin	-----GALT ESQAALVKSSWEF NAN IPKH THRFFILVLEIAPAAKD LFSFLKG TSE
	* : : : * . : : * : * : .
Human beta	PDAVMGN PKVKAH GKKVLGAFSDGLAHLDN----L KGT FATLSELHCDKLHVD PENFRL
Horse beta	PGAVMGN PKVKAH GKKVLHSFGEGVHHLDN----L KGT FAALSELHCDKLHVD PENFRL
Human alpha	---HGS AQVKGH GKKVADALTNAVAHVDD----M PNALSALS SDLHAKLDRVDPVNFKL
Horse alpha	---HGS AQVKAH GKKVGDAITLAVGHLD----L PGALSNLSDL HAKLDRVDPVNFKL
Whale myoglobin	EAEMKASE EDLKKH GVTVLTALGAILKKGH----H EAELKPLAQSHAT KHKIP IKYLEF
Lamprey globin	ADQLKKS ADVRWHAERI I NAVNDAVASMDDT--EKM SMKLRDLSGK HAKSFQVD PQYFKV
Lupin globin	VP--QNN PELOAH AGKVFKLVYEAAIQLQVTGVVVTDATLKNLGSVHVS KG VAD-AHFPV
	. . : * . : . : * . * . : .
Human beta	LGNVLVCVLAHH FGKEFTP PPVQA AYQKVVAGVANALA HKYH-----
Horse beta	LGNVLVVVLAH FGKDFTE LP ELQA SYQKVVAGVANALA HKYH-----
Human alpha	LSHCLLVTLAHL PAEFTPAVHA SLDKFLASV STVLT SKYR -----
Horse alpha	LSHCLLSTLAVH LPNDFTPAVHA SLDKFLSSV STVLT SKYR -----
Whale myoglobin	I SEAI I HVLHSRHPGDFGADAQG AMNKALELFRKDIA AKYKELGYQG
Lamprey globin	LAAVIADTV AAAG---D-----A GFEKIMSMICILLR SAY-----
Lupin globin	VKEAILKTIKEV VGAKWSEELNS AWTIAYDEL AIIVIK KEMNDAA ---
	: : : .

$$\sum_{i=2}^N \sum_{j=1}^{i-1} W_{ij} D_{ij}$$

Time O(L^N)

Human beta	-----VHLT PEEKSAVTALWGKV N-- VDEVGGEALGRLLVVYPWTQR FFESFGDLST
Horse beta	-----VQLS GEEKAAVLALWDKV N-- EEEVGGEALGRLLVVYPWTQR FFDSFGDLSN
Human alpha	-----VLS PADKTNVKAAWGKV GAH AGEYGAEALERMFLSFPTTKT YFPHF-DLS-
Horse alpha	-----VLS AADKTNVKAAWSKV GGH AGEYGAEALERMFLGFPTTKT YFPHF-DLS-
Whale myoglobin	-----VLS EGEWQLVLHVWAKV EAD VAGHGQDILIRLFKSHPETILE KFDRFKHLKT
Lamprey globin	PIVDTGSVAPLS AAEKTKIRSAWAPV YST YETSGVDILVKFFTTPAAQE FFPKFKGLTT
Lupin globin	-----GALT ESQAALVKSSWEF NAN IPKHTHRFFILVLEIAPAAKD LFSFLKGTE

* : : : * . : : * : * : .

Human beta	PDAVMGN PKVKAH GKKVLGAFSDGLAHLDN----L KGT FATLSELHCDKLHVD PENFRL
Horse beta	PGAVMGN PKVKAH GKKVLHSFGEGVHHLDN----L KGT FAALSELHCDKLHVD PENFRL
Human alpha	---HGS AQVKGH GKKVADALTNAVAHVDD----M PNALSALS SDLHAKLDRVDPVNFKL
Horse alpha	---HGS AQVKAH GKKVGDAALTNAVGHLD----L PGALSNLSDL HAKLDRVDPVNFKL
Whale myoglobin	EAEMKASE EDLKKH GVTVLTALGAILKKGH----H EAELKPLAQSHAT KHKIPIKYLEF
Lamprey globin	ADQLKKS ADVRWHAERI I NAVNDAVASMDDT--EKM SMKLRDLSGK HAKSFQVD PQYFKV
Lupin globin	VP--QNN PELOAH AGKVFKLVYEAAIQLQVTGVVVT DATLKNLGSVHVS KGVAD-AHFPV

. . : * . : . : * . * . : .

Human beta	LGNVLVCVLAHH FGKEFTP PVAQ AYQKVVAGVANALA HKYH----
Horse beta	LGNVLVVVLAH HFGKDFPELQA SYQKVVAGVANALA HKYH----
Human alpha	LSHCLLVTLAAH LPAEFTP AVHA SLDKFLASVSTVLT SKYR----
Horse alpha	LSHCLLSTLAVH LPNDFTP AVHA SLDKFLSSVSTVLT SKYR----
Whale myoglobin	ISEAI I HVLHSRHPGDFGADAQG AMNKALELFRKDIA AKYKELGYQG
Lamprey globin	LAAVIADTVAAAG --D----A GFEKIMSMICILLR SAY----
Lupin globin	VKEAILKTIKEV VGAKWSEELNS AWTIAYDELAIVIK KEMNDAA--

: : . : : :

Sequences Time

- 2 1 second
- 3 150 seconds
- 4 6.25 hours
- 5 39 days
- 6 16 years
- 7 2404 years

Time $O(L^N)$

The first automatic MSA?

- David Sankoff, 1973!



Evolution of 5S RNA and the Non-randomness of Base Replacement

SEQUENCES of 5S rRNA have been published for five widely divergent organisms: *Escherichia coli*, *Pseudomonas fluorescens*, yeast, human KB carcinoma, and *Xenopus laevis*¹⁻⁵. The question arises whether sequences for the ancestors of these organisms, represented by X, Y, and Z in Fig. 1, can be reconstructed to any degree of confidence; were this possible, statistical analysis of mutation types would become feasible. Such reconstructions are well known in protein studies⁶ and for tRNAs⁷, but they are somewhat more difficult for 5S rRNA.

- Reconstruct ancestral seqs between pairs
 - Dynamic programming
- Align in steps

NATURE NEV

DAVID SANKOFF
CRISTIANE MOREL
ROBERT J. CEDERGREN

Centre de Recherches Mathématiques and
Département de Biochimie,
Université de Montréal,
Case Postale 6128, Montréal 101

