The Gulbenkian Training Programme in Bioinformatics (Since 1999)

Pedro Fernandes, Organiser

G PB

# ELB18S

# **Entry Level Bioinformatics**

05-09 November 2018

(Second 2018 run of this Course)

# **Basic Bioinformatics Sessions**

Practical 2: Pairwise Sequence Alignment

Sunday 4 November 2018

# **Sensitive Pairwise Alignment**

option.

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

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?

ſ	Nucleotide Protein				
	Enter Query Se	quence Needleman-Wunsch	alignmen	t of two r	nucleotide sequences 😡
	Enter accession nu	mber, gi, or FASTA sequence 😡		Clear	Query subrange 😡
					From
					То
	Or uplead file				
	or, upload me	Browse pax6_genomic.fasta	Θ		
	Job litle	Enter a descriptive title for your BLAST search	Θ		
	Enter Subject S	equence			
	Enter accession nu	mber, gi, or FASTA sequence 😡		Clear	Subject subrange 😡
-	M77844				From
					То
	Or, upload file	Browse No file selected.	Θ		
:	Align	Show results in a new window			
C	Algorithm paramet	ers			
,	Scoring Para	meters			
	Match/Mismatch Scores	2,-3 🔹 😡			
	Gap Costs	Existence: 5 Extension: 2 🔹 🥹			
	Align				
	Aligh	Show results in a new window			

Global Align

Compare two sequences across their entire span (Needleman-Wunsch)



**Basic Bioinformatics - A Practical User Introduction** 

2 of 23

Practical 1: Pairwise Alignment			04 November	2018
Move down to the textual alignment. There are some	Query	661	CGCTGGCGTGGATATTAAGGAAAGTTAGCGCCTGCCTGAGCACCCTCTTTTCTTATCATT	720
weird little bits and pieces at the front of the alignment	Sbjct	1	tate	4
which defy logic I decide not to dwell on these to much	Query	721	GACATTTAAACTCTGGGGCAGGTCCTCGCGTAGAACGCGGCTGTCAGATCTGCCACTTCC	780
beyond noting that the mRNA has some odd bases at the	Sbjct			
front	Query	781	CCTGCCGAGCGGCGGTGAGAAGTGTGGGAACCGGCGCTGCCAGGCTCACCTGCCTCCCCG	840
nont.	Sbjct			
Also, I have faith that the alignment you look at yields	Query	841		900
the highest alignment score, but equally. I doubt most	Ouerv	901	CLITITIC GCT GCC 6 ACT GCT GT CLC A A A T C A A A GC C C G C C C C A GT G G C C C G G G G G G G G G G G G G	960
people would have chosen to throw these odd bases	Sbict			
about with guite such abandon! <b>People</b> are best!	Query	961	CTTGATTTTTGCTTTTAAAAGGAGGCATACAAAGATGGAAGCGAGTTACTGAGGGAGG	1020
	Sbjct	5	GA	6
You can just see evidence of the little patches of whimsy	Query	1021	TAGGAAGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	1080
in the <b>Dot Matrix View</b> .	Sbjct	7	₩	8
	Ouerv	24541	TCTTTCAGAGTTTGAGAGAACCCATTATCCAGATGTGTTTGCCCGAGAAAGACTAGCAGC	24600
Moving down there are a series of far more convincing	Sbjct	1045	Agtttgagagaacccattatccagatgtgtttgcccgagaaagactagcagc	1096
near perfect alignments	Query	24601	ÇAAAATAGATCTAÇÇTGAAGÇAAGAATAÇAGGTACCGAGAGACTGTGCAGTTTCACACTT	24660
near perfect anglinnents.	Sbjct	1097	CAAAATAGATCTACCTGAAGCAAGAATACAGGTA	1130
	Query	24661	TGTGATTCATACCATTTGTCTTTCCTAGAGACAGAGGTGCTTGTACAGAGTACTATTTAT	24720
You must know what these aligned regions represent by	Sbjct			
now?	Query	24721	TTATAGGACTAATATAATAAAAAGGTTCAGTCTGCTAAATGCTCTGCTGCCATGGGCGTG	24780
	Sbjct			
	Query	24781	GGGAGGGCAGCAGTGGAGGTGCCAAGGTGGGGCTGGGCT	24840
But, just in case:	Ouerv	24841		24900
	Shict	1131		1167
	Query	24901	AAAAACTGAGGAATCAGAGAAGACAGGCCAGCAACACACCTAGTCATATTCCTATCAGCA	24960
What do you suppose these regions represent?	Sbjct	1168	AAAAACTGAGGAATCAGAGAAGACAGGCCAGCAACACCCTAGTCATATTCCTATCAGCA	1227
	Query	24961	GTAGTTTCAGCACCAGTGTCTACCAACCAATTCCACAACCCACCACACCGGGTAATTTGA	25020
How many are there and do they correspond nicely to	Sbjct	1228	ĠŦĂĠŦŦŦĊĂĠĊĂĊĊĂĠŦĠŦĊŦĂĊĊĂĂĊĊĂĂŦŦĊĊĂĊĂĂĊĊĊĂĊĂĊĊĠĠ	1278
the lines of the Det Metrix View?	Query	25021	AATACTAATACTACGAATCAATGTCTTTAAACCTGTTTGCTCCGGGCTCTGACTCTCACT	25080
the lines of the Dot Matrix view?	Sbjct	25.003		25140
	Query	1270		25140
How many exone would you say this mena has?	Ouerv	25141		25200
The many exons would you say this mina has?	Sbict	1310		1369
	Query	25201	ÇAŢĢĢÇAAAŢAAÇCŢĢĊĊŢĄŢĢĊĄAGTAAGTGCGGCTGGTGGTGGCCTGCATAACCCAGG	25260
If one was to forgive the strange "bits" at the start	Sbjct	1370		1394
would you say <b>blast</b> seems to have done a reasonable	Query	25261	CCCCAGAGAAGTGAGGAGTGGCTCAGGGCCTGCGGACCTCATTGGCTGTGTCTGCACCCT	25320
ioh here?	Sbjct			
job nere:	Query	25321	TGAGAGCTTTTCGCACTACAGTGATTGGCTTGACCAGTCAAGTCGGAGACAGTCAATCCC	25380
	Sbjct			
I think I would				
- unite 1 1100000.				
	Query	28561	TTTTTGTAAACCTATAAATTTGTATTCCATGTCTGTTTCTCAAAGGGAATATCTACATGG	28620
	SDJCt	20621	CTATTICATCCACTTCACCACTCATTICCCCCCCCCCCC	20600
	Sbict	1547		1582
The final alignment section even has a <b>PolyA Tail</b> !	Query	28681	TCCCGGAAGTGAACCTGATATGTCTCAATACTGGCCAAGATTACAGTAAAAAAAA	28740
Or does it? How you you interpret the run of $\Delta s$ at the	Sbjct	1583	тсссобаабтбаасстбататбтстсаатастббссаабаттасабтаааааааа	1642
and af the final aven?	Query	28741	AAAAAAAAAAAGGAAAGGAAATATTGTGTTAATTCAGTCAG	28800

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. I fact, pretty though it appears, the alignment blast has generated is not as entirely logical as it might first seem. For example, consider:

Query

Sbjct 1643 Å

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

end of the final exon?

