UNIVERSITY OF **Southampton**

# Visualisation of genomic data

Dr Jane Gibson
J.Gibson@soton.ac.uk
6th February 2017

---

## Analysis workflow

UNIVERSITY OF **Southampton**

## Lecture Outline

- Why visualise?

- What can you visualise?

- Visualisation tools

- Focus on Integrative Genome Viewer (IGV)

## Learning Objectives

By the end of this session you should…

- Understand the importance of data visualisation and different types of visualisations

- Be aware of several available tools to view alignment/variant data

- Be familiar with the layout of the IGV software, how to upload your data and navigate to regions of interest

- Gain the information you will need in order to complete the practicals

# Why visualise?

## Why visualise the data?

- From initial exploration of the data…
  - "I just want to have a *look* at it"

- …to publication of the data
  - "I need to *illustrate* a point"



'A picture paints a thousand words'

- To check that an algorithm has dealt with the data as expected

- View our own data *in the context* of the information already known

- Just like in the public genome viewers showing the human reference genome eg. ensembl, UCSC

- Integrate annotation of various types, and visualise all relative to a particular genome 'map'

# Lets just look at the files?

.fastq files
(~3.8GB x2)

Sequencing info

Nucleotide sequence

**Files sizes are for exome (50x)- Whole human genome (30x) BAM >10 times bigger**

.SAM file (32GB)
[Binary .BAM file (6GB)]

- The data are too big…
  - You can look at *some* of the file
  - The files can be very wide and very long

- V. difficult to make sense of lists of numbers, letter, symbols and get the important information out –not intuitive

- Want to be able to go straight to data for a specific region of interest

.VCF file (95MB)
[annotated variant file (65MB)]

INFO meta-information

FILTER meta-information

FORMAT meta-information

Optional: FORMAT field specifying data type
+ Per-sample genotype data

---

# Types of biological visualisation

# Genome Browsers

- *"A graphical interface for display of information from a biological database of genomic data"*

- Early days of the Human Genome Project (1999-2001) a need to visualise of genomic data, as a way of making it available to the research community in a more user friendly format than a database

- Public genome viewers showing the human reference genome
  - eg. ensembl, UCSC Genome Browser, (*NCBI mapviewer/Genome data viewer*)

- Many more specific to particular genomes (plants, animals…) or particular projects/datasets 1000 genomes, COSMIC, ENCODE project

- <u>Interactive</u> browsers, move around and investigate data

- <u>Linear</u> representation of genome ~3bn bp, left to right, <u>one strand</u>, lots of tracks below all lined up

- Customisable view; squishable, expandable, add or remove tracks, add your own tracks

*Introduced to ensembl browser in Omic Techniques module (uploaded bam and vcf files to ensembl)*

---

UNIVERSITY OF
Southampton

# What kind of data can we visualise?

# Input files

Epigenomics
Microarrays
NGS alignments
RNA-Seq
mRNA, CNV, Seq

- BAM
- BED
- BedGraph
- bigBed
- bigWig
- Birdsuite Files
- broadPeak
- CBS
- CN
- Custom File Formats
- Cytoband
- FASTA
- GCT
- genePred
- GFF/GTF
- GISTIC
- Goby
- GWAS
- IGV
- LOH
- MAF (Multiple Alignment Format)
- MAF (Mutation Annotation Format)
- Merged BAM File
- MUT
- narrowPeak
- PSL
- RES
- SAM
- Sample Info (Attributes) file
- SEG
- SNP
- TAB
- TDF
- Track Line
- Type Line
- VCF
- WIG
- chrom.sizes

- NGS data, we want to visualise the alignments (how good is the coverage?) and the variants (how good is the variant?)

- Alignments file = Bam (bai), Variants file = vcf (idx)

- Indexing:
  - Often files need to be indexed (IGVtools can do this) to speed up the navigating around
  - Can jump to specific region without reading in whole file from the beginning

| Heatmap | |
| Bar chart | |
| Scatter plot | |
| Line plot | |

---

# Reference genomes

- What's in a name?
- Everything is annotated with respect to a particular reference genome

- Genome builds
  - UCSC = hg18, hg19
  - NCBI = b36
  - Genome Reference Consortium (GRC) = GRCh37
  - Merged = hg38, GRCh38