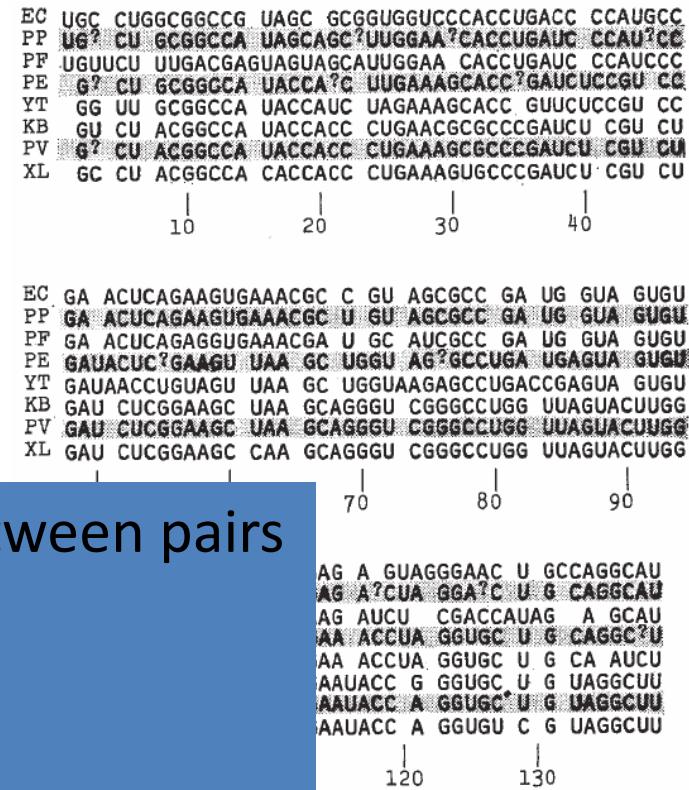


Fig. 2 Alignment of known and reconstructed 5S RNA sequences. From top to bottom: EC, *E. coli*; PP, proto-prokaryote; PF, *P. fluorescens*; PE, proto-eukaryote; YT, yeast; KB, human KB carcinoma; PV, proto-vertebrate; XL, *Xenopus*. Reconstructed sequences are shaded; ?, undetermined base or gap.

Paulien Hogeweg

- Utrecht, The Netherlands
- “Bioinformatics”
 - Hesper B., Hogeweg P. (1970.)
"Bioinformatica: een werkconcept",
Kameleon, 1(6): 28–29.
- Automated MSA 1984
 - Hogeweg, P. and B. Hesper B. (**1984**)
The alignment of sets of sequences and the construction of phyletic trees: an integrated method. *J Mol Evol* 20: 175-186.
 - Align sequences together gradually



Human beta	-----VHLT PEEKSAVTALWGKV N-- VDEVGGEALGRLLVVYPWTQR FFESFGDLST
Horse beta	-----VQLS GEEKAAVLALWDKV N-- EEEVGGEALGRLLVVYPWTQR FFDSFGDLSN
Human alpha	-----VLS PADKTNVKAAWGKV GAH AGEYGAEALERMFLSFPTTKT YFPFH-DLS-
Horse alpha	-----VLS AADKTNVKAAWSKV GGH AGEYGAEALERMFLGFPTTKT YFPFH-DLS-
Whale myoglobin	-----VLS EGEWQLVLHVWAKV EAD VAGHGQDILIRLFKSHPETILE KFDRFKHLKT
Lamprey globin	PIVDTGSVAPLS AAEKTIRSAWAPV YST YETSGVDILVKFFTSTPAAQE FFPKFKGLTT
Lupin globin	-----GALT ESQAALVKSSWEF NAN IPKH THRFFILVLEIAPAAKD LFSFLKG TSE

* : : : * . : : * : * : .

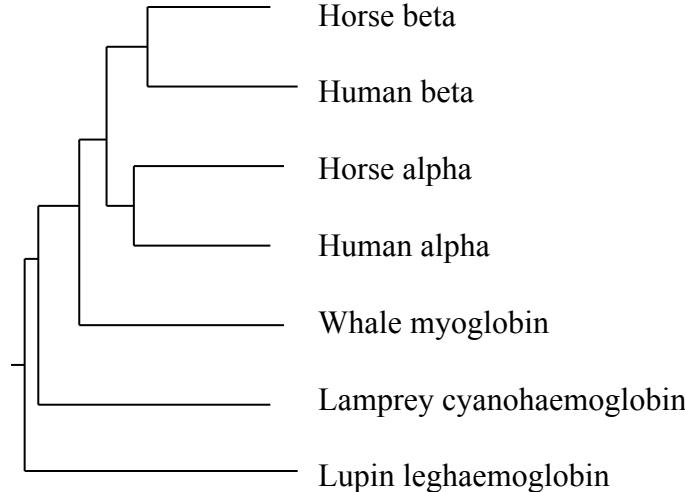
Human beta	PDAVMGN PKVKAH GKKVLGAFSDGLAHLDN----L KGT FATLSELHCDKLHVD PENFRL
Horse beta	PGAVMGN PKVKAH GKKVLHSFGEGVHHLDN----L KGT FAALSELHCDKLHVD PENFRL
Human alpha	---HGS AQVKGH GKKVADALTNAVAHVD----M PNALSALS SDLHAKLDRVDPVNFKL
Horse alpha	---HGS AQVKAH GKKVGDACTLAVGHLD----L PGALSNLSDL HAKLDRVDPVNFKL
Whale myoglobin	EAEMKASE EDLKKHG TVTLTALGAILKKGH----H EAELKPLAQSHAT KHKIP IKYLEF
Lamprey globin	ADQLKKS ADVRWHAERI INAVNDAVASMDDT--EKM SMKLRDLSGKHAKSFQVD PQYFKV
Lupin globin	VP--QNN PELOAHAGKVFKLVYEAA IQLQVTGVVVT DATLKNLGSVHVS KGVAD-AHFPV

. . : * : . : : * . * . : .

Human beta	LGNVLVCVLAHH FGKEFTP PPVQA AYQKVVAGVANALA HKYH----
Horse beta	LGNVLVVVLAH FGKDFTE LPQAS SYQKVVAGVANALA HKYH----
Human alpha	LSHCLLVTLAH LPAEFTPAVHA SLDKFLASVSTVLT SKYR----
Horse alpha	LSHCLLSTLAH LPNDFTPAVHA SLDKFLSSVSTVLT SKYR----
Whale myoglobin	I SEAI I HVLHSR HPGDFGADAQG AMNKALELFRKDIA AKYKELGYQG
Lamprey globin	LAAVIADTVAAAG --D----A GFEKIMSMICILLR SAY----
Lupin globin	VKEAILKTIKEV VGAKWSEELNS AWTIAYDELAIVIK KEMNDAA--

: : . : : :

“Guide Tree”



Progressive Alignment:

- Barton and Sternberg, 1987
- Florence Corpet, 1988
- Feng and Doolittle, 1987
- Jotun Hein, 1989
- Higgins and Sharp, 1988, 1989
- Hogeweg and Hesper, 1984
- Willie Taylor, 1987, 1988

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

-----VHLT**PEEKSAVTALWGKVN**--VDEVG**GEALGRLLVVY**PWT**QRFFESFGDLST**
-----VQLS**GEEKAAVLALWDKVN**--EEE**EVGGEALGRLLVVY**PWT**QRFFDSFGDLSN**
-----VLS**PADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF**-DLS-
-----VLS**AADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHF**-DLS-
-----VLS**EWEQQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKT**
PIVDTGSVAPLS**AAEKT**KIRSAWAPVYSTYETSGVDILVKFTSTPAAQE**FFPKFKGLTT**
-----GALT**ESQAALVKSSWEEFNANIPKH**THRFFILVLEIA**PAAKD**LFSFLKG**TSE**

* : : : * . : : : * : * : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

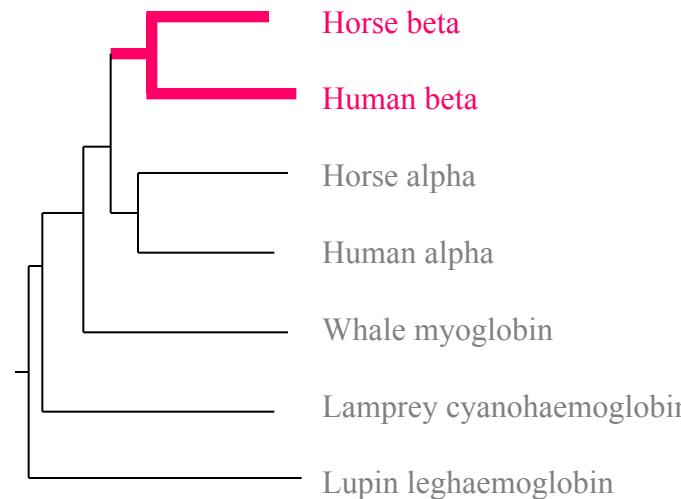
PDAVMGNPKV**KAHGKKVLGAFSDGLAHL**DN----LKGT**FATLSEL**HCDKLHVD**PENFRL**
PGAVMGNPKV**KAHGKKVLHSFGEGVHH**LDN----LKGT**FAALSEL**HCDKLHVD**PENFRL**
---HGS**AQVKGHGKKVADALTNAVAHVDD**----MPNALSALS**SDLH**AHKLRVDPVNFKL
---HGS**AQVKAHGKKVGDA**LTAV**GHLDD**----LPGALS**NLSLH**AHKLRVDPVNFKL
EAEMKASEDLKK**HGTVTL**TALGAILKKGH----HEAE**ELKPLAQSHATKH**KIPIKYLEF
ADQLKKSADVRW**HAERI**INAVNDAVASMDT--EKMS**MKLRDLSGKHAKSFQVDPQYFKV**
VP--QNNPELO**QAHAGKVFKLVYEAA**IQLQVTGVVVT**DATLKNLGSVHVSKGVAD**-AHFPV

. : : * : . : : : * . * : . : .

Human beta
Horse beta
Human alpha
Horse alpha
Whale myoglobin
Lamprey globin
Lupin globin

LGNVLVCVLAHHFGKEFTPPPVQA**AYQKVVAGVANALA**HKYH----
LGNVLVVVLARHFGKDFTPELQA**SYQKVVAGVANALA**HKYH----
LSHCLLV**TLAAHLP**AEFTPAV**HASLDKFLASV**STVLT**TSKYR**----
LSHCLL**STLAVHLP**NDFTPAV**HASLDKFLSSV**STVLT**TSKYR**----
I**SEAI**IIHVL**HSRHPGDFGADAQG**AMNKALE**LFRKDIA**AKYKELGYQG
LAAVIADTV**AAAG**----D----AGFEK**ILMSMICILLRSAY**----
V**KEA**ILKT**IKEVVGAK**SEELNS**AWTIAYDEL**AI**VICKEMNDAA**---

: : : . : . : . : :



Human beta	-----VHLT PEEKSAVTALWGKV N--VDEVGGEALGRLLVVY PWT QRFFESFGDL ST
Horse beta	-----VQLS GEEKAAVLALWDKV N--EEEVGGEALGRLLVVY PWT QRFFDSFGDLSN
Human alpha	-----VLS PADKTNVKAAWGKVGAHAGEYGAELERMFLSFPTTKTYFPHF-DLS-
Horse alpha	-----VLS AADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFPHF-DLS-
Whale myoglobin	-----VLS EWEQWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKT
Lamprey globin	PIVDTGSVAPLS AAEKT KIRSAWAPVYSTYETSGVDILVKFFTSTPAAQE FFPKFKGLT
Lupin globin	-----GALT ESQAALVKSSWEEFNANIPKH THRFFILVLEIAPAAKD LFSFLKG TSE

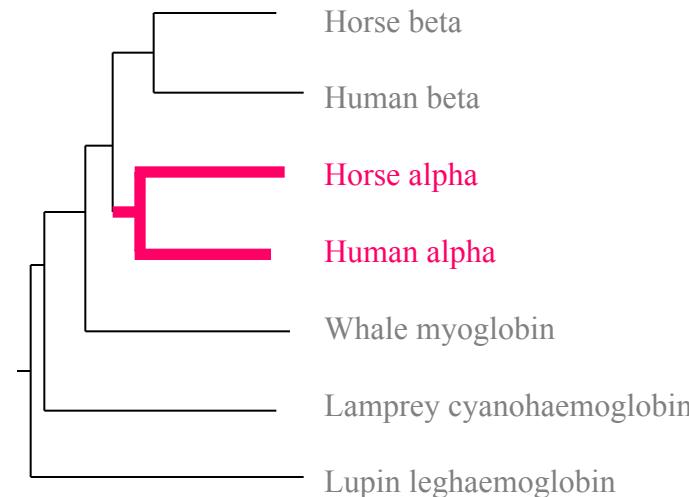
* : : : * . : : * : * : .

Human beta	PDAVMGNPKV KAHGKKVLGAFSDGLA HLDN----LKGTFA T SEL HCDKLHVD PENFRL
Horse beta	PGAVMGNPKV KAHGKKVLHSF GEGVHHLDN----LKGTFA AL SEL HCDKLHVD PENFRL
Human alpha	--- HGS AQVK GHGKKVADALTNAVAH VDD----MPNAL SALS SDL HAKL RVPVNFKL
Horse alpha	--- HGS AQVK KAHGKKVGDALT LAVG HLDD ----LP GALS NLS DLHAKL RVPVNFKL
Whale myoglobin	EAEMKASEDLKK HGTV TVL TALGAIL KKKGH----HEAE LKPLAQSHATKH KIPIKYLEF
Lamprey globin	ADQLKK SADVRW HAERI I NAVNDAVASMDT--EKMS MKLRDLSGKHAKSF QVDPQYFKV
Lupin globin	VP--QNN PELOQAHAGKVFKLVYEAA IQLQVTGVVVT DATLKNLGSVHVSKG VAD-AHF PV

. : : * : . : : : * . : .

Human beta	LGNVLVCV LAHHFGKEFTPPVQ AY QKVVAGVANALA HKYH----
Horse beta	LGNVLVV VVLA RH FGKDFTP ELQAS YQKVVAGVANALA HKYH----
Human alpha	LSHCLL VT LA AHL PAEFTP AVH A SLDK FLASV STVLT SKYR ----
Horse alpha	LSHCLL ST LA V HL PND F TPAVH A SLDK FLSSV STVLT SKYR ----
Whale myoglobin	I SEAI II HVLHSRHPGDFGADAQG AMNK ALELFRKDIA AKY KELGYQG
Lamprey globin	LAAVIAD TVAAG --D----AG FEKILMSMICILL R SAY -----
Lupin globin	V KEAILKT IK EVVGAK SEELNS AWTIAY DEL AI VI KEMNDAA --

: : : . : . . . : :



