



GTPB

The Gulbenkian Training Programme in Bioinformatics
(Since 1999)

Pedro Fernandes, Organiser



ELB19F

Entry Level Bioinformatics

04-08 February 2019

(First 2019 run of this Course)

Basic Bioinformatics Sessions

Practical 2: Pairwise Sequence Alignment

Wednesday 30 January 2019

Sensitive Pairwise Alignment

The purpose of this exercise is to look at some aspects of **Pairwise Sequence Alignment** using the most accurate methods available.

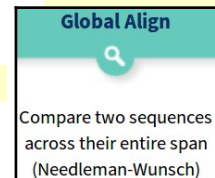
As hopefully has been discussed, sequences can be aligned using a **global** strategy, in which the two sequences being aligned are assumed to be homologous from end to end, or using a **local** approach, in which the sequences are assumed to just have homologous region(s).

Global Pairwise Sequence Comparison

First the **global** approach. In a previous exercise, you already have used the **blast** facility at the **NCBI** to perform crude pairwise alignment. **blast** also offers a sensitive option, so maybe that would be a good place to start.

So, once more to the **NCBI** home page (<http://www.ncbi.nlm.nih.gov/>). From there chose **BLAST** from the

Popular Resources list. Scroll down to the **Specialized searches** section and chose the



option.

A choice of settings for **Nucleotide** or **Protein** alignment is offered. As we are going to investigate the alignment of DNA sequences, the default choice is fine. For the first sequence, browse for the file **pax6_genomic.fasta**, which you created when looking at **Ensembl**. It contains the region of **Chromosome 11** containing the entire **PAX6** gene (with a few extra base pairs either end).

To specify the second sequence, you could load the file **pax6_mrna.fasta**, but just typing the corresponding **Accession** code in the appropriate box seems far more sophisticated, so that is what I chose to do.

Open the **Algorithmic Parameters** section, and see that they are as one might expect. The defaults are fine here as the alignment to be computed is trivial (given the way **blast** will go about the task), so anything not outrageous should work.

Ask to **Show results in a new window** and then click on the **Align** button.

After some significant Rollin' and Tumblin' **blast** will proclaim its lyrical conclusions. First examine the **Dot Matrix View**. This sort of representation has rather gone out of fashion in recent years. A shame, I say, this picture represents such a succinct summary of what should be expected of the textual alignment(s) that are the "real" detailed output of this sort of program.

How would you interpret this picture?

What do the diagonal(ish) lines represent?

What are the gaps in between the lines?

Which axis represents the genomic sequence and which the mrna?



Move down to the textual alignment. There are some weird little bits and pieces at the front of the alignment which defy logic. I decide not to dwell on these too much, beyond noting that the mRNA has some odd bases at the front.

Also, I have faith that the alignment you look at yields the highest alignment score, but equally. I doubt most **people** would have chosen to throw these odd bases about with quite such abandon! **People** are best!

You can just see evidence of the little patches of whimsy in the **Dot Matrix View**.

Query	661	CGCTGGCTGGATATTAAGGAAAGTTAGCGCTGCCTGAGCACCTCTTTCTTATCATT	720
Sbjct	1	-----TATC---	4
Query	721	GACATTTAAACTCTGGGGCAGGTCCTCGCGTAGAACGCGGCTGCAGATCTGCCACTTCC	780
Sbjct		-----	
Query	781	CCTGCCGAGCGGC6GTGAGAAAGTGTGGGAACCGGCGCTGCCAGGCTCACCTGCCTCCCG	840
Sbjct		-----	
Query	841	CCCTCCGCTCCCAAGTAACCGCCCGGGCTCCGGCCCCGGCCGGCTCGGGGCCGCGGGG	900
Sbjct		-----	
Query	901	CCTCTCCGCTGCCAGCAGCTGTGTCCCAAATCAAAGCCGCCCAAGTGGCCCCGGGG	960
Sbjct		-----	
Query	961	CTTGATTTTTGCTTTTAAAGGAGGCATACAAAGATGGAAGCGAGTTACTGAGGGAGGGA	1020
Sbjct	5	-----GA	6
Query	1021	TAGGAAGGGGGTGGAGGAGGACTTGTCTTTGCCGAGTGTCTCTTCTGCAAAAGTAGC	1080
Sbjct	7	TA-----	8

Moving down there are a series of far more convincing near perfect alignments.

You must know what these aligned regions represent by now?

But, just in case:

What do you suppose these regions represent?

How many are there and do they correspond nicely to the lines of the **Dot Matrix View**?

How many exons would you say this mrna has?

If one was to forgive the strange “bits” at the start, would you say **blast** seems to have done a reasonable job here?

I think I would.

The final alignment section even has a **PolyA Tail**!

Or does it? How you you interpret the run of **As** at the end of the final exon?

Query	24541	TCTTTCAGAGTTTGAGAGAACCATTATCCAGATGTGTTTGCCCGAGAAAGACTAGCAGC	24600
Sbjct	1045	-----AGTTTGAGAGAACCATTATCCAGATGTGTTTGCCCGAGAAAGACTAGCAGC	1096
Query	24601	CAAAATAGATCTACCTGAAGCAAGAAATACAGGTACCGAGAGACTGTGCAAGTTTACACTT	24660
Sbjct	1097	CAAAATAGATCTACCTGAAGCAAGAAATACAGGTA-----	1130
Query	24661	TGTGATTACATACATTTTCTTTCTTAGAGACAGAGGTGCTTGTACAGAGTACTATTTAT	24720
Sbjct		-----	
Query	24721	TTATAGGACTAATAATAAAAAAGGTTCAAGTCTGCTAAATGCTCTGCTGCCATGGGCGTG	24780
Sbjct		-----	
Query	24781	GGGAGGGCAGCAGTGGAGGTGCCAAGGTGGGGCTGGGCTCGACGTAGACACAGTGCTAAC	24840
Sbjct		-----	
Query	24841	CTGTCCACCTGATTTCCAGGTATGGTTTTCTAATCGAAGGGCCAAATGGAGAAGAGAAG	24900
Sbjct	1131	-----TGGTTTTCTAATCGAAGGGCCAAATGGAGAAGAGAAG	1167
Query	24901	AAAACTGAGGAATCAGAGAAGACAGGCCAGCAACACCTAGTCATATTTCTATCAGCA	24960
Sbjct	1168	AAAACTGAGGAATCAGAGAAGACAGGCCAGCAACACCTAGTCATATTTCTATCAGCA	1227
Query	24961	GTAGTTTCAGCACCAGTGTCTACCAACCAATTCACAACCCACACACCGGGTAATTTGA	25020
Sbjct	1228	GTAGTTTCAGCACCAGTGTCTACCAACCAATTCACAACCCACACACCGG-----	1278
Query	25021	AATACTAATACTACGAATCAATGTCTTAAACCTGTTTGTCTCGGGCTCTGACTCTCACT	25080
Sbjct		-----	
Query	25081	CTGACTACTGTCAATTTCTTTCCTCAGTTTCTCTCTCACATCTGGCTCCATGTTGGG	25140
Sbjct	1279	-----TTCTCTCTCACATCTGGCTCCATGTTGGG	1309
Query	25141	CCGAACAGACACAGCCCTCACAAACACCTACAGCGCTCTGCCCGCTATGCCAGCTTAC	25200
Sbjct	1310	CCTAACAGACACAGCCCTCACAAACACCTACAGCGCTCTGCCCGCTATGCCAGCTTAC	1369
Query	25201	CATGGCAAATAACCTGCCTATGCAAGTAAGTGGCGGCTGGTGGCTGCATAACCCAGG	25260
Sbjct	1370	CATGGCAAATAACCTGCCTATGCAA-----	1394
Query	25261	CCCCAGAGAAGTGAGAGTGGCTCAGGGCTGCGGACCTATTGGCTGTGTCTGCACCCT	25320
Sbjct		-----	
Query	25321	TGAGAGCTTTTGCCTACAGTATTGGCTTGACCAAGTCAAGTCGGAGACAGTCAATCCC	25380
Sbjct		-----	

Query	28561	TTTTGTAAACCTATAAATTTGATTCCATGTCTGTTTCTCAAGGGAATATCTACATGG	28620
Sbjct		-----	
Query	28621	CTATTTCTTTTCATCCACTTCTAGGACTCATTTCCTGGTGTGTCAGTTCAGTTCAGT	28680
Sbjct	1547	-----ACTCATTTCCCTGGTGTGTCAGTTCAGTTCAGT	1582
Query	28681	TCCCGGAAGTGAACCTGATATGTCTCAATACTGGCCAAGATTACAGTAAAAAAAAAAAA	28740
Sbjct	1583	TCCCGGAAGTGAACCTGATATGTCTCAATACTGGCCAAGATTACAGTAAAAAAAAAAAA	1642
Query	28741	AAAAAAAAAAGGAAAGAAATATTGTGTTAATTCAAGTCAGTCACTATGGGGACACAACAG	28800
Sbjct	1643	A-----	1643
Query	28801	TTGAGCTTTCAGGAAAGAAAGAAATGGCTGTTAGAGCCGCTTCAAGTCTACAATTGTG	28860

Wonderful, but it is not safe to assume that just selecting any service that claims to do a sensitive global pairwise alignment will just work for any pair of sequences. In fact, pretty though it appears, the alignment **blast** has generated is not as entirely logical as it might first seem. For example, consider:

How might the gap around **24,750** in the genomic sequence been positioned more intelligently?