- Other variations
  - Patches; GRCh37.p1 (regional fixes do not change coordinates)
  - Sorting; ordered lexicographically/ karyotypically
  - Chr1 Vs 1
  - Zero/1-based format for specifying locations

| TP53 isoform a Chr17 | Start | End | length |
|---|---|---|---|
| hg18 | 7,512,445 | 7,531,588 | 19,144 |
| hg19 | 7,571,720 | 7,590,868 | 19,149 |
| hg38 | 7,668,402 | 7,687,550 | 19,149 |



| | 1-based | 0-based |
|---|---|---|
| Indicate a single nucleotide | chr1:4-4 G | chr1:3-4 G |
| Indicate a range of nucleotides | chr1:2-4 ACG | chr1:1-4 ACG |
| Indicate a single nucleotide variant | chr1:5-5 T/A | chr1:4-5 T/A |

- 1-based coordinate system
  - Single nucleotides, variant positions, or ranges are specified directly by their corresponding nucleotide numbers
- 0-based coordinate system
  - Single nucleotides, variant positions, or ranges are specified by the coordinates that flank them

https://www.biostars.org/p/84686/
Tutorial: Cheat Sheet For One-Based Vs Zero-Based Coordinate Systems
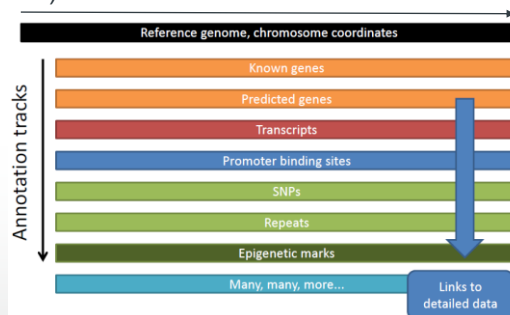Obi Griffith

# Which genome to use?

- Generally use the closest major build, hg18, hg19, hg38
- Check the meta-information/header of your files (sam, vcf)
- Currently in transition between hg19 (2009) and hg38 (2013)
- Archive versions of databases/browsers are available with previous builds
- Conversion tools – eg. UCSC liftOver tool / Assembly Converter @ensembl

When you get it wrong – hg38 bam file viewed on hg19 genome

# Tracks (annotation)

- Sequence – individual nucleotides
- Gene datasets (refGene, Ensembl)
- Databases of known variants
- Sequence conservation
- Repetitive sequence, CpGs
- Phenotype (Omim, COSMIC)
- Many many others…
- Compare samples (trio)
- In context of other peoples data
- In context of other types of data (omics)

Reference genome, chromosome coordinates

Annotation tracks

Known genes
Predicted genes
Transcripts
Promoter binding sites
SNPs
Repeats
Epigenetic marks
Many, many, more…
Links to detailed data

Jon K. Lærdahl, Structural Bioinformatics,
Department of Informatics, University of Oslo

# Data visualisation tools

---

## Tools for visualisation

- **Integrative Genomics Viewer (IGV)**
  - Free
  - Download to PC
  - Windows, Java

Robinson, et al. Integrative Genomics Viewer. Nature Biotechnology 29, 24–26 (2011)
Thorvaldsdóttir,et al Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration. Briefings in Bioinformatics 14, 178-192 (2013).

- **Circos**
  - (http://circos.ca/)
  - Free
  - Download to PC
  - Perl, command-line but possible on windows

Krzywinski, M. *et al*. Circos: an Information Aesthetic for Comparative Genomics. Genome Res (2009) 19:1639-1645

# Tools for visualisation

- **Genome Savant**
  - Free
  - Download to PC
  - Windows, Java
  - Arc view good for visualising structural variants

- **Trackster (Circster) in Galaxy**
  - Free
  - Available via the galaxy server
  - Also available UCSC, ensembl, IGV, IGB



---

# Tools for visualisation

- **UCSC genome browser**
  - Free
  - Web-based
  - Upload your own data track
  - Vast amount of data available through UCSC

- **Ensembl**
  - Free
  - Web-based
  - Upload your own data track
  - Vast amounts of data available through ensembl

## Which tool to use?

- Pros and Cons for each
  - Download loads of annotation from UCSC to view in IGV Or
  - Upload loads of your own files to UCSC?

- Share with colleagues? Available online? Private data?

- How powerful is your PC? Got lots of RAM?

- Easiest to integrate with the external information you want?

- Which do you find easiest?

