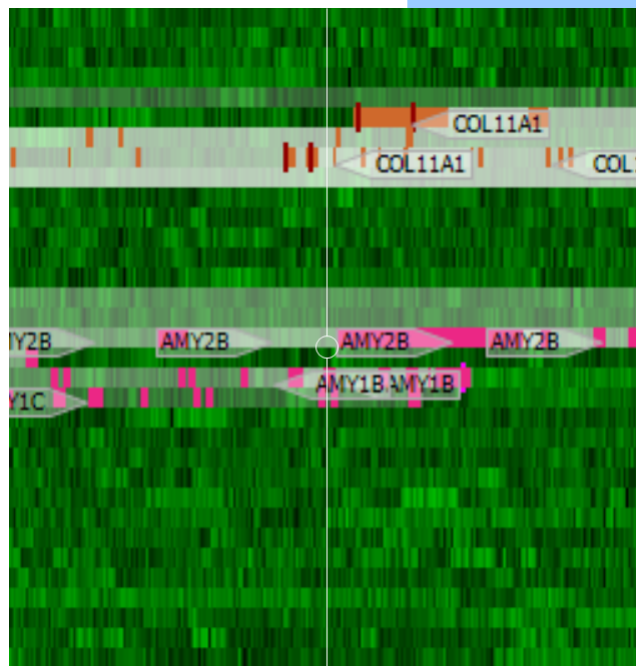# Visual Genome Browser (Beta version)

*This document outlines the Use Cases which led to the development of the offline genome browser software.*
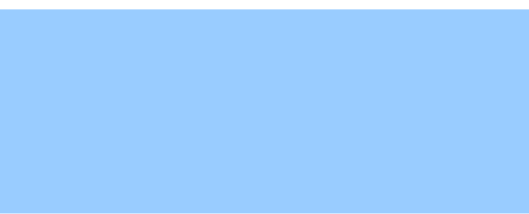
**Heinrich Ferreira**

*29 November 2016*

*heinrichjf@gmail.com*

# Table of Contents

# Genome Browser Concept

## *Introduction:*

I came up with a concept of visualizing the human genome (or the genome of any other organism for that matter) in a 2 dimensional space while playing around with the file formats of genome data provided for free online by the **University of Santa Cruz.** (https://genome.ucsc.edu/cgi-bin/hgGateway).

Being a programmer I wanted to investigate the Human Genome for myself. I wanted to get a feel for the nature of the digital data contained within it. I learned that it contained a code for producing protein machines such as enzymes as well as structural and control proteins such as transcription factors, enhancers and suppressors. I was also intrigued by the repetitive patterns in the genome in regions such as the centromeres and telomeres of chromosomes. The standard way of representing genes on a linear axis simply did not elucidate the structure sufficiently for me. I learnt that the genome had a complex 3D organization in the nucleus of the cell and thought that that there had to be a way to get the a better picture of the structure of the data of the genome.

Coming from the field of GIS (Graphical Information Systems) I knew that one could better represent big amounts of map data in 2 dimensions with multiple layers of information overlaid on top of each other. I wanted to do the same for the genome. ***In no way do I claim that I am the first to think of doing this, there might be other software who have attempted to do just this, but I wanted to get an intuitive "feel" for the data***. I decided to find a way of representing the chromosomes in as compact a way as possible by depicting the linear DNA letters by coloured blocks running from left to right, top to bottom akin to the scan lines of a television screen. If there were any repetitive patterns in the data it should become visible when you use the correct line width. Any repeat occurring at an interval of x bases, will show up when you line up the sequences at a width of x bases.

I needed a metric to us that would sufficiently change throughout different parts of the genome and I decided upon the GC content. I needed to map that metric onto an RGB (Red Green Blue) colour scale and picked the scale of Green colour as it reminded me of GFP (Green Fluorescent Protein) which is often used in fluorescent reporter assays. I found that if I simply mapped nucleotides to pixels, the result simply looked like noise, so I had to apply averaging over bases to get more uniform and visually identifiable "blocks" of a certain colour intensity. Taking the average number G'c and C's in blocks of about 20-50 bases and then mapping it to the 0-255 green intensity provided me with a structure of the genome. Just like origins of replication often occur in areas of higher A's and T's, where the DNA helix is more easily opened up by a helicase for DNA replication to start, I could now distinguish areas of the genome exhibiting distinct GC Content structure.

The best way to explain this is probably to go through a few examples showing how these features are highlighted by visual inspection.

For most of the examples I have obtained data from the UCSC Genome Browser.

https://genome.ucsc.edu/cgi-bin/hgGateway

HG38 **.2bit** download folder: ftp://hgdownload.cse.ucsc.edu/goldenPath/hg38/bigZips/

## *Some examples of distinct GC structure in the genome.*

Here is an extract from the Human X-chromosome.



*Illustration 1: ChrX:152,300,809-153,887,660 at 39672 bases per horizontal line and 29 bases per horizontal pixel.  Green intensity represents average GC content in each pixel.*

The data has a non-repetitive (almost random) structure.  When I map the genes that I downloaded via MySQL from the table data of the UCSC, I find that there is a high concentration of genes in this region.



*Illustration 2: ChrX:152,300,809-153,887,660 at 39672 bases per horizontal line and UCSC genes drawn as blocks of different colours*

I mapped the gene annotations as an overlay on top of the GC Content with small name arrows indicating the direction in which the genes are read during transcription by the RNA Polymerase. In this case I have not separated the exons of the genes, therefore the coloured regions represents the primary transcripts before it is spliced to remove the introns.

A similar **CpG Island output** can be obtained by redrawing the genome (with the "Recreate" option and selecting "CpG Islands" under the "Controls" tab). In the next picture it can be seen that the promoters of genes are mostly found near CpG Islands in the chromosome.
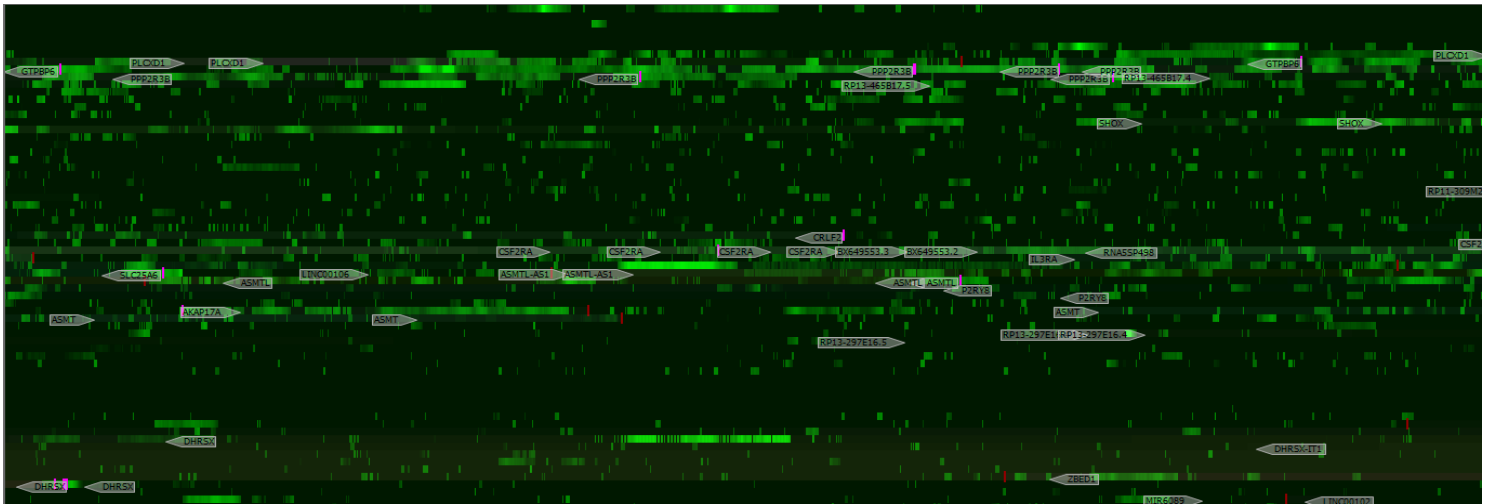
*Illustration 3: ChrX:152,300,809-153,887,660 at 39672 bases per horizontal line showing CpG Islands with gene annotations only indicated with arrows*

Gene annotations data are downloaded by the software using a MySql query from the UCSC table data (`genome-mysql.cse.ucsc.edu`) or obtained directly from GFF3, Genbank, Ensemble, BED, BigBed, or BigWig files. Annotations can also be obtained from the UCSC search results copied from the web.

When you compare the image above with another part of the genome closer to the centromeres, where highly repetitive patterns are seen, the region is devoid of annotated genes.



*Illustration 4: Chr X:59587374 Inside the centromere with highly repetitive alpha repeats*

The picture shows 39672 bases for each horizontal line, and each horizontal pixel represents 29 nucleotides. This means one can get maximum 29 intensities of green. (This value is constantly adjusted per chromosome in order to optimally fit in the information into the fixed pixel width of the display.) All chromosomes are also limited to 32000 vertical display lines (with about 4 vertical pixels per coloured block).

In order to have a closer look at the finer structure of the bases, the Genome Browser provides 2 detail view options for the **DNA View**: (The block at the centre vertically and horizontally always represents the highlighted base position)

- **The default DNA View depicts a continuous extract of bases from the genome at the position of the mouse in the Main Genome View. The genome bases flows from left to right, top to bottom.**

- **An "Aligned" view (representing a further zoom level above the Zoom Gene View) where a cut-out is made for bases spaced every 39672 bases. This allows you to alter the line width using the < and > keys in order to investigate repetitive patterns.**

Each base of the DNA sequence in this region is depicted by a colour:
T = Black
A = Dark Blue
G = Yellow
C = Orange

I deliberately chose dark colours for paired bases with only 2 hydrogen bonds (representing regions which are more easily melted apart) and bright colours for paired bases with 3 hydrogen bonds (representing regions more difficult to separate).

Just using visual observations, one can start to identify the structure of the chromosomes. The structure of the chromosomes has bearing on how densely packed the chromosomes are. Densely packed chromatin / constitutive heterochromatin as is found in the centromeres also tend to stain darker.

When I draw a closer, zoomed in version of this pattern at a width of 171 bases per line, one gets the following close up view (which is only hinted at when you see it in the first image, but upon "closer inspection" becomes clear). One can see a vertical band of blue and black bases, which is probably the bases in the linker section between the nucleosomes, which is more exposed to mutational deamination of C's to U's.



*Illustration 5: ChrX:59801335-59830917 at width 171 bases per line (alpha satellite repeats in the heterochromnatin of the centromeric region)*

Colouring of the bases actually has a bearing on the temperature at which the DNA could be melted apart. I would therefore expect that organisms which live at very high temperatures would have a genome which consisted mostly of brightly coloured bases (such as the Termus aquaticus thermophilic bacteria living in the hot springs of Yellowstone National Park (from which we obtain the TAQ Polymerase used in PCR (Polymerase Chain Reaction) assays) and Deinococcus radiodurans (an extremophilic bacterium highly resistant against DNA damage).

Here is an extract from the genome of D. radiodurans.



*Illustration 6: Extract from Deinococcus radiodurans genome*



*Illustration 7: GC Content of 50400 bases averaged at 20 bases per block of D. radiodurans*

Looking at the GC Content of a larger region of the genome of this bacteria (50400 out of a total of 2,648,638 of its circular positively supercoiled chromosome), where the green light intensity is equated to the average GC content, reveals a genome with high GC content.
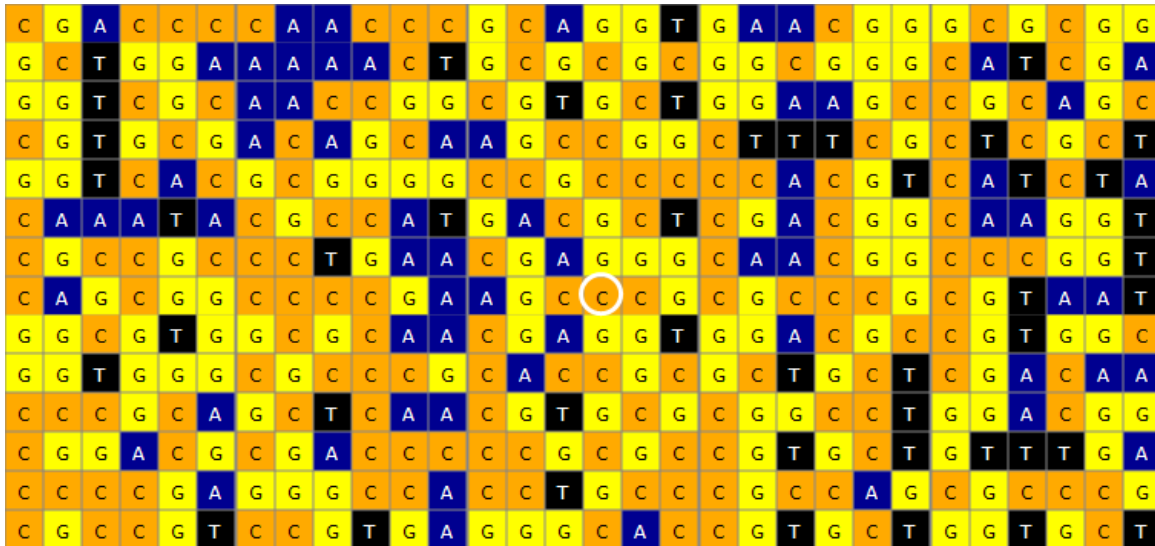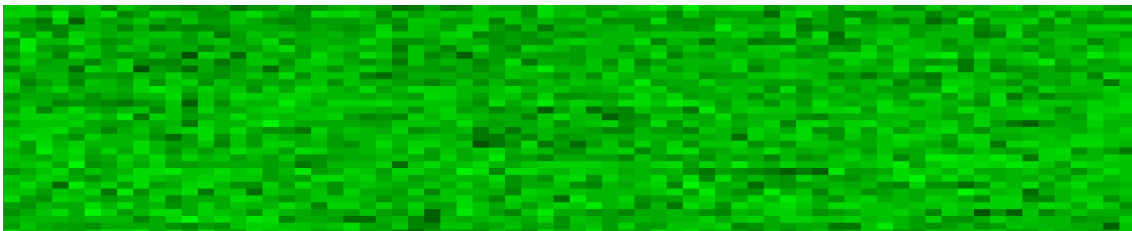
This is clearly an organism suited to replicate its DNA at higher temperatures, and the positively supercoiled DNA helps prevent its DNA from unwinding easily.

When we look at a similar size length of the thermophilic bacteria Termus aquaticus, we get the same result:



*Illustration 8: GC Content of 56,160 bases averaged at 20 bases per block of T. aquaticus*

Doing the same for other bacteria living at much colder temperatures, such as Yersinia enterocolitica, which can live at temperatures of less than 0 degrees, I find a genome with a much lower GC content and a nucleotide bias towards more A's and T's. Its DNA needs to be replicated at lower temperatures, with less energy available to helicases to pull apart the DNA strands. Because there are less energy available at these low temperatures, it would not survive if the DNA strands required a lot of energy to be separated.
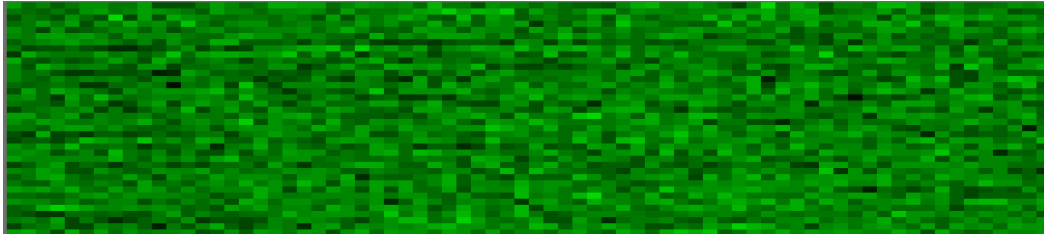
*Illustration 9: GC Content of 56,160 bases averaged at 20 bases per block of Y. enterocolitica*

The lightness of the green squares of the thermophilic bacteria when compared to bacteria which lives at colder temperatures is a striking indication of how organisms are suited for their environment.

Interestingly enough, when you create the same GC Content representation of the pathogenic bacteria causing tetanus (which can survive in soil for extended periods of time), it reveals a genome which is skewed towards a much lower GC content. Suggesting it were well suited to survive in colder temperatures.



*Illustration 10: GC Content of 56,160 bases averaged at 20 bases per block of Clostridium tetani (the cause of tetanis)*

(Gene-centric association analysis for the correlation between the guanine-cytosine content levels and temperature range conditions of prokaryotic species: *It has been found that the environment that a bacteria lives can be predicted with >80% accuracy by looking at the GC Content percentage of its genome.*)
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3024870/

When I use the coloured representation to display the DNA of D. radiodurans, in particular the section of the DNA coding for the DNA-directed RNA-Polymerase in both strains of bacteria, it looks as follows:



Illustration 11: Chr1:923733-925271 of D. radiodurans DNA-directed RNA polymerase beta subunit



Illustration 12: C. tetani of DNA-directed RNA polymerase beta subunit

Interesting how both organisms has to build an enzyme with exactly the same function used for transcription, but utilizing codons picked f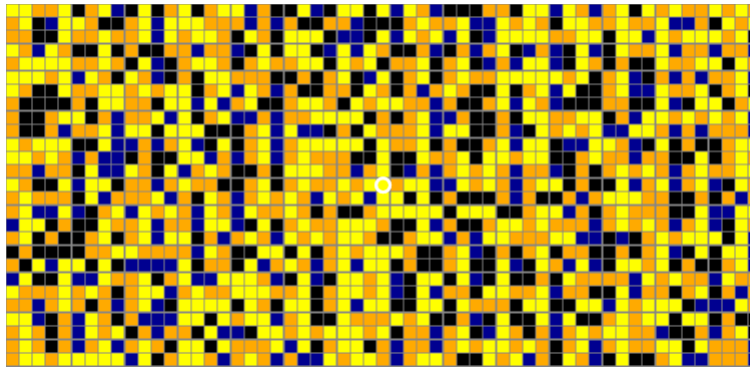rom a different GC content DNA. Interestingly, when you look at all 61 possible codons of the Genetic Code, one finds that for each possible amino acid, there are equal amounts of codons coding for each amino acid that uses A's and T's as there are that uses G's and C'c, and this is divided equally among the amino acids.

I specifically designed this software to allow me to vary the number of bases depicted in each line of the output. Any repetitive pattern could then be discovered easily by simply changing the width until a pattern became evident when you looked at it. The human visual system is sometimes able to discriminate patterns much more easily than sophisticated software.

One such an example is when you look at the alpha satellite repeats found at the centromeres of human chromosomes.

Here is a 2D display of the centromeric region HG38 chr1:122939341-122970111 (Width=170 bases) The CENP-B Box region is highlighted in Magenta.
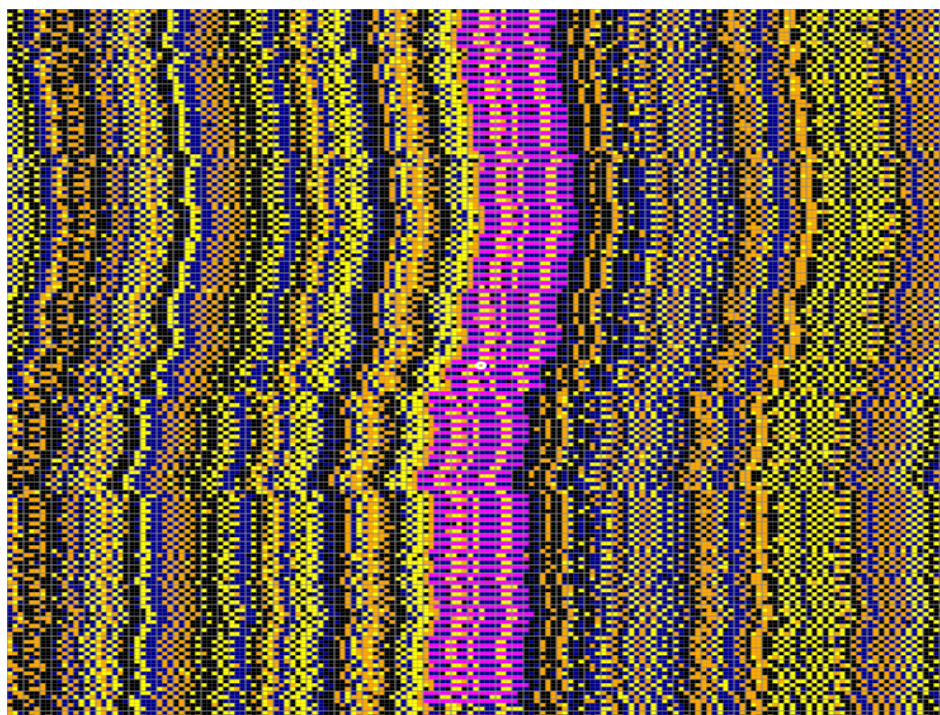


*Illustration 13: chr1:122939341-122970111 at width 170. Clearly showing the minimum width repeat. Wider line widths simply show the repeat multiple times in each row.*

The magenta highlight is part of the pattern search and "probe" functionality of the software which is the equivalent of using fluorescently tagged DNA probes to find specific target sequences via hybridization. However in this software I am using the Boyer-Moore string search algorithm, which I have modified to account for wildcards and multiple possible bases in each sequence position. It also allows you to specify the number of base mismatches it must allow when searching in the genome sequence.

One simply has to enter a search for the **Cenp-B box DNA sequence**:

CTTCGTTGGAAACGGGA

Allowing 2 mismatches it can be found almost 9000 times ONLY in the centromere of the X-Chromosome.

CpG methylation of the CENP-B box reduces human CENP-B binding (http://onlinelibrary.wiley.com/doi/10.1111/j.1432-1033.2004.04406.x/full)

If one wants to search for the sequence, but with either a T or an A in specific positions, it can be put in square brackets:

C*[TA][TA]*CG*[TA][TA]*GG*[AT][AT][AT]*CGGG*[AT]*

Any DNA sequence entered into the **Search field** such as this will be highlighted in the positive and negative (reverse complement) direction in the DNA View.

Additionally, the **"Find Pattern"** button can be used to search the entire genome through multiple selected chromosomes in parallel (threads).

The **Main Genome View** then displays these findings as differently coloured blocks in the **view** as

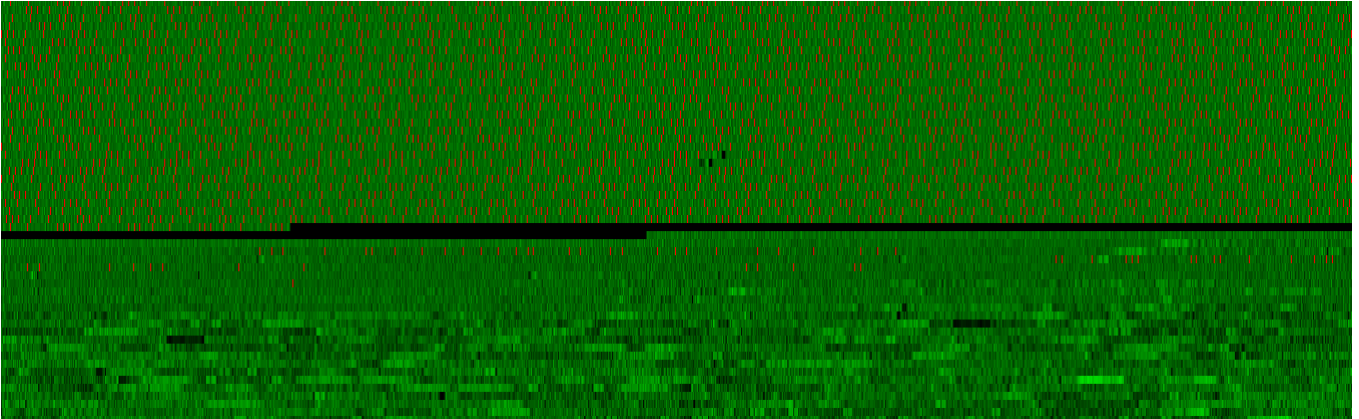is illustrated below. One can clearly see where the centromere ends.



*Illustration 14: Depiction of all the CENP-B box locations near the end of the X-Chromosome centromere.*

From the literature it is known that the DNA binding proteins which are responsible for locating the CENP-B box DNA sequence within the centromeric DNA are part of the protein complex which form part of the kinetochore assembly in mammalian chromosomes. The 170-171 bases is known from literature to represent the number of bases between subsequent nucleosomes (containing the CENP-A variant of the H3 histone) around which the DNA is wrapped to produce the heterochromatin in the centromeres.

Any repetitive pattern is easily visualised by incrementally changing the bases per line until at 170-171 a pattern immediately emerges.

**Looking at exactly the same region, but using a different line width of 80 bases, reveals no identifiable pattern**, which indicates that it is easy to miss these patterns unless you vary the line width at a specific location until repetitive bases start to line up, which is easily spotted by inspection.
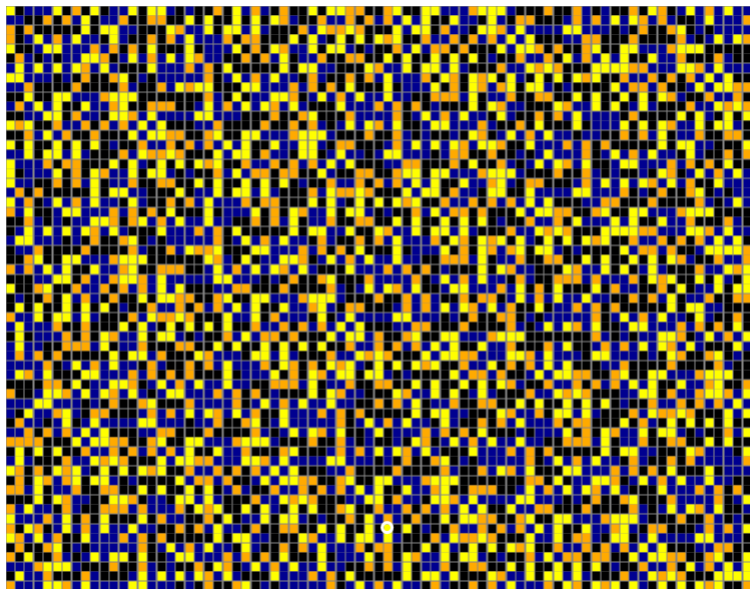


*Illustration 15: chr1:123097253-123105973 at width 80.*
*At most line width there are no discernible patterns.*

Other non-centromeric regions of the genome also contains patterns, such as this one found at 114 bases width on chromosome 2.
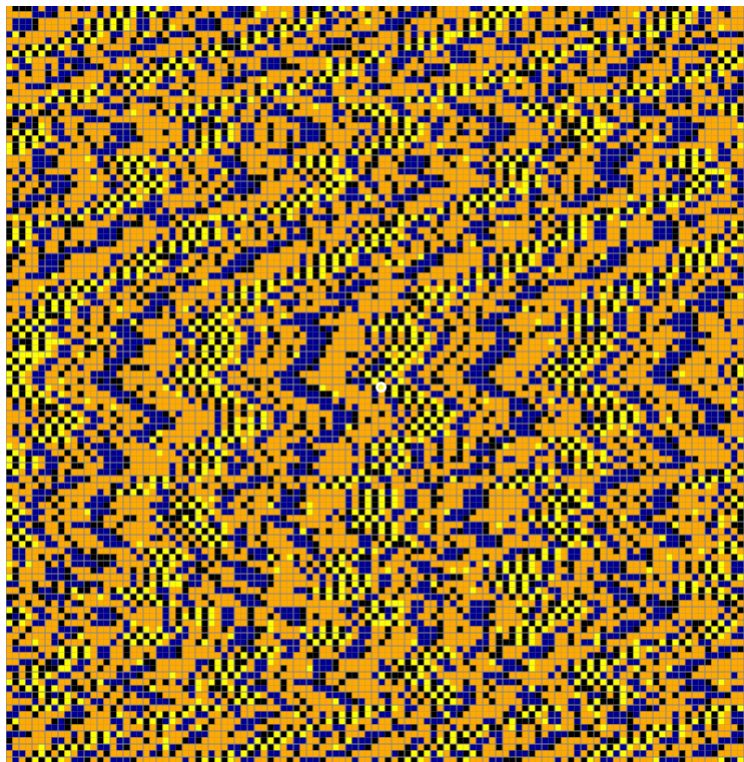


*Illustration 16: chr2:1524735-1538073 at width 114*

The above repeat is found in an intron of the thyroid peroxidase (TPO) gene.

