

# Searching sequence databases

**MBV-INF4410**

**Thursday 10 September 2009**

**Torbjørn Rognes**

**Centre for Molecular Biology and Neuroscience (CMBN), Rikshospitalet  
& Department of Informatics, University of Oslo  
torognes@ifi.uio.no**

# Overview of the presentation

- An example showing how useful database searches can be
- Searching sequence databases
- A walk-through of the BLAST search service
  - BLAST variants
  - About sequence databases
- What is a good match?
- What is homology?
- Significance of alignments
- Problematic sequences

# One example of how useful database searches can be

- The protein AlkB was discovered in the bacterium *E.coli* more than 20 years ago.
- It was known that it protected the bacterium when subjected to DNA-alkylation agents.
- No enzymatic activity was found.
- Perhaps some co-factors were missing?
- In 2001, a bioinformatics paper was published that shed light on the problem. Many similar sequences were found using advanced sequence similarity searches ...

<http://genomebiology.com/2001/2/3/research/0007.1>

Research

**The DNA-repair protein AlkB, EGL-9, and Iprecan define new families of 2-oxoglutarate- and iron-dependent dioxygenases**

L Aravind and Eugene V Koonin

Address: National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA.

Correspondence: L Aravind. E-mail: [aravind@ncbi.nlm.nih.gov](mailto:aravind@ncbi.nlm.nih.gov)

# Example...

Alignment showing conserved amino acids among many sequences

```
CAS_Scla_322266 -----YHRLQPNYMLACSRADHE-----RTAATLVASVRK---70---VTEAVYLEPG-DLLIVDNF-----RTHARTPFSPRWGDKRWLHRVYIRT 302\
IPNS_Bn_124825 -----LITVLYVQ-----SNVQNLQVETAA-----GYQDIADDT-GYLINCGSYMAHLTNNYKAPHRKQVNVN----AERQSLPFFVNL 288|
FLAS_Pet_421946 -----YIILVLP-----NEVQGLQVFKDG-----HWYDVKIYIP-ALIVHIGDQVBIILNGKYKYSVKEHTIVNK----DKTRMSWVFLFEP 309|
LDOX_Pet_1730108 -----ALTFILH-----NMVPLGLQLFYEG-----QWVTAKCVPM-SIIMHIGDTIIBILSNGKYKYSILHRGVVVK----EKVRFSSWAIFCEP 311|
Srg_At_479047 -----GLTVLMQV-----NDVEGLQIKKDG-----KWPVVKPLRN-AFIVNIGDVLBIITNGTYRSIEBERGVVNS----EKERLSIATFHNV 309|
EFE_Le_398992 -----GIILLLFD-----DKVSGLLQLLKD-----QWIDVPPMRH-SIVVNLGDQLBEVITNGKYKYSVLRVIAQT---DGTFRMSLASFYMP 253| Small
Ga20Ox_Sot_10800976 -----SLTILHQ-----DSVSGLVQVMDN-----QWRISISPNLS-AFVVNIGDTFMALSNGRYKSCLRRAVVNN---KTRPKSLAFFLCP 317| molecule
PA0147_Pa_9945977 -----CVTLLYQ-----DAAGGLQVQNRQG-----EWIDAPPIDG-TFVVNIGDMARWSDRYRSTERRVISP---GVHRYSMPPFAEP 274|
PA4191_Pa_9950401 -----LITLLHQ-----DAIGGLQVTRTPQ-----GWLEAPPIDG-SFVGNLGDMLERMTGGLYRSTERRVARNTS---GRDRLSLFLFFDP 277|
ISP7_Sp_729862 -----ALTLMSQ-----DNVKGLEILLDPVSN-----CFLSVSPAPG-ALIANLGDIMAILTNNRYKSSMIRVCNNS---GSDRYTIPFFLQG 353|
SPCC1494.01_Sc_7491815 -----SITLLFQ-----RDAAGLEIRPPNFVKDM---DWIKVNVQPD-VVLVNIADMLQFWTSGLKLRSTVRRVRIIDPG---VKTRQTIAYFVTP 267|
DACC5_Ly1_769809 -----IVELLITQPCP-----NGFVSLQVEIDG-----RFVEVPPRPG-CVVVFCGSIAPLVSDBGKIKAPQRRVVS-PGA4-GSNRTSSVLFLLRP 268/

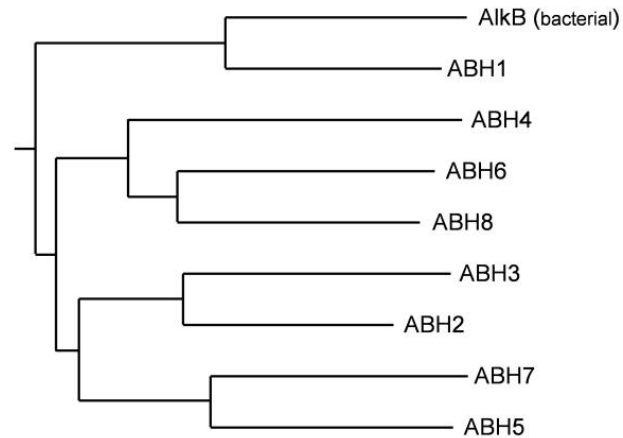
RRPO_SHVX_548840 IYPKG-----NKILTVNAA-----GSGTFSI-----KCAKGE-TTLNLEDGD-YFQMPSGFQETHRKNVA----VTPRLSITFRSTV 743\
POL_ASPV_487652 IYDIN-----HQVLTVNSY-----GDAIFCI-----ECLGSGF-EIPLSGPQ-MLLMPFGFQKSHRSGIKSP---SKGRISLTLFRLLK 853|
POL_BSV_409711 IFMRG-----APVHTVSM-----GNADFGT-----ECAAGR-QYTLTRGNVQFTMPSGFQETHRKNVNT---TAGRVSITFRLLK 841| RNA
RRPO_PMW_139137 CYPKG-----HQVLTINHS-----GCLTQI-----ACQKGA-SITMGGD-YLSPVGFQESHKHAISNT---TGGRVSLTRFRCTV 690| viral
POL_GLV_1154656 IFEKD-----SKILTVCIQ-----GDCBFFF-----RCATGET-GFYMEAPK-QFMMPDGFQSNHVEATREBC---TPGRISAIFRRAK 772| AlkB
Pol_GVA_1405615 CYLPG-----GSVVTVNLH-----GDAATFEVK-----ENQSGKIEKKELHDGD-VYVMSPGMQQTTHKRVVTSH---TDGRCSITLNRKT 738| homologa
RRPO_ACLSV_1710717 CYDD-----DEILTINVV-----GDAKFHT-----TC-HGE--IIDLRQGD-BILMPGGYQKMNKRAVEVA---SEGRTSVTLRVHK 836/
T13L16.2_At_2708738 FL-----RPFCTISFL-----SECNILEGNSNLKVE-----GPGDFSGY-SIPLPVGS-VLVLKNGADVARKVCPAV---FTKRISITFRKMD 420\
T19K4.220_At_3036813 FL-----RPFCTVSL-----SECNILEGNSNLKVL-----GPGDFSGY-SIPLPVGS-VLVLKNGADVARKVCPAV---FTKRISITFRKMD 403|
At2g48080_At_4249414 -----QPISLVL-----SESTMVFGHRLGVD-----NDGNFRGL-TLPLKEGS-LLVMRGNADMARVVCPS---FNKRVAITFFFLK 351|
AK000315.1_Hs_7020317 IFE-----RPIVSVSFF-----SDCALCFGCKFQFK-----PIRVSEPVLSLFPVRRGS-VTVLSGYAADBITHCIRPQDI---KERRAVILRKR 270|
CG17807_Dm_7291441 AFL-----DPILSLSLQ-----SLVVMDFRRG-----DDQV-QVRLPRRS-LLMSGEARVDWBEIRPKHID13RGRKTSITFRRLR 325| Eukaryotic
CG6144_Dm_7297712 FH-----PIIISTISG-----AHTVLEFVKREDTTETBAGDQTTREVLV-KLLLEPRS-LLILKDTLYTDYLAISSETSBD24RSPKISLTIKRV 213| Family of
CG4036_Dm_7297561 IWGERVTVNC-----LGDVLTLT--PYEVQSGKYNLDLVAIYDELLAP-LLTDDQLATFEGKVLRIPIPNLS-LIVLYGPARYQFERSVLRDQV---QERFVCAVREFT 278| AlkB
FLJ2001_Hs_38923019 LWGERLVSLNL-----LSPVLSMC-----REAPGSLLLCSAPSAPEALVDSVIAPSRVLCQSEVAIPLPARS-LLVLTGAARHQWKAIAHRRHI---BARFVCTVFRRLS 274| paraloga
C14B1.10_Ce_6580210 AFD-----DPIVSISLL-----SDVVMDFKD-----GANSARIAPVLLKARS-LCLIQGESRYRWKRGIVNRRYD10RQTFVSLTLRKR 343|
SPAP8A3.02c_Sp_7491301 FGDG-----VAIFGFLSN-----TIMIFTRPE-----LFLKE--KIRLEKGS-LLMSGTARYDWRBEIPFRAGD12RSQRLSITMRRII 219|
L3377.4_Lm_9989036 VYD-----DIFAI CSLG-----SNCLLRFVH-----VQNGBEL-DVMVPDRS-VYIMSGEARYVYVFWLFPV---BAQFSLVFRRSI 193/
MTC1237.14c_Mtu_2052134 RGSTEDTM-----VAIVSLGAT-----RVFALRP-----RGRGSLRLPLAHDG-LLVMGGSCQRTTFHVPKTSAP--TGRVSVIQFRPRD 203\
ALKB_Cc_2055386 ADPR-----FPLLSISLQ-----DIAVFRIGG-----VNRKDTFRSLRLAAGD--VCRLLGPARLAREGVDRIPLPG6-GGGRIINLTLRAR 190|
ALKB_Bc_113638 PDLR-----APIVSVSLG-----LEAIFQFGG-----LKRNDPLKRLLEHGD--VVVWGSBSRLFYGGIQLPKAG5-IDCFYNTLFRQAG 213| Classic
ALKB_Scoe_8894829 RTD-----APVSVLSLQ-----DTCVFRFGN-----PETRTRPTDTELRSGD--LFVFGGSRSLAYGSRVHPG7-LRGRINLTLRVSQ 215| AlkB
ALKB_At_4835778 ADWS-----KPIVMSLSG-----CKAIFLLGGK-----SKDDPEHAMVLRSGD--VVLWAGEARBCEGILLHFQL34KTSRININIRQVF 354|
ALKB_Sp_3080529 BDLT-----LPLILSLWG-----LDCIVLIGTE-----SRSEKPS-ALRLHSGD--VVIMTGSRKARHGKHC---SFYKLYLSQLIA 272|
ALKB_Hs_2134723 LDHS-----KPLLSFSFG-----QSAIFLLGGL-----QRDEAPP-PFMHSGD--IMIMSGFSRLLEAIFRVLPFN39KTA RVNMA RQVL 272/
Consensus (85%): .....sh.h.....s..h.....s..h.....H.s.....+h.h..b...
```

# Example...

- By comparing *E.coli* AlkB to other sequences in the database it was found that AlkB had some features in common with more well-known enzymes
- Based on these similarities the following was suggested regarding AlkB:
  - That AlkB is a dioxygenase
  - That the enzyme is Iron(II) dependent
  - That the enzyme is 2-oxo-glutarate dependent
  - That AlkB repairs alkylated bases through a form of oxidation
  - That the enzyme could demethylate RNA as well (not just DNA)
  - That there were eukaryotic counterparts of the protein
- All of this was later verified in the lab and resulted in three publications in *Nature*.

# Example...

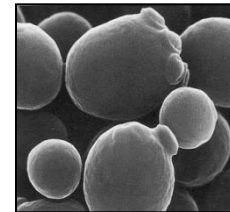
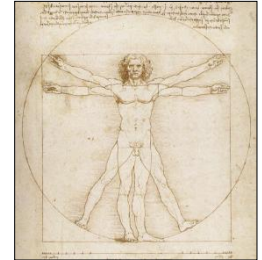
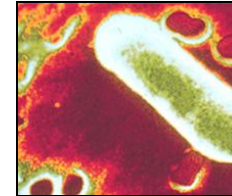
- By further sequence analysis 3 AlkB-like sequences were found in humans:
  - ALKBH1
  - ALKBH2
  - ALKBH3
- And by even more advanced analysis another 5 homologs were found in humans:
  - ALKBH4
  - ALKBH5
  - ALKBH6
  - ALKBH7
  - ALKBH8



- The function of these 8 enzymes are now being studied in detail. Some of them may be related to human diseases.

# Genomes are a huge source of information

- More than 900 completely sequenced genomes available – an enormous source of information. More than 3 000 other genomes underway...
- More than 250 000 000 000 basepairs sequenced to date \*
- Database sizes are growing exponentially – doubling in about 18 months
- Searching sequence databases for a similar sequence is fundamental in many types of analyses in bioinformatics
- Searching a sequence database with a new amino acid or nucleotide sequence allow us to find out more about:
  - Gene function
  - Conserved and probably important residues
  - 3D structure of a protein
  - Distribution of the gene among species
  - Gene structure
  - Chromosomal localisation
- Save time in the lab!
- Database searching is highly compute intensive and is probably the task consuming the largest amount of computing time within bioinformatics.



\* Source: <ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt>

# Searching sequence databases

- Goal: Identify which sequences in a database are significantly similar to a given DNA, RNA or protein sequence.
- How: The query sequence is compared (aligned) with each of the database sequences, and the amount of similarity is determined for each database sequence.

## Example:



Query sequence:

acgatcgattagcca

Database sequences:

Identical (trivial):

acgatcgattagcca

Very similar (easy):

acga**c**cgat**g**agcca

Similar (moderate):

a**t**ga**c**ggat**g**agc**g**a

Very diverged (hard):

a**t**ga**c**gggat**g**agc**g**a



NCBI/ BLAST Home

BLAST finds regions of similarity between biological sequences. more...

Learn more about how to use the new BLAST design

BLAST Assembled Genomes

Choose a species genome to search, or list all genomic BLAST databases.

- Human Mouse Rat Arabidopsis thaliana Oryza sativa Bos taurus Danio rerio Drosophila melanogaster Gallus gallus Pan troglodytes Microbes Apis mellifera

Basic BLAST

Choose a BLAST program to run.

- nucleotide blast Search a nucleotide database using a nucleotide query Algorithms: blastn, megablast, discontinuous megablast protein blast Search protein database using a protein query Algorithms: blastp, psi-blast, phi-blast blastx Search protein database using a translated nucleotide query tblastn Search translated nucleotide database using a protein query tblastx Search translated nucleotide database using a translated nucleotide query

News

New Human and Mouse pre-indexed databases

Human and mouse genomic + transcript megablast searches now use a faster, indexed algorithm that typically reduces run time by two thirds, as compared with standard megablast.

2007-09-04 10:55:00

More BLAST news...

Tip of the Day

How to save custom search pages.

So you have made a few BLAST searches and after adjusting the database, organism limits and maybe a few Algorithm Parameters you arrive at what you think is a good search strategy. Do you want to have to fiddle with pull down menus or remember all the changes you made the next time to want to run a similar search? Now you can use

# Search program variants

Query	Database	Comparisons	FASTA	BLAST	Description
Nucleotide	Nucleotide	Nucleotide (2)	fasta (fastn)	blastn	Compares directly both strands (forward and reverse complement) of the nucleotide query sequence to the nucleotide sequences in the database.
Amino acid	Amino acid	Amino acid (1)	fasta (fastp)	blastp	Compares the amino acid query sequence with the amino acid sequences in the database.
Amino acid	Nucleotide	Amino acid (6)	tfasta, tfastx, tfasty	tblastn	Translates the database nucleotide sequences into all six frames and compares the resulting amino acid sequences with the amino acid query sequences. tfasty allows intra-codon substitutions and frameshifts.
Nucleotide	Amino acid	Amino acid (6)	fastx, fasty	blastx	Translates the nucleotide query sequence into all six frames and compares the resulting amino acid sequences with the amino acid sequences in the database. fasty allows intra-codon substitutions and frameshifts.
Nucleotide	Nucleotide	Amino acid (36)	-	tblastx	Translates both the query nucleotide sequence and the database nucleotide sequences into all six frames and compares the resulting amino acid sequences with each other.

### Enter Query Sequence

Enter accession number, gi, or FASTA sequence [Clear](#)

```
MLDLFADAEFPWQEPPLAAGAVILRRFAFNAAEQLIRDINDVASQSPFRQMVTTCGYTMSVA
MTNCGHLGWTTHRQCYLYSPIDPQTNKWPWPAMPQSFHNLQORAATAAGYPDFQPDACLIN
RYAPGAKLSLHQDRDEPDLR&PIVSVSLGLPAIFQFCGLKRNDPLKRLLLEHGDVVVWGC
ESRLFYHGIQPLKACFHPLTIDCRYNLTFRQ&GKKE
```

Query subrange [Clear](#)

From

To

Or, upload file

Job Title

Enter a descriptive title for your BLAST search [?](#)

### Choose Search Set

Database

Swissprot protein sequences([swissprot](#)) [?](#)

Organism  
Optional

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. [?](#)

Entrez Query  
Optional

Enter an Entrez query to limit search [?](#)

### Program Selection

Algorithm

- blastp (protein-protein BLAST)
  - PSI-BLAST (Position-Specific Iterated BLAST)
  - PHI-BLAST (Pattern Hit Initiated BLAST)
- Choose a BLAST algorithm [?](#)

**BLAST**

Search **database swissprot** using **Blastp (protein-protein BLAST)**

Show results in a new window

[Algorithm parameters](#)

Note: Parameter values that differ from the default are highlighted in yellow

# BLAST databases (protein)

- nr:** All non-redundant GenBank CDS translations + RefSeq Proteins + PDB + SwissProt + PIR + PRF
- refseq:** RefSeq protein sequences from NCBI's Reference Sequence Project.
- swissprot:** Last major release of the SWISS-PROT protein sequence database (no updates).
- pat:** Proteins from the Patent division of GenPept.
- pdb:** Sequences derived from the 3-dimensional structure from Brookhaven Protein Data Bank.
- env\_nr:** Protein sequences from environmental samples.

