

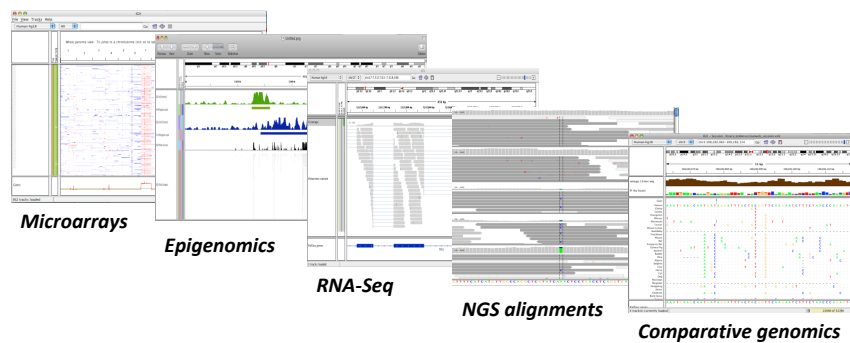
Visualization of Data

Many slides courtesy Broad/MIT.

What is IGV



A desktop/server application
for integrated visualization
of multiple data types and annotations
in the context of the genome



Motivation

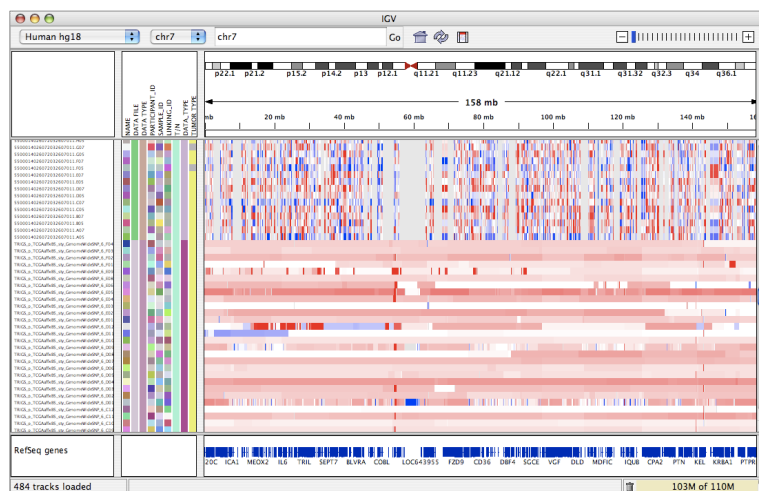


- Easily view investigator-generated datasets alongside publically available data
- Support integration of diverse data types and sample attribute information
- Handle large datasets