---

# IGV

# Getting started with IGV

- Launching IGV

- Selecting relevant genome
  - Choose 'more…' if not present in drop down menu

- Loading your data (alignments)
  - Load from file…

- The index file must have the same file name and must reside in the same directory as the bam file
  - mysample.bam (choose this one)
  - mysample.bam.bai



Figure 6. The Integrative Genome Viewer (IGV)

# Navigating

- You will not see your reads when you first load up your data
- Zoom in to your region of interest
  1. Select Chr from drop down menu
  2. Type in location (chr1 or Chr1:100,000-200,000) or locus (gene symbol)
  3. Highlight regions on genomic ruler
  4. Scroll left and right by dragging with the mouse (Home/End, left/right arrows)
- Different views at different resolutions (zooms)

# Genes and sequence

- Sequence
  - Click to get 3 frame translation
  - Click arrow to reverse the strand

- Genes
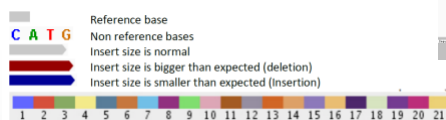  - Multiple transcripts – click to expand
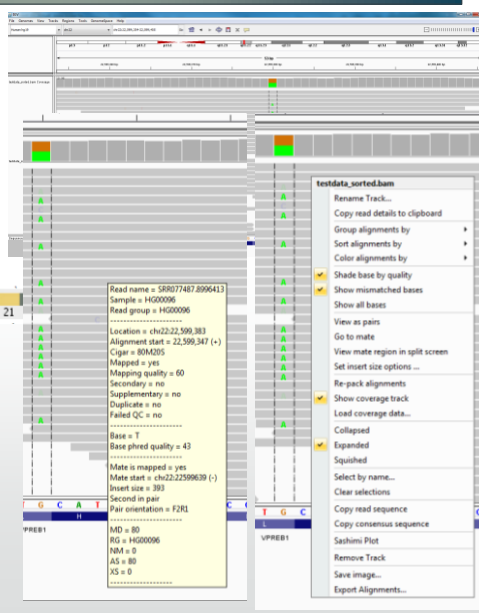  - Positive/negative strand



# Viewing alignments

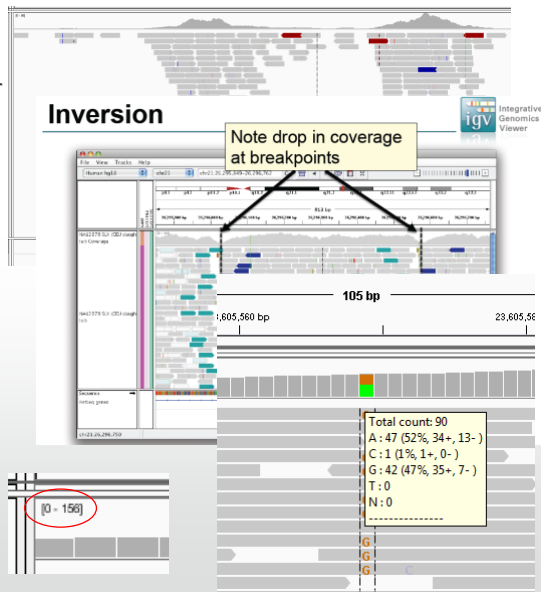- Zoom in to see grey bars (arrow end shows direction)

- Low quality bases are 'faded'

- Many of the thresholds can be altered in preferences

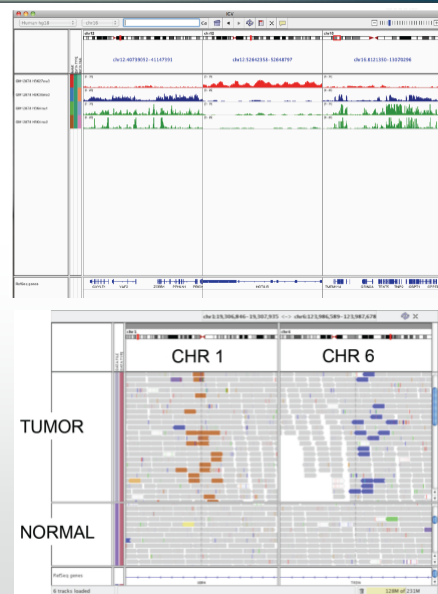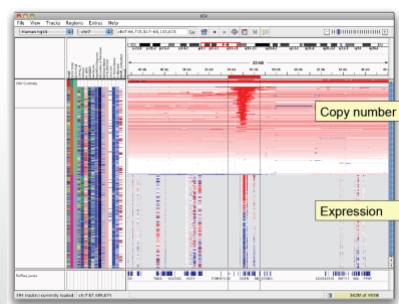- Hover for more information, right click to get more viewing options

# Coverage track

- Histogram across the top

- Coverage profile different for Exome/genome

- Does my gene have good coverage/depth?

