



IX International Course of Massive Data Analysis FOR GENOMICS



Overview

Interface

Loading Data

Browsing the Data

Sessions

Exercises

Overview

Interface

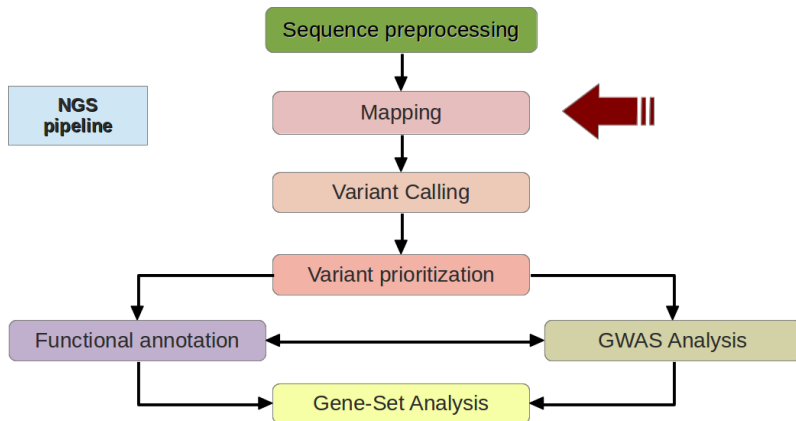
Loading Data

Browsing the Data

Sessions

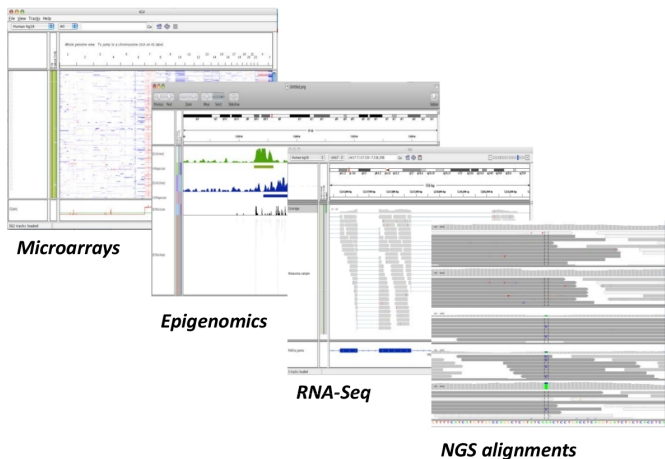
Exercises

Where are we?



Why IGV?

- ▶ IGV is an integrated visualization tool of large data types



Why IGV?

- ▶ Integrate different data types simultaneously
- ▶ View large datasets easily
- ▶ Fast navigation
- ▶ Run it locally on desktop
- ▶ Easy to use interface

Why IGV?

- ▶ Large-Scale projects using IGV

- ▶ The Cancer Genome Atlas

<http://cancergenome.nih.gov/>

- ▶ Multiple Myeloma Research Consortium

<http://themmrc.org/>

- ▶ 1000 Genomes Project

<http://www.1000genomes.org/>

<http://www.broadinstitute.org/igv/>

The screenshot shows the IGV website homepage. On the left is a navigation sidebar with links for Home, Downloads, Documents, Hosted Genomes, FAQ, IGV User Guide, File Formats, Release Notes, Credits, and Contact. Below the sidebar is a search box and the Broad Institute logo. The main content area features a large banner image of the IGV interface with the text "Integrative Genomics Viewer". Below the banner are sections for "What's New" (listing updates from December 2012 and April 2012), "Overview" (describing IGV as a high-performance visualization tool), "Downloads" (with a download icon and instructions), "Citing IGV" (providing citation information for publications), and "Funding" (listing funding sources like the National Cancer Institute and the Starr Cancer Consortium). At the bottom right, there is a logo for GBPA (Genomics & Bioinformatics Platform of Andalusia) and a note that IGV is participating in the GenomeSpace initiative.