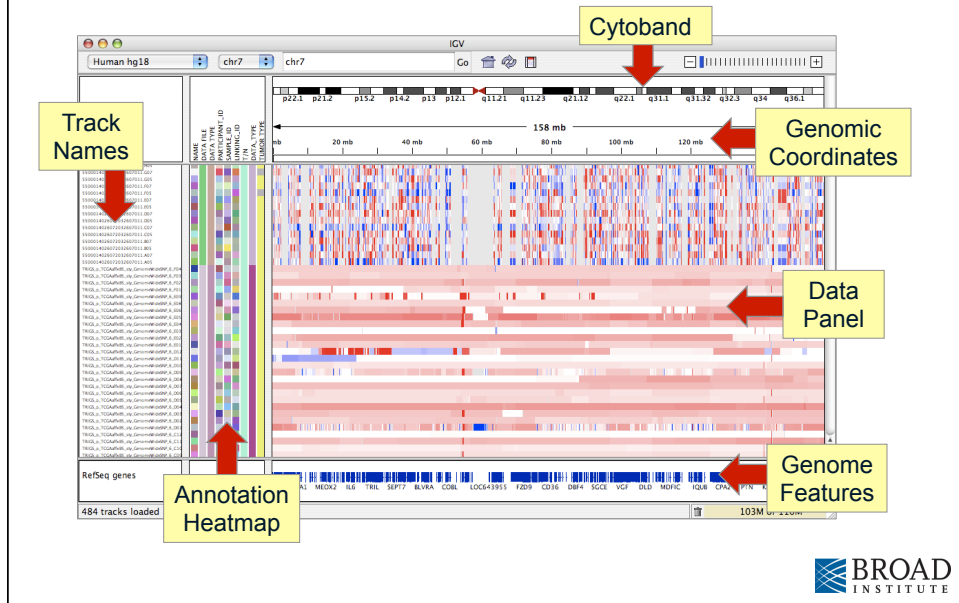


IGV layout

Expression and copy number data

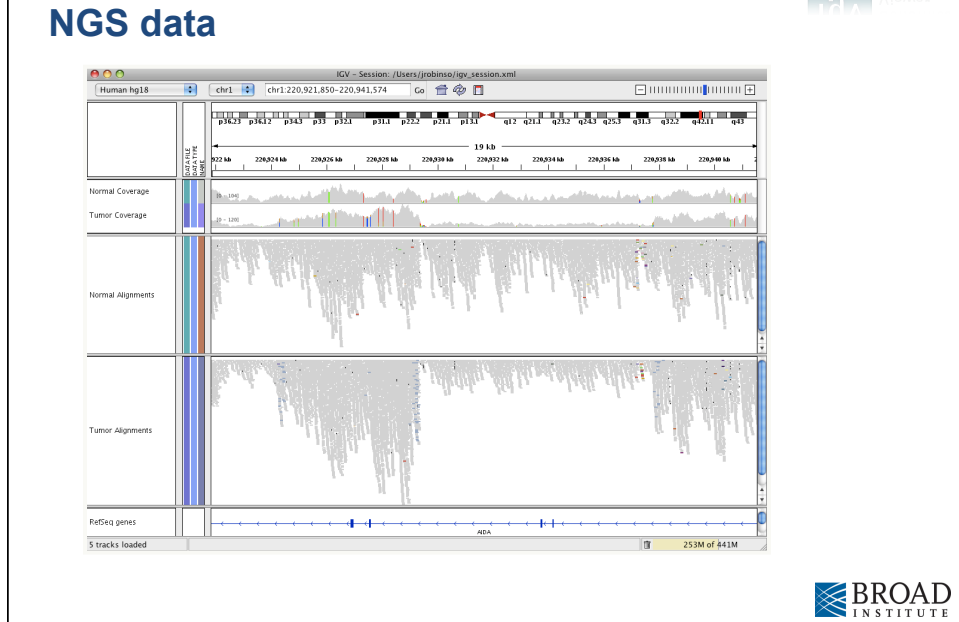


IGV layout



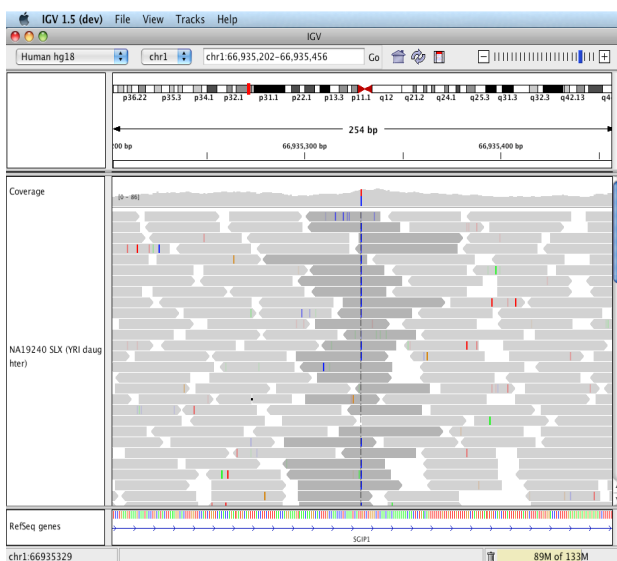
IGV layout

NGS data



IGV layout

NGS data



UI basics



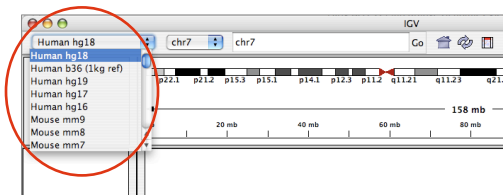
- Selecting a reference genome
- Loading data
- Navigating through the data
- Setting track attributes



Selecting a reference genome



- Select one of the hosted genomes from the pull-down menu



- For more information see www.broadinstitute.org/igv/Genomes
- You can import other genomes if you have the sequence data



Loading data



Types of data

- Any data tied to genomic coordinates
- Genome annotations
- Sample attributes/annotations

File formats

- Many different file formats supported
- See www.broadinstitute.org/igv/FileFormats



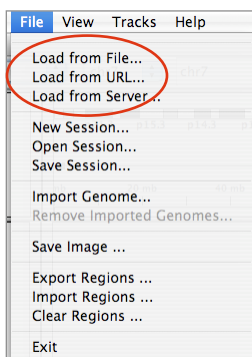
Tracks



- Two generic types:
 - data (continuous valued data)
 - annotation (features)
- Specialized types include
 - alignments
 - mutations
 - multiple alignments
- Type is defined by file format, and can be overridden by the user
- IGV uses type to determine
 - initial placement in a panel
 - display options and options for other track attributes



Loading data



#1 : Load local file

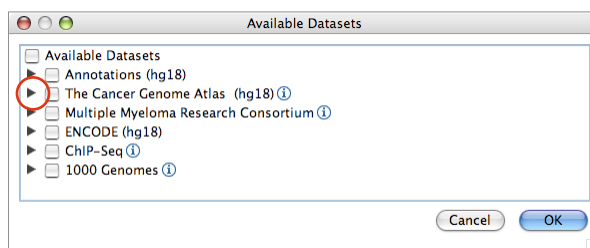
#2 : Load from URL

#3 : Load from server

(Broad IGV data server,
other data server)



“Load from server” menu



What you see depends on :

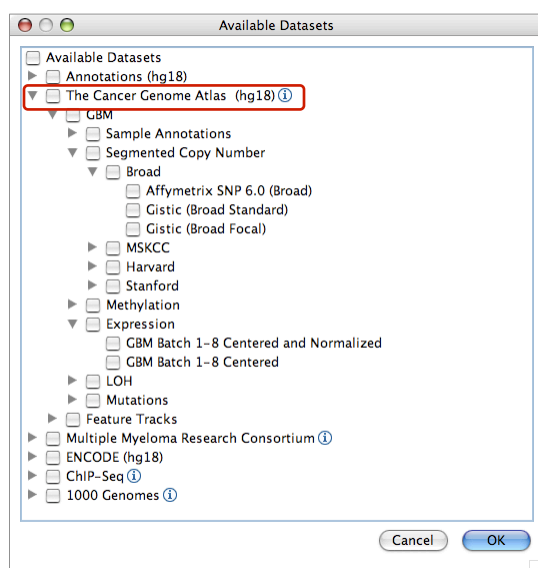
- (1) which server you selected – default is Broad server
- (2) which reference genome you've selected