Human beta	-----VHLT PEEKSAVTALWGKVN --VDEVG GEALGRLLVVY PWT QRFFESFGDLST
Horse beta	-----VQLS GEEKAAVLALWDKVN --EEEVG GEALGRLLVVY PWT QRFFDSFGDLSN
Human alpha	-----VLS PADKTNVKAAWGKVGAHAGEYGAELERMFLSFPTTKTYFPHF-DLS
Horse alpha	-----VLS AADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFPHF-DLS
Whale myoglobin	-----VLS EWEQWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKT
Lamprey globin	PIVDTGSVAPLS AAEKT KIRSAWAPVYSTYETSGVDILVKFFTSPAAQE FFPKFKGLTT
Lupin globin	-----GALT ESQAALVKSSWEEFNANIPKH THRFFILVLEIAPAAKD LFSFLKG TSE

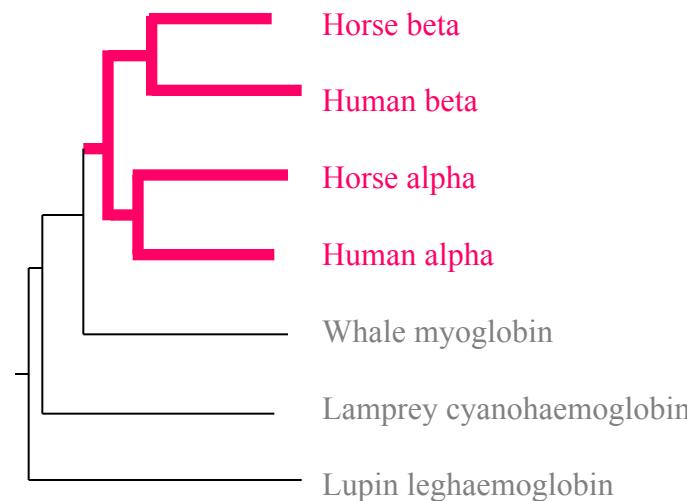
* : : : * . : : : * : * : .

Human beta	PDAVMGNPKV KAHGKKVLGAFSDGLA HLDN----LKGT FATLSELHCDKLHVD PENFRL
Horse beta	PGAVMGNPKV KAHGKKVLHSFGE GVHLDN----LKGT FAALSELHCDKLHVD PENFRL
Human alpha	---HGSAQVK GHGKKVADALTNAVA HVD----MPNALSALS DLHAKL RVPVNFKL
Horse alpha	---HGSAQVK AHGKKVGDA LTAV GH LDL--LPGALS NLSDLH AHKLRVPVNFKL
Whale myoglobin	EAEMKASEDLKK HGTV TL TA LGAILKKGH----HEAE ELKPLAQSH ATKH KIPI KYLEF
Lamprey globin	ADQLKK SADVRW HAERI I NAVNDAVASMDT--EKMS MKLRDLSGKHAKSF QVDPQYFKV
Lupin globin	VP--QNNPELO QAHAGKVFKLV YEAA I QLQVTGVVVT DATLKNLGSVHVS KGVAD-AHF PV

. : : * : . : : : * . : .

Human beta	LGNVLVCVLAHHFGKEFTPPVQ AY QKVVAGVANALA HKYH----
Horse beta	LGNVLVVVLARHFGKDFTPELQ ASY QKVVAGVANALA HKYH----
Human alpha	LSHCLLVTLA AHLPAEFTP AVH ASLDKFLASV STVLT SKYR ----
Horse alpha	LSHCLLSTLAVHLPNDFTPAVH ASLDKFLSSV STVLT SKYR ----
Whale myoglobin	I SEAI IIHVLHSRHPGDFGADAQGAMNKALELFRKDIA AKY KELGYQG
Lamprey globin	LAAVIADTV AAG--D----AG FEKILMSMICILLRSAY ----
Lupin globin	V KEAILKTIKEVVGAK SEELNS AWTIAYDEL AIVIK KEMNDAA --

: : : . : . : . : :



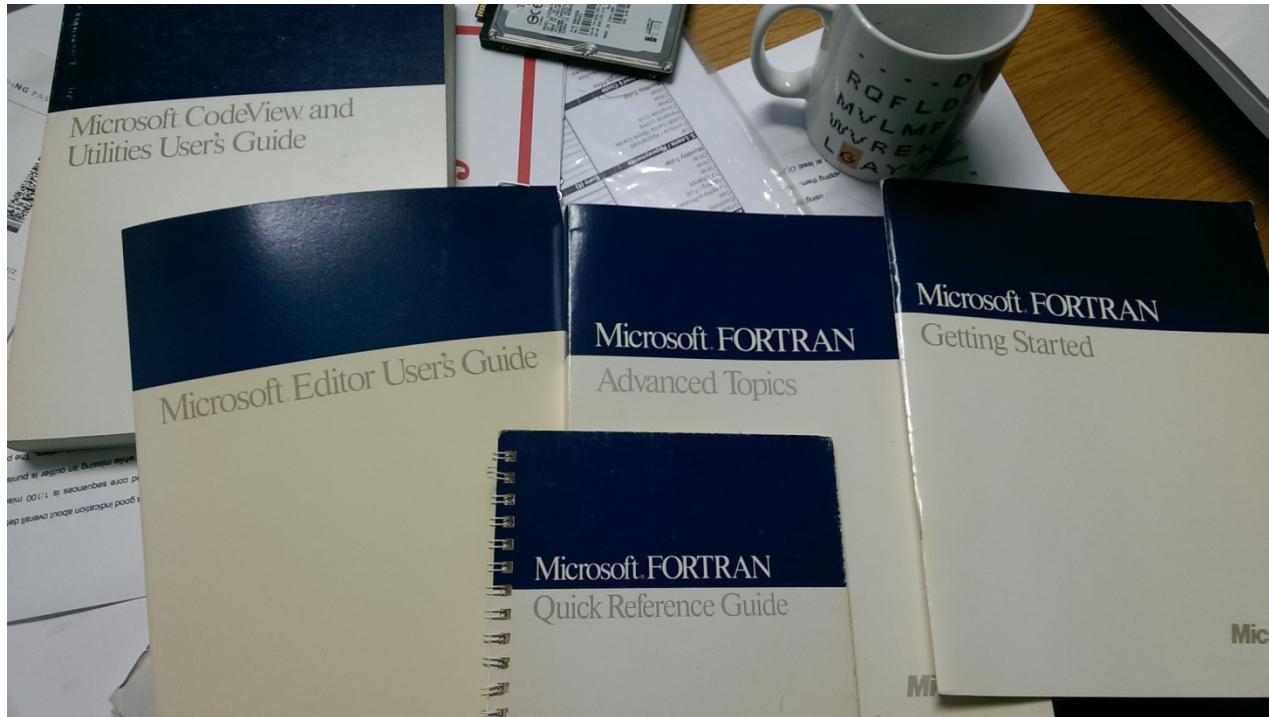
1988!

- Paul Sharp lab.
 - TCD
- No internet (much)
- IBM PC
 - MS-DOS
 - 640kb memory
 - 0-20MB hard disk



Fortran!