# BLAST databases (nucleotides)

<b>nr:</b>	All GenBank + RefSeq Nucleotides + EMBL + DDBJ + PDB sequences (excluding HTGS0,1,2, EST, GSS, STS, PAT, WGS). No longer "non-redundant".
<b>refseq_rna:</b>	RNA entries from NCBI's Reference Sequence project
<b>refseq_genomic:</b>	Genomic entries from NCBI's Reference Sequence project
<b>est:</b>	Database of GenBank + EMBL + DDBJ sequences from EST Divisions
<b>est_human:</b>	Human subset of est.
<b>est_mouse:</b>	Mouse subset.
<b>est_others:</b>	Non-Mouse, non-Human subset of est.
<b>gss:</b>	Genome Survey Sequence, includes single-pass genomic data, exon-trapped sequences, and Alu PCR sequences.
<b>htgs:</b>	Unfinished High Throughput Genomic Sequences: phases 0, 1 and 2 (finished, phase 3 HTG sequences are in nr)
<b>pat:</b>	Nucleotides from the Patent division of GenBank.
<b>pdb:</b>	Sequences derived from the 3-dimensional structure from Brookhaven Protein Data Bank
<b>month:</b>	All new or revised GenBank + EMBL + DDBJ + PDB sequences released in the last 30 days.
<b>dbsts:</b>	Database of GenBank+EMBL+DDBJ sequences from STS Divisions .
<b>chromosome:</b>	A database with complete genomes and chromosomes from the NCBI Reference Sequence project..
<b>wgs:</b>	A database for whole genome shotgun sequence entries.
<b>env_nt:</b>	Nucleotide sequences from environmental samples, including those from Sargasso Sea and Mine Drainage projects.

Note: Parameter values that differ from the default are highlighted in yellow

Algorithm parameters

General Parameters

**Max target sequences**  Select the maximum number of aligned sequences to display

**Short queries**  Automatically adjust parameters for short input sequences

**Expect threshold**

**Word size**

Scoring Parameters

**Matrix**

**Gap Costs**

**Compositional adjustments**

Filters and Masking

**Filter**  Low complexity regions

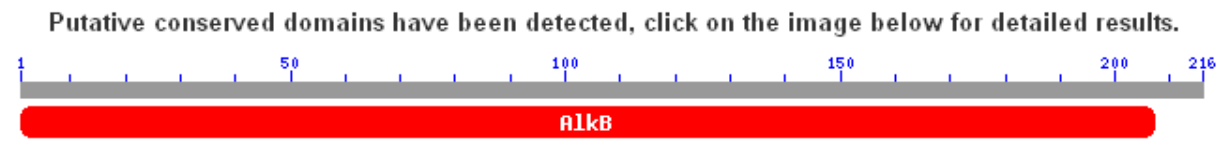
**Mask**  Mask for lookup table only  
 Mask lower case letters

**BLAST**

Search **database swissprot** using **Blastp (protein-protein BLAST)**  
 Show results in a new window

NCBI/BLAST/blast/Formatting Results - JE920YVS012 [\[Formatting options\]](#)

### Job Title: Protein sequence(216 letters)



Request ID	JE920YVS012
Status	Searching
Submitted at	Tue Oct 30 08:32:00 2007
Current time	Tue Oct 30 08:32:07 2007
Time since submission	

This page will be automatically updated in 12 seconds

NCBI/BLAST/blastp/Formatting Results - JE920YVS012 [\[Reformat these Results\]](#) [\[Edit and Resubmit\]](#) [Sign in above to save your search strategy]

**Job Title: Protein sequence(216 letters)** [▶ Show Conserved Domains](#)

BLASTP 2.2.17 (Aug-26-2007)

**Reference:**  
Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schäffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.

**Reference:**  
Schäffer, Alejandro A., L. Aravind, Thomas L. Madden, Sergei Shavirin, John L. Spouge, Yuri I. Wolf, Eugene V. Koonin, and Stephen F. Altschul (2001), "Improving the accuracy of PSI-BLAST protein database searches with composition-based statistics and other refinements", Nucleic Acids Res. 29:2994-3005.

RID: JE920YVS012

**Database:** Non-redundant SwissProt sequences  
259,844 sequences; 98,643,457 total letters

If you have any problems or questions with the results of this search please refer to the [BLAST FAQs](#)  
[Taxonomy reports](#)

**Query=**  
Length=216

[Distribution of 9 Blast Hits on the Query Sequence](#)

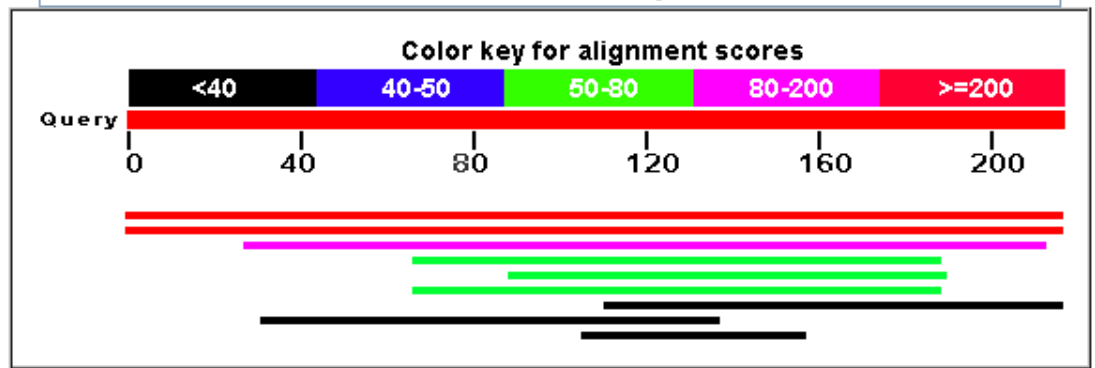
Mouse-over to show define and scores, click to show alignments

**Color key for alignment scores**



### Distribution of 9 Blast Hits on the Query Sequence

Mouse-over to show define and scores, click to show alignments



[Distance tree of results](#) **NEW**

Sequences producing significant alignments:

Sequence ID	Description	Score (Bits)	E Value
<a href="#">sp P05050.1 ALKB_ECOLI</a>	Alpha-ketoglutarate-dependent dioxygen...	450	2e-126
<a href="#">sp P37462 ALKB_SALTY</a>	Alpha-ketoglutarate-dependent dioxygenas...	363	3e-100
<a href="#">sp 005725 ALKB_CAUCR</a>	Alpha-ketoglutarate-dependent dioxygenase a	139	6e-33
<a href="#">sp Q9SA98 ALKBH_ARATH</a>	Alkylated DNA repair protein alkB homolog	75.5	2e-13
<a href="#">sp Q13686 ALKB1_HUMAN</a>	Alkylated DNA repair protein alkB homolog	69.3	1e-11
<a href="#">sp 060066 ALKBH_SCHPO</a>	Alkylated DNA repair protein alkB homolog	68.2	2e-11
<a href="#">sp Q5UR03 YL905_MIMIV</a>	Uncharacterized protein L905	36.2	0.085
<a href="#">sp Q5UR02 CHLE_MIMIV</a>	Probable cholinesterase precursor (Acylchol	30.8	4.0
<a href="#">sp Q8N661 TM86B_HUMAN</a>	Transmembrane protein 86B	30.0	7.6

### Alignments

- Get selected sequences
- Select all
- Deselect all
- Distance tree of results

> [sp|P05050.1|ALKB\\_ECOLI](#) Alpha-ketoglutarate-dependent dioxygenase alkB (Alkylated DNA repair protein alkB)  
 Length=216

> [sp|P05050.1|ALKB\\_ECOLI](#) Alpha-ketoglutarate-dependent dioxygenase alkB (Alkylated DNA repair protein alkB)  
Length=216

Score = 450 bits (1157), Expect = 2e-126, Method: Composition-based stats.  
Identities = 216/216 (100%), Positives = 216/216 (100%), Gaps = 0/216 (0%)

Query	1	MLDLFADAEPWQEPLAAGAVILRRFAFNAAEQLIRDINDVASQSPFRQMVTTPGGYTMSVA	60
		MLDLFADAEPWQEPLAAGAVILRRFAFNAAEQLIRDINDVASQSPFRQMVTTPGGYTMSVA	
Sbjct	1	MLDLFADAEPWQEPLAAGAVILRRFAFNAAEQLIRDINDVASQSPFRQMVTTPGGYTMSVA	60
Query	61	MTNCGHLGWTTTHRQGYLYSPIDPQTNKPWPAMPQSFHNLQRAATAAGYPDFQPDACLIN	120
		MTNCGHLGWTTTHRQGYLYSPIDPQTNKPWPAMPQSFHNLQRAATAAGYPDFQPDACLIN	
Sbjct	61	MTNCGHLGWTTTHRQGYLYSPIDPQTNKPWPAMPQSFHNLQRAATAAGYPDFQPDACLIN	120
Query	121	RYAPGAKLSLHQDKDEPDLRAPIVSVSLGLPAIFQFGGLKRNPLKRLLEHGDVVVWGG	180
		RYAPGAKLSLHQDKDEPDLRAPIVSVSLGLPAIFQFGGLKRNPLKRLLEHGDVVVWGG	
Sbjct	121	RYAPGAKLSLHQDKDEPDLRAPIVSVSLGLPAIFQFGGLKRNPLKRLLEHGDVVVWGG	180
Query	181	ESRLFYHGIQPLKAGFHPLTIDCRYNLTFRQAGKKE	216
		ESRLFYHGIQPLKAGFHPLTIDCRYNLTFRQAGKKE	
Sbjct	181	ESRLFYHGIQPLKAGFHPLTIDCRYNLTFRQAGKKE	216

> [sp|P37462|ALKB\\_SALTY](#) Alpha-ketoglutarate-dependent dioxygenase alkB (Alkylated DNA repair protein alkB)  
Length=216

Score = 363 bits (932), Expect = 3e-100, Method: Composition-based stats.  
Identities = 172/216 (79%), Positives = 193/216 (89%), Gaps = 0/216 (0%)

Query	1	MLDLFADAEPWQEPLAAGAVILRRFAFNAAEQLIRDINDVASQSPFRQMVTTPGGYTMSVA	60
		MLDLFAD PWQEPLA GAV+LRRFAF AA+ L+ DI VASQSPFRQMVTTPGGYTMSVA	
Sbjct	1	MLDLFADEAPWQEPLAPGAVVLRRAFAFRAAQSLLDDIGFVASQSPFRQMVTTPGGYTMSVA	60
Query	61	MTNCGHLGWTTTHRQGYLYSPIDPQTNKPWPAMPQSFHNLQRAATAAGYPDFQPDACLIN	120
		MTNCG LGWTT R GY Y+ DP T+KPWPA+P SF ++C++AA AAGY FQPDACLIN	
Sbjct	61	MTNCGALGWTTDRHGVCYAVRDPLTDKWPALPLSFASVCRQAAIAAGYASFQPDACLIN	120
Query	121	RYAPGAKLSLHQDKDEPDLRAPIVSVSLGLPAIFQFGGLKRNPLKRLLEHGDVVVWGG	180
		RYAPGAKLSLHQDKDEPDLRAPIVSVSLG+PA+FQFGGL+R+DP++R+LLEHGD+VVVWGG	
Sbjct	121	RYAPGAKLSLHQDKDEPDLRAPIVSVSLGVPVAVFQFGGLRRSDPIQRILLEHGDIVVWGG	180
Query	181	ESRLFYHGIQPLKAGFHPLTIDCRYNLTFRQAGKKE	216
		ESRLFYHGIQPLKAGFHP+T + RYNLTFRQA +KE	
Sbjct	181	ESRLFYHGIQPLKAGFHPMTGEFRYNLTFRQAAEKE	216

> [sp|005725|ALKB\\_CAUCR](#) Alpha-ketoglutarate-dependent dioxygenase alkB homolog  
Length=220

Score = 139 bits (351), Expect = 6e-33, Method: Composition-based stats.

> [sp|Q13686|ALKB1 HUMAN](#) **G** Alkylated DNA repair protein alkB homolog 1  
Length=389

Score = 69.3 bits (168), Expect = 1e-11, Method: Composition-based stats.  
Identities = 33/101 (32%), Positives = 55/101 (54%), Gaps = 0/101 (0%)

Query 89 WPAMPQSFHNLQRAATAAGYPDFQPDACLINRYAPGAKLSLHQDKDEPDLRAPIVSVSL 148  
+ P L ++ A A G+ DF+ +A ++N Y + L +H D+ E D P++S S  
Sbjct 189 YTPFPSDLGFLSEQVAAAACGFEDFRAEAGILNYYRLDSTLGIHVDRSELDHSPKLLSFSF 248

Query 149 GLPAIFQFGGLKRNPLKRLLEHGDVVVWGGESRLFYHGI 189  
G AIF GGL+R++ + + GD+++ G SRL H +  
Sbjct 249 GQSAIFLLGGLQRDEAPTAMFMHSGDIMIMSGFSRLLNHAV 289

> [sp|O60066|ALKBH SCHPO](#) Alkylated DNA repair protein alkB homolog  
Length=273

Score = 68.2 bits (165), Expect = 2e-11, Method: Composition-based stats.  
Identities = 37/123 (30%), Positives = 63/123 (51%), Gaps = 1/123 (0%)

Query 67 LGWITHRQGYLYSPIDPQTNKPWPAMPQSFHNLQRAAT-AAGYPDFQPDACLINRYAPG 125  
L W T + Y ++ + P P+ + ++ + + ++ +A ++N Y+PG  
Sbjct 135 LRWVTLGEQYDWTITKEYPDPSPKSPGPFKDLGDFVEKVVKESTDFLHWKAEAAIVNFYSPG 194

Query 126 AKLSLHQDKDEPDLRAPIVSVSLGLPAIFQFGGLKRNPLKRLLEHGDVVVWGGESRLF 185  
L S H D+ E DL P++S+S+GL I+ G R++ L L GDVV+ G SR  
Sbjct 195 DTLSAHIDESEEDLTLPLISLSMGLDCIYLIGTESRSEKPSALRLHSGDVVIMTGTSRKA 254

Query 186 YHG 188  
+HG  
Sbjct 255 FHG 257

> [sp|Q5URO3|YL905 MIMIV](#) Uncharacterized protein L905  
Length=210

Score = 36.2 bits (82), Expect = 0.085, Method: Composition-based stats.  
Identities = 30/113 (26%), Positives = 55/113 (48%), Gaps = 11/113 (9%)

Query 111 DFQPDACLINRYAPGAKLSLHQDKDEPDLRAPIVSVSLGLPAIFQFGGLKRNPL---KR 167  
D +PD ++N Y PG L H D+ + + I+ +SLG I +F +N P+ K+  
Sbjct 88 DQKPDQIIVNEYKPGEGLPKPHFDRKDY-QQNVIIIGLSLGSSTIMEF---YKNKPIPEKKK 143

Query 168 LLEHGDVVVWGGESRLFY-HGIQPLK---AGFHPLTIDCRYNLTFRQAGKKE 216  
+ + + + ++R + HCI P K + + R ++TFR K++  
Sbjct 144 IYIPRSLYIIKDDARYIWKHGIPPRKYDEINGKKIPRETRISITFRNVIKEK 196

> [sp|Q5URO2|CHLE MIMIV](#) **G** Probable cholinesterase precursor (Acylcholine acylhydrolase)  
Length=579

NCBI Entrez Protein  
My NCBI  
[Sign In] [Register]

Search Protein for [ ] Go Clear  
Limits Preview/Index History Clipboard Details  
Display GenPept Show 5 Send to [ ]  
Range: from begin to end Features: [x] CDD + Refresh

I: [Q13686](#). Reports Alkylated DNA rep...[gi:12643239] BLink, Conserved Domains, Links

[Comment](#) [Features](#) [Sequence](#)

LOCUS Q13686 389 aa linear PRI 10-JUL-2007  
DEFINITION Alkylated DNA repair protein alkB homolog 1.  
ACCESSION Q13686  
VERSION Q13686.2 GI:12643239  
DBSOURCE swissprot: locus ALKB1\_HUMAN, accession [Q13686.2](#);  
class: standard.  
extra accessions: Q8TAU1, Q9ULA7  
created: Dec 1, 2000.  
sequence updated: Dec 1, 2000.  
annotation updated: Jul 10, 2007.  
xrefs: [X91992.1](#), [CAA63047.1](#), [AC008044.4](#), [AAF01478.1](#), [BC025787.1](#),  
[AAH25787.1](#), [S64736](#)  
xrefs (non-sequence databases): UniGene:Hs.94542,  
Ensembl:ENSG00000100601, KEGG:hsa:8846, H-InvDB:HIX0011855,  
HGNC:17911, MIM: [605345](#), ArrayExpress:Q13686,  
GermOnline:ENSG00000100601, RZPD-ProtExp:IOH12000,  
RZPD-ProtExp:T4267, GO:0006307, InterPro:IPR004574,  
TIGRFAMs:TIGR00568  
KEYWORDS .  
SOURCE Homo sapiens (human)  
ORGANISM [Homo sapiens](#)  
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;

KEYWORDS .