Next, try aligning the same two sequences with another program (implementing the same algorithm) at the **EBI**.

Global Alignment
 Global alignment tools create an end-to-end alignment of the sequences to be aligned. There are separate forms for protein or nucleotide sequences.

Needle (EMBOSS)
 EMBOSS Needle creates an optimal global alignment of two sequences using the Needleman-Wunsch algorithm.

Protein Nucleotide

Go to the **Pairwise Sequence Alignment EBI** page:

(<http://www.ebi.ac.uk/Tools/psa/>).

Select the **Nucleotide** option for the **Global Alignment** program **Needle**.

Needle implements the best global pairwise algorithm exactly.

Load up the first sequence from **pax6_genomic.fasta**.

Load up the second sequence from **pax6_mrna.fasta**.

Click on the **More options** button to see what parameters you can set. They should be as you might expect. The defaults are fine for the first run.

Click on the **Submit** button to get **Needle** into action.

Pairwise Sequence Alignment (NUCLEOTIDE)
 EMBOSS Needle reads two input sequences and writes their optimal global sequence alignment to file.

This is the form for nucleotide sequences. Please go to the **protein** form if you wish to align protein sequences.

STEP 1 - Enter your nucleotide sequences

Enter or paste your first nucleotide sequence in any supported format:

Or, upload a file: pax6_genomic.fasta

AND

Enter or paste your second nucleotide sequence in any supported format:

Or, upload a file: pax6_mrna.fasta

STEP 2 - Set your pairwise alignment options

MATRIX: DNAfull | GAP OPEN: 10 | GAP EXTEND: 0.5 | OUTPUT FORMAT: pair

END GAP PENALTY: false | END GAP OPEN: 10 | END GAP EXTEND: 0.5

STEP 3 - Submit your job

Be notified by email (Tick this box if you want to be notified by email when the results are available)

```

pax6_genomic 22301 TTCTTGTAAACACAATGTGGCCCGCTGCACGCCCTCAAGAGAATC-----C 22344
M77844.1      1  -----TATCGATAAGT 11
pax6_genomic 22345 TTTTGTGTCCCGCCTCATTGTAGCCTCAAAAT-TCTGCCACGAAAGTT 22393
M77844.1     12  TTTT-----TTATGT-----CAATCTCTG----- 34
pax6_genomic 22394 TGCCAACGCTCCTGCCAGGAGTTTAATAGTTTCCCTTACTCGGGGC 22443
M77844.1     35  -----TCTCCT-TCCCAGGAATCTGAGGATTGCTTACACAC----- 71
pax6_genomic 22444 ATTGTGCAGCGCTGAAAAGCAGCCCTCGCTATTCAAGTGTGGTGTCA 22493
M77844.1     72  -----CAACCCAGCAA-CATCC-----GTGGAGA 94
pax6_genomic 22494 ---TCTCAATAG-ATCTCCAAGGGCCCATATGGTGCCAGTCCGATGA 22538
M77844.1     95  AAACCTTCACCGCACTCC----- 114
pax6_genomic 22539 ATCCGCTGTTTAAATGGGGAGAAAGTTGGGTTTTAAACAT----- 22582
M77844.1    115  -----TTTAAA-----ACACGT--CATTCAAACCAATTGTGGT 146
pax6_genomic 22583 -TTCAA-----AGTTCTGAAAAGATCCC-----ACT----- 22608
M77844.1    147  CTTCAAGCAACAACAGCAGCACAAAACCCCAACCAACAAAACCTCTTG 196
  
```

Well! Nothing like as convincing as the alignment **blast** produced!

Alignment does not even begin until over **22,300** base pairs along the genomic sequence. Even then it is not convincing, as in **wrong**, if we accept the results already obtained from **blast** as a fair approximation of the truth.

```

pax6_genomic 23746 TTACCTTGGGAATGTTTGGTGA--GGCTGTCCGGATATAATGCTCTTG 23792
M77844.1      804  -----ATGT-----TGAACGGCGAGCCGG-----AAGC---TG 829
pax6_genomic 23793 GAGTTTAAGACTACACCGCCCT-TTTGGAGGCTCCAAGTTAATCC-- 23839
M77844.1     830  GGG-----CACCG--CCCTGGTTGG-----TATCCGG 855
pax6_genomic 23840 AAATTTCTCTTAC---CATCTATTCTTTTGTTCAGATGGCTGCCAG 23885
M77844.1     856  GGACTTCGGTGCCAGGGCAACCTA-----CGCAAGATGGCTGCCAG 896
pax6_genomic 23886 CAACAGGAAGGAGGAGAGAATAACCACTCATCAGTTCCAACGGAGA 23935
M77844.1     897  CAACAGGAAGGAGGAGAGAATAACCACTCATCAGTTCCAACGGAGA 946
pax6_genomic 23936 AGATTCAGATGAGGCTCAAATGCGACTTCAGCTGAAGCGGAAGCTGCAA 23985
M77844.1     947  AGATTCAGATGAGGCTCAAATGCGACTTCAGCTGAAGCGGAAGCTGCAA 996
pax6_genomic 23986 GAAATAGAACATCCTTTACCAAGAGCAAATGAGGCCCTGGAGAAAGTT 24035
M77844.1     997  GAAATAGAACATCCTTTACCAAGAGCAAATGAGGCCCTGGAGAA---- 1042
pax6_genomic 24036 GATAGAGTTTTTCAAAGTAGAGAAGCAGTAATCAAAGTAAATGCCACAT 24085
M77844.1    1043  ----- 1042
pax6_genomic 24086 CTTCAAGTAAAGAGCTAAATTTAGCCAGGCCCTTGCATAGAAGAAATG 24135
M77844.1    1043  ----- 1042
  
```

There are some well aligned regions after genomic position **24,500**.

```

pax6_genomic 24836 CTAACCTGTCCACCTGATTTCCAGGTATGGTTTCTAATCGAAGGCCA 24885
M77844.1     1125  -----CAGGTATGGTTTCTAATCGAAGGCCA 1152
pax6_genomic 24886 AATGGAGAAGAGAGAAAACCTGAGGAATCAGAGAAGACAGGCAGCAAC 24935
M77844.1     1153  AATGGAGAAGAGAGAAAACCTGAGGAATCAGAGAAGACAGGCAGCAAC 1202
pax6_genomic 24936 ACACCTAGTCATATTCCTATCAGCAGTAGTTTCAGCACCAGTGTCTACCA 24985
M77844.1     1203  ACACCTAGTCATATTCCTATCAGCAGTAGTTTCAGCACCAGTGTCTACCA 1252
pax6_genomic 24986 ACCAATTCCACAACCCACCACCGGTAATTTGAAATACTAATACTACG 25035
M77844.1     1253  ACCAATTCCACAACCCACCACCG----- 1277
pax6_genomic 25036 AATCAATGTCTTAAACCTGTTTGTCTCCGGGCTCTGACTCTCACTGTAC 25085
M77844.1     1278  ----- 1277
pax6_genomic 25086 TACTGTCAATTTCTTCTGGCCCTCAGTTTCTCTTCCACATCTGGCTCCATG 25135
M77844.1     1278  -----GTTTCTCTTCCACATCTGGCTCCATG 1304
pax6_genomic 25136 TTGGGCCGAACAGACACAGCCCTCAAAAACCTACAGCGCTCTGCGGCC 25185
M77844.1     1305  TTGGGCCGAACAGACACAGCCCTCAAAAACCTACAGCGCTCTGCGGCC 1354
pax6_genomic 25186 TATGCCAGCTTCCACATGGCAAATAACCTGCCTATGCAAGTAAAGTCGG 25235
M77844.1     1355  TATGCCAGCTTCCACATGGCAAATAACCTGCCTATGCAAGTAAAGTCGG 1394
pax6_genomic 25236 CTGGTGGTGGCTGCATAACCCAGGCCAG--AGAAGTGAGGAGTGGCT 25283
M77844.1     1395  -----CC-----CCAGTCCCAGCCAG----- 1413
pax6_genomic 25284 CAGGCCCTCGGACCTCAT-----TGGCTGTGCTG--CACCTTGAGAG 25326
M77844.1     1414  -----CCT-----CCTCATACTCTGCTG--CTGCCACC-----AG 1444
  
```

Then a resumption of chaos after **25,230** or so.

