

Moderní metody analýzy genomu

Bioinformatika II

Karol Pál

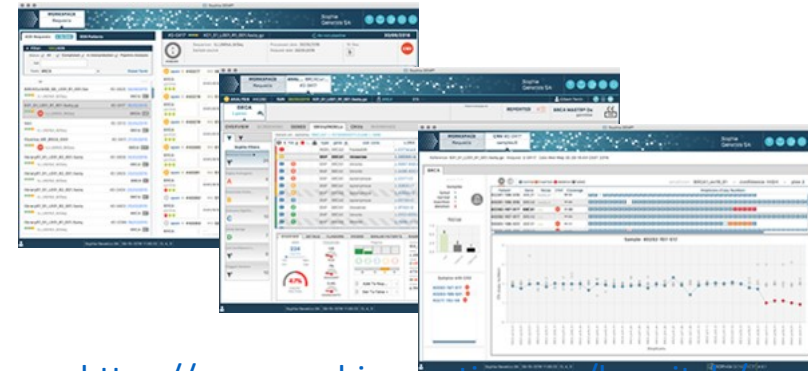
(Šárka Pospíšilová Research Group - Centre for Molecular Medicine
CEITEC)

Bioinformatics

- Use of computer to analyze and catalogue biological data
- Interdisciplinary field
- Algorithms for calculating local alignment 1970's
- New applications in NGS

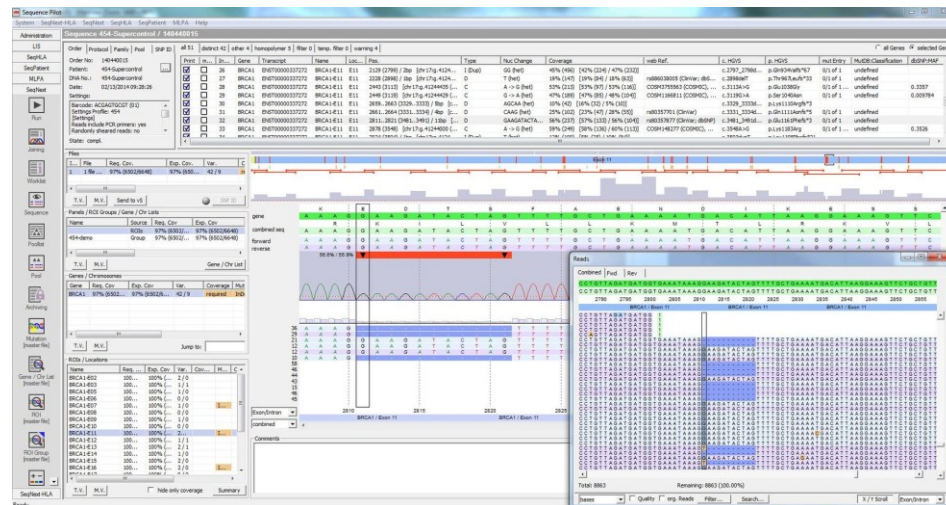
Commercial solutions

- - Expensive
- + Support
- Diagnostics



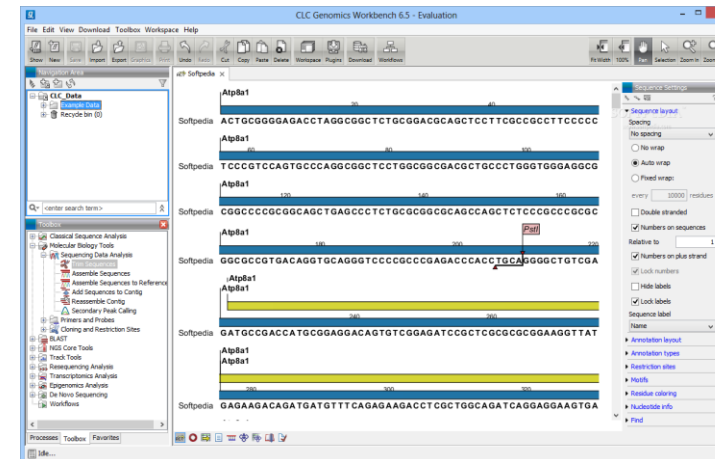
<https://www.sophiagenetics.com/hospitals/sophia-ddm/sophia-ddmr-details.html>

JSI Sequence pilot



<https://www.jsi-medisys.de/products/sequence-pilot/seqnext/>

CLC genomics workbench



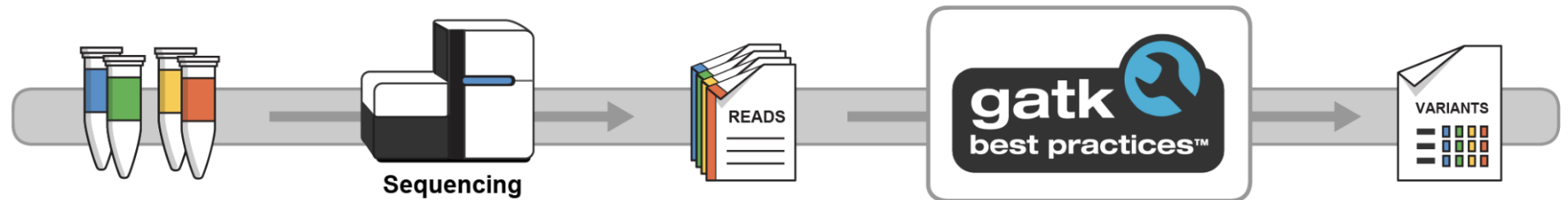
<https://www.qiagenbioinformatics.com/products/clc-genomics-workbench>

Open Source

- The basic steps are the same
- Community driven development by researchers
- Solutions from Open Source often implemented in commercial software
- - Expensive Bioinformatician
- Research

Genome Analysis Toolkit

Variant Discovery in High-Throughput Sequencing Data



Pipelines



Pipelines

- Workflow consisting of several steps
 - Data preprocessing
 - Quality Control
 - Detect Variation
 - Compare with known data sources
 - Present/Visualize results

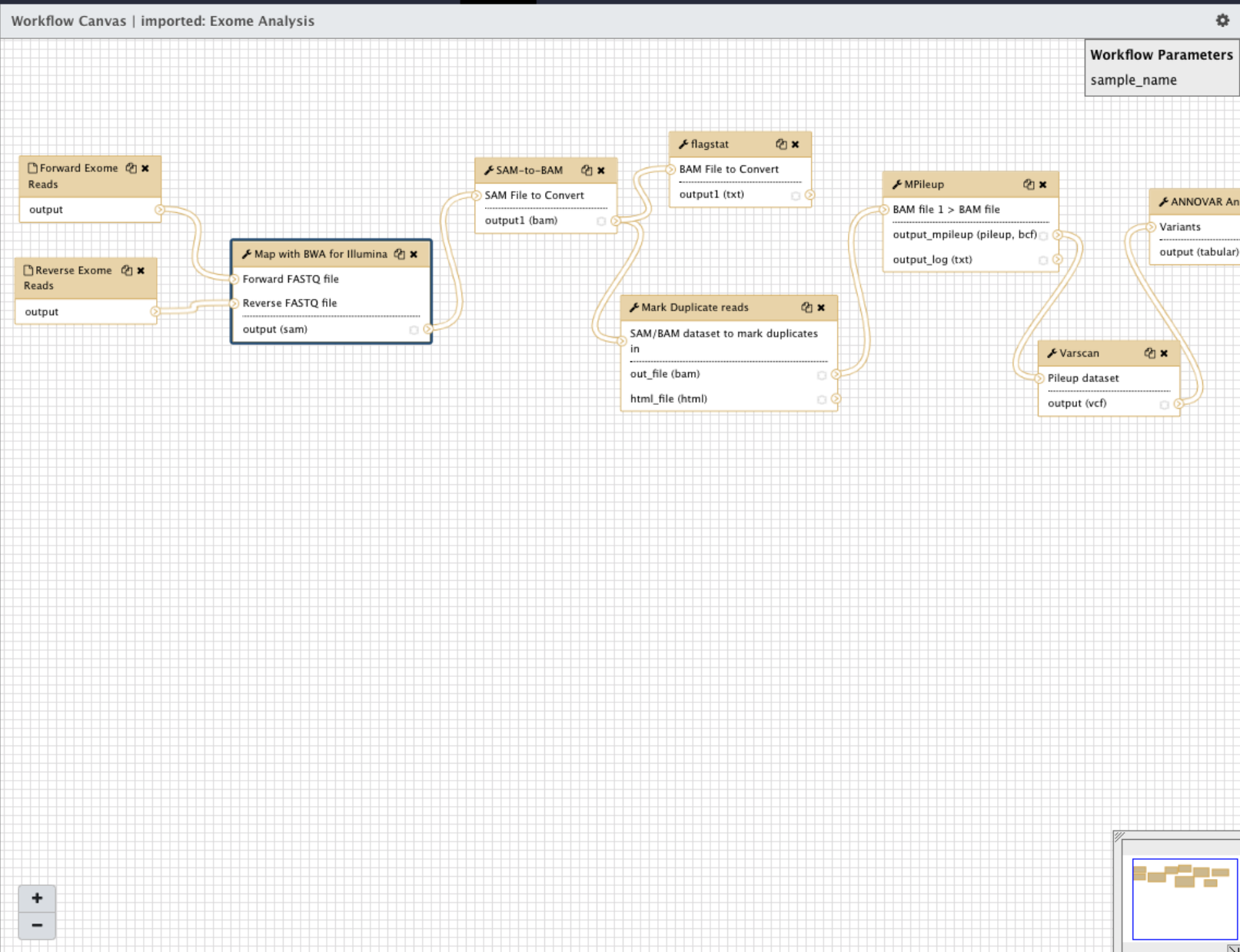
Pipelines

- Workflow consisting of several steps
- Each step can be a separate program in a different programming language
 - Perl – Old schoolers
 - R – Statisticians
 - C/C++ – Programmers
 - Python – Human beings

Pipeline manager

- Combine individual steps into one "package"
- Manage reusability
- Project organization
- Logs
- Examples
 - Snakemake
 - Bcbio
 - Galaxy (graphical interface) <https://usegalaxy.org/>
 - Bash

- Tools
- search tools
- Inputs
 - Get Data
 - Send Data
 - Lift-Over
 - Collection Operations
 - Text Manipulation
 - Datamash
 - Convert Formats
 - Filter and Sort
 - Join, Subtract and Group
 - Fetch Alignments/Sequences
 - NGS: QC and manipulation
 - NGS: DeepTools
 - NGS: Mapping
 - NGS: RNA Analysis
 - NGS: SAMtools
 - NGS: BamTools
 - NGS: Picard
 - NGS: VCF Manipulation
 - NGS: Peak Calling
 - NGS: Variant Analysis
 - NGS: RNA Structure
 - NGS: Du Novo
 - NGS: Gemini
 - NGS: Assembly
 - NGS: Chromosome Conformation
 - NGS: Mothur
 - Operate on Genomic Intervals
 - Statistics
 - Graph/Display Data
 - Phenotype Association
 - BEDTools
 - Genome Diversity



Details

Map with BWA for Illumina (Galaxy Version 1.2.3)

Workflow Parameters
sample_name

Label

Add a step label.

Annotation

Add an annotation or notes to this step. Annotations are available when a workflow is viewed.

Will you select a reference genome from your history or use a built-in index?

Use a built-in index

Select a reference genome

Human (Homo sapiens) (b37):...

Is this library mate-paired?