SOURCE Homo sapiens (human)

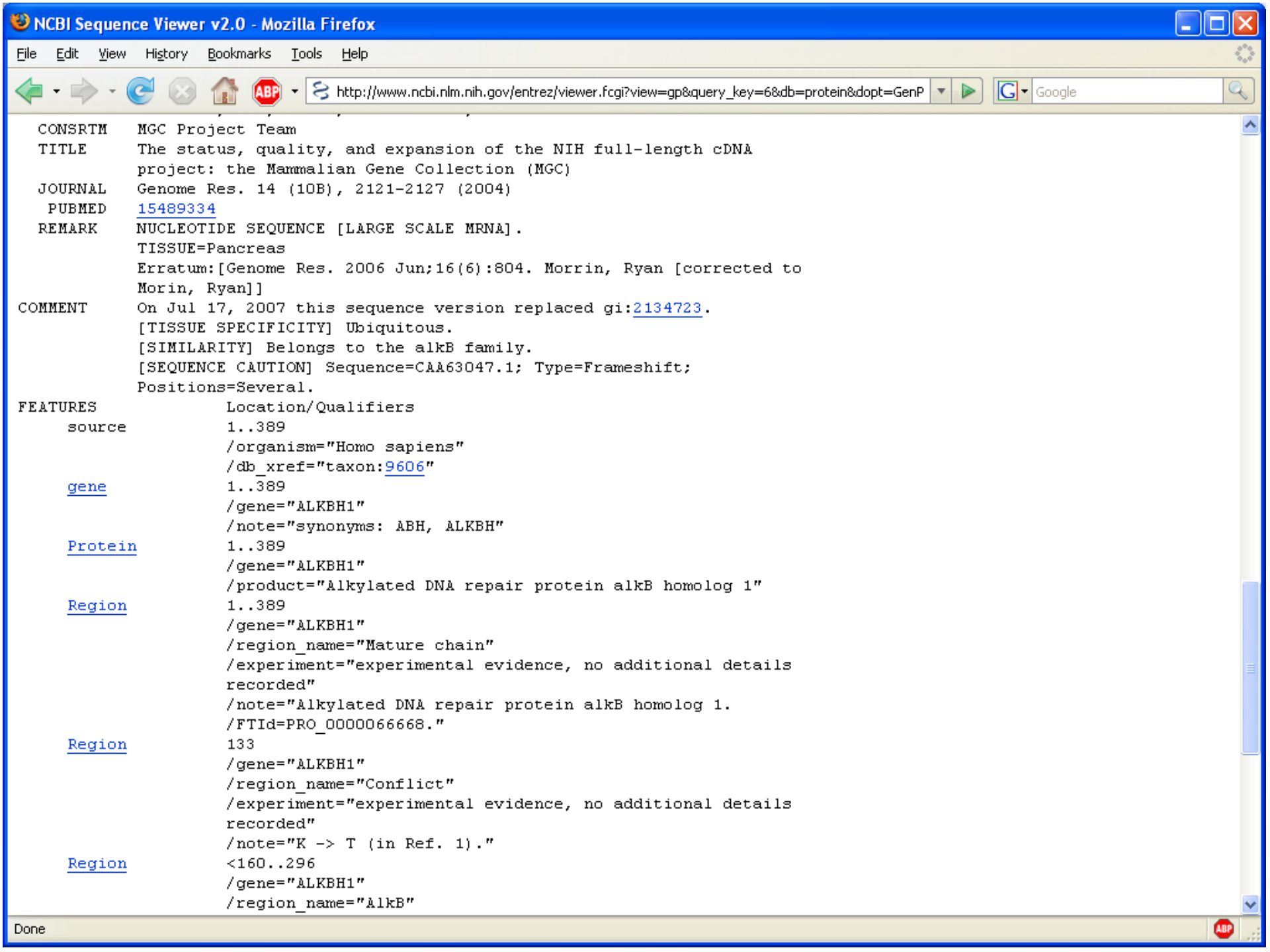
ORGANISM Homo sapiens
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
Catarrhini; Hominidae; Homo.

REFERENCE 1 (residues 1 to 389)
AUTHORS Wei, Y.F., Carter, K.C., Wang, R.P. and Shell, B.K.
TITLE Molecular cloning and functional analysis of a human cDNA encoding
an Escherichia coli AlkB homolog, a protein involved in DNA
alkylation damage repair

JOURNAL Nucleic Acids Res. 24 (5), 931-937 (1996)
PUBMED 8600462
REMARK NUCLEOTIDE SEQUENCE [MRNA].
TISSUE=Synovial sarcoma

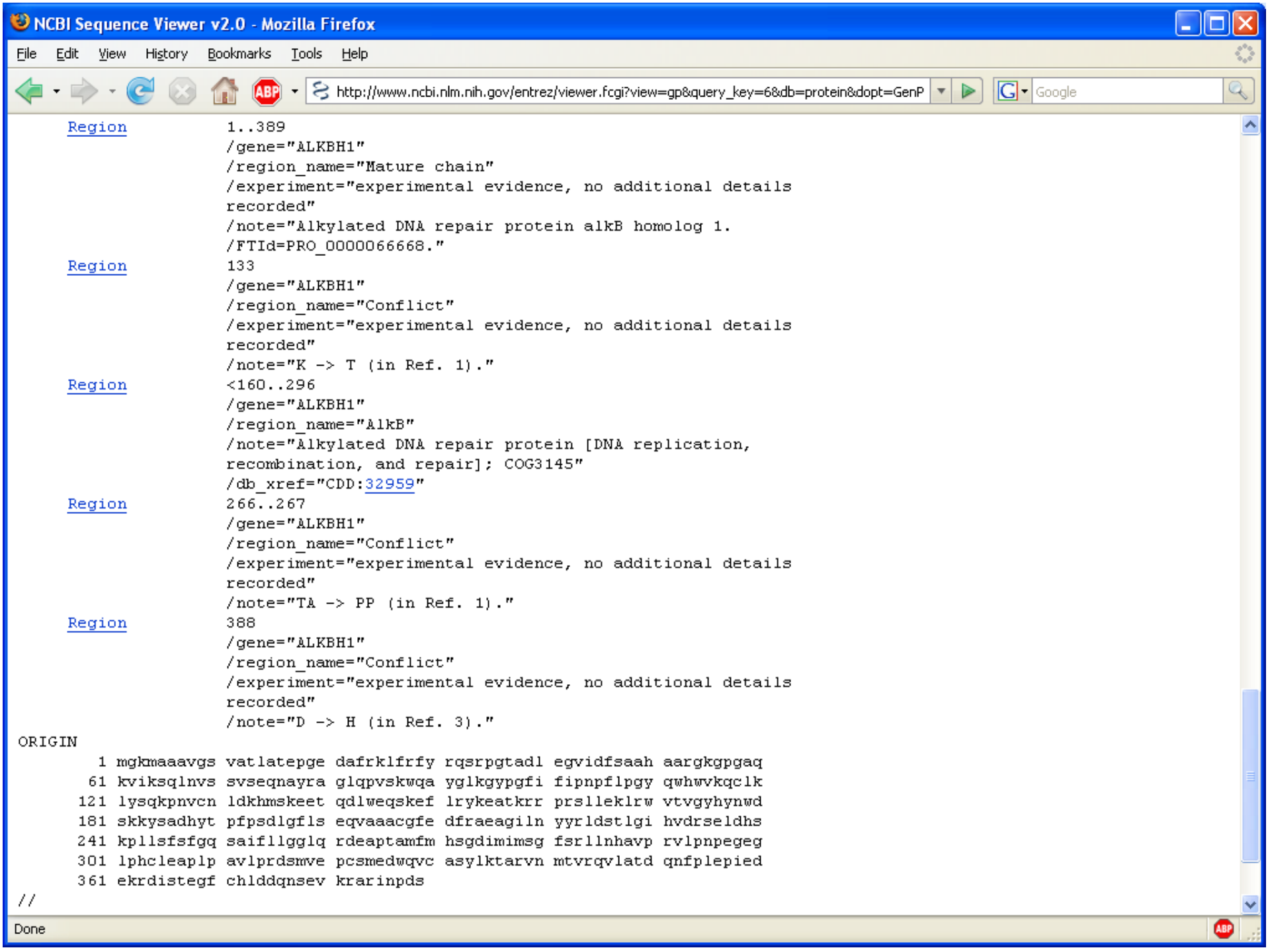
REFERENCE 2 (residues 1 to 389)
AUTHORS Heilig, R., Eckenberg, R., Petit, J.L., Fonknechten, N., Da Silva, C.,
Cattolico, L., Levy, M., Barbe, V., de Berardinis, V., Ureta-Vidal, A.,
Pelletier, E., Vico, V., Anthouard, V., Rowen, L., Madan, A., Qin, S.,
Sun, H., Du, H., Pepin, K., Artiguenave, F., Robert, C., Cruaud, C.,
Bruls, T., Jaillon, O., Friedlander, L., Samson, G., Brottier, P.,
Cure, S., Segurens, B., Aniere, F., Samain, S., Crespeau, H., Abbasi, N.,
Aiach, N., Boscus, D., Dickhoff, R., Dors, M., Dubois, I., Friedman, C.,
Gouyvenoux, M., James, R., Madan, A., Mairey-Estrada, B., Mangenot, S.,
Martins, N., Menard, M., Oztas, S., Ratcliffe, A., Shaffer, T.,
Trask, B., Vacherie, B., Bellemere, C., Belser, C., Besnard-Gonnet, M.,
Bartol-Mavel, D., Boutard, M., Briez-Silla, S., Combette, S.,
Dufosse-Laurent, V., Ferron, C., Lechaplais, C., Louesse, C.,
Muselet, D., Magdelenat, G., Pateau, E., Petit, E.,
Sirvain-Trukniewicz, P., Trybou, A., Vega-Czarny, N., Bataille, E.,
Bluet, E., Bordelais, I., Dubois, M., Dumont, C., Guerin, T.,
Haffray, S., Hammadi, R., Muanga, J., Pellouin, V., Robert, D.,
Wunderle, E., Gauguier, G., Roy, A., Sainte-Marthe, L., Verdier, J.,
Verdier-Discala, C., Hillier, L., Fulton, L., McPherson, J.,
Matsuda, F., Wilson, R., Scarpelli, C., Gyapay, G., Wincker, P.,
Saurin, W., Quetier, F., Waterston, R., Hood, L. and Weissenbach, J.

TITLE The DNA sequence and analysis of human chromosome 14
JOURNAL Nature 421 (6923), 601-607 (2003)
PUBMED 12508121
REMARK NUCLEOTIDE SEQUENCE [LARGE SCALE GENOMIC DNA].



CONSRTM MGC Project Team  
TITLE The status, quality, and expansion of the NIH full-length cDNA project: the Mammalian Gene Collection (MGC)  
JOURNAL Genome Res. 14 (10B), 2121-2127 (2004)  
PUBMED [15489334](#)  
REMARK NUCLEOTIDE SEQUENCE [LARGE SCALE MRNA].  
TISSUE=Pancreas  
Erratum:[Genome Res. 2006 Jun;16(6):804. Morrin, Ryan [corrected to Morin, Ryan]]  
COMMENT On Jul 17, 2007 this sequence version replaced gi:[2134723](#).  
[TISSUE SPECIFICITY] Ubiquitous.  
[SIMILARITY] Belongs to the alkB family.  
[SEQUENCE CAUTION] Sequence=CAA63047.1; Type=Frameshift; Positions=Several.

FEATURES Location/Qualifiers  
source 1..389  
/organism="Homo sapiens"  
/db\_xref="taxon:[9606](#)"  
[gene](#) 1..389  
/gene="ALKBH1"  
/note="synonyms: ABH, ALKBH"  
[Protein](#) 1..389  
/gene="ALKBH1"  
/product="Alkylated DNA repair protein alkB homolog 1"  
[Region](#) 1..389  
/gene="ALKBH1"  
/region\_name="Mature chain"  
/experiment="experimental evidence, no additional details recorded"  
/note="Alkylated DNA repair protein alkB homolog 1."  
/FTId=PRO\_0000066668."  
[Region](#) 133  
/gene="ALKBH1"  
/region\_name="Conflict"  
/experiment="experimental evidence, no additional details recorded"  
/note="K -> T (in Ref. 1)."  
[Region](#) <160..296  
/gene="ALKBH1"  
/region\_name="alkB"



Region 1..389  
 /gene="ALKBH1"  
 /region\_name="Mature chain"  
 /experiment="experimental evidence, no additional details recorded"  
 /note="Alkylated DNA repair protein alkB homolog 1.  
 /FTId=PRO\_0000066668."

Region 133  
 /gene="ALKBH1"  
 /region\_name="Conflict"  
 /experiment="experimental evidence, no additional details recorded"  
 /note="K -> T (in Ref. 1)."

Region <160..296  
 /gene="ALKBH1"  
 /region\_name="AlkB"  
 /note="Alkylated DNA repair protein [DNA replication, recombination, and repair]; COG3145"  
 /db\_xref="CDD:32959"

Region 266..267  
 /gene="ALKBH1"  
 /region\_name="Conflict"  
 /experiment="experimental evidence, no additional details recorded"  
 /note="TA -> PP (in Ref. 1)."

Region 388  
 /gene="ALKBH1"  
 /region\_name="Conflict"  
 /experiment="experimental evidence, no additional details recorded"  
 /note="D -> H (in Ref. 3)."

ORIGIN

```

1 mgkmaaaavgs vatlatepge dafrklirfy rqsrrpqtadl egvidfsaah aargkpggaq
61 kviksqliavs svseqnayra glqpvskwqa yglkgypgfi fipnpflpgy qwhwvkqclk
121 lysqkpnvcn ldkhmskeet qdlweqskef lrykeatkrr prsllleklrw vtvgyhynwd
181 skkysadhyt pfpstdlgfls eqvaaacgfe dfraeagiln yyrldstlgi hvdrseidhs
241 kpillsfsfgq saifllgglq rdeaptamfm hsgdiminsg fsrllnhavp rvlpnpegeg
301 lphcleaplp avlprdsmvv pcsmedwqvc asylktarvn mtvrqvlvtd qnfplepied
361 ekrdistegf chlddqnvsef krarinpds

```

//

# Common alignment scoring system

- Substitution score matrix
  - Score for aligning any two residues to each other
  - Identical residues have large positive scores
  - Similar residues have small positive scores
  - Very different residues have large negative scores
- Gap penalties
  - Penalty for opening a gap in a sequence (Q)
  - Penalty for extending a gap (R)
  - Typical gap function:  $G = Q + R * L$ , where L is length of gap
  - Example: Q=11, R=1

**BLOSUM62 amino acid substitution score matrix**

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

```

E.c. Alka 127 SVAMAAKLTARVAQLYGERLDDFPE--YICFPTPQRLAAADPQA-LKALGMPLKRAEALI 183
      ++|   +  |+ | +| ||   +  |  ||+ | ||  + +| |+ ||+  ||  +
H.s. OGG1 151 NIARITGMVERLCQAFGPRLIQLDDVITYHGFPSLQALAGPEVEAHLRKLGLGY-RARYVS 209

E.c. Alka 184 HLANAALE-----GTLPM TIPGDVEQAMKTLQTFPGIGRWTANYFAL 225
      | | ||   |   | +| | |  ||+|  |+  |
H.s. OGG1 210 ASARAILEEQGLAWLQQLRESSYEEAHKALCILPGVGTKVADCICL 256
    
```

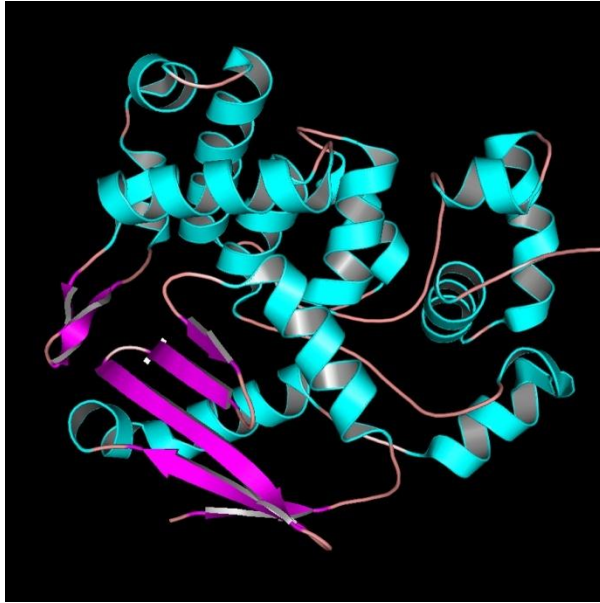


# Amino acid substitution score matrix

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

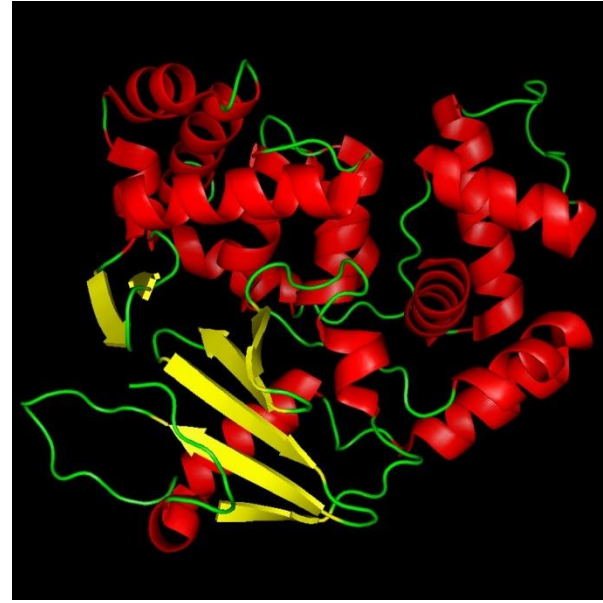
BLOSUM62

# Structure and sequence alignment



*E. coli* AlkA

Hollis *et al.* (2000) *EMBO J.* **19**, 758-766 (PDB ID 1DIZ)



Human OGG1

Source: Bruner *et al.* (2000) *Nature* **403**, 859-866 (PDB ID 1EBM)

