



Earlham
Institute

Earlham Institute summer school on bioinformatics

25-29 July 2016

Basic Bioinformatics Sessions

Practical 4: Multiple Sequence Alignment

Thursday 21 July 2016

Multiple Sequence Alignment

Here we will look at some software tools to align some protein sequences. Before we can do that, we need some sequences to align. I propose we try all the human **homeobox** domains from the well annotated section of **UniprotKB**. Getting the sequences is a trifle clumsy, so concentrate now! There used to be a much easier way, but that was made redundant by foolish people intent on making the future ever more tricky!!

So, begin by going to the home of **Uniprot**:

<http://www.uniprot.org/>

Choose the **Advanced** option of the **Search** button.

First specify that you are only interested in **Human** proteins. To do this, set the first field to **Organism [OS]** and **Term** to **Human [9606]**.

Set the second field selector to **Reviewed** and the corresponding **Term** to **Yes** (that is, choose to find only **SwissProt** entries).

Click on the **+** button to request a further field selection option. Set the new field to **Function**. Set the type of **Function** to **DNA binding**. Set the **Term** selection to **Homeobox**.

The screenshot shows the Uniprot search criteria form. It includes a search bar at the top right. The criteria are: Organism [OS] Human [9606], AND Reviewed Yes, AND Function DNA binding, Term Homeobox, Length range 50-70, Evidence Any assertion method. There are buttons to add, remove, and clear criteria, and a search button at the bottom right.

From previous investigations, you should be aware that a **Homeobox** domain is **generally 60** amino acids in length. To avoid partial and/or really weird **Homeobox** proteins, set the **Length** range settings to recognise only **homeoboxes** between **50** and **70** amino acids long.

Leave the **Evidence** box as Any assertion method, one does not wish to be too fussy! Address the **Search** button with authority to get the search going.

Entry	Entry name	Protein names	Gene names	Organism	Length
P52952	NKX25_HUMAN	Homeobox protein Nkx-2.5	NKX2-5 CSX,NKX2.5,NKX2E	Homo sapiens (Human)	324
P26367	PAX6_HUMAN	Paired box protein Pax-6	PAX6 AN2	Homo sapiens (Human)	422
Q99697	PITX2_HUMAN	Pituitary homeobox 2	PITX2 ARP1,RGS,RIEG,RIEG1	Homo sapiens (Human)	317
Q99801	NKX31_HUMAN	Homeobox protein Nkx-3.1	NKX3-1 NKX3.1,NKX3A	Homo sapiens (Human)	234
P49639	HXA1_HUMAN	Homeobox protein Hox-A1	HOXA1 HOX1F	Homo sapiens (Human)	335
Q92786	PROX1_HUMAN	Prospero homeobox protein 1	PROX1	Homo sapiens (Human)	737
Q01860	PO5F1_HUMAN	POU domain, class 5, transcription ...	POU5F1 OCT3,OCT4,OTF3	Homo sapiens (Human)	360
P43699	NKX21_HUMAN	Homeobox protein Nkx-2.1	NKX2-1 NKX2A,TITF1,TTF1	Homo sapiens (Human)	371
Q01826	SATB1_HUMAN	DNA-binding protein SATB1	SATB1	Homo sapiens (Human)	763
Q15475	SIX1_HUMAN	Homeobox protein SIX1	SIX1	Homo sapiens (Human)	284

A fine miscellany of sequences will assemble upon you screen. Most seem to declare themselves in possession of a **Homeobox** or two (including **PAX6_HUMAN**), so I suggest a declaration of success.

Now save the entire list into a file using the [Download](#) button. Set the download to **uncompressed**. Make sure you have **all** sequences selected and that **Text** (i.e. **EMBL** or **SwissProt**) format selected. Press the **Go** button and do whatever it takes to ensure your results end up in a file residing on your **Desktop** called:

human_homeobox_proteins.emb

```
ID NKX25_HUMAN Reviewed; 324 AA.
AC P52952; A8K3K0; B4DNB6; E9PBU6;
DT 01-OCT-1996, integrated into UniProtKB/Swiss-Prot.
DT 01-OCT-1996, sequence version 1.
DT 06-JUL-2016, entry version 173.
DE RecName: Full=Homeobox protein Nkx-2.5;
DE AltName: Full=Cardiac-specific homeobox;
DE AltName: Full=Homeobox protein CSX;
DE AltName: Full=Homeobox protein NK-2 homolog E;
GN Name=NKX2-5; Synonyms=CSX, NKX2.5, NKX2E;
OS Homo sapiens (Human).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
OC Catarrhini; Hominidae; Homo.
OX NCBI_TaxID=9606;
RN [1]
RP NUCLEOTIDE SEQUENCE [MRNA] (ISOFORM 1).
RC TISSUE=Heart;
RX PubMed=8900537;
RA Turbay D., Wechsler S.B., Blanchard K.M., Izumo S.;
RT "Molecular cloning, chromosomal mapping, and characterization of the
RT human cardiac-specific homeobox gene hCsx.";
RL Mol. Med. 2:86-96(1996).
```

Take a swift look at the file you have just created. Your neat list of **Human Homeobox** sequences will have transformed into a flood of many **SwissProt** format **UniProtKB** entries. Ugly, but what is required.

Search (Control F) for the term **DNA_BIND**.

It should occur many times (at least once per sequence) in the Feature Tables and most often refer to a **Homeobox** region.

In the **DNA_BIND** Feature Table entries, the position of the **Homeoboxes** are recorded and will be used by the next program to isolate the sequence of the **Homeoboxes**.

```
FT CHAIN 1 374 Pre-B-cell leukemia transcription factor
FT 4.
FT /FTid=PRO_0000049241.
FT Homeobox; TALE-type.
FT {ECO:0000255|PROSITE-ProRule:PRU00108}.
FT V -> I (in dbSNP:rs8108180).
FT /FTid=VAR_059355.
FT M -> V (in dbSNP:rs8108981).
FT /FTid=VAR_059356.
FT T -> M (in a colorectal cancer sample;
FT somatic mutation; dbSNP:rs376647012).
FT {ECO:0000269|PubMed:16959974}.
FT /FTid=VAR_036439.
FT CONFLICT 368 368 I -> T (in Ref. 1; BAG53471).
FT {ECO:0000305}.
FT SQ SEQUENCE 374 AA; 40854 MW; B9CE8BE93D0B7ABC CRC64;
MAAPPRPAPS PPAPRRLDTS DVLQQIMAIT DQSLDEAQR KHALNCHRMK PALFSVLCEI
KEKTVVSIRG IQDEDPPDAQ LLRLDNMLLA EGVCRPEKRG RGGAVARAGT ATPGGCPNDN
SIEHSDYRAK LSQIRQIYHS ELEKYEQACR EFTTHVTNLL QEQSRMRPVS PKEIERMVGA
IHGKFSAIQM OLKQSTCEAV MTLRSRLDA RRKRRNFSKQ ATEVLNEYFY SHLNPNYPSE
EAKEELARKG GLTISQVSNW FGNKRIRYKK NMGKFQEEAT IYTGKTAVDI TEVGVPGNHA
SCLSTPSSGS SGPFPLPSAG DAFLTLRTLA SLQPPPGGC LQSAQGSWQ GATPQPATAS
PAGDPSINS STSN
//
```

Now to extract from the whole protein sequences you have saved in a file, the sequences of just the **Homeobox** domains. One way of doing this (possibly not the best), is to use an **EMBOSS** package program called **extractfeat**. This can be found in many places, including the Bioinformatics server at **Wageningen** in the Netherlands. Go to:

<http://emboss.bioinformatics.nl/>

EDIT

[aligncopy](#)
[aligncopypair](#)
[biosed](#)
[codcopy](#)
[cutseq](#)
[degapseq](#)
[descseq](#)
[entret](#)
[extractalign](#)
[extractfeat](#)

Find the program **extractfeat** (in the **EDIT** section), and set it going.

Use the **Choose File** button to **upload** the **SwissProt** format sequences from **UniProtKB** that you saved in the file **human_homeobox_proteins.emb**.

Set **Type of feature to extract field** to **DNA_BIND** (Make sure you remove the “*”).

Set **Value of feature tags to extract** to **Homeobox*** (Make sure you append the “*” to ensure hits with, for example “homeoboxes”).

Set the **Output sequence format** to **SwissProt** (Fasta would do, but **SwissProt** retains more annotation).

Click on the **Run extractfeat** button to start **extractfeat** going. Many sequences of **60** amino acids (or so) in length will leap into view.

Input section

Select an input sequence. Use one of the following three fields:

1. To access a sequence from a database, enter the USA here:
2. To upload a sequence from your local computer, select it here: human_homeobox_proteins.emb
3. To enter the sequence data manually, type here:

Additional section

Amount of sequence before feature to extract

Amount of sequence after feature to extract

Source of feature to display

Type of feature to extract

Sense of feature to extract
(default is 0 - any sense, 1 - forward sense, -1 - reverse sense)

Minimum score of feature to extract

Maximum score of feature to extract

Tag of feature to extract

Value of feature tags to extract

Output section

Output introns etc. as one sequence?

Append type of feature to output sequence name?

Feature tag names to add to the description

Output sequence format

Run section

Email address:

If you are submitting a long job and would like to be informed by email when it finishes, enter your email address here.

OUTPUT FILE [outseq](#)

```
ID NKX25 HUMAN 138 197 Reviewed; 60 AA.
DE [DNA_contact] Homeobox protein Nkx-2.5 (Cardiac-specific homeobox) (Homeobox protein CSX) (Homeobox protein NK-2 homolog E)
SQ SEQUENCE 60 AA; 7514 MW; 16EE564D071E5E8A CRC64;
RRKPRVLFSQ AQVYELERRF KQQRYSAPR RDQLASVLKL TSTQVKIWFQ NRRYCKKRQR
//
ID PAX6 HUMAN 210 269 Reviewed; 60 AA.
DE [DNA_contact] Paired box protein Pax-6 (Aniridia type II protein) (Oculorhombin)
SQ SEQUENCE 60 AA; 7447 MW; 075C194DB9F33ED9 CRC64;
LQNRNRSFTQ EQIEALEKEF ERTHYPDVFA RERLAADIDL PEARIQVWFS NRRAKWRREE
//
ID PITX2 HUMAN 85 144 Reviewed; 60 AA.
DE [DNA_contact] Pituitary homeobox 2 (ALL1-responsive protein ARP1) (Homeobox protein PITX2) (Paired-like homeodomain transcription factor 2) (R1EG bicoid-related homeobox transcription factor) (Solurshin)
SQ SEQUENCE 60 AA; 7622 MW; 49CF61CFC17E1E0E CRC64;
QRRQRTHFTS QQLQLEATF QRNRYPDMST REEIAVWNL TEARVRVWFK NRRAKWRKRE
//
ID NKX31 HUMAN 124 183 Reviewed; 60 AA.
DE [DNA_contact] Homeobox protein Nkx-3.1 (Homeobox protein NK-3 homolog A)
SQ SEQUENCE 60 AA; 7339 MW; F665B481E2E574BB CRC64;
QKRSRAAFSH TQVIELERKF SHQKYSAPR RAHLAKNLKL TETQVKIWFQ NRRYKTKRKR
//
ID HXA1 HUMAN 229 288 Reviewed; 60 AA.
DE [DNA_contact] Homeobox protein Hox-A1 (Homeobox protein Hox-1F)
SQ SEQUENCE 60 AA; 7365 MW; 53E2BC59B06F544E CRC64;
PNAVRTNFTT KQLTELEKEF HFKNYLTRAR RVEIAASLQL NETQVKIWFQ NRRMKQKKRE
//
```

Right click the **outseq** button and select **Save Link as...** . Do whatever it takes to save all your **Homeobox** domains into a file residing on your **Desktop** called:

homeobox_human.emb

Finally, we have some sequences with which to investigate the multiple sequence alignment programs.