History and Usage

- ▶ First release was in August 2008
- ▶ Current version: 2.2
- ▶ Open source and freely available

Overview

Interface

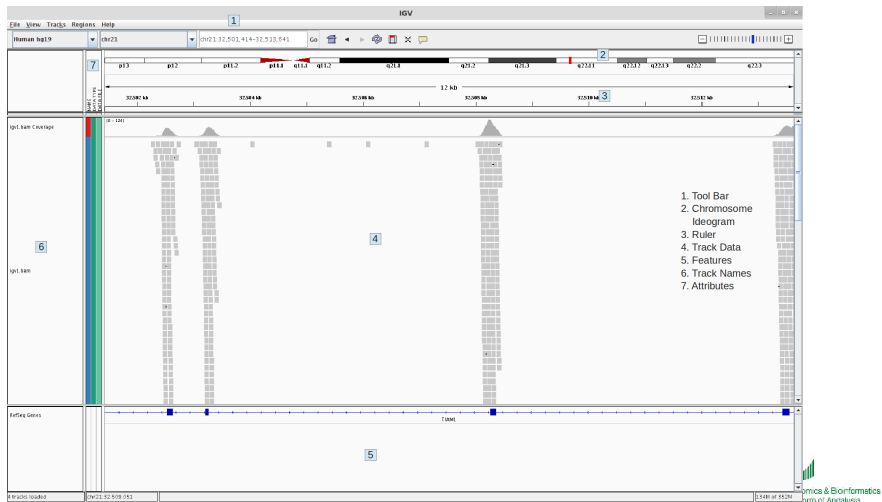
Loading Data

Browsing the Data

Sessions

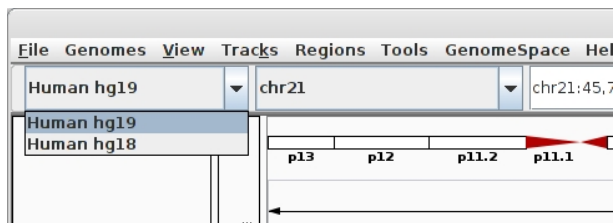
Exercises

Main Window



Genome drop-down box		Loads a genome
Chromosome drop-down box		Zooms to a chromosome
Search box		Displays a genomic location
Whole genome view		Zooms to whole genome view
		Moves backward and forward
Refresh		Refreshes the display
Define a region		Defines a region of interest
		Fits all tracks to the window
		Toggles on/off pop-up information
Zoom slider		Moves to a chromosome

Available Genomes



- ▶ Human, Mouse, *S. cerevisiae*, *C. Elegans*, *D. melanogaster*, and others

<http://www.broadinstitute.org/software/igv/Genomes>

Overview

Interface

Loading Data

Browsing the Data

Sessions

Exercises

- ▶ General characteristics
 - ▶ Any data related to genome coordinates
 - ▶ Sample annotation/attributes
 - ▶ Genome annotations

Data types

- ▶ IGV supports multiple file formats

Source Data	Recommended File Formats
Sequence alignment data	SAM (must be sorted/indexed) BAM format (must be indexed)
Genome annotations	GFF or GFF3 format, BED format
ChIP-Seq, RNA-Seq	TDF format. Use igvtools to generate a binary read count. Load the resulting TDF file into IGV.
Any numeric data	IGV format, TAB format WIG format
Gene expression data	GCT format RES Format

Loading a BAM file

- ▶ BAM format: SAM binary. Reduces disk space and access time.
 - ▶ For each read, provides the position(s) where it maps and information about the alignment.
 - ▶ BAM files need to be indexed (samtools). SAM files need to be sorted by start position and indexed.

Index an example bam file

```
samtools index /home/biouser/mda13/mqc-igv/igv1.bam  
ls -la /home/biouser/mda13/mqc-igv/igv1.*
```