How many convincingly aligned regions did you see?

How many did you expect?

Clearly, this alignment is not correct. Can you explain why?

I assume you have all read the lucid answers to the question above? If so, I am confident you will agree that there are **3** ways to get an answer, similar to that generated by **blast**, from the tools offered at the **EBI**. They are:

- Make gap penalties so cheap that **Needle** will have no excuse to avoid gaps where they are needed. This works if you use a gap opening penalty of **1.0** (the lowest allowed by the web interface) and a gap extension penalty of **0.0**, allowed by the program *but not by the EBI web interface!!* The lowest value the web interface allows is **0.0005**, which really should be sufficiently small, but provably is not. The most important question being “*Why would a web interface restrict a program's capabilities other than to prevent excessive resource use?*”. I have no answer for that one, I will just petulantly include some extra low gap alignments (made without a web interface) in your **Backup_Results** directory and retire with self righteous hauteur! Note that making gaps completely free (i.e. both gap **opening** and **extension** equal to **0.0**) will not work at all! **needle** would simply match each base of the mRNA with the next identical base of the genomic sequence until it runs out of letters. You could do this from the command line, but it would clearly not make sense.

Actually, using gap penalties to suit huge gaps that are really introns, will only work when the exons are so similar (as here) that any gap penalties will work for their alignment. Generally, you need to pick gap penalties to optimise exon alignment. So this is a very horrible way to “fix” the situation anyway.

- Tell **Needle** to penalise the gaps it puts at either end of the alignment in the same way it penalises gaps it puts in the middle. By default, end gaps are free!! Which is not very logical here. This is possible using the website.
- Use **Stretcher**, which uses essentially the same algorithm as **Needle**, except, it also applies a bit of common sense (**heuristics**, if you like). **Stretcher** takes a look at the sequences before it starts to do any serious computation. It identifies any “*good regions*” (all **12** exon matches in this case) and then says “*OK, I am definitely having those, how best can I deal with the rest?*”. In essence, **Stretcher** does a quick **Dot Matrix View** before it starts and so only goes to work when it has a pretty good idea what the answer should look like. It works in this case, but not always. **Stretcher** is faster than **Needle** but does not necessarily generate the highest scoring alignment. **Stretcher** works in a fashion far closer to the way a human would work, which has to be good! Well, usually anyway.

So, try the **Needle** with penalised **End Gaps** approach by returning to the **Needle** launch page from your results. You should find the two sequences are still selected, so you should only have to click on **More Options** again and change the **END GAP PENALTY** field from **false** to **true**.

STEP 2 - Set your pairwise alignment options			
MATRIX	GAP OPEN	GAP EXTEND	OUTPUT FORMAT
DNAfull	10	0.5	pair
END GAP PENALTY	END GAP OPEN	END GAP EXTEND	
true	10	0.5	

Click on the **Submit** button and **Needle** will be on the road again.

How many matching regions are there this time? Is the count **now** roughly as you would expect?

Stretcher (EMBOSS)

EMBOSS Stretcher uses a modification of the Needleman-Wunsch algorithm that allows larger sequences to be globally aligned.

[Protein](#) [Nucleotide](#)

Finally, check that **Stretcher** works as expected.

Go again to the **Pairwise Sequence Alignment EBI** page (<http://www.ebi.ac.uk/Tools/psa/>).

From there, select the **Nucleotide** option for the **Global Alignment** program **Stretcher**.

Load up the sequences exactly as for **Needle**.

Take a look at the parameters and see there is nothing unexpected hiding there.

Set **Stretcher** sequence rope stretching.

How do you feel about the results this time?

How do you think **blast** achieve the correct results without any fuss?

Tools > Pairwise Sequence Alignment > EMBOSS Stretcher

Pairwise Sequence Alignment (NUCLEOTIDE)

EMBOSS Stretcher calculates an optimal global alignment of two sequences using a modification of the classic dynamic programming algorithm which uses linear space.

This is the form for nucleotide sequences. Please go to the protein form if you wish to align protein sequences.

STEP 1 - Enter your nucleotide sequences

Enter or paste your first nucleotide sequence in any supported format:

Or, upload a file: pax6_genomic.fasta

AND

Enter or paste your second nucleotide sequence in any supported format:

Or, upload a file: pax6_mrna.fasta

STEP 2 - Set your pairwise alignment options

MATRIX	GAP OPEN	GAP EXTEND	OUTPUT FORMAT
DNAfull	16	4	pair

STEP 3 - Submit your job

Be notified by email (Tick this box if you want to be notified by email when the results are available)

Pairwise Sequence Comparison using Specialised Software

None of the alignments generated thus far have been entirely correct.

By persuading the general global alignment software to treat huge gaps (i.e. the introns) in some sort of special manner, a reasonable answer was obtained. However, the general software could not know that something more than just **Substitutions** and **Indels** were at issue here. Consequently, it stood no chance of dealing with the intron/exon boundaries sensibly.

The solution is not to fiddle around with the parameters of the general tools. Aligning **mRNAs** with **Genomic** sequence is simply not “*General Alignment*”. It is an example of a problem that is sufficiently particular to require specialised software for an optimal solution.

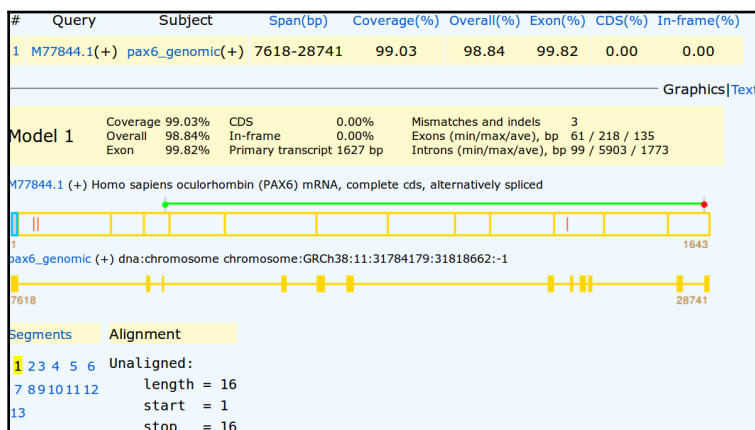
There is a program in the **EMBOSS** package (the same collection of programs as **Needle** and **Stretcher**), called **est2genome**, which is specifically designed for the alignment of cDNA/mRNA and genomic sequences. **est2genome** (and similar programs) may assume much more about the sequences to be aligned than can a general purpose alignment program. Gaps representing introns can be placed far more accurately if they are **known** to represent introns. Programs such as **est2genome** seek the highly conserved bases that occur at intron/exon boundaries, **C/T** rich intronic regions, **polyA** regions and **Stop/Start** codons to assist its detection of exons and gene structures.

est2genome is a fine program, but the option offered at the **NCBI** in America does the same job, I think, somewhat more nicely. The **NCBI** program is called **splign**. To investigate, go to the home of **splign** at:

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

Click on the **Online** button. In the **Genomic** section, **Browse** to upload **pax6_genomic.fasta**.

In the **cDNA** section, paste the sequence **pax6_mrna.fasta**. Where **cDNA** and **Genomic** sequences share exons that are nearly identical, **splign** uses the comparison algorithm **megablast** (default). Where exons are less similar (e.g. when the **cDNA** and **Genomic** sequences are from different organisms) the more sensitive option **discontinuous megablast**, is a better choice¹. Note the option to compare your **cDNA** with a **Whole genome** (including Human). Today, the default options are fine. Click the **Align** button.



Your results will appear showing the cDNA split into **12** sections (the predicted exons) corresponding to **12** regions of the genomic sequence indicated by yellow rectangles. A **13th** region of **16** base pairs is displayed and declared to be **unaligned**. These are the **16** mystery base pairs at the start of this particular mRNA that **Needle** and **Stretcher** had trouble treating sensibly also. I wonder what they are?

Any theories?

¹ Why this is so will be considered later when we look at the database searching program **blast**.

Click on the first exon section of the cDNA display.

Here there shows two **substitutions**. These were also apparent in the successful **blast**, **Needle** and **Stretcher** alignments. You might have spotted them?

Though these are in a non-coding region, they could easily still be very significant. However, for the purposes of this exercise, let us assume they are not.

The **Start** (green) and **Stop** (red) codons delimiting the **CoDing Sequence (CDS)** are illustrated by the bar above the cDNA display.