E. c.	AlkA	127	SVAMAAKLTARVAQLYGERLDDFPE--YICFPTPQRLAAADPQA-LKALGMPLKRAEALI	183
			++  +   +     +     +     + +     +     +	
H. s.	OGG1	151	NIARITGMVERLCQAFGPRLIQLDDVTYHGFPSLQALAGPEVEAHLRKLGLGY-RARYVS	209
E. c.	AlkA	184	HLANAALE-----GTLPMTIPGDVEQAMKTLQTFPGIGRWTANYFAL	225
			+         +     +	
H. s.	OGG1	210	ASARAILEEQGGLAWLQQLRESSYEEAHKALCILPGVGTKVADCICL	256

# Significance of alignments

- Even random sequences may reach a high score when aligned optimally, so when is a sequence alignment significant?
- Statistical methods compare the score of a match with the distribution of alignment scores found using random sequences
- The most commonly used indicator of significance:  
E-value = Expect value = expected number of random matches at least as good as this one (with at least this alignment score)
- Some other indicators of significance (less accurate):
  - Percentage of identical residues
  - Percentage of similar residues
  - Raw alignment score

# Repeats and low complexity regions

- Repeats and low complexity regions constitute more than one third of the human genome.
- Highly locally biased composition occurs in regions of many proteins and in DNA. E.g. structural proteins in hair.
- Low complexity regions may give rise to high alignment scores – but are usually biologically uninteresting
- They can (and should usually) be masked using programs like RepeatMasker, DUST or SEG before a database search is carried out. The sequence in each region is then replaced by Ns or Xs.
- Examples:
  - interspersed repeats:
    - Short interspersed elements (SINEs)
    - Long interspersed elements (LINEs)
  - simple repeats (microsatellites)
    - usually 1 to 7 nucleotides are repeated a large number of times
    - E.g. ...AGAGAGAGAGAGAGAGAG...
    - E.g. ...CCGCCGCCGCCGCCGCCGCCG...
  - low complexity regions,
    - Protein example: PPCDPPPPPKDKKKKDDGPP
    - DNA example: AAATAAAAAAATAAAAAAT

# BLAST online resources

- NCBI BLAST website  
<http://www.ncbi.nlm.nih.gov/BLAST/>
- NCBI tutorial on BLAST  
<http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/information3.html>
- NCBI Handbook, Chapter 16, BLAST  
<http://www.ncbi.nlm.nih.gov/books/bv.fcgi?rid=handbook.chapter.ch16>
- BLAST FAQ  
[http://www.ncbi.nlm.nih.gov/blast/blast\\_FAQs.shtml](http://www.ncbi.nlm.nih.gov/blast/blast_FAQs.shtml)
- Wikipedia on BLAST  
<http://en.wikipedia.org/wiki/BLAST>

# **Finding families of proteins**

# Overview

- Searching with a family of proteins
  - Back to the example
  - PSI-BLAST example
  - Sequence patterns
  - Sequence profiles and logos
  - Iterated searches and PSI-BLAST
  - Advanced PSI-BLASTing

# Back to the example...

How are all these sequences found? Ordinary BLAST is not enough...

CAS_scla_322266	-----YHRLQPNYVMLACSRADHE-----	RTAATLVASVRK--70--	VTEAVYLEEG-DLLIVDNF-----	RTTARTPFSPRWGDKRLLHRVVI	302 \	
IPNS_Bn_124825	-----LITVLYIQ-----	SNVQNLQVETAA--	GYQDIBADDT-GYLINCGSYMAHLTNNYKAPIDRVKVVN	-----AE	RQSLPPFVNL 288	
FLAS_Pet_421946	-----YITLILVP-----	NEVQGLQVFKDG--	HWYDVKIYFN-ALIVHIGDQVEILSNGKYSVYRITTVNK	---	DKTRMSWPVLEP 309	
LDOX_Pet_1730108	-----ALTFILH-----	NMVFGLQLFYEG--	QWVTAKCVFN-SIIMHIGDTIILSNGKYSILRQGVN	---	BKVRFSWAIFCEP 311	
Srg_At_479047	-----GLITVLMQV-----	NDVBSGLQIKKDG--	KWVPVKLPN-AFIVNIGDVLEIITNGTYRSIEIRGVVNS	---	EKERLSIATFHNV 309	
EFE_Le_398992	-----GIIILLFQD-----	DKVBSGLQLLKDE--	QWIDVPPMRH-SIVVNLGDLQLEVITNGKYSVLRVIAQT	---	DGTRMSLASFYNP 253	
Ga200x_Sot_10800976	-----SITLILHQ-----	DSVBSGLQVFMND--	QWRISPNLS-AFVNVIGDTFMALSNGRYKSCLRAVVNN	---	KTPRKLAFPLCP 317	
PA0147_Pa_9945977	-----CVTLLLYQ-----	DARGGLQVQNRQG--	EWIDAPPIDG-TFVNVIGDMMARWSNDRYRSTPFRVISP	---	GVHRYSMPPFAEP 274	
PA4191_Pa_9950401	-----LLTLLHQ-----	DAIGGLQVTRPQ--	GWLEAPPIDG-SFVNLGDLBRMTGGLYRSTPFRVARTS	---	GRDRLSFPLFFDP 277	
ISP7_sp_729862	-----ALTLMSQ-----	DNVKGLEILLDPVSN--	CFLSVSPAPG-ALIANLGDIMAILTNNRYKSSMERVCNNS	---	GSDRTYIPFFLQG 353	
SPCC1494.01_sc_7491815	-----SITLLFQ-----	RDAAGLEIRPPNFVKDM--	DWIKVNVQSD-VVLVNIADMLQFWTSGKLRSTVFRVIDPG	---	VKTRQTIAYVTP 267	
DACCS_Lyl_769809	-----IVSLILQTPCP-----	NGFVSLQVBI DG--	RFVEVPPRPG-CVVVFCGSIAPLVSDGKIKAPQFRVNS	PGA4--	GSNRTSSVLELRP 268/	
RRPO_SHVX_548840	IYPKG-----	NKILTVNAA-----	GSSTFGI-----	KCAKGE-TTLNLEDGD-YFQMPSGFQETHKRVVA	----	VTPRSLSTFRSTV 743 \
POL_ASPV_487652	IYDIN-----	HQVLTVNS-----	GDAIFCI-----	ECLGSGF-EIPLSGPQ-MLLMPFGQKSHRGIKSP	----	SKGRISLTFRLTA 853
POL_BSV_409711	IFMRG-----	APVHTVSHD-----	GNADFGT-----	ECAGR--QYTLRGNVQFTMPSGQETHKRAVNT	----	TAGRVSYTFRLLA 841
RRPO_FMV_139137	CYPKG-----	HQVLTINHS-----	GBOLTQI-----	ACQKGA-SITMGEGD-YLSPVGFQESHKRAVNT	----	TGGRVSLTFRECTV 690
POL_GLV_1154656	IFEKD-----	SKILTVCIQ-----	GCDFRE-----	RCATGET-GFYMEAPK-QFMMPDGFQSNHVAVREB	----	TPGRI SATFRAX 772
Pol_GVA_1405615	CYLPG-----	GSVVTVNLH-----	GDATFEVK-----	ENQSGKIEKELHGD-VYVMGPMQQTHKRVYTS	----	TDGRCSITLANKT 738
RRPO_ACLSV_1710717	CYDD-----	DEILTINVV-----	GCAXFHT-----	TC-HGE--IIDLRQGD-EILMPGGYQKMNKRAVEVA	----	SEGRTSVTLRVHK 836/
TL3L16.2_At_2708738	FL-----	RPCTISFL-----	SECDILFGSNLKVE-----	GPDFGSGY-SIPLPVGSLVFLNGGADVAKCVPVAV	----	PTKRSITFRKMD 420 \
T19K4.220_At_3036813	FL-----	RPCTVSFL-----	SECNILFGSNLKVL-----	GPDFGSGY-SIPLPVGSLVFLNGGADVAKCVPVAV	----	PTKRSITFRKMD 403 \
At2g48080_At_4249414	-----	QPISTLVL-----	SESTMVFGHRLGVD-----	NDGNFRGSL-TLPLKEGS-LLVMRGNADMARVWCPS	----	PNKRVAITFFILK 351
AK000315.1_Hs_7020317	IFE-----	RPIVSVSFF-----	SDSALCFGCKQFK-----	PIRVSEVLSLPRRGS-VTVLSGYAADBITCIRPQDI	----	KERAVIILAKTR 270
CG17807_Dm_7291441	AFL-----	DPILSLSLQ-----	SDVMDFRG-----	DDQV-QVLRPRS-LLLMSGARYDWTGIRPKHID13RKR	-----	TKRSLTFRRLR 325
CG6144_Dm_7297712	FH-----	PIISTISTG-----	AHTVLEFRKREDTTETREAGDQTTRELVF	-----	LKLEPRS-LLLKDLYTDYLAISBTSBD24RSPRISLTI	RNVV 213
CG4036_Dm_7297561	IWGBRVVTVNC-----	LGDSVLTLT--PYEVQSGKYNLDLVA	SYDELLAP-LLTDDQLATFEGKVLRI	PMPNLS-LIVLYGPARYQFESVLR	-----	QERVCVAYREFT 278
FLJ2001_Hs_38923019	LWGBRLVSLNL-----	LSPTVLSMC-----	REAPGSLLLCPSA	RAPEALVDSVIAPSRVLCQVEVAIPLPARS-LLVLTGAARHQWKA	-----	AIHRRHI--BARVCVTFRELS 274
C14B1.10_Ce_6580210	AFD-----	DPIVSISSL-----	SDVMEFKD-----	GANSARIAPVLLKARS-LCLIQESRYRWEK	-----	GIVNRYKD10RQTRVSLTKIR 343
SPAP8A3.02c_sp_7491301	FGDG-----	VAIFSLSN-----	TIMFTPE-----	LKLS--KIRLEKGS-LLLMSGTARYDWTGIRPKHID12RQ	-----	RSLVTMERRI 219
L3377.4_Lm_9989036	YVD-----	DIPAI CSLG-----	SNCLLRVH-----	VQNGSEL-DVMVPRDS-VYIMSGPARYVYFMYL	-----	EAQRSLVFRRSI 193
MTCI237.14c_Mtu_2052134	RGSTEDTM-----	VAIVSLGAT-----	RVFALRP-----	RGRGSLRPLAAGD-LLVMGSGCQR	-----	TFEAVPKTSAP--TGRVSIQFRPD 293 \
AlkB_Cc_2055386	ADPR-----	FPLLSISLG-----	DTAVFRIGG-----	VNRKDPTRSLRLASGD--VCRLLGARLAF	-----	STDRILPG6-GGGRNILTLRAR 190
AlkB_Bc_113638	PDLR-----	APIVSVSLG-----	LFAIFQGG-----	LKRNDPLKRLLEHGD--VVVWGSESLFY	-----	GIQPLKAG5-IDCRNLTFRQAG 213
AlkB_Scoe_8894829	RTD-----	APVVSLSG-----	DTCVFRFG-----	PETRTFRPYTDELRSGD--LFVFGPSRLAY	-----	GVPRVHPG7-LRGRNLTLRVSG 215
AlkB_At_4835778	ADWS-----	KPIVSMISG-----	CKAIFLLGK-----	SKDDPPHMYLRSGD--VVLMAGBARECF	-----	GNLLHFQL34KTSRININIQVF 354
alkb_sp_3080529	BDLT-----	LPLISLSMG-----	LDGIYLGTE-----	SRSEKPS-ALRLHSGD--VVIMTGTSRKAF	-----	SKHC--SFKLYLSLIA 272
AlkB_Hs_2134723	LDHS-----	KPLLSFSG-----	QSAIFLLGGL-----	QRDEAPP-EMFMSGD--IMIMSGSRLLN	-----	GVPRVLPN39KTA
Consensus (85%):	.....sh.h.....	.....s.....h.....	.....s.....h.....	.....H.s.....	.....+h.h..b...	

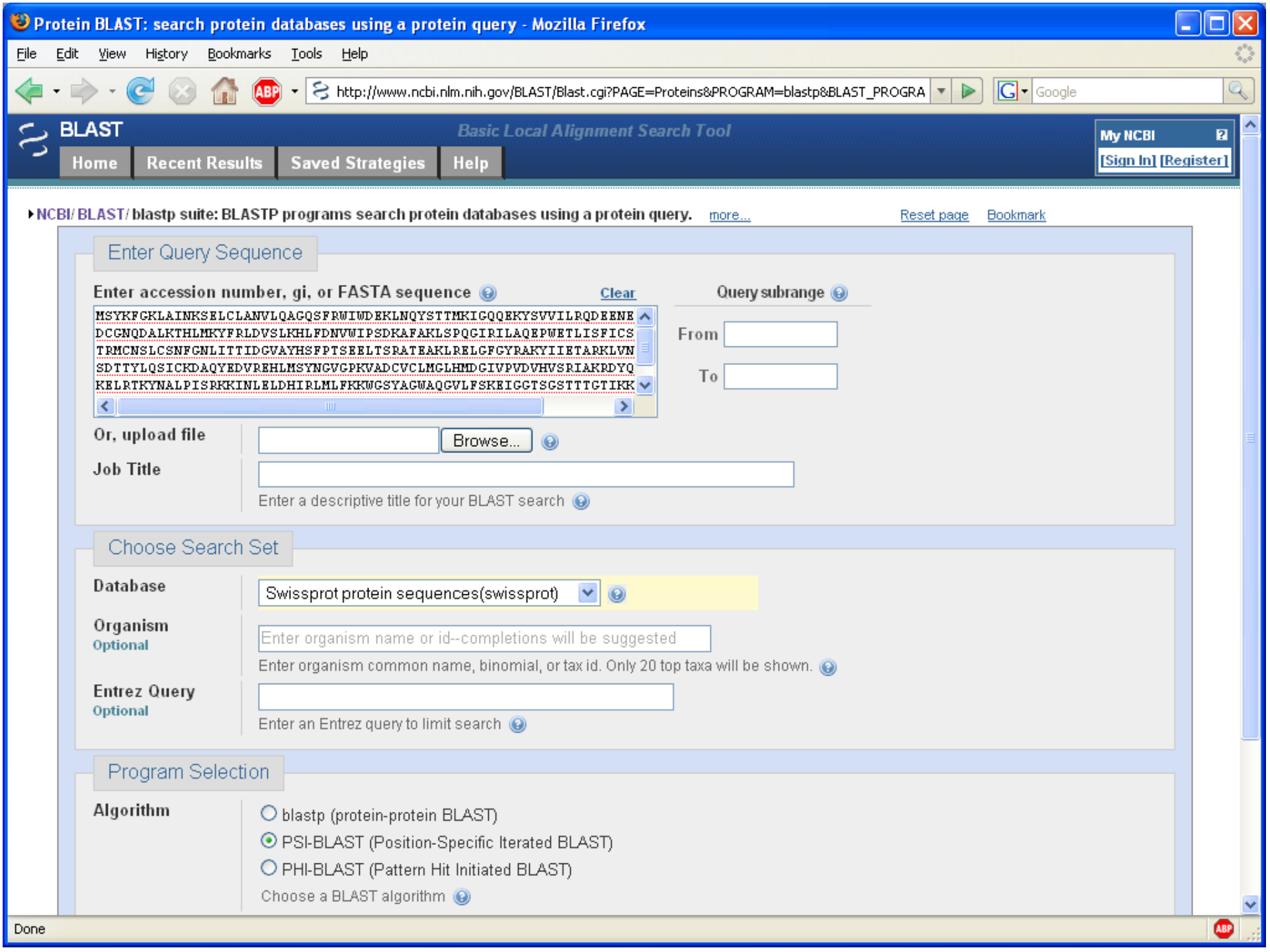


# Excerpt from the AlkB paper

## Results and discussion

### The 2OG-Fe(II) dioxygenase protein superfamily: classification and functional prediction

The Non-redundant Protein Sequence Database (NCBI) [21] was searched using the PSI-BLAST program [22] run to convergence, with a profile-inclusion threshold of 0.01 and AlkB protein sequences from various organisms as queries. In addition to the AlkB orthologs, these searches retrieved from the database, with statistically significant expectation (e) values, several other more distant homologs of AlkB, including uncharacterized eukaryotic proteins and fragments of the polyproteins of plant RNA viruses from the carla-, tricho- and potexvirus families. Examples of homologs found include: *Leishmania* L3377.4, iteration 5, e-value =  $8 \times 10^{-7}$ ; *Drosophila* CG17807, iteration 3, e-value =  $4 \times 10^{-6}$ ; papaya mosaic virus, iteration 3, e-value =  $2 \times 10^{-4}$ . Further iterations of the search using each of the detected proteins as a new query resulted in the detection of several more eukaryotic proteins, including EGL-9 and leprecan, several uncharacterized bacterial proteins and prolyl and lysyl hydroxylases. Finally, another iteration of database searches initiated with the sequences of bacterial proteins, typified by *E. coli* YbiX, resulted in the unification of these proteins with plant dioxygenases such as leucoanthocyanidin oxidase and gibberellin-20 oxidase. In this context, it should be noted that the DNA sequence encoding YbiX is highly similar to that of the



# BLAST

## Basic Local Alignment Search Tool

My NCBI  
[Sign In](#) [Register](#)

- [Home](#)
- [Recent Results](#)
- [Saved Strategies](#)
- [Help](#)

NCBI/BLAST/blastp suite: BLASTP programs search protein databases using a protein query. [more...](#) [Reset page](#) [Bookmark](#)

### Enter Query Sequence

Enter accession number, gi, or FASTA sequence [Clear](#)

```
MSYKFGKLAINKSELCLANVLQAGQSFRIWDEKLNQYSTTMKIQQEKYSVVILRQDEENE
DCGNQDALRTHLMRYFRDVSLSKHLFDNVWIPSDKAFKLSPOGIRILAQEPWETLISFICS
TRHCNSLCSNFGNLIITIDGVAYHSFPTSEELTSRATFAKLRRELGFQYRAKYIETARKLVM
SDTTYLQSIKDAQYEDVREHLMSYNGVGPVADCVCLMGLHMDGIVPVDVHVSRIAKRDYQ
KELRTKYNALPISRFKINLELDHIRLMLFKRWGSYACWAQGVLFPSKEIGGTSGSTTTGTIRK
```