1643

28801 TTGAGCTTTCAGGAAAGAAAGAAAAATGGCTGTTAGAGCCGCTTCAGTTCTACAATTGTG 28860

Next, try aligning the same two sequences with another same algorithm) at the <b>EBI</b>	er program (1	implement	ing the	GIODAI AIIGNMENT Global alignment tools create an end-to-end alignment	ment c				
				the sequences to be aligned. There are separate fo protein or nucleotide sequences.	orms fo				
Jo to the Pairwise Sequence Alignment EBI page:		D		Needle					
(http://www.ebi.a	ac.uk/Tools/p	osa/).		EMBOSS Needle creates an optimal global alig of two sequences using the Needleman-Wunso	gnmen ch				
Select the Nucleotide option for the Global Alignment	program Nee	<mark>dle</mark> .		algorithm. <u>A Protein A Nucleotide</u>					
Pairwise Sequence Alignment (NUCLEOTIDE)	eedle implem	ents the be	est globe	al nairwise algorithm exac	•tlv				
This is the form for nucleotide sequences. Please go to the protein form if you wish to align protein sequences.			51 <u>5</u> 1000		uy.				
STEP 1 - Enter your nucleotide sequences	bad up the firs	st sequenc	e from	pax6_genomic.fasta.					
Enter or paste your first nucleotide sequence in any supported format:	bad up the sec	cond seque	ence fro	m pax6_mrna.fasta.					
C	Click on the More options button to see what parameters you								
ca	n set. They sl	hould be a	as you r	night expect. The defaults	s ar				
Or, upload a file: Browse pax6_genomic.fasta	the for the first	run.							
Enter or paste your second nucleotide sequence in any supported format:	lick on the Su	bmit butt	on to ge	et Needle into action.	223				
		M77844.1	1	.    . 	1				
		pax6_genomic 2 M77844 1	2345 TTTTGTTC     .   12 TTTTTT	GTCCGCGCTCATTGTAGCCTCAAAAT-TCTGCCCACGAAAGTT	2239				
Or, upload a file: Browse pax6_mrna.fasta		pax6_genomic 2	2394 TGCCAACO	SCTCCTGCCCCAGGAGTTTAATAGTTTCCCTTACTCGCGGGGC	2244				
STEP 2 - Set your pairwise alignment options		M77844.1	351	ICTCCT-TCCCAGGAATCTGAGGATTGCTCTTACACAC					
MATRIX GAP OPEN GAP EXTEND OUTPUT FORMAT		M77844.1	72CA	ACCCAGCAACATCC	224				
END GAP PENALTY END GAP OPEN END GAP EXTEND		pax6_genomic 2	2494TCTC	CAATAG-ATCTCCAAGGGCCCATATGGTGGCCAGTGCCGATGA	225				
		pax6_genomic 2	2539 ATCCGCCT	IGTTTAAATGGGGGAGAAAGTTGGGGTTTTAAAACAT	225				
STEP 3 - Submit your job Be notified by email ( <i>Tick this box if you want to be notified by email when the results are available</i> )		M77844.1	115	. .     .   .    TTTAAAACACCGTCATTTCAAACCATTGTGGT	14				
Submit		pax6_genomic 2	2583 -TTCAA       147 CTTCAAGO	AGTICCIGAAAAGAICCCACI    .   .    CAACAACAGCAGCACAAAAAAACCCCCAACCAAACAAAACTCTTG	2260				
Well Nothing like as convincing as the alignment <b>blast</b> i	oroduced	pax6_genomic 2	3746 TTACCTTC	GGAATGTTTTGGTGAGGCTGTCGGGATATAATGCTCTTG	2379				
wen: Nothing like as convincing as the anglinent blast		pax6_genomic 2	3793 GAGTTTA	AGACTACCAGGCCCCT-TTTGGAGGCTCCAAGTTAATCC	238				
Alignment does not even begin until over 22,300 base	e pairs along	M77844.1	830 GGG	······CACCCGCCCTGGTTGG-·····TATCCGG	238				
he genomic sequence. Even then it is not convincing, as	in <u>wrong</u> , if	M77844.1	.   . 856 GGACTTCC		8				
we accept the results already obtained from <b>blas</b>	t as a fair	pax6_genomic 2	3886 CAACAGGA	AAGGAGGGGGAGAGAATACCAACTCCATCAGTTCCAACGGAGA	239				
		pax6_genomic 2	3936 AGATTCAG	5ATGAGGCTCAAATGCGACTTCAGCTGAAGCGGAAGCTGCAAA	239				
		M77844.1	947 AGATTCAG	GATGAGGCTCAAATGCGACTTCAGCTGAAGCGGAAGCTGCAAA	9! 240				
There are some well aligned regions after genomic positi	ion <b>24,500</b> .	M77844.1	997 GAAATAG	AACATCCTTTACCCAAGAGCAAATTGAGGCCCTGGAGAACT	104				
		pax6_genomic 2	4036 GATAGAG1	ITTTTCAAAGTAGAGAAGCAGTAAATCAAAGTAAATGCCACAT	2408				
		pax6_genomic 2	4086 CTTCAGTA	ACAAAGAGCTAAATTTAGCCAGGGCCCTTTGCATAGAAGAATG	241				
		M77844.1	1043		248				
		M77844.1	1125	CAGGTATGGTTTTCTAATGAAGGCCA	11				
		pax6_genomic 2	4886 AATGGAGA	AGAGAAGAAAAACTGAGGAATCAGAGAAGACAGGCCAGCAAC	249				
Then a resumption of chaos after <b>25,230</b> or so.		pax6_genomic 2	4936 ACACCTAG	TCATATTCCTATCAGCAGTAGTTTCAGCACCAGTGTCTACCA	249				
		M77844.1	1203 ACACCTAG	GTCATATTCCTATCAGCAGTAGTTTCAGCACCAGTGTCTACCA	12				
		pax6_genomic 2	4986 ACCAATTC          1253 ACCAATTC	CACAACCACCACACCGGGTAATTTGAAATACTAATACTACG                    CACAACCCACCACACCG	12				
		pax6_genomic 2	5036 AATCAATC	GTCTTTAAACCTGTTTGCTCCGGGCTCTGACTCTCACTCTGAC	250				
Iow many convincingly aligned regions did you see?		M77844.1 pax6_genomic 2	1278 5086 TACTGTCA	ATTTCTCTTGCCCTCAGTTTCCTCCTTCACATCTGGCTCCATG	251				
		M77844.1	1278	GTTTCCTCCTTCACATCTGGCTCCATG	13				
How many did you expect?		paxo_genomic 2 M77844.1	1305 TTGGGCCG	JAALAGALALAGULU I LALAAALALU I ACAGEGETETGEEGEE 	251 13				
		pax6_genomic 2	5186 TATGCCCA	AGCTTCACCATGGCAAATAACCTGCCTATGCAAGTAAGTGCGG	252				
learly this alignment is not correct. Can you evplain wi	<b>a</b> v?	m/7844.1 pax6_genomic 2	1355 TATGCCCA 5236 CTGGTGGT	ITTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTT	13 252				
Jeany, uns angument is not correct. Can you explain wi	цу (	M77844.1	1395	.          CCCCCAGTCCCCAGCCAGA	14				
		pax6_genomic 2	5284 CAGGGCCT	IGCGGACCTCATTGGCTGTGTCTGCACCCTTGAGAG	253				

#### **Practical 1: Pairwise Alignment**

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 <u>but not by the EBI web interface!</u> 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 appr	broach by returning to the <b>Needle</b> launch page from your results.
You should find the two sequences are still	STEP 2 - Set your pairwise alignment options
selected, so you should only have to click on More	GAP OPEN GAP EXTEND OUTPUT FORMAT
<b>Options</b> again and change the END GAP	DNAfull v 10 v 0.5 v pair v
PENALTY field from false to true.	END GAP PENALTY END GAP OPEN END GAP EXTEND

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?