#	Query	Subject	Span(bp)	Coverage(%)	Overall(%)	Exon(%)	CDS(%)	In-frame(%)
1	M77844.1(+)	pax6_genomic(+)	7618-28741	99.03	98.84	99.82	0.00	0.00

Model 1	Coverage	CDS	Overall	In-frame	Mismatches and Indels	3
	99.03%	0.00%	98.84%	0.00%	Exons (min/max/ave), bp	61 / 218 / 135
	99.82%	0.00%	99.82%	0.00%	Introns (min/max/ave), bp	99 / 5903 / 1773

M77844.1 (+) Homo sapiens oculorhombin (PAX6) mRNA, complete cds, alternatively spliced

Segments Alignment

```

1 2 3 4 5 6
7 8 9 10 11 12
13
17 TTTTATTGTCAATCTCTGTCTCCTTCCCAGGAATCTGAGGATTGCTCTACACACCAACCCAGCAACATC
7618 TTTTATTGTCAATCTCTGTCTCCTTCCCAGGAATCTGAGATTGCTCTACACACCAACCCAGCAACATC
87 CGTGGAGAAAACCTCTCACCAGCAACTCCTTTAAACACCGTCATTTCAAACCATTGGTCTTCAAGCAA
7688 CGTGGAGAAAACCTCTCACCAGCAACTCCTTTAAACACCGTCATTTCAAACCATTGGTCTTCAAGCAA
157 CAACAGCAGCACAAAAACCCCAACCAAAACAAAACCTTTGACAGAAGCTGTGACAACCAAGAAAGGATGCC
7758 CAACAGCAGCACAAAAACCCCAACCAAAACAAAACCTTTGACAGAAGCTGTGACAACCAAGAAAGGATGCC
227 TCATAAAG.....
7828 TCATAAAGGTGAG
    
```

Click on the exon including the green **Start** codon (the **3rd**).

#	Query	Subject	Span(bp)	Coverage(%)	Overall(%)	Exon(%)	CDS(%)	In-frame(%)
1	M77844.1(+)	pax6_genomic(+)	7618-28741	99.03	98.84	99.82	0.00	0.00

Model 1	Coverage	CDS	Overall	In-frame	Mismatches and Indels	3
	99.03%	0.00%	98.84%	0.00%	Exons (min/max/ave), bp	61 / 218 / 135
	99.82%	0.00%	99.82%	0.00%	Introns (min/max/ave), bp	99 / 5903 / 1773

M77844.1 (+) Homo sapiens oculorhombin (PAX6) mRNA, complete cds, alternatively spliced

Segments Alignment

```

1 2 3 4 5 6
7 8 9 10 11 12
13
312 .....AGCCCATATTCGAGCCCGTGGAAATCCCGCGGCCCAACGAGCCAGCATGCAGAACA.....
12196 AACAGAGCCCATATTCGAGCCCGTGGAAATCCCGCGGCCCAACGAGCCAGCATGCAGAACAGTAA
373 .
12266 G
    
```

The first coding exon is now displayed with translation of the mRNA where appropriate.

The statistics at the top of the display include the claim that there are **3** discrepancies (**Mismatches** and **Indels**) between the **cDNA** and **Genomic** sequences.

Two of these are the **substitutions** we have already seen in the first exon of the cDNA. The third is indicated by the red bar in the **10th** exon of the cDNA display.

Click on the **10th** exon section of the cDNA display.

The third difference, a substitution, should be clear to see. Given it changes the coded protein, this substitution is likely to be the most significant.

Irritatingly, in the extreme! **splign** only translates the mRNA. So one has to work to discover the alternative suggested by the Genomic sequence.

Vital if we were really doing this seriously, but for an exercise, it is fine to relax. I do not intrude on real life much and **it**, largely, leaves **me** untouched in grateful response.

#	Query	Subject	Span(bp)	Coverage(%)	Overall(%)	Exon(%)	CDS(%)	In-frame(%)
1	M77844.1(+)	pax6_genomic(+)	7618-28741	99.03	98.84	99.82	0.00	0.00

Model 1	Coverage	CDS	Overall	In-frame	Mismatches and Indels	3
	99.03%	0.00%	98.84%	0.00%	Exons (min/max/ave), bp	61 / 218 / 135
	99.82%	0.00%	99.82%	0.00%	Introns (min/max/ave), bp	99 / 5903 / 1773

M77844.1 (+) Homo sapiens oculorhombin (PAX6) mRNA, complete cds, alternatively spliced

Segments Alignment

```

1 2 3 4 5 6
7 8 9 10 11 12
13
1279 .....TTTCCTCTCACATCTGGCTCCATGTTGGGCTTAACAGACACAGCCCTCACAAAACACCTACAGC
25105 CTCAGTTTCCTCTCACATCTGGCTCCATGTTGGGCTTAACAGACACAGCCCTCACAAAACACCTACAGC
1344 GCTCTGCCGCTATGCCAGCTTACCATGGCAAATAACCTGCCTATGCAA.....
25175 GCTCTGCCGCTATGCCAGCTTACCATGGCAAATAACCTGCCTATGCAAGTAAAG
    
```

What is the amino acid corresponding to the mutated position in the **Genomic** sequence?

What are the **Genomic** and **mRNA** base positions corresponding to the mutation at amino acid position **33**?

Sensitive Local Pairwise Sequence Comparison

Finally, a swift look at sensitive local pairwise sequence alignment. You have already used **blast** to do a local pairwise alignment in the last Practical, when you aligned the two human genomic sequencing contigs that covered the **PAX6** location in **Chromosome 11**. **blast** did not use a sensitive approach however, nothing subtle was required for that particular alignment.

For a more accurate alignment, return to the **Pairwise Sequence Alignment EBI page** (<http://www.ebi.ac.uk/Tools/psa/>).

From there, select the **Nucleotide** option for the **Local Alignment program Matcher**.

Water or **LALIGN** would also be fine options, but I declare the nucleotide option of **Matcher** to be choice of the day.

Local Alignment

Local alignment tools find one, or more, alignments describing the most similar region(s) within the sequences to be aligned. There are separate forms for protein or nucleotide sequences.

Water (EMBOSS)

EMBOSS Water uses the Smith-Waterman algorithm (modified for speed enhancements) to calculate the local alignment of two sequences.

[Protein](#) [Nucleotide](#)

Matcher (EMBOSS)

EMBOSS Matcher identifies local similarities between two sequences using a rigorous algorithm based on the LALIGN application.

[Protein](#) [Nucleotide](#)

LALIGN

LALIGN finds internal duplications by calculating non-intersecting local alignments of protein or DNA sequences.

[Protein](#) [Nucleotide](#)

Tools > Pairwise Sequence Alignment > EMBOSS Matcher

Pairwise Sequence Alignment (NUCLEOTIDE)

EMBOSS Matcher identifies local similarities in two input sequences using a rigorous algorithm based on Bill Pearson's lalign application, version 2.0u4 (Feb. 1996).

This is the form for nucleotide sequences. Please go to the [protein](#) form if you wish to align protein sequences.

STEP 1 - Enter your nucleotide sequences

Enter or paste your first nucleotide sequence in any supported format:

Or, upload a file: pax6_genomic.fasta

AND

Enter or paste your second nucleotide sequence in any supported format:

Or, upload a file: pax6_mrna.fasta

STEP 2 - Set your pairwise alignment options

MATRIX	GAP OPEN	GAP EXTEND	ALTERNATIVES MATCHES	OUTPUT FORMAT
DNAfull	16	4	1	pair

STEP 3 - Submit your job

Be notified by email (Tick this box if you want to be notified by email when the results are available)

Load up the **Genomic** and **mRNA** sequences as you did for **Needle**.

Click on the **More options** button to see what parameters you can set. They should be as you might expect. The defaults are fine for the first run.

Click on the **Submit** button to get **Matcher** into Matchbox mode.

After due consideration of all the possibilities, **Matcher** will enrich your screen with its conclusions.

But, only one alignment? A good one, covering the highest scoring region of all those considered, but it cannot be the whole story, which must tell the tale of **12** exons! Here is but one.

In common with most local alignment programs, by default **Matcher** will only show you the single best local alignment between two sequences.

A good reason to have a **Dot Matrix View** to inform one of roughly what to expect, which is not one miserable alignment in this case.

pax6_genomic	16871	CACTTCCCCTAT---GCAGGTGTCCAACGGATGTGTGAGTAAAATTCGG	16917
M77844.1	485	534
pax6_genomic	16918	GCAGGTATTACGAGACTGGCTCCATCAGACCCAGGGCAATCGGTGGTAGT	16967
M77844.1	535	584
pax6_genomic	16968	AAACCGAGAGTAGCGACTCCAGAAGTTGTAAGCAAATAGCCAGTATAA	17017
M77844.1	585	634
pax6_genomic	17018	GCGGGAGTGCCCGTCCATCTTTGCTTGGGAAATCCGAGACAGATTACTGT	17067
M77844.1	635	684
pax6_genomic	17068	CCGAGGGGGTCTGTACCAACGATAACATACCAAGCGTAAGTTCATTGAGA	17117
M77844.1	685	734
pax6_genomic	17118	ACA--TCTGCCCTCCCTGCC	17135
M77844.1	735	754

Of course, it is also miserable biologically! **Matcher** fails to align the exons accurately for all the same reasons that the **Needle** failed to represent the *biological* reality.