- MS Fortran
for MS-DOS
- No recursion
- No memory
management
(much)



Clustal 1-3

- Fast k-tuple pairwise alignments (Clustal 1)
- UPGMA (Clustal 2)
- Align larger and larger alignments (Clustal 3)
 - Fast approximate word based alignments



Gene
Volume 73, Issue 1, 15 December 1988, Pages 237-244

CLUSTAL: a package for performing multiple sequence alignment on a microcomputer

Desmond G. Higgins , Paul M. Sharp

Show more

[https://doi.org/10.1016/0378-1119\(88\)90330-7](https://doi.org/10.1016/0378-1119(88)90330-7)

[Get rights and content](#)

Abstract

An approach for performing multiple alignments of large numbers of amino acid or nucleotide sequences is described. The method is based on first deriving a phylogenetic tree from a matrix of all pairwise sequence similarity scores, obtained using a fast pairwise alignment algorithm. Then the multiple alignment is achieved from a series of pairwise alignments of clusters of sequences, following the order of branching in the tree. The method is sufficiently fast and economical with memory to be easily implemented on a microcomputer, and yet the results obtained are comparable to those from packages requiring mainframe computer facilities.



[Previous article in issue](#)



[Next article in issue](#)

Keywords

Cluster analysis; phylogenetic tree; protein secondary structure; RNA secondary structure; globin; 5S RNA; dendrogram

Abbreviations

aa, amino acid(s); ASCII, American Standard Code for Information Interchange; nucleotide(s); S, sedimentation constant; UPGMA, Unweighted Pair-Group Arithmetic Averages

Register to receive
based on you

[Register now](#)

2-Seq

Alignment?

- $O(L^2)$
- $L = 1000?$
 - 1-4 MB memory
- Myers and Miller, 1988

Optimal alignments in linear space

Eugene W. Myers^{1,2} and Webb Miller²

Abstract

Space, not time, is often the limiting factor when computing optimal sequence alignments, and a number of recent papers in the biology literature have proposed space-saving strategies. However, a 1975 computer science paper by Hirschberg presented a method that is superior to the new proposals, both in theory and in practice. The goal of this paper is to give Hirschberg's idea the visibility it deserves by developing a linear-space version of Gotoh's algorithm, which accommodates affine gap penalties. A portable C-software package implementing this algorithm is available on the BIONET free of charge.

Introduction

Consider the following problem. Given sequences $A = a_1a_2\dots a_M$ and $B = b_1b_2\dots b_N$, find a set of 'evolutionary operations' that converts A to B and minimizes the sum of the operations' costs. The allowed operations are (i) replace one symbol by another, (ii) delete k consecutive symbols, or (iii) insert k consecutive symbols. In addition, the problem statement requires that every symbol of A must be either replaced or deleted. Replacement costs are specified by a table, w , where $w(a,b)$ gives the cost of replacing a by b . Note that a symbol of A is effectively left unedited if it is replaced by itself at no cost, i.e. $w(a,a) = 0$. Two non-negative constants, g and h , specify an affine function, $\text{gap}(k) = g + hk$, for the cost of a k -symbol indel (insertion or deletion). Informally, opening up a gap costs g and each symbol in the gap costs h .

The problem is often formulated as maximizing the similarity score of an alignment, rather than minimizing the difference score of a conversion. A bonus $\sigma(a,b)$ is added for every aligned pair (a,b) and a 'gap penalty' $q + rk$ is subtracted for every k -symbol gap. This formulation is converted to a difference problem by the transformations

$$\begin{aligned}w(a,b) &= \sigma_{\max} - \sigma(a,b) \text{ for all pairs } (a,b) \\g &= q \\h &= r + \frac{1}{2}\sigma_{\max}\end{aligned}$$

¹Department of Computer Science, University of Arizona, Tucson, AZ 85721, USA

²Department of Computer Science, The Pennsylvania State University, University Park, PA 16802, USA

where $\sigma_{\max} = \max_{(a,b)} \sigma(a,b)$ (Smith *et al.*, 1981). Thus, to produce an alignment that maximizes the similarity score, first apply these transformations and then run the program described in this paper with the resulting w , g and h . If the minimum conversion score is C , then the corresponding maximum alignment score is $\frac{1}{2}(M + N)\sigma_{\max} - C$.

Gotoh (1982) gave an algorithm that solves such problems in $O(MN)$ time. If only the minimum cost is desired, then it is easy to implement the algorithm in $O(N)$ space, where N can be taken as the shorter sequence length. If one also desires a set of operations attaining the minimum cost, then straightforward implementations need $O(MN)$ space. In practice, this space requirement often limits the method's applicability, and several papers (Taylor, 1984; Watanabe *et al.*, 1985; Altschul and Erickson, 1986; Gotoh, 1986, 1987) have presented strategies that reduce space consumption by constant factors. These papers fail to note that Hirschberg (1975) showed how to produce an optimal conversion or alignment in $O(N)$ space. When only a single optimal alignment of A and B is desired, Hirschberg's approach is superior to the others. For example, in one megabyte of memory, our program based on Hirschberg's method can align two sequences of length 62 500. Altschul and Erickson (1986) propose keeping 7 bits for each of MN entries, so the limit for their method is $7N^2 \leq 8 \times 10^6$, or $N < 1070$. Moreover, any program that packs and unpacks bits or uses disk storage is doomed to be slow and, probably, non-portable.

$O(MN)$ -space methods permit the construction of all optimal alignments. However, the number of alignments that attain the minimum cost is often astronomical, in part because a brute force enumeration lists many arrangements whose differences are insignificant to the user. Moreover, when one is searching for a particular 'biologically meaningful' arrangement, it may be necessary to consider slightly sub-optimal alignments (Waterman, 1983; Waterman and Byers, 1985). One alternative to explicitly constructing all optimal alignments is to modify our linear-space program to produce 'left-most' and 'right-most' optimal alignments that delineate the range of possibilities. In any case, it is important to understand that a single optimal alignment can be found in far less space than is needed to record 'traceback' information for finding all optimal alignments.

Hirschberg's original presentation treats a simpler alignment problem, known as the longest common subsequence problem, where $w(a,b) = 1$ if $a \neq b$, $w(a,a) = 0$, and $\text{gap}(k) = k$. However, the approach is quite general. To the best of our