Query subrange [?](#)

From

To

Or, upload file

[Browse...](#) [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

### Choose Search Set

Database

Swissprot protein sequences(swissprot) [?](#)

Organism  
Optional

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. [?](#)

Entrez Query  
Optional

Enter an Entrez query to limit search [?](#)

### Program Selection

Algorithm

- blastp (protein-protein BLAST)
- PSI-BLAST (Position-Specific Iterated BLAST)
- PHI-BLAST (Pattern Hit Initiated BLAST)

Choose a BLAST algorithm [?](#)



Sequences with E-value BETTER than threshold

Sequences producing significant alignments:

			Score (Bits)	E Value		
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp P53397 OGG1 YEAST</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">741</a>	0.0	<b>G</b>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp 070249 OGG1 RAT</a>	N-glycosylase/DNA lyase [Includes: 8-oxogu...	<a href="#">194</a>	4e-49	<b>G</b>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp 008760 OGG1 MOUSE</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">194</a>	4e-49	<b>G</b>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp 015527 OGG1 HUMAN</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">192</a>	2e-48	<b>G</b>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q9V3I8 OGG1 DROME</a>	N-glycosylase/DNA lyase (dOgg1) [Include...	<a href="#">153</a>	1e-36	<b>G</b>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp 027397 OGG1 METHH</a>	Probable N-glycosylase/DNA lyase [Includ...	<a href="#">90.5</a>	1e-17	<b>G</b>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q9SJQ6 ROS1 ARATH</a>	Protein ROS1 (Repressor of silencing 1) ...	<a href="#">43.9</a>	0.001	<b>G</b>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp 049498 DML3 ARATH</a>	DEMETER-like protein 3	<a href="#">42.4</a>	0.004	<b>G</b>

Run PSI-Blast iteration 2

Sequences with E-value WORSE than threshold

<input type="checkbox"/>	<a href="#">sp Q9SR66 DML2 ARATH</a>	DEMETER-like protein 2	<a href="#">40.8</a>	0.009	
<input type="checkbox"/>	<a href="#">sp Q4UK93 END3 RICFE</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">35.8</a>	0.30	
<input type="checkbox"/>	<a href="#">sp Q8LK56 DME ARATH</a>	Transcriptional activator DEMETER (DNA gl...	<a href="#">35.4</a>	0.36	<b>G</b>
<input type="checkbox"/>	<a href="#">sp P18479 POLG ZYMVC</a>	Genome polyprotein [Contains: P1 protein...	<a href="#">35.4</a>	0.36	
<input type="checkbox"/>	<a href="#">sp P75506.1 Y269 MYCPN</a>	UPF0144 protein MG130 homolog	<a href="#">33.5</a>	1.3	
<input type="checkbox"/>	<a href="#">sp Q10630 ADA MYCTU</a>	Putative regulatory protein ada (Regulato...	<a href="#">33.5</a>	1.6	
<input type="checkbox"/>	<a href="#">sp Q58030 Y613 METJA</a>	Putative endonuclease MJ0613	<a href="#">33.1</a>	2.0	
<input type="checkbox"/>	<a href="#">sp Q63072 BST1 RAT</a>	ADP-ribosyl cyclase 2 precursor (Cyclic AD...	<a href="#">32.7</a>	2.9	<b>G</b>
<input type="checkbox"/>	<a href="#">sp Q64277 BST1 MOUSE</a>	ADP-ribosyl cyclase 2 precursor (Cyclic ...	<a href="#">32.3</a>	3.0	<b>G</b>
<input type="checkbox"/>	<a href="#">sp Q40002 VATA HORVU</a>	Vacuolar ATP synthase catalytic subunit ...	<a href="#">32.0</a>	4.0	
<input type="checkbox"/>	<a href="#">sp Q89330 POLG ZYMVR</a>	Genome polyprotein [Contains: P1 protein...	<a href="#">32.0</a>	4.8	
<input type="checkbox"/>	<a href="#">sp Q84934 POLG PPVSK</a>	Genome polyprotein [Contains: P1 protein...	<a href="#">31.6</a>	5.2	
<input type="checkbox"/>	<a href="#">sp P17767 POLG PPVRA</a>	Genome polyprotein [Contains: P1 protein...	<a href="#">31.2</a>	7.4	
<input type="checkbox"/>	<a href="#">sp P13529 POLG PPVD</a>	Genome polyprotein [Contains: P1 proteina...	<a href="#">31.2</a>	7.4	

Run PSI-Blast iteration 2

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:			Score	E
			(Bits)	Value
<input checked="" type="checkbox"/>	<a href="#">sp P53397 OGG1 YEAST</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">617</a>	2e-176
<input checked="" type="checkbox"/>	<a href="#">sp 008760 OGG1 MOUSE</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">462</a>	1e-129
<input checked="" type="checkbox"/>	<a href="#">sp 070249 OGG1 RAT</a>	N-glycosylase/DNA lyase [Includes: 8-oxogu...	<a href="#">459</a>	7e-129
<input checked="" type="checkbox"/>	<a href="#">sp 015527 OGG1 HUMAN</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">450</a>	4e-126
<input checked="" type="checkbox"/>	<a href="#">sp Q9V3I8 OGG1 DROME</a>	N-glycosylase/DNA lyase (dOggl) [Include...	<a href="#">417</a>	5e-116
<input checked="" type="checkbox"/>	<a href="#">sp 027397 OGG1 METHH</a>	Probable N-glycosylase/DNA lyase [Includ...	<a href="#">247</a>	6e-65
<input checked="" type="checkbox"/>	<a href="#">sp Q9SJQ6 ROS1 ARATH</a>	Protein ROS1 (Repressor of silencing 1) ...	<a href="#">114</a>	7e-25
<input checked="" type="checkbox"/>	<a href="#">sp 049498 DML3 ARATH</a>	DEMETER-like protein 3	<a href="#">104</a>	6e-22
<input checked="" type="checkbox"/>	<a href="#">sp Q9SR66 DML2 ARATH</a>	DEMETER-like protein 2	<a href="#">88.5</a>	4e-17
<input checked="" type="checkbox"/>	<a href="#">sp Q8LK56 DME ARATH</a>	Transcriptional activator DEMETER (DNA gl...	<a href="#">85.8</a>	3e-16
<input checked="" type="checkbox"/>	<a href="#">sp P37878 3MGA BACSU</a>	DNA-3-methyladenine glycosylase (3-methy...	<a href="#">52.7</a>	2e-06
<input checked="" type="checkbox"/>	<a href="#">sp P73715 END3 SYNY3</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">46.9</a>	2e-04
<input checked="" type="checkbox"/>	<a href="#">sp Q9WYK0 END3 THEMA</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">44.9</a>	5e-04
<input checked="" type="checkbox"/>	<a href="#">sp Q10630 ADA MYCTU</a>	Putative regulatory protein ada (Regulato...	<a href="#">44.9</a>	5e-04
<input checked="" type="checkbox"/>	<a href="#">sp P39788 END3 BACSU</a>	Probable endonuclease III (DNA-(apurinic...	<a href="#">43.4</a>	0.002
<input checked="" type="checkbox"/>	<a href="#">sp P46303 UVEN MICLU</a>	Ultraviolet N-glycosylase/AP lyase (UV-e...	<a href="#">42.3</a>	0.003

Run PSI-Blast iteration 3

Sequences with E-value WORSE than threshold

<input type="checkbox"/>	<a href="#">sp Q92383 MAG1 SCHPO</a>	DNA-3-methyladenine glycosylase 1 (3-met...	<a href="#">41.1</a>	0.007
<input type="checkbox"/>	<a href="#">sp Q58030 Y613 METJA</a>	Putative endonuclease MJ0613	<a href="#">41.1</a>	0.008
<input type="checkbox"/>	<a href="#">sp P54137 NTH1 CAEEL</a>	Probable endonuclease III homolog (DNA-(...	<a href="#">39.6</a>	0.020
<input type="checkbox"/>	<a href="#">sp P44319 END3 HAEIN</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">39.6</a>	0.023
<input type="checkbox"/>	<a href="#">sp Q8KA16 END3 BUCAP</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">39.6</a>	0.024
<input type="checkbox"/>	<a href="#">sp Q68W04 END3 RICTY</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">38.8</a>	0.040
<input type="checkbox"/>	<a href="#">sp Q4UK93 END3 RICFE</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">38.8</a>	0.041
<input type="checkbox"/>	<a href="#">sp 083754 END3 TREPA</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">38.0</a>	0.070
<input type="checkbox"/>	<a href="#">sp Q92GH4 END3 RICCN</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">38.0</a>	0.073
<input type="checkbox"/>	<a href="#">sp Q35980 NTHL1 MOUSE</a>	Endonuclease III-like protein 1	<a href="#">37.6</a>	0.093
<input type="checkbox"/>	<a href="#">sp P78549 NTHL1 HUMAN</a>	Endonuclease III-like protein 1	<a href="#">37.2</a>	0.11

Sequences with E-value BETTER than threshold

Sequences producing significant alignments:			Score	E	
			(Bits)	Value	
<input checked="" type="checkbox"/>	<a href="#">sp P53397 OGG1 YEAST</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">575</a>	<a href="#">1e-163</a>	
<input checked="" type="checkbox"/>	<a href="#">sp O70249 OGG1 RAT</a>	N-glycosylase/DNA lyase [Includes: 8-oxogu...	<a href="#">464</a>	<a href="#">2e-130</a>	
<input checked="" type="checkbox"/>	<a href="#">sp O08760 OGG1 MOUSE</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">456</a>	<a href="#">6e-128</a>	
<input checked="" type="checkbox"/>	<a href="#">sp O15527 OGG1 HUMAN</a>	N-glycosylase/DNA lyase [Includes: 8-oxo...	<a href="#">437</a>	<a href="#">5e-122</a>	
<input checked="" type="checkbox"/>	<a href="#">sp Q9V3I8 OGG1 DROME</a>	N-glycosylase/DNA lyase (dOggl) [Include...	<a href="#">386</a>	<a href="#">9e-107</a>	
<input checked="" type="checkbox"/>	<a href="#">sp O27397 OGG1 METHH</a>	Probable N-glycosylase/DNA lyase [Includ...	<a href="#">319</a>	<a href="#">1e-86</a>	
<input checked="" type="checkbox"/>	<a href="#">sp P37878 3MGA BACSU</a>	DNA-3-methyladenine glycosylase (3-methy...	<a href="#">235</a>	<a href="#">3e-61</a>	
<input checked="" type="checkbox"/>	<a href="#">sp P73715 END3 SYNY3</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">147</a>	<a href="#">1e-34</a>	
<input checked="" type="checkbox"/>	<a href="#">sp P39788 END3 BACSU</a>	Probable endonuclease III (DNA-(apurinic...	<a href="#">145</a>	<a href="#">2e-34</a>	
<input checked="" type="checkbox"/>	<a href="#">sp Q9S3Q6 ROS1 ARATH</a>	Protein ROS1 (Repressor of silencing 1) ...	<a href="#">144</a>	<a href="#">5e-34</a>	
<input checked="" type="checkbox"/>	<a href="#">sp Q8LK56 DME ARATH</a>	Transcriptional activator DEMETER (DNA gl...	<a href="#">138</a>	<a href="#">3e-32</a>	
<input checked="" type="checkbox"/>	<a href="#">sp P46303 UVEN MICLU</a>	Ultraviolet N-glycosylase/AP lyase (UV-e...	<a href="#">137</a>	<a href="#">6e-32</a>	
<input checked="" type="checkbox"/>	<a href="#">sp Q10630 ADA MYCTU</a>	Putative regulatory protein ada (Regulato...	<a href="#">133</a>	<a href="#">2e-30</a>	
<input checked="" type="checkbox"/>	<a href="#">sp O49498 DML3 ARATH</a>	DEMETER-like protein 3	<a href="#">118</a>	<a href="#">6e-26</a>	
<input checked="" type="checkbox"/>	<a href="#">sp Q9SR66 DML2 ARATH</a>	DEMETER-like protein 2	<a href="#">117</a>	<a href="#">9e-26</a>	
<input checked="" type="checkbox"/>	<a href="#">sp Q9WYK0 END3 THEMA</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">101</a>	<a href="#">6e-21</a>	
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp P63541 END3 MYCBO</a>	Endonuclease III (DNA-(apurinic or apyri...	<a href="#">95.0</a>	<a href="#">4e-19</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp P44319 END3 HAEIN</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">90.4</a>	<a href="#">1e-17</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q9CB92 END3 MYCLE</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">86.2</a>	<a href="#">2e-16</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp P0AB83.1 END3 ECOLI</a>	Endonuclease III (DNA-(apurinic or apy...	<a href="#">84.2</a>	<a href="#">8e-16</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q92GH4 END3 RICCN</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">82.7</a>	<a href="#">2e-15</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q68W04 END3 RICTY</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">81.5</a>	<a href="#">5e-15</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q4UK93 END3 RICFE</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">81.1</a>	<a href="#">7e-15</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp O05956 END3 RICPR</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">80.0</a>	<a href="#">1e-14</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q89AW4 END3 BUCBP</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">78.1</a>	<a href="#">5e-14</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp Q8KA16 END3 BUCAP</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">71.1</a>	<a href="#">6e-12</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp O83754 END3 TREPA</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">69.6</a>	<a href="#">2e-11</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp P54137 NTH1 CAEEL</a>	Probable endonuclease III homolog (DNA-(...	<a href="#">67.7</a>	<a href="#">7e-11</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp P04395.1 3MG2 ECOLI</a>	DNA-3-methyladenine glycosylase 2 (DNA...	<a href="#">63.8</a>	<a href="#">1e-09</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp P57219 END3 BUCAI</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">63.8</a>	<a href="#">1e-09</a>
<b>NEW</b>	<input checked="" type="checkbox"/>	<a href="#">sp O92383 MAG1 SCHPO</a>	DNA-3-methyladenine glycosylase 1 (3-met...	<a href="#">61.9</a>	<a href="#">4e-09</a>

NEW	<input checked="" type="checkbox"/>	<a href="#">sp P17802.1 SNO2_BUCAI</a>	DNA 3-methyladenine glycosylase 2 (DNA...	<a href="#">53.8</a>	<a href="#">1e-09</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp P57219 END3_BUCAI</a>	Endonuclease III (DNA-(apurinic or apyrimid	<a href="#">63.8</a>	<a href="#">1e-09</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q92383 MAG1_SCHPO</a>	DNA-3-methyladenine glycosylase 1 (3-met...	<a href="#">61.9</a>	<a href="#">4e-09</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q58030 Y613_METJA</a>	Putative endonuclease MJ0613	<a href="#">61.5</a>	<a href="#">6e-09</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp P78549 NTHL1_HUMAN</a>	Endonuclease III-like protein 1	<a href="#">59.2</a>	<a href="#">3e-08</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q2KID2 NTHL1_BOVIN</a>	Endonuclease III-like protein 1	<a href="#">58.4</a>	<a href="#">5e-08</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp P17802.1 MUTY_ECOLI</a>	A/G-specific adenine glycosylase	<a href="#">56.1</a>	<a href="#">2e-07</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q94468 MAG2_SCHPO</a>	Probable DNA-3-methyladenine glycosylase...	<a href="#">53.8</a>	<a href="#">1e-06</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q05869 MUTY_SALTY</a>	A/G-specific adenine glycosylase	<a href="#">53.4</a>	<a href="#">1e-06</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q58829 Y1434_METJA</a>	Putative endonuclease MJ1434	<a href="#">52.3</a>	<a href="#">3e-06</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q35980 NTHL1_MOUSE</a>	Endonuclease III-like protein 1	<a href="#">52.3</a>	<a href="#">4e-06</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp P31378 NTG1_YEAST</a>	DNA base excision repair N-glycosylase 1...	<a href="#">51.5</a>	<a href="#">5e-06</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q8K926 MUTY_BUCAP</a>	A/G-specific adenine glycosylase	<a href="#">51.5</a>	<a href="#">6e-06</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q10159 MYH1_SCHPO</a>	A/G-specific adenine DNA glycosylase	<a href="#">51.1</a>	<a href="#">7e-06</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q09907 END3_SCHPO</a>	Endonuclease III homolog (DNA-(apurinic ...	<a href="#">51.1</a>	<a href="#">8e-06</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp P57617 MUTY_BUCAI</a>	A/G-specific adenine glycosylase	<a href="#">50.7</a>	<a href="#">1e-05</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp P29588 GTMR_METTF</a>	G/T mismatches repair enzyme (Thymine-DN...	<a href="#">49.2</a>	<a href="#">3e-05</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q9UIF7 MUTYH_HUMAN</a>	A/G-specific adenine DNA glycosylase (MutY	<a href="#">46.5</a>	<a href="#">2e-04</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q08214 NTG2_YEAST</a>	DNA base excision repair N-glycosylase 2	<a href="#">45.3</a>	<a href="#">5e-04</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q8R5G2 MUTYH_RAT</a>	A/G-specific adenine DNA glycosylase (MutY h	<a href="#">44.5</a>	<a href="#">7e-04</a>	<b>G</b>
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q89A45 MUTY_BUCBP</a>	A/G-specific adenine glycosylase	<a href="#">44.5</a>	<a href="#">7e-04</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp P44320 MUTY_HAEIN</a>	A/G-specific adenine glycosylase	<a href="#">43.4</a>	<a href="#">0.002</a>	
NEW	<input checked="" type="checkbox"/>	<a href="#">sp Q9XAI4 RECR_STRCO</a>	Recombination protein recR	<a href="#">41.9</a>	<a href="#">0.005</a>	

Run PSI-Blast iteration 4

Sequences with E-value WORSE than threshold

<input type="checkbox"/>	<a href="#">sp Q99P21 MUTYH_MOUSE</a>	A/G-specific adenine DNA glycosylase (MutY	<a href="#">41.1</a>	<a href="#">0.008</a>	<b>G</b>
<input type="checkbox"/>	<a href="#">sp P22134 MAG_YEAST</a>	DNA-3-methyladenine glycosylase (3-methyl...	<a href="#">40.7</a>	<a href="#">0.009</a>	<b>G</b>
<input type="checkbox"/>	<a href="#">sp Q65UP1 RUVB_MANSM</a>	Holliday junction ATP-dependent DNA helicase	<a href="#">39.5</a>	<a href="#">0.022</a>	
<input type="checkbox"/>	<a href="#">sp P21269 CCA1_YEAST</a>	tRNA nucleotidyltransferase, mitochondri...	<a href="#">39.5</a>	<a href="#">0.023</a>	<b>G</b>
<input type="checkbox"/>	<a href="#">sp Q82EQ7 RECR_STRAW</a>	Recombination protein recR	<a href="#">38.4</a>	<a href="#">0.054</a>	
<input type="checkbox"/>	<a href="#">sp Q29876 OGG1_ARCFU</a>	Probable N-glycosylase/DNA lyase [Includ...	<a href="#">36.8</a>	<a href="#">0.13</a>	
<input type="checkbox"/>	<a href="#">sp P74925 MUTL_THEMA</a>	DNA mismatch repair protein mutL	<a href="#">35.3</a>	<a href="#">0.42</a>	
<input type="checkbox"/>	<a href="#">sp Q9Z512 UVRC_STRCO</a>	UvrABC system protein C (Protein uvrC) (...	<a href="#">34.9</a>	<a href="#">0.52</a>	
<input type="checkbox"/>	<a href="#">sp Q49Y80 RUVB_STAS1</a>	Holliday junction ATP-dependent DNA helicase	<a href="#">34.1</a>	<a href="#">0.83</a>	<b>G</b>

# Searching with a family of proteins

Instead of searching with a simple sequence, we can search with a family of proteins, represented by a model.