So, what can one do but try again! By returning to the **Matcher** launch page from your results. You should find the two sequences are still selected, so you should only have to click on **More Options** again and set the **ALTERNATIVE MATCHES** field **20**.

STEP 2 - Set your pairwise alignment options

MATRIX	GAP OPEN	GAP EXTEND	ALTERNATIVES MATCHES	OUTPUT FORMAT
DNAfull	16	4	20	pair

Actually, as you know there are only **12** exons. And that some might well be close enough to be included in the same alignment, you do not need to go as high as **20**. However, the web interface restricts choice (**WHY!?**) such that this is the most sensible cautious choice.

pax6_genomic	24856	TCCAGGTATGGTTTTCTAATCGAAGGGCCAAATGGAGAAGAGAAGAAAA	24905
M77844.1	1123	TACAGGTATGGTTTTCTAATCGAAGGGCCAAATGGAGAAGAGAAGAAAA	1172
pax6_genomic	24906	CTGAGGAATCAGAGAAGACAGGCCAGCAACACACCTAGTCATATTCCTAT	24955
M77844.1	1173	CTGAGGAATCAGAGAAGACAGGCCAGCAACACACCTAGTCATATTCCTAT	1222
pax6_genomic	24956	CAGCAGTAGTTTCAGCACCAGTGTCTACCAACCAATTCACAACCCACCA	25005
M77844.1	1223	CAGCAGTAGTTTCAGCACCAGTGTCTACCAACCAATTCACAACCCACCA	1272
pax6_genomic	25006	CACCGGTAATTTGAAATACTAATACTACGAATCAATGTCTTTAAACCTG	25055
M77844.1	1273	CACCGG-----	1278
pax6_genomic	25056	TTTGCTCCGGGCTCTGACTCTCACTCTGACTACTGTCTATTCTCTTGCC	25105
M77844.1	1279	-----	1278
pax6_genomic	25106	TCAGTTTCCTCCTTCACATCTGGCTCCATGTTGGGCCGAACAGACACAGC	25155
M77844.1	1279	----TTTCCTCCTTCACATCTGGCTCCATGTTGGGCCGAACAGACACAGC	1324
pax6_genomic	25156	CCTCACAACACCTACAGCGCTCTGCCCTATGCCAGCTTCACCATGG	25205
M77844.1	1325	CCTCACAACACCTACAGCGCTCTGCCCTATGCCAGCTTCACCATGG	1374
pax6_genomic	25206	CAAATAACCTGCCTATGCAA	25225
M77844.1	1375	CAAATAACCTGCCTATGCAA	1394

Click on the **Submit** button and **Matcher** will trust and obey.

At the top of your output will be some nice believable local alignments, some involving more than one exon.

Matcher tries to make each alignment as long as it can, stopping only when, to stretch the alignment any further would involve the alignment score decreasing due to the necessity for gap penalties.

```

#-----
#
# Aligned sequences: 2
# 1: pax6_genomic
# 2: M77844.1
# Matrix: EDNAFULL
# Gap_penalty: 16
# Extend_penalty: 4
#
# Length: 46
# Identity:      31/46 (67.4%)
# Similarity:   31/46 (67.4%)
# Gaps:         1/46 ( 2.2%)
# Score: 83
#
#-----
pax6_genomic 11618 ACAGTTTGACTGAGCCCTAGATGCATGTGTTTTT-CCTGAGAGTGA 11662
M77844.1    1043  AGAGTTTGAGAGAACCATTATCCAGATGTGTTTGCCCGAGAAAGA 1088

#-----
#
# Aligned sequences: 2
# 1: pax6_genomic
# 2: M77844.1
# Matrix: EDNAFULL
# Gap_penalty: 16
# Extend_penalty: 4
#
# Length: 58
# Identity:      39/58 (67.2%)
# Similarity:   39/58 (67.2%)
# Gaps:         6/58 (10.3%)
# Score: 83
#
#-----
pax6_genomic 2554 GCTGGACGCCACCCGGCGCCAGA--GCCGGG--CTGAGGAGCGGGGTC 2598
M77844.1    425  GCCGGACTCCACCCGGCAGAAGATTGTAGAGCTAGCTCAC-AGCGGGGCC 473

pax6_genomic 2599 TGGCCGGG 2606
M77844.1    474  CGCCGTG 481
    
```

Go to far down the list of alignments and you will realise what a literal interpretation **Matcher** has of its duties.

You asked for **20** alignments?

So here are the best **20** alignments and it is entirely up to you to decide where “silly” begins.

Not too difficult in this case I suggest.

Why do you suppose your aligned exons are not presented in the correct positional order?

THE END

DPJ – 2019.01.30

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 background 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 Global Alignment:

What do you suppose these regions represent?

Exons

Or does it? How you you interpret the run of As at the end of the final exon?

Summary:

Well, whatever they are they cannot be a **PolyA Tail** as they exist both as part is the **mRNA** *AND* the **Genomic** sequence!

As you assuredly know already, **Polyadenylation** (the addition of a **poly(A) tail** to a messenger RNA) is part of the process that produces mature messenger RNA (**mRNA**). So the As of a **poly (A)** tail occur only at the end of the **messenger RNA**, *NOT* in the genomic sequence!

So, I would suppose the As in question are the **3' UnTranslated Region (UTR)**, or at least part of it.

Full Answer:

This **mRNA** was born in **1991**, as can be confirmed by a quick glance at its **Genbank** annotation.

```

REFERENCE 1 (bases 1 to 1643)

AUTHORS Ton,C.C.T., Hirvonen,H., Miwa,H., Well,M.M., Monaghan,P.,
Jordan,T., van Heyningen,V., Hastle,N.D., Meijers-Heijboer,H.,
Drechsler,M., Royer-Pokora,B., Collins,F.S., Swaroop,A.,
Strong,L.C. and Saunders,G.F.
TITLE Positional cloning and characterization of a paired box- and
homeobox-containing gene from the aniridia region
JOURNAL Cell 67 (6), 1059-1074 (1991)
PUBMED 1684738
    
```

mRNA sequences of this era quite often were submitted with incomplete **UTRs**.

The absence of a **polyA_site Feature** further suggests the As at the end of **M77844** are not a complete **3' UTR**.

```

misc feature 1017..1166
/gene="PAX6"
/gene_synonym="aniridia"
/note="Region: homeobox"

ORIGIN
1 tatcgataag tttttttttt attgtcaatc tctgtctcct tcccaggaat ctgaggattg
61 ctcttacaca ccaaccagc aacatccgtg gagaaaactc tcaccagcaa ctctcttaaa
121 acaccgtcat ttcaaacatt tgtgtcttc aagcaacaac agcagcacia aaaaccccaa
181 ccaaacaaaa ctcttgacag aagctgtgac aaccagaaag gatgctcat aaaggggaa
241 gactttaact aggggcccgc agatgtgtga ggcctttat tgtgagagt gacagacatc
301 cgagattcca gagcccata ttcgacccc gtggaatccc ggcgcccaca gccagagcca
361 gcctgcagaa cagtcacagc ggaatgaatc agctcgggtg tctctttgac aacggcggg
421 cactgccgga ctccaccggg cagaagattg tagagctaac tcacagcggg gcccgccgtg
481 gcgacatttc ccgaattctg caggtgtcca acggatgtgt gagtaaaatt ctggccaggt
541 attacagagc tggctccatc agaccaggg caatcgggtg tagtaaacgg agagttagca
601 ctccagaagt tgaagcaaaa atagcccagt ataagcggga gtgccctgcc atctttgctt
661 gggaaaatcc agacagatta ctgtccgagg gggtctgtac caacgataac ataccaagcc
721 tgtcatcaat aaacagagtt ctctgcaacc tggctagcga aaacacaacag atggggcag
781 atggcagta tgaataacta aggatgttga acgggcagac cggagcctgg gccaccgcc
841 ctggttgata tccgggactc tgggtccagg gcaacctac gcaagatgac tccaagcaad
901 aggaaggagg gggagaaat accaactcca tcagttccaa cggagaagat tcaatgagg
961 ctcaaatgag acttcaactg aagcggaaag tcaaaagaaa tagaacatcc ttacccaag
1021 agcaaatgga gcccctggag aaagattttg agagaaccca ttatccagat gtgtttgcc
1081 gagaagaact agcagcmeta atagatctac ctgaagcaag aatacagata tggttttcta
1141 atcgaaggcc caaatgaga agagaagaaa aactgaggaa tcagaagaaa caggccagca
1201 acacacctag tcatattctc atcagcagta gtttcagcac cagttctcac caaccaattc
1261 caaacccac cacaccggtt tctctcttca catctggctc catgttgggc ctaacagaca
1321 cagccctcac aaacacctac agcctctctc cgcctatgcc cagtttacc atggcaataa
1381 acctgcctat gcaaccacca gtcccagccc agacctcttc atactctgc atgtctccca
1441 ccagcccttc ggtgaatggg cggagttatg atactacac ccccacat atgcagacac
1501 acatgaacag tcagccaagt ggcacctcgg gcaccactc aacaggactc atttccctg
1561 gtgtgtcagt tccagttcaa gtcccggaa gtgaacctga tatgttcaa tactggccaa
1621 gattacagta aaaaaaaaaa aaa
//
    
```

In the following example of a more recent (**2018.02**) **mRNA GenBank** entry, there is a **polyA_site** at the end of the final exon (highlighted, and implying a complete **3' UTR**) with the **polyA** itself included as a part of the sequence.