	Tools > Pairwise Sequence Alignment > EMBOSS Stretcher					
Finally check that <b>Stretcher</b> works as expected	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.					
Go again to the Pairwise Sequence Alignment EBI	This is the form for nucleotide sequences. Please go to the protein form if you wish to align protein sequences.					
page (http://www.ebi.ac.uk/Tools/psa/).	STEP 1 - Enter your nucleotide sequences					
	Enter or paste your first nucleotide sequence in any supported format:					
From there select the Nucleotide option for the Global						
Alignment program Stretcher						
Angument program Stretcher.						
Lood we the genuerous expectives for Needle						
Load up the sequences exactly as for freedie.	Or, upload a file: Browse pax6_genomic.fasta					
	AND Enter or paste your second nucleotide sequence in any supported format:					
Take a look at the parameters and see there is nothing						
unexpected hiding there.						
Set Stretcher sequence rone stretching						
set stretener sequence tope stretening.	Or, upload a file: Browse pax6 mrna.fasta					
How do you feel about the regults this time?						
How do you leef about the results this time?	STEP 2 - Set your pairwise alignment options					
	MATRIX GAP OPEN GAP EXTEND OUTPUT FORMAT					
How do you think <b>blast</b> achieve the correct results	UNAUII					
without any free?	STEP 3 - Submit your job					
without any fuss?	Be notified by email (Tick this box if you want to be notified by email when the results are available)					
	Submit					

# 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.

**cDNA** In the 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 choice<sup>1</sup>. Note the option to compare your cDNA with a Whole genome (including Human). Today, the default options are fine. Click the Align button.

#	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			
Mo	odel 1	Coverage 99.03% Overall 98.84% Exon 99.82%	CDS 0 In-frame 0 Primary transcript 1	.00% Misma .00% Exons 627 bp Intro	atches and ind (min/max/av ns (min/max/a	lels 3 re), bp 61 rve), bp 99	/ 218 / 135 / 5903 / 17	773			
M77844.1 (+) Homo sapiens oculorhombin (PAX6) mRNA, complete cds, alternatively spliced											
1 pax	6_genomic (	+) dna:chromosome	chromosome:GRCh38	8:11:31784179:3	1818662:-1			1643			
ŀ											
76	18							28741			
Se	gments	Alignment									
1	23456	Unaligned:									
7	89101112	length = 1	L6								
13		start = 1	16								
		stop = 1	16								



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 13<sup>th</sup> 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?

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.

										04 P	lovemb	er 2018
#	Query	Subject	t S	Span(bp)	Coverage(%)	Overall(%)	Exon(%)	CDS(%)	In-frame(%	)		
1	M77844.1(+)	pax6_genom	nic(+) 76	18-28741	99.03	98.84	99.82	0.00	0.00			
		Coverage 9	9.03%	CDS	0.00	0%	Misr	natches an	d indels 3			
M	odel 1	Overall 9 Exon 9	8.84% 9.82%	In-frai Primai	me 0.00 ry transcript 162	)% 7 bp	Exo Intr	ns (min/ma ons (min/n	ix/ave), bp=6: hax/ave), bp 99	218 / 135 ) / 5903 / 1773		
M7	7844.1 (+) Hom	o sapiens oculor	nombin (PA	(X6) MRNA,	complete cos, alt	ernatively spi	cea					
П							1		T			
1									1643			
pa:	x6_genomic (+)	dna:chromosom	e chromoso	me:GRCh38	:11:31784179:3	1818662:-1						
76	18								28741			
Se	gments A	lignment										
1	<mark>2</mark> 3456											
7	89101112	17 TT	ГТТАТТ <sup>,</sup>	GTCAAT	стстбтстс	CTTCCCA	GGAATC	T G A G <b>G</b> A	TTGCTCT	TACACACC	A A C C C A G (	CAACATC
13	3	7618 11	ΓΤΤΑΤΤ:	GTCAAT	стстбтстс	CTTCCCA	GGAATC	TGAG <mark>A</mark> A	TTGCTCT	CACACACC	AACCCAG	CAACATC
		07.00	STOCAG		CT CA C CA C C	AACTOCT	T T A A A A		CATTICA		COTOTA	
		87 00										
		7688 0	GT GGA G	AAAACT	CTCACCAGO	AACTCCT	ТТАААА	CACCG	CATTTCA	AACCATTG	гостстто	CAAGCAA
		157 C/	ACAGO	AGCACA			CAAAAC	тстте		TGTGACAA		GATGCC
		1		11111		111111	11111	11111	111111			
		7758 C/	AACAGC	AGCACA	AAAAACCCC	AACCAAA	СААААС	тсттби	CAGAAGC	TGTGACAA	CCAGAAA	GGATGCC
		227 T	CATAAA	G								
		1		I								
		7828 T	CATAAA	GGTGAG								

							Click on the exon including the green Start
# Query	Subject	Span(bp) Cov	verage(%) Overall(%	) Exon(%) CDS(%	6) In-frame(%)	)	codon (the <b>3</b> <sup>rd</sup> )
1 M//844.1(+	pax6_genomic(+)	/618-28/41	99.03 98.84	99.82 0.00	0.00		
							The first adding even is new displayed with
Model 1	Coverage 99.03% Overall 98.84%	CDS In-frame	0.00% 0.00%	Mismatches Exons (min/r	and indels 3 max/ave), bp 61	/ 218 / 135	The first coding exon is now displayed with
	Exon 99.82%	Primary tra	anscript 1627 bp	Introns (min	/max/ave), bp 99	/ 5903 / 1773	translation of the mRNA where appropriate.
M77844.1 (+) Hon	mo sapiens oculorhombin	(PAX6) mRNA, comp	plete cds, alternatively s	liced			
П							The statistics at the top of the display include
1 pax6_genomic (+)	) dna:chromosome chrom	osome:GRCh38:11:3	31784179:31818662:-1		1643		the claim that there are <b>3</b> discremancies
-							(Mismatches and Indels) between the
7618					28741		(Mismatches and Inders) between the
Segments	Alignment						cDNA and Genomic sequences.
1 23 <mark>4</mark> 5 6	312 /	GCCCCATATT	CGAGCCCCGTGGA	ATCCCGCGGCG	CCCAGCCAG	M Q N GAGCCAGCATGCAGAACA	
13	12196 AACAGA						Two of these are the substitutions we have
	12150 / 100						already seen in the first exon of the cDNA
	373 .						The third is indicated by the red her in the
	12266 G						The unit is indicated by the red bar in the
	12200 0						<sup>1</sup> 10 <sup>m</sup> exon of the cDNA display.