Click on the  for more information about the data source

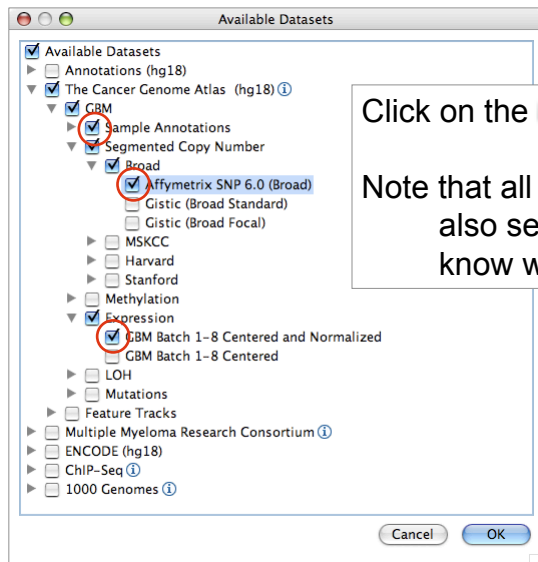
Click on the  to expand the sub-menus



“Load from server” menu



“Load from server” menu

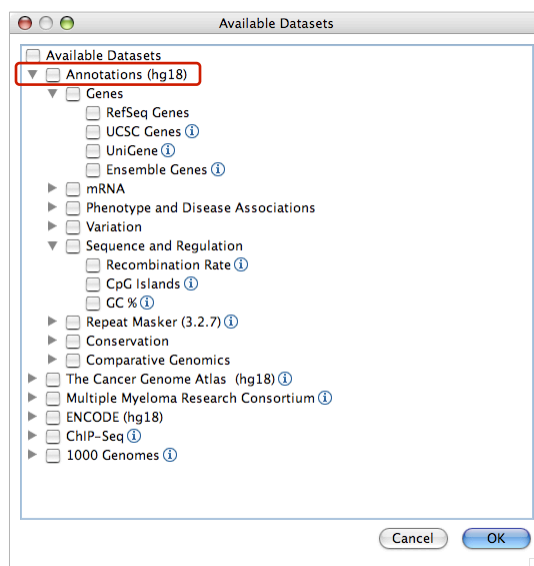


Click on the ☐ to select datasets

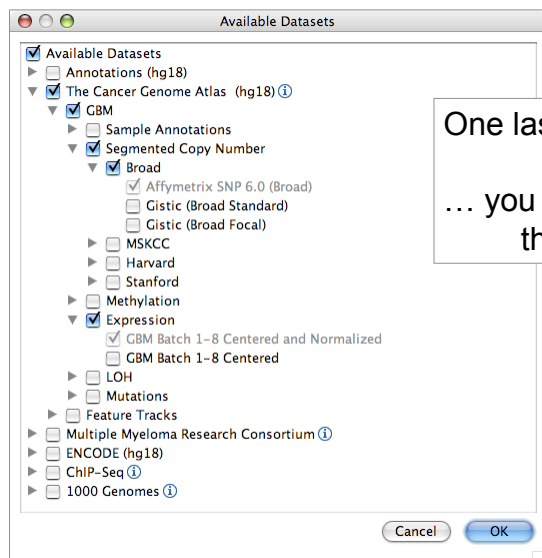
Note that all nested datasets are also selected – make sure you know what you’ve selected



“Load from server” menu



“Load from server” menu



One last thing ...

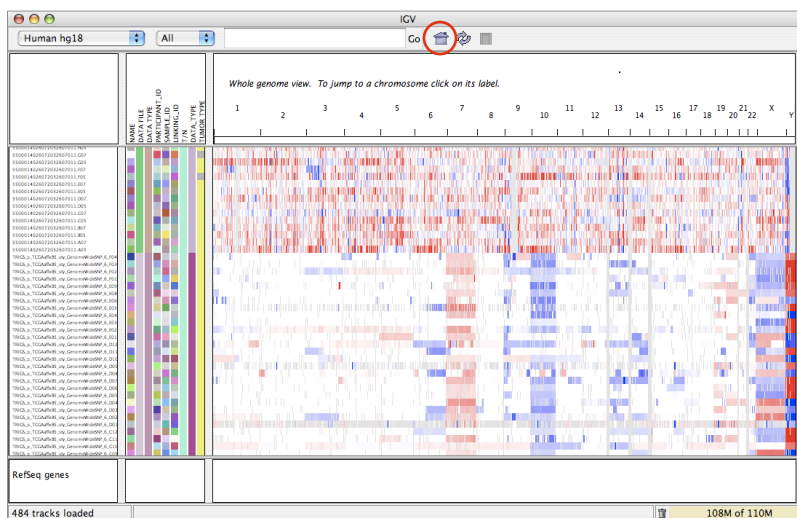
... you cannot **unload** using the checkboxes



Navigating through the data



Whole genome view



Navigating through the data



Zooming in to the chromosome level

Select chromosome from menu

Click on chromosome number

Whole genome view. To jump to a chromosome click on its label.