Models for the representation of a family of protein sequences:

- Set of sequences
- Consensus sequence
- Patterns: Simplified "regular expressions"
- Profiles: position-specific scoring matrices (PSSMs) based on probabilities of amino acid substitutions (Gribskov *et al.* 1987)
- Hidden Markov models (HMMs): probabilistic model for linear sequences (Haussler *et al.* 1993)

A good multiple alignment of the sequences in the family is essential for most of these models.



# Pattern example

Entry for alkylbase DNA glycosylases (AlkA) in PROSITE:

G-I-G-x-W-[ST]-[AV]-x-[LIVMFY](2)-x-[LIVM]-x(8)-[MF]-x(2)- [ED]-D

Syntax:

x denotes any amino acid

[ST] denotes alternative amino acids, e.g. serine (S) or threonine (T)

(2) denotes that the previous part should occur 2 times

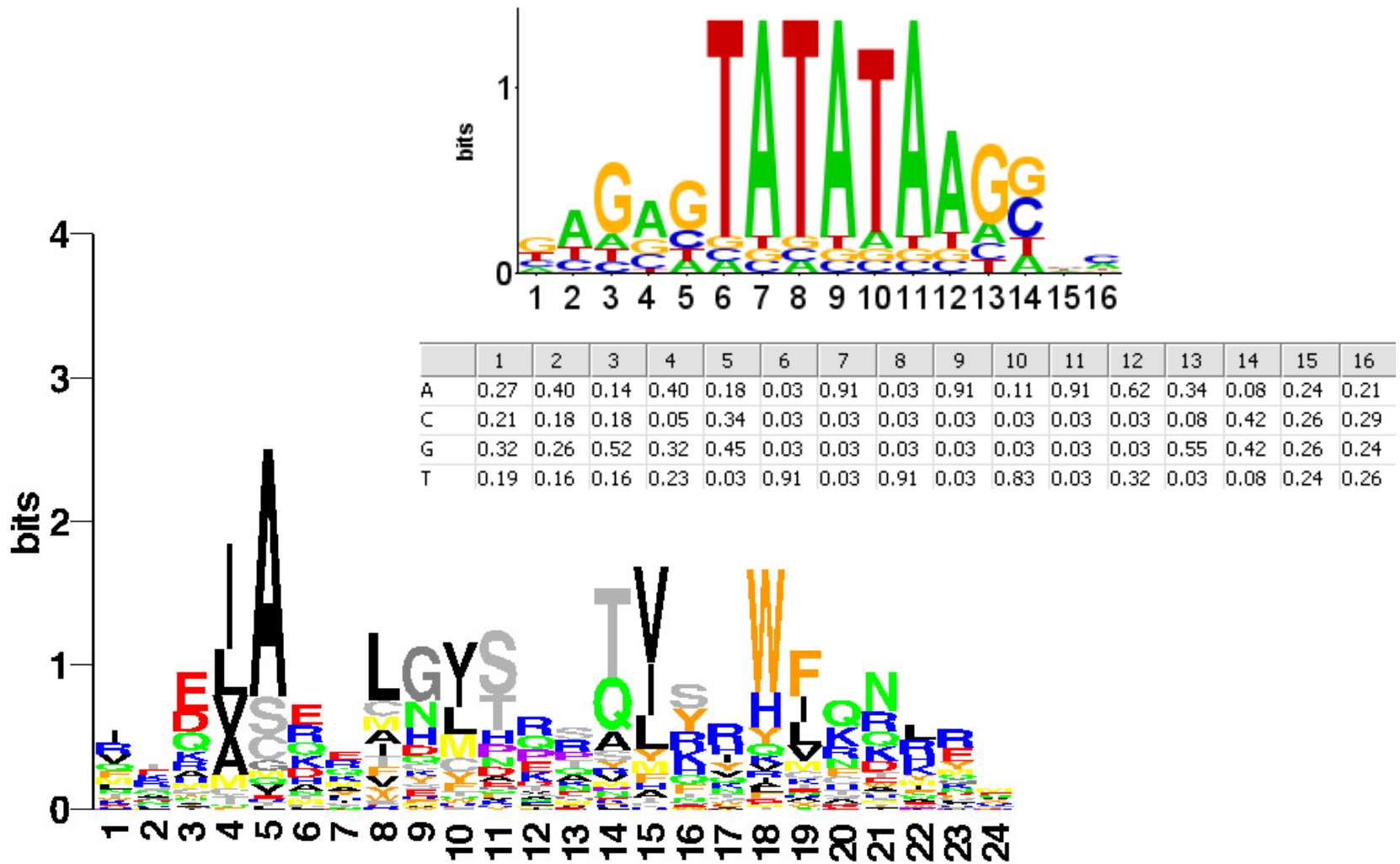
- Derived from a multiple alignment of a family of proteins
- Generated automatically (e.g. by PRATT) or manually
- Collected in PROSITE database ([www.expasy.org/prosite](http://www.expasy.org/prosite))



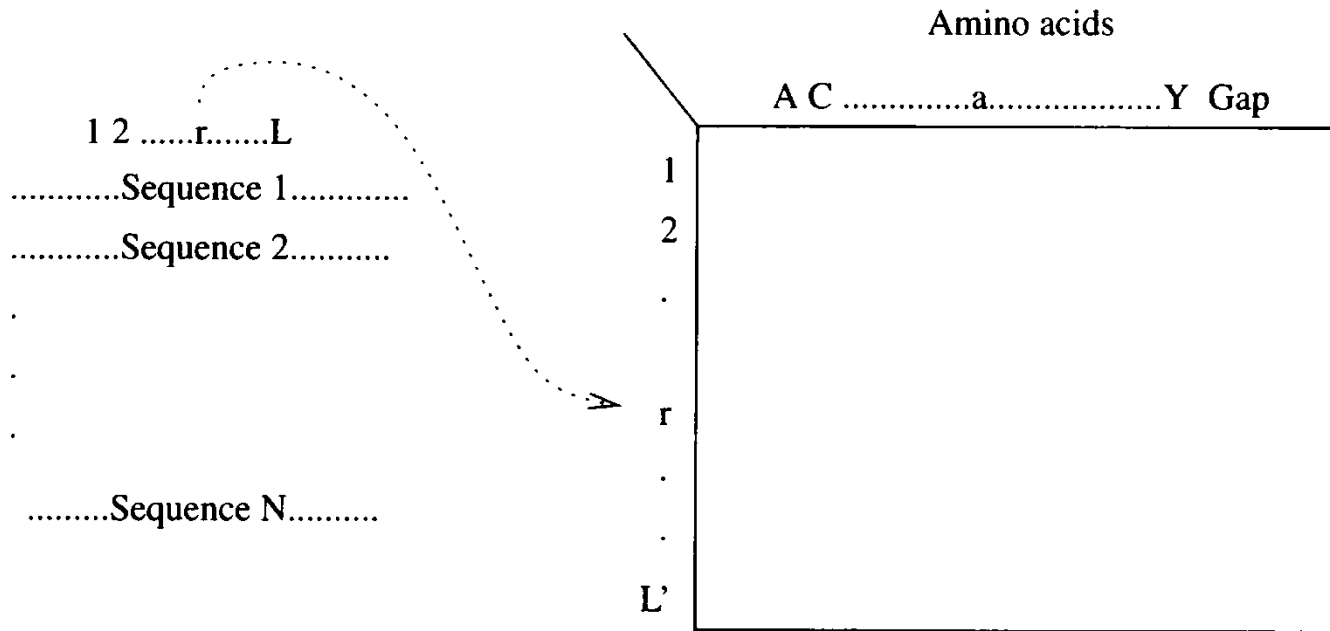
# Sequence profiles (PSSMs)

- Position-specific scoring matrices
- Based on a multiple alignment of proteins in a family
- A matrix of  $21 \times L$  cells, where  $L$  is the length of the alignment (21 for the 20 amino acids + gap)
- Scores in each cell are calculated as a weighted average of the scores from a substitution score matrix (e.g. BLOSUM62) for matching a certain amino acid with each of the amino acids present in the proteins in a specific position in the multiple alignment.
- Sequences are weighted in order to reduce the effect of many similar sequences.

# DNA and protein sequences logos



# Calculating PSSMs



**Figure 6.1** The figure shows the connection between a block of a multiple alignment and a derived profile. Positions  $1 \dots r \dots L$  are positions in the profile. Gap specifies gap penalties for the gap in profile positions. ( $L'$  can be different from  $L$  if some of the columns are left out.)

# Profile example

Position-specific scoring matrix (PSSM) derived from a family of ATP binding RNA helicases ("DEAD" box family)

Protein	Sequence
rhle_ecoli	EILVLD <b>DEAD</b> RMLDMGFIHDI
dbp2_schpo	TYLVLD <b>DEAD</b> RMLDMGFEPQI
dbp2_yeast	TYLVLD <b>DEAD</b> RMLDMGFEPQI
dbpa_ecoli	NTLVMD <b>DEAD</b> RMLDMGFSDAI
rm62_drome	TYLVLD <b>DEAD</b> RMLDMGFEPQI
p68_human	TYLVLD <b>DEAD</b> RMLDMGFEPQI
rh1b_ecoli	QVVVLD <b>DEAD</b> RMVDLGFIKDI
yn21_caeel	KFLIM <b>DEAD</b> RILNMDFEVEL
yhm5_yeast	KFLVMD <b>DEAD</b> RLLDMEFGPVL
me31_drome	RILVLD <b>DEAD</b> KLLSLDFQGML
drs1_yeast	EILVMD <b>DEAD</b> RMLEEGFQDEL
if4a_rabit	KMFVLD <b>DEAD</b> EMLSRGFKDQI
if41_human	KMFVLD <b>DEAD</b> EMLSRGFKDQI
yk04_yeast	RYIVLD <b>DEGD</b> KLMELGFDETI
...	...

Pos	Cons	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y	Gap
1	K	4	-11	16	17	-21	6	9	-7	27	-12	0	16	7	16	20	12	8	-7	0	-18	100
2	F	-3	3	-14	-10	29	-10	-3	20	-9	26	19	-7	-12	-13	-10	-6	5	14	10	20	100
3	L	-5	-19	-21	-14	42	-16	-8	34	-12	51	40	-15	-12	-10	-15	-12	-3	32	16	17	100
4	V	10	11	-11	-11	15	8	-17	66	-11	46	34	-17	4	-12	-17	-6	11	83	-43	-4	100
5	L	-4	-37	-24	-14	53	-23	-11	39	-10	73	67	-19	-14	-4	-14	-19	-4	39	17	11	100
6	D	30	-50	150	100	-100	70	40	-20	30	-50	-40	70	10	70	0	20	20	-20	-110	-50	100
7	E	30	-60	100	150	-70	50	40	-20	30	-30	-20	50	10	70	0	20	20	-20	-110	-50	100
8	A	121	24	25	25	-41	58	-8	0	0	-9	0	17	41	16	-24	33	33	16	-66	-25	100
9	D	30	-50	150	100	-100	70	40	-20	30	-50	-40	70	10	70	0	20	20	-20	-110	-50	100
10	R	-7	-16	7	10	-23	-7	17	-9	36	-14	7	8	10	18	49	5	0	-10	35	-25	100
11	M	-1	-28	-19	-10	33	-16	-12	31	1	60	62	-15	-10	-2	-1	-14	-1	31	-2	2	100
12	L	-6	-35	-26	-16	61	-26	-10	41	-15	74	64	-20	-17	-7	-20	-20	-5	39	24	18	100
13	D	12	-9	41	32	-30	23	11	-7	12	-17	-12	23	7	20	3	19	10	-7	-27	-19	100
14	M	5	-15	-3	1	8	-3	-3	13	6	25	30	-3	0	5	6	-1	2	13	-1	-6	100
15	G	26	4	31	23	-25	54	-5	-8	-1	-16	-9	17	11	10	-10	23	17	8	-41	-27	100
16	F	-31	-6	-63	-44	96	-36	-4	44	-45	77	31	-31	-45	-50	-32	-19	-19	13	82	90	100
17	E	7	-7	10	13	-7	5	4	7	10	3	5	6	2	10	3	6	5	5	-9	-7	100
18	D	12	-12	31	27	-27	19	12	-6	12	-12	-8	17	12	21	4	8	9	-3	-30	-19	100
19	D	8	-14	25	22	-18	11	13	0	9	-3	-1	13	5	24	5	3	8	0	-24	-12	100
20	I	0	-5	-8	-7	27	-10	-9	43	-7	36	29	-10	-8	-7	-11	-6	6	35	-9	5	100

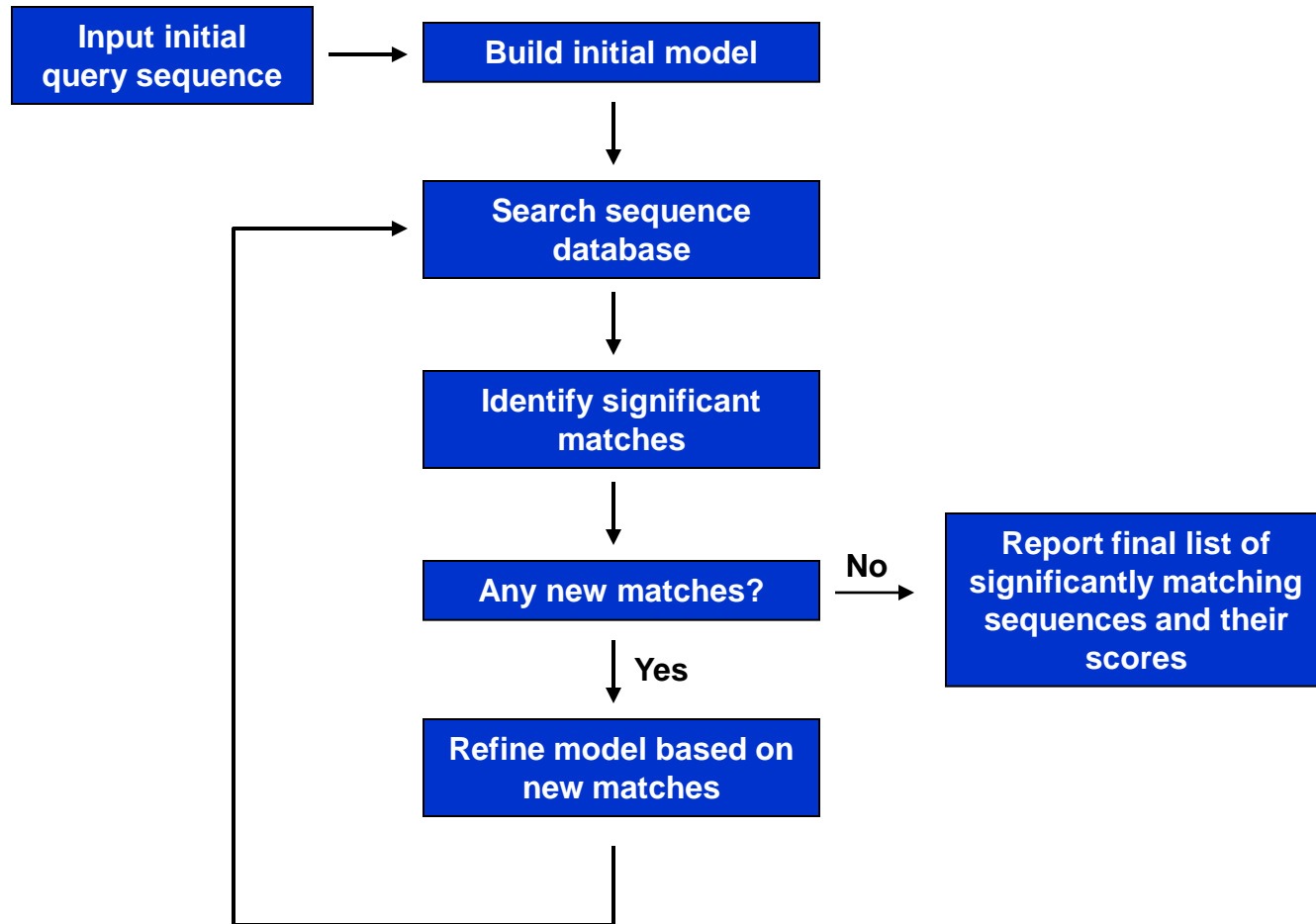
Source: Gribskov and Veretnik (1996) *Meth.Enz.* **266**, 198-212