# Click on the **10<sup>th</sup>** 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.

LI		que, y	Subject	Span(bp)	coverage( /0)	Overall(70)	Ex011( 70)	CD3( 70)	in nunc( ///	
	1	M77844.1(+	<pre>-) pax6_genomic(+)</pre>	7618-28741	99.03	98.84	99.82	0.00	0.00	
	Мо	del 1	Coverage 99.03% Overall 98.84% Exon 99.82%	CDS In-fra Prima	0.0 me 0.0 iry transcript 162	0% 0% 27 bp	Misi Exo Intr	matches ai ns (min/m ons (min/r	nd Indels 3 ax/ave), bp 61 max/ave), bp 99	l / 218 / 135 / 5903 / 1773
<i>r</i>	M77	844.1 (+) Hor    6_genomic (+ 8	) dna:chromosome chror	n (PAX6) mRNA,	complete cds, al	ternatively spl	iced		1643	
	Seg	ments	Alignment							
,	1 2 7 8 13	23456 3910 <mark>11</mark> 12	1279 25105 CTCAG	V S S TTTCCTCC IIIIIIII TTTCCTCC	F T S TTCACATCI IIIIIIIII TTCACATCI	G S IGGCTCCA IIIIIIII IGGCTCCA	M L TGTTGG IIIIII TGTTGG	G L GGCC <b>T</b> A III I GGCC <b>G</b> A	T D T ACAGACAC IIIIIIII ACAGACAC	TALTNTYS CAGCCCTCACAAACACCTACAGC
,			A L 1344 GCTCT       25175 GCTCT	. P P GCCGCCTA           GCCGCCTA	M P S TGCCCAGCI           TGCCCAGCI	F T M TTCACCAT	1 A M GGCAAA         GGCAAA	N TAACC        TAACC	L P M TGCCTATO	Q GCAA         GCAAGTAAG

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?



region.

04 November 2018 Subject Ouerv Coverage(%) Overall(%) Exop(%) M77844.1(+) pax6\_genomic(+) 7618-28741 99.03 98.84 99.82 0.00 0.00 verage 99.03% Mismatches and Indels 3 Exons (min/max/ave), bp 61 / 218 / 135 Introns (min/max/ave), bp 99 / 5903 / 1773 1odel 1 98.84% 99.82% Click on the last exon section in the cDNA display. You should now see the final exon of ic (+) dna:chro me chromosome:GRCh38:11:31784179:31818662:-1 the cDNA with the Stop codon and polyA Alianment 1 2 3 4 5 6 GACTCATTTCCCCTGGTGTGTCAGTTCCAGTTCAAGTTCCCGGAAGTGAACCTGATATGTCTCA 1546 7 89101112 TTCCCCTGGTGTGTCAGTTCCAGTTCCAGTTCCCGGAAGTGAACCTGATAT 28639 1611 28709

#	Query	Subject	Span(bp)	) (	Coverage(%)	Overall(	%) Exon(%)	) CDS(%)	In-frame(%)			1
1	M77844.1(+)	pax6_genomic(+)	7618-287	41	99.03	98.84	99.82	0.00	0.00			
											Graphics	
#	Query	Subject	Idty	Len	Q.Start	Q.Fin	S.Start	S.Fin	Туре	Detai	ls	
+1	M77844.1	pax6_genomic	-	16	1	16	-	-	<l-gap></l-gap>	-		
+1	M77844.1	pax6_genomic	0.991	218	17	234	7618	7835	CA <exon>GT</exon>	M39RI	48RM169	
+1	M77844.1	pax6_genomic	1	77	235	311	11738	11814	AG <exon>GC</exon>	M77		<b>.</b>
+1	M77844.1	pax6_genomic	1	61	312	372	12201	12261	AG <exon>GT</exon>	M61		Fina
+1	M77844.1	pax6_genomic	1	131	373	503	15829	15959	AG <exon>GT</exon>	M131		text
+1	M77844.1	pax6_genomic	1	216	504	719	16887	17102	AG <exon>GT</exon>	M216		10AC
+1	M77844.1	pax6_genomic	1	166	720	885	17807	17972	AG <exon>GT</exon>	M166		
+1	M77844.1	pax6_genomic	1	159	886	1044	23875	24033	AG <exon>GT</exon>	M159		
+1	M77844.1	pax6_genomic	1	83	1045	1127	24549	24631	AG <exon>GT</exon>	M83		
+1	M77844.1	pax6_genomic	1	151	1128	1278	24861	25011	AG <exon>GT</exon>	M151		
+1	M77844.1	pax6_genomic	0.991	116	1279	1394	25110	25225	AG <exon>GT</exon>	M33R	182	
+1	M77844.1	pax6_genomic	1	151	1395	1545	27803	27953	AG <exon>GT</exon>	M151		
+1	M77844.1	pax6_genomic	1	98	1546	1643	28644	28741	AG <exon></exon>	M98		

lly, click on the **Text** link to view the al summary of the **splign** results.

exon

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

Where is the **3<sup>rd</sup>** substitution in the mRNA?

Where is the  $3^{rd}$  substitution in the Genomic Sequence?



2

 $<sup>\</sup>overline{2}$ The original label for this very nice graphic is:

This figure shows two "sequence logos" which represent sequence conservation at the 5' (donor) and 3' (acceptor) ends of human introns. The region between the black vertical bars is removed during mRNA splicing. The logos graphically demonstrate that most of the pattern for locating the intron ends resides on the intron. This allows more codon choices in the protein-coding exons. The logos also show a common pattern "CAGIGT", which suggests that the mechanisms that recognize the two ends of the intron had a common ancestor. See R. M. Stephens and T. D. Schneider, "Features of spliceosome evolution and function inferred from an analysis of the information at human splice sites", J. Mol. Biol., 228, 1124-1136, (1992).