Paired-end

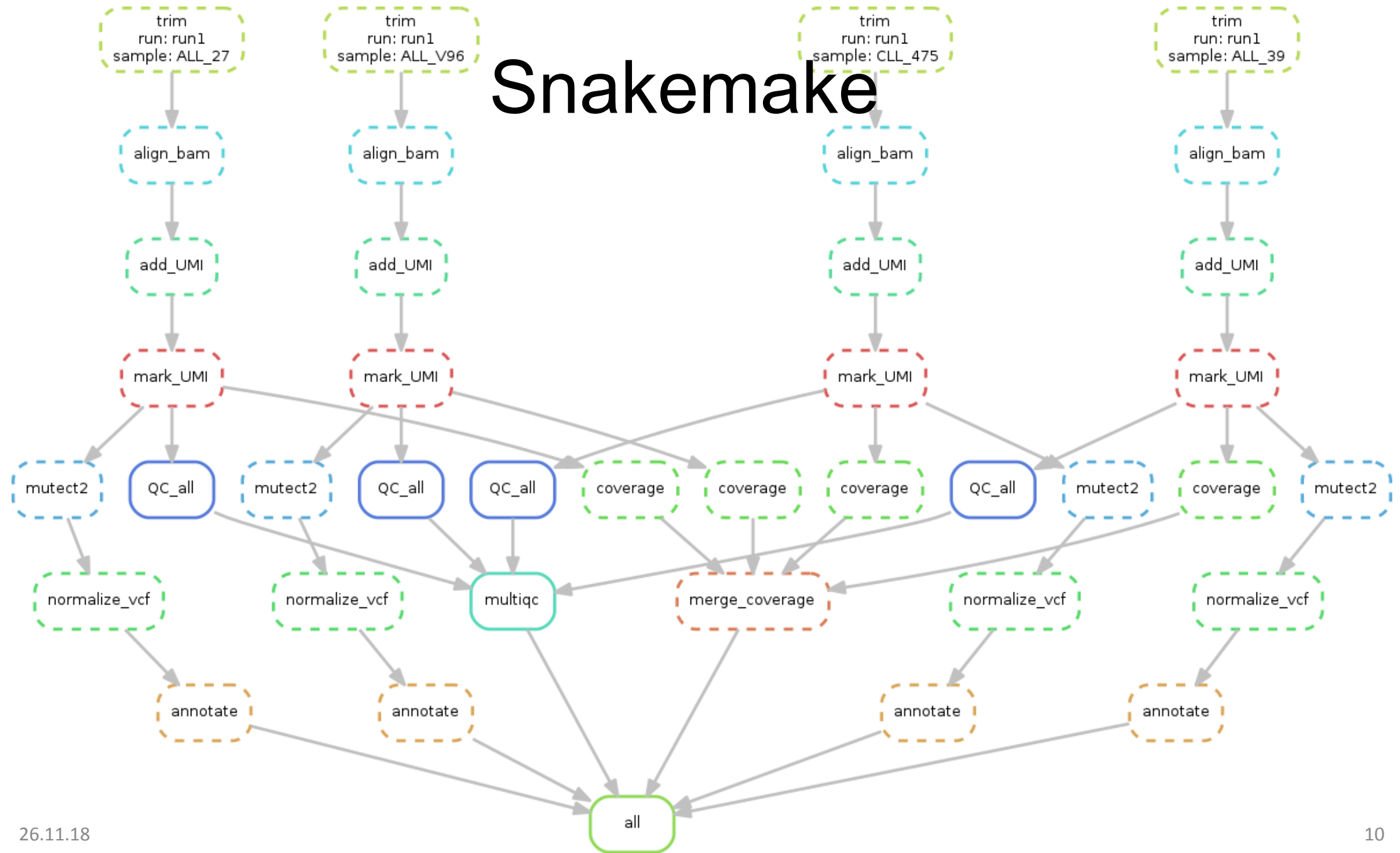
Forward FASTQ file

Data input 'input1' (fastqsanger or fastqillumina)
FASTQ with either Sanger-scaled quality values (fastqsanger) or Illumina-scaled quality values (fastqillumina)

Reverse FASTQ file

Data input 'input2' (fastqsanger or fastqillumina)
FASTQ with either Sanger-scaled quality values (fastqsanger) or Illumina-scaled quality values (fastqillumina) 9

Snakemake



A note on bash

- A command line environment for controlling the computer
- Scripting language
- Fast, built in command for text manipulation
- Working with text files too big for a text editor

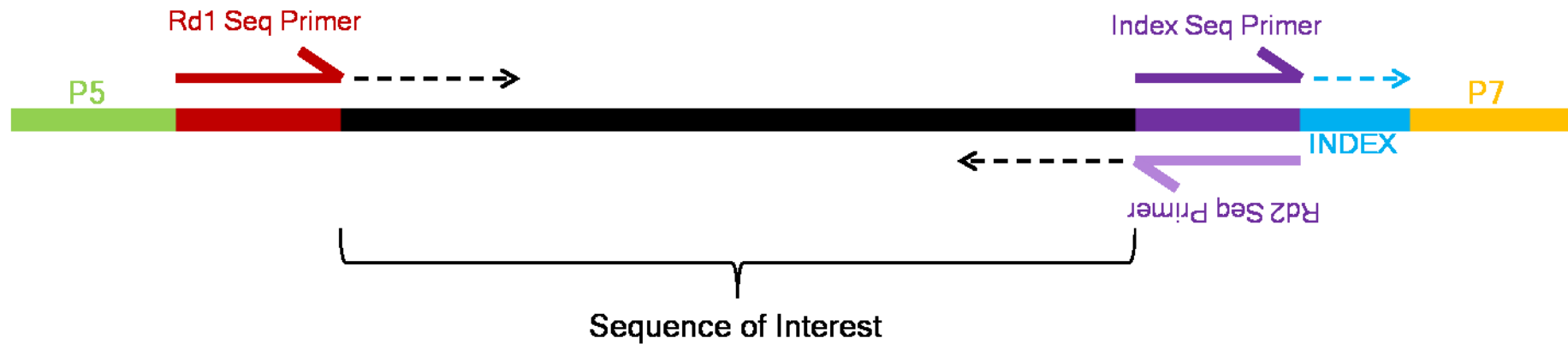
A note on bash

- A command line environment for controlling the computer
- Scripting language
- Fast, built in command for text manipulation
- Working with text files too big for a text editor

```
$ zcat JJ1462_dia.vardict.filt.sorted.vcf.gz | grep 'PASS' -c
935
$
```

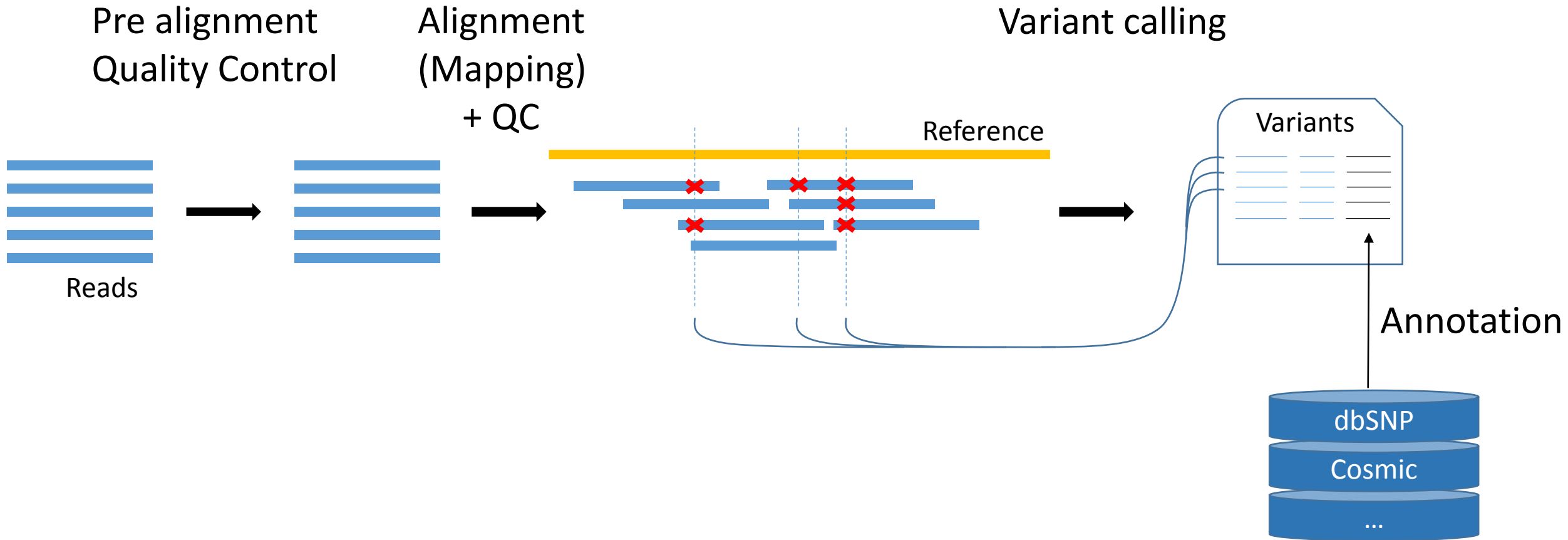
Recap – Sequencing reads

STRUCTURE DETAILS

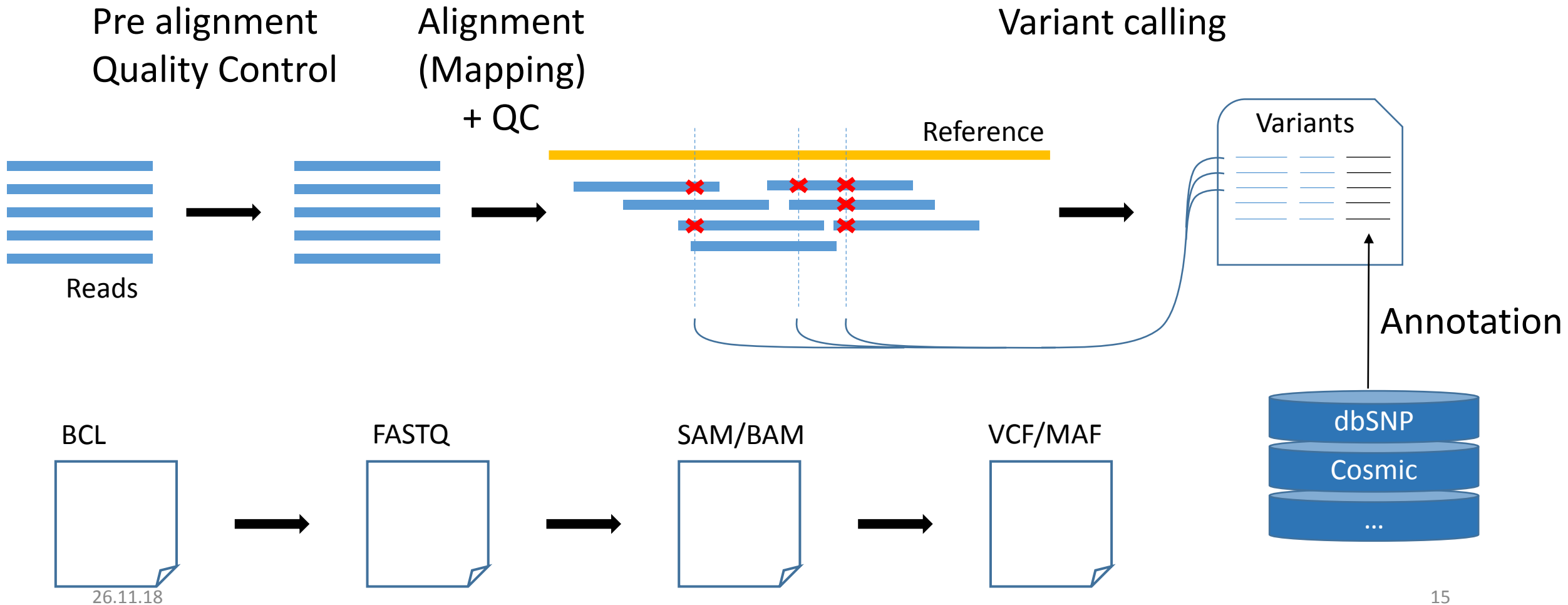


<http://nextgen.mgh.harvard.edu/CustomPrimer.html>

Data analysis pipeline (DNaseSeq)

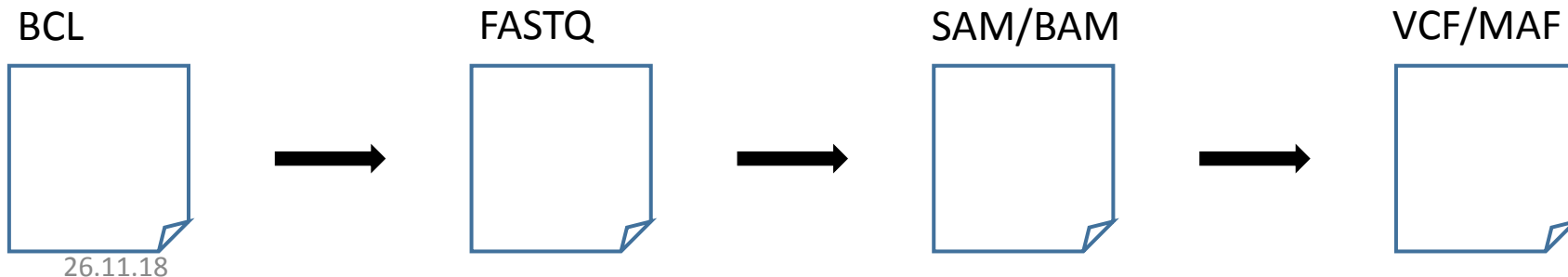


Data analysis pipeline



Quality Control

- Different steps
 - Reads summary statistics
 - Aligners
 - Post alignment statistics
- Different tools
 - Different kinds of outputs



MultiQC

Aggregate results from bioinformatics analyses across many samples into a single report

MultiQC searches a given directory for analysis logs and compiles a HTML report. It's a general use tool, perfect for summarising the output from numerous bioinformatics tools.



Introduction to MultiQC (1:19)

Installing MultiQC (4:33)

Running MultiQC (5:21)

Using MultiQC Reports (6:06)

 [GitHub](#)

 [Python Package Index](#)

 [Documentation](#)

 [56 supported tools](#)

 [Publication / Citation](#)

 [Get help on Gitter](#)

[Quick Install](#)

```
pip install multiqc # Install
multiqc .           # Run
```

[pip](#) [conda](#) [manual](#)

Need a little more help? [See the full installation instructions.](#)

MultiQC

Aggregate results from bioinformatics analyses across many samples into a single report

MultiQC searches a given directory for analysis logs and compiles a HTML report. It's a general use tool, perfect for summarising the output from numerous bioinformatics tools.