The point that I wish to make is: It is possible to use the Visual Genome Browser's ability to represent the DNA sequence in a 2D matrix in order to locate and visualise many inherent structural features of the genome which is not always evident by just looking at a one dimensional representation of it.
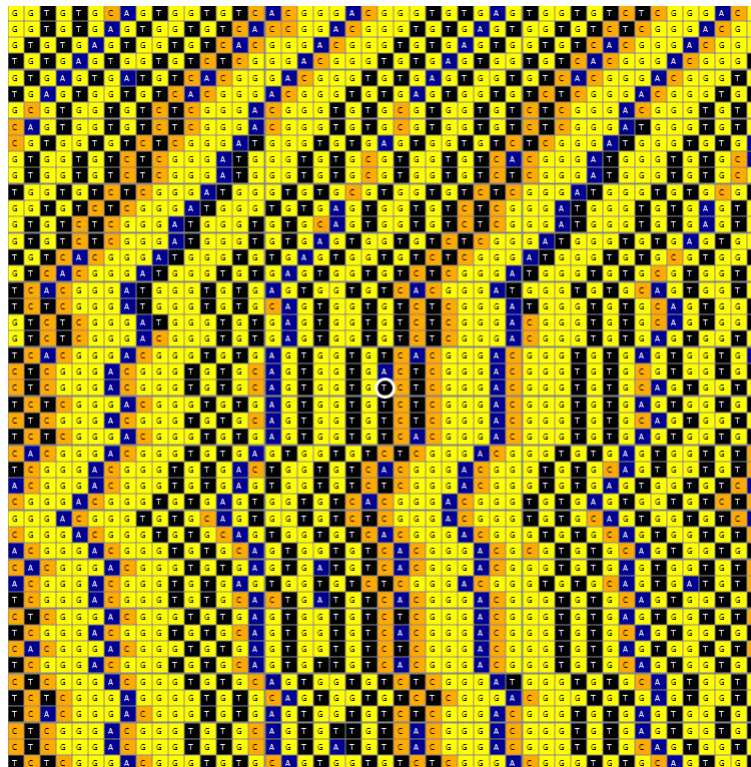


*Illustration 17: Pattern at chr16:857351-859703 with width 47 bases per line*

Other DNA sequence structures that can be identified are the telomeric sequences of chromosomes which consists of short repeats determined by the RNA template which is part of the DNA Telomerase enzyme responsible for lengthening the ends of chromosomes to overcome the Hayflick limit whereby the number of cell divisions are limited due to the shortening of the ends of linear chromosomes during replication.
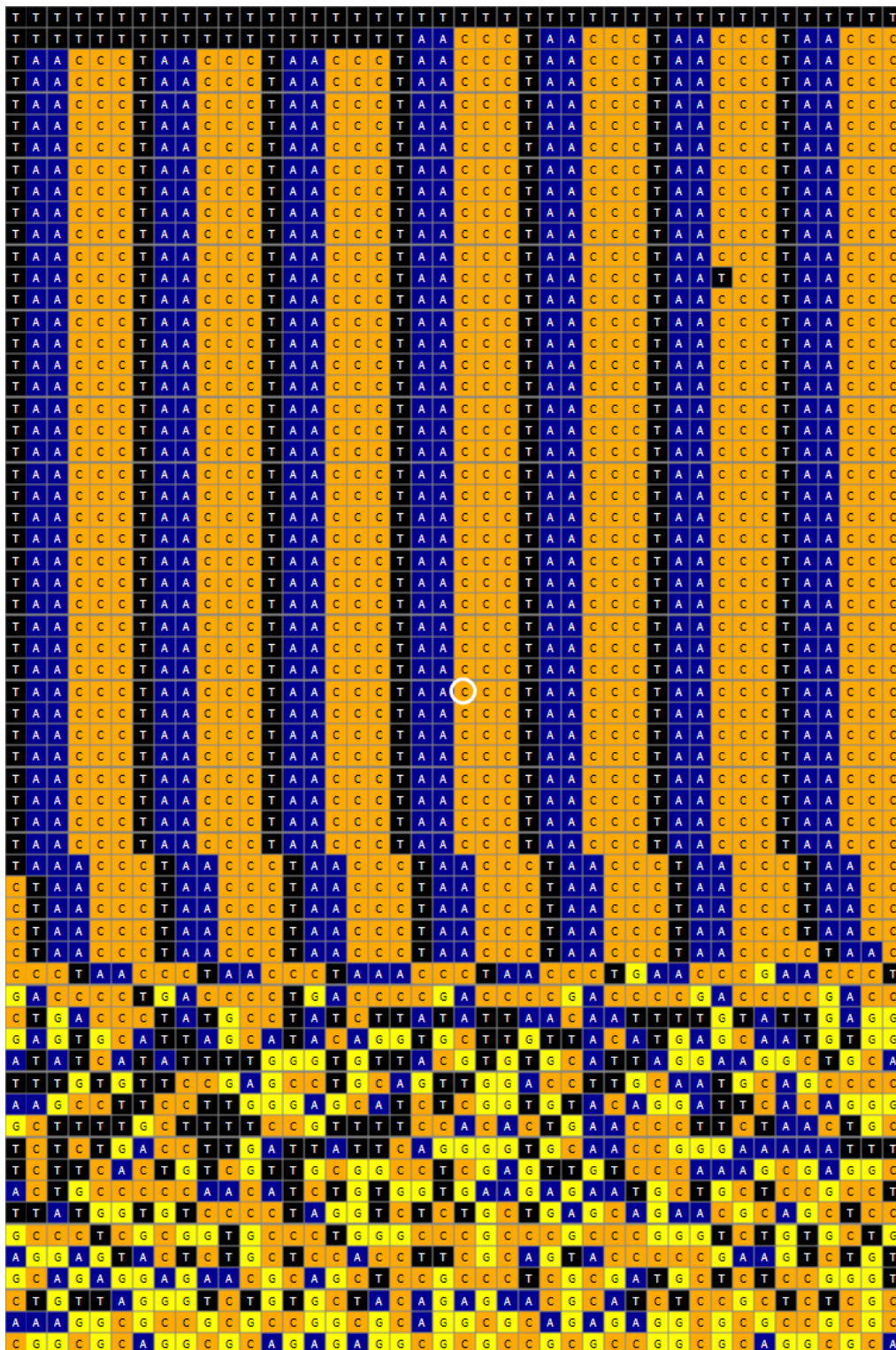


*Illustration 18: chr5:9941-12587 at width 42*

At a line width of 42 (7 times the telomerase template width of 6) one can clearly see the highly repetitive TAACCC repeats at the end of chromosome 5.

Looking at the opposite end of chromosome X, one can see a similar pattern, but this time reverse complemented.



*Illustration 19: chrX:156029066-156031046 at width 36*

The reverse complement GGGTTA of TAACCC can be found at the 3' end of the X Chromosome, this time repeating 6 times due to the width of 36. These patterns are known to be found in the RNA template which is part of the Telomerase protein-RNA complex.

## Possible sources of genomes and sequence data for the Genome Browser

At its core, the Visual Genome Browser stores the genomic sequence data in the **.2bit** sequence format, which is the most compact (yet still uncompressed) format that you can use to represent DNA data. In essence this means that 2 bits of each byte is used to store a single DNA letter. Since you can represent $2^2 = 4$ possibilities using 2 bits. This means one can fit in $8/2 = 4$ DNA letters in each byte (8 bits) of data. The complete version human reference genome (HG38) can therefore by represented by 796 MB of information.

All of the genome data used by the Visual Genome Browser are maintained in a Data Folder which contains the genomes of different organisms in sub-folders. When you want to add a new organism's genome to the browser, you simply copy the downloaded **.2bit** genome into this Data Folder. The Browser will then scan this folder at start-up for all the **.2bit** files and add it to the list of genomes.

The UCSC maintains the gene annotation data in a MySQL database which has to be queried by the genome browser in order to download the gene annotations which are overlaid on the genome structure.

The software also scans the Data Folder for other file types which may contain DNA sequence data, such as FASTA, GENBANK and ENSEMBL browser files. These files need to be put in folders prefixed with the Text : "Fasta", "Genbank" or "Embl". It typically examines the headers of these files to determine if it is of the specified type. Fasta files generally do not contain annotation data, but this can be loaded separately from **BED**, **BigBED** or **GFF3** files after the DNA sequence is displayed. The files located in each of these specially prefixed folders are then combined into a single **.2bit** files for quick and uniform access (one **.2bit** file for each folder involved). If any files gets added or removed from these folders, the .2bit and the annotation data (from Genbank and Ensembl files) is also rebuilt. Another way to obtain annotated sequences is to enter the Genbank accession number into a field and click a download button. Single Genbank, or a range of Genbank files are then downloaded from the

Url:

https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=nucleotide&id={0}&rettype=gbwithparts

Each **.2bit** file corresponds to a single "folder", which may contain a list of possible sequences or "chromosomes". When you select a specific folder, you are able to search through all possible fields of all chromosomes or sequences of that folder. You can then go and pick a specific chromosome to display in the 2 dimensional **Main Genome View**. The browser loads all annotations for this sequence into a memory data structure (range tree) which allows for very quick overlap checks between genomic ranges in order to overlay the correct annotations on top of the GC Content display of the structure.
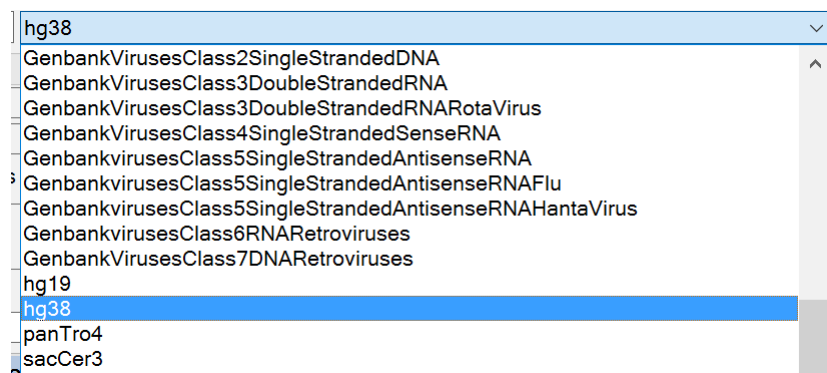


*Illustration 20: Folder List*

When you select a folder, the sequence or chromosome list is populated with sizes for the sequences



*Illustration 21: Sequence or chromosome list*



# The Different Views

Clicking on "DRAW/REDRAW" will load the 2bit data and display it in the **Main Genome View.**

If you also want to load the gene overlay, you click "LOAD GENES". The chromosomes selected in the sequence list always determines which sequences' annotations are loaded or searched for genes/patterns. The selected chromosome is then loaded into the Main Genome Browser screen.



*Illustration 22: Main Genome Browser Screen*

All available annotation data are kept in indexed binary files making them searchable and allowing fast lookup of genes from a single text input (similar to the Google predictive search, which gives you suggestions as you start typing the gene names).
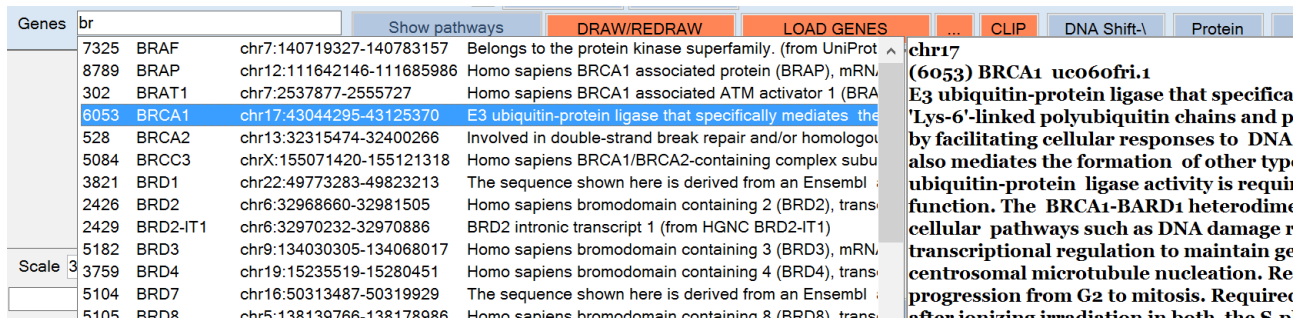


*Illustration 23: Predictive gene lookup*

This is the default field which has the focus when the software starts up, allowing the user to start a gene lookup immediately. (This field can also be reached by pressing **Ctrl-G**)

As you change the selection, a pop up window shows pertinent information regarding the selected gene. When you press enter, the software will locate the gene and highlight it in the corresponding navigable views.

In addition to the quick lookup, one can also do an exhaustive search through all or selected fields in the annotation data. When (**Selected Chr**) is chosen, only the checked sequences will be searched, otherwise all the index will be used to search through all of the sequences in the current folder.
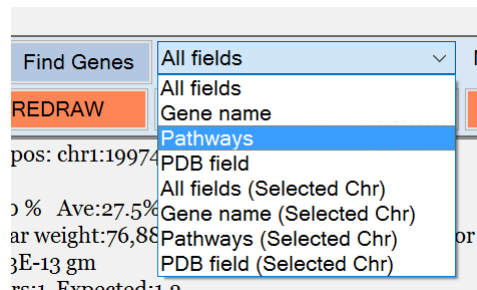


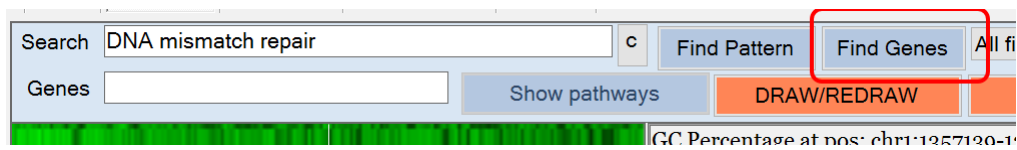*Illustration 24: Fields that can be searched*



*Illustration 25: Searching all fields*

The gene results are then given in the **Gene Search Results** tab as shown in the following picture.

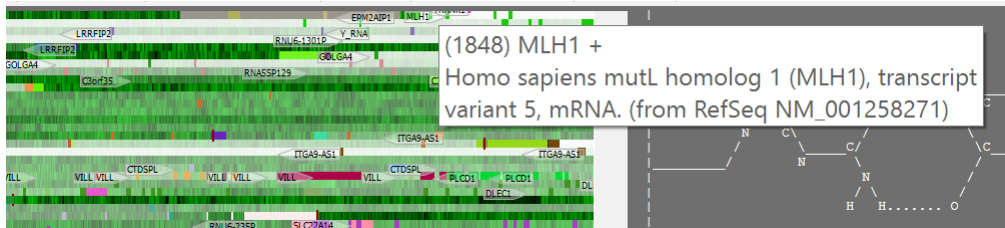When you double-click on any gene, the software highlights the selected gene in the **Main Genome View**.

*Illustration 26: Selected gene in Main Genome View*


*Illustration 27: Search options*



| Main Genome View | Chromosomes View | Gene Search Results (66) | Pathways | Information | Sequencing Gel | Protein | Splicing Graph | Controls |

⊞ MLH1 + (1842) chr3 (bases=57,587) (exons=19) (residues=757) Homo sapiens mutL homolog 1 (MLH1), transcript variant 1, mRNA. (from RefSeq NM_000249)
⊞ MLH1 + (1843) chr3 (bases=57,493) (exons=18) (residues=516) Homo sapiens mutL homolog 1 (MLH1), transcript variant 6, mRNA. (from RefSeq NM_001258273)
⊟ MLH1 + (1848) chr3 (bases=57,272) (exons=17) (residues=679) Homo sapiens mutL homolog 1 (MLH1), transcript variant 5, mRNA. (from RefSeq NM_001258271)
   chr3:36993573-37050844
   Homo sapiens mutL homolog 1 (MLH1), transcript variant 5, mRNA. (from RefSeq NM_001258271)
   MLH1
   (57,272 bases)
   name = uc062ibj.1
   chrom = chr3
   strand = +
   txStart = 36993572
   txEnd = 37050844
   cdsStart = 36993572
   cdsEnd = 37050653
   exonCount = 17
   spDisplayID = H0Y818_HUMAN
   geneSymbol = MLH1
   refseq = NM_001258271
   protAcc = NM_001258271
   pathway = Mismatch repair=hsa03430
   pathwaymapid = hsa03430
   summary = This gene was identified as a locus frequently mutated in hereditary nonpolyposis colon cancer (HNPCC). It is a human homolog of the E.
   coli DNA mismatch repair gene mutL, consistent with the characteristic alterations in microsatellite sequences (RER+phenotype) found in HNPCC.
   Alternative splicing results in multiple transcript variants encoding distinct isoforms. Additional transcript variants have been described, but
   their full-length natures have not been determined.[provided by RefSeq, Nov 2009].
   PDB =
⊞ MLH1 + (1849) chr3 (bases=57,071) (exons=20) (residues=516) Homo sapiens mutL homolog 1 (MLH1), transcript variant 7, mRNA. (from RefSeq NM_001258274)
⊞ MLH1 + (1851) chr3 (bases=57,066) (exons=18) (residues=516) Homo sapiens mutL homolog 1 (MLH1), transcript variant 4, mRNA. (from RefSeq NM_001167619)
⊞ MLH1 + (1852) chr3 (bases=56,909) (exons=19) (residues=516) Homo sapiens mutL homolog 1 (MLH1), transcript variant 3, mRNA. (from RefSeq NM_001167618)
⊞ MLH1 + (1854) chr3 (bases=57,041) (exons=19) (residues=659) Homo sapiens mutL homolog 1 (MLH1), transcript variant 2, mRNA. (from RefSeq NM_001167617)
⊞ MLH3 - (4371) chr14 (bases=37,769) (exons=13) (residues=1454) Homo sapiens mutL homolog 3 (MLH3), transcript variant 1, mRNA. (from RefSeq NM_001040108)
⊞ MLH3 - (4372) chr14 (bases=37,764) (exons=12) (residues=1430) Homo sapiens mutL homolog 3 (MLH3), transcript variant 2, mRNA. (from RefSeq NM_014381)

*Illustration 28: Gene Annotation Results*

*Illustration 29: Zoom Gene View*



The **Zoom Gene View** *provides you with a scalable magnified view of the 2D genomic region in question.*

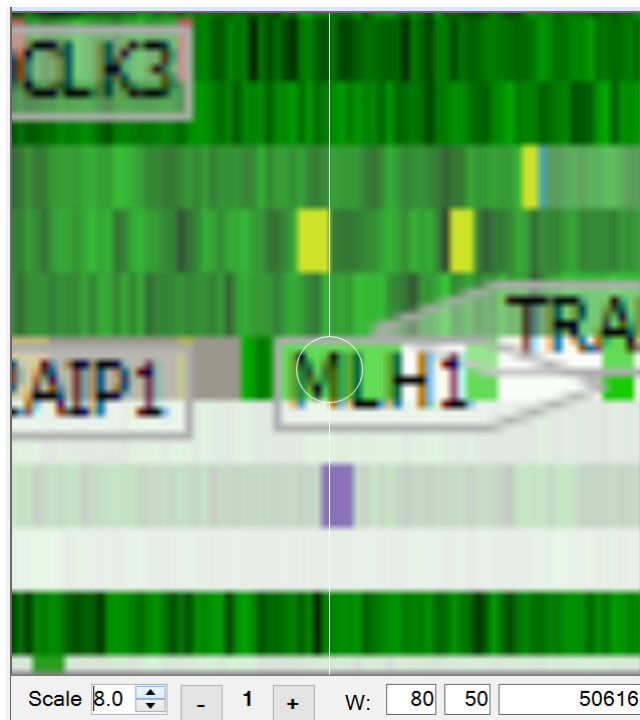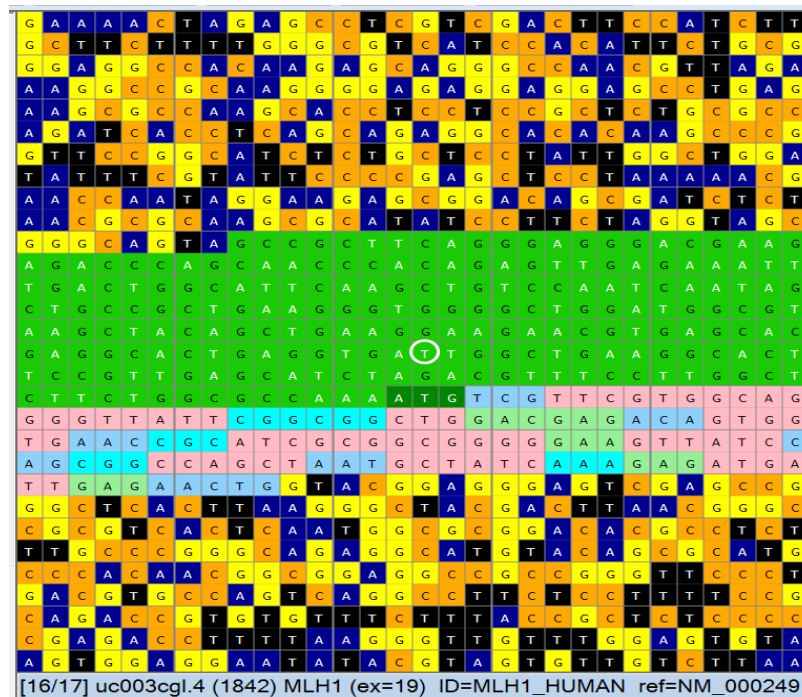When the magnification is increased, it provides you with a more close up view:



*Illustration 30: 8x Magnified Zoom Gene View*

*DNA View* shows the finer detail at sequence level.



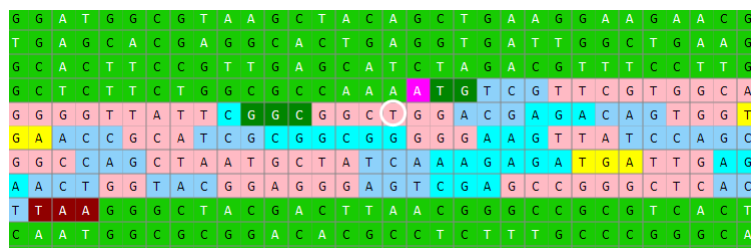*Illustration 31: DNA View showing the **MLH1** gene's 5'UTR*

This is where you can really get acquainted with the primary structure of the gene DNA sequence.

In the DNA view above you can see the 5' Untranslated Region (UTR), the first exon, the start codon (ATG) and part of the protein coding sequence.

Navigation around the genome is achieved via the cursor keys of the keyboard, the mouse via a dragging motion or via touch on a machine with a touch display. All display windows are **navigationally locked**, meaning when you pan or scroll one window, the other zoomed widows moves accordingly. This allows for intuitive operation around the chromosomes of the genome.

## Copying DNA sequences

By moving the mouse in the **DNA View**, the base position at the mouse is displayed in a field on the left. By **double clicking** anywhere in the **DNA View**, a Magenta block is inserted. Moving the mouse will now show you the distance between the mouse location and the magenta base position.



The field displays the **Centre position, the length in bases from the magenta marker, the base offset (r), the x-difference and the y-difference in position. (as well as the chromosome)**

Double clicking again will copy the bases from the magenta base to the current base to the information tab **as well as to the windows clipboard**.

```
chr3:36993548-36993694
```

```
ATGTCGTTCGTGGCAGGGGTTATTCGGCGGCTGGACGAGACAGTGGTGAACCGCATCGCGGCGGGGGAAGTTATCCAGCGGCCAGCTAATGCTATCAAAG
AGATGATTGAGAACTGGTACGGAGGGAGTCGAGCCGGGCTCACTTAA
```

## *Sequencing Gel View*

And, if the sequence is the right length, a "simulated gel electrophoresis" will be shown in the "**Sequencing Gel**" tab:



The Gel View is really **just because is looks cool and reminds us of how far genome sequencing has come**. (It tries to fit in the number of selected bases based on the gel density and length, but sometimes it does not find the proper parameters and does not display it).

# Using the Visual Genome Browser views and keyboard keys

The software has been written with the purpose of making navigating genomes in 2D as easy as possible. All views are navigationally locked to each other allowing a user to move the other views by dragging any one of the Main Genome View, DNA View or Zoomed Gene views.

The views can be moved by **the mouse**, **the keyboard arrow keys** or **by touch** on a computer with a touch display such as the Microsoft Surface Pro (on which it was developed).

Many of the options are also available by keyboard shortcuts:

## *Keyboard and mouse shortcuts*

Mouse wheel

- Move in Zoom DNA View

Mouse left button

- Drag and move views

Mouse centre button

- Jump between Main Genome View and DNA View to correct genome position

Mouse right button

- Gene context menu (If no gene shown, menu for current genome position)

Ctrl

- Press and hold to temporarily enable tooltips

- When moving the mouse in the Main Genome View, will lock the vertical position and only follow the horizontal position of the mouse in the other windows. This makes it easier to track the mouse from left to right on the same genome sequence.

- Holding down the CTRL key while double clicking in the chromosomes list, will select only the most important chromosomes for the Human genomes

Shift

- Increase arrow key navigation speed by 10 in the Main Genome View

- Holding down the SHIFT key while double clicking in the chromosomes list, will select all of the sequences

Alt

- When pressed will slow down the speed of the mouse cursor in the Main Genome View

F1

- Show gene info screen from where you can create favourites

F2

- Show History

F3

- Show Favourites

F4

- Jump to last history or favourite selected

F5

- Show Protein view for current gene or redraw protein view

F6

- Show context menu for current gene

F7

- Move mouse between Main Genome View and Zoom DNA View at correct position in genome

F8

- Collapse/Expand middle panel

F9

- Collapse/Expand bottom panel

F11

- Toggle Full Screen View

F12

- Toggle between current and last expression cell type and redraw genome

Page Down/Up

- Jump to next/previous gene/filtered position

Arrow keys

- Navigate in Main Genome View or Zoom DNA View (Shift key increases move speed by 10)

+ or -

- When mouse cursor in Main Genome View, Zoom DNA View, Splicing View, Zoomed Genes View
  Cycle through different overlapping genes/transcripts at current position

- When showing All Chromosomes view
  Change expression level threshold redraw expression overlay

< or >

- Change number of bases per line in "Aligned" Zoom DNA View

- When mouse cursor in Main Genome View, Zoom DNA View, Splicing View, Zoomed Genes View, Protein View
  Jump to next/previous exon when appropriate context menu is selected

- If Exon selection not active it will select next/previous ALLELE when showing Variants

- When showing All Chromosomes view
  Change expression level threshold redraw expression overlay

[ OR ]

- When Main Genome View tab selected and mouse cursor in Main Genome View, Zoom DNA View, Zoomed Genes View

- Change bases per line in Zoom DNA View

- When Protein View tab selected and mouse cursor in Protein view
  Change amino acids per line in Protein View

Shift-/

- When Mouse cursor in Main Genome View or Zoom DNA View
  Copy DNA bases from DNA View as text to clipboard

Ctrl-F

- Give focus to Gene/Pattern Search field allowing you to search through gene fields or look for patterns

Ctrl-G

- Give focus to Gene quick lookup combo allowing you to type and look for genes

```
Ctrl-S
```

- Save Favourites

```
Ctrl-A
```

- Select all text in most text boxes

```
Windows Menu key
```

- Show gene context menu for current position of mouse in Main Genome View

## *The buttons on the left of the Sequence tab allows you to collapse/expand the panels to suit your preference*

It is sometimes very difficult to fit everything one wants to display on the same view. This is why the software provides buttons to collapse either the middle or bottom panels, or even go into full-screen mode, in order to better view the genomes.