484 tracks loaded

108M of 110M

BROAD INSTITUTE

Navigating through the data



Chromosome view

Click on chromosome number

484 tracks loaded

103M of 110M

BROAD INSTITUTE

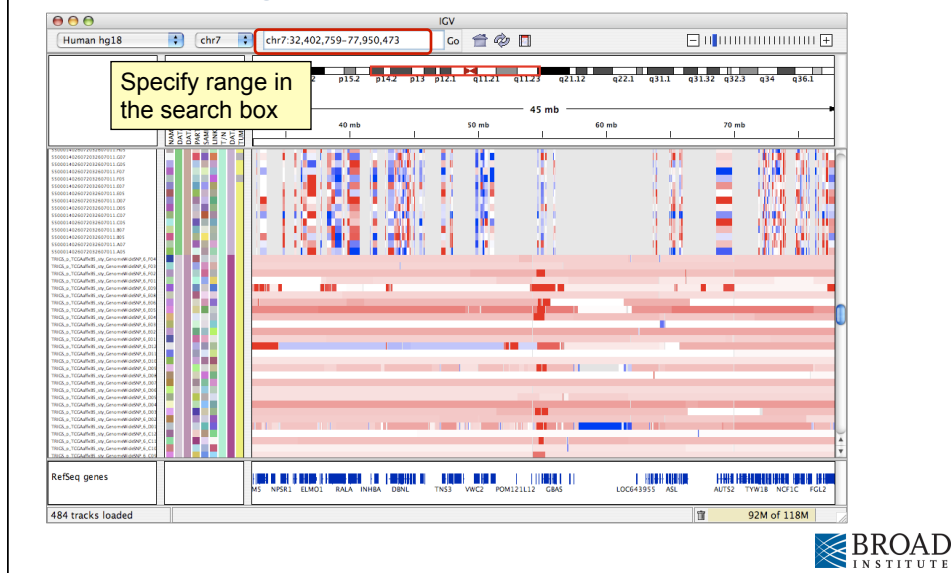
Navigating through the data

Zooming further in



Navigating through the data

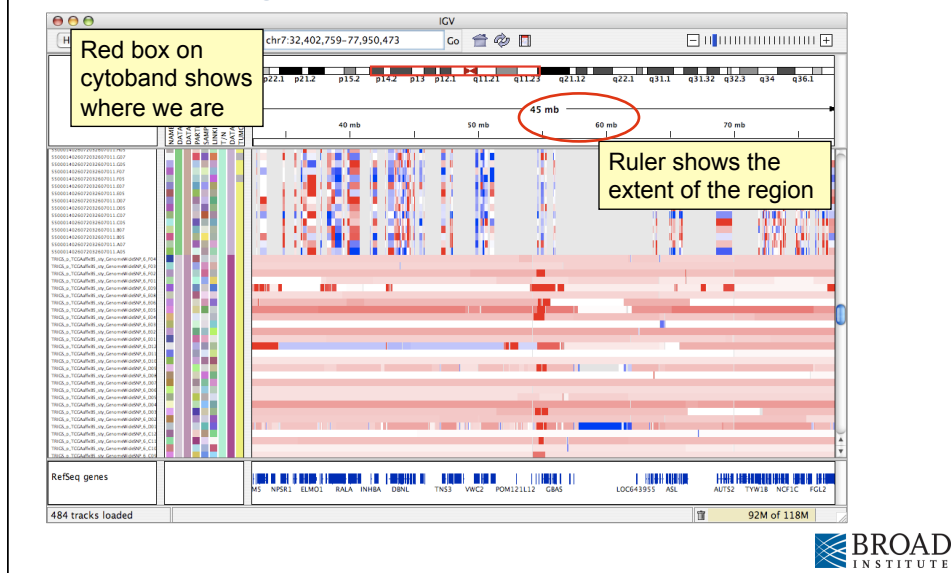
Zooming further in



Navigating through the data



Zooming further in



Navigating through the data



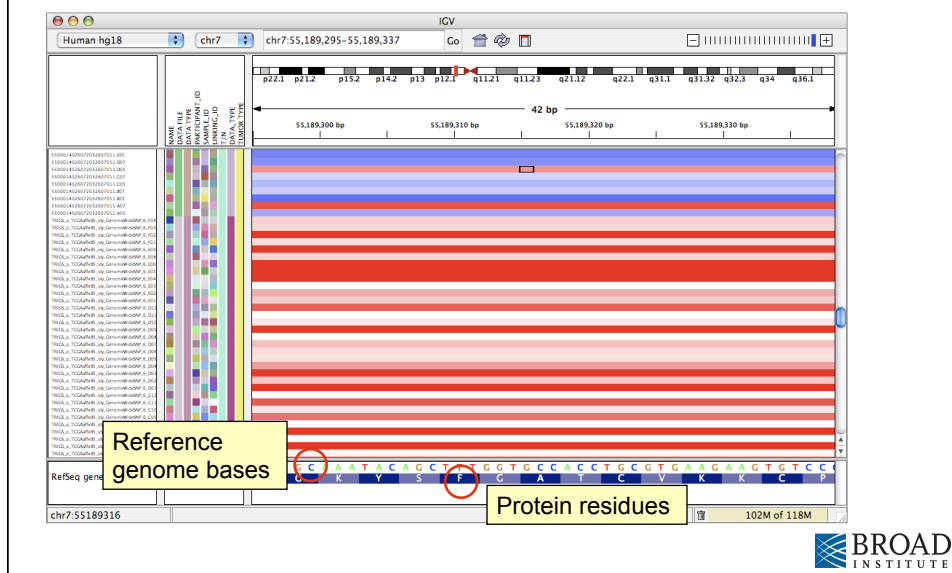
Scroll or jump to location at same zoom level



Navigating through the data



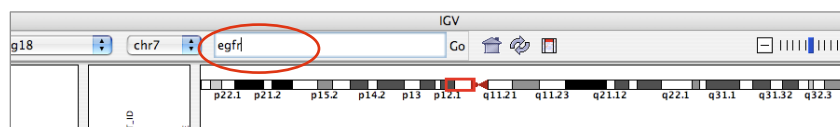
Zoomed in to base pair view



Navigating through the data



Jump to feature



- Enter name of feature in search box
 - With or without zoom (View > Preferences > General)
- Click on a feature track (e.g. gene track, BED, GFF)
 - Ctrl+F = jump forward to next feature
 - Ctrl+B = jump backward to previous feature



Setting track attributes



Right-click popup menu

The screenshot shows the IGV interface with a track selected. A right-click popup menu is displayed over the track, listing various options for modifying the track's appearance and data. The menu includes options for the type of graph, data range, track settings, and track removal.

- Type of Graph
 - Heatmap (selected)
 - Bar Chart
 - Scatterplot
 - Line Plot
- Data Range
 - Set Data Range...
 - Log scale
 - Autoscale
 - Set Heatmap Scale...
- Track Settings
 - Rename Track
 - Change Track Color (Positive Values)
 - Change Track Color (Negative Values)
 - Change Track Height
 - Remove Tracks

The background shows a genomic track with a heatmap visualization of data across a region of chromosome 7 (chr7:19,349,624-64,897,338). The track is labeled "TRIGS_p_TCGAafixB5_sty_GenomeWideSNP_6_E05_223158".

Setting track attributes



Multiple tracks

The screenshot shows the IGV interface with multiple tracks loaded. Two callouts provide instructions on how to select multiple tracks:

- Select multiple tracks by clicking on track names : Shift-click / Ctrl-click** (indicated by a red arrow pointing to the track list on the left).
- Select multiple tracks by clicking on color in annotation heatmap** (indicated by a red arrow pointing to the heatmap visualization).

The background shows a genomic track with a heatmap visualization of data across a region of chromosome 7 (chr7:49,624-64,897,338). The track is labeled "TRIGS_p_TCGAafixB5_sty_GenomeWideSNP_6_E05_223158". The track list on the left shows 757 tracks loaded.