Introduction to MultiQC (1:19)


Installing MultiQC (4:33)

Running MultiQC (5:21)


Using MultiQC Reports (6:06)

 [GitHub](#)

 [Python Package Index](#)

 [Documentation](#)

 [56 supported tools](#)

 [Publication / Citation](#)

 [Get help on Gitter](#)

[Quick Install](#)

```
pip install multiqc # Install
multiqc .           # Run
```

[pip](#) [conda](#) [manual](#)

Need a little more help? [See the full installation instructions.](#)

MultiQC: Supported Tools

Pre-alignment tools

Alignment tools

Post-alignment tools

Quality Control

Skewer	Skewer is an adapter trimming tool specially designed for processing next-generation sequencing (NGS) paired-end sequences.
SortMeRNA	SortMeRNA is a program tool for filtering, mapping and OTU-picking NGS reads in metatranscriptomic and metagenomic data.
Trimmomatic	Trimmomatic is a flexible read trimming tool for Illumina NGS data
Bismark	Bismark is a tool to map bisulfite converted sequence reads and determine cytosine methylation states.
Bowtie 1	Bowtie 1 is an ultrafast, memory-efficient short read aligner.
Bowtie 2	Bowtie 2 is an ultrafast and memory-efficient tool for aligning sequencing reads to long reference sequences.
BBMap	BBMap is a suite of pre-processing, assembly, alignment, and statistics tools for DNA/RNA sequencing reads.
HiCUP	HiCUP (Hi-C User Pipeline) is a tool for mapping and performing quality control on Hi-C data.
HISAT2	HISAT2 is a fast and sensitive alignment program for mapping NGS reads (both DNA and RNA) to reference genomes.
Kallisto	kallisto is a program for quantifying abundances of transcripts from RNA-Seq data.
Salmon	Salmon is a tool for quantifying the expression of transcripts using RNA-seq data.
STAR	STAR is an ultrafast universal RNA-seq aligner.
TopHat	TopHat is a fast splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to mammalian-sized genomes.
Bamtools	BamTools provides both a programmer's API and an end-user's toolkit for handling BAM files.
Bcftools	BCFtools is a set of utilities that manipulate variant calls in the Variant Call Format (VCF) and its binary counterpart BCF.
BUSCO	BUSCO assesses genome assembly and annotation completeness with Benchmarking Universal Single-Copy Orthologs.
Conpair	Conpair estimates concordance and contamination for tumour-normal pairs
Disambiguate	Disambiguation algorithm for reads aligned to two species (e.g. human and mouse genomes) from Tophat, Hisat2, STAR or BWA mem.

Quality Control

MultiQC Example Reports RNA-Seq Whole-Genome Seq Bisulfite Seq HI-C MultiQC_NGI

MultiQC
v1.3

General Stats

- featureCounts
- STAR
- Cutadapt
- FastQC
- Sequence Quality Histograms
- Per Sequence Quality Scores
- Per Base Sequence Content
- Per Sequence GC Content
- Per Base N Content
- Sequence Length Distribution
- Sequence Duplication Levels
- Overrepresented sequences
- Adapter Content

26.11.18

MultiQC

A modular tool to aggregate results from bioinformatics analyses across many samples into a single report.

Report generated on 2017-11-03, 14:21 based on data in:
`/Users/ewels/GitHub/MultiQC_website/public_html/examples/rna-seq`

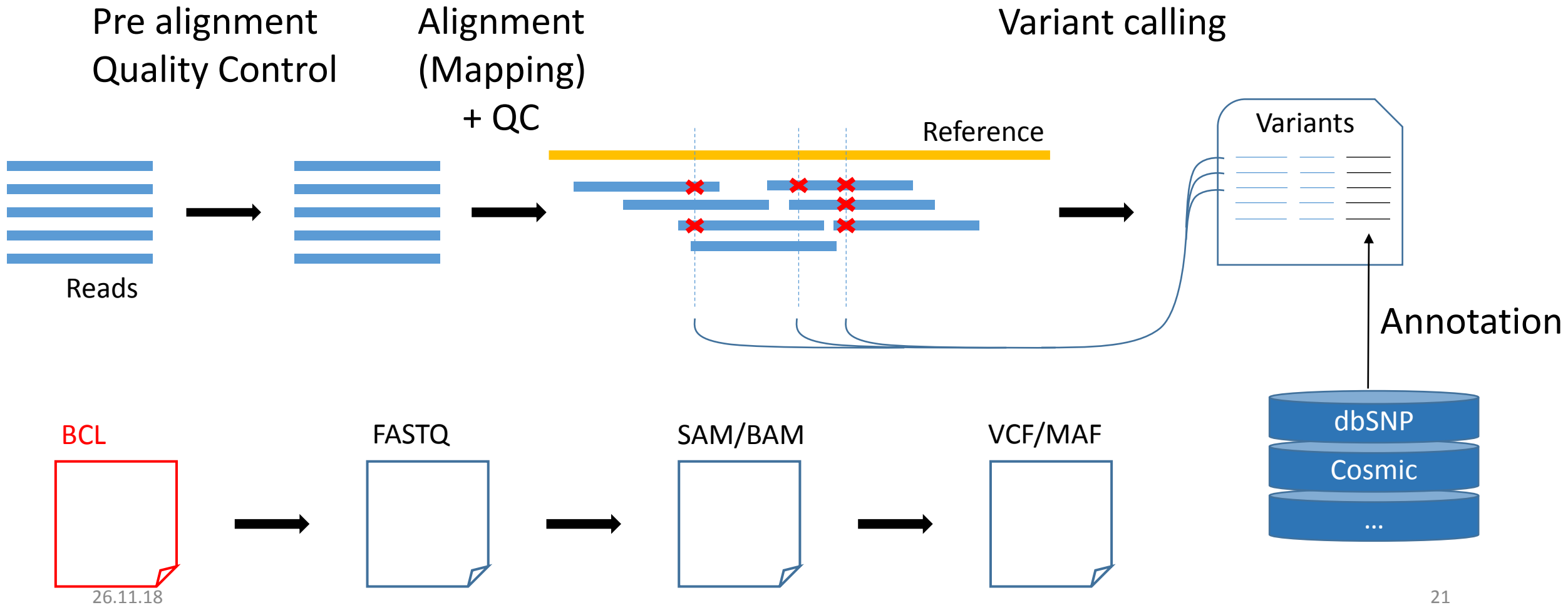
Welcome! Not sure where to start? [Watch a tutorial video](#) (6:06) [don't show again](#) ✕

General Statistics

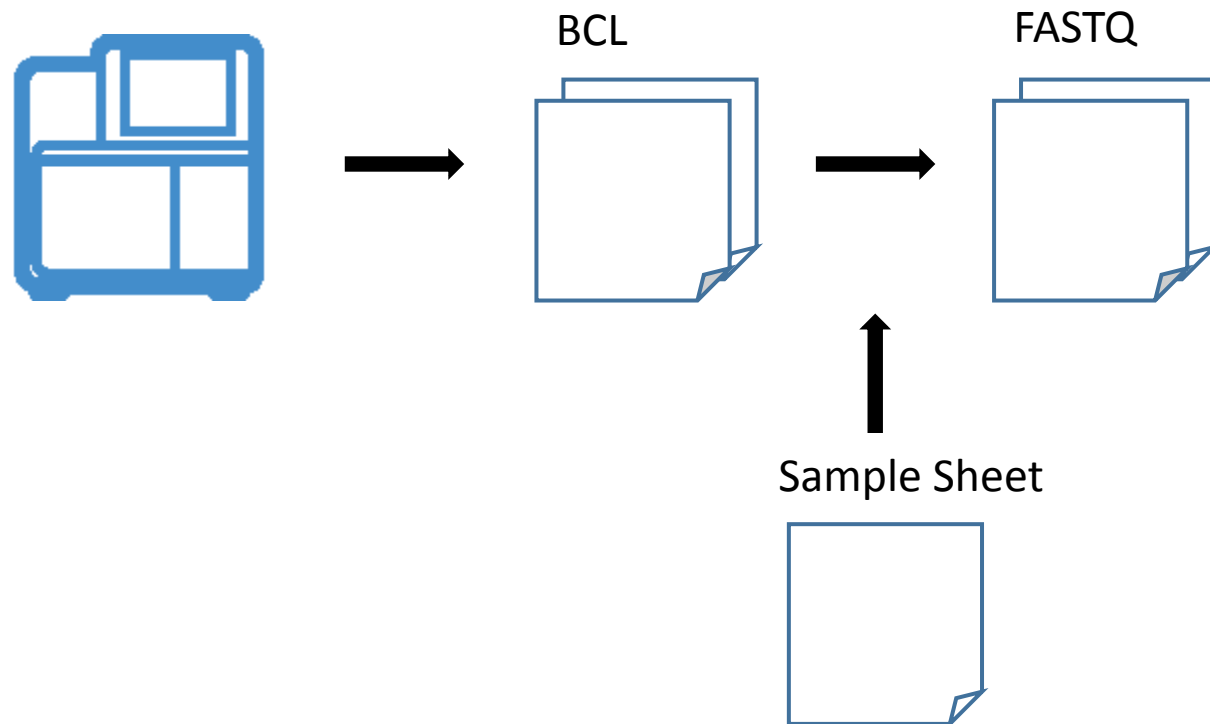
[Copy table](#) [Configure Columns](#) [Plot](#) Showing 8/8 rows and 8/10 columns.

Sample Name	% Assigned	M Assigned	% Aligned	M Aligned	% Trimmed	% Dups	% GC	M Seqs
SRR3192396	67.5%	71.9	93.7%	97.8	4.0%	78.9%	51%	104.4
SRR3192397	66.6%	63.0	94.7%	87.1	3.5%	77.2%	49%	92.0
SRR3192398	50.9%	36.5	88.2%	58.7	5.0%	55.3%	47%	66.6
SRR3192399	52.3%	42.3	88.2%	65.6	5.0%	57.4%	47%	74.3
SRR3192400	70.3%	63.4	77.3%	73.4	7.2%	74.1%	45%	94.9

Data analysis pipeline



BCL to FASTQ



- BCL - raw sequencing output
- Convert to FASTQ format
- Split into sample files
- May be automated

Sample sheet

[Header]

IEMFileVersion,4
Experiment Name,Exom.20171013
Date,9.10.2017
Workflow,GenerateFASTQ
Application,FASTQ Only
Assay,TruSeq LT
Description,
Chemistry,Default

[Reads]

[Settings]

ReverseComplement,0

[Data]

Sample_ID,Sample_Name,Sample_Plate,Sample_Well,I7_Index_ID,index,Sample_Project,Description,
BRN01077_normal,,,,AD002,CGATGT,,
BRN01404_normal,,,,AD007,CAGATC,,
BRN01503_normal,,,,AD019,GTGAAA,,

Sample sheet

[Header]

IEMFileVersion,4

Experiment Name,Exom.20171013

Date,9.10.2017

Workflow,GenerateFASTQ

Application,FASTQ Only

Assay,TruSeq LT

Description,

Chemistry,Default

[Reads]

[Settings]

ReverseComplement,0

[Data]

Sample_ID,Sample_Name,Sample_Plate,Sample_Well,I7_Index_ID,index,Sample_Project,Description,

BRN01077_normal,,AD002,CGATGT,,

BRN01404_normal,,AD007,CAGATC,,

BRN01503_normal,,AD019,GTGAAA,,

Raw reads – bcl2fastq

MultiQC output

General Statistics

Copy table

Configure Columns

Plot

Showing 4/4 rows and 3/3 columns.