*Illustration 32: With the TOP PANEL prominent*

*Illustration 33: With the BOTTOM PANEL prominent*

*Illustration 34: Hiding the BOTTOM PANEL*

The Full-screen view allows you to see only the **Main Genome View, Protein View, Splicing View, Sequence Gel, Information, Search results or Pathway tree**. or any of the views obtained from the tabs.



*Illustration 35: FULLSCREEN View*

## CpG Sequences near the start of genes

When you want to look at a specific pattern such as DNA transcription factor binding sites or CpG sites, one can enter a pattern in the **Search field** such as **CG** and it will be highlighted in the **DNA View**. Notice how many **CpG di-nucleotides** can be found near the 5'UTR of the previous gene.



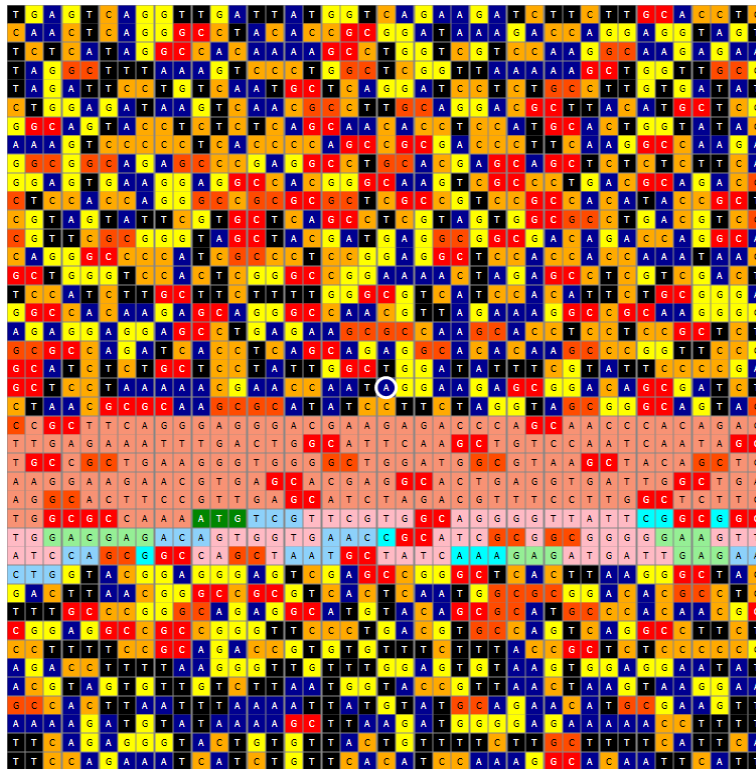The small **C** button can be used to toggle the pattern to its reverse complement.



*Illustration 36: CHR3:36992431-36994111 near gene MLH1 at width 41*

Genes which are constitutively switched on usually have a region of high CpG concentration near the transcription start site, to which transcriptional regulator proteins bind which either prevent RNA Polymerase II from binding or, which recruits epigenetic factors responsible for epigenetically switching off genes by modifying methylation tags on DNA or histones. Genes may also require DNA binding factors and activators to recruit the transcription machinery to the start of transcription.

## Promoter and terminator sequences

The software is capable of matching specific promoter and terminator sequence of genes in order to investigate transcriptional regulatory elements. Take for example the gene IL12A at chr3:159995404-159996019.



Promoter sequences are highlighted in Magenta while terminator sequences are highlighted in dark red.
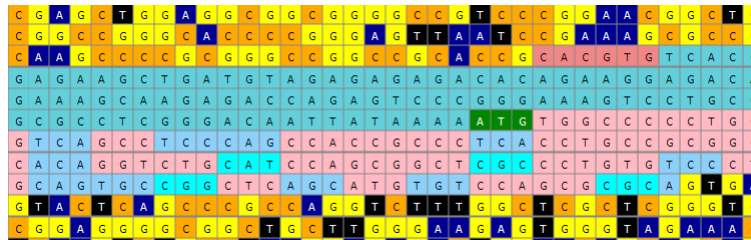


*Illustration 37: Gene IL12A with CACGTG promoter sequence shown in pink. (chr3:159995404-159996019)*
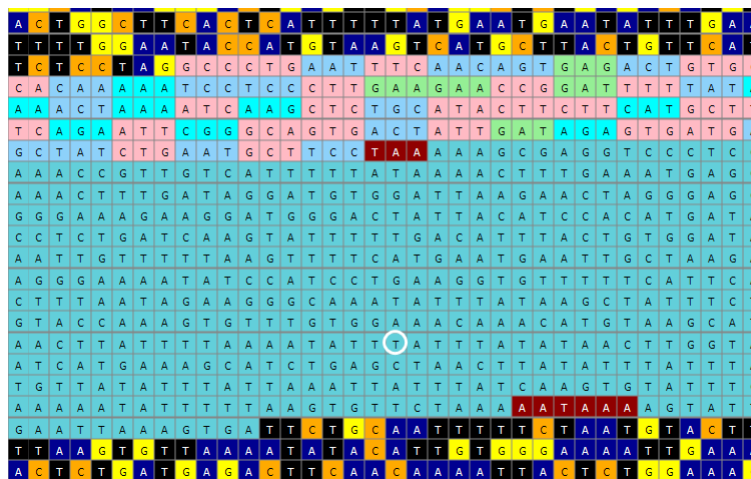


*Illustration 38: The terminator sequence AATAAA at the end of the 3' UTR of gene IL12A*

## *Chromosomes View*

When one select the **Chromosomes View**, a birds-eye view is provided of exactly where on all the chromosomes the search results for the search term: "**dna mismatch repair**" can be found. The currently displayed selection is additionally highlighted with a moving rectangle display as shown in the picture below.
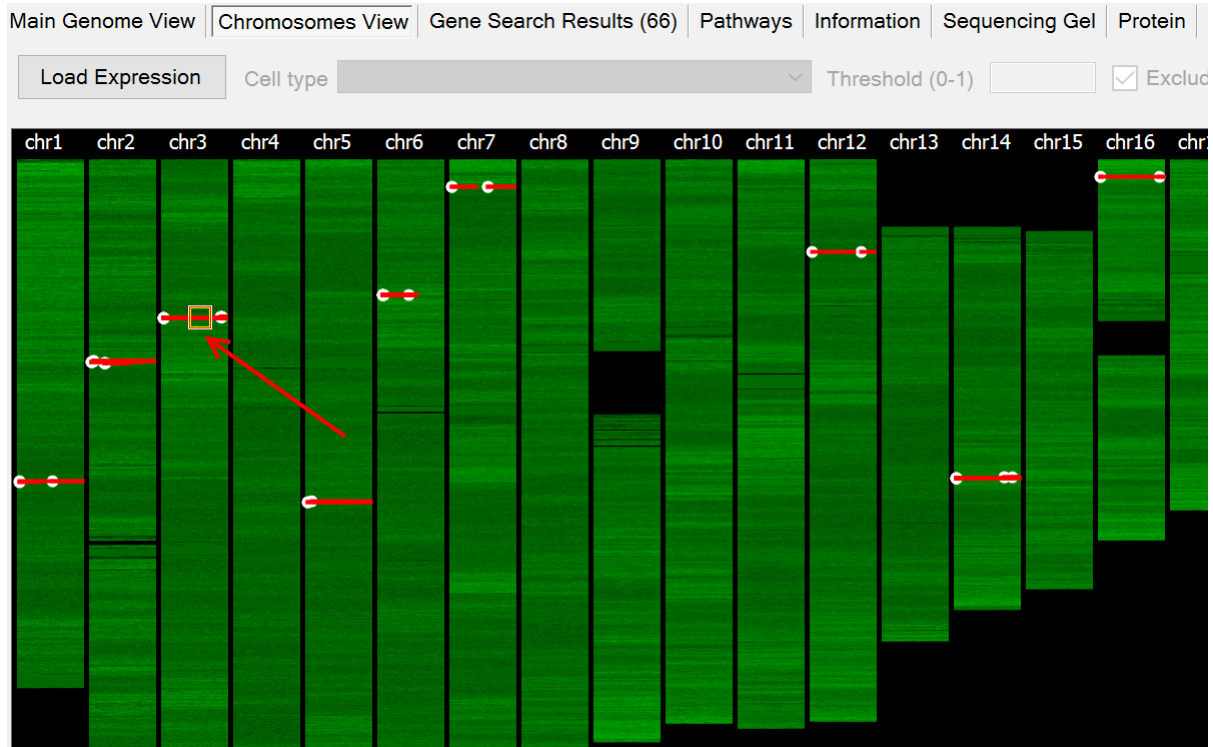


*Illustration 39: Chromosomes View (providing a true birds-eye view)*

Again, hovering over the **Chromosomes View** above will display the corresponding region in the **DNA View** and *clicking on the view*, will navigate down to the specific region and swap back to the **Main Genome View**.

This means one is able to drill down from the **"Big Picture"** (**Chromosomes View**) – to –
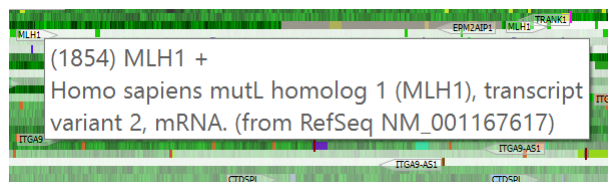
the **Main Genome View** – to –>

the **DNA View** – to –>

the **Molecule View** (where a comic depiction of the molecules in question is shown)

When the mouse pointer is moved across the **Main Genome View**, the chromosomal position is shown (which can be copied and pasted in the online UCSC genome browser)



One can also enter coordinates and then click on "**Go position**" to jump to a specific sequence position. An info tooltip is provided as you hover across the genes in the **Main Genome View**.

Normally, the **Chromosomes View** would be used to provide you with a quick way to see where the genes associated with a specific pathway is located across all of the chromosomes, or to display gene search results obtained from the UCSC genome browser.  As an example, lets do a search for all the tRNA's found in the human genome.  There are more than 500 tRNA's encoded in the human genome, corresponding to 48 anti-codons.  By copying the search results from the UCSC browser and clicking the **CLIP** button, the clipboard contents is parsed in order to plot all the tRNAs at the appropriate positions on the HG19 (in this case) version of the human genome sequence.

## UCSC Genes

TRNA_Pseudo (uc021ofs.1) at chr1:7990339-7990408 - transfer RNA pseudogene (anticodon ???)
TRNA_Asn (uc021ogp.1) at chr1:16847080-16847153 - transfer RNA Asn (anticodon GTT)
TRNA_Asn (uc021ogq.1) at chr1:16858893-16858966 - transfer RNA Asn (anticodon GTT)
TRNA_Gly (uc021ogs.1) at chr1:16872434-16872504 - transfer RNA Gly (anticodon CCC)
TRNA_Pseudo (uc021ogt.1) at chr1:16874160-16874232 - transfer RNA pseudogene (anticodon CAC)
TRNA_Gly (uc021ogw.1) at chr1:17004766-17004836 - transfer RNA Gly (anticodon CCC)
TRNA_Val (uc021ogx.1) at chr1:17006501-17006573 - transfer RNA Val (anticodon CAC)
TRNA_Pseudo (uc021ogz.1) at chr1:17052061-17052133 - transfer RNA pseudogene (anticodon CAC)
TRNA_Gly (uc021oha.1) at chr1:17053780-17053850 - transfer RNA Gly (anticodon CCC)
TRNA_Pseudo (uc021ohb.1) at chr1:17180900-17180971 - transfer RNA pseudogene (anticodon CTG)
TRNA_Pseudo (uc021ohd.1) at chr1:17186693-17186765 - transfer RNA pseudogene (anticodon CAC)
TRNA_Gly (uc021ohe.1) at chr1:17188416-17188486 - transfer RNA Gly (anticodon CCC)
TRNA_Asn (uc021ohg.1) at chr1:17201958-17202031 - transfer RNA Asn (anticodon GTT)
TRNA_Asn (uc021ohh.1) at chr1:17216172-17216245 - transfer RNA Asn (anticodon GTT)
TRNA_Pseudo (uc021olx.1) at chr1:39970195-39970267 - transfer RNA pseudogene (anticodon CTT)
TRNA_Lys (uc021onv.1) at chr1:55423542-55423614 - transfer RNA Lys (anticodon CTT)
TRNA_Cys (uc021opz.1) at chr1:93981834-93981906 - transfer RNA Cys (anticodon GCA)
TRNA_Arg (uc021oqb.1) at chr1:94313129-94313213 - transfer RNA Arg (anticodon TCT)
TRNA_Pseudo (uc021oqx.1) at chr1:108496275-108496345 - transfer RNA pseudogene (anticodon ???)
TRNA_Asn (uc021otf.1) at chr1:143690028-143690101 - transfer RNA Asn (anticodon GTT)
TRNA_Asn (uc021otg.1) at chr1:143879832-143879905 - transfer RNA Asn (anticodon GTT)
TRNA_Asn (uc021oty.1) at chr1:144301611-144301684 - transfer RNA Asn (anticodon GTT)
TRNA_Asn (uc021otz.1) at chr1:144308614-144308687 - transfer RNA Asn (anticodon GTT)
TRNA_Asn (uc021ouc.1) at chr1:144481840-144481913 - transfer RNA Asn (anticodon GTT)
TRNA_Asn (uc021oud.1) at chr1:144488843-144488916 - transfer RNA Asn (anticodon GTT)
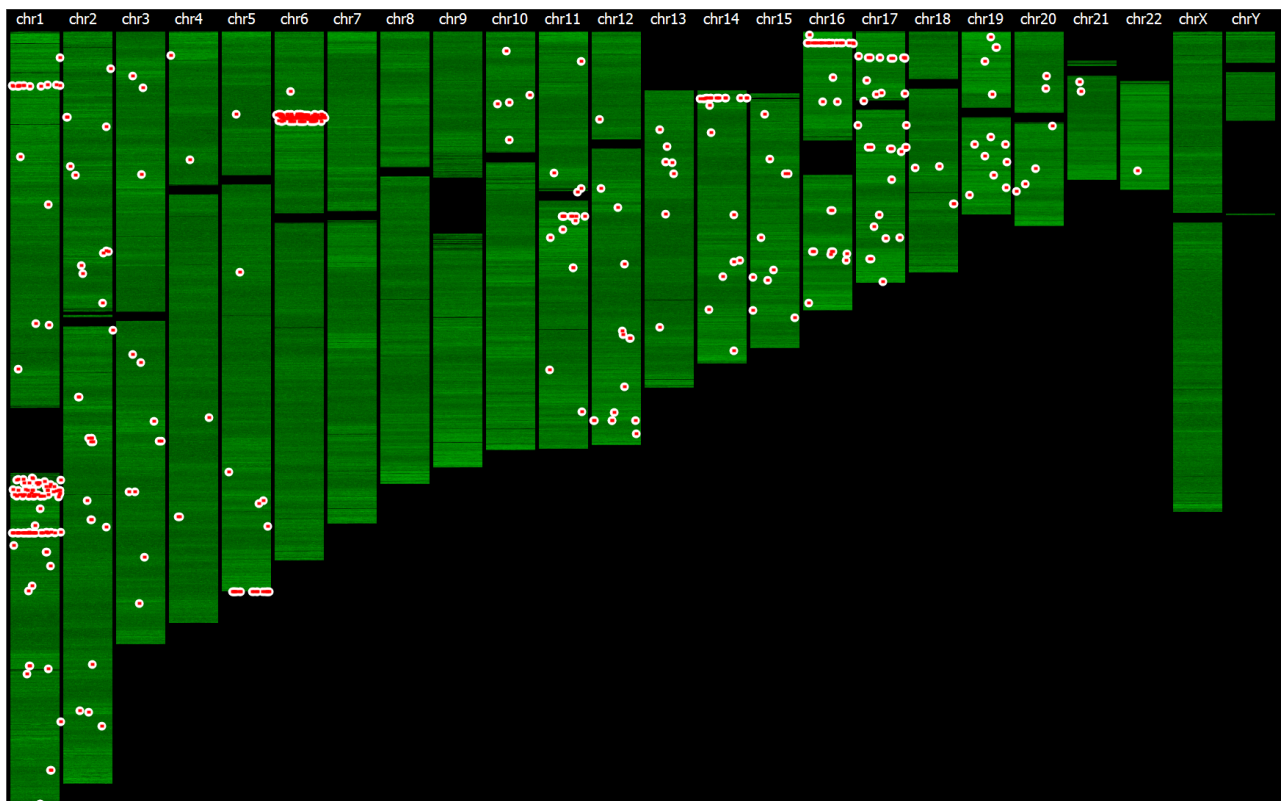TRNA_Gln (uc021oug.1) at chr1:144839436-144839507 - transfer RNA Gln (anticodon CTG)
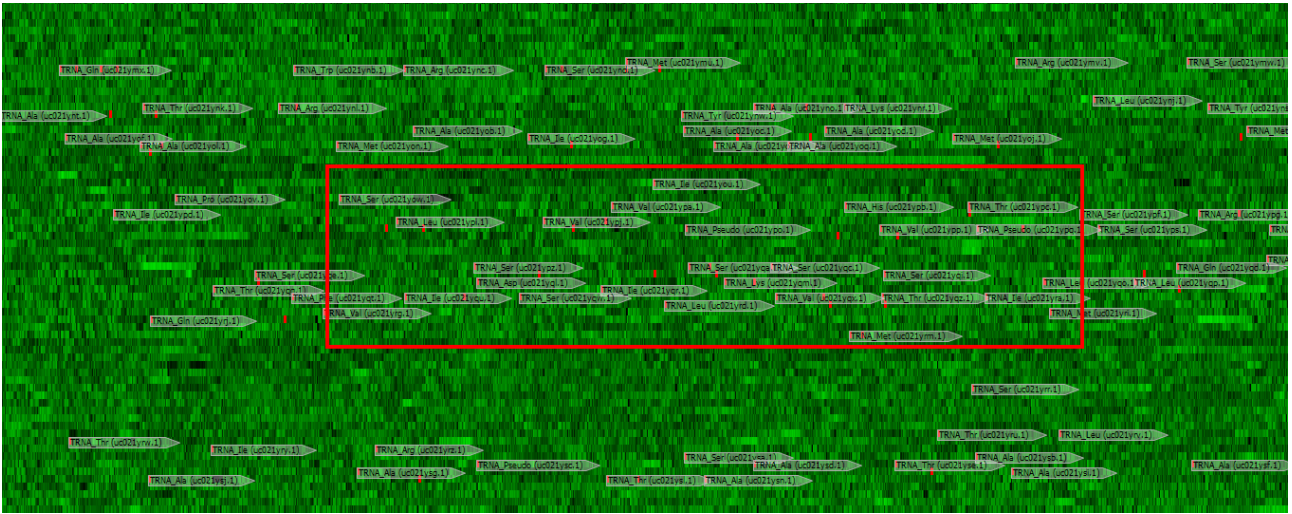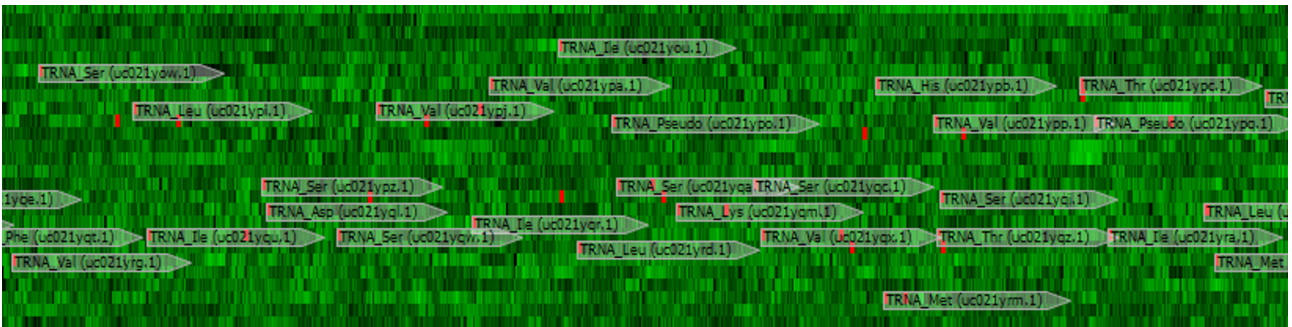


*Illustration 40: Chromosomes view showing all the positions of tRNAs in the genome*
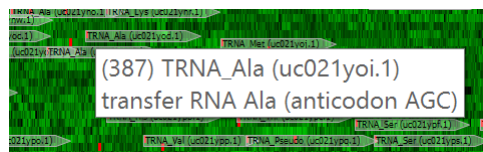
## *The Main Genome View*

When you now click on the band of red dots on chromosome 6, you can see their distribution clearer.
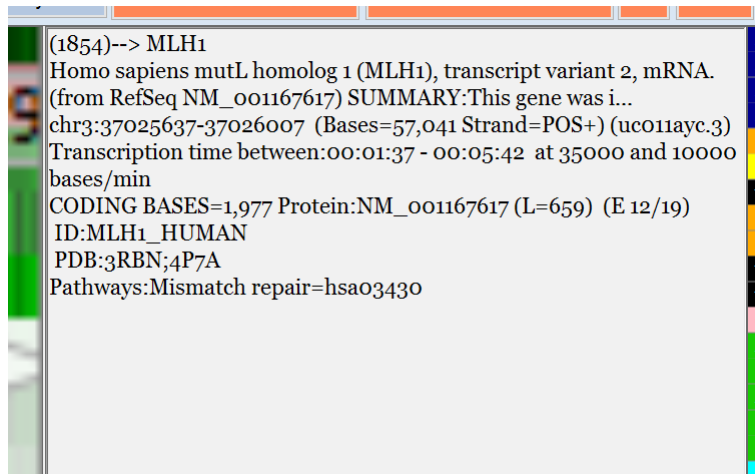


With the highlighted area magnified...



Hovering the mouse over any of these genes will bring up a tool tip:



(387) TRNA_Ala (uc021yoi.1)
transfer RNA Ala (anticodon AGC)

## *The Information Display window*

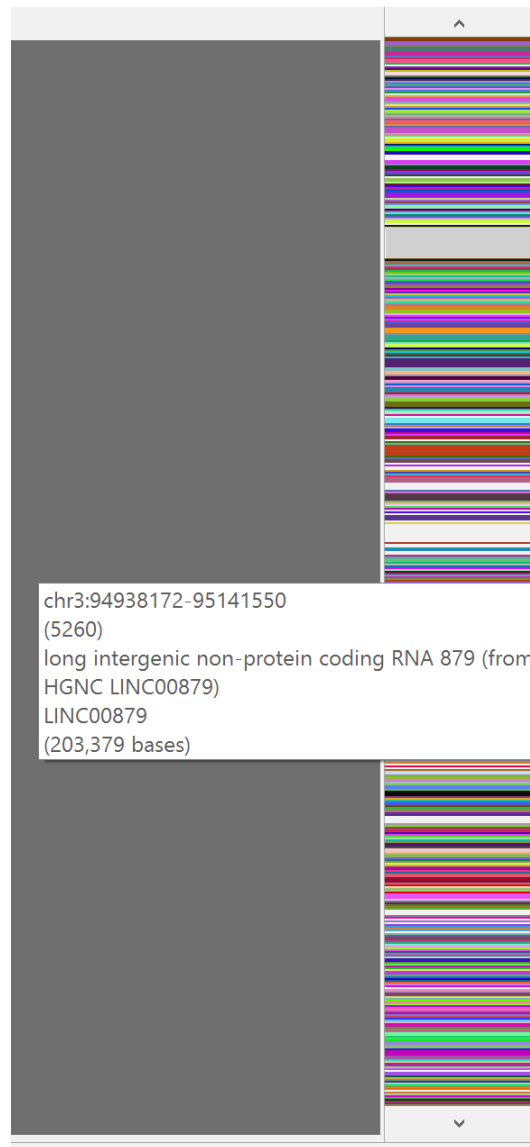The **Information Display** in the top centre displays relevant information for the specific gene.



*Illustration 41: Information Display*

## *The enhanced Chromosome scrollbar*

The chromosome scrollbar on the right provides another quick way of finding gene annotations.



*Illustration 42: Chromosome scrollbar.*
*Hovering provides an info tooltip*

## Filtering display results

Probably one of the **most useful** features of the Visual Genome Browser is the ability to filter the amount of annotations displayed. This can be done via the **Filter Input** field at the top right.
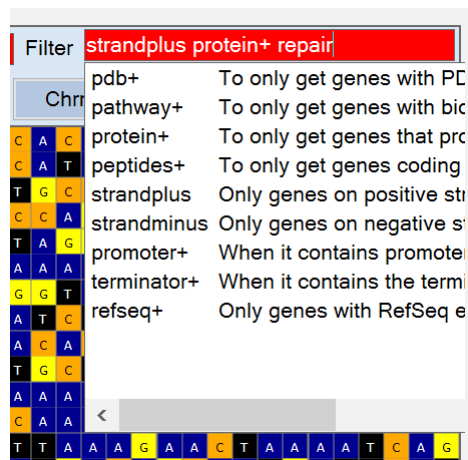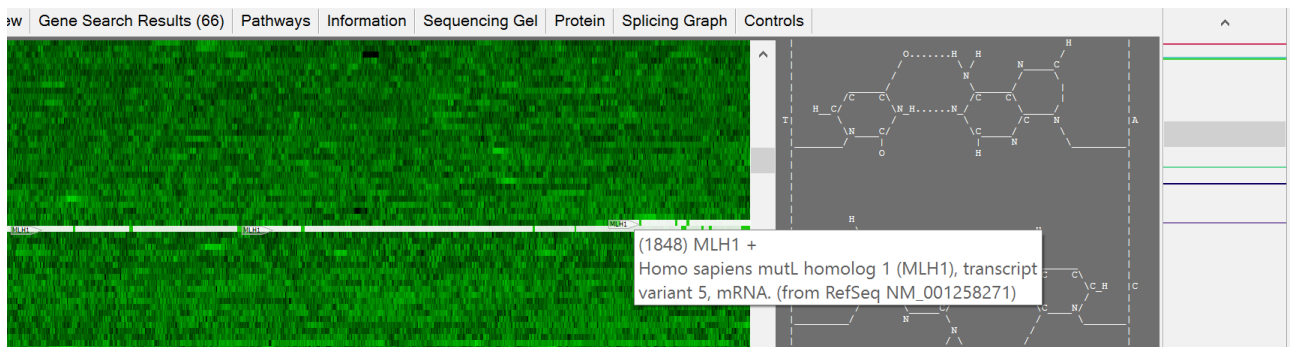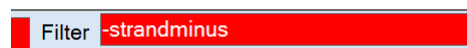

*Illustration 43: Filter Input field*



When you click "DRAW/REDRAW", **all** browser views are updated to reflect this new filter criteria, including the scrollbar. This makes it much easier to see the **forest for the trees (genes)**.