Take a look at the file you have created. You should have many human **homeobox** domains in **SwissProt** format, looking rather as they did in your browser window. Happily **ClustalX**, the first multiple alignment program to be investigated, accepts multiple sequence **SwissProt** format files as input.

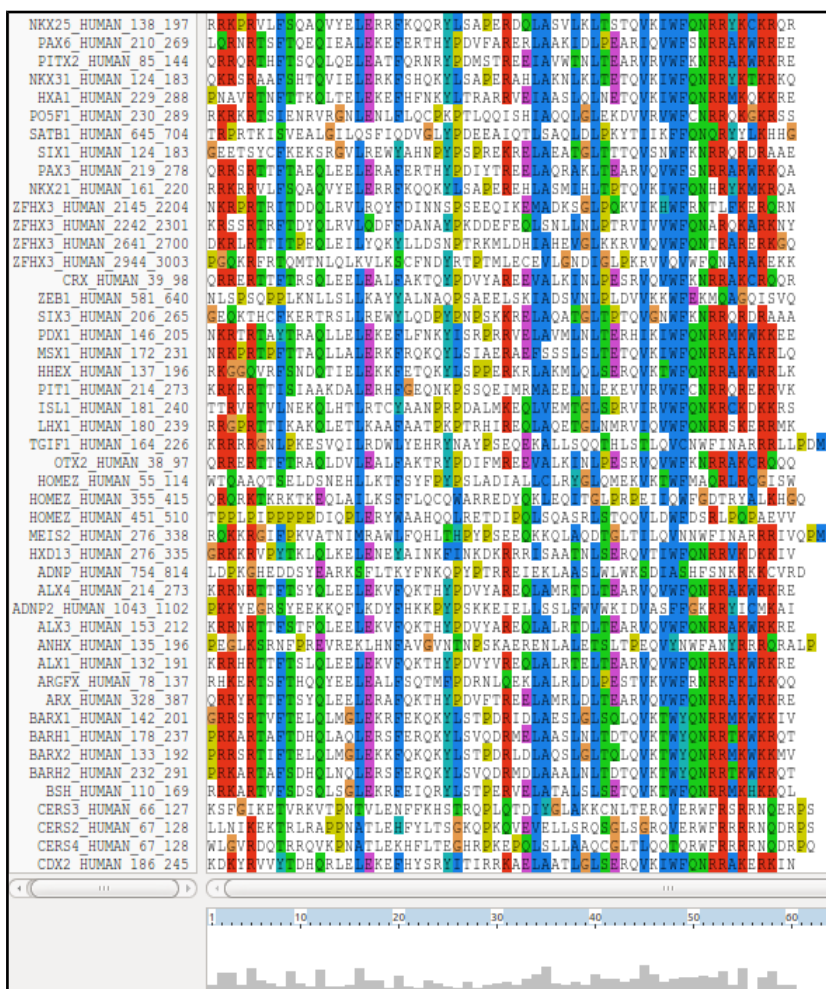
ClustalX is a part of the mostly widely known family of **Multiple Sequence Alignments (MSA)** programs, originating in the **1980s**. Until relatively recently, it was the only real option. **ClustalX** still has merit, although it lacks some of the sophistication of more recent programs. **ClustalX** runs on effectively all workstations and has a nice **Graphical User Interface (GUI)**. A good place for us to start. It is installed on your workstations.

Start up the program **ClustalX**¹. The **ClustalX** **Graphical User Interface (GUI)** will regally mount your screen.

Select **Load Sequences from the File** pull down menu and load your file of **homeobox** domains.

The sequences will arrange themselves colourfully. Many of the **homeoboxes** are similar enough to look convincing even before alignment. Note the “Manhattan skyline” under the sequences indicating the varying degrees of conservation.

Font: 10 You might like to increase the **Font** size from the minute default setting designed for Hawks and Eagles, to something more comfortable. **24** works tolerably well for me.



Pairwise Parameters

OK

Fast-Approximate

Slow/Accurate Pairwise Parameters Fast/Approx Pairwise Parameters

Gap Penalty [1-500]: 3

K-Tuple Size [1-2]: 1

Top Diagonals [1-50]: 5

Window Size [1-50]: 5

From the **Alignment** pull down menu, go to the **Alignment parameters** menu and select **Pairwise Alignment Parameters**. Just for a moment, change the setting from **Slow-Accurate** to **Fast-Approximate**. Bring the corresponding parameters into view by clicking on **Fast/Approx Pairwise Parameters** tab².

Hopefully, we will have discussed the way **ClustalX** (and similar multiple alignment tools) work. Intuitively, it should not make a lot of difference how the initial pairwise comparison stage is conducted. However, it very often does.

Specifically for this set of proteins, as well as generally, **ClustalX** will give a noticeably better alignment if the initial pairwise alignment stage is done carefully. Accordingly, reverse your whimsical setting change by moving back from **Fast-Approximate** to **Slow-Accurate**.

¹ Of course, you could run **Clustal** from websites all over the world if you wished. Specifically, it is available at the Bioinformatics server at **Wageningen**. Try it if you have time. You get the same results but will, sadly, lose the pretty interface.

<http://www.bioinformatics.nl/tools/clustalw.html>

The **EBI** no longer offer basic **Clustal** any longer.

² The **Fast-Approximate** algorithm is essential that which the database searching program **fasta** employs. Assuming we have discussed how **fasta** (or **blast**) works, it should require little further explanation here.

Click on the **Slow/Accurate Pairwise Parameters** tab for a final look at the default parameters to be used. The **Slow-Accurate** option is essentially a version of **Global Alignment** algorithm we will have discussed previously. Hopefully, all the parameter options will therefore be familiar to you.

I will assume both sets of parameters at least seem familiar? If not please ask. The default **Slow/Accurate Pairwise Parameters** you now have in view are fine. Click the **OK** button to dismiss the **Pairwise Parameters** window.

Before proceeding, save the **homeobox** sequences in **FASTA** format, which will better suit the other **MSA** programs we will try. Do this by selecting **Save sequences as...** from the **File** pull down menu. Deselect **CLUSTAL format**, select **FASTA format**.

Change the default file output file name to **homeobox_human_full**

Click **OK**. A file called **homeobox_human_full.fasta** will be created. Take a look to check it is as you would expect.

Strangely, saving your sequences in **FASTA** format convinces **clustalx** that it should now output its alignments in **FASTA** format. To prevent this, select **Output Format Options** from the **Alignments** pull down menu. Deselect **FASTA format** and select **CLUSTAL format**. Click **OK**.

From the **Alignment** pull down menu, select **Do Complete Alignment**. Accept the default names for output files and click on the **OK** button. **ClustalX** will start to think deeply and eventually come up with it view of how the **homeobox** domains should be aligned.

Not a bad first try. From an entirely non scientific, cosmetic, viewpoint, the ragged ends offend a trifle, as does the gap just before position 30!



In reality, these features might be very interesting, but here I go for pretty!

So, just to investigate what is possible, select all the **homeobox** sequences that are causing the gap around position **30** by clicking on their names (quite a lot of them I fear). Hold the **Ctrl** key down to allow multiple selection.

All selected, go to the **Edit** pull down menu and select **Cut Sequences**. Then select **Remove Gap-Only columns** from the **Edit** pull down menu. Nasty gap gone ... along with all scientific credibility, but ... never mind.

You could recompute the alignment from scratch with the reduced sequence set ending up with the same answer. Just for the sake of it, select **Select All Sequences** from the **Edit** pull down menu. Then select **Remove All gaps** from the **Edit** pull down menu and confirm your intentions. You are now back where you started, but without the sequences that mess up the alignment intolerably!