Only the highlighted **A** at position **1284**, which is the **polyA_site**, will occur in the **Genomic** sequence.

```

polyA_site 1284
/gene="LDHC"
/gene_synonym="CT32; LDH3; LDHX"

ORIGIN
1 cgtgcgtgtc tcgagtcgca cggagggcaa cgtcgcagc gcttagcgc tcaactgtgc
61 ttgggtgatt tttctgggtt cacttctgtg ccttccttca aaggtggtgc tttgtccctg
121 tgggtcatct gtactgattg cgcaaacgaa agcatttgtt ctccaaatgt caactgtcaa
181 ggagcagcta attgagaagc taattgagga tgatgaaaac tcccagtgta aaattactat
241 tgttgaact ggtgccgtag gcatggctgt tgctattagt atcttactga aggatttggc
301 tgatgaact gcccttgttg atgttgcatt ggacaaactg aaggagaaaa tgatggatct
361 tcagcatggc agtcttttct ttagtacttc aaagattact tctgaaaaag attacagttg
421 atctycaaac tccagaatag ttattgtcac agcaggtgca aggcagcagg agggagaaac
481 tcgcttgccc ctggctcaac gtaatgtggc tataatgaaa tcaatcattc ctgccatagt
541 ctattatag cctgattgta aaattctgtg ttttcaaat ccagtgata tttgacata
601 tatagtctgg aagataagtg gcttacctgt aactctgtga attggaagtg gttgtaactc
661 agactctgcc cgittccggt acctaattgg agaaaagttg ggtgtccacc ccacaagctg
721 ccatggttgg attattggag aacatggtga tcttagtgtg cccttagtga gtggggtgaa
781 tgttctggtt gttgctctga agactctgga ccctaaatta ggaacggatt cagataagga
841 acactggaaa aatatccata aacaagtat tcaaatgtcc tatgaaatta tcaagctgaa
901 ggggtatacc tcttgggcta ttgactctgt tgtgatgatt ctggtagatt ccattttgaa
961 aaatcttagg agagtgacc cagtttccac catggttaag ggattatgtg gaataaaaa
1021 agaactcttt ctcaatgac cttgtgtctt gggcggaaat ggtgtctcag atgttgtgaa
1081 aattaaactg aattctgagg aggaggccct ttcaagaag agtgcagaaa cactttggaa
1141 tttcaaaaag gatcctaata tttaaattaa agccttctaa ttttccactg tttggagaac
1201 agaagatagc agcgtgtgta ttttaatttt tgaagatttt ttcatttgat ctttaaaaaa
1261 taaaaacaaa ttggagacct gtgcaaaaa aaaaaaaaaa aaaaaaaaaa aaaaaaaaa
//
    
```


How many convincingly aligned regions did you see?

4

How many did you expect?

12, as that was how many **blast** found, not including the silly ones at the beginning.

The 4 that were found correspond the illustrated 4 diagonal lines grouped together in the **Dot Matrix View** made by **blast**.

Clearly, this alignment is not correct. Can you explain why?



This alignment algorithm only wishes to maximise an alignment score. It sees ALL the high scoring exon regions, however, as the gaps between many of the exons (introns that is) are so long that the penalties for representing them correctly are greater than the gain achieved by the inclusion the extra exons in the alignment. Arithmetically, it is better to align all the exons either side of the 4 exons that were aligned sensibly, in the biologically improbably fashion shown. Arithmetically the best alignment, biologically ridiculous!

This behaviour is exaggerated because this program regards the enormous gaps in has suggested at the start and end of the alignments as “free”. Some global alignment programs (including this one if you ask politely, as you will see) offer the option of penalising the ends gaps in the same way as for internal gaps. Normally, not penalising end gaps is sensible as it allows for the sequences to have slightly different lengths. In this case, penalising end gaps will result in a far better alignment.

Had you used **stretcher** (also offered by the **EBI**) you would have got a much improved answer in this case (but not necessarily in generally). This is because **stretcher** works in a way far closer to the way an informed human might think. **stretcher** does not mindlessly insist of the highest alignment score. Instead, it looks for all the high scoring regions (i.e. all the exons) and then computes the best way to link them together. The result is a far more convincing alignment, but not the arithmetically best scoring answer.

How many matching regions are there this time?

Were you to trawl though your textual output carefully (or simply take my immaculate word for it), you would find 12 perfectly (or nearly so) aligned regions, implying 12 exons.

To be pedantic, the nicely aligned regions do not match the exons exactly (as has been discussed), but well enough to claim definite evidence for the number of exons. 12 is good enough for me.

Is the count now roughly as you would expect?

Yes, exactly the same as **blast** predicted in the first place. More exons than 17 might have been a surprise as that is how many the gene record for **PAX6** at the **NCBI** suggested. Any given transcript may have less than 17 exons or exactly 17 exons, but not more than 17 exons if the heroes of the **NCBI** are not mistaken.

How do you think **blast** achieve the correct results without any fuss?

The only way **blast** could have got the right answer, as it did, would be to use one of the strategies listed previously. **blast** did not use the horrible idea of making gaps super cheap! Not only is that a disgustingly dirty trick, but **blast** actually declares that it is using quite sensible gap penalties.

Leaving **penalising end gaps** and/or using the same sort of heuristics employed by **stretcher**. I would strongly suspect **blast** uses a **stretcher** approach. After all, **blast** has clearly already identified all the “promising regions” in order to construct its **Dot Matrix View**. Also the **stretcher** strategy is similar to that of all **blast** searches (discussed in the next Practical). Finally, **blast** is often used to align very long DNA sequences to detect very strongly similar large regions. This is exactly what the faster (if less pure) **stretcher** approach is all about.

From your investigations comparing mRNA/cDNA with genomic DNA:

What is the amino acid corresponding to the mutated position in the **Genomic** sequence?

```

T S G S M L G L T D T A L T N
ACATCTGGCTCCA TGTTGGGCCTAACAGACACAGCCCTCACAAAC
|||||
ACATCTGGCTCCA TGTTGGGCCGAAACAGACACAGCCCTCACAAAC
    
```

The top sequence is the mRNA. **splign** is kind enough to explicitly inform us that the “mutated” codon, **CTA**, will be expressed as **Leucine**.

So, why not translate the **Genomic** sequence also **splign**?! Easy enough to look up. But I resent having to do so!

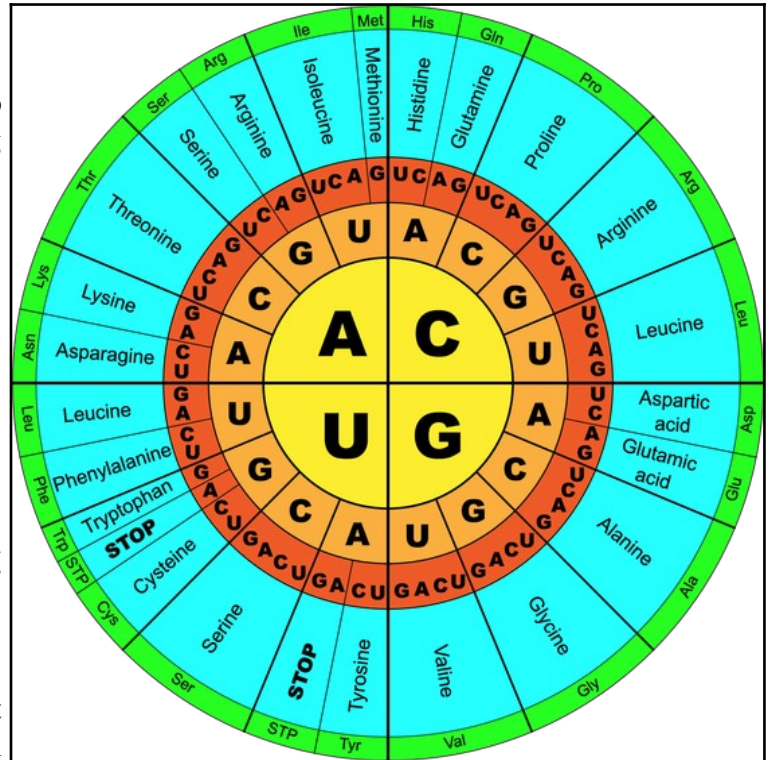
From this rather beautiful representation of the **Genetic Code**, I conclude:

mRNA CTA → Leucine (L)
 Genomic CGA → Arginine (R)

I checked, and this does not appear to be a substitution that is associated with any “interesting” phenotype.