Finally, a swift look at sensitive local pairwise sequence al	ignment. You have	e already used blast to do a local
pairwise alignment in the last Practical, when you aligned the	two human genon	nic sequencing contigs that covered
the PAX6 location in Chromosome 11. blast did not use a se	nsitive approach ho	owever, nothing subtle was required
for that particular alignment.		Local Alignment
		Local alignment tools find one, or more, alignments
For a more accurate alignment return to the Pairwise Sec	uence Alignment	describing the most similar region(s) within the sequences to be aligned. There are separate forms for protein or
EBI page (http://www.ebi.ac.uk/Tools/psa/)		nucleotide sequences.
		Water 🛿 (FMBOSS)
	• ,	EMBOSS Water uses the Smith-Waterman algorithm
From there, select the Nucleotide option for the Local Al Matcher	ignment program	(modified for speed enhancements) to calculate the local alignment of two sequences.
		🔌 Protein 🔌 Nucleotide
Water or LALIGN would also be fine options, but I decl	are the nucleotide	Matcher 🕘 (EMBOSS)
option of <b>Matcher</b> to be choice of the day.	EMBOSS Matcher identifies local similarities between two sequences using a rigorous algorithm based on the LALIGN application.	
		A Protein A Nucleotide
Tools > Pairwise Sequence Alignment > EMBOSS Matcher		LALIGN Ø
PairWise Sequence Alignment (NUCLEOTIDE) EMBOSS Matcher identifies local similarities in two input sequences using a rigorous algorithm based on Bill Pearson's lalign		LALIGN finds internal duplications by calculating
application, version 2.0u4 (Feb. 1996).		non-intersecting local alignments of protein or DNA
This is the form for nucleotide sequences. Please go to the protein form if you wish to align protein sequences.		sequences.
STEP 1 - Enter your nucleotide sequences		Protein Viddeotide
Enter or paste your first nucleotide sequence in any supported format:	Load up the Gen	omic and mRNA sequences as you
	did for Needle.	1
Or, upload a file: Browse pax6_genomic.fasta		
AND Enter or paste your second nucleotide sequence in any supported format:		
	Click on the M	ore options button to see what
	parameters you ca	in set They should be as you might
	expect The defaul	ts are fine for the first run
	enpeet. The defud	
or, uplease a line. Browse paxo_ninna.idsia		
STEP 2 - Set your pairwise alignment options ALTERNATIVES		
MATRIX     GAP OPEN     GAP EXTEND     MATCHES     OUTPUT FORMAT       DNAfull     1     1     1     pair     1		
STEP 3 - Submit your job	Click on the Sul	mit button to get Matcher inte
Be notified by email (Tick this box if you want to be notified by email when the results are available)	Matchhov mode	onne outton to get whatcher into
Submit	matchoox 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	pax6_genomic	16871 CACTTCCCCTATGCAGGTGTCCAACGGATGTGTGAGTAAAATTCTGG	16917
highest scoring region of all those considered, but it	M77844.1	485 CATTTCCCGAATTCTGCAGGTGTCCAACGGATGTGTGAGTAAAATTCTGG	534
cannot be the whole story, which must tell the tale of 12	pax6_genomic	16918 GCAGGTATTACGAGACTGGCTCCATCAGACCCAGGGCAATCGGTGGTAGT	16967
exons! Here is but one.	M77844.1	535 GCAGGTATTACGAGACTGGCTCCATCAGACCCAGGGCAATCGGTGGTAGT	584
	pax6_genomic	16968 AAACCGAGAGTAGCGACTCCAGAAGTTGTAAGCAAAATAGCCCAGTATAA	17017
In common with most local alignment programs, by	M77844.1	585 AAACCGAGAGTAGCGACTCCAGAAGTTGTAAGCAAAATAGCCCAGTATAA	634
default Matcher will only show you the single best local	pax6_genomic	17018 GCGGGAGTGCCCGTCCATCTTTGCTTGGGAAATCCGAGACAGATTACTGT	17067
alignment between two sequences.	M77844.1	635 GCGGGAGTGCCCGTCCATCTTTGCTTGGGAAATCCGAGACAGATTACTGT	684
	pax6_genomic	17068 CCGAGGGGGTCTGTACCAACGATAACATACCAAGCGTAAGTTCATTGAGA	17117
A good reason to have a <b>Dot Matrix View</b> to inform	M77844.1	685 CCGAGGGGGTCTGTACCAACGATAACATACCAAGCGTGTCATCAATAAAC	734
one of roughly what to expect, which is not one	pax6 genomic	17118 ACATCTGCCCTCCCTGCC 17135	
miserable alignment in this case.	M77844.1	.     . . . .  735 AGAGTTCTTCGCAACCTGGC 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.

Practical 1: Pairwise Alignment

Sensitive Local Pairwise Sequence Comparison

04 November 2018

Practical 1: Pairwise Alignment					04 Novembe	er 2018
So, what can one do but try again! By returning to the	e Matcher	launcl	n page from	your results	. You should fin	nd the
two sequences are still selected, so you should ster	2 - Set your pairwise alig	gnment optio	ns	-		
only have to click on More Options again and set MATE	IX	GAP OPE	EN GAP EXTEND	ALTERNATIVES MATCHES	OUTPUT FORMAT	
the ALTERNATIVE MATCHES field 20.	full	• ] 16	▼ ] [ 4	• 20	▼ ) pair	•
Actually, as you know there are only 12 exons. And	pax6_genomic	24856	TCCAGGTATGGTTTTC  .	CTAATCGAAGGGCCA/	AATGGAGAAGAGAAGAAAAAA	24905
that some might well be close enough to be included	M77844.1	1123	TACAGGTATGGTTTTC	TAATCGAAGGGCCA	AATGGAGAAGAAGAAGAAAAA	1172
in the same alignment, you do not need to go as high	pax6_genomic	24906	CTGAGGAATCAGAGA	AGACAGGCCAGCAAC	ACACCTAGTCATATTCCTAT	24955
as 20. However, the web interface restricts choice	M77844.1	1173	CTGAGGAATCAGAGAA	AGACAGGCCAGCAAC	ACACCTAGTCATATTCCTAT	1222
(WHY!?) such that this is the most sensible cautious	pax6 genomic	24956	CAGCAGTAGTTTCAG		ACCAATTCCACAACCCACCA	25005
choice.	M77844 1	1223				1272
		25000				25055
	paxo_genomic	25000		ATACTAATACTACG	AATCAATGICIIIAAACCIG	20000
	M77844.1	1273	CACCGG			1278
Click on the Submit button and Matcher will trust	pax6_genomic	25056	тттостссооостсто	GACTCTCACTCTGAC	TACTGTCATTTCTCTTGCCC	25105
and obey.	M77844.1	1279				1278
	pax6_genomic	25106	тсадтттсстссттс	ACATCTGGCTCCATG	TTGGGCCGAACAGACACAGC	25155
	M77844.1	1279	TTTCCTCCTTCA	ACATCTGGCTCCATG	.	1324
At the top of your output will be some nice believable	pax6 genomic	25156	ССТСАСАААСАССТАС		TATGCCCAGCTTCACCATGG	25205
local alignments some involving more than one exon	M77044 1	1225				1274
	11//844.1	1972			TATUCCCAUCTICACCATO	1574
	pax6_genomic	25206	CAAATAACCTGCCTAT	FGCAA 25225		
Matahan trias to make each alignment as long as it	M77844.1	1375	CAAATAACCTGCCTAT	GCAA 1394		
<b>Wratcher</b> tries to make each alignment as long as it						

can, stopping only when, to stretch the alignment any further would involve the alignment score deceasing due to the necessity for gap penalties.



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

THE END DPJ – 2018.11.04

# 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. <u>BUT</u>, it is not intended for printing or for reading serially, so I submit, being long and wordy does not matter. Feel free to disagree.

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** <u>*AND*</u> 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**, <u>**NOT**</u> 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	<mark>Cell 67 (6), 1059-1074 (1991)</mark>
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.

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.