One can also **filter out** certain genes, by adding a minus in front of the terms in the filter:



Entering **–strandminus** will filter out all genes on the negative strand of the chromosome.

Entering something like **–protein+** (will filter out any genes that are coding for proteins and only display non-coding genes), because **protein+** is the search term for genes which has an ORF or coding sequence.

## *Examining genes using the different views*

Because there are multiple overlapping transcripts at most positions, it is possible to use the **+/-** buttons to filter the **DNA View** in order to only see a specific splice variant.
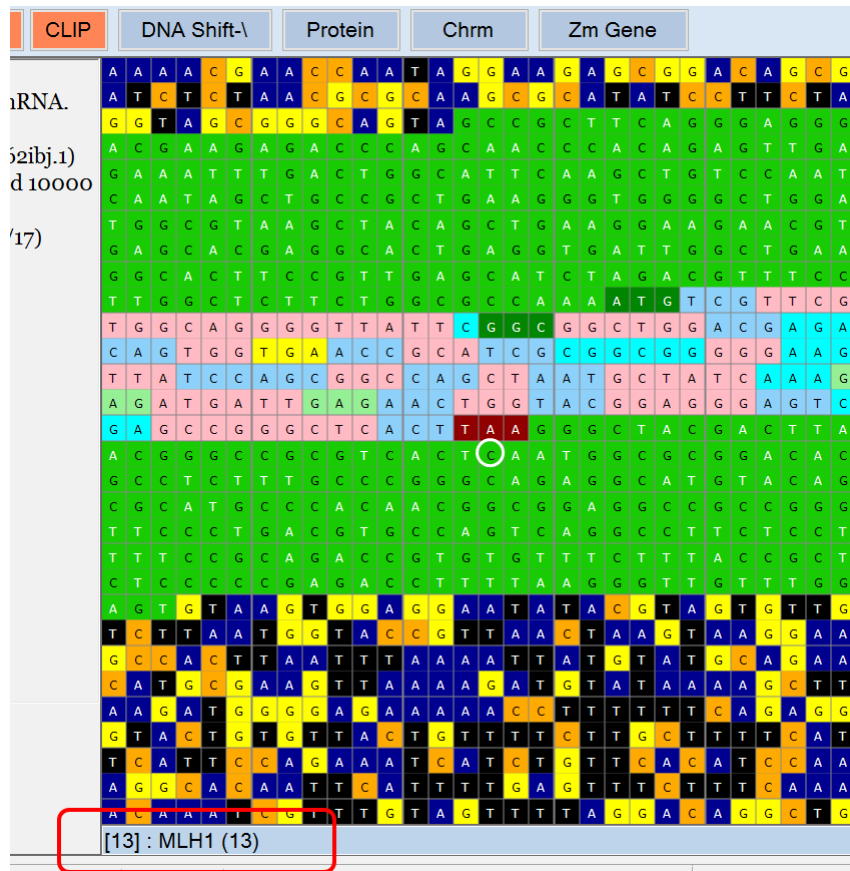


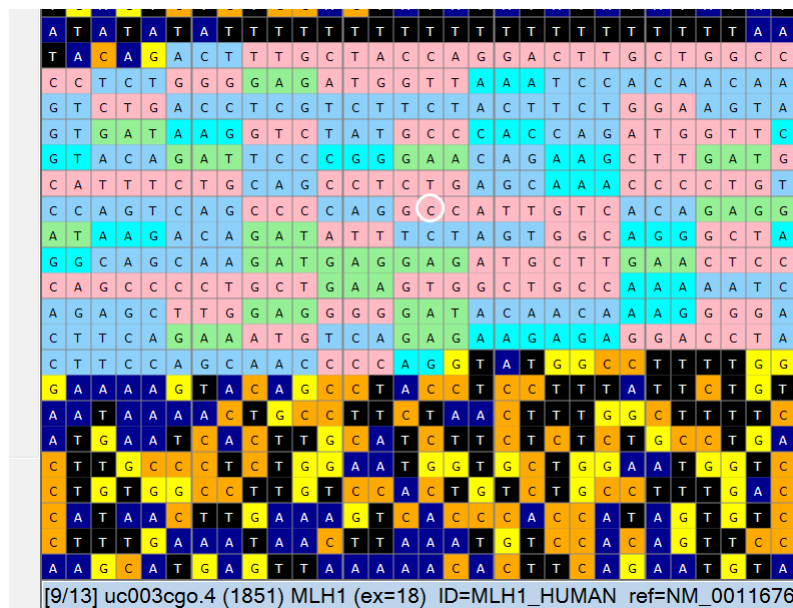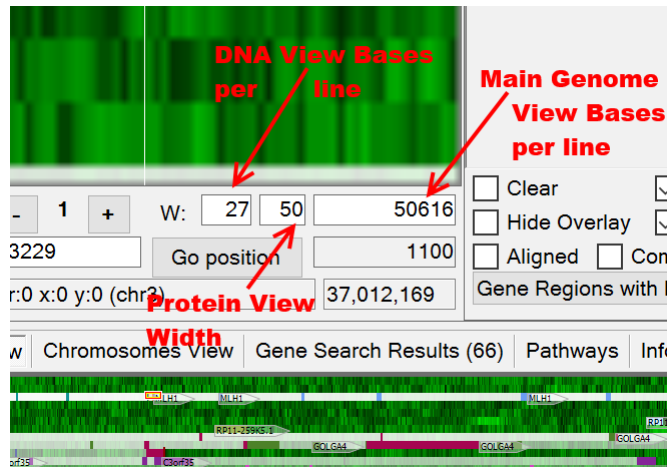*Illustration 44: Use plus/minus to select specific overlapping gene at same position.*



*Illustration 45: Different transcript selected*

You can change the number of bases per line in the display by pressing **[ or ]** or by changing the width in the text field:



The **DNA View**'s number of bases per line is correspondingly altered. When the **Codons** checkbox is selected, the coding sequences are displayed using colours which correspond to the hydrophobic/polar character of the amino acids coded by these codons. In the picture below one can also observe the bases responsible for demarcating the start and ends of introns which are removed by the Spliceosome when it produces the mature RNA transcript.



*Illustration 46: DNA View with Splice signal bases highlighted*

Codons which are **pink** code for amino acids which are hydrophobic and tend to be on the interior of the resultant folded protein or embedded in lipid membranes, whilst amino acids which are polar

(from the **blue and green** codons) tend to be on the exterior of the folded protein. The colour of the codons reveals much about the properties of a protein. We would therefore expect that a protein like **Aquaporin** (which is a transmembrane protein used to transport water molecules across the lipid membranes), would consist of a lot of hydrophobic amino acids on its outside, with polar amino acids on its interior to interact with the polar water molecules.

This can be seen in the following **Protein Sequence View**, which is obtained by **double-clicking** on a gene in the **Main Genome View** or by selecting "**Show Protein Sequence**" in the context menu.



By changing the number of amino acids in each row, it is possible to line up the amino acids at the start of each exon, indicated by red blocks in the protein view. It is interesting how similar amino acids line up vertically (almost like a crossword puzzle). Also notice how many hydrophobic amino acids can be found in the protein.

Again, notice how the displays are linked. When you hover the mouse on the **Protein View**, the **DNA View** will navigate to the corresponding codon.

## Protein Statistics

You can obtain the gene context menu by right clicking on the protein view. If you select "**Show Protein Statistics**" you will see why this membrane protein is considered a mostly hydrophobic protein.
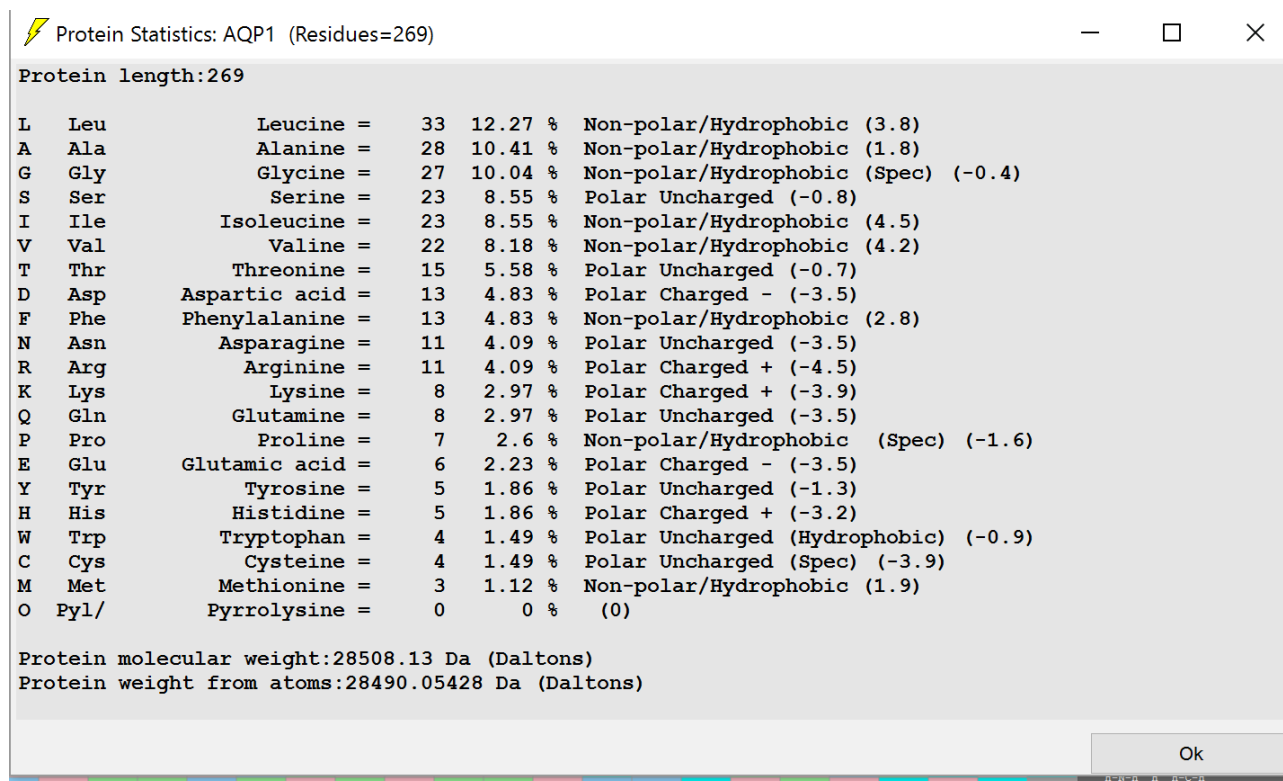
```
Protein Statistics: AQP1  (Residues=269)                          —    □    ✕

Protein length:269

L    Leu              Leucine =    33  12.27 %   Non-polar/Hydrophobic (3.8)
A    Ala              Alanine =    28  10.41 %   Non-polar/Hydrophobic (1.8)
G    Gly              Glycine =    27  10.04 %   Non-polar/Hydrophobic (Spec) (-0.4)
S    Ser               Serine =    23   8.55 %   Polar Uncharged (-0.8)
I    Ile            Isoleucine =    23   8.55 %   Non-polar/Hydrophobic (4.5)
V    Val               Valine =    22   8.18 %   Non-polar/Hydrophobic (4.2)
T    Thr             Threonine =    15   5.58 %   Polar Uncharged (-0.7)
D    Asp         Aspartic acid =    13   4.83 %   Polar Charged - (-3.5)
F    Phe         Phenylalanine =    13   4.83 %   Non-polar/Hydrophobic (2.8)
N    Asn            Asparagine =    11   4.09 %   Polar Uncharged (-3.5)
R    Arg              Arginine =    11   4.09 %   Polar Charged + (-4.5)
K    Lys                Lysine =     8   2.97 %   Polar Charged + (-3.9)
Q    Gln             Glutamine =     8   2.97 %   Polar Uncharged (-3.5)
P    Pro               Proline =     7    2.6 %   Non-polar/Hydrophobic  (Spec) (-1.6)
E    Glu         Glutamic acid =     6   2.23 %   Polar Charged - (-3.5)
Y    Tyr              Tyrosine =     5   1.86 %   Polar Uncharged (-1.3)
H    His             Histidine =     5   1.86 %   Polar Charged + (-3.2)
W    Trp            Tryptophan =     4   1.49 %   Polar Uncharged (Hydrophobic) (-0.9)
C    Cys              Cysteine =     4   1.49 %   Polar Uncharged (Spec) (-3.9)
M    Met            Methionine =     3   1.12 %   Non-polar/Hydrophobic (1.9)
O    Pyl/           Pyrrolysine =     0     0 %     (0)


Protein molecular weight:28508.13 Da (Daltons)
Protein weight from atoms:28490.05428 Da (Daltons)


                                                                      Ok
```
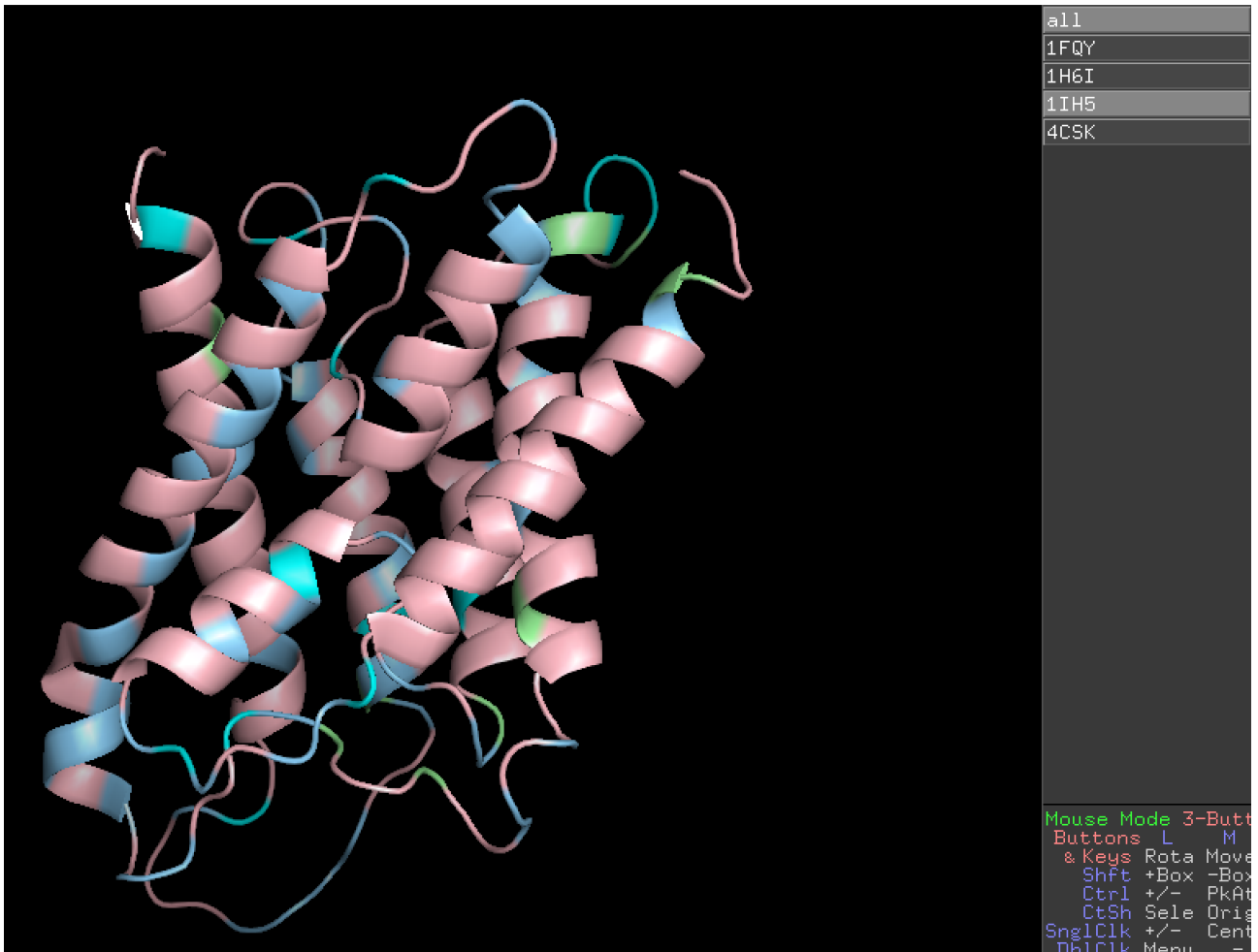
*Illustration 47: The protein statistics shows important statistical values for each amino acid in the protein sequence*
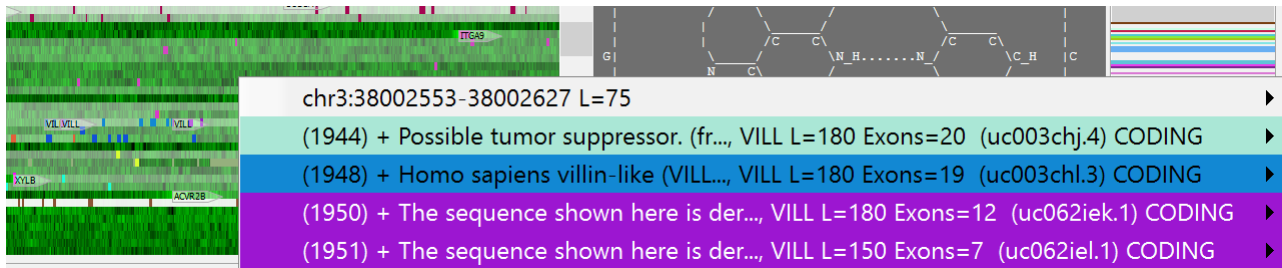
## Launching PDB Crystal structures in Pymol

Selecting "Show PDB Structure in Pymol" will download a few PDB files from the Protein Database and colour code the amino acids according to the same colour legend. It will keep these PDB files in a disk cache in order that it does not need to be downloaded next time.

The 3D visualization software "PyMol" will then be launched in order to display the crystal structures determined for these proteins.  Other visualization software will later also be supported.

# The Gene Context Menu

When you want to get more information for the overlapping genes, you can right-click on any of the windows or trees in order to get a context menu which is coloured according to the genes below the mouse cursor, allowing you to easily select the most appropriate gene. A "software generated" **unique** number is always shown in brackets. **This will help you to identify a unique gene annotation, no matter in which window it is selected from**.



The context menu provides a list of context specific options which can be performed on each gene. The context menu is available from the following windows:

- **Main Genome View** (when right clicking at any location in the genome or on gene regions)
- **DNA View**
- **Protein Sequence View**
- **Gene Search Results tree**
- Biochemical **Pathways tree**

*The following options are available in the gene context menu:*

Move mouse to other window (F7)
Copy position to clipboard
Show info in window (F1)
Locate chromosome/sequence
Open in UCSC Browser
Retrieve UCSC Gene Information
Retrieve Encode Dnase Sites
Jump to first exon (use < > to move between exons)
Copy transcribed DNA to clipboard + Comp
Copy spliced transcript to clipboard + Comp
Copy Splice Info to clipboard
Copy only coding sequence to clipboard + Comp
Copy Protein to clipboard + Comp
Show gene homology in browser
Show biochemical pathway in browser
Show PDB Structure in Pymol
Show Protein Sequence
Show Protein Statistics
Show Protein Sequence with PDB Secondary Structure
Show Refseq Info at UCSC
Show in Human Protein Atlas
Show Splicing Graph
Info

*Illustration 48: Context menu options*

## Move mouse to other window (F7)

When you want to quickly navigate between the **Main Genome View** and the **DNA View**, this option allows you to quickly jump between the two windows. It can also be achieved by pressing F7 or using the *middle button of the mouse*.
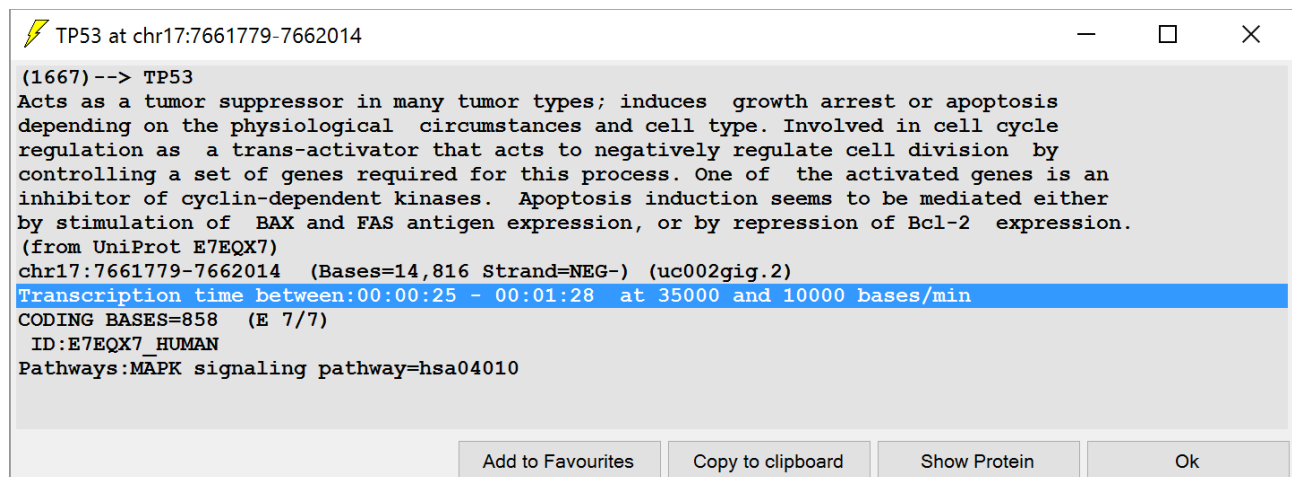
## Copy position to clipboard

The Visual Genome Browser uses the 1-based position format also used in the online UCSC browser. This option will simply copy the genome position below the cursor to the clipboard eg. **chr3:159995734-159996019**

## Show Info in window (F1)

The **Information Display** in the top centre of the main screen displays information of the gene or annotation currently below the mouse cursor or which is selected in the centre of the cross hairs of the **DNA View** or the **Zoom Gene View.**

This option will open up a separate window showing this information, and also gives you the option of adding the selected gene to the **Favourites**.

⚡ TP53 at chr17:7661779-7662014       — ☐ ✕

```
(1667)--> TP53
Acts as a tumor suppressor in many tumor types; induces  growth arrest or apoptosis
depending on the physiological  circumstances and cell type. Involved in cell cycle
regulation as  a trans-activator that acts to negatively regulate cell division  by
controlling a set of genes required for this process. One of  the activated genes is an
inhibitor of cyclin-dependent kinases.  Apoptosis induction seems to be mediated either
by stimulation of  BAX and FAS antigen expression, or by repression of Bcl-2  expression.
(from UniProt E7EQX7)
chr17:7661779-7662014   (Bases=14,816 Strand=NEG-) (uc002gig.2)
Transcription time between:00:00:25 - 00:01:28  at 35000 and 10000 bases/min
CODING BASES=858   (E 7/7)
 ID:E7EQX7_HUMAN
Pathways:MAPK signaling pathway=hsa04010
```

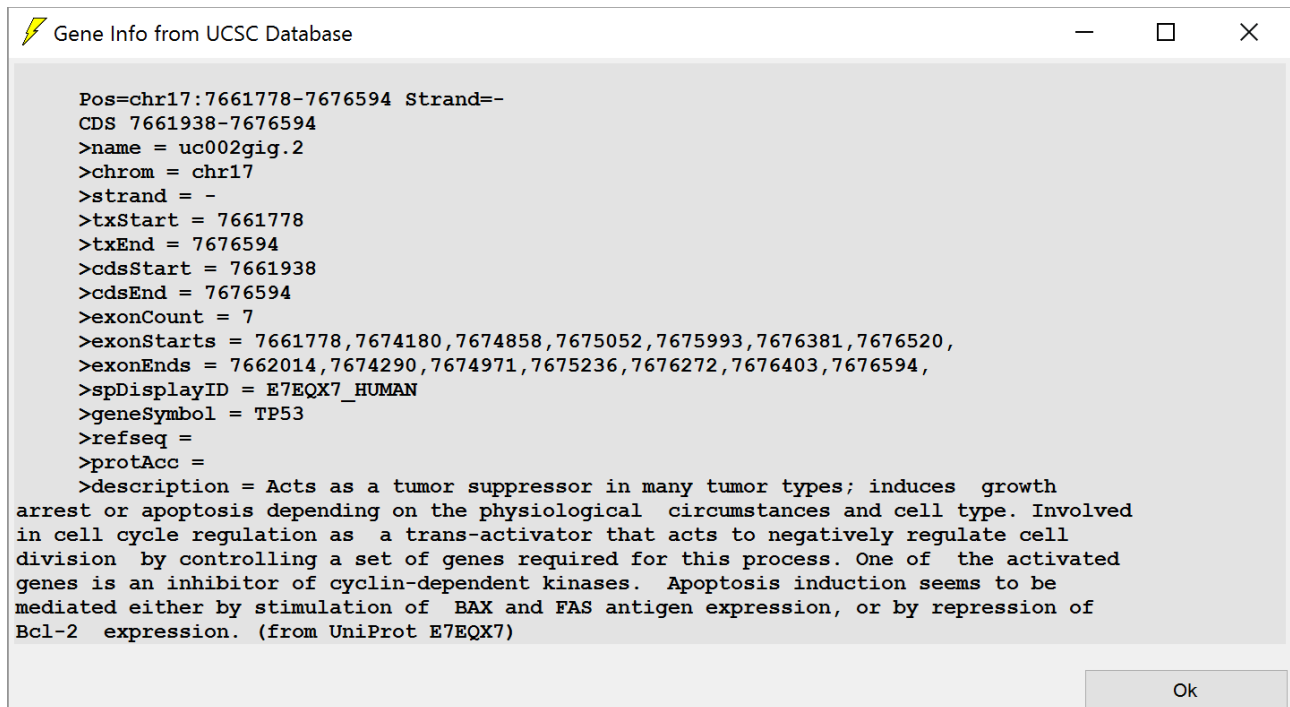| | Add to Favourites | Copy to clipboard | Show Protein | Ok |

## Locate chromosome/sequence

When one has searched for genes matching a specific search criteria and you are looking at in a tree containing genes from many chromosomes, this option allows you to quickly locate the gene's current chromosome in the **Chromosomes/Sequences list**.

## Open in UCSC Browser

This option will open the selected gene position in the UCSC Genome Browser

## *Retrieve UCSC Gene Information*

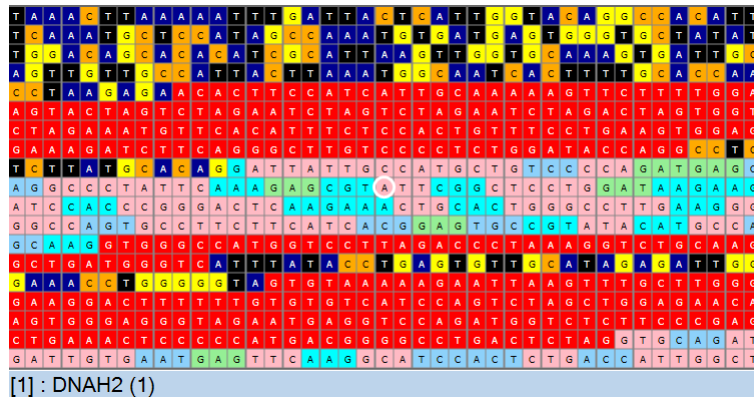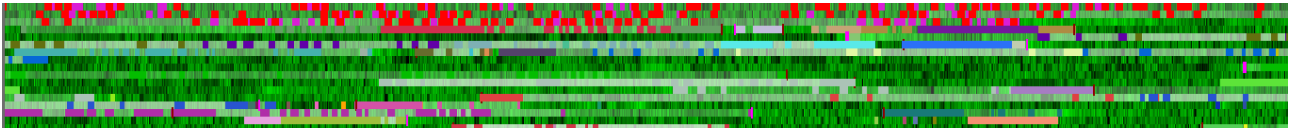This option will retrieve the MySQL query fields from the UCSC Tables and display it in a window



It returns all of the same information which were downloaded for the indexed gene annotations.