There is no real reason why it should. We did not pause to find out anything about the mRNA downloaded from the **NCBI**, The annotation is particularly unrevealing by itself (it is in **Backup_Files** if you really want to check).

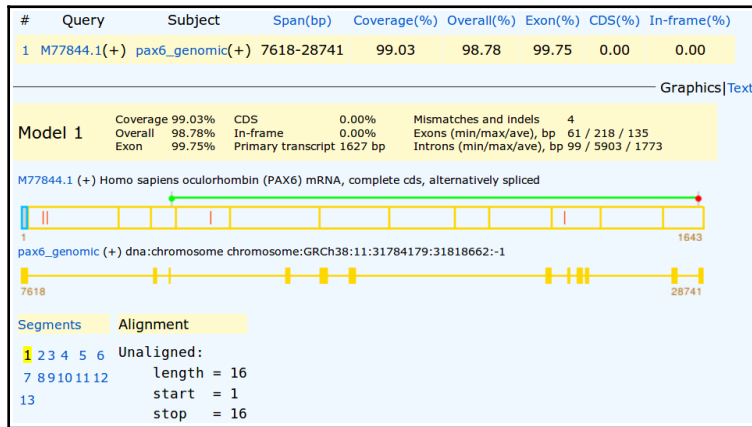
Let us simply assume it is a benign **Accepted Point Mutation (PAM)**. Yes indeed, that feels comfortable. Not so very tricky this Science stuff after all what!



What are the **Genomic** and **mRNA** base positions corresponding to the mutation at amino acid position **33**?

Remember the **Natural variation** at amino acid position **33**? You looked at it in passing during the course of the first exercise. It is a major cause of **Aniridia**. An **Alanine** mutated to a **Proline** at the end of a **Helix** vital to the **DNA Binding** function of the **PAX6** protein.

Natural variant	(VAR_008694)	29	I → S in AN.	1 Publication
Natural variant	(VAR_003811)	29	I → V in AN.	1 Publication
Natural variant	(VAR_008695)	33	A → P in AN.	1 Publication
Natural variant	(VAR_008696)	37 - 39	Missing in AN.	1 Publication
Natural variant	(VAR_008697)	42	I → S in AN; mild.	1 Publication
Natural variant	(VAR_008698)	43	S → P in AN.	1 Publication
Natural variant	(VAR_003812)	44	R → Q in AN.	1 Publication



splign shows alignments for all exons and from those alignments the answer to this question is thus clearly available. To make finding the right spot in the alignment to study easier, I ran **splign** again with an edited version of the **mRNA** (saved as **pax6_mrna_edited.fasta** amongst your cheat files) against the same **Genomic** sequence. Had there been a suitable **mRNA** sequence in the databases, I would have used it for the exercise, but there is not.

You should be able to clearly see the extra mutation is in the **5th** segment.

Focussing on the **5th** segment, the substitution is clear. Using the same methods as were used for the previous question, it is easy to confirm that the variation at amino acid position **33** amounts to:

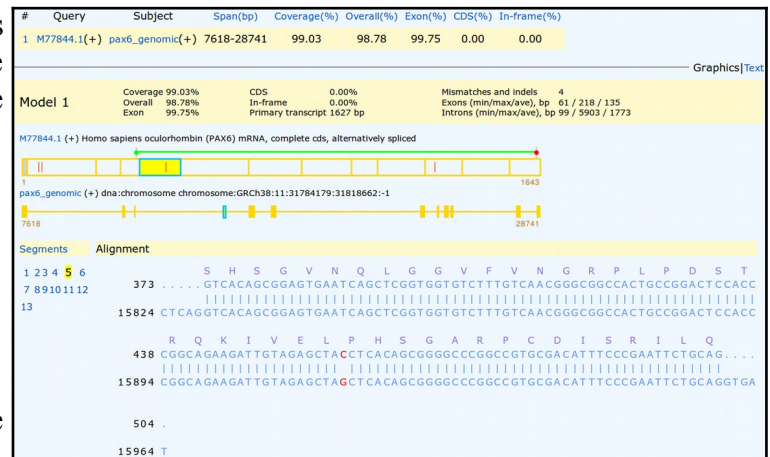
Affected Patient protein:

CCT → Proline (P)

Canonical protein:

GCT → Alanine (A)

Squinting madly, you can also discover that the variation base positions are:



Affected Patient mRNA: Base position 459 → C

Wild Type Genomic DNA: Base position 15915 → G

In case you were wondering, chasing these values around is a little more than tragic pedantry. You will need this information later when you investigate **Primer Design**. No need to take notes, I will remind you of what you need when the time comes. Here I just want to show how the values could be determined, if you had to. Not difficult, just tedious!

How do you interpret the **Details** column for exons 1 and 10?

Summary:

The **Details** column shows the alignments of each exon in a compressed format described in the **splign** documentation as illustrated.

11. Alignment transcript	Alignment transcript represents full details of the alignment in a form of a string composed of characters 'M', 'R', 'I' and 'D' where each character corresponds to an elementary command (Match, Replace, Insert or Delete) needed to transform the query segment into the subject segment. The string is encoded with RLE.
--------------------------	---

The majority of the exon alignments are trivial.

#	Query	Subject	Span(bp)	Coverage(%)	Overall(%)	Exon(%)	CDS(%)	In-frame(%)
1	M77844.1(+)	pax6_genomic(+)	7618-28741	99.03	98.84	99.82	0.00	0.00

[Graphics](#)|[Text](#)

#	Query	Subject	Idty	Len	Q.Start	Q.Fin	S.Start	S.Fin	Type	Details
+1	M77844.1	pax6_genomic	-	16	1	16	-	-	<L-Gap>	-
+1	M77844.1	pax6_genomic	0.991	218	17	234	7618	7835	CA<exon>GT	M39RM8RM169
+1	M77844.1	pax6_genomic	1	77	235	311	11738	11814	AG<exon>GC	M77
+1	M77844.1	pax6_genomic	1	61	312	372	12201	12261	AG<exon>GT	M61
+1	M77844.1	pax6_genomic	1	131	373	503	15829	15959	AG<exon>GT	M131
+1	M77844.1	pax6_genomic	1	216	504	719	16887	17102	AG<exon>GT	M216
+1	M77844.1	pax6_genomic	1	166	720	885	17807	17972	AG<exon>GT	M166
+1	M77844.1	pax6_genomic	1	159	886	1044	23875	24033	AG<exon>GT	M159
+1	M77844.1	pax6_genomic	1	83	1045	1127	24549	24631	AG<exon>GT	M83
+1	M77844.1	pax6_genomic	1	151	1128	1278	24861	25011	AG<exon>GT	M151
+1	M77844.1	pax6_genomic	0.991	116	1279	1394	25110	25225	AG<exon>GT	M33RM82
+1	M77844.1	pax6_genomic	1	151	1395	1545	27803	27953	AG<exon>GT	M151
+1	M77844.1	pax6_genomic	1	98	1546	1643	28644	28741	AG<exon>	M98

For example:

For **Exon 2**, **splign** informs us **M77**, meaning “There are **77** bases aligned and they all **Match** perfectly”.

For **Exon 4**, **splign** informs us **M131**, meaning “There are **131** bases aligned and they all **Match** perfectly”.

The only **2** interesting entries are those where there are some disagreements. That is, the entries for **Exons 1** and **5**, which, following the documentation, I translate thus:

Exon 1 – M39RM8RM169

An alignment of **218** bases, the first **39** of which **Match** perfectly (**M39**), there then follows an **Replacement (R)**, a further **8 Matched bases(M8)**, a second **Replacement (R)** all finished off with **169 Matched bases (M169)**.

Exon 10 – M33RM82

An alignment of **116** bases, the first **33** of which **Match** perfectly (**M33**), there then follows a **Replacement (R)** and a further **82 Matched bases(M82)**.

Its a pity there are no **Insertions (I)** and **Deletions (D)**, but this was the best **mRNA** I could find.

Full Answer:

A point of pedantry to commence. From a different example, which included **InDels**, I got the display illustrated.