Sample Name	Total Reads	Mb Yield \geq Q30	% Perfect Index
BRNO062_tumor	54 513 947.0	7 943.6	100.0%
BRNO0047_tumor	52 169 492.0	7 596.2	100.0%
BRNO1503_tumor	49 439 468.0	7 199.7	100.0%
undetermined	8 933 116.0	1 024.4	0.0%

Raw reads – bcl2fastq

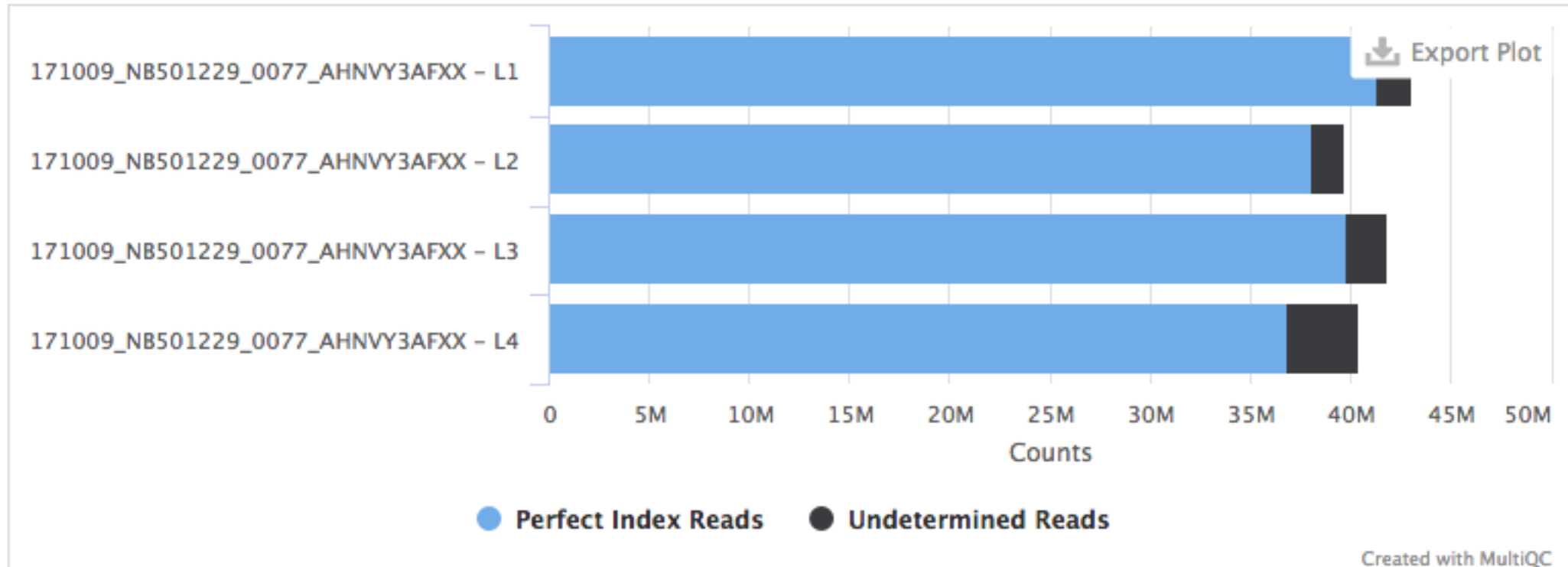
Clusters by lane

Number of reads per lane (with number of perfect index reads)

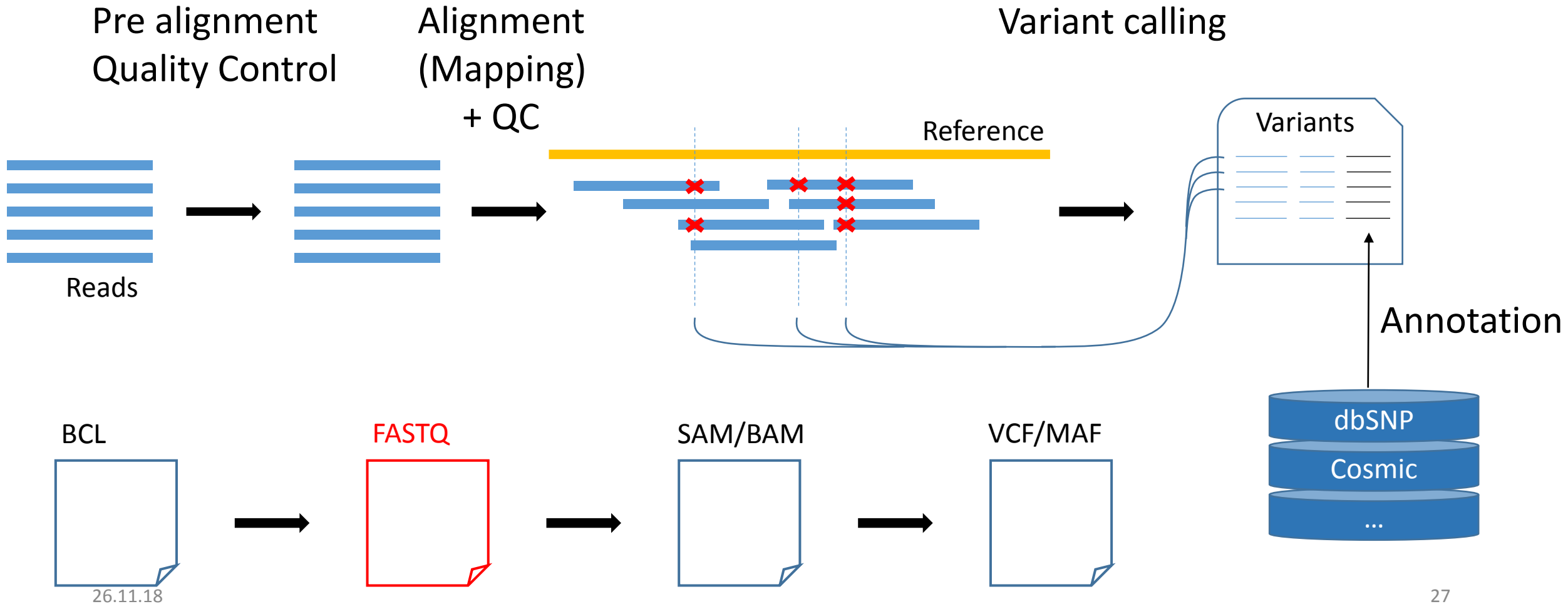
[Help](#)

MultiQC output

Counts Percentages



Data analysis pipeline

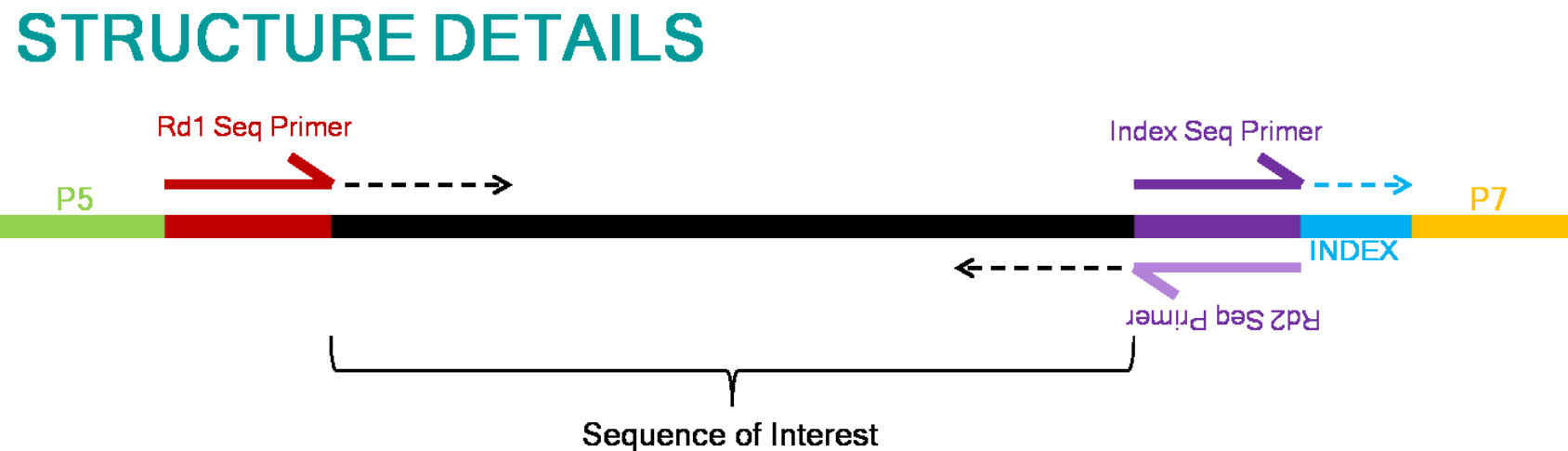


FastQC

- Summary statistics
- Two modes
 - Stand alone program
 - Command line (output can be integrated to MultiQC)
- Input: Fastq or BAM file

Trimming

- Adaptors
- Low quality ends of reads
- Tools:
 - Cutadapt
 - Trimmomatic

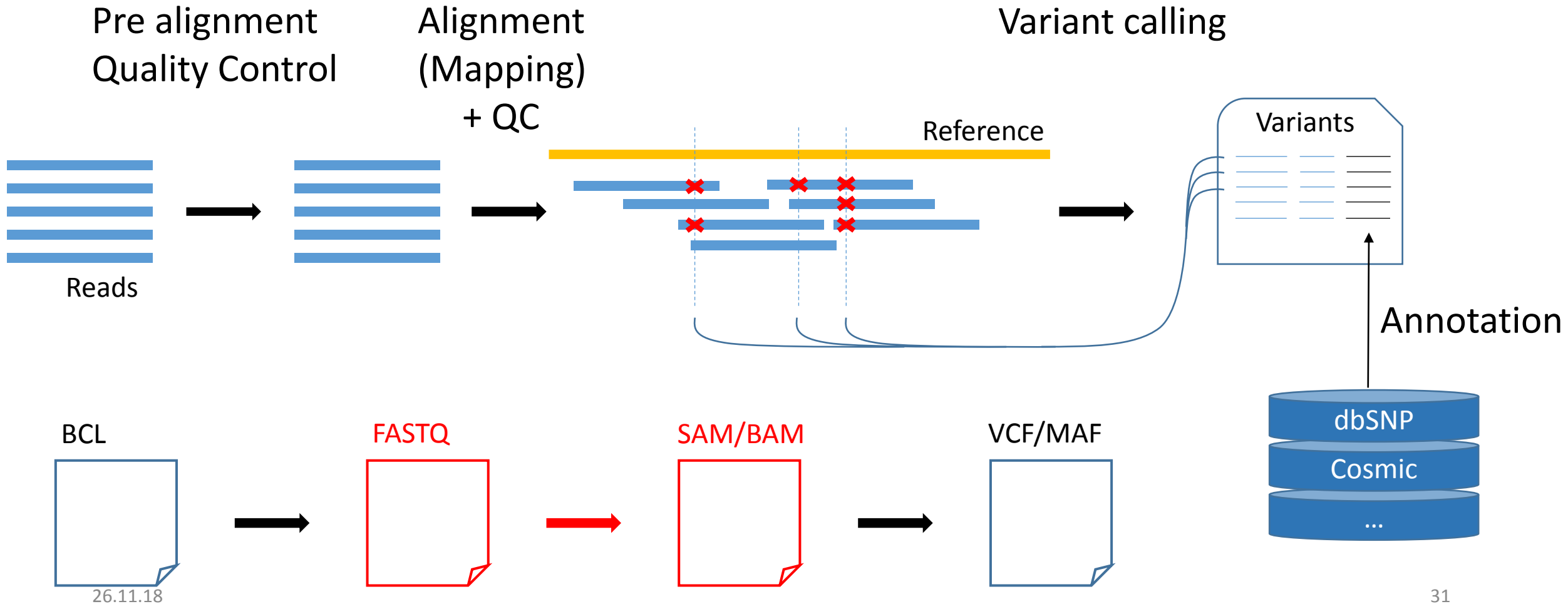


Trimming



```
> cutadapt
-a AGATCGGAAGAGC \
-A AGATCGGAAGAGC \
-o BR_0296_I.trimmed.1.fastq.gz \
-p BR_0296_I.trimmed.2.fastq.gz \
BR_0296_I.R1.fq.gz BR_0296_I.R2.fq.gz
```

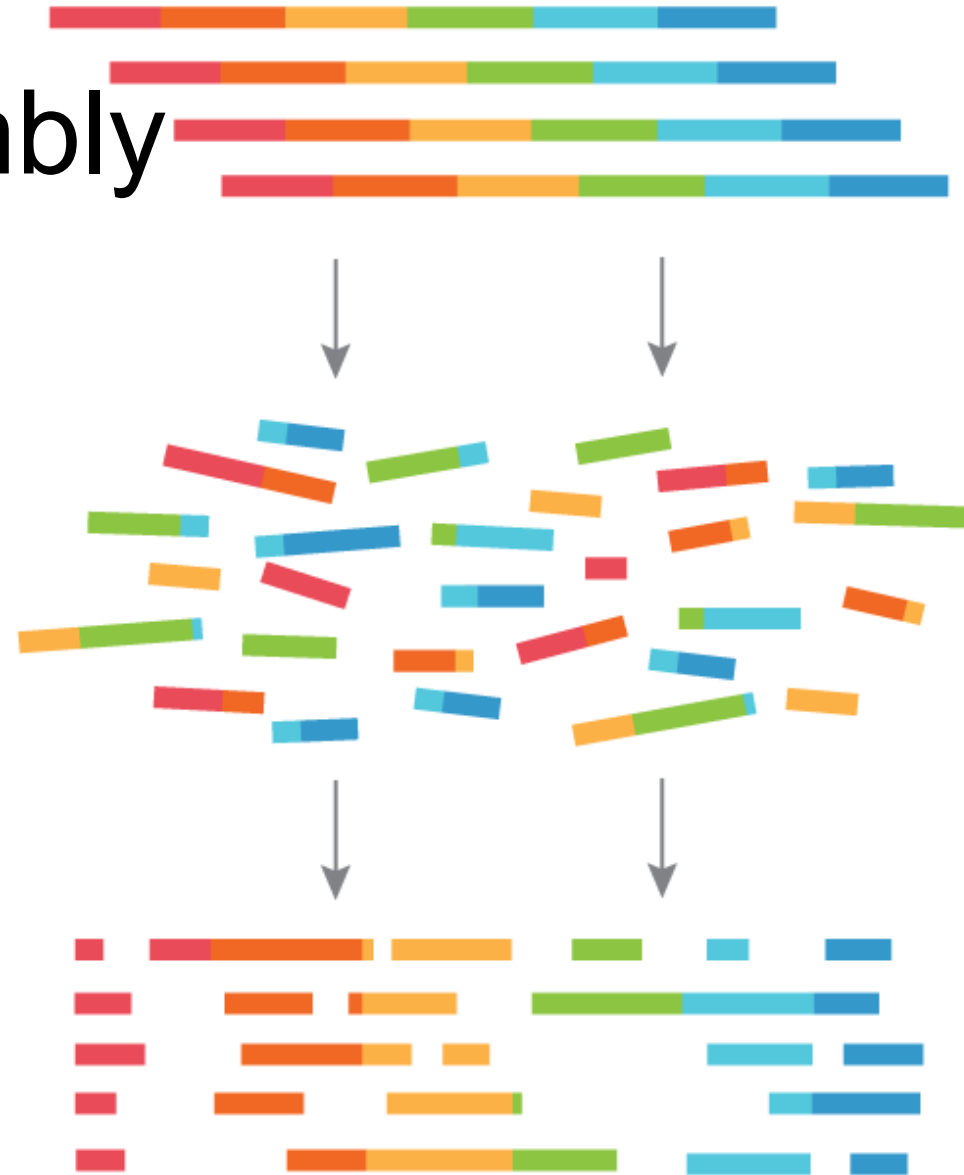
Recap – Data analysis pipeline



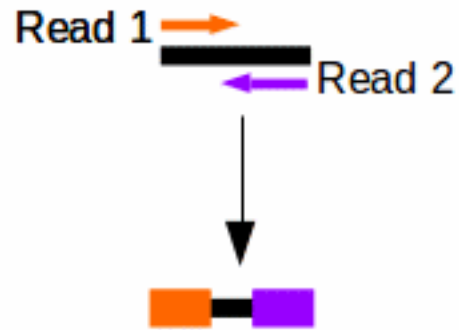
DNA

- De Novo Assembly
 - Create a new reference
 - Find structural variants
- Map to an existing reference
 - Alignment (BWA)
- Map against several references
 - Blast