Gene annotations where downloaded from the UCSC database using the following SQL query:

SELECT HG38.knownGene.name, HG38.knownGene.chrom, HG38.knownGene.strand, HG38.knownGene.txStart,

HG38.knownGene.txEnd, HG38.knownGene.cdsStart, HG38.knownGene.cdsEnd,

HG38.knownGene.exonCount, HG38.knownGene.exonStarts, HG38.knownGene.exonEnds,

HG38.kgXref.spDisplayID, HG38.kgXref.geneSymbol, HG38.kgXref.refseq,

HG38.kgXref.protAcc, HG38.kgXref.description,

HG38.keggMapDesc.description as pathway, HG38.keggPathway.mapID as pathwaymapid,

hgFixed.refSeqSummary.summary as summary

FROM HG38.knownGene

inner join HG38.kgXref on HG38.knownGene.name = HG38.kgXref.kgID

left outer join HG38.keggPathway on HG38.knownGene.name = HG38.keggPathway.kgID

left outer join HG38.keggMapDesc on HG38.keggPathway.mapID = HG38.keggMapDesc.mapID

left outer join hgFixed.refSeqSummary on HG38.kgXref.refseq = hgFixed.refSeqSummary.mrnaAcc

where HG38.knownGene.chrom like 'chrX'

ORDER BY HG38.knownGene.txStart;
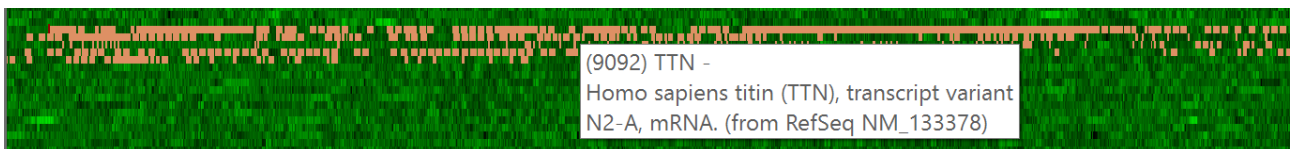
## Retrieve Encode Dnase Hypersensitivity Sites

This is currently only available for this software when using **HG19**. It downloads the DNAse 1 Hypersensitivity sites for a region around the selected position and it is highlighted in red.





[1] : DNAH2 (1)

## Jump to first exon (use < > to move between exons)

It is sometimes helpful to be able to navigate quickly from exon to exon in the **DNA View**. This is where this option comes in handy. When selected, it immediately jumps to the first exon of the gene. After this is selected you can use the **< and > keys** in order to quickly jump through the list of exons in the gene.

Take the TTN (Titin gene) for example. Setting the intron transparency to 100% and filtering out all genes not containing the software based unique no (9092), cleans up the display to allow you to only see the single Titin transcript consisting of 312 exons.



(9092) TTN -
Homo sapiens titin (TTN), transcript variant
N2-A, mRNA. (from RefSeq NM_133378)

Pressing > repeatedly navigates forward through the exons up to exon 312 and

pressing < navigates backwards to exon 1.

The TTN gene is transcribed from the negative strand and the RNA Polymerase II enzyme moves backwards through the chromosome to correctly read the bases.

Notice that the splice signals GT....AG are reverse complemented as CT....AC
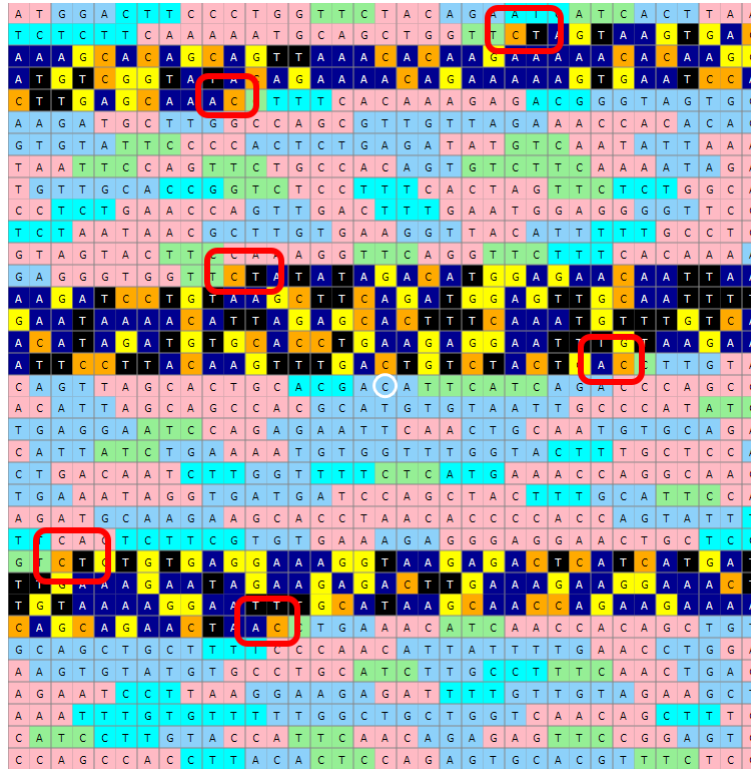


*Illustration 49: chr2:178723267-178724491 within TTN gene at width 35*

The splice signals which are part of the introns are highlighted with red rectangles. The codons have to be read in the reverse direction. It is actually possible to conveniently complement the bases by selecting the checkbox option "**Complement**" when the gene in question is encoded on the reverse strand as is the case here. (Genes on the positive strand will also show correctly)



*Illustration 50: chr2:178723267-178724491 with bases complemented (Splice bases are highlighted)*

The codons can now be read backwards, but are showing the correct complemented base.

It seems like, although the lengths of introns may vary a lot, the lengths of the exons tend to fall in a narrow range of lengths. What is interesting to me is how the codons of a protein coding gene can be scattered out over large distances, but still join up (even in the middle of codons a lot of the
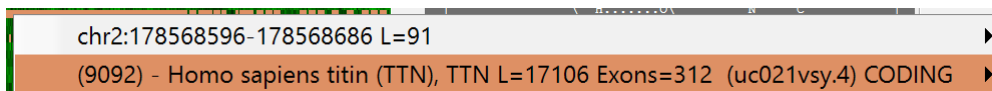


*Illustration 51: chr2:178723267-178724491 with bases complemented (Broken up codons are highlighted)*

time).

In the picture above you can see the same exons (Exons 69-72 of a total of 312) with the codons which span more than one exon. Only after the exons are correctly spliced together are these codons able to make sense in terms of the protein coding region.

It is like a puzzle which only makes sense once you have correctly assembled the individual pieces. I sometimes wonder how these original genes got split up into exons and then got incorporated with the correct splice signals (GT...AG) in at sometimes distant positions among intronic regions. Still, when you assemble the exons, the codons suddenly match up and can code for proteins.

The software provides the following information for this gene, which produces the biggest protein in the human body:

chr2:178568596-178568686 L=91 ▶

(9092) - Homo sapiens titin (TTN), TTN L=17106 Exons=312  (uc021vsy.4) CODING ▶

**Here is an example of the kind of information displayed for a gene:**

(9092)--> TTN

**Homo sapiens titin (TTN)**, transcript variant N2-A, mRNA. (from **RefSeq NM_133378**)

SUMMARY:This gene encod...

**chr2:178530241-178535849  (Bases=281,433 Strand=NEG-)** (uc021vsy.4)

**Transcription time between:00:08:02 - 00:28:08  at 35000 and 10000 bases/min**

**CODING BASES=100,272** Protein:NM_133378 (**L=33424**)  (**E 307/312**)

 ID:TITIN_HUMAN


**PDB**:1BPV;1G1C;1NCT;1NCU;1TIT;1TIU;1TKI;1TNM;1TNN;1WAA;1YA5;2A38;2BK8;2F8V;2ILL;2J8H;2J8O;2NZI;2RQ8;2WP3;2WWK;2WWM;2Y9R;3B43;3KNB;3LCY;3LPW;3PUC;3Q5O;3QP3;4C4K;4JNW;4O00;4QEG;4UOW

**Pathways**:*Hypertrophic cardiomyopathy* (HCM)=hsa05410

This gene encodes a large abundant protein of striated muscle. The product of this gene is divided into two regions, a N-terminal I-band and a C-terminal A-band. The I-band, which is the elastic part of the molecule, contains two regions of tandem immunoglobulin domains on either side of a PEVK region that is rich in proline, glutamate, valine and lysine. The A-band, which is thought to act as a protein-ruler, contains a mixture of immunoglobulin and fibronectin repeats, and possesses kinase activity. An N-terminal Z-disc region and a C-terminal M-line region bind to the Z-line and M-line of the sarcomere, respectively, so that a single titin molecule spans half the length of a sarcomere. Titin also contains binding sites for muscle associated proteins so it serves as an adhesion template for the assembly of contractile machinery in muscle cells. It has also been identified as a structural protein for chromosomes. Alternative splicing of this gene results in multiple transcript variants. Considerable variability exists in the I-band, the M-line and the Z-disc regions of titin. Variability in the I-band region contributes to the differences in elasticity of different titin isoforms and, therefore, to the differences in elasticity of different muscle types. Mutations in this gene are associated with familial hypertrophic cardiomyopathy 9, and autoantibodies to titin are produced in patients with the autoimmune disease scleroderma. [provided by RefSeq, Feb 2012].


## *Copy transcribed DNA to clipboard + Comp*

This option copies the raw DNA bases to the clipboard (and the Comparison list on the Main tab) Bases are copied as FASTA output and in the correct direction (in other words: negative strand transcription will appear in the forward direction as it is transcribed).  Intronic bases is shown in lower case and exons in upper case.

```
>TTN ; Homo sapiens titin (TTN), transcript variant N2-A, mRNA. (from RefSeq NM_133378) ; Transcript length = 306160 (uc021vsy.4 Full
gene position = chr2:178525989-178807421)

GCAGTCGTGCATTCCCAGCCTCGCCTCGGGTGTAGGGATTGCATAGAAAAGCAAAACTACACAGTCTTGACTGTGTAGTTTTGTTTTTAGGATTAGAGGC

TCACCGATTCATGTCGGAGATGGTCAGAAAAACCAACTCTCCATAGGACGTCGTTTCAGAAGCAACCTTGGGCTTAGTCCCACCCTTTTTAGGCACTCTT

GAGAAATCAGGTGAGCAGCAACTTTTCTTATTTTAATAGTAGAGTCACTTTCTTATTGATCTCTGGCTCTGGATTATTTGTGTGTGCTAAATGCCTATGT

ACAAGTGTTTGTATGACTATTTCTTGCCAGTGCTGACCCATTATTTAAATCTGCTTGTTAAATTCTATTTTAAAATGTCTGAACATTTTTCCTTAAAAAG

TTGAAAGCTTTATTATATCAGCTAAATGATTCCATGATTCAAATTGTTTAAGTTCACTTAATATTCTAAAAGTTGAGGTCTGAAATGTTGAGAGAAATGT

TTACTTCTTTTCTAGACTCTCTTTGCCTGTTGAGAATGTTCAAGGTTATATTACTCAATGATAAAGTTAAGCACACCAAGAATAGCTCTGATAAAATGCA

TGTGGTATCAGCTACACCCTGATGACTTTTTAAGGGAATCAACTGCTGAGGTATCAGTTGCAGGACAGTGACAGAACACTTGCCTGCTGTCAGTCAACCT

TGGAGAGTTTAGGAGAGTTTATTTTAAGAGCATGACATTTTAGACGCATACCTTCAGCCCCAGTTAGACAAAGGCCTGTGAAGAGCAAGTGCTAAGGCAA

TAAAATTTACCTAGACAATAGGAATGTAAGATCAATTTAAGCTGGTTAATTGCAACATAGTTGGAGATTTGGAATAAAATGATATAATTGTTTTCTAAAA

TATTCAGAATCTCTTCATGCATTTTAACTCAGCACACAGTCCAATTATATATTTTGAGATAAATGTTGACACTAAGAAGGGTAAGAAATAATGTTTAAAG
```

GAAAGGTTTCTTACTTCTAAATACACTTTGATTTCTCATCTTCCTAGAATTTGATTACTAATTGTCAATGCCTGTAAATACATAGATATATAATTAGGTA
TTCTATGTCTAAGAGGACAGCTATTTTTAGATATACTAACTTATAAACAGACAAAGTTAGACACAGAAAGTTGTTCATCCCAAAGATAATATGCTACAAA
AAGGAGTTTCTGACAAAAAAAATTTATGAAAGCTTTTGGATATTATCTGATTTTTCCAAATATAAATTTTATTGTCTCAGAGAATGCTGGTTTCATTTTT
TTAACTAAGTAGAAGGTTAcagtaaggaacagaattatcactaagaagattattgctccatgaagtttagtaaacaaagaatttaatgaaaaatagtaa
tgaagtgattgacaattcagaaatgttctatttcagaacagaatagtcatgttgaataattaggtatttatgcatcagaaaattcttcacataatgtctc
tgaaactgataggcaaagtactttcttgcacttaacaaaatgacttaactgcaaattttatattcgtttctttgatgttgaatattttactttatgaaag
cataggcctttcacttcatctacattttacaccagccttccattattcaaagtttgcaaattgttccatattaagaactcaagccagcatttttctttc
agaactataataaactttggcctgttcactaaaatttaaattttctgaatatttacttcagcagtaacaacatttactgcacagttacagtatgaagaa
atgttaagaaatagaaagcaccgtatgccgtctacaaaaagcttgcatatttgataaggcttgccaaagcttctccagtcaggaaagagattacaaagtt
cagatttttacttctggcaactattttaagtaatcagataatgtctgaatcaatgtcaactttataaatatagagtgaacgaaaaatattggggagaagg
aaaatttttaaaataattttcacttactaagctgatgcagaggtttatttattttaatttttttttttttttttttttttgagacagagtctcattctgtca
cccaggctggagtgtggtggcaccatctttgctcactgcaacctccacctcctgggtttaagcaattctcatgcctcagcctcccaagtagctgggacta
cagacatgggccaccacgcctggctaattttttgtatttgtagtagagacagggggtttcaccatgttagccaaactggtctcaaactcctgacctcaagta
atcctcccgcctcggcctcctaaagtgctgggattaccagcatgagccaccacgccaggcccagaggtttattttacatttatcacatgtcttgtctctc
taattacacttcactaaccctgcattgctgtggttttcttcttcaactctgaacctggttttccaccattggtaagactgtatattccctctttgaccac
tagggttcattccttgttctgtttggctaatacttttttccatttcccttacacttgctactttgggttaggctttaagtcaactcaatgatgaactat
gcaaaactgttcaagccacaaaagagaagacctgcccatcactctaggcccacagatgacctatggagcaatccatttggagaaacgtgtgtctctgcta
tctgcaaagcagctccagagtgacccttagctgggacaccctaatttatttctctcttcttttttcagAGTGCCTAGAAAGATGACAACTCAAGCACCGACG
TTTACGCAGCCGTTACAAAGCGTTGTGGTACTGGAGGGTAGTACCGCAACCTTTGAGGCTCACATTAGTGGTAAGCTCACACATTCACACTTTTGTTTTT
TTTTCCTTTGCCTCCTCTCCAGTAAGTTAACGTTGCTGCAGGACTTGACGCCAAGTTTAAGCCCTGCTTTCACTTCGGAAAATTAGTCCAACACTATGGA
...

## *Copy spliced transcript to clipboard + Comp*

This option will copy the mature transcript (after splicing) to the clipboard and and Comparison list. This means it will include the 5'UTR + Exons + 3'UTR

## *Copy only coding sequence to clipboard + Comp*

This option will copy only the bases used by the ribosome to produce proteins. In other words the coding sequence.

eg.

```
>TTN ; Homo sapiens titin (TTN), transcript variant N2-A, mRNA. (from RefSeq NM_133378) ; CDS length = 100272 (uc021vsy.4 Full gene
position = chr2:178525989-178807421)
```

**ATGA**CAACTCAAGCACCGACGTTTACGCAGCCGTTACAAAGCGTTGTGGTACTGGAGGGTAGTACCGCAACCTTTGAGGCTCACATTAGTGGTTTTCCAG

TTCCTGAGGTGAGCTGGTTTAGGGATGGCCAGGTGATTTCCACTTCCACTCTGCCCGGCGTGCAGATCTCCTTTAGCGATGGCCGCGCTAAACTGACGAT

CCCCGCCGTGACTAAAGCCAACAGTGGACGATATTCCCTGAAAGCCACCAATGGATCTGGACAAGCGACTAGTACTGCTGAGCTTCTCGTGAAAGCTGAG

ACAGCACCACCCAACTTCGTTCAACGACTGCAGAGCATGACCGTGAGACAAGGAAGCCAAGTGAGACTCCAAGTGAGAGTGACTGGAATCCCTACACCTG

….

TTCCATCTGATATCAGCATTGATGAAGGCAAAGTTCTAACAGTAGCCTGTGCTTTCACGGGTGAGCCTACCCCAGAAGTAACATGGTCCTGTGGTGGAAG

AAAAATCCACAGTCAAGAACAGGGGAGGTTCCACATTGAAAACACAGATGACCTGACAACCCTGATCATCATGGACGTACAGAAACAAGATGGTGGACTT

TATACCCTGAGTTTAGGGAATGAATTTGGATCTGACTCTGCCACTGTGAATATACATATTCGATCCATT<span style="color:red">TAA</span>

## *Copy Protein to clipboard + Comp*

This option will **use the appropriate Genetic Code table** selected in the "Controls" tab to translate the coding sequence into a protein and copy that amino acid letters as a FASTA string to the clipboard.

```
>TTN ; Homo sapiens titin (TTN), transcript variant N2-A, mRNA. (from RefSeq NM_133378) ; Protein length = 33423 (uc021vsy.4 Full
gene position = chr2:178525989-178807421)
```

**M**TTQAPTFTQPLQSVVVLEGSTATFEAHISGFPVPEVSWFRDGQVISTSTLPGVQISFSDGRAKLTIPAVTKANSGRYSLKATNGSGQATSTAELLVKAE

TAPPNFVQRLQSMTVRQGSQVRLQVRVTGIPTPVVKFYRDGAEIQSSLDFQISQEGDLYSLLIAEAYPEDSGTYSVNATNSVGRATSTAELLVQGEEEVP

AKKTKTIVSTAQISESRQTRIEKKIEAHFDARSIATVEMVIDGAAGQQLPHKTPPRIPPKPKSRSPTPPSIAAKAQLARQQSPSPIRHSPSPVRHVRAPT

PSPVRSVSPAARISTSPIRSVRSPLLMRKTQASTVATGPEVPPPWKQEGYVASSSEAEMRETTLTTSTQIRTEERWEGRYGVQEQVTISGAAGAAASVSA

…

SGKYTIKAKNFRGQCSATASLMVLPLVEEPSREVVLRTSGDTSLQGSFSSQSVQMSASKQEASFSSFSSSSASSMTEMKFASMSAQSMSSMQESFVEMSS

SSFMGISNMTQLESSTSKMLKAGIRGIPPKIEALPSDISIDEGKVLTVACAFTGEPTPEVTWSCGGRKIHSQEQGRFHIENTDDLTTLIIMDVQKQDGGL

YTLSLGNEFGSDSATVNIHIRS**I**

## *Copy Splice Info to clipboard*

This option will copy important length information related to splicing to the clipboard.

It looks as follows for the TTN gene:

TTN ; Homo sapiens titin (TTN), transcript variant N2-A, mRNA. (from RefSeq NM_133378) ; uc021vsy.4 (Full gene position = chr2:178525989-178807421)

**Transcript Length: 281433**

**Exon bases length: 101518**

**Intron bases length: 179915**


**Exon 1/312 length:210**

**Intron 1/311 length:2556**

Exon 2/312 length:104

Intron 2/311 length:2210

Exon 3/312 length:204

Intron 3/311 length:1455

Exon 4/312 length:288

Intron 4/311 length:484

Exon 5/312 length:86

Intron 5/311 length:93

Exon 6/312 length:245

Intron 6/311 length:4234

Exon 7/312 length:331

Intron 7/311 length:370

Exon 8/312 length:153

Intron 8/311 length:857

Exon 9/312 length:138

Intron 9/311 length:1206

Exon 10/312 length:126

Intron 10/311 length:1226

…

…

Exon 310/312 length:154

Intron 310/311 length:525

Exon 311/312 length:303

**Intron 311/311 length:138**

**Exon 312/312 length:1319**


Exon+Intron 1 length:2766

Exon+Intron 2 length:2314

Exon+Intron 3 length:1659

…

## Show gene homology in browser

This option **opens the NIH website** with the Known Gene name eg. TTN, BRCA1, TP53 etc. in order to find homologs fo this gene in other organisms.

**HomoloGene:130650. Gene conserved in Euteleostomi**

**Genes**
*Genes identified as putative homologs of one another during the construction of HomoloGene.*

**Proteins**
*Proteins used in sequence comparisons and their conserved domain architectures.*

TTN, *H.sapiens*
titin

NP_001254479.1
35991 aa

LOC703527, *M.mulatta*
titin-like

XP_002808058.1
33365 aa

TTN, *C.lupus*
titin

XP_535982.5
35162 aa

TTN, *B.taurus*
titin

XP_002685306.2
34369 aa

Ttn, *M.musculus*
titin

NP_035782.3
33467 aa

It effectively calls: https://www.ncbi.nlm.nih.gov/homologene/?term=TTN

When you now click on the alignments highlighted in red, you can get sequence alignments between the organisms.

**Sequence Alignment** ☐ include c

Reformat | Format: Hypertext ▼ | Row Display: up to 10 ▼ | Color Bits: 2.0 bit ▼ | Type Selection: the most diver

```
               10        20        30        40        50        60        70        80
         ....*....|....*....|....*....|....*....|....*....|....*....|....*....|....*....|
3MFR_A        26 YELCEVIGKGAFSVVRRCINRE-TGQQFAVKIVDVAKFTSSPglstEDLKREASICHML---------KHP------HIV  89
gi 297474669 220 YEVLKVIGKGSFGQVVKAYDHK-VHQHVALKMVRNEKR--FH----RQAAEEIRILEHLr-------kQDK------DNT 279
gi 120537328 150 YEIDSLIGKGSFGQVVKAYDRV-EQEWVAIKIIKNKKAF--L----NQAQIEVRLLELMnk-----hdTEM------KYY 211
gi 19113931   11 LTDIRHLTDGTISEVFVGERKN-SKKLYVIKVQGLVFKR-PP----HDAMRGKLILESI---------GHP------HIE  69
gi 151945999  39 VTNHNSLGDGNFSVVKECMNIH-TKDLYAMKLIKKQTVKNKI----QLIQREFDLLRSIsekirdmekKNE------HSL 107
gi 25146830  147 YEVLEFLGKGTFGQVVKAWKKG-TSEIVAIKILKKHPS--YA----RQGQIEVSILSRLsne----nsEEF------NFV 209
gi 19075761  159 YIVQSNLGKGMFSTVVSALDRN-RNQTFAIKIIRNNEVM--Y----KEGLKEVSILERLqaa---dreGKQ------HII 222
gi 291409415 674 YNVYGYTGQGVFSNVVRARDNArANQEVAVKIIRNNEL--MQ----KTGLKELEFLKKLnda--dpdDKF------HCL 738
gi 6322320   369 YLVLDILGQGTFGQVVKCQNLL-TKEILAVKVVKSRTE--YL----TQSITEAKILELLnq------kIDPt-----NKH 430
3KVW_A        99 YEVLKVIGKGSFGQVVKAYDHK-VHQHVALKMVRNEKR--FH----RQAAEEIRILEHL---------RKQdkdntmNVI 162

               90       100       110       120       130       140       150       160
         ....*....|....*....|....*....|....*....|....*....|....*....|....*....|....*....|
3MFR_A        90 ELLETYSSD----------GMLYMVFEFMDGa--DLCFEIVKRAdagfV--YSEAVASHYMRQILEALRYCHDN--NIIH 153
gi 297474669 280 MNVIHMLENft------frNHICMTFELLSM---NLYELIKKNK----FqgFSLPLVRKFAHSILQCLDALHKN--RIIH 344
gi 120537328 212 IVHLKRHFMf--------rNHLCLVFEMLSy---NLYDLLRNTN----FrgVSLNLTRKFAQQMCTALLFLATPelSIIH 276
gi 19113931   70 RIVDSFIDNe---------aGSVYLITSFKSF--VLSDVMDE---------ISIDTKCKIVLQISSALEYLEKH--GILH 127
gi 151945999 108 DIFEGHHHIlqlfdyfetaDNIVLITQLCQKg--DLYEKIVENQ----Cl-DLETQVTSYCACLVSVLEFLHSQ--GIVH 178
gi 25146830  210 RAFECFNHK----------SHTCLVFEMLEQ---NLYDFLKQNK----FmpLPLNAIRPILFQVLTALLKLKSL--GLIH 270
gi 19075761  223 HYERHFMHK----------NHLCMVFEMLSLnlrDILKKFGRNV----G--LSIKAVRLYAYQMFMALDLLKQC--NVIH 284
gi 291409415 739 RLFRHFYHK----------QHLCLVFEPLSMnlrEVLKKYGKDV----G--LHIKAVRSYSQQLFLALKLLKRC--NILH 800
gi 6322320   431 HFLRMYDSFvh-------kNHLCLVFELLSN---NLYELLKQNK----FhgLSIQLIRTFTTQILDSLCVLKES--KLIH 494
3KVW_A       163 HMLENFTFR----------NHICMTFELLSM---NLYELIKKNK----FqgFSLPLVRKFAHSILQCLDALHKN--RIIH 223
```

## Show biochemical pathway in browser

Many of the genes in the UCSC table data, especially genes coding for enzymes involved in biochemical pathways, contains reference numbers to the KEGG biochemical pathway maps.

When you now select "Show biochemical pathway in browser", you will be taken to the KEGG pathway where this gene is found:

http://www.kegg.jp/pathway/hsa04010+TP53



From here you can click on the p53 signalling pathway to drill down further.

Take note that this only provides a link to online information and is not part of this software.

Searching or filtering for genes with the **pathway+** search term, will result in all the genes containing pathway information. This can also be obtained for all loaded genes (after clicking the "**LOAD GENES**" button).

Any time you looking at genes on the genome, you can click the "**Show pathways**" button. This will load all displayed genes (which have pathways) into a pathway tree, which will group the genes below their corresponding biochemical pathway.

For example: Loading the genes for the current sequence and then clicking **"Show pathways"** after entering the search term TP53, will load only the genes related to TP53 in the pathway tree.





From here you can:

- double click on the genes to navigate to them

- double click on the pathway to filter only the genes containing those pathways and navigate to the first one of those

- Right click on the gene to get the gene context menu in order to execute all the usual actions.

## Show Protein Sequence

All genes which code for proteins, are marked in the context menus with **CODING** at the end.

This indicates that the UCSC annotation contains the coding sequence or ORF position information indicating where translation begins and ends.