The exon was reported as: **M53IM5IM43**

This implies that the choice of **Insertion (I)** or **Deletion (D)** is made to describe the type of variation required to transform the **cDNA (Query)** sequence into the **genomic (Subject)**. Hence the two **InDels** displayed here are considered to be **Insertions**.

```

1 CAGAGGTCAGGCTTCGCTAATGGGCCAGTGAGGAGCGGTGGAGGCGAGGCCGG - CGCCG - CACACACACA
|||||
7245 CAGAGGTCAGGCTTCGCTAATGGGCCAGTGAGGAGCGGTGGAGGCGAGGCCGGGCGCCGGCACACACACA
    
```

Not that it is a vital issue, but I would have thought the other way around was more logical? That is, to consider the **genomic** sequence as the **reference** against which a particular **mRNA** might vary. In other words, what we see here would surely be more relevantly recorded as “This **mRNA/cDNA** has two **Deletions** relative to the **genomic** sequence which, presumably, attempts to represent the norm in the general population”? Just the reflection of an irretrievable pedant, but I am right, nevertheless!!!

In the documentation (see illustration in the **Summary** answer) it enigmatically states “The string is encoded with **RLE**.”. Just in case, **RLE** stands for **Run-length encoding** which is succinctly defined by **Wikipedia**. In a nutshell, it is a very simple form of data compression that recognizes that:

XX

can be compressed to:

60X

which has to be very effective for any data that has runs of identical characters of significant length. This is certainly the case here where one would expect long stretches of **Ms** in most alignments. Of course, life would get tricky if the data included numeric characters, but that is not an issue here⁴.

I think it worth mentioning, that this way of representing an alignment is a simplification of **CIGAR** format⁵. This format is used for **SAM** (Sequence Alignment Map) and **BAM** (Binary Alignment Map, exactly the same as **SAM**, except compressed) files. You will be engulfed in **SAM/BAM** files if you ever do any **Next Generation Sequencing (NGS)**⁶.

CIGAR: CIGAR string. The CIGAR operations are given in the following table (set '*' if unavailable):

Op	BAM	Description
M	0	alignment match (can be a sequence match or mismatch)
I	1	insertion to the reference
D	2	deletion from the reference
N	3	skipped region from the reference
S	4	soft clipping (clipped sequences present in SEQ)
H	5	hard clipping (clipped sequences NOT present in SEQ)
P	6	padding (silent deletion from padded reference)
=	7	sequence match
X	8	sequence mismatch

- H can only be present as the first and/or last operation.
- S may only have H operations between them and the ends of the CIGAR string.
- For mRNA-to-genome alignment, an N operation represents an intron. For other types of alignments, the interpretation of N is not defined.
- Sum of lengths of the M/I/S/=/X operations shall equal the length of SEQ.

So, straight from the **SAM/BAM Format Specification** I copy the table of **CIGAR** enlightenment.

Note, in particular, the extended range of **Operators** and the different meaning associated with the operator '**M**'. The operators '=' and 'X' are such that any '**M**' is either an '=' or and 'X' but never both. Which leaves one pondering when one might use '**M**' in preference to either an '=' or an 'X'?

4 The **Wikipedia** article shows how this complication might be overcome.
 5 There may or may not be some justification for calling the format **CIGAR**, but if there is, I have no idea what it might be.
 6 **NGS** is also referred to as **High Troughput Sequencing (HTS)**, which, on the whole, I think is a more meaningful name.

Compare the predicted **splign** intron/exon boundaries with the conservation suggested by the logo?

What deviation(s) from the model suggested by the logo can you see?

You may have gathered, I rather like this logo, although I rather think it is leading me to make the same point a trifle to often?

The logo is in almost **100%** agreement with the predictions of **splign**.

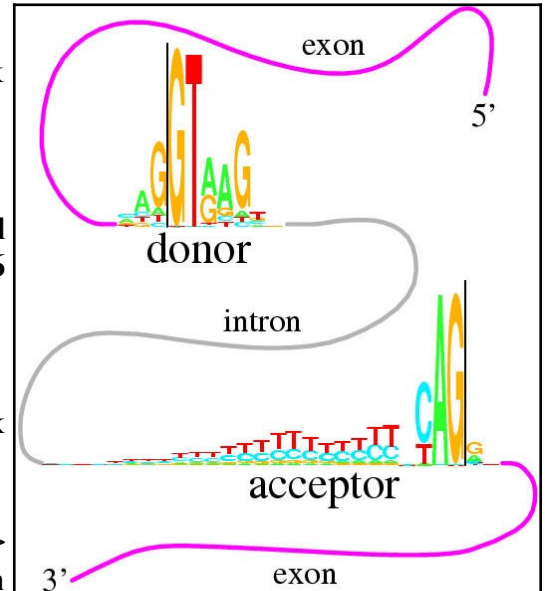
As you will have noted previously, when looking at the **Ensembl** predictions of exons locations of a similar transcript of the **PAX6** human gene (previous Practical), there is a single exception.

Type
<L-Gap>
CA<exon>GT
AG<exon>GC
AG<exon>GT
AG<exon>GT
AG<exon>GT
AG<exon>GT
AG<exon>GT
AG<exon>GT
AG<exon>GT
AG<exon>GT
AG<exon>

The easiest way to show this in the **splign** output is to look at the **splign** text output again.

The **Type** column records the type of all the <exon> alignments it predicts. It also records **2 flanking intron base pairs**.

It is clear that the only time the **splign** prediction deviates from the model suggested by the logo is at the end of the **2nd** exon. Here there is **GC** rather than **GT**. Well, nothing is perfect!



From your investigations of **Local Alignment**:

Why do you suppose your aligned exons are not presented in the correct positional order?

To **Matcher**, the logical order in which to present the alignments is that governed by quality rather than position. So, the highest scoring alignment, rather than the first exon alignment, will be at the top of the list. I think this is generally logical. Once again, the program **splign**, knowing it was looking for an ordered set of exons, was more specifically logical.

DPJ – 2019.01.30

Discussion Points and Casual Questions arising from the Instructions Text.**Notes:*****Work in progress I fear.***

The intention is to provide a full consideration of some issues skimmed over in the exercise proper.

If you are attending a “supervised” presentation of the exercise, I would hope to have conducted a live discussion of all these issues to an extent that reflects:

- the depth that seems appropriate
- the time available
- the degree to which the issues seem to match the interests of the class
- how many of you are awake

Here, I hope to write out very full answers were such a response exists. Accordingly, I suggest you will not need to read much of many of these discussions. There will be much detail of interest to rather few of you. Possibly a bit self indulgent, but I wish to make a note of all the background I have discovered while writing these exercises.

In a nutshell, the exercises are trying to make very general points avoiding too much detail. Nevertheless, I record the detail outside the main exercise text, just in case it might be of interest. Some of the answers to the “**Casual Questions**” are exceedingly trivial. Some of the “**Discussion Points**” are exceedingly long and rambling. You have been warned.

How would you interpret this picture?

What do the diagonal(ish) lines represent?

What are the gaps in between the lines?

Which axis represents the genomic sequence and which the mRNA?

The **Genomic** sequence is represented by the longer **X-Axis**. The **mRNA** is represented by the shorter **Y-Axis**. The two sequences are not represented in strict proportion, but the **Genomic** axis is sufficiently longer than the **mRNA** axis to feel and look intuitively correct.

The sloping lines represent the **Exons** that comprise this **mRNA**. The sloping lines are not at **45** degrees because the **Genomic** sequence is longer than the **mRNA**.



Considered together they cover the whole length of the **mRNA** (except for a few mystery bases at the start).

They represent regions of the **Genomic** sequence (still **Exons**) that are separated by gaps of varying length which are, of course, the **Introns**.

All terribly simple, and I am sure you worked all this out for yourself. However, a fine excuse for yet another beautiful picture.

How many aligned regions are there and do they correspond nicely to the lines of the **Dot Matrix View**?

How many exons would you say this mRNA has?

Well, looking only at the **Dotplot**, I would estimate **12 Exons**. Of course, that would be a dangerous prediction as the resolution of the picture might disguise some very small **Introns**. However, after counting the aligned regions and coming again to a count of **12** (ignoring the silly bit at the start), exactly corresponding to the evidence of the **Dotplot**, I would predict **12 Exons** with confidence.

If one was to forgive the strange “bits” at the start, would you say **blast** seems to have done a reasonable job here?

Yes indeed!

How do you feel about the results this time?

The results generated by **stretcher**, that is.

Well, they are effectively the same as were generated by **blast**. Both **blast** and **stretcher** produce credible alignments whereas **needle** (with default settings) generates a nonsense. On the face of it, rather strange as **needle** is the most exacting of the three options.

Any theories?

Concerning the few wayward bases at the start of the mRNA.

I cannot help you here? Maybe some sequencing artefact? It is a sequence of some antiquity after all.

DPJ – 2019.01.30