Setting track attributes



Global attributes

Tracks > Fit Data to Window
Tracks > Set Track Height

484 tracks loaded

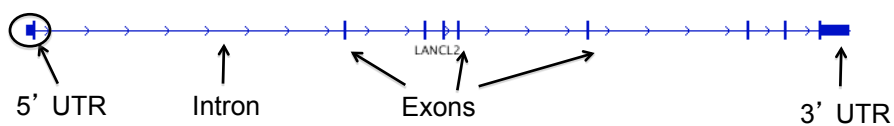
85M of 135M

BROAD INSTITUTE

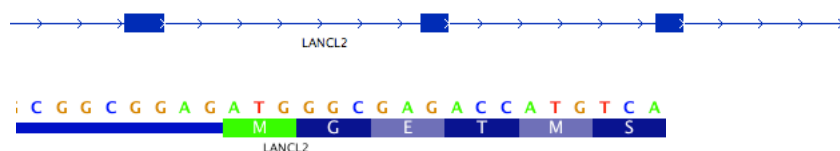
Annotation track



Gene representation



Zoomed in views



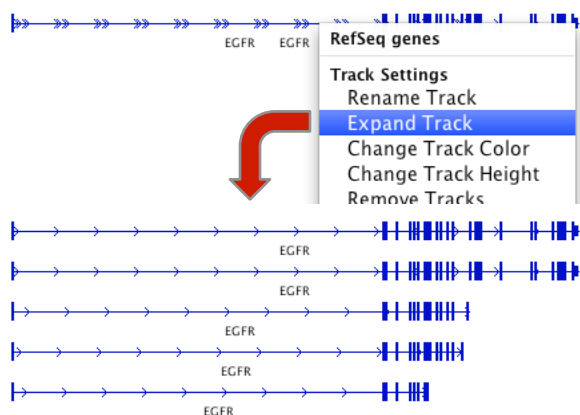
Annotation display mode



1. Features are drawn in a single row, by default



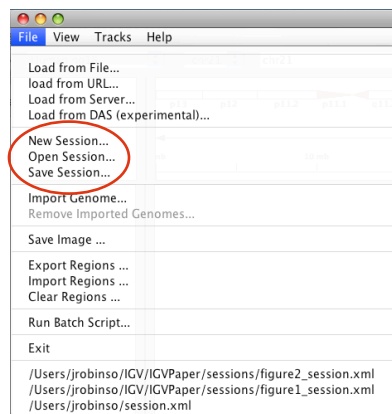
2. Expand the track using the popup menu



Sessions



- Save current state of IGV to a named session file.
- Use to
 - restore the same state
 - share session with colleagues



File Formats



File formats



- Annotation File Formats
- Data File Formats
- Track Line
- Genomes and FASTA Files



Annotation File formats



- BED - UCSC standard format. Useful for displaying any feature type from simple blocks to genes.

<http://genome.ucsc.edu/FAQ/FAQformat.html#format1>

- GFF – Two variants, GFF2 and GFF3. Can also be used all feature types, tends to be more verbose and slower to parse than BED. File sizes can be significantly larger.

<http://www.sequenceontology.org/gff3.shtml>

•Note: BED file coordinates are “zero-based half-open”. This means an interval spanning the first base is represented as 0-1. GFF files are “one-based open”. An interval spanning the first base is represented as 1-1. This difference is responsible for many off-by-one bugs.



Data File formats



•Single Track Formats

- WIG – for fixed or variable step data with fixed spans

<http://genome.ucsc.edu/goldenPath/help/wiggle.html>

- BEDGraph – similar to BED format

<http://genome.ucsc.edu/goldenPath/help/bedgraph.html>



UCSC track line



A track line can be used to control many aspects of the track display such as graph type, color, and scale.

Can be used with wig, bed, gff, igv, cn, and gct files.

Line begins with “track” for wig and bed, “#track” for other formats.

Track line consists of key=value pairs, separated by a single space

Example:

track name=“my custom track” graphType=bar color=255,0,0



Viewing NGS Data



Next-generation sequencing

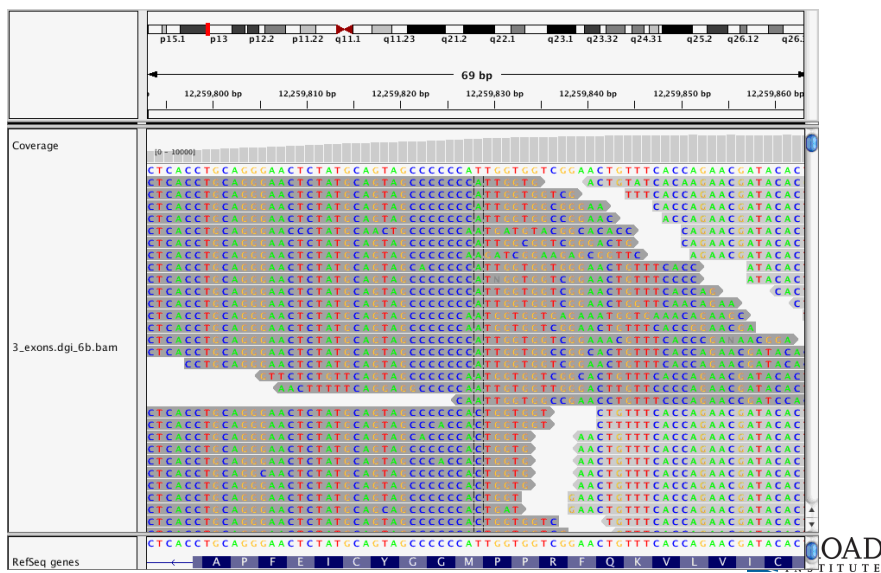


The size of NGS datasets presents many challenges, including:

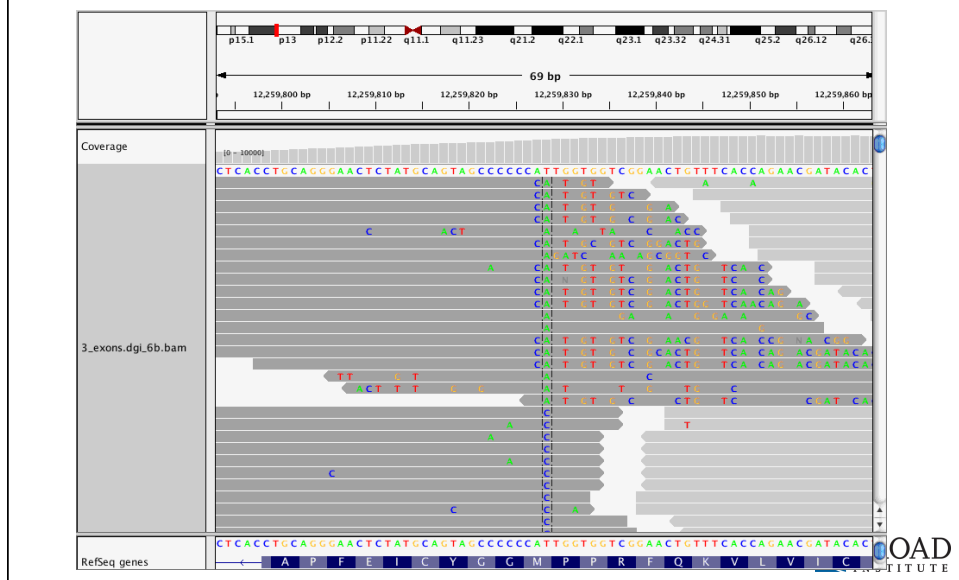
- Implementation
 - Managing terabyte size files with modest compute resources (desktop computers).
- Visual design
 - Highlight events of interest
 - Deemphasize irrelevant details
 - Avoid information overload



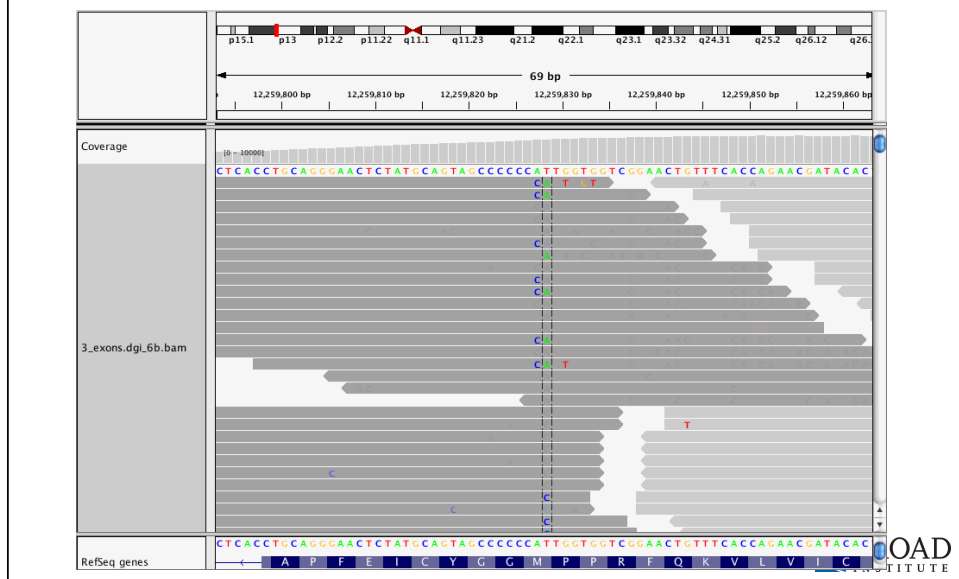
Aligned reads – all bases



Aligned reads - mismatches



Aligned reads – base quality



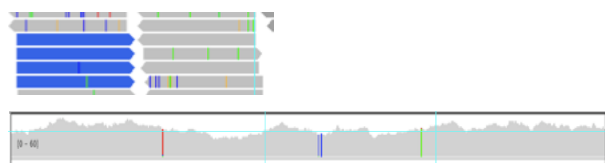
Vary view by resolution scale



Whole chromosome -- calculated summary data, e.g. coverage.



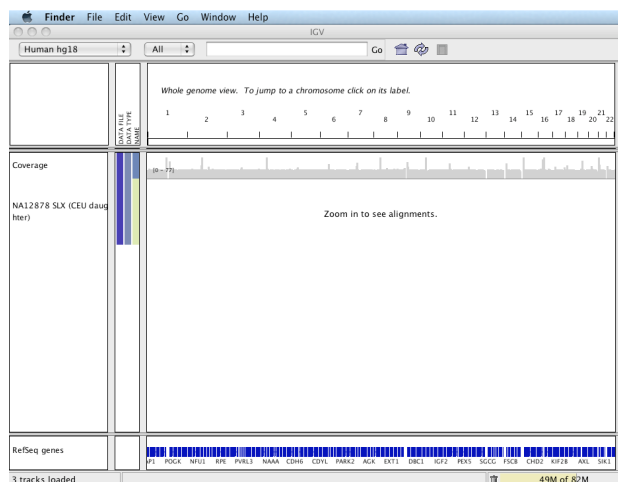
~ 50-100 kb -- putative rearrangements, SNPs