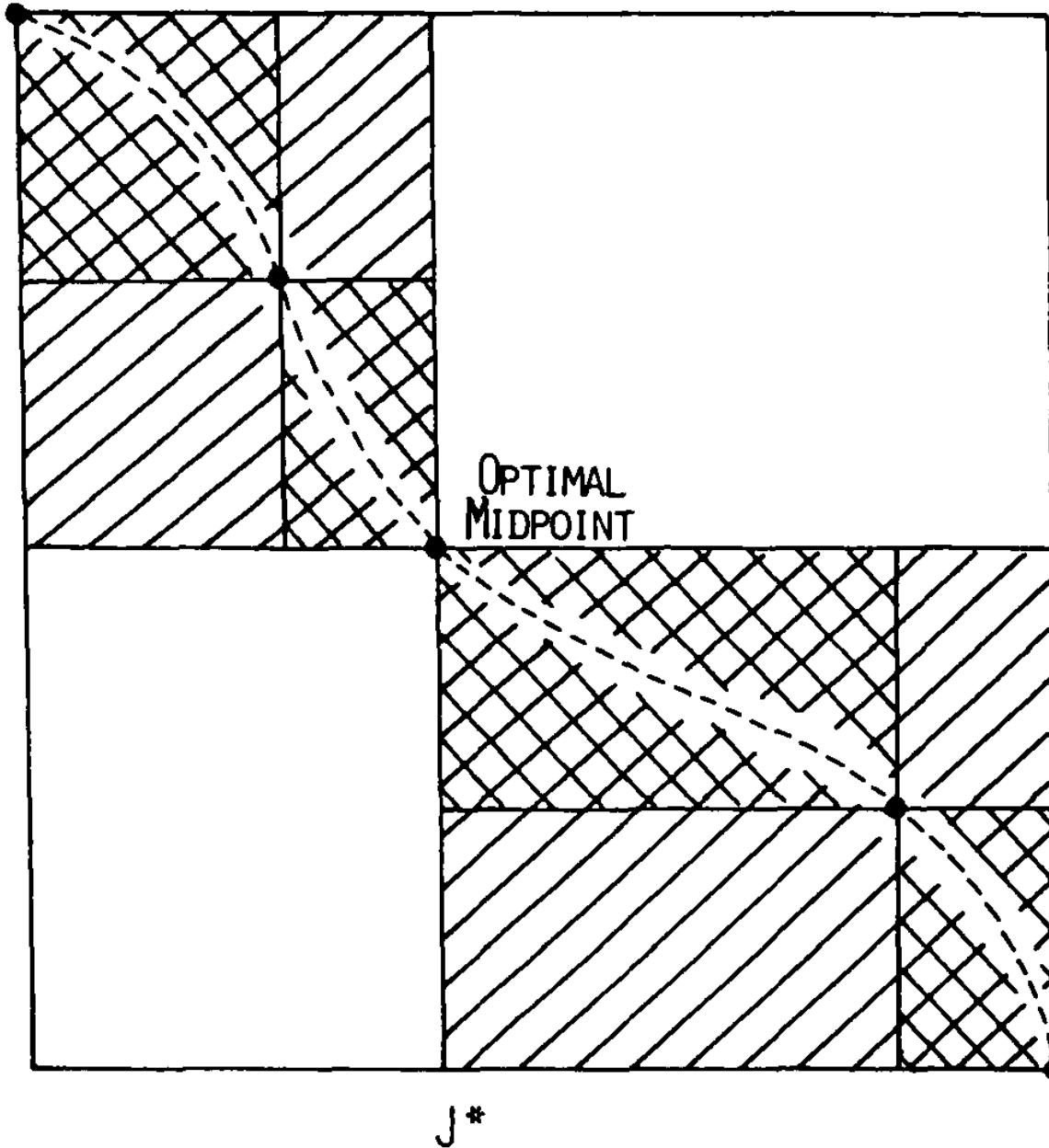


Fig. 2. Splitting the problem into sub-problems.

Clustal 4

- Full progressive alignment
- On PC
- Fast
- Accurate

CABIOS COMMUNICATIONS

Vol. 5, no. 2, 1989
Page 151–153

Fast and sensitive multiple sequence alignments on a microcomputer

Desmond G. Higgins* and Paul M. Sharp

Abstract

A strategy is described for the rapid alignment of many long nucleic acid or protein sequences on a microcomputer. The program described can handle up to 100 sequences of 1200 residues each. The approach is based on progressively aligning sequences according to the branching order in an initial phylogenetic tree. The results obtained using the package appear to be as sensitive as those from any other available method.

Introduction

In the recent literature on biological sequence analysis, at least a dozen methods for performing multiple alignments of nucleic acid or protein sequences have been described [e.g. Bains (1986), Sobel and Martinez (1986), Barton and Sternberg (1987), Feng and Doolittle (1987), Santibanez and Rohde (1987), Taylor (1987)]. The motivation for this effort has been the need for the automatic alignment of three or more sequences for the purposes of evolutionary or structural comparisons or for attempting to demonstrate similarity between sets of sequences. In this paper, we describe a strategy which we believe offers the best combination of speed and sensitivity available for any multiple alignment method. We offer a program which can perform multiple alignments of up to 100 sequences of maximum length 1200 residues on a microcomputer in a reasonable amount of time. We judge the program to be

for each residue in one set compared against each residue in the second set. Any gaps introduced into either set of sequences are scored as single gaps. The main difficulty in using this approach on a microcomputer arises from the excessive memory requirements of the Needleman and Wunsch (1970) method—memory usage is proportional to the square of the average sequence length.

In a previous paper (Higgins and Sharp, 1988) we described a strategy for the very rapid multiple alignment of large numbers of sequences on a microcomputer. This method also comprised a progressive approach, using the fast, but approximate, two-sequence alignment method of Wilbur and Lipman (1983). While this approach is extremely rapid and economical with core memory, it works well only for closely related sequences. We did not consider using the exactly optimal method of Needleman and Wunsch (1970) for the progressive alignments because of the excessive memory requirements. However, a recent paper by Myers and Miller (1988) demonstrates how to achieve exactly optimal alignments of two sequences where memory usage varies only linearly with sequence length, without making use of bit packing or secondary disk storage. Thus, a progressive series of alignments of larger and larger groups of sequences, using the method of Myers and Miller (1988) for each alignment, is the key to the current approach.

System

Distribution?

- Mainly by post
- The EMBL Mail Server

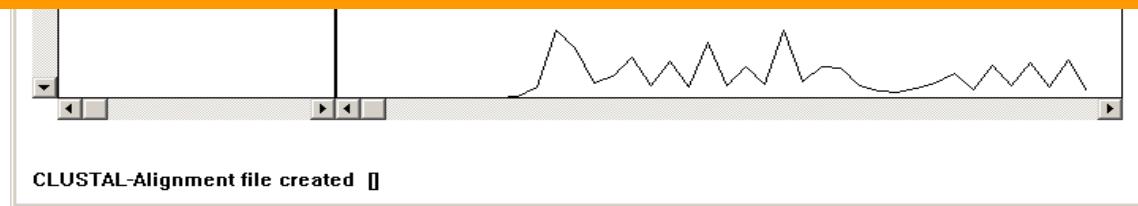
System

The program described in this paper was written in standard FORTRAN 77 and compiled using the Microsoft FORTRAN compiler, version 4.0. Program performance was tested on an IBM AT compatible microcomputer, running at 10 MHz with no maths coprocessor, 640 kbytes of memory and a hard disk. This program (CLUSTAL4) is an extension to the package described in Higgins and Sharp (1988). Copies of the executable files, documentation and test data files will be sent on request. Please send three 5.25 inch floppies formatted to 360 kbytes, or one high density 5.25 inch floppy formatted to 1.2 Mbytes.