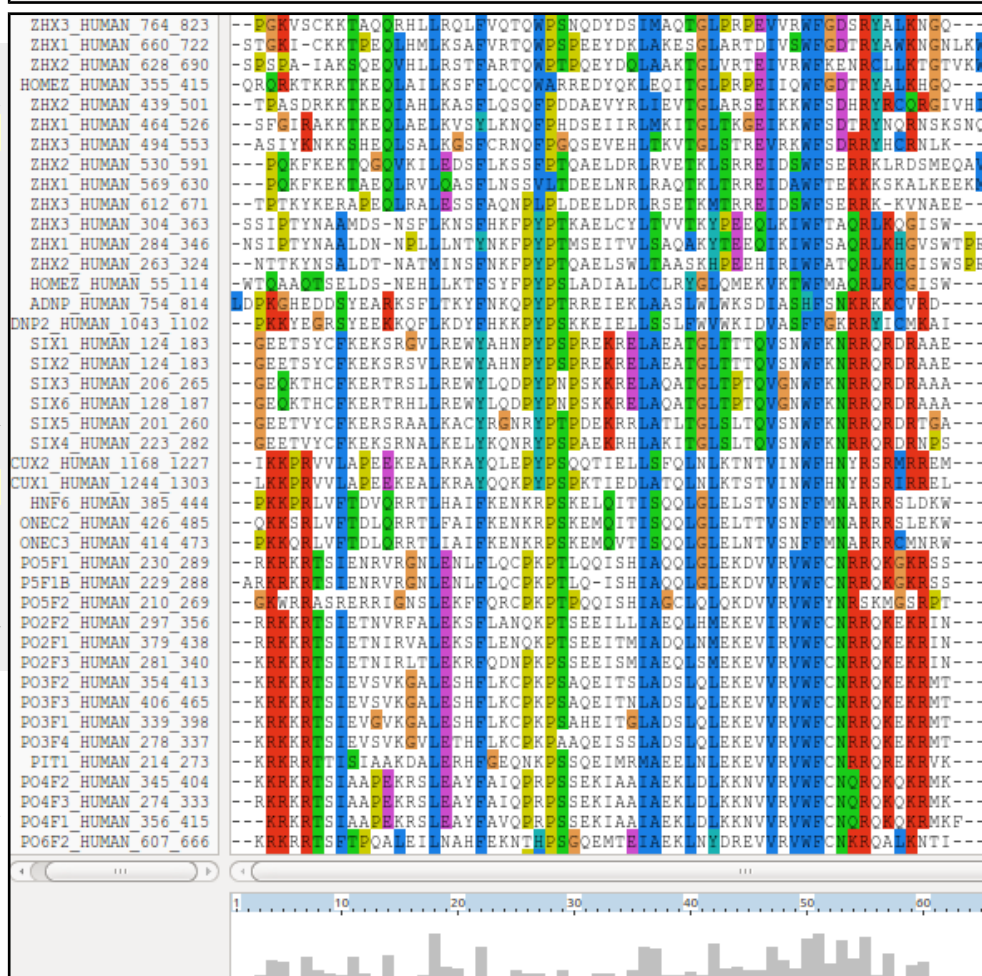
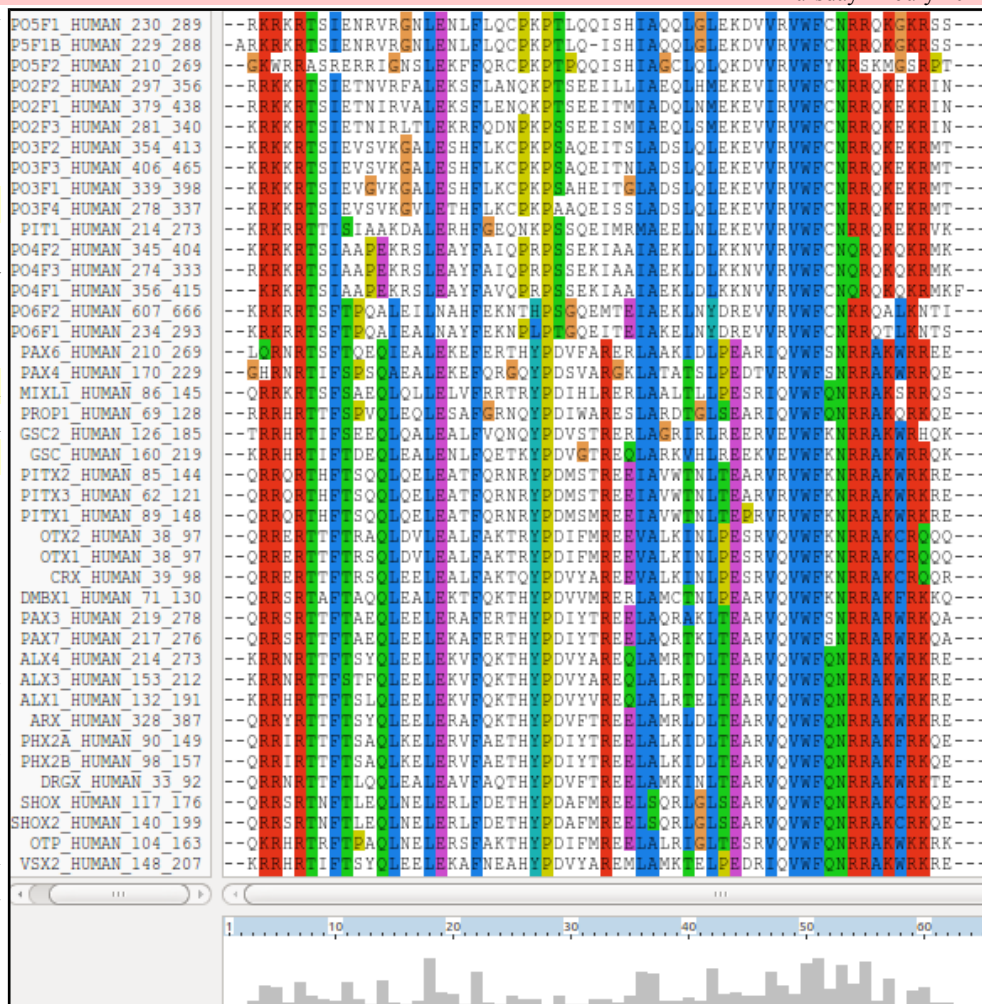
Save your filtered set of sequences. From the **File** pull down menu select **Save Sequences as...** Choose **FASTA format only**. This time, create a file with the default name:

homeobox_human.fasta

The full original set of sequences was saved in a differently named file, as a precaution. I am convinced the sequences eliminated would not align convincingly with any of the tools we have at hand. Let us lose them! Press the **OK** button.

From the **Alignment** pull down menu, select **Output Format Options** and then select **CLUSTAL format only**.

From the **Alignment** pull down menu, select **Do Complete Alignment**. Accept the default names for the output files. This will overwrite your previous efforts, but no matter. More deep thought. Well, I got back to where I was, no gaps around position **30** but still with ragged ends!



It is difficult to prove you have exactly the same alignment as previously as the order of the **MSA** will be different. This order being determined by the pairwise comparison stage of the **clustalx** MSA computation.

The **Prosite** motif database uses patterns to represent protein features (in addition to HMMs). The pattern for a **homeobox** is the ever memorable:

```
[LIVMFYFG] - [ASLVR] -x (2) - [LIVMSTACN] -x- [LIVM] - {Y} -x (2) - {L} - [LIV] - [RKNQESTAIY] -
[LIVFSTNKH] -W- [FYVC] -x- [NDQTAH] -x (5) - [RKNAIMW]
```

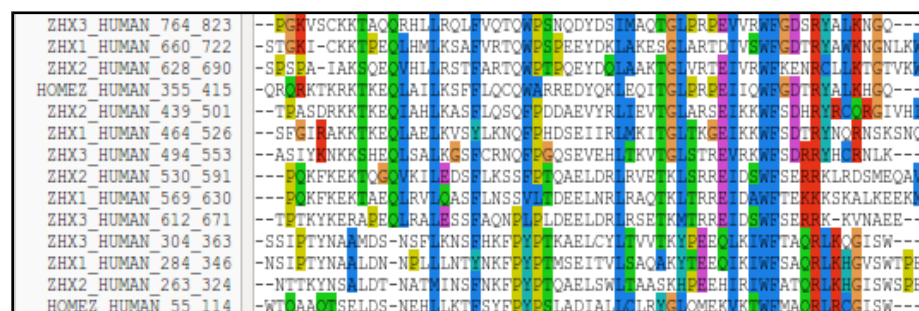
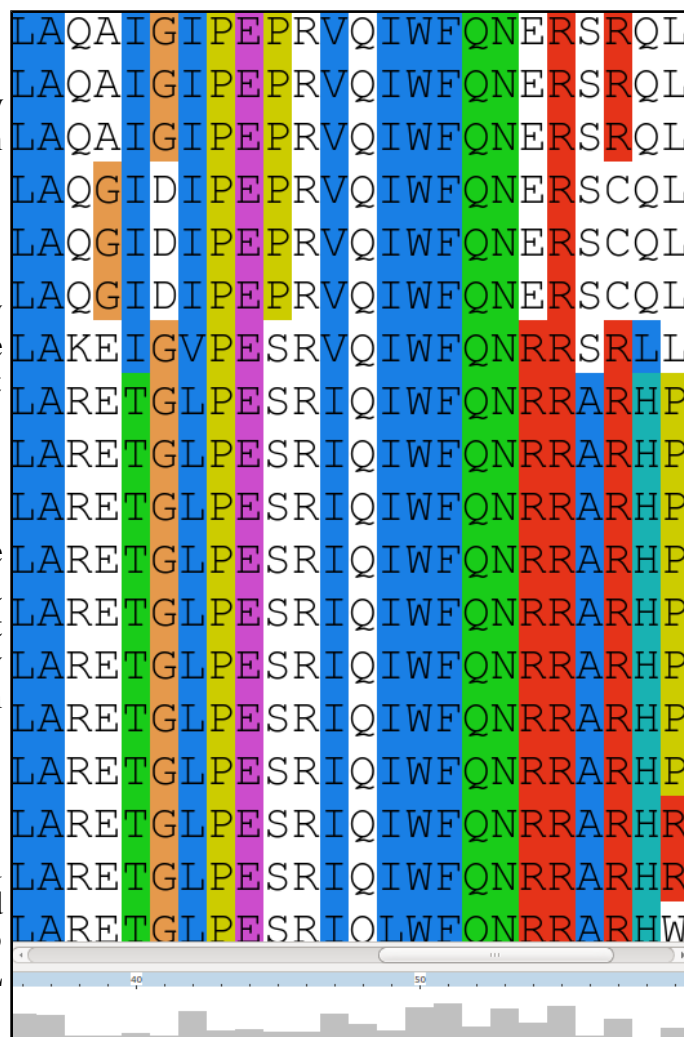
Any speculations as to how this might be interpreted?

This pattern corresponds to positions 36 to 59 in my alignment. See that the “Manhattan Skyline” is encouraging in the parts of this region that matter.

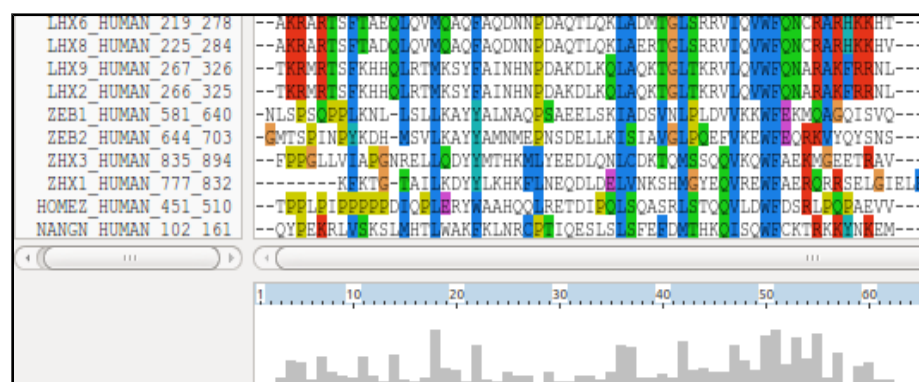
Note that the profile **Tryptophan**, in position 50, is very consistent, but not quite **100%** as suggested by the **Prosite** pattern³. The **W** was even conserved in the sequences that were cosmetically removed.

Position 52 is not conserved (“-x-”) according to the **Prosite** pattern. In the alignment segment offered here, it looks like a pretty consistent **Q**. However, the “**Manhattan skyline**” at this position is very low, suggesting that the sequences in view might not be typical of the whole alignment set. Which, upon checking they are not!

Looking through this alignment, I get the feeling I could design a better, stricter pattern for the region between 36 and 59. Possibly true, but remember the pattern in **Prosite** aims to represent the conservation of **Homeobox** domains in **ALL** organisms. Here we have only sequences from **Human**.



Of course, things are not quite so convincing throughout. If you look at the top and bottom few sequences, you will see that **ClustalX** had its moments of uncertainty.



Note, however, the consistent **W** in position 50 despite the surrounding crumble.

³ From the “**Manhattan Skyline**”, you can see the conservation is less than **100%**. Less conserved than the **F** that immediately follows in fact? Look at your alignment, the “**Manhattan Skyline**” does not seem to reflect reality? The **W** is very well conserved, although the scoring matrices would regard any deviation from **W** as serious? I need to find out more about how the **Skyline** is computed.

Now to show existence of some **msa** program options available on the web. There are many. They are available from a number of server sites. An obvious place to start has to be the **EBI** page dedicated to **MSA**. **Go to:**

<http://www.ebi.ac.uk/Tools/msa/>

Offered here is a selection of popular, current generation **MSA** tools. Each is accompanied by advice to guide the choice of tool to best fit the circumstances. Each tool is provided with a link to its **Launch** interface. All the **Launch** interfaces are very consistent. Once you have run one of the **MSA** options, you should have no trouble running any of the others.