poly	A_site	1284							
-		/gene="LDHC"							
		/gene synonym="CT32; LDH3; LDHX"							
ORIGIN									
1	cgtgcgtgtc	tcgagtcgca	cggagggcaa	ccgtcgacgg	gcttagcgcc	tcaactgtcg			
61	ttggtgtatt	tttctggtgt	cacttctgtg	ccttccttca	aaggtggtgc	tttgtccctg			
121	tgggtcatct	gtactgattg	cgccaagcaa	agcatttgtt	ctccaaatgt	caactgtcaa			
181	ggagcagcta	attgagaagc	taattgagga	tgatgaaaac	tcccagtgta	aaattactat			
241	tgttggaact	ggtgccgtag	gcatggcttg	tgctattagt	atcttactga	aggatttggc			
301	tgatgaactt	gcccttgttg	atgttgcatt	ggacaaactg	aagggagaaa	tgatggatct			
361	tcagcatggc	agtcttttct	ttagtacttc	aaagattact	tctggaaaag	attacagtgt			
421	atctgcaaac	tccagaatag	ttattgtcac	agcaggtgca	aggcagcagg	agggagaaac			
481	tcgccttgcc	ctggtccaac	gtaatgtggc	tataatgaaa	tcaatcattc	ctgccatagt			
541	ccattatagt	cctgattgta	aaattcttgt	tgtttcaaat	ccagtggata	ttttgacata			
601	tatagtctgg	aagataagtg	gcttacctgt	aactcgtgta	attggaagtg	gttgtaatct			
661	agactctgcc	cgtttccgtt	acctaattgg	agaaaagttg	ggtgtccacc	ccacaagctg			
721	ccatggttgg	attattggag	aacatggtga	ttctagtgtg	cccttatgga	gtggggtgaa			
781	tgttgctggt	gttgctctga	agactctgga	ccctaaatta	ggaacggatt	cagataagga			
841	acactggaaa	aatatccata	aacaagttat	tcaaagtgcc	tatgaaatta	tcaagctgaa			
901	ggggtatacc	tcttgggcta	ttggactgtc	tgtgatggat	ctggtaggat	ccattttgaa			
961	aaatcttagg	agagtgcacc	cagtttccac	catggttaag	ggattatatg	gaataaaaga			
1021	agaactcttt	ctcagtatcc	cttgtgtctt	ggggcggaat	ggtgtctcag	atgttgtgaa			
1081	aattaacttg	aattctgagg	aggaggccct	tttcaagaag	agtgcagaaa	cactttggaa			
1141	tattcaaaag	gatctaatat	tttaaattaa	agccttctaa	tgttccactg	tttggagaac			
1201	agaagatagc	aggctgtgta	ttttaaattt	tgaaagtatt	ttcatttgat	ctttaaaaaa			
1261	taaaaacaaa	ttggagacct	gtg <mark>a</mark> aaaaaa	aaaaaaaaaa	aaaaaaaaaa	aaaaaaa			
11									

mise reactive										
			/gene="PAX6"							
			/gene_synonym="aniridia"							
			/note="Region: homeobox"							
ORIGIN										
	1	tatcgataag	ttttttttt	attgtcaatc	tctgtctcct	tcccaggaat	ctgaggattg			
	61	ctcttacaca	ccaacccagc	aacatccgtg	gagaaaactc	tcaccagcaa	ctcctttaaa			
1	121	acaccgtcat	ttcaaaccat	tgtggtcttc	aagcaacaac	agcagcacaa	aaaaccccaa			
1	181	ccaaacaaaa	ctcttgacag	aagctgtgac	aaccagaaag	gatgcctcat	aaagggggaa			
2	241	gactttaact	aggggcgcgc	agatgtgtga	ggccttttat	tgtgagagtg	gacagacatc			
З	801	cgagatttca	gagccccata	ttcgagcccc	gtggaatccc	gcggccccca	gccagagcca			
3	861	gcatgcagaa	cagtcacagc	ggagtgaatc	agctcggtgg	tgtctttgtc	aacgggcggc			
4	21	cactgccgga	ctccacccgg	cagaagattg	tagagctagc	tcacagcggg	gcccggccgt			
4	181	gcgacatttc	ccgaattctg	caggtgtcca	acggatgtgt	gagtaaaatt	ctgggcaggt			
5	541	attacgagac	tggctccatc	agacccaggg	caatcggtgg	tagtaaaccg	agagtagcga			
6	501	ctccagaagt	tgtaagcaaa	atagcccagt	ataagcggga	gtgcccgtcc	atctttgctt			
6	61	gggaaatccg	agacagatta	ctgtccgagg	gggtctgtac	caacgataac	ataccaagcg			
7	21	tgtcatcaat	aaacagagtt	cttcgcaacc	tggctagcga	aaagcaacag	atgggcgcag			
7	781	acggcatgta	tgataaacta	aggatgttga	acgggcagac	cggaagctgg	ggcacccgcc			
8	341	ctggttggta	tccggggact	tcggtgccag	ggcaacctac	gcaagatggc	tgccagcaac			
g	901	aggaaggagg	gggagagaat	accaactcca	tcagttccaa	cggagaagat	tcagatgagg			
9	961	ctcaaatgcg	acttcagctg	aagcggaagc	tgcaaagaaa	tagaacatcc	tttacccaag			
10	21	agcaaattga	ggccctggag	aaagagtttg	agagaaccca	ttatccagat	gtgtttgccc			
10	81	gagaaagact	agcagccaaa	atagatctac	ctgaagcaag	aatacaggta	tggttttcta			
11	141	atcgaagggc	caaatggaga	agagaagaaa	aactgaggaa	tcagagaaga	caggccagca			
12	201	acacacctag	tcatattcct	atcagcagta	gtttcagcac	cagtgtctac	caaccaattc			
12	261	cacaacccac	cacaccggtt	tcctccttca	catctggctc	catgttgggc	ctaacagaca			
13	821	cagccctcac	aaacacctac	agcgctctgc	cgcctatgcc	cagcttcacc	atggcaaata			
13	881	acctgcctat	gcaaccccca	gtccccagcc	agacctcctc	atactcctgc	atgctgccca			
14	41	ccagcccttc	ggtgaatggg	cggagttatg	atacctacac	cccccacat	atgcagacac			
15	501	acatgaacag	tcagccaatg	ggcacctcgg	gcaccacttc	aacaggactc	atttcccctg			
15	61	gtgtgtcagt	tccagttcaa	gttcccggaa	gtgaacctga	tatgtctcaa	tactggccaa			
16	521	gattacagta	aaaaaaaaaa	aaa						
17										

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

Sunday 4 November 2018

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

**blast** has positioned a gap in this region merely to maximize the overall alignment score. There is more than one way of achieving this simple goal. However, if it were to be recognized that the gap to be positioned was to represent an intron, then one of the arithmetically equivalent options becomes far more attractive than the others. This "best" option is not the one chosen by **blast**, which is forgiveable as **blast** had nor reason to expect an intron and was not written to understand the properties of introns anyway.

The alignment chosen for this region by blast was:

Genomic	24601	CAAAATAGATCTACCTGAAGCAAGAATACAGGTACCGAGAGACTGTGCAGTTTCACACTT	24660
mRNA	1097	CAAAATAGATCTACCTGAAGCAAGAATACAGGTA	1130
		• • •	
Genomic	24781	GGGAGGGCAGCAGTGGAGGTGCCAAGGTGGGGCTGGGCT	24840
mRNA			
Genomic	24841	CTGTCCCACCTGATTTCCAGGTATGGTTTTCTAATCGAAGGGCCAAATGGAGAAGAGAAG	24900
mRNA	1131	TITITITITITITITITITITITITITITITITITITI	1167

Shifting the gap 3 places to the left neither changes the size of the gap nor the perfection of the alignment either side of the gap and so does not affect the alignment score.

However, it does mean the gap begins with an **GT** and ends with a **AG** which is what one might expect if it were known that the gap represented an intron. I include the beautiful **Intron/Exon** logo. As you might gather, I rather like this one.



So, if **blast** was a little better informed, the improved alignment would have been:

Genomic	24601	CAAAATAGATCTACCTGAAGCAAGAATACAG <mark>GT</mark> ACCGAGAGACTGTGCAGTTTCACACTT	24660
mRNA	1097	CAAAATAGATCTACCTGAAGCAAGAATACAG	1130
		• • •	
Genomic	24781	GGGAGGGCAGCAGTGGAGGTGCCAAGGTGGGGCTGGGCT	24840
mRNA			
Genomic	24841	CTGTCCCACCTGATTTCCACGAGGGCCAAATGGAGAAGAGAAG	24900
mRNA	1131	GTATGGTTTTCTAATCGAAGGGCCAAATGGAGAGAGAGAAGAGAGAAG	1167

This is the alignment that one might expect from any program customized to align **mRNA** with **Genomic** sequence, as you will see in the fullness of time.