Clustal

- Clustal1-Clustal4
 - 1988, Paul Sharp, Dub
- Clustal V 1992
 - EMBL Heidelberg,
 - Rainer Fuchs
 - Alan Bleasby
- Clustal W, Clustal X 1994-2007
 - EBI, University College Cork
 - Toby Gibson, EMBL, Heidelberg
 - Julie Thompson, ICGEB, Strasbourg
- Clustal W and Clustal X 2.0 2007
 - University College Dublin
- >130k citations on Google Scholar

“Guide Trees”
Bootstrapped NJ trees as output
With
Kimura 2-parameter distance for DNA or
Kimura empirical distance for protein



LOWRY, OH; ROSEBROUGH, NJ; FARR, AL; RANDAL
LAEMMLI, UK
BRADFORD, MM

<http://www.nature.com/news/the-top-100-papers-1.16224>

SANGER, F; NICKLEN, S; COULSON, AR
CHOMCZYNSKI, P; SACCHI, N
TOWBIN, H; STAEHELIN, T; GORDON, J
LEE, CT; YANG, WT; PARR, RG
FOLCH, J; LEES, M; STANLEY, GHS
BECKE, AD
THOMPSON, JD; HIGGINS, DG; GIBSON, TJ
KAPLAN, EL; MEIER, P
ALTSCHUL, SF; GISH, W; MILLER, W; MYERS, EW; L
Sheldrick, George M.
MURASHIGE, T; SKOOG, F
Altschul, SF; Madden, TL; Schaffer, AA; Zhang, JH; Zha
FOLSTEIN, MF; FOLSTEIN, SE; MCHUGH, PR
Perdew, JP; Burke, K; Ernzerhof, M
SOUTHERN, EM
BLIGH, EG; DYER, WJ
SAITOU, N; NEI, M
COX, DR
SHANNON, RD
Otwinowski, Z; Minor, W
Livak, KJ; Schmittgen, TD
BECKE, AD
DUBOIS, M; GILLES, KA; HAMILTON, JK; REBERS, PA
REYNOLDS, ES
WEBER, K; OSBORN, M
Thompson, JD; Gibson, TJ; Plewniak, F; Jeanmougin, F
CHIRGIN, JM; PRZYBYLA, AE; MACDONALD, RJ; R

Top 30 papers on ISI up to 2015

Since Clustal...

- >100 new programs in past 15 years

- 1999: BENCHMARK

Balibase, Oxbench, Prefab, Sabre

- T-Coffee (and R-Coffee, M-Coffee, 3D-Coffee)

- Cedric Notredame
<http://www.tcoffee.org>

- MAFFT

- Kazutaka Katoh, 2001
 - <http://mafft.cbrc.jp/alignment/software/>



- MUSCLE

- Bob Edgar, ISMB 2004
 - <http://www.drive5.com/muscle>

- PROBCONS

- Tom Do, Michael Brudno, Serafim Batzoglou, ISMB 2004

MSAProbs

Yongchao Liu, Bertil Schmidt, Douglas L. Maskell (2010)

"MSAProbs: multiple sequence alignment based on pair hidden Markov models and partition function posterior probabilities".
Bioinformatics, 26(16): 1958 -1964

Consistency!

Two kinds of programs

- Fast (<10,000 seqs)
 - Clustal W
 - MAFFT
 - mafft --partree >>10,000 seqs
 - MUSCLE
 - Kalign
- Accurate (100s seqs)
 - T-Coffee
 - ProbCons
 - MAFFT
 - MSAProbs

Prank and Pagan

Accurate?

- Protein structure comparisons
- Simulation
- Protein structure versus phylogeny
- Accurate analysis of indels/closely related sequences/analysing sites and selection?
 - Use e.g. Prank

So: What is the Problem?

- What if $N \gg 100,000$?
- e.g. SSU rRNA
 - www.arb-silva.de
 - **4,346,367** seqs
- e.g. ABC transporters
 - PFAM May 2015
 - ABC_tran PF00005
 - **150k** seqs
- Metagenomics

GENOME 10K.
Unveiling animal diversity

Genome 10K Project

To understand how life has evolved and become what it is today, we must sequence every vertebrate genome. This project will help us understand how life would change if we stopped our sequencing efforts. The group of scientists involved in this project is dedicated to coordinating efforts in tissue specimen collection that will lay the groundwork for a large-scale sequencing and analysis project.

•Sequence 10,000 vertebrate genomes! =>5,000,000 protein kinases, GPCRs

Join us

Co-directors

David Haussler
haussler@soe.ucsc.edu
831-459-1477
CBSE/ITI
UC Santa Cruz
1156 High Street
Santa Cruz, CA 95064

Stephen J. O'Brien
stephen.obrien@nih.gov
301-846-1296
National Cancer Institute
Laboratory of Genomic Diversity
Frederick, MD 21702

Oliver A. Ryder
orvder@sandiegozoo.org

Location: Chaminade Resort
Santa Cruz, CA

G10K meeting
March 16-18, 2011
Registration open until
January 31, 2011
First-come, first-served
Location: Chaminade Resort
Santa Cruz, CA

Accomplishments

► G10K announces first
101 species for
sequencing

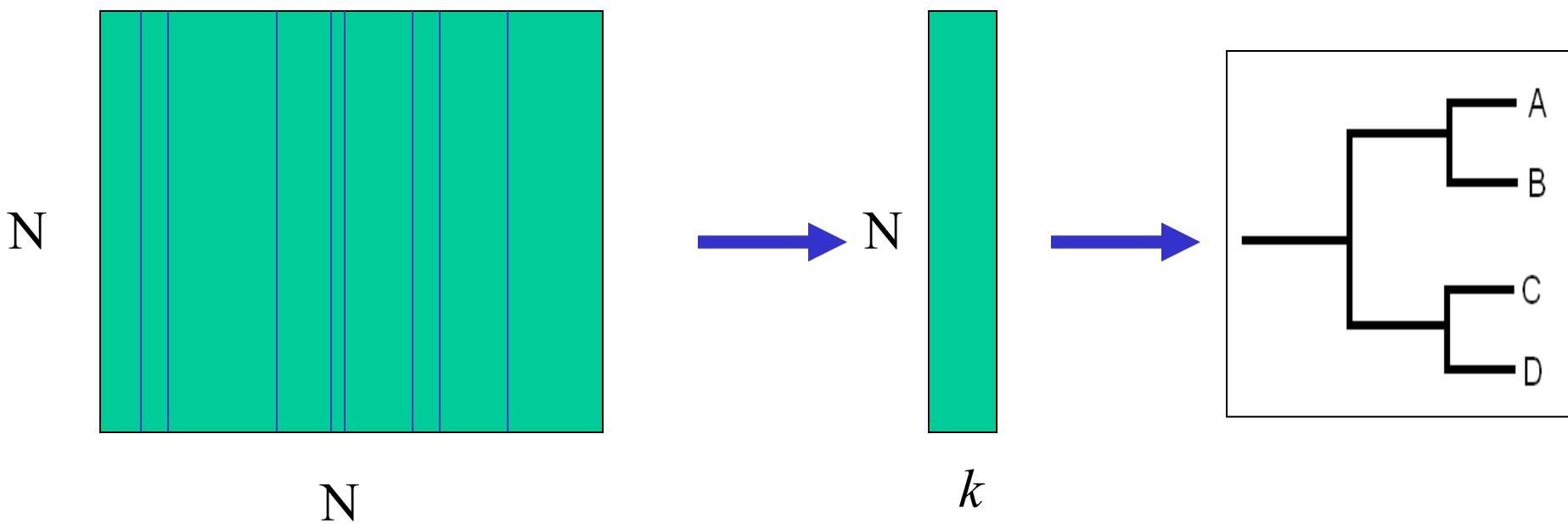


Big Alignments?

- Once a guide tree is made:
 - Multiple alignment: for N seqs, fixed len L
 - Time $O(N)$
- To make a guide tree:
 - Time and memory: $O(N^2)$
- Maximum alignment size 10-20k seqs.

mBED

k seeds



Blackshields G, Sievers F, Shi W, Wilm A, Higgins DG. (2010)
Sequence embedding for fast construction of guide trees for multiple sequence alignment.
Algorithms Mol Biol. 14;5:21.