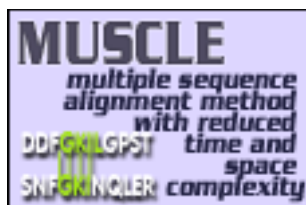
Clustal Omega ? New MSA tool that uses seeded guide trees and HMM profile-profile techniques to generate alignments. Suitable for medium-large alignments. Launch Clustal Omega	MUSCLE ? Accurate MSA tool, especially good with proteins. Suitable for medium alignments. Launch MUSCLE
Kalign ? Very fast MSA tool that concentrates on local regions. Suitable for large alignments. Launch Kalign	MView ? Transform a Sequence Similarity Search result into a Multiple Sequence Alignment or reformat a Multiple Sequence Alignment using the MView program. Launch MView
MAFFT ? MSA tool that uses Fast Fourier Transforms. Suitable for medium-large alignments. Launch MAFFT	T-Coffee ? Consistency-based MSA tool that attempts to mitigate the pitfalls of progressive alignment methods. Suitable for small alignments. Launch T-Coffee
	WebPRANK The EBI has a new phylogeny-aware multiple sequence alignment program which makes use of evolutionary information to help place insertions and deletions. Try it out at WebPRANK .

Here I intend to align again the human **homeboxes** with just one of the tools on offer. Then take a quick look at how the machine generated multiple alignment can be manually edited using **Jalview**, a program that is installed on your workstation and that you have already used as an alignment viewer when investigating **Pfam** and also **Jpred**.

Then I will invite you to try a few of the other options for yourself and see that they do not all produce the same alignment! Differences reflect not only the parameters selected, which we will have discussed, but also the particular objectives of the program selected. For example, a multiple protein sequence alignment optimal for investigating conservation of protein structure might well not be identical to one best representing protein evolution.

Used to align the **Homeobox** sequences used in this exercise, I do not expect you will see much difference between the outputs of any of these options. They will all work sufficiently on such a simple data set.

The program whose use I choose to describe carefully, leading on to a short **Jalview** exercise is **MUSCLE**. I choose thus as **MUSCLE** is now the first choice of most of the people with whom I work. Also popular are **Clustal Omega**, **MAFFT** and, for **phylogeny**, **WebPRANK**.



So the plan now is to use **MUSCLE**⁴ to align again the **homeobox** sequences previously aligned with **clustalX**. **MUSCLE** works in a way similar to **clustalX** but it takes rather more care in the generation of the **Guide Tree** used to control the order of pairwise construction of the final multiple alignment⁵. Particularly for more difficult alignments, **MUSCLE** should do a better job than **clustalX**. The alignment you will generate here will certainly be different. I leave you to judge for yourselves whether it is better.

Start by requesting to [Launch MUSCLE](#).

Use the [Browse...](#) button to upload the file containing the **FASTA** format **homeobox** sequences, **homeobox_human.fasta**. This file should not included the sequences with a mess around position 30.

STEP 1 - Enter your input sequences

Enter or paste a set of sequences in any supported format:

Or upload a file: [Browse...](#) homeobox_human.fasta

STEP 2 - Set your Parameters

OUTPUT FORMAT: [ClustalW](#)

The default settings will fulfill the needs of most users and, for that reason, are not visible.

[More options...](#) (Click here, if you want to view or change the default settings.)

Take a look at the

Set your Parameters section of the page. I find the claim that “*The default settings will fulfill the needs of most users and, for that reason, are not visible*” a little strange? What about the users who are not in the

category “*most*”? I want all the control over the programs I use that their creators deemed it sensible to make available⁶?

The default settings behind the [More options...](#) button are not those that affect the computation of the **MSA**. I confess myself confused at the lack of any meaningful options to consider? I was expecting at least the **gap open** and **gap extension penalty** options (available elsewhere, including **Wageningen**), plus a way to change the **scoring matrix**. I have inquired why things are as they are (most recently **2016.04.17**). No practical issue here, as I intended to suggest the defaults whatever they were. **Look** at the range of settings for the **OUTPUT TREE** parameter. **none** is indeed the thinking persons choice, but ... one or the other (but not both?) of the **Guide Trees** that **MUSCLE** will compute can be saved if you wish⁷. You may also set the **OUTPUT ORDER** to **aligned** or ... **aligned**?

STEP 2 - Set your Parameters

OUTPUT FORMAT: [ClustalW](#)

OUTPUT TREE: [none](#) OUTPUT ORDER: [aligned](#)

ClustalW

Pearson/FASTA

ClustalW

ClustalW (strict)

HTML

GCG MSF

Phylip interleaved

Phylip sequential

There are a number of **OUTPUT FORMATS** offered. For a quick glance at your results, both **ClustalW** or **HTML** are fine. Here I suggest it would be nice to generate an output that can be downloaded and viewed in **Jalview**⁸. The default **ClustalW** or **Pearson/FASTA** serve for this purpose. As **ClustalW** looks more like an alignment in the web page, I choose **ClustalW**⁹.

How do the options for the **OUTPUT TREE** relate to the output files of **ClustalX** and the difference between the way that **ClustalX** and **muscle** work?

Comment on how one might choose between the range of options offered for the aligned parameter?

⁴ Available from a variety of websites in addition to the **EBI**, including the Bioinformatics server at **Wageningen**: <http://www.bioinformatics.nl/tools/muscle.html>

⁵ As discussed, superficially at least, previously. I hope.

⁶ I have asked the **EBI** about their policy (the same for all the locally provided **MSA** options). Discussion is ongoing (**2016.04.20**).

⁷ A useful option if you thought it possible you might want to rerun **MUSCLE** with different parameter setting for the stages after the **Guide Tree(s)** are generated. The same possibilities exist for **ClustalX**. Of course, utterly pointless if it is impossible to control the relevant parameters so I really cannot see the point of any of the **More options** section? I am open to elucidation from all/any sources.

⁸ A widely used **java** alignment editor and viewer.

⁹ But feel free to try the others. **HTML** is the default at **Wageningen**. The **Phylip** formats are the best if you are going to analyse your output further with the phylogeny programs of the **PHYLP** package.

After considering these enigmas, or before if you prefer, Click on the **Submit** button and sit back to admire **muscle** in action.

The alignment that is computed is, superficially at least, similar to that offered by **ClustalX**.

The alignment is irritatingly split into two sections. A nice extra parameter might have been “How wide would you like your alignment to be”? A problem with the format rather than the program, to be fair.

At the very bottom of the page, **muscle** whines:

PLEASE NOTE: Showing colors on large alignments is slow.

So click the **Show Colors** button at the top of the page and try to live with the pain of such gross Trans-Atlantic inept spelling in a European site!!! Good Grief! They get everywhere!!

Well, an improvement I suppose? Colours are very useful (even slow ones) in the interpretation of alignments. Various colour schemes are used to clarify the message of alignments. Colouring can indicate shared amino acid properties not immediately evident when the letter representations differ.

But any decoration available here is far short of what can be achieved with **Jalview**, so click on the **Download Alignment File** button to save you alignment in a file on your **Desktop** called:

homeobox_human_muscle.aln

```

ARX_HUMAN_328_387      --QRRYR-TTFTSYQLEELERAFQKTHYPDVFTRLEAMRLDLTEARVQVWFQNRRAKWR
ALX1_HUMAN_132_191     --KRRHR-TTFTSLQLEELKVFQKTHYPDVVREQLALRTELTEARVQVWFQNRRAKWR
ALX4_HUMAN_214_273     --KRRNR-TTFTSYQLEELKVFQKTHYPDVYAREQLAMRTDLTEARVQVWFQNRRAKWR
ALX3_HUMAN_153_212     --KRRNR-TTFTSYQLEELKVFQKTHYPDVYAREQLALRDLTEARVQVWFQNRRAKWR
ISL1_HUMAN_181_240     --TTRVR-TVLNEKQLHLRTCYAANPRPDALMKEQLVEMTGLSPRVIRVWFQNKRCCKDK
ISL2_HUMAN_191_250     --TTRVR-TVLNEKQLHLRTCYAANPRPDALMKEQLVEMTGLSPRVIRVWFQNKRCCKDK
LHX9_HUMAN_267_326     --TKRMR-TSFKHHQLRTMKSIFYAINHNPDADKDLKLAQKTGLTKRVLQVWFQNARAKFR
LHX2_HUMAN_266_325     --TKRMR-TSFKHHQLRTMKSIFYAINHNPDADKDLKLAQKTGLTKRVLQVWFQNARAKFR
LHX6_HUMAN_219_278     --AKRAR-TSFTAELQVMQAQFAQDNNPDQATLQKLADMTGLSRRVIQVWFQNCRAHKK
LHX8_HUMAN_225_284     --AKRAR-TSFTAELQVMQAQFAQDNNPDQATLQKLADMTGLSRRVIQVWFQNCRAHKK
ZFHX3_HUMAN_2641_2700  --DKRLR-TTITPEQLEILYQKYLSDSNPTRKMLDHIAREVGLKRRVVQVWFQNTRARER
ZFHX4_HUMAN_2560_2619  --DKRLR-TTITPEQLEILYQKYLSDSNPTRKMLDHIAREVGLKRRVVQVWFQNTRARER
ZFHX2_HUMAN_1857_1916  --DKRLR-TTILPEQLEILYRWYMQSDSNPTRKMLDCLSEEVGLKRRVVQVWFQNTRARER
ZFHX2_HUMAN_2065_2124  --QRRYR-TQMSLQKIMKACYEAYRTPTMQEVLGEEIGLPRKRVVQVWFQNARAKEK
ZFHX3_HUMAN_2944_3003  --QRRYR-TQMSLQKIMKACYEAYRTPTMQEVLGEEIGLPRKRVVQVWFQNARAKEK
ZFHX4_HUMAN_2884_2943  --HKRFR-TQMSNLQKLVKACFSYRTPTMQEVLGEEIGLPRKRVVQVWFQNARAKEK
LHX1A_HUMAN_195_254    --PKRPR-TILTQORRAFKASFEVSSKPCRKRVRETAAETGLSVRVVQVWFQNRRAKMK
LHX1B_HUMAN_219_278    --PKRPR-TILTQORRAFKASFEVSSKPCRKRVRETAAETGLSVRVVQVWFQNRRAKMK
LHX1_HUMAN_180_239     --RRGPR-TTIAKQLETLKAAFAATPKPTRHIREQLAQETGLNMRVIQVWFQNRRAKMK
LHX5_HUMAN_180_239     --RRGPR-TTIAKQLETLKAAFAATPKPTRHIREQLAQETGLNMRVIQVWFQNRRAKMK
LHX3_HUMAN_157_216     --AKRPR-TTITAKQLETLKSAINTSPKPARHVREQLSSETGLDMRVVQVWFQNRRAKMK
LHX4_HUMAN_157_216     --AKRPR-TTITAKQLETLKSAINTSPKPARHVREQLSSETGLDMRVVQVWFQNRRAKMK
                        :      :      :      :      :
HOMEZ_HUMAN_451_510    EVV---
ZHX3_HUMAN_835_894     RAV---
ZHX1_HUMAN_777_832     LGIELF
HOMEZ_HUMAN_55_114     ISW---
ZHX2_HUMAN_263_324     ISWSPE
ZHX3_HUMAN_304_363     ISW---
ZHX1_HUMAN_284_346     VSWTPE
ZEB2_HUMAN_644_703     SNS---
ZEB1_HUMAN_581_640     SVQ---
NANGN_HUMAN_102_161    KEM---
ZHX1_HUMAN_569_630     LKEEKM
ZHX2_HUMAN_530_591     SMEQAV
ZHX3_HUMAN_612_671     AEE---
ZHX2_HUMAN_439_501     RGIVHI
ZHX3_HUMAN_494_553     NLK---
ZHX1_HUMAN_464_526     NSKSNQ
HOMEZ_HUMAN_355_415    HGQ---