# Amino acid substitution score matrix

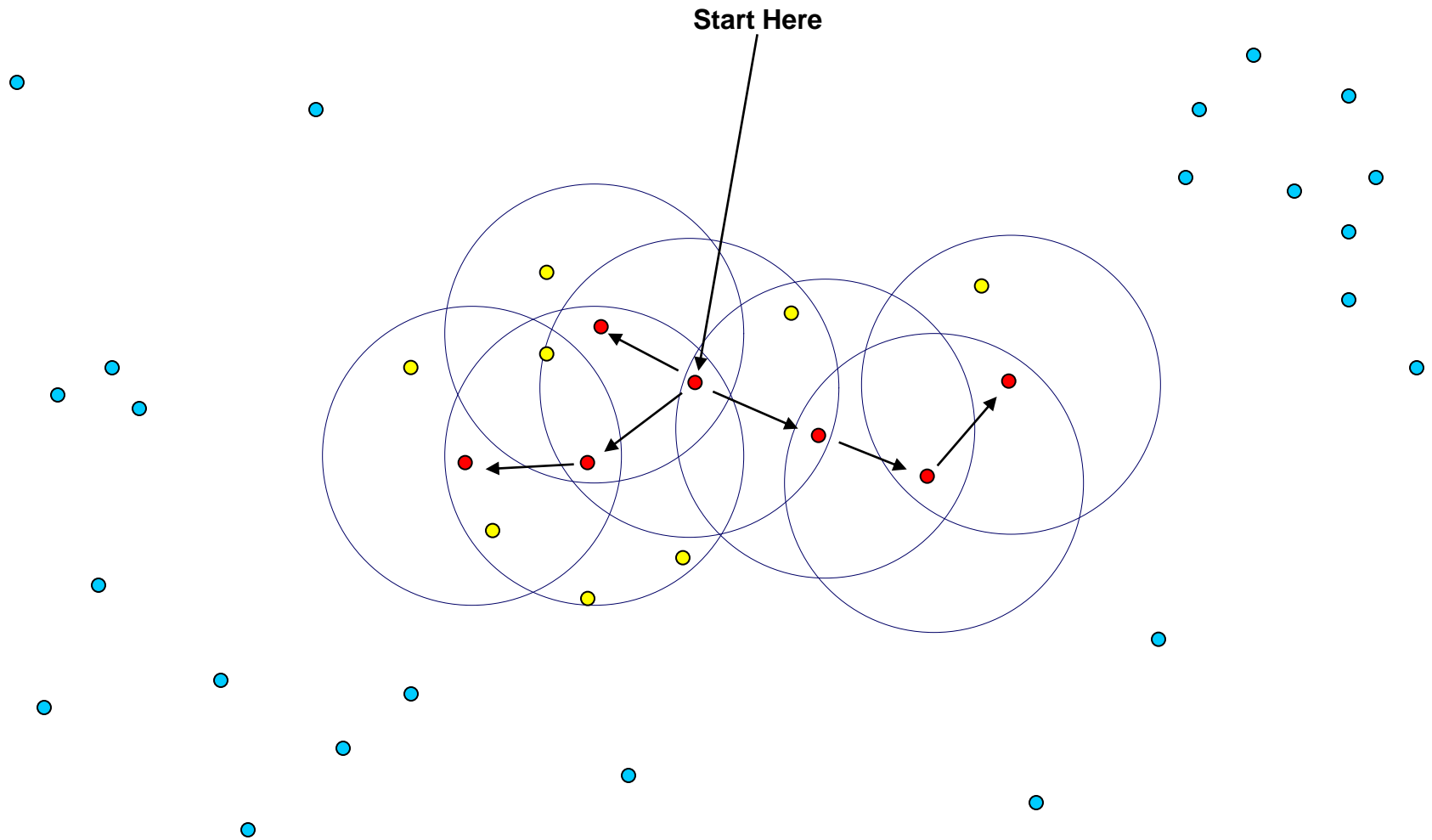
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-1	-2	-1	1	0	-3	-2	0
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	2	-1	-3	-2	-1	-1	-3	-2	-3
N	-2	0	6	1	-3	0	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2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

BLOSUM62

# Iterated searches

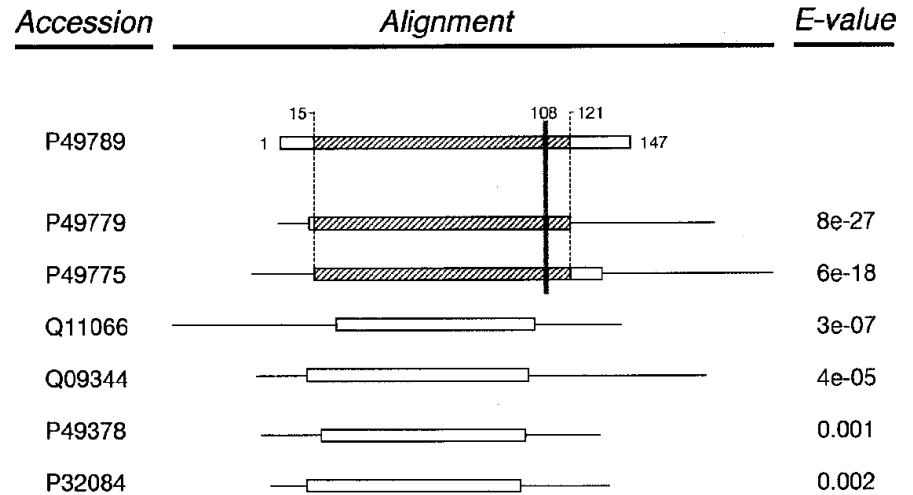


# Iterated searches in protein space



# PSI-BLAST: Calculating the profile

- Search database and collect clearly significant matches (e.g.  $E < 0.002$ )
- Create multiple sequence alignment (MSA) after each iteration
- Simplistic MSA based on pairwise alignment to the original query sequence, ignore gaps
- Only one copy of very similar sequences (98% identical) are used, other sequences are weighted



Source: Altschul *et al.* (1997) *Nucleic Acids Res.* **25**, 3389-3402



# Advanced PSI-Blasting

PSI-BLAST search with 5 rounds and inclusion threshold 0.002:

```
blastpgp -d nr -i query.fsa -JT -sT -j 5 -h 0.002
```

Same PSI-BLAST search, but save profile:

```
blastpgp -d nr -i query.fsa -JT -sT -j 5 -h 0.002 -C pssm
```

Run search with a saved profile:

```
blastpgp -d nr -i query.fsa -JT -sT -R pssm
```

Run search in translated nucleotide database (tblastn-like) with saved protein profile:

```
blastall -p psitblastn -d nt -i query.fsa -JT -FF -R pssm
```

Recommended options for blastpgp: -JT -sT

Recommended options for blastall: -JT -FF

# Using your own alignment

Run PSI-BLAST search based on own multiple alignment:

```
blastpgp -d nr -i query.fsa -JT -sT -j 5 -h 0.002 -B alignment
```

The alignment file:

```
mdap12itam  AETESPYQELQG-QRPEVYSDLNTQ
rdap12itam  AETESPYQELQG-QRPEVYSDLNTQ
hdap12itam  TETESPYQELQG-QRSDVYSDLNTQ
ssdap12itam AETESAYQELQG-QRSDVYSDLNTQ
btdap12itam PETESPYQELQG-QRTDVYSDLNTQ
ggnfam1itam PPPSPVYDCLDS-QQVEVYSVLKNN
ccnilt1itam VDPDQIYTELNASRQSDVYQSLRTD
btzetaitam2 NPNEVVYNELRKDKMAEAYSEIGMK
rzetaitam2  NPHEVVYNELRKDKMAEAYSEIGMK
mzetaitam2  NPQEGVYNALQKDKMAEAYSEIGTK
ggzetaitam2 NPHDTVYSSLQKDKMGEAYSEIGKK
mgraspitam  APGEELYAALEDYHPAELYRALAVS
```

# Literature

## PSI-BLAST paper

- *Gapped BLAST and PSI-BLAST: a new generation of protein database search programs*  
Altschul SF et al. (1997)  
**Nucleic Acids Research**, 25, 3389-3402.  
<http://nar.oupjournals.org/cgi/content/abstract/25/17/3389>



## AlkB paper

- *The DNA-repair protein AlkB, EGL-9, and Iprecan define new families of 2-oxoglutarate- and iron-dependent dioxygenases*  
Aravind L, Koonin EV (2001)  
**Genome Biology**, 2(3):RESEARCH0007.  
<http://genomebiology.com/2001/2/3/RESEARCH/0007>



## Recommended reading

- *Getting the most from PSI-BLAST*  
David T Jones and Mark B Swindells (2002)  
**Trends in Biochemical Sciences**, 27, 161-164.  
[http://dx.doi.org/10.1016/S0968-0004\(01\)02039-4](http://dx.doi.org/10.1016/S0968-0004(01)02039-4)



# Multiple sequence alignment

# Overview

- What is multiple sequence alignment?
- Examples
- An objective function – sum of pairs score
- Exhaustive algorithms
- Heuristic algorithms
  - Progressive alignment
    - CLUSTAL W
  - Iterative alignment
    - MUSCLE
  - MAFFT
- Aligning protein-coding DNA sequences

# What is a multiple alignment?

- Extension of pairwise alignments to three or more sequences
- Usually global alignments – entire sequences included
- Indicates common conserved residues in all or most sequences – usually important for function / activity
- Indicates accepted residues in the different positions
- Indicates positions where gaps are more likely
- Basis for construction of phylogenetic trees
- Basis for sequence motifs and profiles
- Essential for evolutionary studies and phylogenetics

# Example

```

CAS_Scla_322266 -----YHRLQPNYVMLACSRADHE-----RTAATLVASVRK---70---VTEAVYLEPG-DLLIVDNF-----RTDARTPFSPRWGDKRWLHRVYIRT 302\
IPNS_Bn_124825 -----LITVLYQ-----SNVQNLQVETAA-----GYQDIADDT-GYLINCGSYMAHLTNNYKAPIDHRKQVNVN----AERQSLPFFVNL 288\
FLAS_Pet_421946 -----YIILVLP-----NEVQGLQVFKDG-----HWYDVKIYIP-ALIVHIGDQVIBLSNGKYKYSVKEHTIVNK----DKTRMSWVFLFEP 309\
LDOX_Pet_1730108 -----ALTFILH-----NMVPLGLQLFYEG-----QWVTAKCVPM-SIIMHIGDTIBILSNGKYKYSILHRGVVVK----EKVRFPSWAIFCEP 311\
Srg_At_479047 -----GLTVLMQV-----NDVEGLQIKKDG-----KWPVVKPLRN-AFIVNIGDVLBIITNGTYRSIBERGVVNS----EKERLSIATFHNV 309\
EFE_Le_398992 -----GIILLLQD-----DKVSGLLQLKDE-----QWIDVPPMRH-SIVVNLGDQLBEVITNGKYKYSVLRVIAQT---DGTMSLSASFYMP 253\
Ga20Ox_Sot_10800976 -----SLTILHQ-----DSVSGLVQVMDN-----QWRISISPNLS-AFVVNIGDTFMALSNGRYKSCIDHRVVMN---KTRPKSLAFFLCP 317\
PA0147_Pa_9945977 -----CVTLLYQ-----DAAGGLQVQNRQG-----EWIDAPPIDG-TFVVNIGDMMARWENDRYRSTEHRTISPR---GVHRYSMPPFAEP 274\
PA4191_Pa_9950401 -----LITLLHQ-----DAIGGLQVTRTPQ-----GWLEAPPIDG-SFVCNLGDMLEBRMTGGLYRSTEHRRARNTS---GRDRLSLFLFFDP 277\
ISP7_Sp_729862 -----ALTLMSQ-----DNVKGLEILDPEVSN-----CFLSVSPAPG-ALIANLGDIMAILTNNRYKSSMIRVCNNS---GSDRYTIPFFLQG 353\
SPCC1494.01_sc_7491815 -----SITLLFQ-----RDAAGLEIRPPNFVKDM---DWIKVNVQPD-VVLVNIADMLQFWTSGKLRSTVHRVIDPG---VKTRQTIAYFVTP 267\
DACS_Ly1_769809 -----IVLILQTPCP-----NGFVSLQVEIDG-----RFVEVPPRGD-CVVVFCGSIAPLVSDDGKIKAPQHRVVS-PGA4-GSNRTSSVLFLLRP 268\

RRPO_SHVX_548840 IYPKG-----NKILIVNAA-----GSGTFSI-----KCAKGE-TTLNLEDGD-YFQMPSGFQETHRKNVA----VTPRLSITFRSTV 743\
POL_ASPV_487652 IYDIN-----HQVLTVNSY-----GDAIFCI-----ECLGSGF-EIPLSGPQ-MLLMPFGFQKSHRSGIKSP---SKGRISLTLFRLLK 853\
POL_BSV_409711 IFMRG-----APVHVSMD-----GNADFGT-----ECAAGR-QYTLTRGNVQFTMPSGFQETHRKNVNT---TAGRVSYTFRRLA 841\
RRPO_PMW_139137 CYPKG-----HQVLTINHS-----GCLTQI-----ACQKGA-SITMGGD-YLSPVGFQESHKAPVNT---TKGRVSLTRFRCTV 690\
POL_GLV_1154656 IFEKD-----SKILIVCIQ-----GDCBFFF-----RCATGET-GFYMEAPK-QFMMPDGFQSNHVEATREBC---TPGRISAIFRRAK 772\
Pol_GVA_1405615 CYLPG-----GSVVIVNLH-----GDAATFEVK-----ENQSGKIEKKELHDGD-VYVMSPGMQQTTHKRVVTSH---TDGRCSITLNRNK 738\
RRPO_ACLSV_1710717 CYDD-----DEILTINVV-----GDAKFHT-----TC-HGE--IIDLRQGD-BILMPGGYQKMNKAWEVA---SEGRTSVTLRVHK 836\
T13L16.2_At_2708738 FL-----RPFCTISFL-----SECDILEGSLNKVSE-----GPGDFSGSY-SIPLPVGS-VLVLKNGADVARKVCPAV---FTKRSITLFRKMD 420\
T19K4.220_At_3036813 FL-----RPFCTVSL-----SECNILFGSLNKVL-----GPGDFSGSY-SIPLPVGS-VLVLKNGADVARKVCPAV---FTKRSITLFRKMD 403\
At2g48080_At_4249414 -----QPISLVL-----SESTMVFGHRLGVD-----NDGNFRGL-TLPLKES-LLVMRGNADMARVVCPS---FNKRVAITFFFLK 351\
AK000315.1_Hs_7020317 IFE-----RPIVSVSFF-----SDCALCFGCKFQFK-----PIRVSEVLSLFPVRRGS-VTVLSGYAADBITHCIRPQDI---KERRAVILRKR 270\
CG17807_Dm_7291441 AFL-----DPILSLSLQ-----SLVVMDFRRG-----DDQV-QVRLPRRS-LLMSGEARVDWBEIRPKHID13RGRKTSITFRRLR 325\
CG6144_Dm_7297712 FH-----PIIISTISG-----AHEVLEFVKREDTITTEBAGDQTTREVLV-KLLEPRS-LLILKDTLYTDYLAISSETSBD24RSPKISLTIKRV 213\
CG4036_Dm_7297561 IWGERVVTVNC-----LGSVLTLLT--PYEVQSGKYNLDLVAIYDELLAP-LLTDDQLATFEGKVLRIKMPNLS-LIVLYGPARYQFHSVLRREDV---QERFVCAVREFT 278\
FLJ2001_Hs_38923019 LWGERLVSLLN-----LSPVLSLWC-----REAPGSLLLCSAPSAEPAEALVDSVIAPSRVLCQSEVAIPLPARS-LLVLTGAARHQWKAIAHRRHI---BARFVCTVFRRLS 274\
C14B1.10_Ce_6580210 AFD-----DPIVSISLL-----SDVVMDFKD-----GANSARIAPVLLKARS-LCLIQGESRYRWKRGIVNRRYD10RQTFVSLTLRKR 343\
SPAP8A3.02c_Sp_7491301 FGDG-----VAIFGFLSN-----TIMIFTRPE-----LKLK--KIRLEKGS-LLMSGTARYDWRBEIPFRAGD12RSQRLSITMRRII 219\
L3377.4_Lm_9989036 VYD-----DIFAI CSLG-----SNCLLRFVH-----VQNGBEL-DVMVPPRS-VYIMSGEARVYVYFHWLPV---BAQFSLVFRRSI 193\
MTC1237.14c_Mtu_2052134 RGSTEDTM-----VAIVSLGAT-----RVFALRP-----RGRGSLRLPLAHD-LLVMGGSQRTTFHVPKTSAP--TGRVSVIQFRPRD 203\
ALKE_Cc_2055386 ADPR-----FPLLSISLG-----DIATVIGG-----VNRKIDPTRSLRLAGSD--VCRLLGPARLARBGGDRILPG6-GGGRIINLTLRRAR 190\
ALKE_Bc_113638 PDLR-----APIVSVSLG-----LEAIFQFGG-----LKRNDPLKRLLEHGD--VVVWGSBSRLFYGGIQLKAG5-IDCFYNLTFRQAG 213\
ALKE_Scoe_8894829 RTD-----APVSVSLG-----DTCVFRFGN-----PETRTRPYDTELRSGD--LFVFGGSRSLAYGTPRVHPG7-LRGRINLTLRVS 215\
ALKE_At_4835778 ADWS-----KPIVGSLSG-----CKAIFLLGGK-----SKDDPEHAMLRSGD--VVLWAGEARBCEGILLHFQL34KTSRININIRQV 354\
ALKE_Sp_3080529 BDLT-----LPLLSLSG-----LDCIVLIGTE-----SRSEKPS-ALRLHSGD--VVIMTGSRKARHGKHC---SFYKLYLSQLIA 272\
ALKE_Hs_2134723 LDHS-----KPLLSFSFG-----QSAIFLLGGL-----QRDEAPP-PFMHSGD--IMIMSGFSRLLMBAFVRLFN39KTA RVNMA RQVL 272\
Consensus (85%): .....sh.h.....s..h.....s..h.....H.s.....+h.h..b...