One example

| chr2:190691236-190691326 L=91 | ▶ |
| --- | --- |
| (9593) + Homo sapiens NGFI-A binding pr…, NAB1 L=2522 Exons=10 (uc002usb.4) CODING | ▶ |
| (9596) + NGFI-A binding protein 1 (EGR1…, NAB1 L=982 Exons=7 (uc002usc.4) CODING | ▶ |
| (9597) + The sequence shown here is der…, NAB1 L=1029 Exons=6 (uc061qsg.1) CODING | ▶ |

and another

| chr3:36993776-36993850 L=75 | ▶ |
| --- | --- |
| (1849) + Homo sapiens mutL homolog 1 (M…, MLH1 L=115 Exons=20 (uc003cgn.5) CODING | ▶ |

When you select the **"Show Protein Sequence"** menu option, the software will perform essentially the same process that happens in a cell. It will join all the bases in the exons (like the Spliceosome) and only extract the coding sequence, after which it will perform protein translation as it would happen in the Ribosome.

The Genetic Codon table it will use is the one selected below the "Controls" tab:



The amino acid primary structure sequence will then be displayed in the **Protein** tab.

This action can be obtained from the context menu, or by double clicking on the gene in the **Main Genome View**. When you have highlighted the gene using cursor keys (which works in both the **Main Genome View** or the **DNA View**) you can also press **F5** which will always draw or redraw the current protein view. If you want to search for specific peptide signalling sequences in the protein, you can enter the amino acid letters in the **Search field** and it will be highlighted in the **Protein View**.

Something else to observe is that the **Protein View** and **DNA View** is navigationally linked. When you hover the mouse over a specific amino acid, the **DNA** View navigates and centres on the corresponding codon, which is coloured with the same colour when the "**Codons**" checkbox is checked. The display is scrollable when the entire sequence does not fit.



You can also change the number of amino acids in each line, by using the **[** and **]** keys, or by changing it in the text box:



When the amino acids per line becomes small enough, the display also starts to display the amino acid index (given that the option Show Amino Acid Numbers is enabled in the settings), which you also get when you hover over any amino acid in the display .

When you are centring the **DNA View** on a different codon, and the **Protein View** is currently displaying the exact same gene's protein, then the corresponding amino acid is highlighted in the **Protein View** as follows:



When there are multiple overlapping genes in the **DNA View** you can step through them by using the + **and** − keys.

When you want to copy and paste a specific range of amino acids you can copy them by double clicking on the first and last amino acid and the result will be copied to the **Information** tab and clipboard. Both the DNA coding sequence and the amino acid sequence will be shown (but only for the region YOU selected).

```
>Amino acids 1-514 (Len=515) Bases=37014478->37050648 for gene: MLH1 on + strand
MNGYISNANYSVKKCIFLLFINHRLVESTSLRKAIETVYAAYLPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLLGSNSSR
MYFTQTLLPGLAGPSGEMVKSTTSLTSSSTSGSSDKVYAHQMVRTDSREQKLDAFLQPLSKPLSSQPQAIVTEDKTDISSGRARQQDEEMLELPAPAEVA
AKNQSLEGDTTKGTSEMSEKRGPTSSNPRKRHREDSDVEMVEDDSRKEMTAACTPRRRIINLTSVLSLQEEINEQGHEVLREMLHNHSFVGCVNPQWALA
QHQTKLYLLNTTKLSEELFYQILIYDFANFGVLRLSEPAPLFDLAMLALDSPESGWTEEDGPKEGLAEYIVEFLKKKAEMLADYFSLEIDEEGNLIGLPL
LIDNYVPPLEGLPIFILRLATEVNWDEEKECFESLSKECAMFYSIRKQYISEESTLSGQQSEVPGSIPNSWKWTVEHIVYKALRSHILPPKHFTEDGNIL
QLANLPDLYKVFERC

>DNA Bases=37014478->37050648 (Bases=1545/1545)   for gene: MLH1 on + strand
ATGAATGGTTACATATCCAATGCAAACTACTCAGTGAAGAAGTGCATCTTCTTACTCTTCATCAACCATCGTCTGGTAGAATCAACTTCCTTGAGAAAAG
CCATAGAAACAGTGTATGCAGCCTATTTGCCCAAAAACACACACCCATTCCTGTACCTCAGTTTAGAAATCAGTCCCCAGAATGTGGATGTTAATGTGCA
CCCCACAAAGCATGAAGTTCACTTCCTGCACGAGGAGGAGCATCCTGGAGCGGGTGCAGCAGCACATCGAGAGCAAGCTCCTGGGCTCCAATTCCTCCAGG
ATGTACTTCACCCAGACTTTGCTACCAGGACTTGCTGGCCCCTCTGGGGAGATGGTTAAATCCACAACAAGTCTGACCTCGTCTTCTACTTCTGGAAGTA
GTGATAAGGTCTATGCCCACCAGATGGTTCGTACAGATTCCCGGGAACAGAAGCTTGATGCATTTCTGCAGCCTCTGAGCAAACCCCTGTCCAGTCAGCC
CCAGGCCATTGTCACAGAGGATAAGACAGATATTTCTAGTGGCAGGGCTAGGCAGCAAGATGAGGAGATGCTTGAACTCCCAGCCCCTGCTGAAGTGGCT
GCCAAAAATCAGAGCTTGGAGGGGGATACAACAAAGGGGACTTCAGAAATGTCAGAGAAGAGAGGACCTACTTCCAGCAACCCCAGAAAGAGACATCGGG
AAGATTCTGATGTGGAAATGGTGGAAGATGATTCCCGAAAGGAAATGACTGCAGCTTGTACCCCCCGGAGAAGGATCATTAACCTCACTAGTGTTTTGAG
TCTCCAGGAAGAAATTAATGAGCAGGGACATGAGGTTCTCCGGGAGATGTTGCATAACCACTCCTTCGTGGGCTGTGTGAATCCTCAGTGGGCCTTGGCA
CAGCATCAAACCAAGTTATACCTTCTCAACACCACCAAGCTTAGTGAAGAACTGTTCTACCAGATACTCATTTATGATTTTGCCAATTTTGGTGTTCTCA
GGTTATCGGAGCCAGCACCGCTCTTTGACCTTGCCATGCTTGCCTTAGATAGTCCAGAGAGTGGCTGGACAGAGGAAGATGGTCCCAAAGAAGGACTTGC
TGAATACATTGTTGAGTTTCTGAAGAAGAAGGCTGAGATGCTTGCAGACTATTTCTCTTTGGAAATTGATGAGGAAGGGAACCTGATTGGATTACCCCTT
CTGATTGACAACTATGTGCCCCCTTTGGAGGGACTGCCTATCTTCATTCTTCGACTAGCCACTGAGGTGAATTGGGACGAAGAAAAGGAATGTTTTGAAA
GCCTCAGTAAAGAATGCGCTATGTTCTATTCCATCCGGAAGCAGTACATATCTGAGGAGTCGACCCTCTCAGGCCAGCAGAGTGAAGTGCCTGGCTCCAT
TCCAAACTCCTGGAAGTGGACTGTGGAACACATTGTCTATAAAGCCTTGCGCTCACACATTCTGCCTCCTAAACATTTCACAGAAGATGGAAATATCCTG
CAGCTTGCTAACCTGCCTGATCTATACAAAGTCTTTGAGAGGTGT
```

## Show PDB Structure in PYMOL

Many of the gene annotations from the UCSC table data comes with the Protein Database (PDB) codes for protein coding genes. These PDB files contains the 3D Crystallography structure obtained for these genes.

When you have the open source Molecular visualization software PYMOL installed, the Visual Genome Browser will firstly download a set number of PDB files referenced in the gene annotations, create a hydrophobic colour scheme and then launch the PDB files in PYMOL.

"PyMOL is an open-source, user-sponsored, molecular visualization system created by Warren Lyford DeLano and commercialized initially by DeLano Scientific LLC, which was a private software company dedicated to creating useful tools that become universally accessible to scientific and educational communities. It is currently commercialized by Schrödinger, Inc. PyMOL can produce high-quality 3D images of small molecules and biological macromolecules, such as proteins."

When one selects the TBP (TATA-box binding protein, which is a transcription factor responsible for recruiting RNA Polymerase II to genes containing the TATA sequence promoters) this option shows the Protein structure determined for this protein in PYMOL on the following page.



*Illustration 52: TBP (TATA Box binding protein crystal structure)*

## Show Protein sequence with PDB Secondary Structure

I was looking for a way to incorporate the secondary structure of a protein (alpha helix, beta sheet, loop etc.) into the 2D protein view and I realised that I could try to do an alignment between the Protein Sequence View amino acids and the amino acids in the PDB files from the Protein Data Bank and then extract the secondary structure information from the amino acids in the PDB files and then associate that information with the corresponding amino acids in the Protein Sequence view.  In this was I did not have to resort to trying to predict the secondary structure using molecular angles etc., but I could extract it directly from the crystallography structures.

http://www.rcsb.org/pdb/home/home.do

You obtain this by selecting **Show Protein Sequence with PDB Sequence** Alignment. That results in a protein display which looks as follows (Notice the correspondence between the first amino acids of the crystal structure and the highlighted amino acids SGIV in the image below):



The colours represent different secondary structures:
Blue = alpha helix
Green = Loop
Yellow = Beta sheet
Gray = Residues for which no alignment could be found with the PDB (This also gives us an indication for which part of the protein the crystal structure was determined)

## Red markers for amino acids at start of exons

The Red blocks are markers for the amino acids which are at the start of each exon.  This is sometimes useful to have in order to see how subsequent and similar protein coding sequences, which are part of the same protein coding gene, can be found so far apart and separated by many intron bases. The display of amino acids at the start of exons as red blocks can also be switched off in the Settings. The fact that the amino acids at the start of each exon is marked in red, will allow you to get some insight into the different segments of the protein which might be spliced out for alternatively spliced transcripts.

I referred to these patterns in proteins on my blog:
http://splicejunction.blogspot.com.au/2015/08/protein-scrabble.html

## Show Splicing Graph

I wanted to get a way to "graphically visualise" the introns and exons of a transcribed gene as there can sometimes be hundreds of exons which are spliced together to produce the final mature RNA transcript.

I came up with the concept of a **Splicing Graph**. I knew from the theory on the Spliceosome, that proteins complexed with RNA is able to recognize the splice signals (GT....AG) and then correctly form a loop (Lariat) which is cut out before the exons are joined or ligated together.

I discussed this on my blog: http://splicejunction.blogspot.com.au/2015/09/spliced-genes-natures-model-hobby-kit_18.html

I decided to draw the introns which "loop out" as vertical lines which represent these loops **to scale** on the vertical axis of these graphs. Each line therefore represented a loop going downwards up to half the length of the removed intron and half the length back again. **The vertical lines therefore mimicked the real loop which is formed during splicing.** In a similar way, the lengths of the **exons** are represented **to scale** on the horizontal axis of the graph. In this way you can get a feeling of the length spacing of the exons, while simultaneously showing the differences in the intron lengths as vertical lines.

When I create a **Splicing Graph** of the TTN gene transcript : RefSeq NM_133378
ucsc Id: uc021vsy.4



*Illustration 53: Splicing Graph for the TTN gene with a total of 312 exons and 311 Intron which are represented by the vertical lines*

When I do the same for a gene for Collagen, which has less exons (66 in total) I get the following Splicing Graph.

The spacing of the vertical lines are determined by the lengths of the exons between them.



*Illustration 55: Splicing Graph of the gene for collagen type V alpha 1 (COL5A1). Intron loops are on the vertical axis while exons are the spacing between the vertical lines on the horizontal axis.*

*Illustration 56: Protein View for COL5A1 gene at width 18 residues*

*If one is rather interested in the exon lengths*, it is possible to change the splicing view setting in order to rather show the exon lengths.  For exactly the same collagen protein here is a picture showing the exon lengths:



When splicing graph settig is changed, ths length represents the EXON length.

*Illustration 57: Splicing Graph of the gene for collagen type V alpha 1 (COL5A1). In this case EXONS are on the vertical and horizontal axes.  This provides a quick way to visualise the exons which makes up the coding sequence.*

The setting can be changed in the **Settings tab**.

For example, look how the marked amino acids at the start of each exon form a pattern in the **TTN** protein. The distribution of the exon lengths seems fall in a relatively narrow band throughout a big part of the protein. This is a block representation of only part of the whole protein.
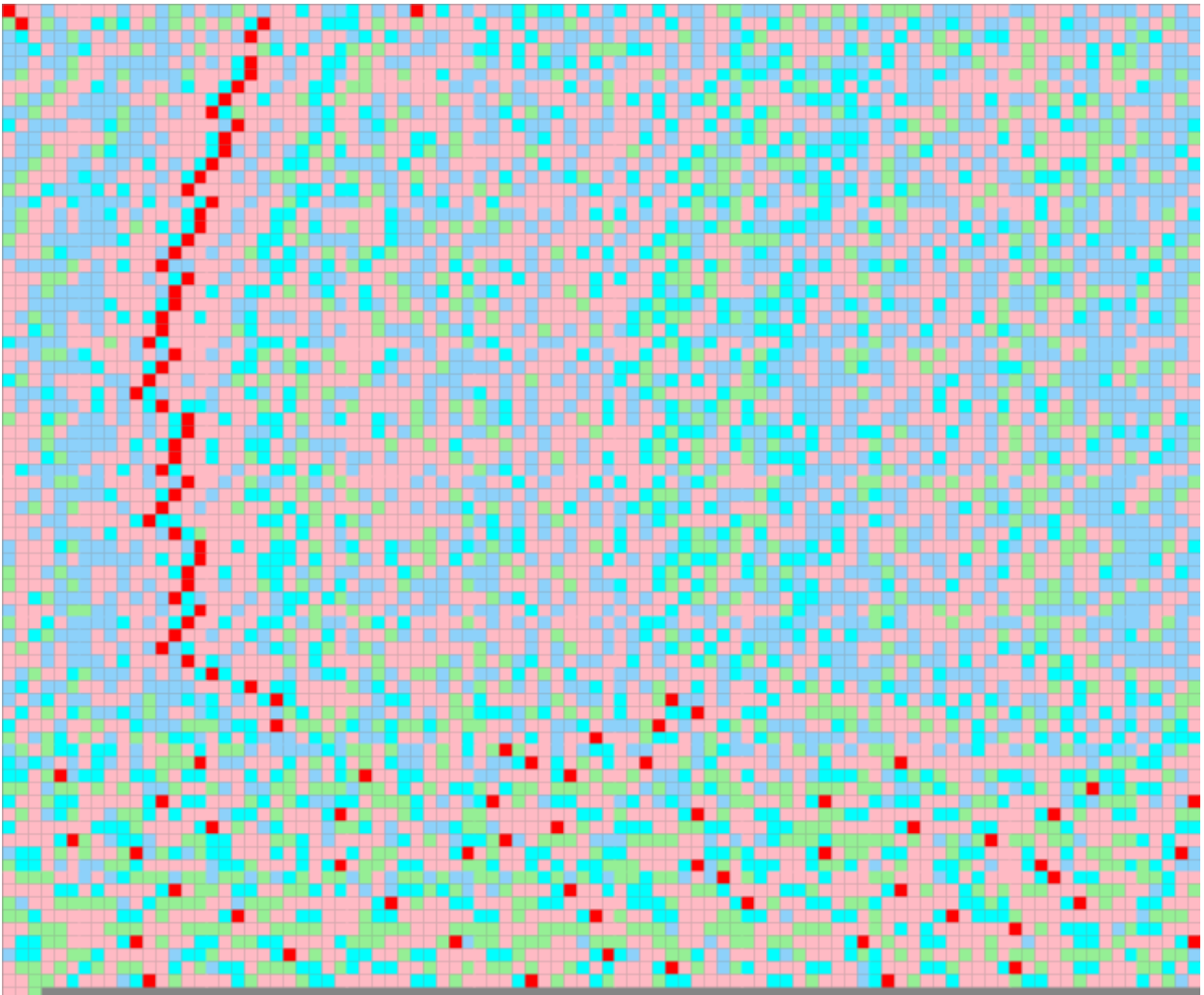


*Illustration 58: Patterns in the protein view made by amino acids at the start of exons. (For the Titin TTN gene on chromosome 2)*

The **TITIN** gene codes for a large abundant protein in the human body. The gene is quite probably the largest gene in the human genome. The protein effectively acts as **shock absorbers** for muscle fibres (sarcomeres) and needs to be structurally very long. In the picture above, remember that the pink areas represent **hydrophobic** residues, green negatively charged polar residues, blue uncharged polar residues and the cyan positively charged residues.

Here is an extract from the UCSC summary information for this gene:

*"**This gene encodes a large abundant protein of striated muscle**. The product of this gene is divided into two regions, a N-terminal I-band and a C-terminal A-band. ... An N-terminal Z-disc region and a C-terminal M-line region bind to the Z-line and M-line of the sarcomere, respectively, so that a single titin molecule spans half the length of a sarcomere. **Titin also contains binding sites for muscle associated proteins so it serves as an adhesion template for the assembly of contractile machinery in muscle cells.** It has also been identified as a structural protein for chromosomes. ... Mutations in this gene are associated with familial hypertrophic cardiomyopathy 9, and auto-antibodies to **titin** are produced in patients with the autoimmune disease scleroderma."*

# Protein Sequence display of various genes

Next, we look at the primary structure sequence depiction of the collagen COLL11A1 and COL protein product, the amino acids at the start of each exon is marked with red blocks. Interestingly, observe how often the amino acid at the start of each exon is a **Glycine**. This is because the consensus sequence for an intron (in IUPAC nucleic acid notation) is: G-G-[cut]-G-T-R-A-G-T (donor site) ... intron sequence ... Y-T-R-A-C (branch sequence 20-50 nucleotides upstream of acceptor site) ... Y-rich-N-C-**A-G-[cut]-G** (acceptor site).

This means that the first base following the intron is very often equal to **G**. When we refer back to the genetic code, we find that the first amino acids at the start of exons will often be one of the following: **G=Glycine (GGG, GGA, GGC, GGT)**, **A=Alanine (GCG, GCA,GCC,GCT)**, **E=Glutamic Acid (GAG,GAA) , D=Aspartic Acid (GAC,GAT) or V=Valine (GTG,GTA,GTC,GTT).**

```
M L S F V D T R T L L L L A V T L C L A T C Q F F E T V R K G P A G
D R G P R G E R G P P G P P G R D G E D G P T G P P G P P G P P G P
P G L G G N F A A Q Y D G K G V G L G P G P M G L M G P R G P P G A
A G A P G P Q G F Q G P A G E P G E P G Q T G P A G A R G P A G P P
G K A G E D G H P G K P G R P G E R G V V G P Q G A R G F P G T P G
L P G F K G I R G H N G L D G L K G Q P G A P G V K G E P G A P G E
N G T P G Q T G A R G L P G E R G R V G A P G P A G A R G S D G S V
G P V G P A G P I G S A G P P G F P G A P G P K G E I G A V G N A G
P A G P A G P R G E V G L P G L S G P V G P P G N P G A N G L T G A
K G A A G L P G V A G A P G L P G P R G I P G P V G A A G A T G A R
G L V G E P G P A G S K G E S G N K G E P G S A G P Q G P P G P S G
E E G K R G P N G E A G S A G P P G P P G L R G S P G S R G L P G A
D G R A G V M G P P G S R G A S G P A G V R G P N G D A G R P G E P
G L M G P R G L P G S P G N I G P A G K E G P V G L P G I D G R P G
P I G P A G A R G E P G N I G F P G P K G P T G D P G K N G D K G H
A G L A G A R G A P G P D G N N G A Q G P P G P Q G V Q G G K G E Q
G P P G P P G F Q G L P G P S G P A G E V G K P G E R G L H G E F G
L P G P A G P R G E R G P P G E S G A A G P T G P I G S R G P S G P
P G P D G N K G E P G V V G A V G T A G P S G P S G L P G E R G A A
G I P G G K G E K G E P G L R G E I G N P G R D G A R G A P G A V G
A P G P A G A T G D R G E A G A A G P A G P A G P R G S P G E R G E
V G P A G P N G F A G P A G A A G Q P G A K G E R G A K G P K G E N
G V V G P T G P V G A A G P A G P N G P P G P A G S R G D G G P P G
M T G F P G A A G R T G P P G P S G I S G P P G P P G P A G K E G L
R G P R G D Q G P V G R T G E V G A V G P P G F A G E K G P S G E A
G T A G P P G T P G P Q G L L G A P G I L G L P G S R G E R G L P G
V A G A V G E P G P L G I A G P P G A R G P P G A V G S P G V N G A
P G E A G R D G N P G N D G P P G R D G Q P G H K G E R G Y P G N I
G P V G A A G A P G P H G P V G P A G K H G N R G E T G P S G P V G
P A G A V G P R G P S G P Q G I R G D K G E P G E K G P R G L P G L
K G H N G L Q G L P G I A G H H G D Q G A P G S V G P A G P R G P A
G P S G P A G K D G R T G H P G T V G P A G I R G P Q G H Q G P A G
P P G P G P P G P P G V S G G G Y D F G Y D G D F Y R A D Q P R S
A P S L R P K D Y E V D A T L K S L N N Q I E T L L L T P E G S R K N
P A R T C R D L R L S H P E W S S G Y Y W I D P N Q G C T M D A I K
V Y C D F S T G E T C I R A Q P E N I P A K N W Y R S S K D K K H V
W L G E T I N A G S Q F E Y N V E G V T S K E M A T Q L A F M R L L
A N Y A S Q N I T Y H C K N S I A Y M D E E T G N L K K A V I L Q G
S N D V E L V A E G N S R F T Y T V L V D G C S K K K T N E W G K T I
I E Y K T N K P S R L P F L D I A P L D I G G A D Q E F F V D I G P
V C F K
```

*Illustration 59: Pattern formed by amino acid at start of exons for Collagen COL1A2 at 34 amino acids per line*

*Illustration 60: Similar patterns in the protein of Collagen COL11A1 at width 36 residues per line*

*Different representations of the Retinitis pigmentosa GTPase regulator (RPGR) protein (Id = uc004ded.2)*

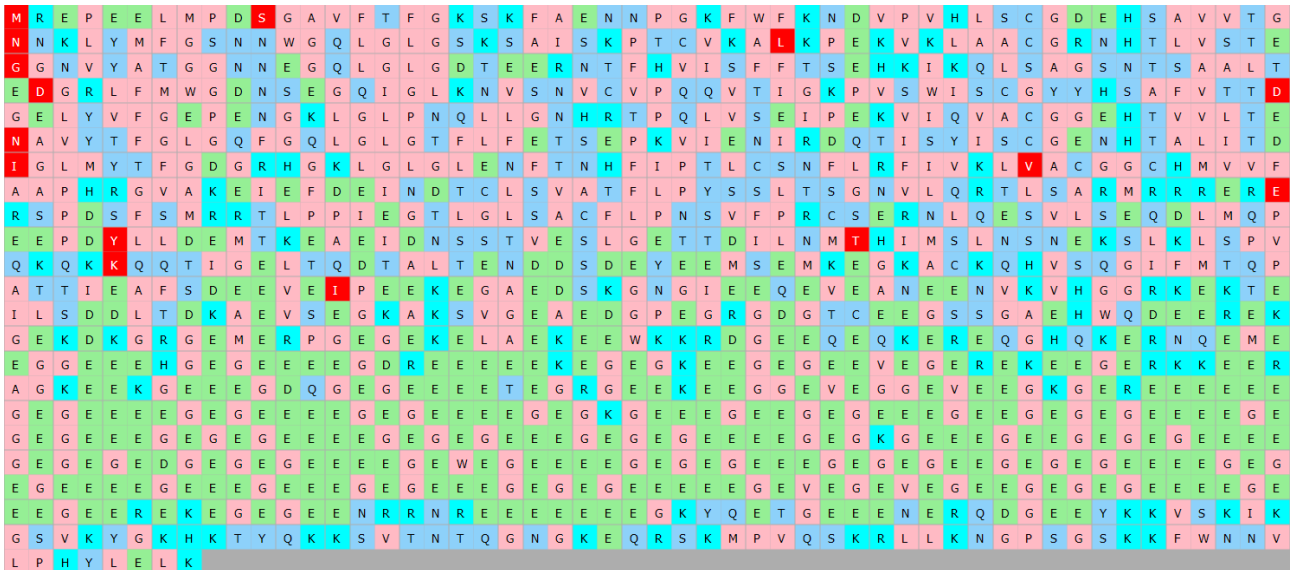*On the Protein domain as colour coded amino acids:*



*Illustration 61: RPGR protein product at chrX:38327340-38327544 at 52 residues width*

*The initial pink and blue parts represent the first part of the protein which is embedded in the membrane*

*In the DNA View with codons colour coded (only the last 2 exons)*



*Illustration 62: The last exon of the negative strand encoded gene on the DNA domain at 52 bases per line*

*Illustration 63: The same RPGR protein product with secondary protein structure colour coded at 52 residues width (Yellow=beta sheet, Blue=alpha helix, Green=Loop)*

The tooltip on the amino acids indicates that the section that was aligned in order to obtain the secondary structure was the PDB crystal structure: 4QAM and the PDB atoms which were matched to the specific amino acid were 2207-2217



When the 3D protein structure is now launched in Pymol with exactly the same colour coding:



Notice how the transmembrane occupying amino acids are mostly hydrophobic (pink) due to the requirement of interacting with the hydrophobic (oily) membrane.

# The Genetic Code - How the colour scheme was chosen

The colour scheme is an indication of how hydrophobic or polar the amino acids are. Protein structures fold in water spontaneously due to the interactions the side chains of the amino acids have with water. Hydrophobic amino acids tend to move away from polar water molecules and therefore tend to be on the inside of folded proteins (OR on the outside when embedded in the hydrophobic lipid membranes of a cell membrane).

Colour coding the primary structure of a protein based on the polarity of the amino acids gives you a better idea of which amino acids of the protein will be on the inside and which will be on the outside, by just looking at the sequence. It is however not always as simple and one also have to look at the 3D secondary structure such as alpha helices, beta sheets and loops.

When this colour coding is extended from the protein to the DNA domain by doing the same for the corresponding codons, it allows you to better understand which part of the protein is coded by different exons when the coding sequence is assembled.

It is also known that amino acids with similar properties (such as hydrophobicity) can often be substituted for each other without much detrimental effect on the protein folding. A missense mutation in the DNA is when the change in codon base results in a different amino acid. When a single base missense mutation causes the translation of an amino acid with similar hydrophobic properties as the original amino acid it is known as a **conservative missense mutation**.

I therefore decided to arrange the Genetic Codon table in such a fashion that there will always be only a single base change between adjacent codons in the table. In digital communication this is know as a Hamming Code and it is used during error correction coding. Codes are given a redundancy by only choosing a subset of the available codes as valid codes, in other words, with a hamming distance of 3. When single bit changes are introduced in transmission, one can then determine the correct original bit by choosing the code words which are closest in Hamming Distance to the received code.

I ended up with a Codon Table which resembles a "torus" where the codons in the top row is also only one base away from the codons of the bottom row and the same is valid in the left and right directions. A Genetic Code arranged in this way will therefore only code for the directly adjacent amino acids when single base changes are introduced by mutation. Diagonal movement represents 2 base changes (except when looping to the other side of the torus).

After colour coding the amino acids coded for by those codons and adding the hydrophobicity value of the amino acids in the centre, an **interesting "periodic table like"** arrangement emerged where amino acids with similar properties tended to cluster together.