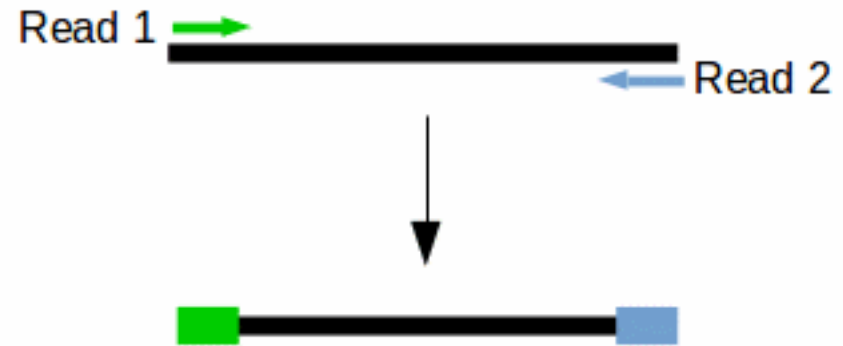
De Novo Assembly



Short-insert paired-end reads



Long-insert paired-end reads (Mate pair)



De novo sequencing



De Novo Assembly

SCIENTIFIC REPORTS







OPEN

***De novo* yeast genome assemblies from MinION, PacBio and MiSeq platforms**

Received: 17 January 2017

Accepted: 8 May 2017

Published online: 21 June 2017

Francesca Giordano¹, Louise Aigrain¹, Michael A Quail¹, Paul Coupland², James K Bonfield¹, Robert M Davies¹, German Tischler³, David K Jackson ¹, Thomas M Keane¹, Jing Li ⁴, Jia-Xing Yue ⁴, Gianni Liti⁴, Richard Durbin ¹ & Zemin Ning¹

Alignment

Consensus contig

ACGCGATTCAGGTTACCACGCGTAGCGCATTACACAGATTAG

Aligned reads



ACGCGATTCAGGTTACCACG
GCGATTCAGGTTACCACGCG
GATTCAGGTTACCACGCGTA
TTCAGGTTACCACGCGTAGC
CAGGTTACCACGCGTAGCGC
GGTTACCACGCGTAGCGCAT
TTACCACGCGTAGCGCATT
ACCACGCGTAGCGCATTACA
CACGCGTAGCGCATTACACA
CGCGTAGCGCATTACACAGA
CGTAGCGCATTACACAGATT
TAGCGCATTACACAGATTAG

Alignment

GCTGATGTGCCGCCTCACTTCGGTGGTGAGGTG Reference sequence

CTGATGTGCCGCCTCACTTCGGTGGT	Short read 1
TGATGTG-CGCCTCACT A CGGTGGTG	Short read 2
GATGTG-CGCCTCACTTCGGTGGTGA	Short read 3
GCTGATGTGCCGCCTCACT A CGGTG	Short read 4
GCTGATGTGCCGCCTCACT A CGGTG	Short read 5

Alignment

BED file

Chr7 127471196 127472363 Pos1 0 +

GCTGATGTGCCGCCTCACTTCGGTGGTGAGGTG Reference sequence

CTGATGTGCCGCCTCACTTCGGTGGT

Short read 1

TGATGTG-CGCCTCACT**A**CGGTGGTG

Short read 2

GATGTG-CGCCTCACTTCGGTGGTGA

Short read 3

GCTGATGTGCCGCCTCACT**A**CGGTG

Short read 4

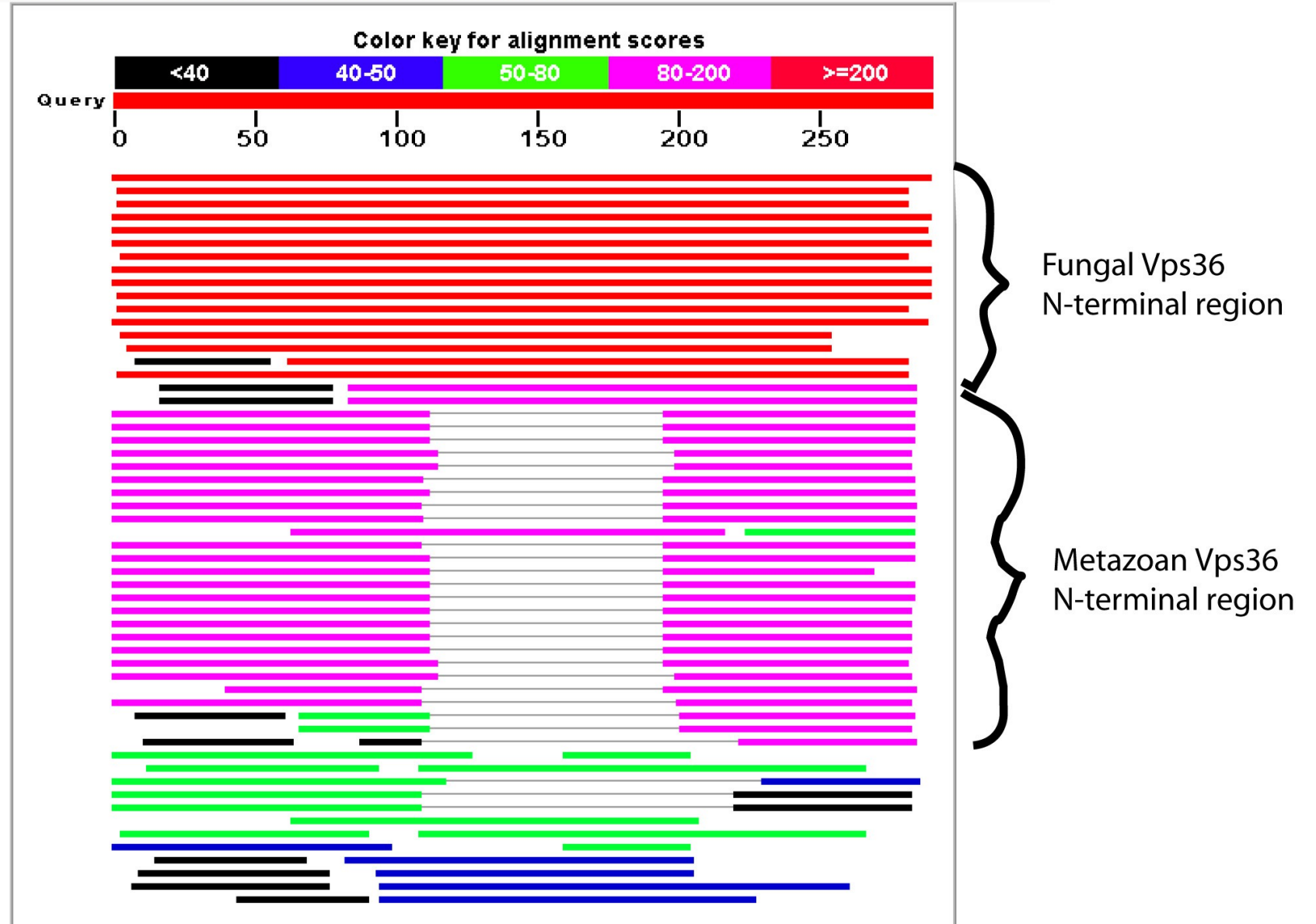
GCTGATGTGCCGCCTCACT**A**CGGTG

Short read 5

Distribution of 440 Blast Hits on the Query Sequence

Blast

Mouse-over to show define and scores, click to show alignments



No alignment?

[Nucleic Acids Res.](#) 2017 Jan 9; 45(1): 39–53.

PMCID: PMC5224470

Published online 2016 Nov 28. doi: [10.1093/nar/gkw1002](https://doi.org/10.1093/nar/gkw1002)

Alignment-free d_2^* oligonucleotide frequency dissimilarity measure improves prediction of hosts from metagenomically-derived viral sequences

[Nathan A. Ahlgren](#),^{1,*†} [Jie Ren](#),^{2,†} [Yang Young Lu](#),² [Jed A. Fuhrman](#),¹ and [Fengzhu Sun](#)^{1,2,3}

[Author information](#) ▶ [Article notes](#) ▶ [Copyright and License information](#) ▶

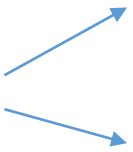
This article has been [cited by](#) other articles in PMC.

Alignment to human genome

- GRCh37 vs hg19 released 2009

Fasta format:

Unique
sequence name



```
>seq1  
ACGTCGTG  
>seq2 additional info  
TCGCAGCG
```

Alignment to human genome

- GRCh37(NCBI) vs hg19(UCSC) released 2009

```
173390@BioDA-server /m/n/s/0/r/G/seq> cat GRCh37.fa | grep ">"
```

```
>1 dna:chromosome chromosome:GRCh37:1:1:249250621:1  
>2 dna:chromosome chromosome:GRCh37:2:1:243199373:1  
>3 dna:chromosome chromosome:GRCh37:3:1:198022430:1  
>4 dna:chromosome chromosome:GRCh37:4:1:191154276:1  
>5 dna:chromosome chromosome:GRCh37:5:1:180915260:1
```

...

```
>22 dna:chromosome chromosome:GRCh37:22:1:51304566:1  
>X dna:chromosome chromosome:GRCh37:X:1:155270560:1  
>Y dna:chromosome chromosome:GRCh37:Y:2649521:59034049:1  
>MT gil251831106|ref|NC_012920.1| Homo sapiens mitochondrion, complete genome  
>GL000207.1 dna:supercontig supercontig::GL000207.1:1:4262:1  
>GL000226.1 dna:supercontig supercontig::GL000226.1:1:15008:1  
>GL000229.1 dna:supercontig supercontig::GL000229.1:1:19913:1
```

Alignment to human genome

- GRCh37(NCBI) vs hg19(UCSC) released 2009

```
173390@BioDA-server /m/n/s/0/r/G/seq> 173390@BioDA-server /m/n/s/0/r/h/seq> grep ">" hg19.fa
>1 dna:chromosome chromosome:GRCh37:1>chrM
>2 dna:chromosome chromosome:GRCh37:2>chr1
>3 dna:chromosome chromosome:GRCh37:3>chr2
>4 dna:chromosome chromosome:GRCh37:4>chr3
>5 dna:chromosome chromosome:GRCh37:5>chr4
...
...
>22 dna:chromosome chromosome:GRCh37:>chr21
>X dna:chromosome chromosome:GRCh37:X>chr22
>Y dna:chromosome chromosome:GRCh37:Y>chrX
>MT gil251831106|ref|NC_012920.1| Hom>chrY
>GL000207.1 dna:supercontig supercont>chr1_g1000191_random
>GL000226.1 dna:supercontig supercont>chr1_g1000192_random
>GL000229.1 dna:supercontig supercont>chr4_ctg9_hap1
>chr4_g1000193_random
```

Alignment to human genome

GRCh37(NCBI) vs hg19(UCSC) released Feb 2009

VS

GRCh38(NCBI) or hg38(UCSC) released Dec 2013

Alignment to human genome

Heng Li's blog

Archive

Categories

Pages

Tags

GRCh37 GRCh38 **Which human reference genome to use?**

13 November 2017

TL;DR: If you map reads to GRCh37 or hg19, use `hs37-1kg`:

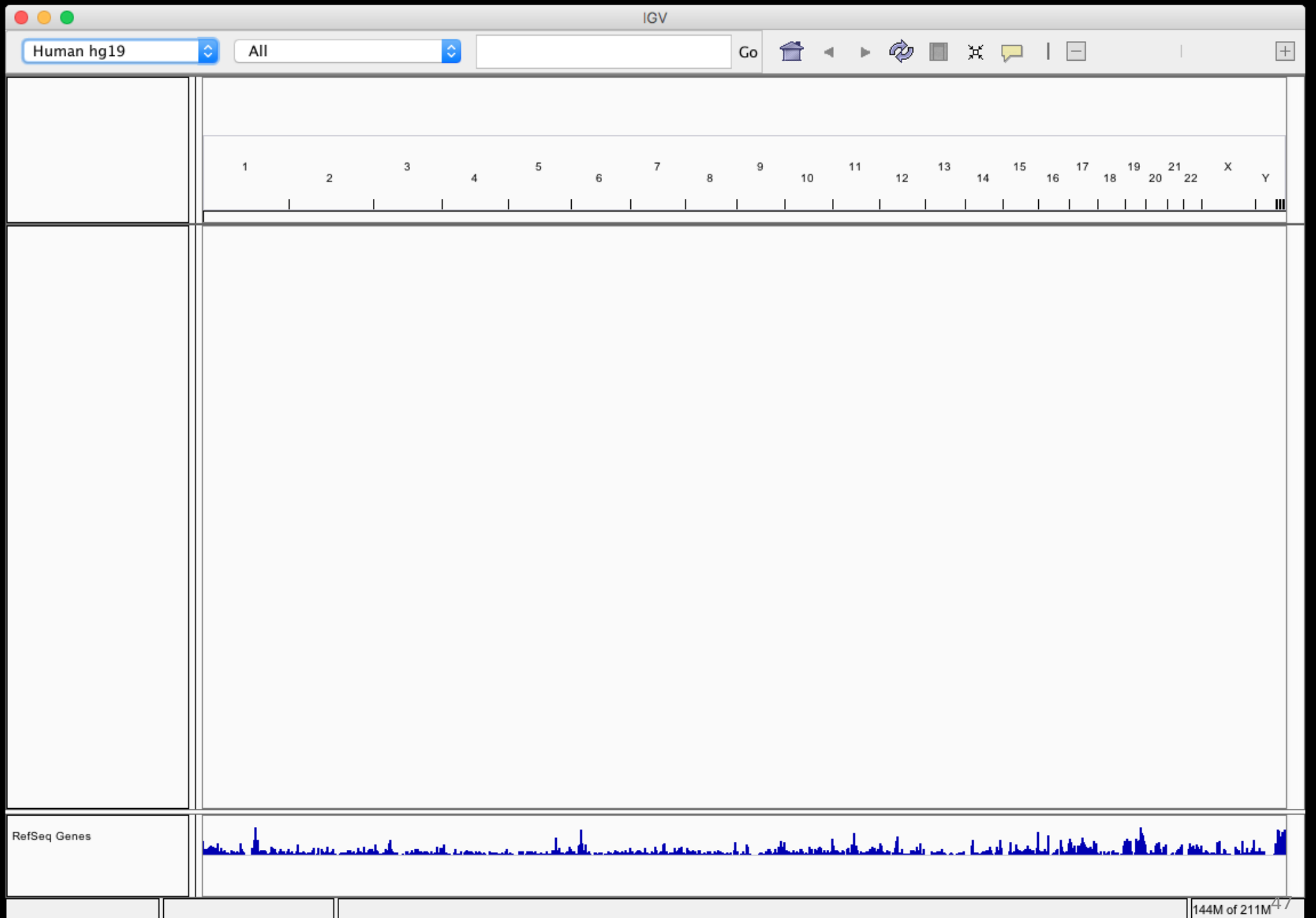
```
ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp/technical/reference/human_g1k_v37.fasta.gz
```

If you map to GRCh37 and believe decoy sequences help with better variant calling, use `hs37d5`:

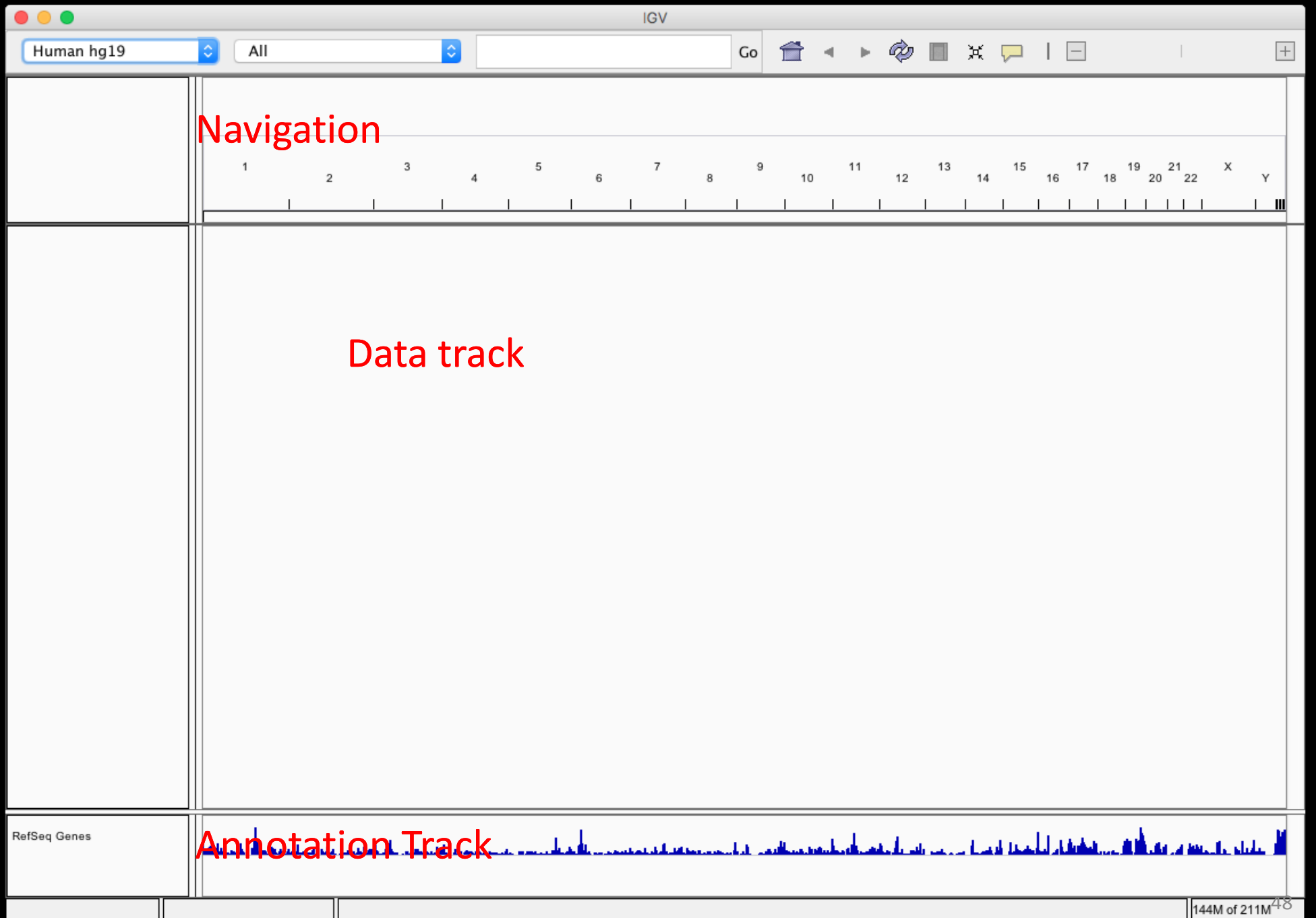
Alignment to genome

NCBI/UCSC applies also to the mouse genome
GRCm38/mm10

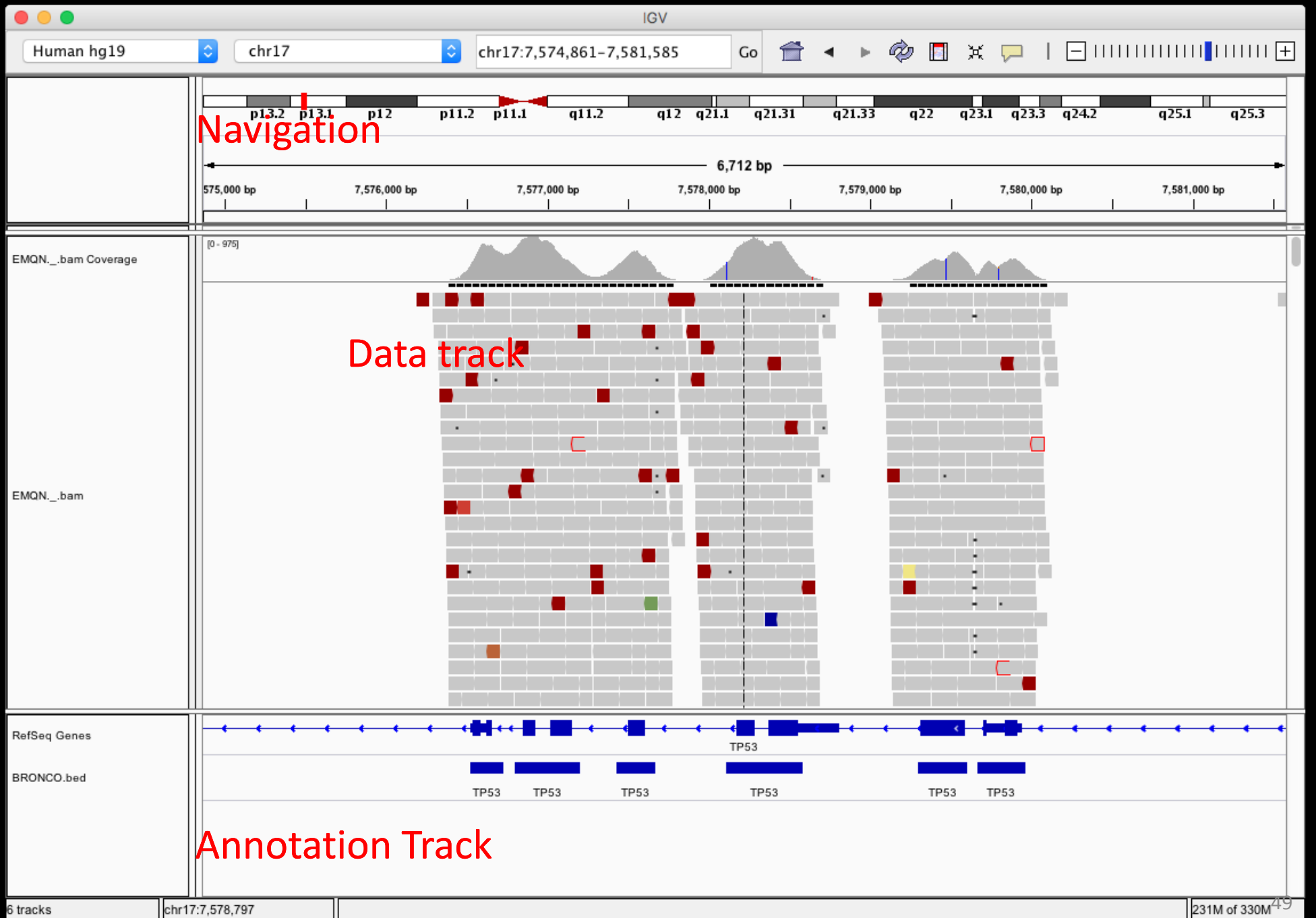
IGV



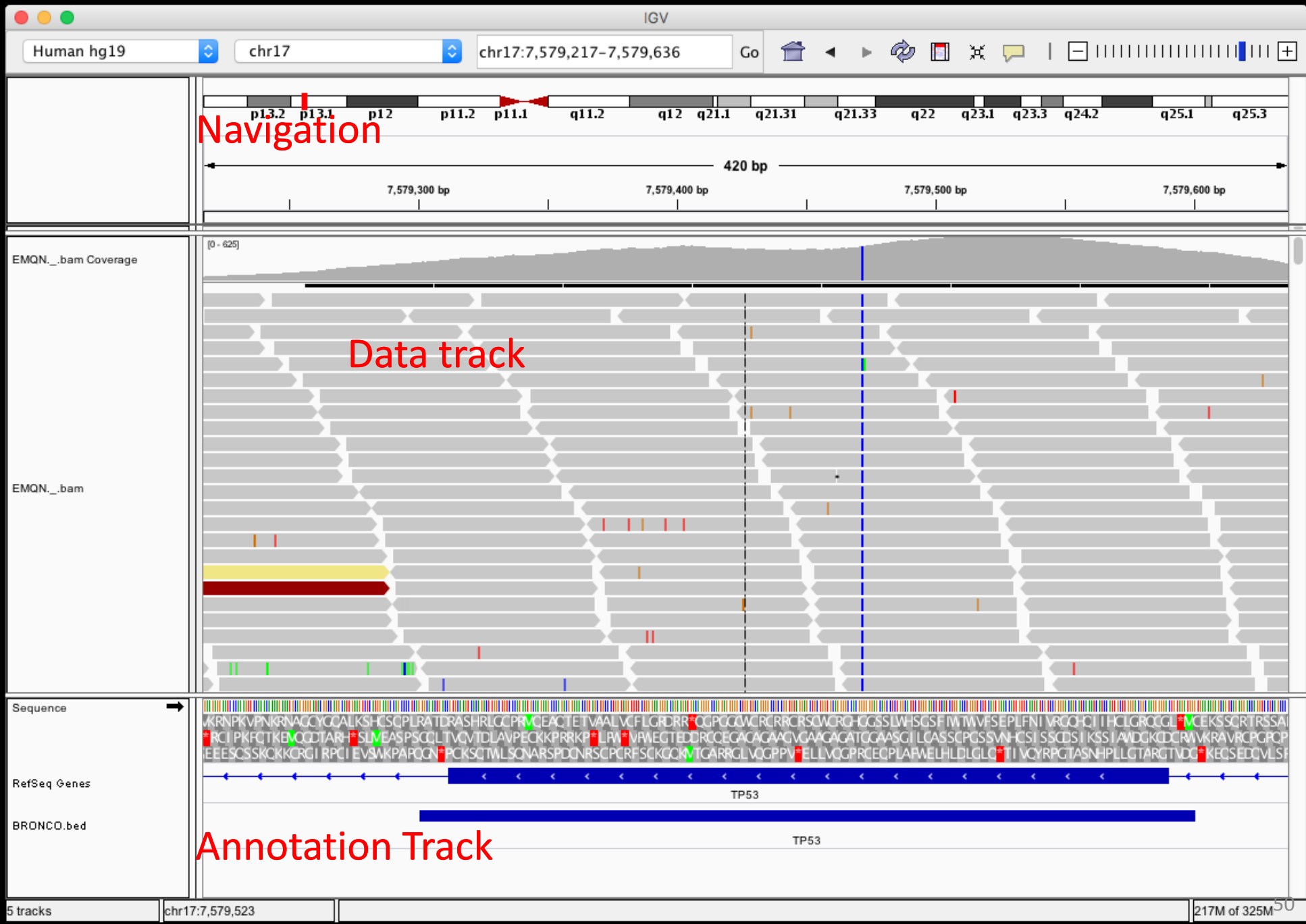
IGV



IGV



IGV

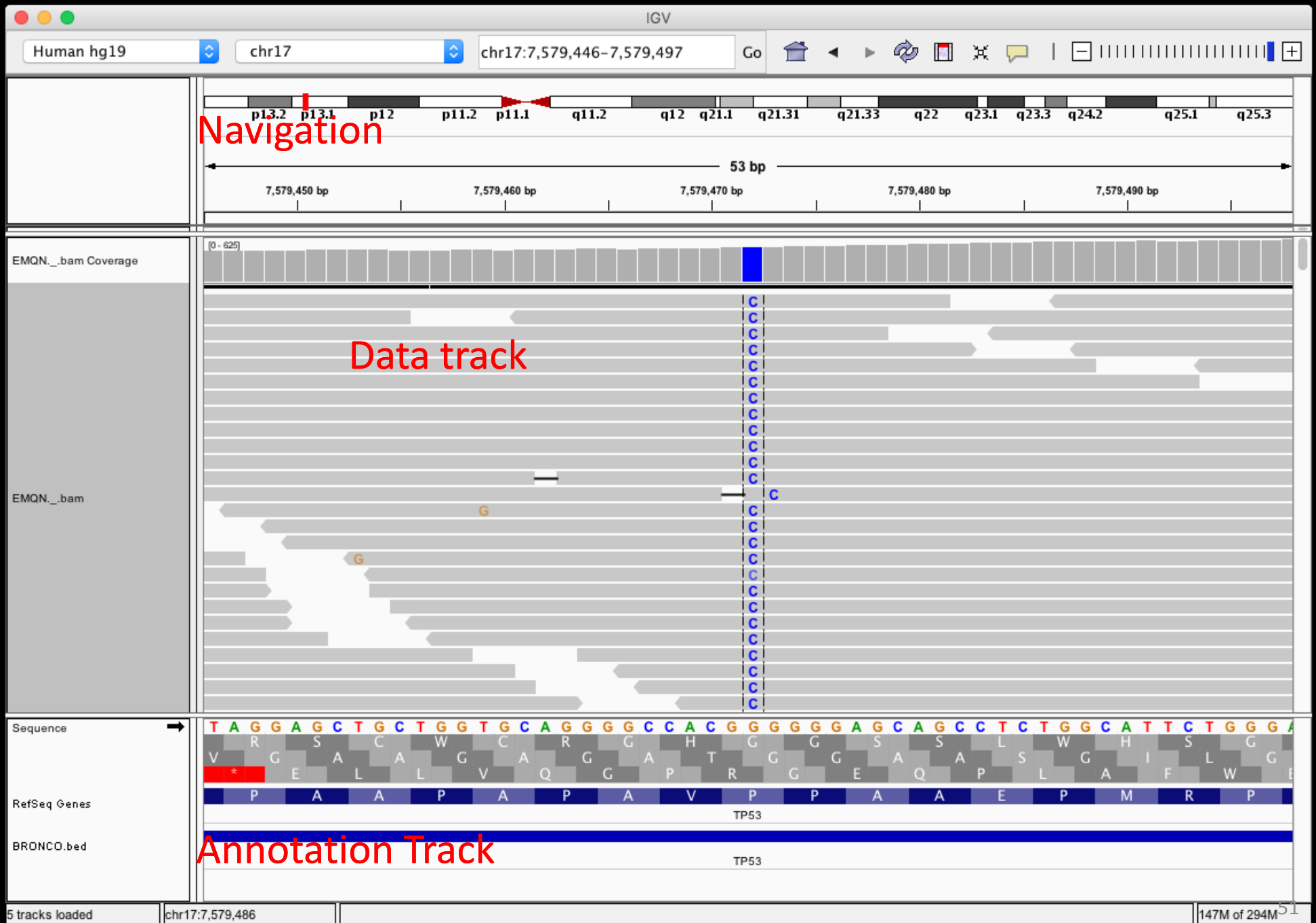


Navigation

Data track

Annotation Track

IGV



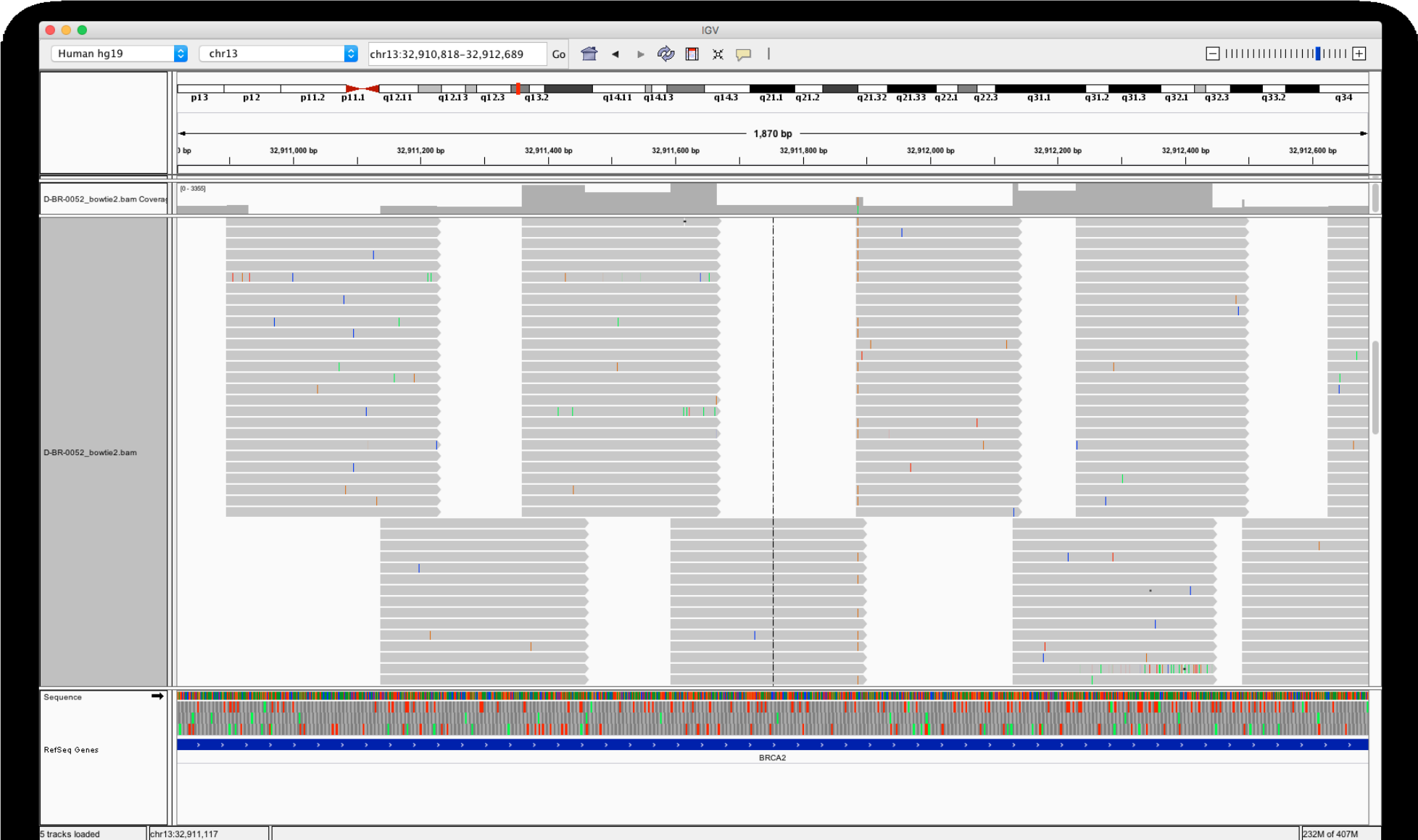
Alignment QC - Coverage statistics

- How many reads are aligned?
- How even is the overall coverage
- Average insert size
- How many reads come from the region of interest
 - On/Off target reads
 - Bed file – defines region of interest
- What is the average coverage
- How many % of target bases have at least X coverage

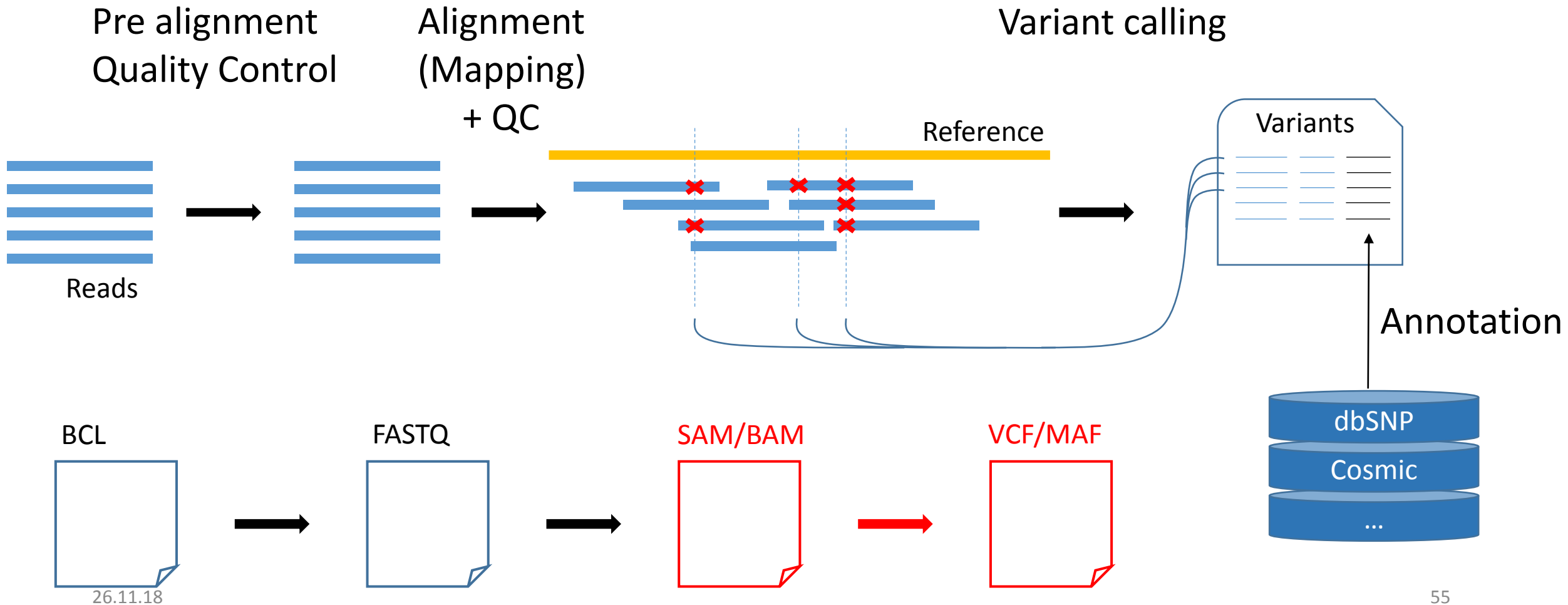
Alignment – Coverage statistics Picard-Tools