```

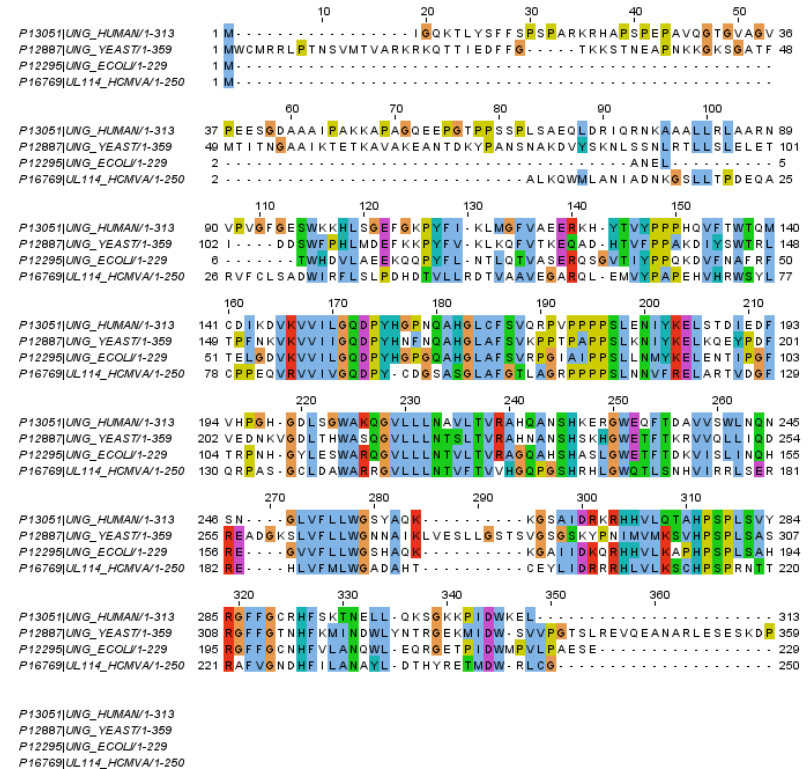
# A lot of software...

- MSA – full DP
- DCA – divide and conquer
- Clustal W / Clustal X - progressive
- T-Coffee - progressive
- DbClustal - progressive
- Poa - progressive
- PRALINE - progressive
- PRRN - iterative
- MUSCLE - iterative
- Dialign, Dialign2 – blocks-based
- Match-Box – blocks-based
- MAFFT – various techniques
- ProbCons – probabilistic
- ...



# Jalview demo/example

- Jalview is a multiple sequence alignment editor
- [www.jalview.org](http://www.jalview.org)
- Can run the algorithms Clustal W, MUSCLE and MAFFT from within the program
- Very useful for making nice colorful figures



# Finding the best MSA

- How should we design a score model for MSA that leads us to the best biological alignment?
- Given a scoring model, how do we find the alignment that achieves the highest score?

# Sum-of-pairs (SP) score

- Objective function indicating multiple alignment quality
- The sum-of-pairs score (SP score) is a simple and commonly used objective function
- Corresponds approximately to the sum of all pairwise alignment scores
- For the alignment  $A$  of  $m$  sequences  $s^1$  til  $s^m$  we have the sum-of-pairs score  $S(A)$ :

$$S(\mathcal{A}) = \sum_{i=1}^{m-1} \sum_{j=i+1}^m S(\bar{s}^i, \bar{s}^j).$$

- $S(a,b)$  is the pairwise score of  $a$  and  $b$ , and  $\bar{s}^i$  is the projection of  $s^i$ , that is,  $s^i$  with inserted gaps
- This leads to an alignment that optimizes all pairwise alignments with equal weight on all pairs.
- Columns with gaps in all sequences are ignored.
- The SP score does not necessarily give the biologically best alignment

# Calculating the sum-of-pairs score

M	Q	P	I	L	L	L
M	L	R	-	L	L	-
M	K	-	I	L	L	L
M	P	P	V	L	I	L

$$\text{score}(k) = S(P,R) + S(P,-) + S(P,P) + S(R,-) + S(R,P) + S(-,P)$$

score for  
column  $k = 3$

We have  $S(-,-) = 0$

$$\text{Total score} = \text{score}(1) + \text{score}(2) + \dots + \text{score}(N)$$

# Exhaustive algorithms

- Identifies the mathematically optimal alignment that is guaranteed to have the maximum sum-of-pairs score
- Dynamic programming
- Smith-Waterman generalized to N dimensions
- MSA
- (DCA)
- Computationally very demanding
- Only practical for a few short sequences
- Impossible for many and/or long sequences

# Progressive algorithms

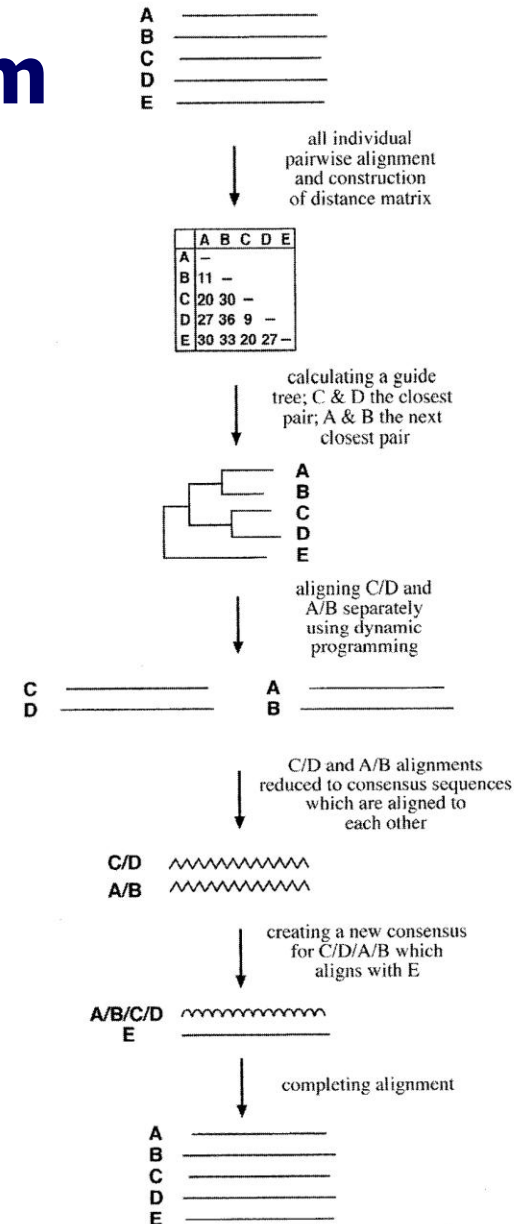
- Clustal W
- T-Coffee
- DbClustal
- Poa
- PRALINE

# Progressive alignment

- Start with one sequence and add progressively more sequences into the alignment
- Align one sequence with an alignment, or align to alignments with each other
- Disadvantage: Early errors in alignment may be fatal and cannot be fixed later
- The choice of which sequences to align and in which order is usually based on clustering or phylogenetic analysis and a guide tree.

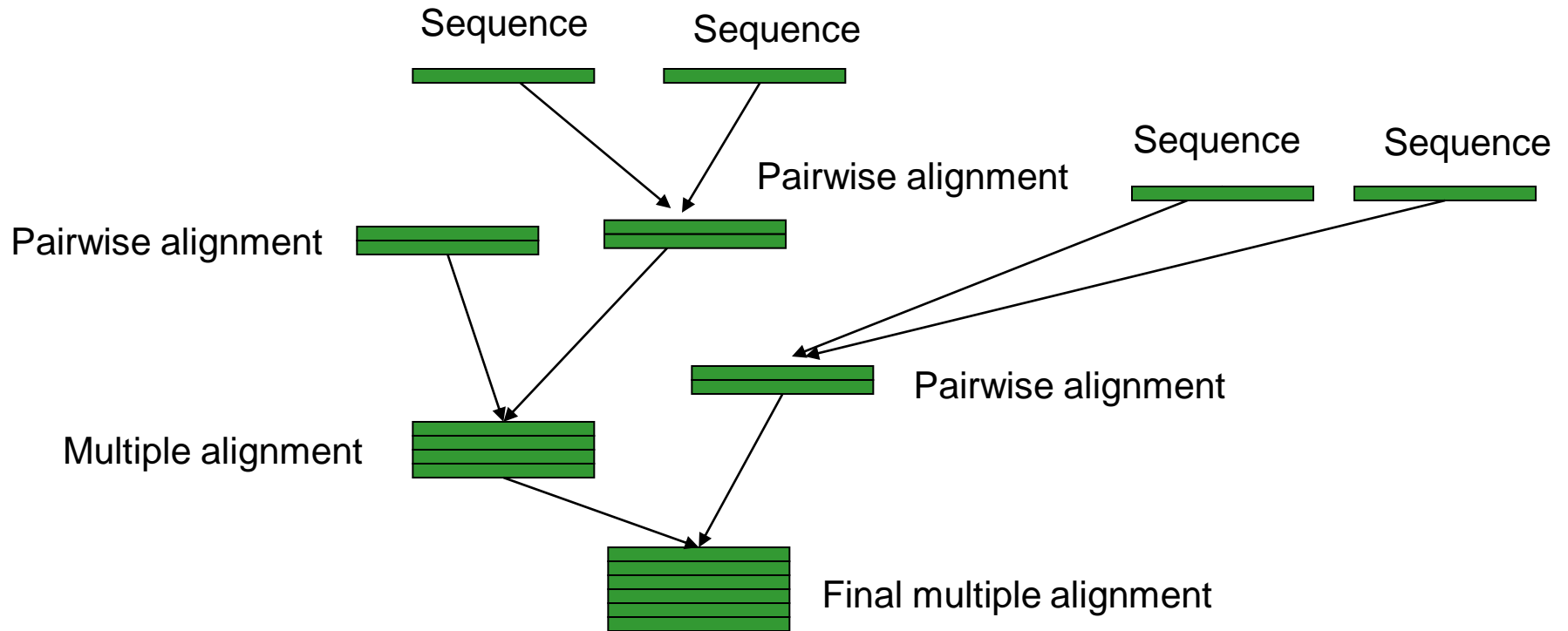
# Progressive alignment algorithm

1. Pairwise alignment of all  $n(n-1)/2$  pairs of sequences. Construct distance matrix from scores.
2. Group sequences progressively based on pairwise distances and produce "guide tree"
3. Based on guide tree, start by realigning the two closest sequences and treat them as one "consensus" sequence in later steps
4. Gradually realign all pairs of sequences according to guide tree, always starting with the closest sequences, until all sequences are aligned
5. Repeat step 4 until all sequences are aligned.





# Profile alignment



# Clustal W

- One of the most commonly used and well-known tools for multiple sequence alignment
- Now somewhat outdated and surpassed by other tools.
- Many versions, including one with a graphical user interface (Clustal X).
- Paper: Thompson JD, Higgins DG, Gibson TJ. (1994) CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. Nucleic Acids Res. 22, 4673-4680.
- <http://www.ebi.ac.uk/clustalw/>

# Iterative alignment

- Iteration is used to gradually increase the quality of the alignment by repeating some steps in the algorithm until a good alignment is reached
- Reduces/eliminates the effect of the order of aligning sequences used in progressive alignment

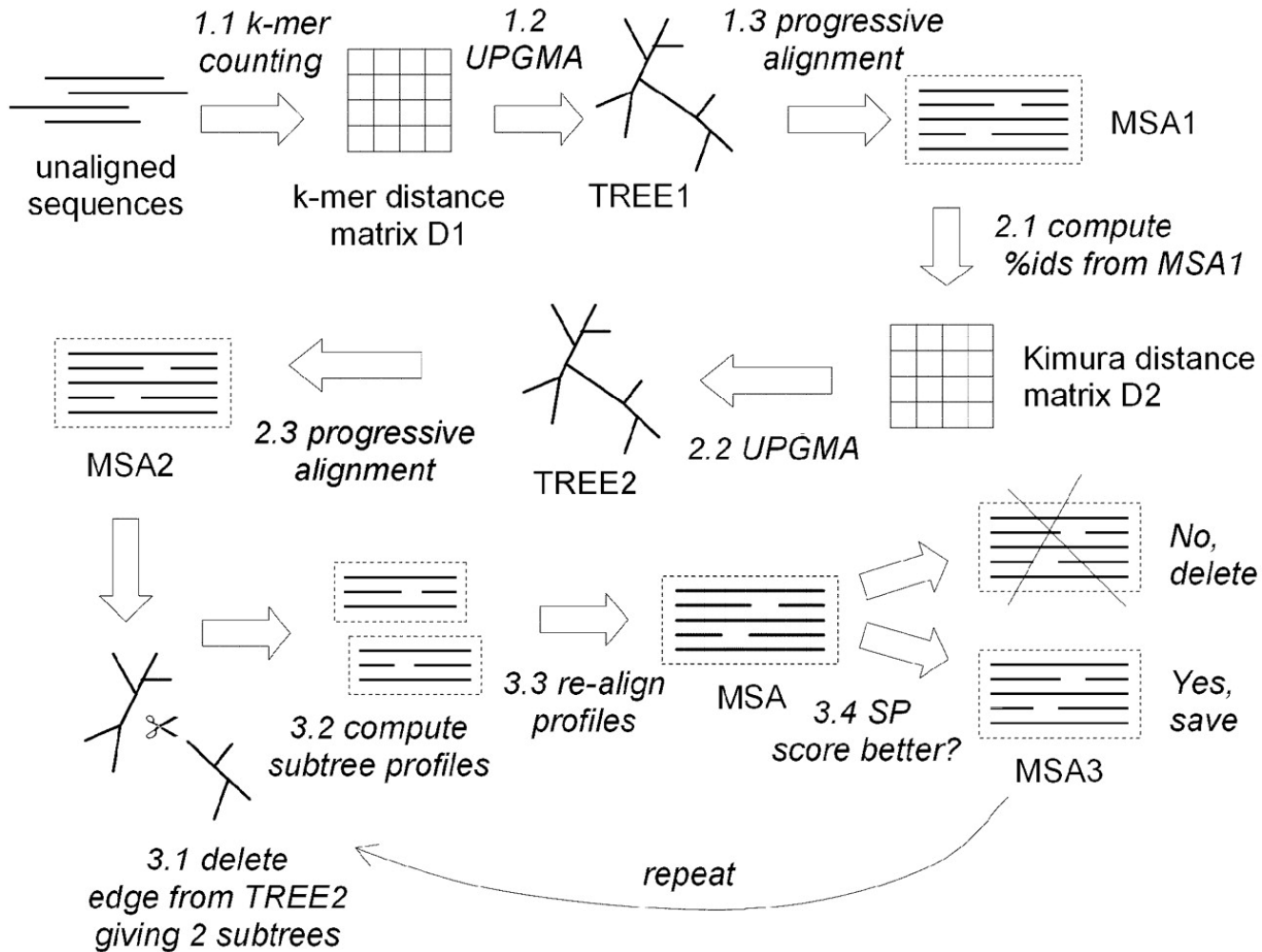
Programs:

- PRRN
- Muscle

# MUSCLE

- MUSCLE = Multiple Sequence Comparison by Log Expectation
- Very high quality of alignments
- Much faster than Clustal W
  
- Improvements:
  - Bases initial sequence distances on k-mer counting
  - Progressive alignment based on a new and better score function
  - Iterated partitioning of the tree and realignment
  
- Edgar RC (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.*, 32, 1792-1797.

# MUSCLE - Algorithm

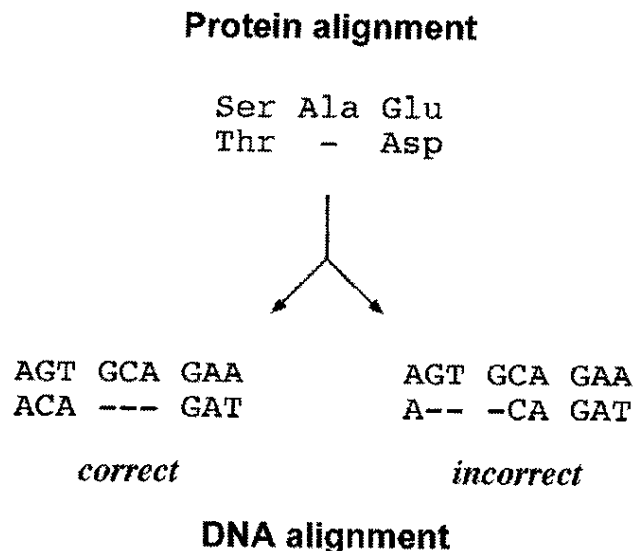


# MAFFT

- MAFFT is one of the most accurate MSA programs currently available
- Uses a combination of approaches
- Improved objective function
- <http://align.bmr.kyushu-u.ac.jp/mafft/software/>

# Alignment of protein-coding DNA sequences

- Protein-coding DNA sequences should be aligned at the amino acid level
- Alignment at the DNA level may introduce unwanted gaps within codons
- Translate DNA sequences to protein before alignment, then translate back if necessary
- RevTrans and PROTA2DNA are multiple sequence alignment programs that align sequences by translating them to proteins



**Figure 5.5:** Comparison of alignment at the protein level and DNA level. The DNA alignment on the left is the correct one and consistent with amino acid sequence alignment, whereas the DNA alignment on the right, albeit more optimal in matching similar residues, is incorrect because it disregards the codon boundaries.