Mouel Allsweis			Sullua	ly 4 November 2016
How many convincingly aligned regions did you se	e?			
4	1,643			
How many did you expect?				
<b>12</b> , as that was how many <b>blast</b> found, not including the silly ones at the beginning.				U
The <b>4</b> that were found correspond the illustrated <b>4</b> diagonal lines grouped together in the <b>Dot Matrix View</b> made by <b>blast</b> .		/	1	
Clearly, this alignment is not correct. Can you explain why?	El Query_27637	10 K 15 F	с 20 К	25 K 28,741

This alignment algorithm only wishes to maximise an alignment score. It sees <u>*ALL*</u> 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 that 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.

#### Model Answers

# 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 THE top S ACATCTGGCTCCATGTTGGGCCTAACAGACACAGCCCTCACAAAC	equence is the mRNA. <b>splign</b> is kind enough to explicitly s that the "mutated" codon, <b>CTA</b> , will be expressed a
So, why not translate the <b>Genomic</b> sequence also <b>splign</b> ?! Easy enough to look up. But I resent having to do so! From this rather beautiful representation of the	Anginine de UC A G UC A
Genetic Code, I conclude: mRNA CTA $\rightarrow$ Leucine (L)	Asparagine G A A C U Aspartia
Genomic CGA $\rightarrow$ Arginine (R)	Eleucine A U U G A A A A A A A A A A A A A A A A
I checked, and this does not appear to be a substitution that is associated with any "interesting" phenotype.	TINT STOP UGACUGACUGACUGACUGACUGACUGACUGACUGACUGAC
There is no real reason why it should. We did not pause to find out anything about the mRNA downloaded from the <b>NCBI</b> . The annotation is	Sor STP Tyr Val

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!

Model Answers		Sunday 4 November 2018
What are the Genomic and mRNA base positions correspondin	g to the mutation	at amino acid position <b>33</b> ?
Demonstra the Network contraction of antiparticle and	Natural variant i (VAR_008694)	29 I $\rightarrow$ S in AN. <b>4</b> 1 Publication <b>•</b>
Remember the Natural variation at amino acid position	Natural variant i (VAR_003811)	29 I $\rightarrow$ V in AN. $\clubsuit$ 1 Publication $\checkmark$
<b>33</b> ? You looked at it in passing during the course of the first	Natural variant i (VAR_008695)	33 A $\rightarrow$ P in AN. $@$ 1 Publication $\checkmark$
exercise. It is a major cause of Aniridia. An Alanine	Natural variant i (VAR 008696)	37 - 39 Missing in AN. 🛛 1 Publication 🚽

Natural variant (VAR\_008697)

Natural variant i (VAR\_008698)

exercise. It is a major cause of Aniridia. An Alanine Natural variant<sup>1</sup> (VAR\_008696) mutated to a **Proline** at the end of a **Helix** vital to the **DNA** Binding function of the PAX6 protein.

#	Query		Subject	Spar	ı(bp)	Coverage(%)	Overall(%)	Exon(%)	CDS(%)	In-frame(%	b)
1	M77844.1(	+) pax6	_genomic(	+) 7618-	28741	99.03	98.78	99.75	0.00	0.00	
										— Graphics	Text
Mo	odel 1	Coverage Overall Exon	99.03% 98.78% 99.75%	CDS In-frame Primary trar	0 0 nscript 1	.00% Mism .00% Exon 627 bp Intro	natches and in is (min/max/a ons (min/max/	dels 4 ve), bp 61 ave), bp 99	/ 218 / 13 / 5903 / 1	5 773	
M7	7844.1 (+) He	omo sapie	ns oculorho	mbin (PAX6)	mRNA,	complete cds, al	ternatively sp	iced			
								1			
1 pax	6_genomic (·	+) dna:chr	omosome o	hromosome	GRCh38	8:11:31784179:3	1818662:-1			1643	
761	18	-		-	•					28741	
Se	gments	Alignm	ent								
<mark>1</mark> 7 13	23456 89101112	Unalig le st st	ined: ingth = art = iop =	16 1 16							

Natural variant i (VAR\_003812) 44 R  $\rightarrow$  Q in AN.  $\blacksquare$  1 Publication  $\neg$ 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.

42 I  $\rightarrow$  S in AN; mild.  $\checkmark$  1 Publication

43 S  $\rightarrow$  P in AN.  $\clubsuit$  1 Publication  $\checkmark$ 

You should be able to clearly see the extra mutation is in the 5<sup>th</sup> segment.

Focussing on the 5<sup>th</sup> 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<sup>3</sup> amounts to:

# Affected Patient protein:

CCT  $\rightarrow$  Proline (P)

**Canonical protein:** 

GCT  $\rightarrow$  Alanine (A)

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

#	Query	Subje	ect Spa	an(bp) Co	verage(%)	Overall(%)	Exon(%)	CDS(%)	In-frame(%)	
1	M77844.1(-	+) pax6_geno	mic(+) 7618	-28741	99.03	98.78	99.75	0.00	0.00	
										Graphics Tex
м	odel 1	Coverage Overall Exon	99.03% 98.78% 99.75%	CDS In-frame Primary	0.0 0.0 transcript 16	00% 00% 27 bp	M Er In	ismatches cons (min/i itrons (min	and Indels max/ave), bp I/max/ave), bp	4 61 / 218 / 135 99 / 5903 / 1773
M	77844.1 (+) Ho	omo sapiens ocul	orhombin (PAX	5) mRNA, com	plete cds, all	ternatively spl	iced			
	11						1			
1 pa	x6_genomic (+	-) dna:chromoso	me chromosom	e:GRCh38:11	:31784179:3	1818662:-1			1643	
•								<b>—</b>		
76	18								28741	
S	gments	Alignment								
1	234 <mark>5</mark> 6 89101112	373	S GTC	H S G ACAGCGG	V N	Q L TCAGCTC	G G GGTGGT	V F GTCTT	V N TGTCAACG	G R P L P D S T GGCGGCCACTGCCGGACTCCACC
1	3	15824	CTCAGGTC	ACAGCGG	AGTGAA	TCAGCTC	GGTGGT	GTCTT	TGTCAACG	GGCGGCCACTGCCGGACTCCACC
		438	R Q K CGGCAGAA	I V GATTGTA	E L GAGCTA	P H CCTCACA	s g GCGGGG	A R CCCGG	P C D CCGTGCGA	I S R I L Q CATTTCCCGAATTCTGCAG
		15894	CGGCAGAA	GATTGTA	GAGCTA	GCTCACA	GCGGGG	cccee	ссбтбсба	CATTTCCCGAATTCTGCAGGTGA
		504								
		15964	т							

**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!

Proving beyond reasonable doubt that substitution is exactly at amino acid position 33 requires a little more counting, dividing by 3 and subtracting the number you first thought of. For now, just trust me? I really am more honest than I look.

### Summary:

The <b>Details</b> column shows the alignments		Alignment transcript represents full details of the alignment in a form of a string
of each exon in a compressed format	11. Alignment	composed of characters 'M', 'R', 'I' and 'D' where each character corresponds to an
described in the <b>splign</b> documentation as	transcript	elementary command (Match, Replace, Insert or Delete) needed to transform the query
illustrated		segment into the subject segment. The string is encoded with RLE.