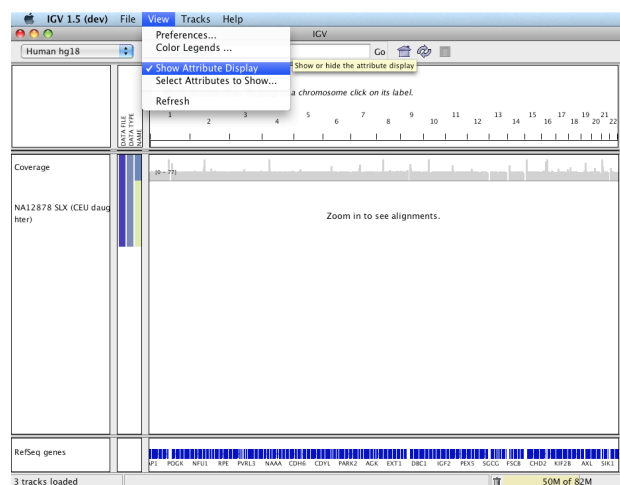
~ 500 bp -- bases



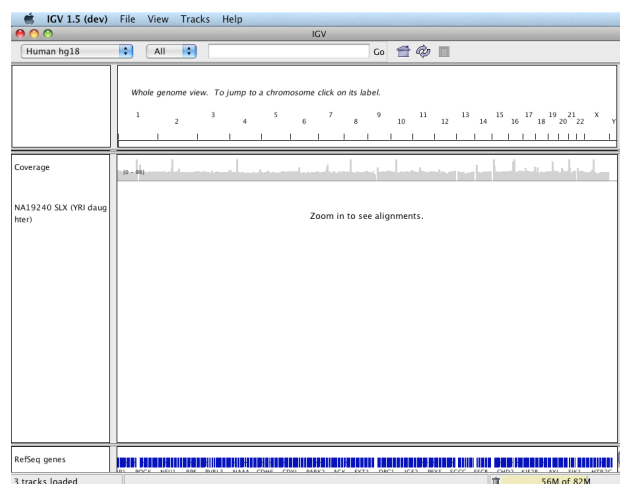
Viewing NGS data



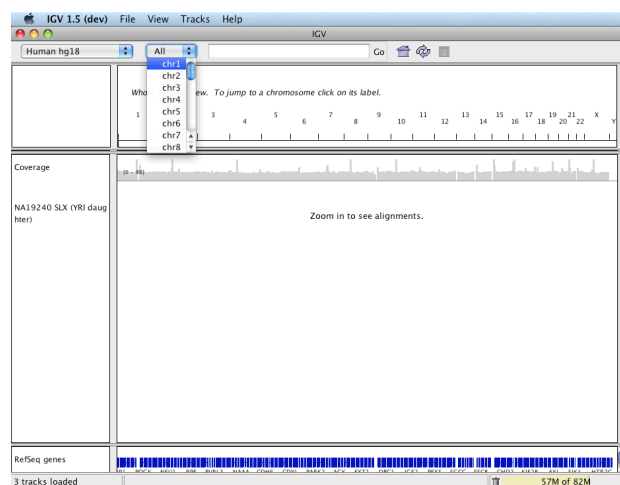
Viewing NGS data



Viewing NGS data



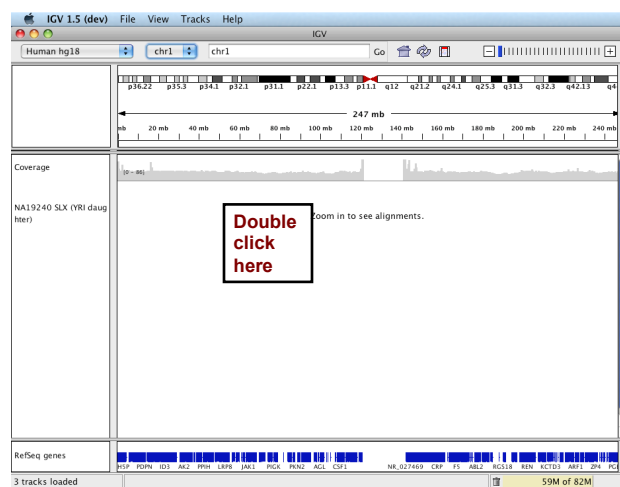
Viewing NGS data



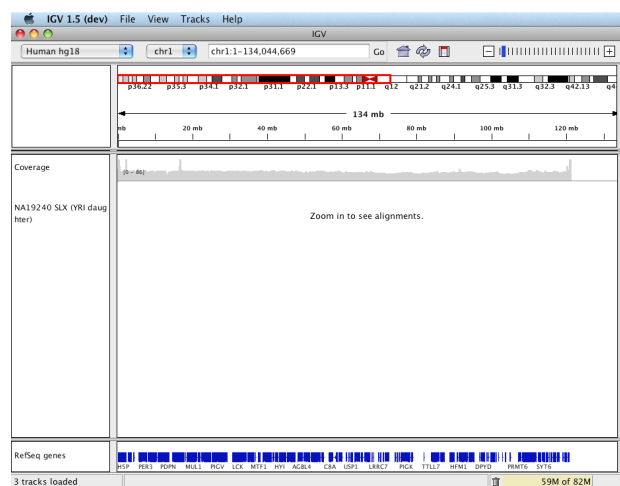
Viewing NGS data



Viewing NGS data



Viewing NGS data



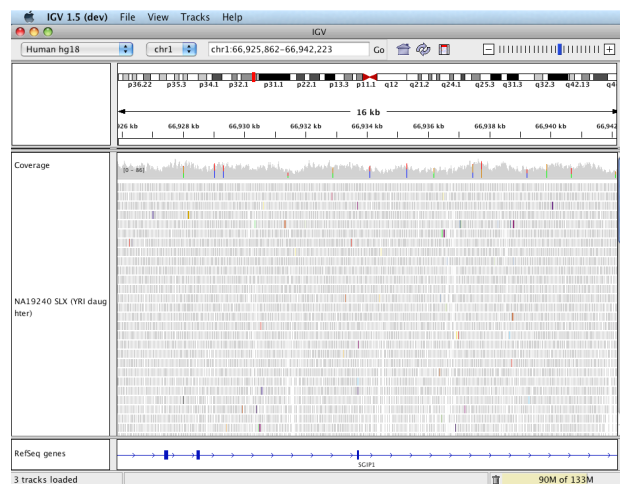
Viewing NGS data



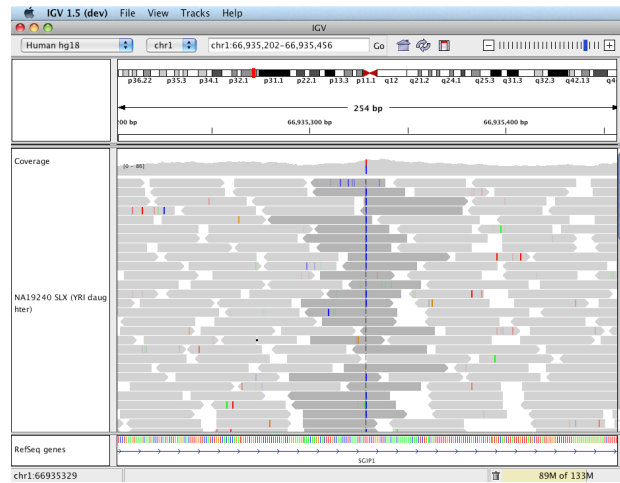
Click
here



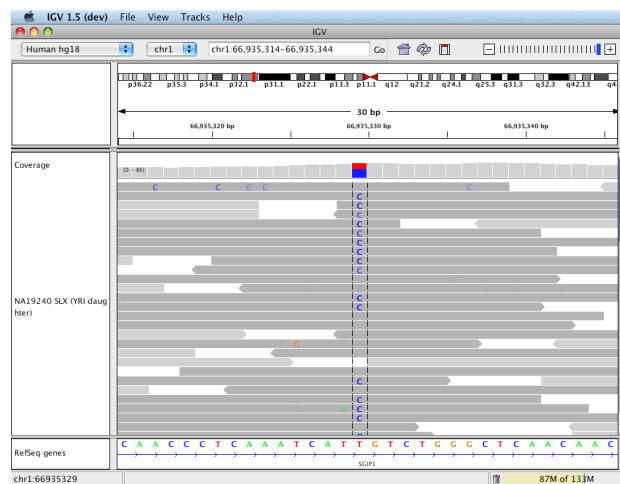
Viewing NGS data



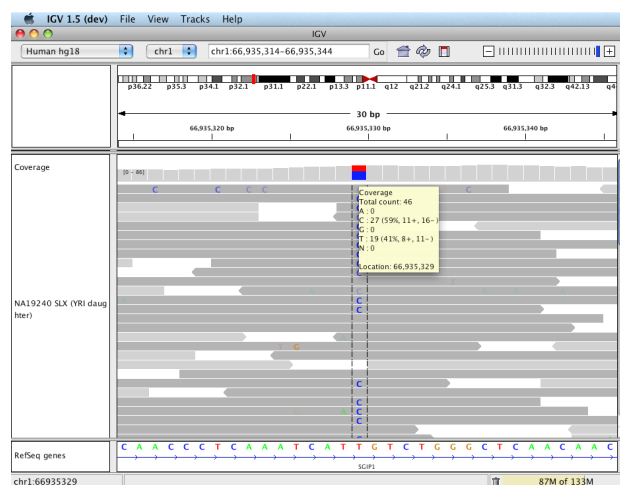
Viewing NGS data



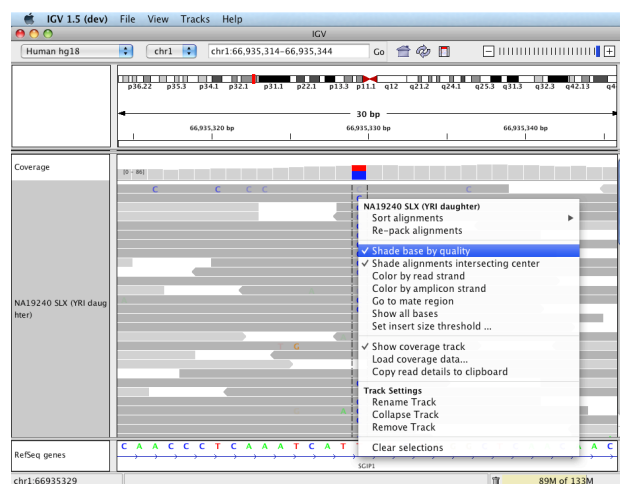
Viewing NGS data



Viewing NGS data



Viewing NGS data



IGVTools



IGVTools



IGVTools is a set of utilities for preparing large files for efficient display.

tile: converts a sorted data input file to a binary tiled data (.tdf) file.
Supported input file formats: .wig, .cn, .snp, .igv, .gct

count: computes average alignment or feature density for over a specified window size across the genome.

Supported input file formats: .sam, .bam, .aligned, .sorted.txt, .bed

sort: sorts the input file by start position. Supported input file formats: .cn, .igv, .sam, .aligned, and .bed.

index: creates an index file for an input ascii alignment file.
Supported input file formats: .sam, .aligned, .sorted.txt



IGVTools tile



The *tile* utility converts large ascii data files into tiled data format (.tdf) files. TDF files have the following advantages