```


Jalview can be easily installed under all commonly used operating systems and run locally. For these exercises, I attempt to use services available freely from the **INTERNET** wherever possible, so let us run **Jalview** from the web here by first going to:

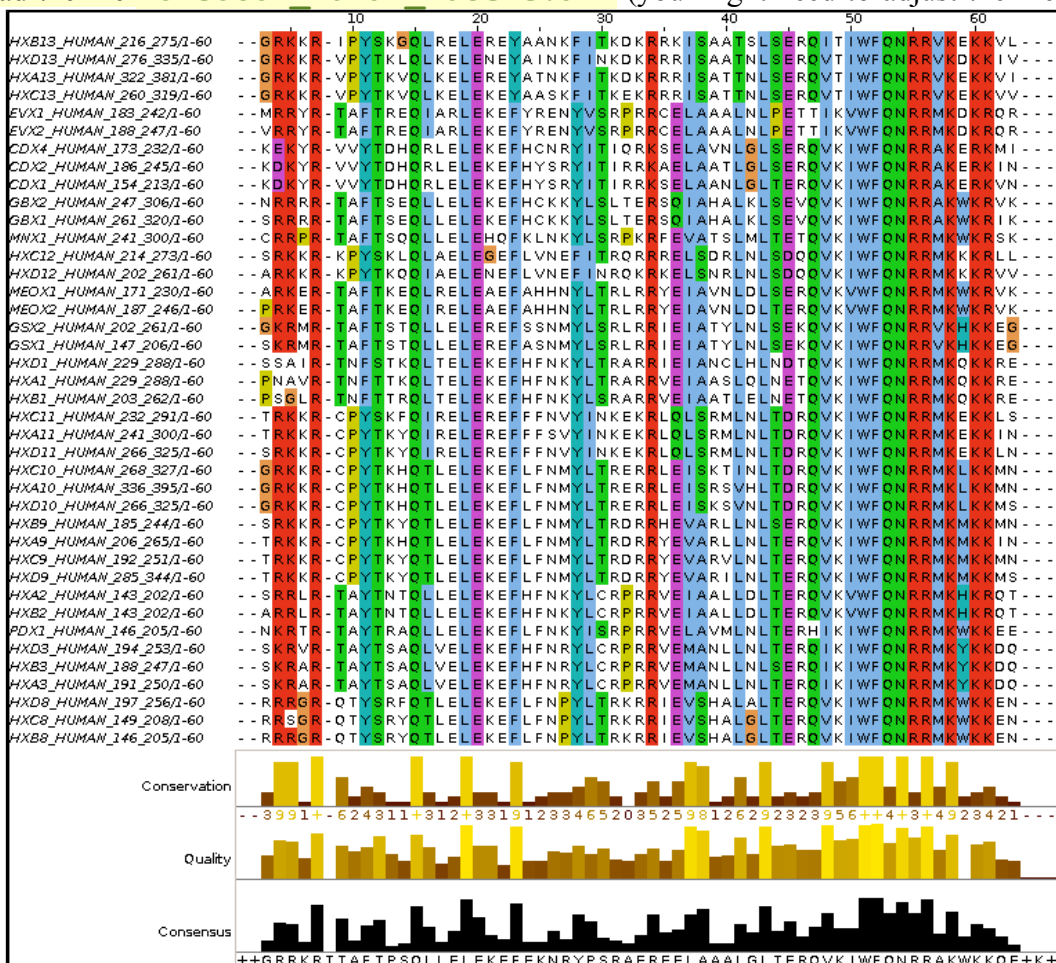
<http://www.jalview.org/>

and selecting the **Launch Jalview Desktop** link at the top of the page. Close down all the example outputs **Jalview** sees fit to show you on start up. From the **File** pull down menu choose **from File** from the **Input Alignment** option. Locate and load the file **homeobox_human_muscle.aln** (you might need to adjust the file name filter to include **.aln** files).

The default view is a trifle bland. Try a few of the options from the **Colour** pull down menu.

You could try the default colour scheme used by **ClustalX**, for example.

Now the **MUSCLE** and massaged **ClustalX** alignments look even more similar! In the nicely aligned regions at least.



There are many **Jalview** features that merit investigation. Have a look around if you have time. In particular, **Jalview** will compute simple phylogenetic trees for you employing a number of methods (**Calculate Tree** from the **Calculate** pull down menu). Try it, but be aware this is only sensible if you were very sure of your alignment (and have a few less and more meaningfully selected sequences maybe?).

Jalview is made by the same group as produce **Jpred** (an extremely effective **Secondary Structure Prediction** system). You could send your alignment for **Secondary Structure Prediction** via the **Web Service** pull down menu, if you wished.

A very important purpose of **Jalview** is to allow users to edit alignments as well as just to view them. For example, hold down the **Shift** key, click and hold on any amino acid at the edge of a gap, slide left and right and see that you can introduce and/or alter the position of gaps. It is very important to be able to edit alignments generated by even the best of programs. As I hope has been made clear, the alignment algorithms are crude. If you know something about the sequences you are aligning it is very reasonable to suppose you can improve upon the computer's alignments. **Jalview** tries to make this possibility easy. Look through some of the other **Edit** pull down menu options, it does not matter how much you mangle your alignment, you can always make another one.

Finally, take a look at the **Jalview** "Manhattan Skyline" for the highly conserved **W** at position 51. This seems better quality than **clustalX** managed? I am not sure how one can make further comment without knowing what parameters were used. Is there really an improvement? If so, is it due to the improved algorithm or more appropriate choice of parameters? Impossible to tell as the parameters used for **MUSCLE** are not revealed.

In my alignment, the **W** at position **51** was at position **50**, according to **clustalx**. This slippage to the right is due to **MUSCLE** introducing an extra gap, inspired by just one sequence at position **8**. Is this sensible? No idea ... exactly when it might be good idea to investigate the effect of lighter gap penalties?

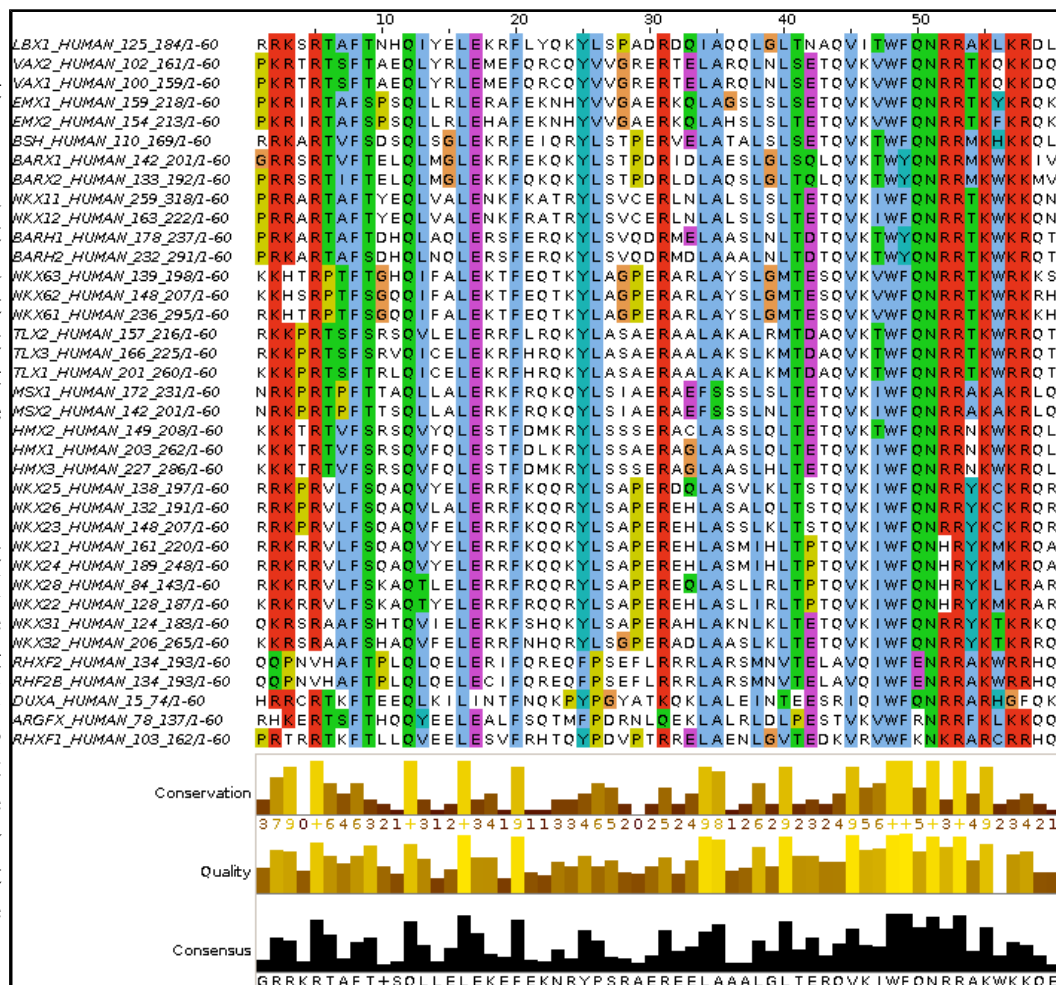
```

10
ZHX3_HUMAN - - SSIPT-YNAA
ZHX1_HUMAN - - NSIPT-YNAA
ZEB2_HUMAN - - GMTSP- INPY
ZEB1_HUMAN - - NLSPS-QPPL
NANGN_HUM - - QYPEK-RLVS
ZHX1_HUMAN - - POKFK- - EKT
ZHX2_HUMAN - - POKFK- - EKT
ZHX3_HUMAN - - TPTKY- KERA
ZHX2_HUMAN - - TPASD- RKKT
ZHX3_HUMAN - - ASIYK- NKKK
ZHX1_HUMAN - - SFGIR- AKKT
HOMEZ_HUM - - QRRKTKRKT
ZHX2_HUMAN - - SPSPA- IAKS
ZHX3_HUMAN - - PGKVS- CKKT
ZHX1_HUMAN - - STGKI- CKKT
ADNP_HUMANLDPKGHE- DDSY
ADNP2_HUMA - - PKKYE- GRSY

```

You can also **Select** and **Cut** sequences in a way similar to that you employed with **clustalx**. I could not resist it! I removed all the ugly sequences that caused the gaps at the start and finish of the alignment (just select their names and then select **Cut** or **Delete** from the **Edit** menu). I achieved the gap-free beautiful alignment illustrated.