Moving horizontally left or right **represents a change in the last "wobble base"** of the codon, which we know is due to the ability of Uracil to pair with either an A or a G in the anti-codon loop of the tRNA. A horizontal base change very often results in a **silent mutation**, which results in a mutated codon still coding for the same amino acid (eg. Alanine (GCC) → Alanine (GCT))

Moving up or down by one base mostly results in a different amino acid with similar properties, in other words: a **conservative missense mutation**.

In order to even further accentuate the Genetic Code's robustness against base changes, I decided to incorporate a matrix which represents how often amino acids are substituted in nature without detrimental effects to resulting proteins. I plotted the corresponding substitution frequency in the corners of the Codon Table.

| | Gly | Ala | Val | Leu | Ile | Met | Cys | Ser | Thr | Asn | Gln | Asp | Glu | Lys | Arg | His | Phe | Tyr | Trp |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Gly** | | | | | | | | | | | | | | | | | | | |
| **Ala** | 58 | | | | | | | | | | | | | | | | | | |
| **Val** | 10 | 37 | | | | | | | | | | | | | | | | | |
| **Leu** | 2 | 10 | 30 | | | | | | | | | | | | | | | | |
| **Ile** | | 7 | 66 | 25 | | | | | | | | | | | | | | | |
| **Met** | 1 | 3 | 8 | 21 | 6 | | | | | | | | | | | | | | |
| **Cys** | 1 | 3 | 3 | | 2 | | | | | | | | | | | | | | |
| **Ser** | 45 | 77 | 4 | 3 | 2 | 2 | 12 | | | | | | | | | | | | |
| **Thr** | 5 | 59 | 19 | 5 | 13 | 3 | 1 | 70 | | | | | | | | | | | |
| **Asn** | 16 | 11 | 1 | 4 | 4 | | | 43 | 17 | | | | | | | | | | |
| **Gln** | 3 | 9 | 3 | 8 | 1 | 2 | | 5 | 4 | 5 | | | | | | | | | |
| **Asp** | 16 | 15 | 2 | | 1 | | | 10 | 6 | 53 | 8 | | | | | | | | |
| **Glu** | 11 | 27 | 4 | 2 | 4 | 1 | | 9 | 3 | 9 | 42 | 83 | | | | | | | |
| **Lys** | 6 | 6 | 2 | 4 | 4 | 9 | | 17 | 20 | 32 | 15 | | 10 | | | | | | |
| **Arg** | 1 | 3 | 2 | 2 | 3 | 2 | 1 | 14 | 2 | 2 | 12 | 9 | | 48 | | | | | |
| **His** | 1 | 2 | 3 | 4 | | | 1 | 3 | 1 | 23 | 24 | 4 | 2 | 2 | 10 | | | | |
| **Phe** | 2 | 2 | 1 | 17 | 9 | 2 | | 4 | 1 | 1 | | | | 1 | 2 | | | | |
| **Tyr** | | 2 | 2 | 2 | 1 | | 3 | 2 | 2 | 4 | | | 1 | 1 | | 4 | 26 | | |
| **Trp** | | | | 1 | | | | 2 | | | | | | | 3 | | 1 | 1 | |
| **Pro** | 5 | 35 | 5 | 4 | 1 | | 1 | 27 | 7 | 3 | 9 | 1 | 4 | 4 | 7 | 5 | 1 | | |

Notice how the frequency is much higher between adjacent amino acids with similar colour (i.e. hydrophobicity), but noticeably lower when the colour is different. The robustness of the genetic code has to do with finding the optimal arrangement where the non-detrimental substitution frequency is maximal across all codons.

| Glutamic acid (E) | Glutamic acid (E) | Aspartic acid (D) | Aspartic acid (D) | Histidine (H) |
|---|---|---|---|---|
| 2  GAG | GAA  83│83 | GAC  16│16 | GAT  4│4 | CAT |
| (-3.5) Pol- | (-3.5) Pol- | (-3.5) Pol- | (-3.5) Pol- | (-3.2) Pol+ |
| 11  11│11 | 11  11│16 | 16  16│16 | 16  9│1 | 10  10│10 |
| 11  11│11 | 11  16│11 | 16  16│16 | 16  1│9 | 10  10│10 |
| Glycine (G) | Glycine (G) | Glycine (G) | Glycine (G) | Arginine (R) |
| GGG | GGA | GGC | GGT  1│1 | CGT |
| (-0.4) HP | (-0.4) HP | (-0.4) HP | (-0.4) HP | (-4.5) Pol+ |
| 58  58│58 | 58  58│58 | 58  58│58 | 58  5│3 | 7  7│7 |
| 58  58│58 | 58  58│58 | 58  58│58 | 58  3│5 | 7  7│7 |
| Alanine (A) | Alanine (A) | Alanine (A) | Alanine (A) | Proline (P) |
| 5  GCG | GCA | GCC | GCT  35│35 | CCT |
| (1.8) HP | (1.8) HP | (1.8) HP | (1.8) HP | (-1.6) HP |
| 0  37  37│37 | 37  37│37 | 37  37│37 | 37  10│5 | 4  4│4 |
| 37  37│37 | 37  37│37 | 37  37│37 | 37  5│10 | 4  4│4 |
| Valine (V) | Valine (V) | Valine (V) | Valine (V) | Leucine (L) |
| 0  GTG | GTA | GTC | GTT  30│30 | CTT |
| (4.2) HP | (4.2) HP | (4.2) HP | (4.2) HP | (3.8) HP |
| 0  8  66│8 | 66  66│66 | 66  66│66 | 66  1│25 | 17  17│17 |
| 1  8  8│66 | 66  66│66 | 66  66│66 | 66  25│1 | 17  17│17 |
| Methionine (M) | Isoleucine (I) | Isoleucine (I) | Isoleucine (I) | Phenylalanine (F) P |
| 1  ATG  6│6 | ATA | ATC | ATT  9│9 | TTT |
| (1.9) HP | (4.5) HP | (4.5) HP | (4.5) HP | (2.8) HP |

# The Colour Coded (Hamming code arranged) Genetic Codon Table

| | | | | | | | | Row |
|---|---|---|---|---|---|---|---|---|
| Glutamic acid (E) GAG (-3.5) Pol- | Glutamic acid (E) GAA (-3.5) Pol- | Aspartic acid (D) GAC (-3.5) Pol- | Aspartic acid (D) GAT (-3.5) Pol- | Histidine (H) CAT (-3.2) Pol+ | Histidine (H) CAC (-3.2) Pol+ | Glutamine (Q) CAA (-3.5) Pol | Glutamine (Q) CAG (-3.5) Pol | _A_ |
| Glycine (G) GGG (-0.4) HP | Glycine (G) GGA (-0.4) HP | Glycine (G) GGC (-0.4) HP | Glycine (G) GGT (-0.4) HP | Arginine (R) CGT (-4.5) Pol+ | Arginine (R) CGC (-4.5) Pol+ | Arginine (R) CGA (-4.5) Pol+ | Arginine (R) CGG (-4.5) Pol+ | _G_ |
| Alanine (A) GCG (1.8) HP | Alanine (A) GCA (1.8) HP | Alanine (A) GCC (1.8) HP | Alanine (A) GCT (1.8) HP | Proline (P) CCT (-1.6) HP | Proline (P) CCC (-1.6) HP | Proline (P) CCA (-1.6) HP | Proline (P) CCG (-1.6) HP | _C_ |
| Valine (V) GTG (4.2) HP | Valine (V) GTA (4.2) HP | Valine (V) GTC (4.2) HP | Valine (V) GTT (4.2) HP | Leucine (L) CTT (3.8) HP | Leucine (L) CTC (3.8) HP | Leucine (L) CTA (3.8) HP | Leucine (L) CTG (3.8) HP | _T_ |
| Methionine (M) ATG (1.9) HP | Isoleucine (I) ATA (4.5) HP | Isoleucine (I) ATC (4.5) HP | Isoleucine (I) ATT (4.5) HP | Phenylalanine (F) TTT (2.8) HP | Phenylalanine (F) TTC (2.8) HP | Leucine (L) TTA (3.8) HP | Leucine (L) TTG (3.8) HP | _T_ |
| Threonine (T) ACG (-0.7) Pol | Threonine (T) ACA (-0.7) Pol | Threonine (T) ACC (-0.7) Pol | Threonine (T) ACT (-0.7) Pol | Serine (S) TCT (-0.8) Pol | Serine (S) TCC (-0.8) Pol | Serine (S) TCA (-0.8) Pol | Serine (S) TCG (-0.8) Pol | _C_ |
| Arginine (R) AGG (-4.5) Pol+ | Arginine (R) AGA (-4.5) Pol+ | Serine (S) AGC (-0.8) Pol | Serine (S) AGT (-0.8) Pol | Cysteine (C) TGT (-3.9) Pol | Cysteine (C) TGC (-3.9) Pol | Stop (*) TGA | Tryptophan (W) TGG (-0.9) HP | _G_ |
| Lysine (K) AAG (-3.9) Pol+ | Lysine (K) AAA (-3.9) Pol+ | Asparagine (N) AAC (-3.5) Pol | Asparagine (N) AAT (-3.5) Pol | Tyrosine (Y) TAT (-1.3) Pol | Tyrosine (Y) TAC (-1.3) Pol | Stop (*) TAA | Stop (*) TAG | _A_ |
| __G | __A | __C | __T | __T | __C | __A | __G | |

**Legend:**

- Non-polar Hydrophobic (pink)
- Polar Uncharged (blue)
- Polar Charged + (cyan)
- Polar Charged - (green)
- Stop Codons (yellow)

Periodic Genetic Code table organised by Hamming distance

Each block's codon differs only by one base from any adjacent block
The single base changes also wrap from right edge to left edge and bottom edge to top edge like a 3D torus.
Numbers in the corners and sides of blocks represent the frequency with which one amino acid is found substituted
  in nature by the adjacent one for the corresponding protein from different organisms.
The centre value in brackets shows the hydrophobicity of the amino acid.
This table serves to illustrate the robustness of the genetic code against single and double nucleotide mutations.
  Single nucleotide mutations (horizontally or vertically adjacent blocks) or
  Double nucleotide mutations (diagonally adjacent blocks)
    seem to either result in the same amino acid (silent mutations) or amino acids
    with similar properties (conservative missense mutations) or similar molecular structure.

*The Visual Genome Browser's **Genetic Code** tab, contains this Genetic Codon Table.*

## Show RefSeq Info at UCSC

Choosing this option for a gene will open the online summary information for the specific gene. It essentially represents the detail page for the gene, after you have clicked on the gene region in the UCSC genome display.

For the gene RAD51 it provides you with a hub of links leading to other information, such as OMIM (Online Mendelian Inheritance Of Man – providing online information as to what evilnesses are linked to mutations in this gene)
PubMed, etc.

http://genome.ucsc.edu/cgi-bin/hgc?g=refGene&i=NM_001164270

**RefSeq Gene**

### RefSeq Gene RAD51

**RefSeq:** NM_001164270.1   **Status:** Reviewed
**Description:** Homo sapiens RAD51 recombinase (RAD51), transcript variant 3, mRNA.
**CCDS:** CCDS53932.1
**CDS:** 3' complete
**OMIM:** 179617
**Entrez Gene:** 5888
**PubMed on Gene:** RAD51
**PubMed on Product:** DNA repair protein RAD51 homolog 1 isoform 3
**GeneCards:** RAD51
**AceView:** RAD51
**Related GeneReviews disease(s):** mirror (Congenital Mirror Movements)

**Summary of RAD51**

The protein encoded by this gene is a member of the RAD51 protein family. RAD51 family members are highly similar to bacterial RecA and Saccharomyces cerevisiae Rad51, and are known to be involved in the homologous recombination and repair of DNA. This protein can interact with the ssDNA-binding protein RPA and RAD52, and it is thought to play roles in homologous pairing and strand transfer of DNA. This protein is also found to interact with BRCA1 and BRCA2, which may be important for the cellular response to DNA damage. BRCA2 is shown to regulate both the intracellular localization and DNA-binding ability of this protein. Loss of these controls following BRCA2 inactivation may be a key event leading to genomic instability and tumorigenesis. Multiple transcript variants encoding different isoforms have been found for this gene. [provided by RefSeq, Aug 2009].

## Show in Human Protein Atlas

This option takes you to a the human protein atlas' information on the selected protein coding gene.

I wanted to find a way to quickly determine which body tissues a specific protein coding gene are expressed at the highest levels and found that I could link to the **Human Protein Atlas** site from the context menu.

For, example, when it is selected for the human insulin gene (INS) which is on chromosome 11, you are taken to:

http://www.proteinatlas.org/search/INS

This allows you to look at the expression of this gene's product in different tissue types in the human body.

By clicking on the "**Tissue**" column, you can get the following type of information for the insulin protein:



Illustration 64: Human Protein Atlas showing expression level of gene in different tissues

# Examining gene expression levels

From the excellent online course "Epigenetic Control" presented by Dr. Marnie Blewitt one gets to know that although genes in all body cells, essentially contain exactly the same genetic DNA code, genes in different body cells are expressed at vastly different levels. Some genes are constitutively switched off in some body cells, which other housekeeping genes are always on in most body cells (such as DNA repair genes). The expression level of genes in different body cells are controlled by transcription binding proteins which are sometimes only expressed in some body cells (and which bind to groups of genes which are simultaneously activated) or epigenetic mechanisms such as chromatin packaging applied by histone modifying enzymes and chromatin remodelling proteins. Long non-coding RNA also play a role to guide these modifying enzymes to target regions of the genome.

It is also known that genes are sequestered into regions of the nucleus where there are much higher concentrations of transcriptional machinery, i.e. transcription factories, while other regions such as the nuclear lamina represent areas where genes are mostly switched off.

## *Looking at the transcriptome using the Chromosomes View*

In order to map the expression levels of genes at a genome wide level, I decided to map the gene expression level that one can obtain from the **Gene Expression Barcode**.

They provide publicly available data from http://barcode.luhs.org/

I quote from their website: *"The barcode algorithm is designed to estimate which genes are expressed and which are unexpressed in a given microarray hybridization. The output of our algorithm is a vector of ones and zeros denoting which genes are estimated to be expressed (ones) and unexpressed (zeros). We call this a gene expression barcode. "*

I downloaded the file: abc.ntc.GPL570.csv at

http://barcode.luhs.org/index.php?page=transcriptome **(Provides the expression barcode)**

http://www.affymetrix.com/ **(Provides the mapping between UCSC genes and AffyID)**

This provided me with an Affymetrix Microarray Chip gene expression level **between 0 and 1** for 54613 different gene transcripts for 131 different cell types in body tissues. (**in some cases allowing you to compare normal with tumour tissue expression.**)

The following 24 tumor tissue types are included (based on the Bar Code File provided)

breast_lobular_cells:tumor
breast_stroma:tumor
breast:tumor
cervix:tumor
colon:tumor
endometrium:tumor
fallopian_tube:tumor
glioblastoma:tumor
glioma:tumor
head_and_neck_epithelial_cells:tumor
head_and_neck_squamous_cell_carcinoma:tumor
kidney:tumor
liver:tumor
lung:tumor

omentum:tumor
ovary:tumor
pancreas:tumor
peritoneum:tumor
pilocytic_astrocytoma:tumor
pituitary:tumor
rectum_mucosa:tumor
sigmoid_colon_mucosa:tumor
tongue_squamous_cells:tumor
uterus:tumor

By looking at the transcriptome for different cell types, **the goal was to see if a graphical bird's eye view of the genome**, will not perhaps give insight as to what regions of the chromosomes are situated in transcriptional factories.

***The difference is that, as opposed to an Affimetrix chip, in this case the gene expression is shown directly on the chromosomes, allowing clustering of expressed genes to be observed much better.***

I was expecting to see clusters of genes expressed at higher levels in different regions of the chromosomes. I proceeded to plot the expression level of genes on both the **Chromosomes View** as well as the **Main Genome View**.

Normally, genes are coloured according to their gene names. Genes which the same name will be displayed in the same colour. When "Random color" is not selected, then the biochemical pathway is used (where available) to colour genes the same.

When in this (gene expression) mode, the colours of genes are neither determined by their biochemical pathway or their gene names, but the colours are determined by the expression level on the following scale:



This enables you to clearly see which genes are expressed highly in a specific tissue type.

The software maps between the AffyID (Affimetrix microarray Id) and the UCSC gene Id in order to generate a **colour map** for the expressed genes. **You can also specify a threshold or cut-off value between 0 and 1 to only display genes expressed at a level higher than the threshold value.**

The functionality is loaded by clicking on **Load Expression**. This will load the Affimetrix expression Bar Code values for all the selected cell types. The threshold value will only display genes expressed at a level higher than the specified value. The "**Exclude 1**" option allows you to remove genes expressed at a level of 100%, revealing more of the genes which are differentially expressed (due to chromatin differences or cell specific transcription factors). From here on the threshold setting will also be used as a filter on the **Main Genome View** (allowing you to only see highly expressed genes on the main single chromosome view). You can also choose to switch the colour scale on or off. In order to draw the genes on the **Chromosomes View** (which displays a **karyotype-like view of all the human chromosomes**).



As discussed earlier, you can move the mouse over the **Chromosomes View**, which will then display that region of the chromosomes in the **DNA View**. You can also click on the **Chromosomes View** to jump to that position in the **Main Genome View**.

Chromosomes are drawn to scale and the purple regions represent the centromeric regions of the chromosomes.

Lets look at an example of how this could be useful.  I was searching the literature for gene dysregulation in lung cancer and stumbled up the following article:

**Dysregulation of GIMAP genes in non-small cell lung cancer.**

Shiao YM[1], Chang YH, Liu YM, Li JC, Su JS, Liu KJ, Liu YF, Lin MW, Tsai SF.

⊕ Author information

**Abstract**

The GIMAP (GTPase of the immunity-associated protein) gene family includes seven functional members residing on human chromosome 7. GIMAP genes encode GTP-binding proteins that share a unique primary structure and whose function is largely unknown. However, gene ablation studies reveal that Gimap4 plays an important role in regulating the apoptosis of T cells. In a pilot microarray analysis on six cases of non-small cell lung cancer (NSCLC), we discovered that the expression of GIMAP family members, but not the neighboring non-GIMAP genes, was uniformly lower in the tumor tissues, compared to that in the adjacent nontumor tissues. This finding was subsequently confirmed by quantitative PCR assays in a total of twenty NSCLCs, and we found that GIMAP6 and GIMAP8 showed striking reduction of gene expression in the tumors. In contrast, GIMAP8 mRNA level was abnormally elevated in the adjacent nontumor tissues as compared to that in the control lung tissues. Such reciprocal expression of GIMAPs suggests that this unique gene family might contribute to the pathogenesis of and immune reactions to NSCLC.

I then decided to compare the genome wide expression of **normal lung tissue** with that of **lung tissue with cancer:**

The following picture displays the **high** expression of most genes which are expressed in normal lung cells: (The threshold was set at 0.05 which means that most genes expressed where at high levels, except the genes which are not expressed in lung tissue due to epigenetic control)



When I then loaded the same for **lung cancer cells**, it looks as follows:

One can now visually observe the extent of gene dysregulation in lung cancer tissue. Notice how many genes now have lower expression. In order to better compare, the **F12** key will toggle between the current and the last expression setting. (Or by clicking on the **Toggle** button).

***This toggle function will help you to locate region of the genome where there are gene expression differences.***

In the **Epigenetic Control** course (from the University of Melbourne) mentioned earlier I learnt that, in cancer, there is a genome wide increase in DNA methylation at CpG Islands (often situated at gene promoters), which has the effect of turning gene expression off. (This might well be what is happening here).

To further investigate this, one can not type in the GIMAP gene into the Gene Lookup field:



Pressing the DOWN ARROW KEY and scrolling though the GIMAP genes will display each gene's information in the panel on the right. Press ENTER to select and navigate to the gene.

The colours in the **Main Genome View** will now show the genes as purple with the **Normal Lung gene expression filter.**

*Illustration 65: Normal Lung tissue Gene expression display of GIMAP and surrounding genes at chr7:150,083,281-151,683,811*

When you now press **F12**, the display will toggle back to the **Lung Cancer gene expression filter:**



*Illustration 66:  Lung Cancer tissue Gene expression display of GIMAP and surrounding genes at chr7:150,083,281-151,683,811*

If you want to get rid of any non-GIMAP genes which are still in the display, you can add the GIMAP name in the **Filter Field:**







*Illustration 67: Lung Cancer tissue Gene expression display of GIMAP only genes at chr7:150,083,281-151,683,811*

When you hover the mouse over one of these genes (or when you centre the **DNA View** on one of these genes, you will now see the expression level of the gene in the **Information Display** in the centre panel:



```
(7889)--> GIMAP7
Homo sapiens GTPase, IMAP family member 7 (GIMAP7), mRNA.
(from RefSeq NM_153236) SUMMARY:This gene encod...
chr7:150514830-150521073  (Bases=6,244 Strand=POS+)
(uc003whk.4)
Transcription time between:00:00:10 - 00:00:37 at 35000 and 10000
bases/min
CODING BASES=903 Protein:NM_153236 (L=301) Exons:2
 ID:GIMA7_HUMAN
 PDB:3ZJC
Gene Expression in lung:tumor = 0.688 Disp>=0.05, Yes
```

Another context menu feature is to now go to the **Human Protein Atlas** for one of these genes and then select the **Cancer Tissue** option:



From the picture above one can again see how this gene's expression is lower in lung cancer tissue.

The **Chromosomes View** also allows you to export the "Karyotype" to a **png image file**.

# Visualising DNA Methylation from BigWig and WigBed files

## *Getting DNA methylation data from publicly available sources*

The UCSC Genome browser provides the ability to link external data tracks into the browser's display. These tracks include **DNA Methylation** data, which is obtained using bisulfite sequencing.

This can be obtained via the page:

http://genome.ucsc.edu/cgi-bin/hgHubConnect

| | Public Hubs | My Hubs | | |
|---|---|---|---|---|

Enter search terms to find in public track hub description pages:

| | Search Public Hubs |
|---|---|

*Clicking Connect redirects to the gateway page of the selected hub's default assembly.*

| Display | Hub Name | Description | Assemblies |
|---|---|---|---|
| Connect | Roadmap Epigenomics Data Complete Collection at Wash U VizHub | Roadmap Epigenomics Human Epigenome Atlas Data Complete Collection, VizHub at Washington University in St. Louis | hg19 |
| Connect | Cancer genome polyA site & usage | An in-depth map of polyadenylation sites in cancer (matched-pair tissues and cell lines) | hg19 |
| Connect | ENCODE Analysis Hub | ENCODE Integrative Analysis Data Hub | hg19 |
| Connect | miRcode microRNA sites | Predicted microRNA target sites in GENCODE transcripts | hg19 |
| Connect | DNA Methylation | Hundreds of analyzed methylomes from bisulfite sequencing data | [+] hg38, hg19, hg18, mm9, mm10, panTro2... |

When you browse to a specific position (which can can obtain by the context menu – **copy position to clipboard**), you can then go to the position in the UCSC genome.

Tracks with lots of items will automatically be displayed in more compact modes.

| − | DNA Methylation | disconnect | refresh |
|---|---|---|---|

| [Pub] Akalin 2012 | [Pub] Ball 2010 | [Pub] Banovich-2014 | [Pub] Berman 2012 | [Pub] Blattler 2014 | [Pub] ENCODE 2011 |
|---|---|---|---|---|---|
| hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ |
| [Pub] Gao 2015 | [Pub] Gertz 2011 | [Pub] Grimmer_2014 | [Pub] Guo-Human-2014 | [Pub] Hammoud-2014 | [Pub] Hansen 2011 |
| hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ |
| [Pub] Heyn 2012 | [Pub] Heyn 2012 | [Pub] Hodges 2011 | [Pub] Hon 2012 | [Pub] Huang 2014 | [Pub] Komori-Human-2015 |
| hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ |
| [Pub] Kozlenkov 2014 | [Pub] Laurent 2010 | [Pub] Li 2010 | [Pub] Liao-Human-2015 | [Pub] Lister 2009 | [Pub] Lister 2011 |
| hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ |
| | | | Changes in Human Methylome during Differentiation, Laurent 2010 | | |
| [Pub] Lister 2013 | [Pub] Liu 2014 | [Pub] Lowe 2013 | [Pub] Lu 2014 | [Pub] Lund 2014 | [Pub] Ma 2014 |
| full ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ |
| [Pub] Martins 2012 | [Pub] Pacis_2015 | [Pub] Pei 2012 | [Pub] Roadmap 2015 | [Pub] Schlesinger 2013 | [Pub] Schroeder 2010 |
| hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ | hide ▼ |
| [Pub] Schroeder 2013 | [Pub] Takashima-2014 | [Pub] Thompson_2015 | [Pub] Vandiver_2015 | [Pub] Xie 2013 | [Pub] Zaina-2014 |
| hide ▼ | hide ▼ | hide ▼ | hide ▼ | full ▼ | hide ▼ |

Then click on the side bar to get access to the DNA Methylation source settings:



## Description

| Sample | BS rate* | Methylation | Coverage | %CpGs | #HMR | #AMR | #PMD | |
|---|---|---|---|---|---|---|---|---|
| Chimp_HSPC | 0.993 | 0.758 | 5.190 | 0.971 | 40595 | 0 | 2847 | LowCov; Download |
| Human_Neut | 0.000 | 0.000 | 0.000 | 0.000 | 0 | 0 | 0 | Download |
| Chimp_Neut | 0.993 | 0.742 | 7.327 | 0.977 | 49625 | 0 | 3355 | Download |
| Chimp_BCell | 0.993 | 0.728 | 7.458 | 0.981 | 41309 | 0 | 2609 | Download |
| Human_HSPC | 0.000 | 0.000 | 0.000 | 0.000 | 0 | 0 | 0 | Download |
| Human_CD133HSC | 0.992 | 0.793 | 9.262 | 0.960 | 53891 | 0 | 3669 | Download |
| Human_BCell | 0.000 | 0.000 | 0.000 | 0.000 | 0 | 0 | 0 | Download |

You can now download the tracks with DNA Methylation as **BigWig (Containing the methylation valies) and BigBed (containing the annotation names)** files from here.

# Index of /methbase/data/Hodges-Human-2011/Human_CD133HSC/tracks_hg38

| Name | Last modified | Size | Description |
|---|---|---|---|
| Parent Directory | | - | |
| Human_CD133HSC.hmr.bb | 19-Nov-2015 22:59 | 706K | |
| Human_CD133HSC.meth.bw | 19-Nov-2015 22:56 | 240M | |
| Human_CD133HSC.pmd.bb | 19-Nov-2015 22:59 | 94K | |
| Human_CD133HSC.pmr.bb | 19-Nov-2015 22:59 | 373K | |
| Human_CD133HSC.read.bw | 19-Nov-2015 22:59 | 229M | |

Here is a full description of the process to obtain the data:

http://smithlabresearch.org/software/methbase/

## *Displaying the data in the software*

Because of the huge size of these bigwig files, they are not all loaded into memory as that would take too much space. Instead, the BigWig files are queried only for the region of interest using the indexed structure of the file format.

In order to open a BigWig file, first select the correct genome, eg. HG38 which is the default.

Now draw the chromosome of interest by selecting the correct chromosome in the sequence list on the left and then click **DRAW/REDRAW**. This will display the **Main Genome View** for the selected chromosome. Then load the gene annotations by clicking on the **LOAD GENES** button. (This is optional, but will allow you to see the DNA methylation in the context of the genes and gene promoters). Now load the BigWig file as a layer by clicking on the "Browse" button:



Then select the import file type as **BigWig** files:



This file contains the DNA methylome of B-cells.