- Evidence for depth changes associated with structural changes?

- Coloured to reflect alleles – hover for more information

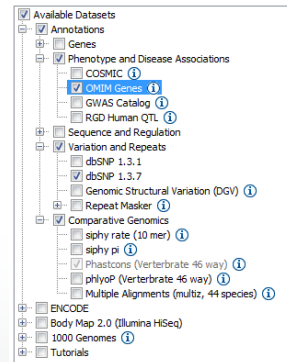- **Note:** scale changes to maximise view



# More complex visualisation

- View sample attributes

- View multiple samples/regions/tissues

# Adding and moving tracks



- Add your own called variants (.vcf)

- Several datasets available (Load from server…)

- Or download your own

- Eg. OMIM genes, or dbSNP 1.3.7

- Can move tracks around to make comparison easier eg. Move dbsnp track to next to vcf track

# Saving your work

- You can save the entire session (.xml), when you restore it, it will be exactly as you left it
  - Saves all your viewing preferences, tracks, location, level of zoom etc.

- Can also save image (.png) of your current view, export visible alignments (.sam), or copy the consensus sequence

- chr22:23915314-23915365
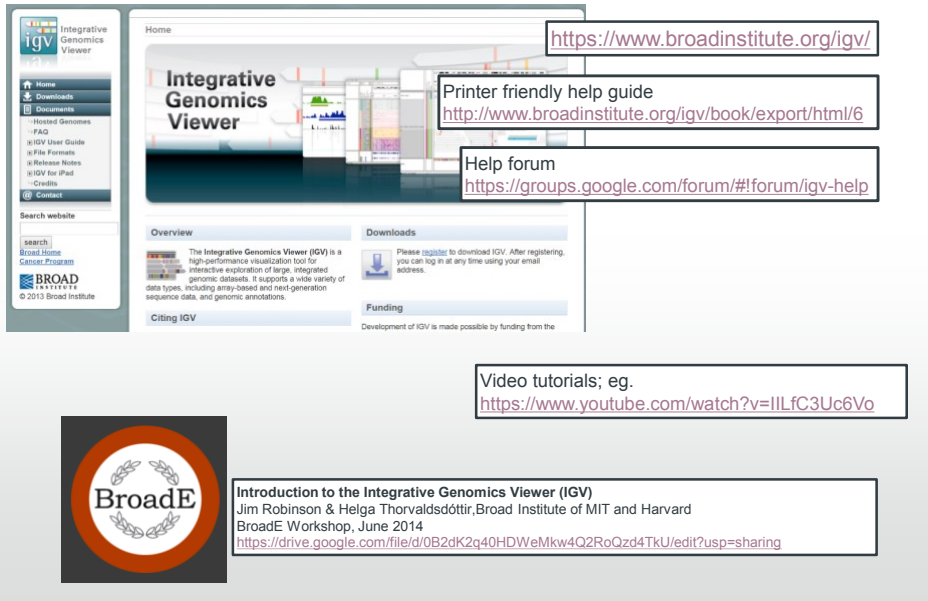  TCTGGTTGACAAAGAGGGTATTTATTKAGGGTTTACTGGGTACAGGGAGAAG

To save a session:

1. Click *File>Save Session*.
2. In the Save Session window, select a directory and session file name and click *OK*.

To restore a saved session:

1. Click *File>Open Session*.
2. In the Open Session window, select a session file and click *OK*. IGV ends the current session and restores the saved session.

## More info/help

https://www.broadinstitute.org/igv/

Printer friendly help guide
http://www.broadinstitute.org/igv/book/export/html/6

Help forum
https://groups.google.com/forum/#!forum/igv-help

Video tutorials; eg.
https://www.youtube.com/watch?v=IILfC3Uc6Vo

**Introduction to the Integrative Genomics Viewer (IGV)**
Jim Robinson & Helga Thorvaldsdóttir, Broad Institute of MIT and Harvard
BroadE Workshop, June 2014
https://drive.google.com/file/d/0B2dK2q40HDWeMkw4Q2RoQzd4TkU/edit?usp=sharing

---

UNIVERSITY OF Southampton

# Next…

## Practical 1: Introduction to Galaxy