Of course, **Jalview** does not compute alignments, so once I had removed all the unfortunate proteins, I had to use an **Edit** option to tidy up my meddling. I used **Remove Empty Columns** to get rid of the gap columns at the start of the alignment. The gaps at the end just melted away once the sequences that supported their presence were removed.



Science is easy! Once you remove the need for honesty that is.

If it could be done slightly more meaningfully, I would suggest you might try some of the other **MSA** tools offered by the **EBI**, to investigate the differences in the alignments computed. Any differences might be due to different parameter selection or differences in the algorithms of the tool you select.

For full control, you really need to download the various tools and run them locally. The **EBI** is not the only site that hides significant parameters from their users.

PSI-BLAST

This program is used to find a comprehensive set of relatives of a protein. First, **BLAST** is used to find closely related proteins. From an alignment of these proteins a general "profile" (a **Position Specific Scoring Matrix - PSSM**) is computed. A **PSSM** is very similar in concept and purpose to an **HMM** profile in that it summarises significant features present in the sequences it represents.

A further search of the protein database is then run using the **PSSM** as a query, and a larger more widely associated group of proteins is found. This larger group is aligned and used to construct another **PSSM**, and the process is repeated until no more significantly matching new sequences can be detected, or the user tires of the whole process.

PSI-BLAST is integrated into the **Secondary Structure Prediction** system **Jpred**. Whenever Jpred is asked to compute structure from a single protein sequence, it will use **PSI-BLAST** to construct an aligned family of protein sequences to enable an improved prediction. An aligned family of proteins is a much better starting point than any single protein sequence.

Similar ideas are used by the domain database **PFAM** to create large alignments of domain regions.

Here we will use **PSI-BLAST** directly from the **NCBI** on the **Paired DOMAIN** of the **PAX6** protein that you saved in a file earlier. It should be possible to detect a large family of **PAX** domains and to eventually multiply align them generating something like the alignment from the **PFAM** database.

To investigate **PSI-BLAST** go first to the **NCBI Home** page at:

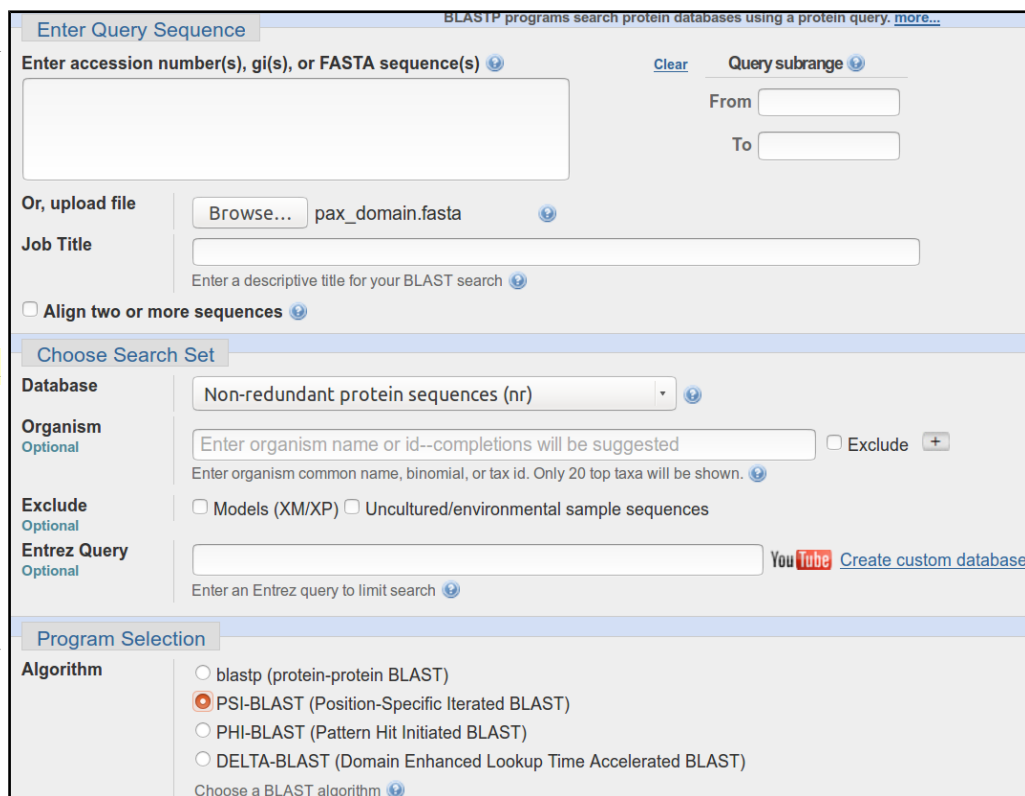
<http://www.ncbi.nlm.nih.gov/>


Click on the **BLAST** option from the **Popular Resources** menu.

Select  from the **Web BLAST** section.

Upload the **PAX6** paired box domain sequence (stored in the file **pax_domain.fasta**) using the appropriate **Browse** button.

Select **PSI-BLAST** from the **Program Selection** section. Leave all the others options at their default settings, particularly the option to search all the proteins available.

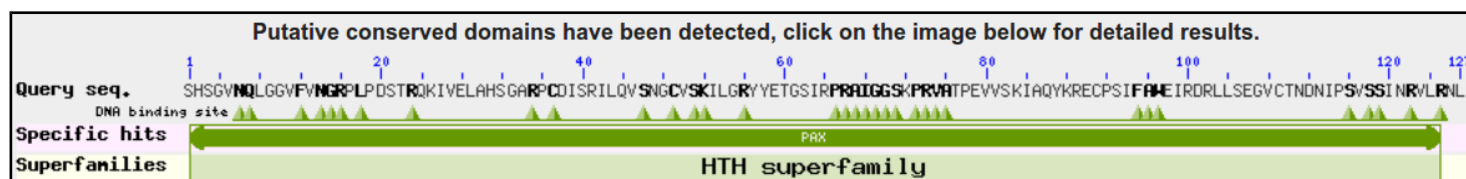


Before you set **PSI-BLAST** going, click on the **Algorithm parameters** link and take a look at the **PSI/PHI/DELTA BLAST** section. Note the option to use a **PSSM** from a previous run of **PSI-BLAST**, potentially on a different database (but with the same query sequence). Accept the default that database entries scoring better than an **Expect Threshold** of **0.005** be offered for inclusion into the **PSSM** of each successive **PSI-BLAST** iteration. Remember the  buttons.

What do you suppose the choice of **Pseudocount** might influence? _____

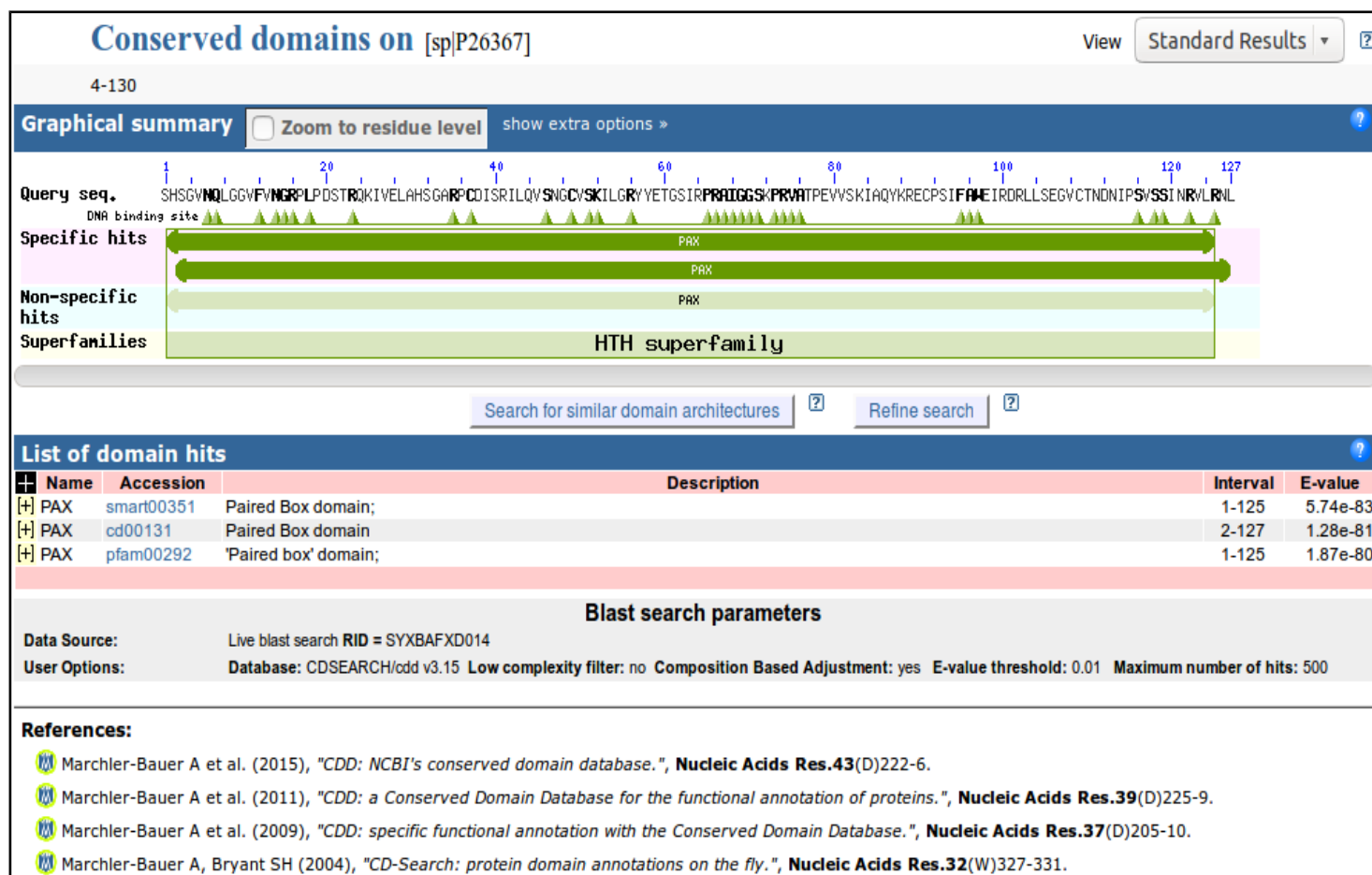
Elect to **Show results in a new window** and then click on the **BLAST** button.

After several moments of deep thought, **PSI-BLAST** will come back with its first set of results, at the top of which is a report that (unsurprisingly) matches have been detected between the query sequence and several domain databases. For more detail, click on the **Conserved Domains** graphic.