Loading a BAM file

Open IGV

```
igv
```

Load an example bam file

- ▶ File → Load from file →
/home/biouser/mda13/mqc-igv/igv1.bam

Overview

Interface

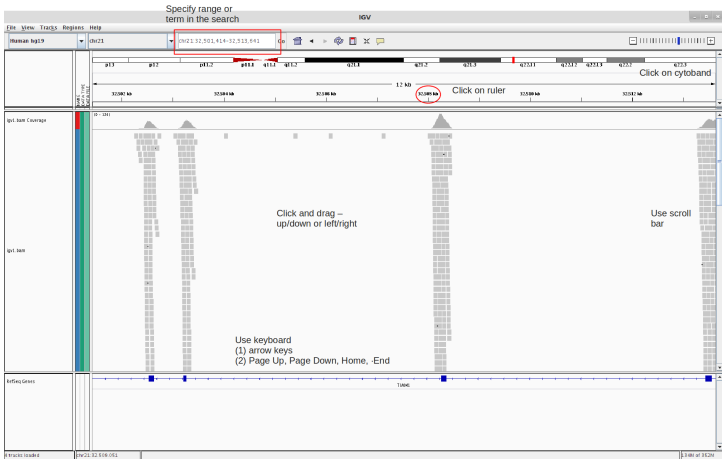
Loading Data

Browsing the Data

Sessions

Exercises

Moving around



Track options

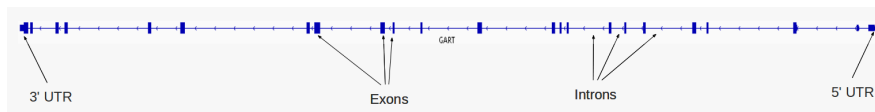
Right click on track

The screenshot displays the IGV interface with a right-click context menu open over the 'igv1.bam' track. The track is currently expanded, showing individual read alignments. Below the track, the 'Sequence' and 'RefSeq Genes' tracks are visible. The context menu includes the following options:

- igv1.bam
- Rename Track...
- Copy read details to clipboard
- Group alignments by ▶
- Sort alignments by ▶
- Color alignments by ▶
- Shade base by quality
- Show mismatched bases
- Show all bases
- View as pairs
- Go to mate
- View mate region in split screen
- Set insert size options ...
- Re-pack alignments
- Show coverage track
- Load coverage data...
- Collapsed
- Squished
- Expanded
- Select by name...
- Clear selections
- Remove Track
- Save image...

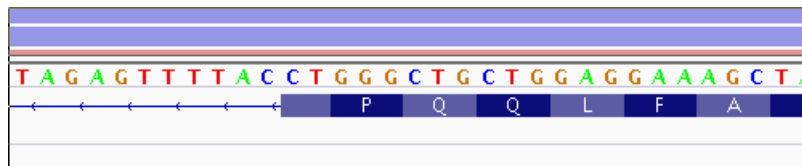
RefSeq track

Type GART in the search box



RefSeq track

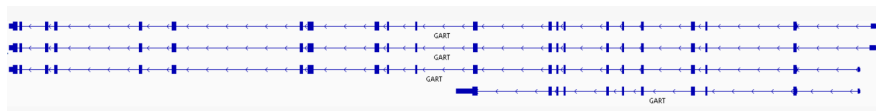
Zoom in to focus on an exon



- ▶ Reference nucleotide sequence
- ▶ Predicted aminoacid sequence

RefSeq track

Right click on the RefSeq track → Expanded



Zoom out to get a general view of the mapping profile

- ▶ Reads are grouped around exons. Is that a coincidence?
No. This bam is part of an exome sequencing experiment.

Loading a BED file

- ▶ Let's load capture regions
- ▶ BED format: to store a list of genomic regions. Text file with the list of regions. Each line contains one region with three required fields separated by tabs: chromosome, start coordinate, end coordinate

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

```
less /home/biouser/mda13/mqc-igv/igv.bed
```

Load a bed file

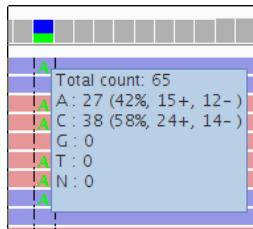
- ▶ File → Load from file → igv.bed
Another track appears: blue boxes indicate target regions

Visualizing variants

Move to 21:48022375

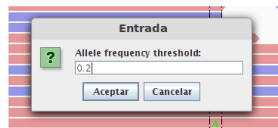
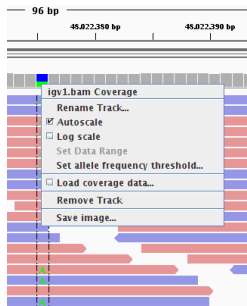
- ▶ Variant: ref:C, alt:A, heterozygosis