Clustal Ω

- Released April 2011
- Scalable
 - mBed
 - Blackshields, et al. (2010) Algorithms in Molecular Biology.
- Accurate
 - Hidden Markov models (HMMs)
 - HHalign
 - Johannes Söding, Munich.
- Re-use old alignments
 - Kevin Karplus
 - UCSC



REPORT

Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega

Fabian Sievers^{1,8}, Andreas Wilm^{2,8}, David Dineen¹, Toby J Gibson³, Kevin Karplus⁴, Weizhong Li⁵, Rodrigo Lopez⁵, Hamish McWilliam⁵, Michael Remmert⁶, Johannes Söding⁶, Julie D Thompson⁷ and Desmond G Higgins^{1,*}

¹ School of Medicine and Medical Science, UCD Conway Institute of Biomolecular and Biomedical Research, University College Dublin, Dublin, Ireland,

² Computational and Systems Biology, Genome Institute of Singapore, Singapore, ³ Structural and Computational Biology Unit, European Molecular Biology Laboratory, Heidelberg, Germany, ⁴ Department of Biomolecular Engineering, University of California, Santa Cruz, CA, USA, ⁵ EMBL Outstation—European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, UK, ⁶ Gene Center Munich, University of Munich (LMU), Muenchen, Germany and

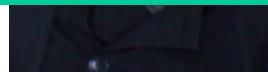
⁷ Département de Biologie Structurale et Génomique, IGBMC (Institut de Génétique et de Biologie Moléculaire et Cellulaire), CNRS/INSERM/Université de Strasbourg, Illkirch, France

⁸ These authors contributed equally to this work

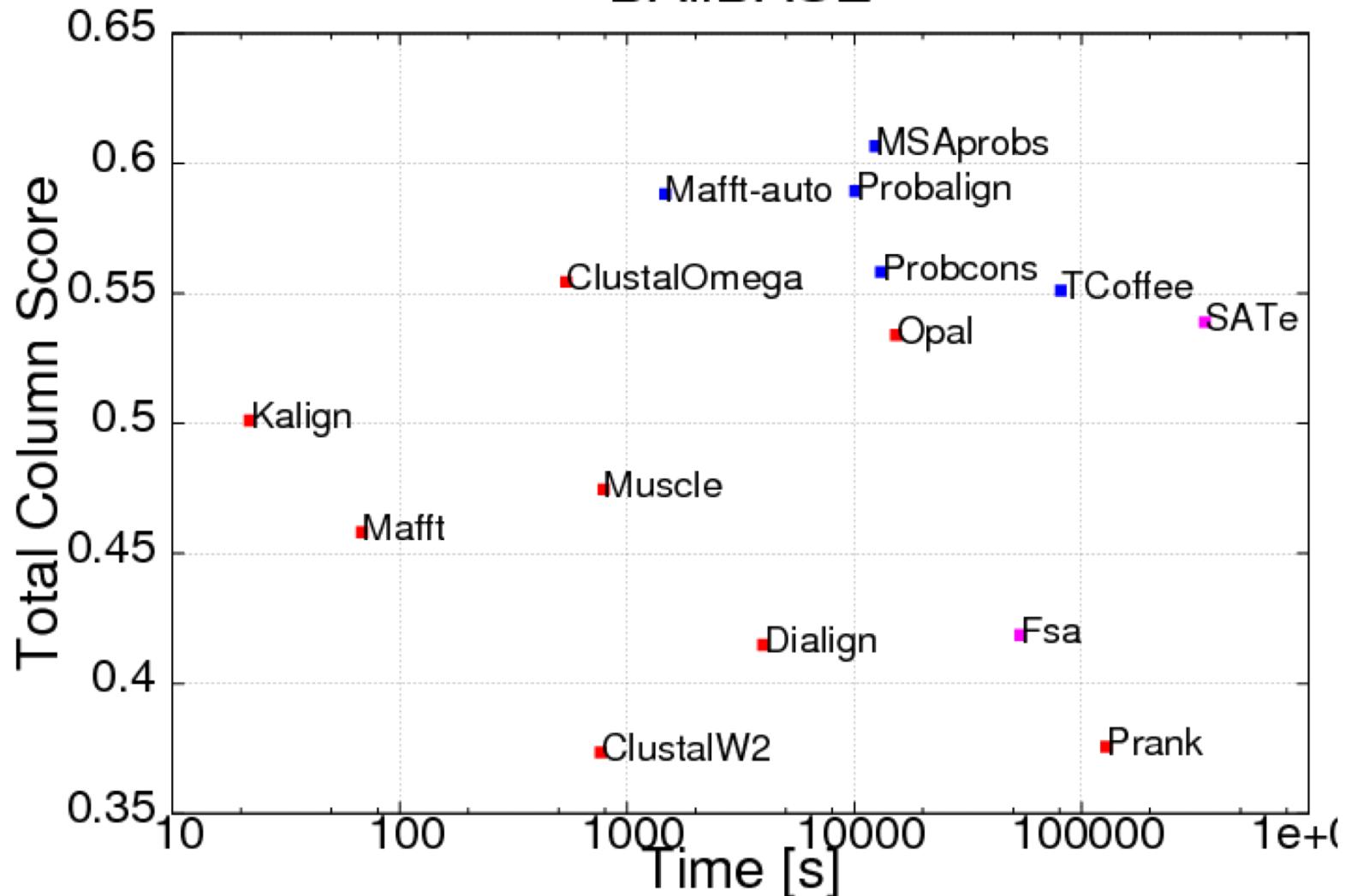
* Corresponding author. UCD Conway Institute of Biomolecular and Biomedical Research, University College Dublin, Belfield, Dublin 4, Ireland. Tel.: + 353 1 716 6833; Fax: + 353 1 716 6713; E-mail: des.higgins@ucd.ie

Received 23.7.11; accepted 6.9.11

• Alignments	BaliBase	% correct	time(s)
– 6			
• More	Clustal Omega	55.4	539
– M	Clustal W	37.4	766
– A	Mafft (default)	45.8	68
• Dow	Muscle	47.5	789
– w	Kalign	50.1	21
	T-Coffee	55.1	81041
	Probcons	55.8	13086
	Mafft (auto/consistency)	58.8	1475
	MsaProbs	60.7	12382



BALiBASE



Clustal Omega Servers

- <http://www.ebi.ac.uk/Tools/msa>
- Download your own copy
 - Command-line only (Unix/Linux style)
 - Visualise with Jalview or SeaView (Manolo Gouy)
 - Proteins or DNA/RNA
 - www.clustal.org

Clustal

Paul Sharp, Edinburgh

Rainer Fuchs, EMBL

Alan Bleasby, Daresbury

Toby Gibson, Ramu Chenna, Nigel Brown, EMBL

Julie Thompson, Francois Jeanmougin, Fred Plewniak, Strasbourg

Paul Mcgettigan, Mark Larkin, Andreas Wilm, Fabian Sievers, UCD

mBED

Gordon Blackshields

Mark Larkin

Clustal Omega

Fabian Sievers, Andreas Wilm, David Dineen, UCD

Johannes Soeding, Michael Remmert, Munich

Rodrigo Lopez, Hamish MacWilliam, Weizhong Li, EBI

Kevin Karplus, UCSC

Current Alignment Group

Fabian Sievers, Quan Le

Kieran Boyce, Gearóid Fox and Peter Jehl