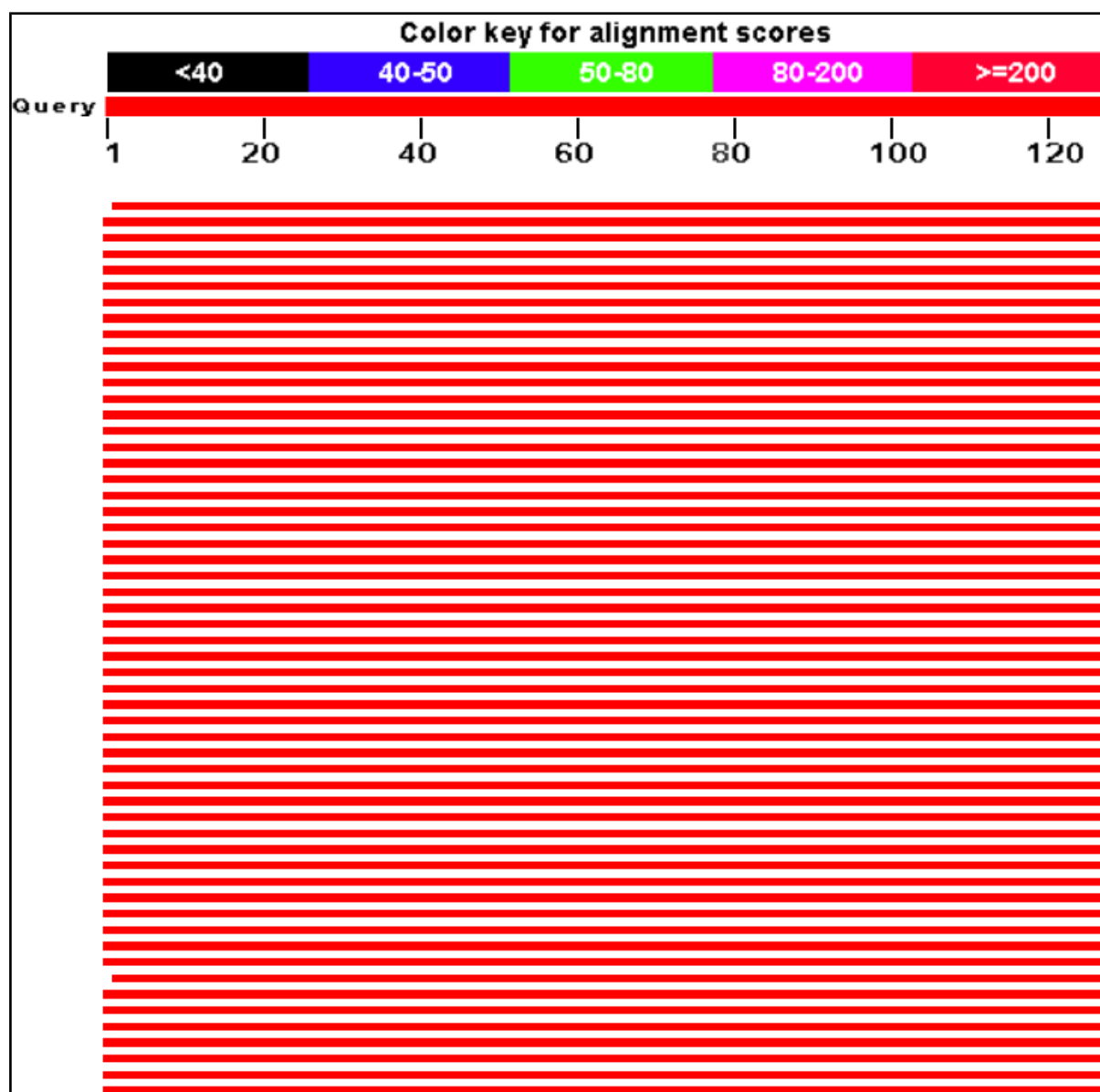


SMART, Pfam and the NCBI Conserved Domains database hits for a **PAX** domain are reported. No surprise here.

There is also a **Superfamilies** (derived from SCOP as briefly mentioned previously) hit recognising that a **PAX** domain, in common with many other domains, includes **Helix-Turn-Helices**.



Moving back to the main **PSI-BLAST** results, you will see that there are many high quality hits covering the whole length of the query sequence.



The best **500** of these are listed.

All the listed hits are selected for inclusion into the **PSSM** for the next iteration. Unless you feel strongly about any particular entry, **leave them all selected**.

Select: All None Selected:0									
Alignments Download GenPept Graphics Distance tree of results Multiple alignment									
	Description	Max score	Total score	Query cover	E value	Ident	Accession	Select for PSI blast	Used to build PSSM
<input type="checkbox"/>	hypothetical protein A6R68_04829 [Neotoma lepida]	257	257	99%	2e-86	100%	OBS66634.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X7 [Protobothrops mucrosquamatus]	262	262	100%	7e-86	100%	XP_015678414.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X4 [Macaca nemestrina]	262	262	100%	4e-85	100%	XP_011722295.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X4 [Macaca mulatta]	262	262	100%	4e-85	100%	XP_014969998.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X2 [Acinonyx jubatus]	263	263	100%	4e-85	100%	XP_014922398.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X4 [Macaca fascicularis]	262	262	100%	4e-85	100%	XP_015289636.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X2 [Ursus maritimus]	263	263	100%	4e-85	100%	XP_008685073.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X7 [Pseudopodoces humilis]	262	262	100%	4e-85	100%	XP_014114466.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	oculorhombin [Homo sapiens]	263	263	100%	4e-85	100%	AAA59962.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	paired box protein Pax-6 [Rattus norvegicus]	263	263	100%	4e-85	100%	NP_037133.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X2 [Fukomys damarensis]	263	263	100%	4e-85	100%	XP_010638711.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X2 [Cavia porcellus]	263	263	100%	4e-85	100%	XP_003464531.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X2 [Miniopterus natalensis]	263	263	100%	4e-85	100%	XP_016061621.1	<input checked="" type="checkbox"/>	
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 [Poecilia reticulata]	261	261	100%	4e-85	98%	XP_017159670.1	<input checked="" type="checkbox"/>	

Download ▾ GenPept Graphics

oculorhombin [Homo sapiens]
Sequence ID: [gb|AAA59962.1](#) Length: 422 Number of Matches: 1

Range 1: 4 to 130 [GenPept](#) [Graphics](#) ▾ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
263 bits(671)	4e-85	Compositional matrix adjust.	127/127(100%)	127/127(100%)	0/127(0%)
Query 1	SHSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYET				60
Sbjct 4	SHSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYET				63
Query 61	GSIRPRAIGGSKPRVATPEVVS KIAQYKRECP SIFAW EIRDRL LSEGVC TNDNIP SVSSI				120
Sbjct 64	GSIRPRAIGGSKPRVATPEVVS KIAQYKRECP SIFAW EIRDRL LSEGVC TNDNIP SVSSI				123
Query 121	NRVLRNL 127				
Sbjct 124	NRVLRNL 130				

Download ▾ GenPept Graphics

paired box protein Pax-6 [Rattus norvegicus]
Sequence ID: [ref|NP_037133.1](#) Length: 422 Number of Matches: 1
[▶ See 3 more title\(s\)](#)

Range 1: 4 to 130 [GenPept](#) [Graphics](#) ▾ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
263 bits(671)	4e-85	Compositional matrix adjust.	127/127(100%)	127/127(100%)	0/127(0%)
Query 1	SHSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYET				60
Sbjct 4	SHSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYET				63
Query 61	GSIRPRAIGGSKPRVATPEVVS KIAQYKRECP SIFAW EIRDRL LSEGVC TNDNIP SVSSI				120
Sbjct 64	GSIRPRAIGGSKPRVATPEVVS KIAQYKRECP SIFAW EIRDRL LSEGVC TNDNIP SVSSI				123
Query 121	NRVLRNL 127				
Sbjct 124	NRVLRNL 130				

Move down to the **Alignments** section of the results and you will see that many of the top hits match the query exactly.

Note that many of the top hits come from the **GenPept** database (roughly equivalent to the **TrEMBL** section of **UniProtKB**).

How might the inclusion of poor quality and duplicated sequences have been minimised?

Download ▾ GenPept Graphics

PREDICTED: paired box protein Pax-6-like isoform X1 [Polistes canadensis]
Sequence ID: [ref|XP_014604192.1](#) Length: 513 Number of Matches: 1

Range 1: 35 to 160 [GenPept](#) [Graphics](#) ▾ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
246 bits(627)	2e-77	Compositional matrix adjust.	118/126(94%)	121/126(96%)	0/126(0%)
Query 2	HSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYETG				61
Sbjct 35	HSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYETG				94
Query 62	SIRPRAIGGSKPRVATPEVVS KIAQYKRECP SIFAW EIRDRL LSEGVC TNDNIP SVSSIN				121
Sbjct 95	SI+PRAIGGSKPRVAT VV+KIA YKRECP SIFAW EIRDRL EGVC NDNIPSVSSIN				154
Query 122	RVLRNL 127				
Sbjct 155	RVLRNL 160				

Download ▾ GenPept Graphics

PREDICTED: paired box protein Pax-6 isoform X2 [Esox lucius]
Sequence ID: [ref|XP_010874535.1](#) Length: 461 Number of Matches: 1

Range 1: 32 to 174 [GenPept](#) [Graphics](#) ▾ Next Match ▲ Previous Match

Score	Expect	Method	Identities	Positives	Gaps
244 bits(623)	2e-77	Compositional matrix adjust.	124/143(87%)	125/143(87%)	17/143(11%)
Query 2	HSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYETG				47
Sbjct 32	HSGVNQLGGVFVNGRPLDPSTRQKIVELAHSGARPCDISRILQVSNCGVSKILGRYYETG				91
Query 48	GCVSKILGRYYETG SIRPRAIGGSKPRVATPEVVS KIAQYKRECP SIFAW EIRDRL LSEG				107
Sbjct 92	GCVSKILGRYYETG SIRPRAIGGSKPRVATPEVVS KIAQYKRECP SIFAW EIRDRL LSEG				151
Query 108	VCTNDNIPS - - - VSSINRVLRNL 127				
Sbjct 152	VCTNDNIPS VSSINRVLRNL 174				

Move down far enough and you will see less perfect matches, some of which involve proteins with the extra **14** amino acids of **isoform 5a** of **PAX6_HUMAN**.


Having browsed your results sufficiently, click on the **Go** button to **Run PSI-Blast iteration 2**. It is at the bottom of the hit list.

<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X1 [Pogonomymex barbatus]	254	254	99%	9e-81	94%	XP_011645781.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	Paired box protein Pax-6 [Papilio machaon]	252	252	99%	9e-81	93%	KPJ13446.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 [Megachile rotundata]	254	254	99%	9e-81	94%	XP_012145708.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 [Wasmannia auropunctata]	254	254	99%	9e-81	94%	XP_011694162.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 [Bombus terrestris]	254	254	99%	9e-81	94%	XP_003396039.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 [Orussus abietinus]	254	254	99%	9e-81	94%	XP_012279617.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X3 [Camponotus floridanus]	255	255	99%	9e-81	94%	XP_011266821.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X1 [Acromymex echinator]	254	254	99%	1e-80	94%	XP_011059463.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X4 [Camponotus floridanus]	254	254	99%	1e-80	94%	XP_011266829.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6-like isoform X3 [Polistes dominula]	254	254	99%	1e-80	94%	XP_015182130.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	Paired box protein Pax-6 [Fukomys damarensis]	255	255	100%	1e-80	100%	KFO26421.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X2 [Cynoglossus semilaevis]	252	252	100%	1e-80	90%	XP_008308818.1	<input checked="" type="checkbox"/>
<input type="checkbox"/>	PREDICTED: paired box protein Pax-6 isoform X4 [Haplochromis burtoni]	251	251	99%	1e-80	90%	XP_014195917.1	<input checked="" type="checkbox"/>

After a few moments, **PSI-BLAST** will return with the results of searching through the database again using the **PSSM** derived from the hits of the first iteration(☒ed). This time the top of the list will be predominantly filled with hits that have already been incorporated into the **PSI-BLAST PSSM**. However, look far enough down the list and you will find some new ones, highlighted yellow.