When we now look at genes more closely, we can examine the DNA methylation patterns across entire genes. It is striking how genes often contain DNA methylation close to their 5'UTR and promoter regions. It is known that special methyl CpG Binding proteins (**such as** MECP2) control the activation and de-activation of genes by binding to DNA methylation close to the promoters of genes.

Here is the **DNA View** of the Major Histocompatibility complex HLA-B gene (encoded on the negative strand)



*Illustration 68: HLA-B (Homo sapiens major histocompatibility complex) chr6:31356167-31356442*

The CpG methylation is shown as colour coded di-nucleotides. It uses the same colour legend as used in the **gene expression** display, where higher DNA methylation is shown as red-purple and lower methylation is shown as green-blue. When one looks at the end of the gene and 3'UTR region, there are much less methylation:

Another gene HLA-G shown with DNA Methylation around its UTR5' region:



*Illustration 69: HLA-G (Homo sapiens major histocompatibility complex) chr6:29826979-29831122*

Then again, we know that other genes such as the gene TP53, which codes for the all important "Guardian of The Genome" P53, always need to be expressed in cells, and we are not expecting any DNA methylation. Here is the **DNA View** of this all important gene, encoded on the negative strand of chromosome 17. (There are almost no DNA methylation)



*Illustration 70: TP53 (Homo sapiens tumor suppressor protein p53), transcript variant 3, chr17:7676521-7676622*

Interestingly, when we look at the WRAP53 gene, JUST NEXT DOOR of TP53, we observe much higher DNA methylation near the 5'UTR of this gene. Again showing that DNA Methylation (which is mitotically heritable via the DNMT1 methyl transferase) is a regulated process which is applied to specific areas of the genome during embryogenesis and gametogenesis in order to switch specific regions on (depending on cell type) and many regions off (such as transposable elements).



*Illustration 71: WRAP53 (Essential component of the telomerase holoenzyme complex) chr17:7686301-7689079*

The WRAP53 gene (which is found just next to TP53) needs to be switched off for most cells to prevent them from producing telomerase which will allow the cells to effectively live forever (due to the ability to indefinitely lengthen their telomeres and thereby circumvent the Hayflick limit of cell division).

WRAP53

*"Essential component of the telomerase holoenzyme complex, a ribonucleoprotein complex essential for the replication of chromosome termini that elongates telomeres in most eukaryotes. In the telomerase holoenzyme complex, it controls telomerase localization to Cajal body."*

It is also possible to load all the DNA methylation onto the **Main Genome View**, by putting a **.all** in the name of the **.bw** file.  (just make sure only gene regions are displayed and not gene labels).

When you move the **DNA View**, a rectangular block is displayed around the region on the **Main Genome View.** When the DNA Methylation is drawn as a layer on top of the GC Content view of the genome, the locations of the **TP53** gene, which is un-methylated and the **WRAP53** gene, which is methylated, is shown in the following image with colour coded DNA methylation:



*Illustration 72: DNA Methylation on Chromosome 17*

In the centre of the magnified view on the right, the methylated region near the WRAP53 gene can be seen.

# Looking at DNA methylation on transposable elements in the human genome

It is also well known that repetitive elements such as transposable elements are "silenced" by DNA methylation to prevent them from jumping to other places in the genome causing genome instability. When I look at some of these elements I find that they are heavily DNA methylated:



*Illustration 73: POGZ (pogo transposable element with ZNF domain) chr1:151440843-151459179*



*Illustration 74: PGBD1 (Homo sapiens piggyBac transposable element derived 1), chr6:28281572-28281918*

*Illustration 75: TIGD2 (tigger transposable element derived 2), chr4:89111500-89111883*

# Ways the software will enable comparison of DNA differences

The Visual Genome Browser software has the following features which will assist in visualising differences in DNA sequences:

- Direct overlay of 2 DNA sequences directly on top of each other (similar to how transparencies of DNA Southern Blots can be overlaid on top of each other in order to spot differences) as depicted in the following picture.



- Using global and linear sequence alignment methods of DNA and amino acid sequences copied via the context menu and other views
- Using the "DNA-probe" method for searching for snippets of DNA throughout the genome sequence or in die displayed **DNA View**.
- Using an brute force search through all protein coding genes and comparing them using the Blosum/Dayhoff substitution matrices.

## *Comparison by overlaying DNA sequences from two different chromosomes or even genomes on top of each other*

Knowing that all of the autosomal chromosomes have homologous regions which have to pair up in order to provide the "molecular pulling force" necessary for the cell cycle to progress into telophase, I was looking for a way to determine how big the para-autosomal regions are which is shared between the X and the Y chromosomes.  I decided that I would overlay the X and Y chromosomes on top of each other: **Everywhere the nucleotide bases matched exactly, I would display the bases in the DNA View in their correct colour, BUT, everywhere there is a mismatch, I would replace the base colour with magenta.** This gave me a simple way to find out up to exactly where the para-autosomal region stretched.  Humans have 2 copies of the genes in this region, one from the X and one from the Y chromosome.

The way this is accomplished in the software is to simply display the X chromosome sequence and then to **double-click** in the text field as indicated in the following picture:



Alternatively you can enter the position of the first chromosome to be compared as: **chrX:1**

Then, you load the second chromosome (which may come from another genome) and simply double click in the adjacent box.



As soon as both fields have a position in them, the comparison will take place automatically in the **DNA View** window:



*Illustration 76: Picture showing the position at exactly chrY:2781480 where the para-autosomal region ends.*

As you move the mouse over the **Main Genome View** on the bottom left, the **DNA View** will update to show where the differences lie. You can slow down the mouse movement by holding the **ALT** key while moving the mouse. Or, you can use the keyboard arrow keys in order to move the mouse in the genome. When you hold the **Shift** key, the keyboard movement will speed up.

Following this method I was able to accurately find the position where the X and Y chromosomes contain exactly the same bases as position **chrY:2 781 480**.

This means that all the DNA bases up to this point is EXACTLY the same up to the very letter.

## Comparing Polio virus strains

Here is an example where I have compared the following 2 strains of the Polio Virus:

Human poliovirus 2 isolate CHN16019c_Sichuan_CHN_2012  complete

Human poliovirus 2 strain Sabine 2 isolate CHN3024_HN_CHN_1999



One can also use this technique to overlap the mitochondrial genomes of different individuals. Here I have compared "Homo sapiens isolate UV1145 mitochondrion" with the HG38 reference mitochondrial genome. By holding the **Ctrl key** while moving the mouse one can actually SLIDE the first sequence over the second one and see how they overlap.



*Illustration 79: With the first sequence at base 1*



*Illustration 77: With the first sequence at base 2*



*Illustration 78: With the first sequence at base 3*

*Say I wanted to determine at which base of its mitochondrial DNA the Horse and the Elephant aligns.*

- First I load the mitochondrial genome of the Horse (educab). I then arbitrarily select a sequence in the centre of the genome (say 4000) and I copy a piece of its DNA by double clicking at the start and end of the section I want to copy: I copy 78 bases and paste it in the **Search Field** as indicated:
  ATAAGCTCACACTGACTAATAATCTGAATCGGATTTGAAATAAATCTACTAGCCAT TATCCCTATCCTAATAAAAAAG

- From now on this becomes my DNA probe.

- I then load the Elelphant's mitochondrial DNA and look in the **DNA View to see if I see any "hybridizations"/matches found. I proceed to increase the number of allowable mismatches until I get a match. eg.23 as indicated**.

- When I now move the mouse across the **DNA View**, I am given a local sequence alignment in the "**Information** tab" as highlighted in the picture. I conclude that there is an alignment between Horse:4000 and Elephant:3969 (In other words: Elephant:1 and Horse:31)

- While having the Elephant sequence open I double click and set Elephant:1
  Using This technique I was able to visually find a region of alignment.




*Illustration 80: With Elelphant:1 aligned with Horse:31*


*Illustration 81: With Elelphant:0 aligned with Horse:31*


*Illustration 82: With Elephan 14 aligned with Horse:31*

## Using the software to look at VCF (Variant Calling Format) files

When you have your genome sequenced, the DNA data is generally not given as a full DNA sequence, but rather as a VCF file representing all the differences between your genome and the reference genome. The raw sequence data consists of FASTQ files containing short string sequences called **reads** from all over your genome. These short **reads** are then aligned with the reference genome resulting in a file called a **BAM** file. This file contains the sequence alignment. This is however not the final output. The data is then fed through an algorithm which produces a file containing a list of all the individual bases where there is a difference between each of your genome's 46 homologous chromosomes (23 which you inherited from your father and 23 inherited from your mother). This means there could possibly be 2 variant records at each base of the reference genome (one from your mother and one from your father). Consecutive changes are however grouped together into a single variant record of the VCF file.

This means that, given your genome's VCF file, it is possible to reconstruct the sequence each of your 46 chromosomes by reading the nucleotide sequence of the reference genome and then replacing the letters with that from the VCF file record. Each chromosome in the reference genome will then produce 2 different versions for each of your homologous pairs. (Homologous means that the 2 chromosomes has mostly corresponding bases but differs approximately every 1000 bases with a different base (single nucleotide polymorphism/SNP) or with a deletion or insertion (INDEL). This means across the entire genome you have about 6 million mostly harmless base changes (3 million from your mother and 3 million from your father), which makes you the unique person which you are.

## Reading compressed VCF files

The software is able to read **compressed VCF files** into the **overlay** display of the **DNA View**. **VCF files** may sometimes contain the genome variants of more than one individual. This means there will be 2 sequences for each of these individuals in the list of layers to choose from.

Each **compressed VCF** will consist of 2 files each:

| | |
|---|---|
| 7z NA12877.vcf.gz | 36,118 KB |
| NA12877.vcf.gz.tbi | 1,521 KB |
| 7z NA12878.vcf.gz | 36,304 KB |
| NA12878.vcf.gz.tbi | 1,525 KB |

*Illustration 83: 2 Different VCF's for 2 different individuals NA12877 and NA12878*

The **.gz.tbi** file represents the **Tabix index file** which contains an index between the genome positions and the blocks of compressed data in the **BGZIP file** with the extension **.gz**

The **BGZIP** file in turn consists of many broken up blocks of VCF records individually compressed using the **gzip** compression algorithm. This allows software to quickly jump to a specific block of VCF data and decompress it into tabular record data.

It is also possible to get a VCF file which contains records for many individuals such as the following one I downloaded from the **1000 Genomes Project**.

ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/

| | |
|---|---|
| 7z ALL.chr22.phase3_shapeit2_mvncall_integrated_v5a.20130502.genotypes.vcf.gz | 209,428 KB |
| ALL.chr22.phase3_shapeit2_mvncall_integrated_v5a.20130502.genotypes.vcf.gz.tbi | 36 KB |

This VCF contains the genome variants for 2504 different individuals (each identified by unique NA_____ codes and consisting of 2 entries for each person representing the chromosomes inherited from father and mother)



*Illustration 84: Subject's chromosome selection in the software)*

At other times, as in the VCF below, the variants of 3 or more individuals are given with a family relationship (sometimes called **trios** because it consist of a father, mother and child's genomes):

ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp/technical/working/20140625_high_coverage_trios_broad

## *Displaying the VCF in the Software*

After selecting the correct release of the reference genome eg. HG38, and then displaying the desired chromosome by clicking "**DRAW/REDRAW**". If you also want to load the genes overlay, click "**LOAD GENES**". Then browse to the gene you want to look at. For example, you can type **PRR35** in the **Gene Entry field** and then quickly jump to a gene at chr16:563256-564376.



*Illustration 85: Homo sapiens proline rich 35 (PRR35), mRNA. (from RefSeq NM_145270)*

Now it is time to load the VCF. Click on the "**Browse**" button and select **.vcf.gz.tbi** file extension.



When displaying VCF data it is always VERY IMPORTANT to make sure that you are displaying the VCF variants using the correct matching reference genome (eg. HG19 or HG38 etc.) which was used to generate the VCF file in the first place. The VCF data contains a snippet of bases from the reference genome which can be used to verify if it is from the correct genome.

Now select the person's genome you want to look at eg. NA12878 (which is the daughter in the trio)

We can now see the variants indicated for the first allele of this specific gene. You can select a different chromosome in the drop down combo box in order to display different alleles (and subjects if there are more than one in the file).



The SNPs (single nucleotide polymorphisms) are indicated in the image below.



*Illustration 86: Variants indicated on gene for Homo sapiens proline rich 35 (PRR35), mRNA.  chr16:563256-564376*

You can quickly switch between alleles by pressing the **< and >** buttons when the mouse if over the **DNA View**.

When you move the mouse over the indicated SNPs, you can also get more information on what kind of variant SNP/INDEL it is.

There are now new options available in the context menu related to the VCF variants which allows you to both copy the changed coding sequence and the changed protein sequence to the clipboard.

There is also an option which will draw the changed protein sequence in the **Protein View**.



*Illustration 87: New variant menu options*

The clipboard options allow you to copy the DNA or Amino acid sequences to the clipboard or the Comparison list, allowing you to do sequence alignment between the sequences.



You can also press the **< and >** buttons to switch between different alleles of this gene while displaying the protein view. This allows you to compare the effect of the different allele's changes on the protein sequence. **Amino acid changes are marked with small magenta rectangles** and when you move the mouse over these changes, the tooltip will display the kind of change such as Conservative/Non-conservative missense or amino acid Insertions. (Deletions are not displayed because they are not available in the changed sequence).

Also, because of the linkage between the views, the **DNA View** will centre on the appropriate codon.

While looking at the protein view you can also shift through the displayed genes by pressing the **Reverse and Forward** buttons:

## When you want to do sequence alignment you can use the Comparison List

Every time you copy a sequence to the clipboard, it also gets put into the **Comparison List**. This list can be found on the **Main Tab**.



All of the following menu options will copy sequences into the comparison view:



Now click on the "**Copy**" button after selecting the desired entries you want to compare:



This will copy the sequences to compare into the fields below it. Now select the correct protein alignment method, eg. Blosum62, check the colors check box and then click on "**Edit Distance**".

This will do a sequence alignment using the selected comparison method and amino acid substitution matrix:

The **BIG / SMALL button** can be used to select between a large and a small display of the alignment letters. The **Colors check box** will colour code the letters based on the same Amino acid polarity scheme that is used elsewhere in the application. Notice how the Blosum62 will often display a colon (:) (partial match for a conservative mis-sense mutation, where the resulting amino acid has similar properties as the original.





*Illustration 88: Global sequence alignment between reference sequence protein and NA12878 variant for gene PRR35*

*Illustration 89: The same comparison can be done on the DNA Coding sequence level*

Sometimes variants causes truncated proteins such as with the RAB40C (UCSC Id = uc059olp.1)

chr16:625899-626083 .

When I look at the NA12878 (2), in other words, the second homologous chromosome of Chr16, I found a truncated protein caused by a deletion of a nucleotide base:



*Illustration 90: HG38 Reference sequence bases : chr16:624937-625233*



*Illustration 91: NA12878 (2) with deletion at chr16:625088. Premature stop codon*

*Illustration 92: Reference sequence gene protein product*



*Illustration 93: Truncated protein caused by single nucleotide deletion (in other words, this is a non-sense mutation causing premature stop codon to be read)*



## Limitations of the software

The Virtual Genome Browser currently do not take into account disruptions in promoters, splice site disruptions or total loss of start or stop codons. It was only intended to give some clues on possible coding sequence disruptions.

# *Finding similar proteins in viruses*

## *Getting the viral genome data*

The software has the ability to look for protein similarity between genes on different sequences or chromosomes from the same or different genomes.

In order to demonstrate how this works, I will show how to search for similarity between different viruses.

The sequence data for the different virus have been obtained by creating text files in folders which starts with "GenBank" for Genbank format files or "Embl" for Ensembl format files.

The software will automatically scan the folders and then add the GenBank files to the per folder **. 2bit** file.

| Name | Size |
|---|---|
| > VennMath  >  Info  >  Genomes  >  GenbankVirusesClass4SingleStrandedSenseRNA | |
| Barley yellow dwarf virus Ker_II isolate K439.txt | 18 KB |
| BovineViralDiarrheaVirus.txt | 24 KB |
| ChikungunyaVirus.txt | 23 KB |
| Citrus tristeza virus.txt | 44 KB |
| Coxsackievirus B1.txt | 15 KB |
| Cricket paralysis virus isolate CrPV_2.txt | 18 KB |
| Cryphonectria hypovirus 1 strain CN280(09280).txt | 23 KB |
| Cucumber mosaic virus  segment RNA1  isolate Palampur ... | 8 KB |
| Cucumber mosaic virus  segment RNA3  isolate Palampur ... | 5 KB |
| Cucumber mosaic virus RdRp  segment RNA2  isolate Pala... | 7 KB |
| Deformed wing virus isolate Chilensis A1.txt | 19 KB |
| Dengue1Virus.txt | 22 KB |
| Drosophila C virus.txt | 18 KB |
| Enterobacteria phage Qbeta.txt | 14 KB |
| Flock house virus isolate TNCL segment RNA1 protein A ... | 7 KB |
| Flock house virus isolate TNCL segment RNA2 protein alp... | 4 KB |
| Foot_and_mouth disease virus _ type O isolate.txt | 17 KB |
| GenbankVirusesClass | 94 KB |
| GenbankVirusesClass | 6 KB |
| HepatitisC.txt | 19 KB |
| MersCoronaVirus.txt | 61 KB |
| NoroVirus.txt | 16 KB |

Type: TXT File
Size: 18.0 KB
Date modified: 2016-09-13 9:38 PM

I systematically went through all viruses in the book "Virus – an illustrated guide to 101 incredible microbes. By Dr Marilyn J Roossinck" and sorted them into sub folders based on the type of genome.

I now have all these viruses at my disposal and it allows me to search for similarities between these viruses.



After selecting a genome/folder, you can download viruses or multiple consecutive sequence data by entering the NCBI accession number in the download field on the **Main Tab.**

**I downloaded the viruses at :** https://www.ncbi.nlm.nih.gov/nuccore





The second field can be used to download genomes of viruses which consists of multiple segments, such as the flu virus.

This will download the GenBank info into a text file an place it into the appropriate folder.

You then simply select the appropriate folder again (**BUT the one prefixed by "Folder:"**) and it will recreate the **.2bit** genome file (containing all of the viruses in the same folder). It will also recreate the search index and quick gene lookup index in order that you can quickly navigate to viral genes.

*Illustration 94: Searching for genes in viruses*



*Illustration 95: Listing all genes in the selected virus*
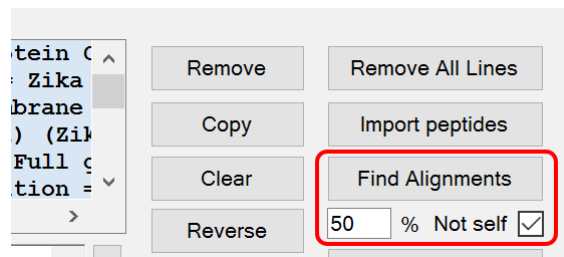
Leaving the **Search Field empty**, and then clicking on **Find Genes** will load all the genes in the **currently selected sequence/virus**.

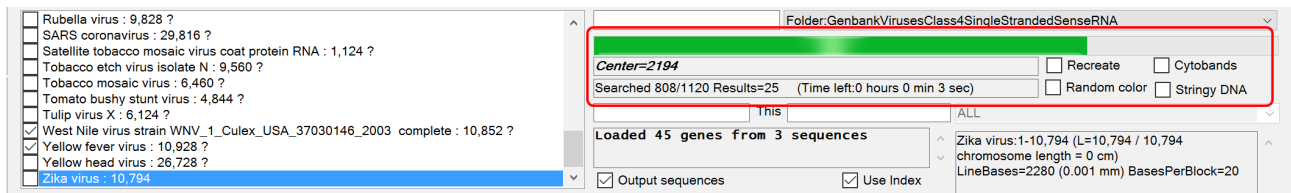## Finding similar proteins between viruses

The first step is to import all the protein coding genes into the **Comparison List**. This is done after the virus/viruses and their genes have already been loaded. Now go to the **Main Tab** and click on the "**Import peptides**" button.
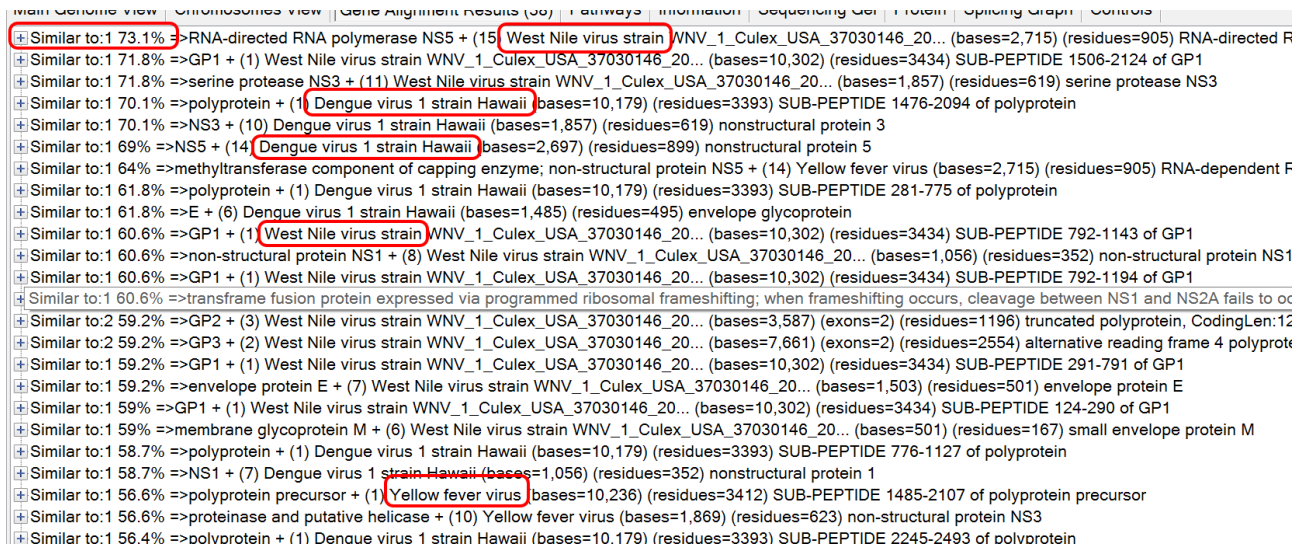


Now proceed to the virus or viruses you want to compare with and load their genes by clicking on "**LOAD GENES**" after selecting the virus sequences you want to compare with:



Then click the "**Find Alignments**" button on the **Main Tab**. The similarity minimum cut-off for proteins can be provided and the "Not self" check box can be used to ignore the virus from being compared against itself (I. e. the Zika Virus)



All the proteins and peptides with similarity is then listed in the Gene Tree:

By double clicking on the gene itself, you will be navigated to the specific gene in the Genome Browser.

By double clicking on the second line for each gene entry, you will be taken to the alignment screen.



Similar to:1 73.1% =>RNA-directed RNA polymerase NS5 + (15) West Nile virus strain WNV_1_Culex_USA_37030146_20... (bases=2,715) (residues=905) RNA-directed RNA polymerase NS5

=>RNA-dependent RNA polymerase NS5 of Zika virus is 73.1% similar to:RNA-directed RNA polymerase NS5 of West Nile virus strain WNV_1_Culex_USA_37030146_2003 complete (DOUBLE

West Nile virus strain WNV_1_Culex_USA_37030146_2003 complete:7648-10362
RNA-directed RNA polymerase NS5
RNA-directed RNA polymerase NS5
(2,715 bases)
name = RNA-directed RNA polymerase NS5
geneSymbol = RNA-directed RNA polymerase NS5
txStart = 7647
txEnd = 10362
cdsStart = 7647
cdsEnd = 10362
strand = +
chrom = West Nile virus strain WNV_1_Culex_USA_37030146_2003 complete
exonCount = 0
parent = ParentPeptide:64-10365:GP1
Similar to:1 71.8% =>GP1 + (1) West Nile virus strain WNV_1_Culex_USA_37030146_20... (bases=10,302) (residues=3434) SUB-PEPTIDE 1506-2124 of GP1
Similar to:1 71.8% =>serine protease NS3 + (11) West Nile virus strain WNV_1_Culex_USA_37030146_20... (bases=1,857) (residues=619) serine protease NS3
Similar to:1 70.1% =>polyprotein + (1) Dengue virus 1 strain Hawaii (bases=10,179) (residues=3393) SUB-PEPTIDE 1476-2094 of polyprotein
=>nonstructural protein NS3 of Zika virus is 70.1% similar to:NS3 of Dengue virus 1 strain Hawaii (DOUBLE CLICK TO VIEW ALIGNMENT)
Dengue virus 1 strain Hawaii:96-10274
SUB-PEPTIDE 1476-2094 of polyprotein
NS3



*Illustration 96: DNA View of RNA directed RNA Polymerase for Yellow Fever virus*

By un-checking the "Codons" the separate genes in the virus genome can more easily be distinguished. When there are multiple genes on top of each other, it is indicated with the line below the DNA View. By pressing + and – you can cycle through the genes in this view.

The alignment view allow you to clearly see the Conservative amino acid substitutions (marked with a colon : as well as the non-conservative subsitutions marked with a space between the similar amino acids which are marked with a vertical line.



*Illustration 97: Protein sequence similarity between Zika and Yellow Fever virus RNA directed RNA polymerase*

All of these alignment results are also output to files in the genomes folder as "Alignments_Date..." in order that you have a record of them.

# Looking at sub-peptides which are obtained due to proteases cleaving larger poly-protein gene products into smaller peptides

In a a group of virus called the **Flaviviridae**, (in fact many other viruses), viral proteins are3 often formed by the mRNA of the virus coding for a single **poly-protein**, which is then cleaved into sub-peptides.

The software is able to detect this in GenBank files and display the protein view appropriately.

When the genes are loaded for the Zika Virus, click on **Find Genes**. This will list all of the genes in the zika virus singl;e stranded RNA genome.  Notice the term: **SUB-PEPTIDES**. This is an indication that there are peptide subdivisions in this virus' genes, which can be appropriately indicated in the protein view.



*Illustration 98: SUB-PEPTIDES available for the Zika Virus*

When you now show the **Protein View** for this specific virus, the different peptide segments which are cleaved are indicated in the following picture.

nonstructural protein NS2B = nonstructural protein NS2B  (1369-1498  L=130)

nonstructural protein NS3 = nonstructural protein NS3   (1499-2115  L=617)

nonstructural protein NS4A = nonstructural protein NS4A   (2116-2242  L=127)

protein 2K = protein 2K  (2243-2265  L=23)

nonstructural protein NS4B = nonstructural protein NS4B   (2266-2516  L=251)

RNA-dependent RNA polymerase NS5 = RNA-dependent RNA polymerase NS5   (2517-3419  L=903)

2715:Glutamine (Gln / Q) Polar Uncharged
CAA = Q or CAG = Q
Sub peptide:2517-3419 L=903 : RNA-dependent RNA polymerase NS5 : RNA-dependent RNA polymerase NS5

1  (C5H8N2O2
Glutamine (Q
[Polar Unchar

2  (C6H12N4O
Arginine (R)
[Polar Charge

3  (C6H12N4O
Arginine (R)
[Polar Charge

4  (C6H7N3O
Histidine (H
[Polar Charge
[Aromatic]

5  (C2H3NO
Glycine (G)
[Special]

6  (C2H3NO
Glycine (G)
[Special]

7  (C2H3NO
Glycine (G)
[Special]

8  (C6H11NO
Leucine (L)
[Non-polar/Hy

Something insightful to do is to examine the amino acids which lie at the transitions between these peptides, in order to determine what the amino acid cleavage motifs are at which the protease enzymes will cut the poly-protein.