The majority of the exon alignments are trivial.

#	Query	Subject	Span(bp	o) (	Coverage(%)	Overall(%	b) Exon(%)	CDS(%)	In-frame(%)	
1	M77844.1(+)	<pre>pax6_genomic(+)</pre>	7618-287	741	99.03	98.84	99.82	0.00	0.00	
										Graphics   Text
#	Query	Subject	Idty	Len	Q.Start	Q.Fin	S.Start	S.Fin	Туре	Details
+1	M77844.1	pax6_genomic	-	16	1	16	-	-	<l-gap></l-gap>	-
+1	M77844.1	pax6_genomic	0.991	218	17	234	7618	7835	CA <exon>GT</exon>	M39RM8RM169
+1	M77844.1	pax6_genomic	1	77	235	311	11738	11814	AG <exon>GC</exon>	M77
+1	M77844.1	pax6_genomic	1	61	312	372	12201	12261	AG <exon>GT</exon>	M61
+1	M77844.1	pax6_genomic	1	131	373	503	15829	15959	AG <exon>GT</exon>	M131
+1	M77844.1	pax6_genomic	1	216	504	719	16887	17102	AG <exon>GT</exon>	M216
+1	M77844.1	pax6_genomic	1	166	720	885	17807	17972	AG <exon>GT</exon>	M166
+1	M77844.1	pax6_genomic	1	159	886	1044	23875	24033	AG <exon>GT</exon>	M159
+1	M77844.1	pax6_genomic	1	83	1045	1127	24549	24631	AG <exon>GT</exon>	M83
+1	M77844.1	pax6_genomic	1	151	1128	1278	24861	25011	AG <exon>GT</exon>	M151
+1	M77844.1	pax6_genomic	0.991	116	1279	1394	25110	25225	AG <exon>GT</exon>	M33RM82
+1	M77844.1	pax6_genomic	1	151	1395	1545	27803	27953	AG <exon>GT</exon>	M151
+1	M77844.1	pax6_genomic	1	98	1546	1643	28644	28741	AG <exon></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 were there are some disagreements. That is, the entries for Exons 1 and 5, which, following the documentation, I translate thus:

# <u>Exon 1 – M39RM8RM169</u>

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**).

### <u>Exon 10 – M33RM82</u>

An alignment of 116 bases, the first 33 of which Match perfectly (M33), there them 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.

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.

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:

#### 

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 here<sup>4</sup>.

I think it worth mentioning, that this way of representing an alignment is a simplification of CIGAR format<sup>5</sup>. This format is used for SAM (Sequence Alignment Map) and BAM (Binary Alignment Map, exactly the same as

SAM, except compressed) files. You will CIGAR: CIGAR string. The CIGAR operations are given in the following table (set '\*' if unavailable): be engulfed in SAM/BAM files if you ever do any Next Generation Sequencing (NGS)<sup>6</sup>.

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

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)
Н	5	hard clipping (clipped sequences NOT present in SEQ)
P	6	padding (silent deletion from padded reference)
=	7	sequence match
X	8	sequence mismatch
can only be p	resent as	the first and/or last operation.
may only have	e H opera	tions between them and the ends of the CIGAR string.
or mRNA-to-g	enome a	lignment, an N operation represents an intron. For other types of all

• Sum of lengths of the M/I/S/=/X operations shall equal the length of SEQ.

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'?

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.

NGS is also referred to as High Troughput Sequencing (HTS), which, on the whole, I think is a more meaningful name.

**Basic Bioinformatics** 

18 of 23

<sup>4</sup> The Wikipedia article shows how this complication might be overcome.

Μ

0

# Where is the **3**<sup>rd</sup> substitution in the mRNA?

Where is the  $3^{rd}$  substitution in the Genomic Sequence?

splign makes one work quite hard to answer this one! Unless I am missing something.

From the alignment of **Exon 10**, the exon including the **3rd R**eplacement, with a bit of squinting, it can be confirmed that the **3<sup>rd</sup> R**eplacement is at:

S S S Ν S D A Т 1279 тостос TCACATCTGGCTCCATGTTGGGCCTAACAGACACAGCCCTCACAAACACCTACAGC 25105 CTCAGTTTCCTCCTTCACATCTGGCTCCATGTTGGGCCGAACAGACACAGCCCTCACAAACACCTACAGC

Base pair position 1,312 of the mRNA

Base pair position 25,143 of the genomic sequence

It might also have been relevant to ask which amino acid position corresponded to the Replacement. To discover this one would need to look at the alignment of **Exon 3**, where the coding begins.

Ν CCCGTGGAATCCCGCGGCC 312 . . . . AG A G A G C C A G C A T G C A G A A C A 

More squinting, and I conclude the A of the ATG representing the initial Methionine of the protein coding region is at position 363. That is, the 5' UTR ends at position 362. So the Replacement is at:

Base position 1312 - 362 = 950 of the protein coding region of the mRNA.

As 950 / 3 is 316 remainder 2, the Replacement is at codon position 2 of the 317<sup>th</sup> amino acid of the protein.

Cannot help thinking that **splign** might have helped a bit more here?

I also reflect that I cannot fully recall why I wanted to know where the mutation was, especially given we have decided to reject any chance that it might be a mutation of consequence. Oh well, some things a man must do, just because they are there to be done!!

Time to move on ... without checking my arithmetic. Bound to be right, I used to be a mathematics teacher you know! Several lifetimes ago.

#### **Postscript:**

After the passage of many months, I now recall why I obsessed as to the position of this amino acid substitution. I wondered if it was in the region of one of the major domains of this protein. If it was, it might increase its chances of being significant?

Well, it is not. In the last exercise, we discovered that:

The Paired-box domain is between positions 4 and 128 (Consensus isoform) or 4 and 142 (isoform 5a).

The Homeo-box domain is between 214 and 266 (Consensus isoform) or 228 and 280 (isoform 5a).

So the Substitution, at position 317, is in a relatively neutral region and so, maybe, less likely to be of great consequence?

Model Answers	Sunday 4 November 2018
Compare the predicted <b>splign</b> intron/exon boundaries with the What deviation(s) from the model suggested by the logo can be	conservation suggested by the logo?
You may have gathered, I rather like this logo, although I it is leading me to make the same point a trifle to often?	rather think exon 5'
The logo is in almost 100% agreement with the predictions	s of <b>splign</b> .
As you will have noted previously, when looking at the	ne Ensembl
human gene (previous Practical), there is a single exception	
Type <l-gap> CA<exon>GT The easiest way to show this in the splign outp</exon></l-gap>	ut is to look
AG <exon>GT</exon>	acceptor
AG <exon>GT AG<exon>GT alignments it predicts. It also records <b>2 flanking</b></exon></exon>	the <exon> intron base 3, exon</exon>
AG <exon>GT</exon>	
AG <exon>GT AG<exon>GT</exon></exon>	
AG <exon>GT AG<exon>GT AG<exon>GT AG<exon> <math>df</math> the end of the <math>2^{nd}</math> exon. Here there is GC rather</exon></exon></exon></exon>	on deviates from the model suggested by the logo is at then <b>GT</b> . Well, nothing is perfect!

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 - 2017.12.23

# **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 if 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 – 2018.11.04**