Once more, click on the **Go** button to **Run PSI-Blast iteration 3**. That is probably enough! It took **8** iterations before there were no more new sequences suggested for inclusion into the **PSMM** when I ran this last, so I advise that you stop here.

Next, move to the just above the **Graphic Summary** and click on the **Multiple alignment** link. You have elected to use the **NCBI** multiple alignment program **Cobalt** to align the **PAX** domain sequences of your final **PSI-BLAST** iteration.

sp P26367 4-130 (127 letters)	
RID	SZYU9GWN01R (Expires on 07-22 06:29 am)
Query ID	Id Query_250067
Description	sp P26367 4-130
Molecule type	amino acid
Query Length	127
Database Name	nr
Description	All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects
Program	BLASTP 2.4.0+ Citation
 No new sequences were found above the 0.005 threshold	
Other reports: Search Summary Taxonomy reports Distance tree of results Related Structures Multiple alignment	

Alignment Parameters	
Gap penalties	-11,-1
End-Gap penalties	-5,-1
CDD Parameters	
Use RPS BLAST	on
Blast E-value	0.003
Find Conserved columns and Recompute	on
Query Clustering Parameters	
Use query clusters	on
Word Size	4
Max cluster distance	0.8
Alphabet	Regular

This can take quite a while. **Cobalt** might even complain it is tired and wants to give up occasionally. If it does, tell it not to be silly!! It will get there in the end. When it does finish, click on the **Alignment parameters** link at the top of the results.

Cobalt reports the parameters it used to make the alignment. It is possible to recompute the alignment with different parameters by using the **Edit and Resubmit** link at the top of the page and then choosing to set **Advanced parameters**. But, maybe not today?

Move past the long list of proteins that have been aligned (the easiest way is to hide the **Descriptions** view).

At the top of the actual alignment, set **View Format to Plain Text** (.... and then hide the **Descriptions** again??), this being the easiest format to understand in a hurry. This might take a while also. I am not sure why? Be patient, it will get therein the end. The alignment will have very ragged ends, but the important region of **120** or so amino acids representing the **PAX** domain is really quite impressive. In particular, the **isoform 5a** insertion is very convincing.

<input checked="" type="checkbox"/>	KRY01441	56	IELAHQGARPCDISRIQL-----VSGNCVSKILC---RYYETGSIRPRAIGGSKPRVATNKVVDKIADYK	117
<input checked="" type="checkbox"/>	ALS19761	36	VELAHSGARPCDISRIQL-----VSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVKIAHYK	97
<input checked="" type="checkbox"/>	XP_011433286	66	VELAHSGARPCDISRIQL-----VSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVKIAQYK	127
<input checked="" type="checkbox"/>	KTG35220	49	VELAHSGARPCDISRIQLTHDDAKVQ-LDNKNVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVGRIAYK	123
<input checked="" type="checkbox"/>	BAA24024	47	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	122
<input checked="" type="checkbox"/>	XP_016061618	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	CAF29075	30	VELAHSGARPCDISRIQLTHADAKVQVLDSENVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_012307695	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_017402601	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_015678407	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_004264009	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_007885968	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_009184622	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_015193788	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVKIAQYK	105
<input checked="" type="checkbox"/>	XP_006975926	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_013814717	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_005530321	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	NP_038655	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	KRZ79048	109	IELAHQGARPCDISRIQL-----VSGNCVSKILC---RYYETGSIRPRAIGGSKPRVATNKVVDKIADYK	170
<input checked="" type="checkbox"/>	XP_006234695	30	VELAHSGARPCDISRIQLTHADAKVQVLDSENVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_010638709	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	CAC80516	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_005064878	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_004851663	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	KRZ09976	62	IELAHQGARPCDISRIQL-----VSGNCVSKILC---RYYETGSIRPRAIGGSKPRVATNKVVDKIADYK	123
<input checked="" type="checkbox"/>	XP_016339217	49	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVGKIAQYK	124
<input checked="" type="checkbox"/>	AAW24017	60	VELAHSGARPCDISRIQL-----VSGNCVSKILA---RYYETGSIRPRAIGGSKPRVATPEVVKIADYK	121
<input checked="" type="checkbox"/>	CBY09679	60	VELAHSGARPCDISRIQL-----VSGNCVSKILA---RYYETGSIRPRAIGGSKPRVATPEVVKIADYK	121
<input checked="" type="checkbox"/>	NP_001091013	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_006140475	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_006195131	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_005530320	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_005001333	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	NP_001595	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105
<input checked="" type="checkbox"/>	XP_014328414	31	VELAHSGARPCDISRIQLQSNCLIKDVCCLFVQVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVGKIAQYK	106
<input checked="" type="checkbox"/>	ALX18491	35	VELAHSGARPCDISRIQL-----VSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATLDVNNKIANYK	96
<input checked="" type="checkbox"/>	XP_014922397	30	VELAHSGARPCDISRIQLTHADAKVQVLDNQNSVSGNCVSKILG---RYYETGSIRPRAIGGSKPRVATPEVVSKEIAQYK	105

Model Answers to Questions in the Instructions Text.

Notes:

For the most part, these “**Model Answers**” just provide the reactions/solutions I hoped you would work out for yourselves. However, sometime I have tried to offer a bit more back ground and material for thought? Occasionally, I have rambled off into some rather self indulgent investigations that even I would not want to try and justify as pertinent to the objective of these exercises. I like to keep these meanders, as they help and entertain me, but I wish to warn you to only take regard of them if you are feeling particularly strong and have time to burn. Certainly not a good idea to indulge here during a time constrained course event!

Where things have got extreme, I am going to make two versions of the answer. One starting:

Summary:

Which has the answer with only a reasonably digestible volume of deep thought. Read this one.

The other will start:

Full Answer:

Beware of entering here! I do not hold back. Nothing complicated, but it will be long and full of pedantry.

This makes the Model answers section very big. **BUT**, it is not intended for printing or for reading serially, so I submit, being long and wordy does not matter. Feel free to disagree.

From your investigations of Multiple Sequence Alignment

How do the options for the **OUTPUT TREE** relate to the output files of **ClustalX** and the difference between the way that **ClustalX** and **muscle** work?

I leave this question here in the hope that one day I will be able to offer a full and sensible answer. First draft answer below.

Essentially, both **ClustalX** and **MUSCLE** work in two stages. First they create **Guide Tree(s)**. Then they create a multiple alignment by pairwise steps ordered by most refined the **Guide Tree**.

ClustalX just computes one based exclusively on the pairwise comparison of its input sequence set.

MUSCLE will create a **Guide Tree** that is the rough equivalent of that computed by **ClustalX**. Then it will offer to refine this **Guide Tree** from computed draft **MSAs** until a user selected maximum number of iterations is met or no further improvement is possible.

ClustalX saves the **Guide Tree** it computes by default. **MUSCLE** offers to save its **Guide Tree** from its first or second refinement iteration.

The purpose of saving the **Guide Tree(s)** to a file is to enable a rerun of the second phase with new parameter settings without having to first recalculate the **Guide Tree**. Of course, as mentioned previously, utterly pointless if there is no way to change the parameters to allow a guide tree to be used as input? but that is the theory.

More investigation by me and expansion of this answer required. Discussion with EBI current (2016.04.20).


Comment on how one might choose between the range of options offered for the aligned parameter?

I cannot ... beyond suggesting it simply does not make sense? Going by what is offered at **Wageningen**, the choice should be between **aligned** and **input order**. i.e. the order of the original set of sequences to be aligned or the order after they have all been compared with each other and arranged into a **Guide Tree** ... or two.

From your investigations of **PSI-Blast**

What do you suppose the choice of **Pseudocount** might influence?

I clicked with confidences upon the link to the help. It opined as illustrated.

Pseudocount	<input type="text" value="0"/>	
Pseduocount parameter. If zero is specified, then the parameter is automatically determined through a minimum length description principle (PMID 19088134). A value of 30 is suggested in order to obtain the approximate behavior before the minimum length principle was implemented.		

I suppose the next step is to read **PMID 19088134**? There is most certainly no elucidation amongst the strangle of words offered here?

The article **Abstract** says:

“Position specific score matrices (PSSMs) are derived from multiple sequence alignments to aid in the recognition of distant protein sequence relationships. The **PSI-BLAST** protein database search program derives the column scores of its **PSSMs** with the aid of **pseudocounts**, added to the observed amino acid counts in a multiple alignment column. In the absence of theory, the number of **pseudocounts** used has been a completely empirical parameter. This article argues that the minimum description length principle can motivate the choice of this parameter. Specifically, for realistic alignments, the principle supports the practice of using a number of **pseudocounts** essentially independent of alignment size. However, it also implies that more highly conserved columns should use fewer **pseudocounts**, increasing the inter-column contrast of the implied **PSSMs**. A new method for calculating **pseudocounts** that significantly improves **PSI-BLAST**'s; retrieval accuracy is now employed by default.”

The article itself, continues in like vein how about we close our eyes and accept the defaults? I would just wonder why the whole thing does not commence with, at least an attempt, to answer the question in the forefront of my inquiry, which is .. “**WHAT, in the current context, IS a pseudocount?**”. I do not believe it is as tricky as they appear to wish us to believe. I will try again later, when my view of the world is less storm infested.

In the meantime I will take comfort in the claim that:

“A new method for calculating **pseudocounts** that significantly improves **PSI-BLAST**'s; retrieval accuracy is now employed by default.”

Jolly good!

How might the inclusion of poor quality and duplicated sequences have been minimised?

At the top of your output is recorded some details of the conditions under which you database search was undertaken. This is a very important step towards making your results reproducible. Not sufficient I would opine. Surely the database versions and a complete record of the parameters used by **blast** are required in order to be able to exactly reproduce a search?

Database Name	nr
Description	All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects
Program	BLASTP 2.4.0+ ▶ Citation

But at least the version of **blast** and the databases that were searched are recorded. The collection of databases searched is rather optimistically called “**nr**”, for non-redundant. A bit of an exaggeration I would think. Surely **PDB** and **SwissProt** overlap a trifle? But let us not be too picky.

The collection of databases that is **nr** includes “*All non-redundant **GenBank CDS translations***” (aka **GenPept**) which, like it European broad equivalent **TrEMBL**, includes some pretty dubious sequences.

I would think that if one wanted to maximise quality and minimise duplication, it would be best to pick just one good quality database. **SwissProt** is the obvious choice. **blast**, in general, and **PSI-BLAST** in particular, allows such a selection.

However, today the objective is not refinement!!! Bloat is good! More the merrier! Never mind the quality, just admire the volume.

DPJ – 2016.07.21