Leave the cursor over that position on the histogram (top of the track)



Visualizing variants

Right click over that position on the histogram



Move to 21:47786817

- ▶ Variant: ref:C, alt:G, homozygosis

Visualizing variants

- ▶ Files containing a variant list can also be loaded: VCF files (Variant Calling Format).

```
less /home/biouser/mda13/mqc-igv/igv.vcf
```

- ▶ VCF format:

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

- ▶ Text file
- ▶ Header including information about the mapping and variant calling processes: set of lines beginning with ##
- ▶ An additional header line beginning with #: contains a table header with column identifiers
- ▶ One line for each variant: chromosome, genomic position, reference and alternative bases

Visualizing variants

- ▶ VCF files need to be indexed before being loaded:

```
igvtools index /home/biouser/mda13/mqc-igv/igv.vcf  
ls -la /home/biouser/mda13/mqc-igv/igv.vcf*
```

Load a VCF file

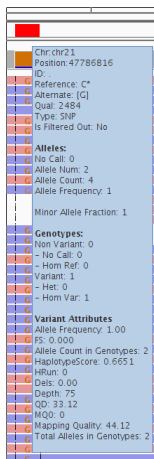
- ▶ File → Load from file → igv.vcf A new track is added

Move to 21:47786817

- ▶ A peak appears indicating the variant position.

Visualizing variants

- ▶ Leave the cursor over that position on the vcf track:



Overview

Interface

Loading Data

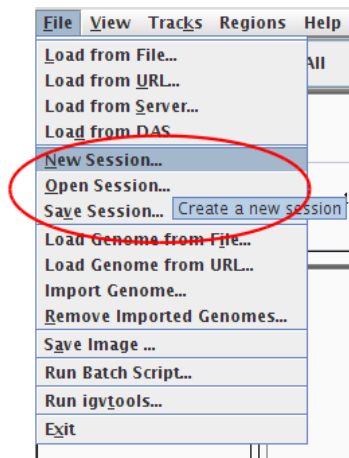
Browsing the Data

Sessions

Exercises

Create/Open/Save Sessions

- ▶ Your current session can be saved



Overview

Interface

Loading Data

Browsing the Data

Sessions

Exercises

Exercise 1

- ▶ Load `/home/biouser/mda13/mqc-igv/igv2.bam`
- ▶ Load `/home/biouser/mda13/mqc-igv/igv1.bam`
- ▶ Move to `chr11:1,016,174-1,018,316`
- ▶ Can you see any difference between both tracks?
- ▶ Is the `igv2.bam` mapping a good or a bad result?
- ▶ Can you think of the reasons that lead to such result?
- ▶ Can you think of any way of improving it?

Exercise 2

- ▶ Load `/home/biouser/mda13/mqc-igv/igv1.bam`
- ▶ Move to `21:47917047`
- ▶ What is happening in our sample in that position?

Exercise 3

- ▶ Load `/home/biouser/mda13/mqc-igv/igv1.bam`
- ▶ Move to `21:26973663`
- ▶ What is happening in our sample in that position?

Exercise 4

- ▶ Load `/home/biouser/mda13/mqc-igv/igv1.bam`
- ▶ Which of these variants would you trust? Why? Which is the sequence change in each case?
 1. 21:47821726
 2. 21:46596230
 3. 21:42848560
 4. 21:47917170

Exercise 5

- ▶ Considering that `/home/biouser/mda13/mqc-igv/igv.bed` contains the target regions of a given sequencing experiment and that `/home/biouser/mda13/mqc-igv/igv1.bam` is the mapping result, how well are the target regions of the gene MX2 covered?
- ▶ Focus on a particular target region and have a look at the histogram on the top of the `igv1.bam` track.
 - ▶ Which is the shape of the histogram? Can you explain this?
 - ▶ Which would be the desired shape?

Exercise 6

- ▶ Considering that `/home/biouser/mda13/mqc-igv/igv.bed` contains the target regions of a given sequencing experiment and that `/home/biouser/mda13/mqc-igv/igv1.bam` is the mapping result, look for unsequenced target regions of the PRMT2 gene.