- 1.Data is indexed for efficient retrieval.
- 2.Data for zoomed out views are preprocessed.
- 3.TDF files are web friendly, large data files can be shared over the web. Only small slices of the file are actually transferred as needed.



IGVTools count

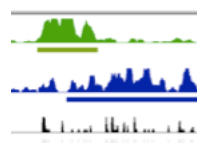


The count command is used to transform alignment files to read density TDF files, e.g. for ChIP-Seq, RNA-Seq, & similar alignment counting experiments.



Alignments

Alignments in bam/sam, .aligned, or bed format.



Read density

"Tiled Data File" indexed and optimized for fast retrieval at multiple resolution scales



IGVTools sort



This utility sorts IGV supported genomic formats by start position.

Example:

```
igvtools sort -m 1000000 -t ~/myTmpDir inputFile.sam outputFile.sorted.sam
```

The sort command uses a combination of memory and disk to handle large files.

-m = maximum # of lines to hold in memory. When this number is exceeded a temporary file is created.

-t = directory used to create temporary files during sorting.



IGVTools index



Used to create an index file for viewing SAM (not BAM) files

Note: to be confused with the *samtools* index, which is used to create an index for BAM files

SAM => igvtools

BAM => samtools

Example: igvtools index inputFile.sam

Result inputFile.sam.sai

The index file must remain with the sam file to be found, IGV just appends .sai to the end.