BAIT_SET	BRNO1077norm	BRNO1404norm	BRNO1503norm
TOTAL_READS	79842182	98157468	106336660
PF_READS	79842182	98157468	106336660
PF_UNIQUE_READS	69127998	87214990	95287834
PCT_PF_UQ_READS	0.865808	0.888521	0.896096
PF_UQ_READS_ALIGNED	68554192	86493216	94875098
PCT_PF_UQ_READS_ALIGNED	0.991699	0.991724	0.995669
ON_BAIT_BASES	4259279846	5132027582	5495198420
NEAR_BAIT_BASES	921266922	1284426047	1344776758
OFF_BAIT_BASES	1112205142	1324505871	1573330165
ON_TARGET_BASES	2765311894	3544263597	3773976192
PCT_SELECTED_BASES	0.823256	0.828896	0.812995
PCT_OFF_BAIT	0.176744	0.171104	0.187005
ON_BAIT_VS_SELECTED	0.822168	0.799823	0.803394
MEAN_BAIT_COVERAGE	93.968208	113.222763	121.235036
MEAN_TARGET_COVERAGE	61.008295	78.193523	83.261441
MEDIAN_TARGET_COVERAGE	52	68	72
PCT_TARGET_BASES_1X	0.971254	0.973571	0.974168
PCT_TARGET_BASES_2X	0.965912	0.969276	0.970293
PCT_TARGET_BASES_10X	0.928313	0.941296	0.945291
PCT_TARGET_BASES_20X	0.858879	0.896376	0.90591
PCT_TARGET_BASES_30X	0.760404	0.835876	0.852135
PCT_TARGET_BASES_40X	0.646449	0.760458	0.783717
PCT_TARGET_BASES_50X	0.530314	0.674041	0.702404

Alignment PCR library



Recap – Data analysis pipeline



DNA Variant calling

- Single Nucleotide Variants (SNV's) + short indels
 - Somatic/Germline
- Copy Number variants (CNV)
- Structural Variants

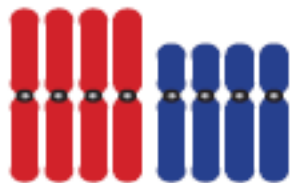
Normal diploid genome



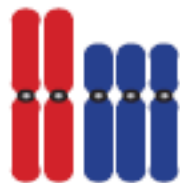
Normal diploid genome



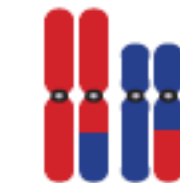
Polyploid



Aneuploid



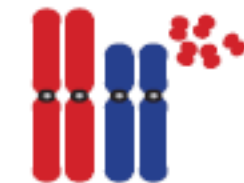
Reciprocal translocation



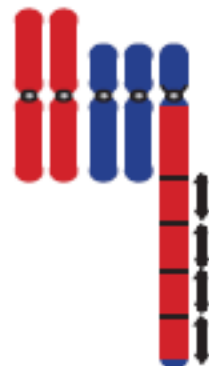
Non-reciprocal translocation



Amplification (double minutes)



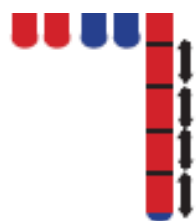
Amplification (HSR)



Amplification (distributed insertions)

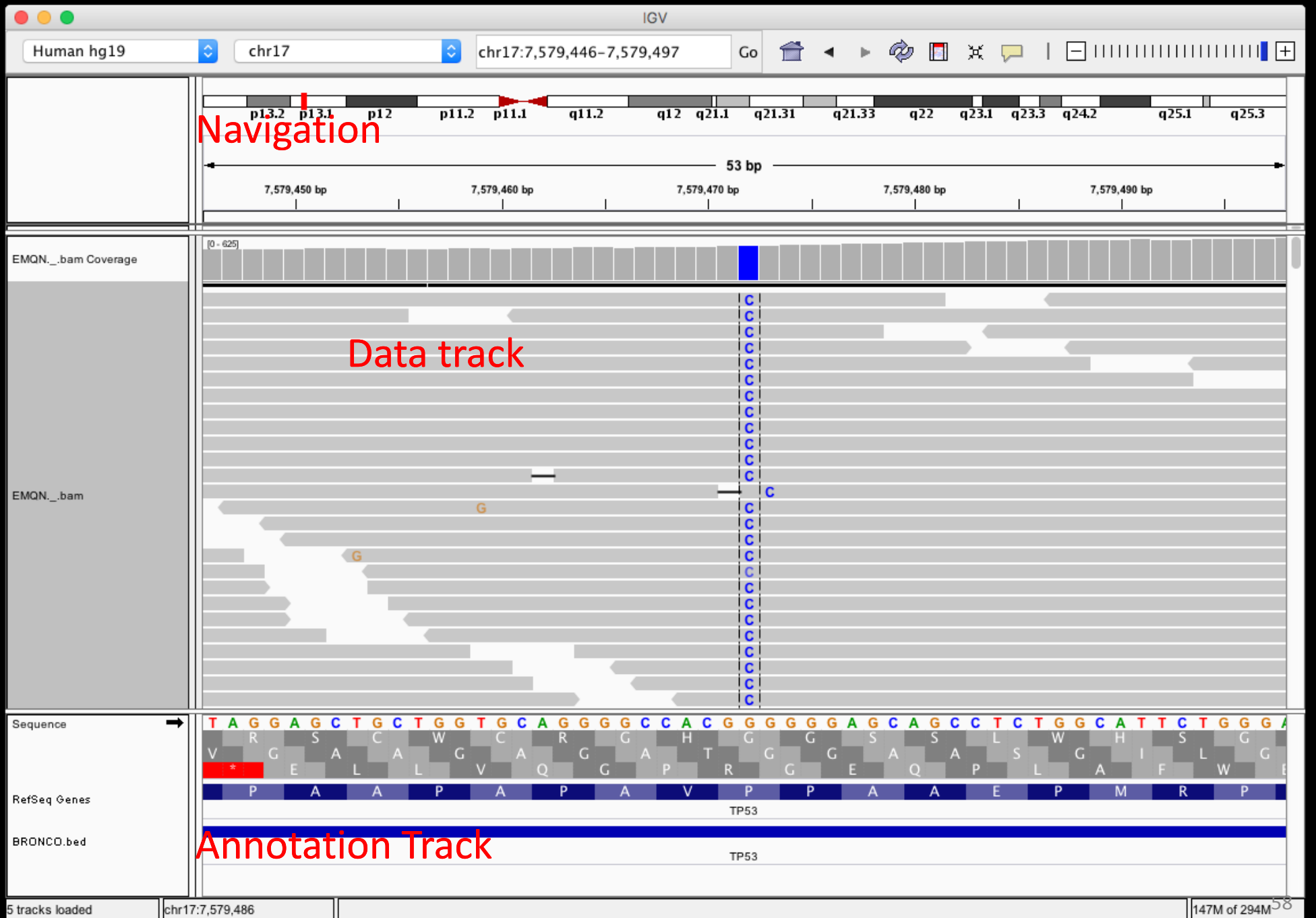


LOH (somatic recombination) LOH (duplication)

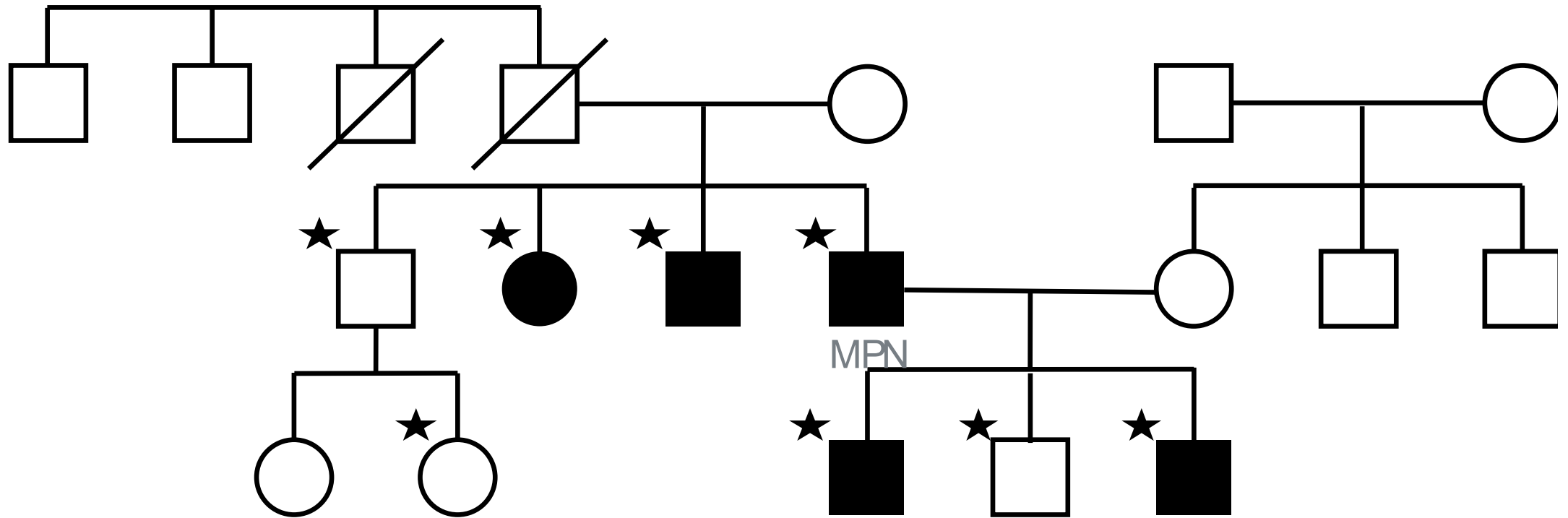


<https://doi.org/10.1038/ng1215>

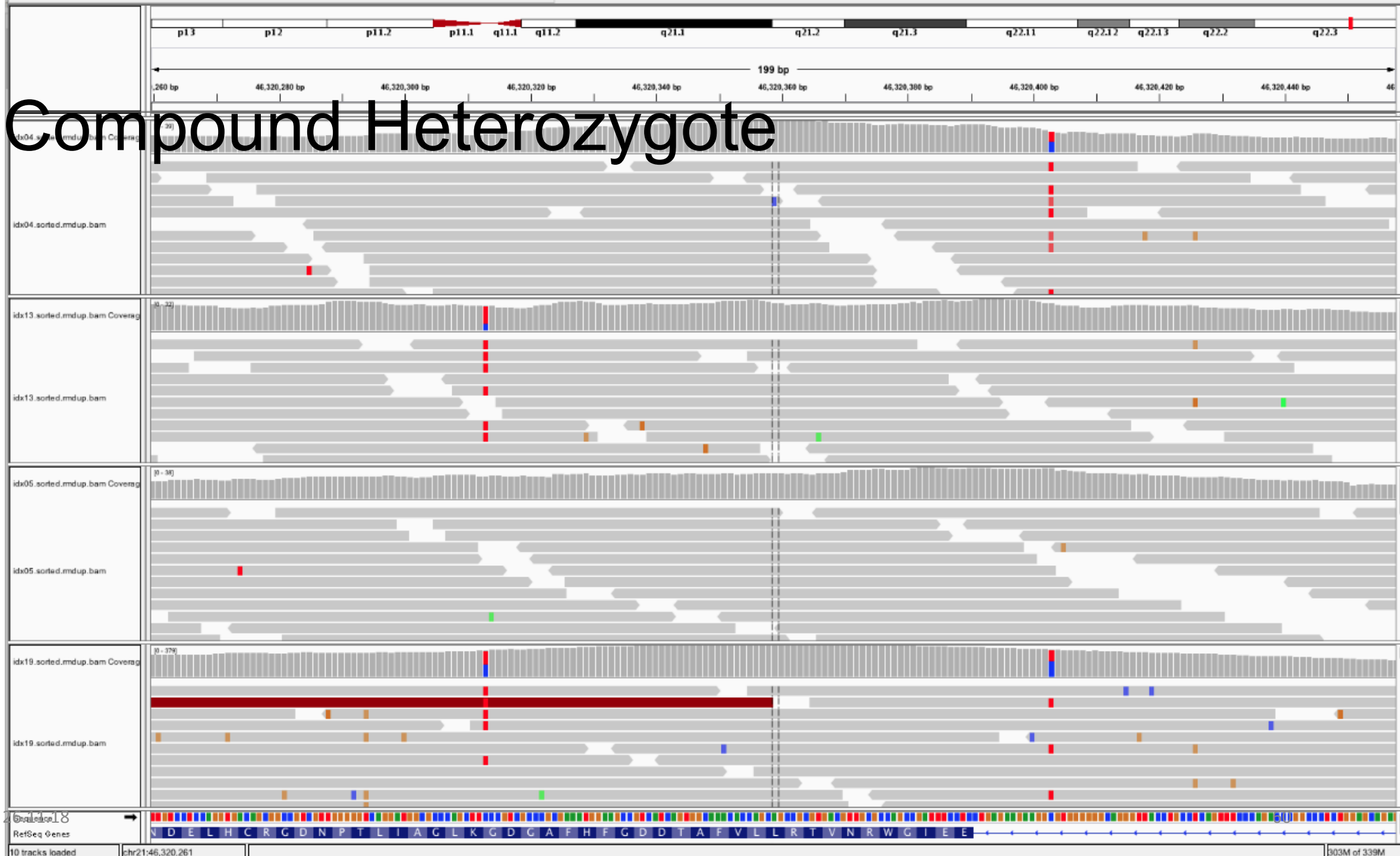
IGV



Disease causing germline mutations



Compound Heterozygote



chr21:46,320,261

10 tracks loaded

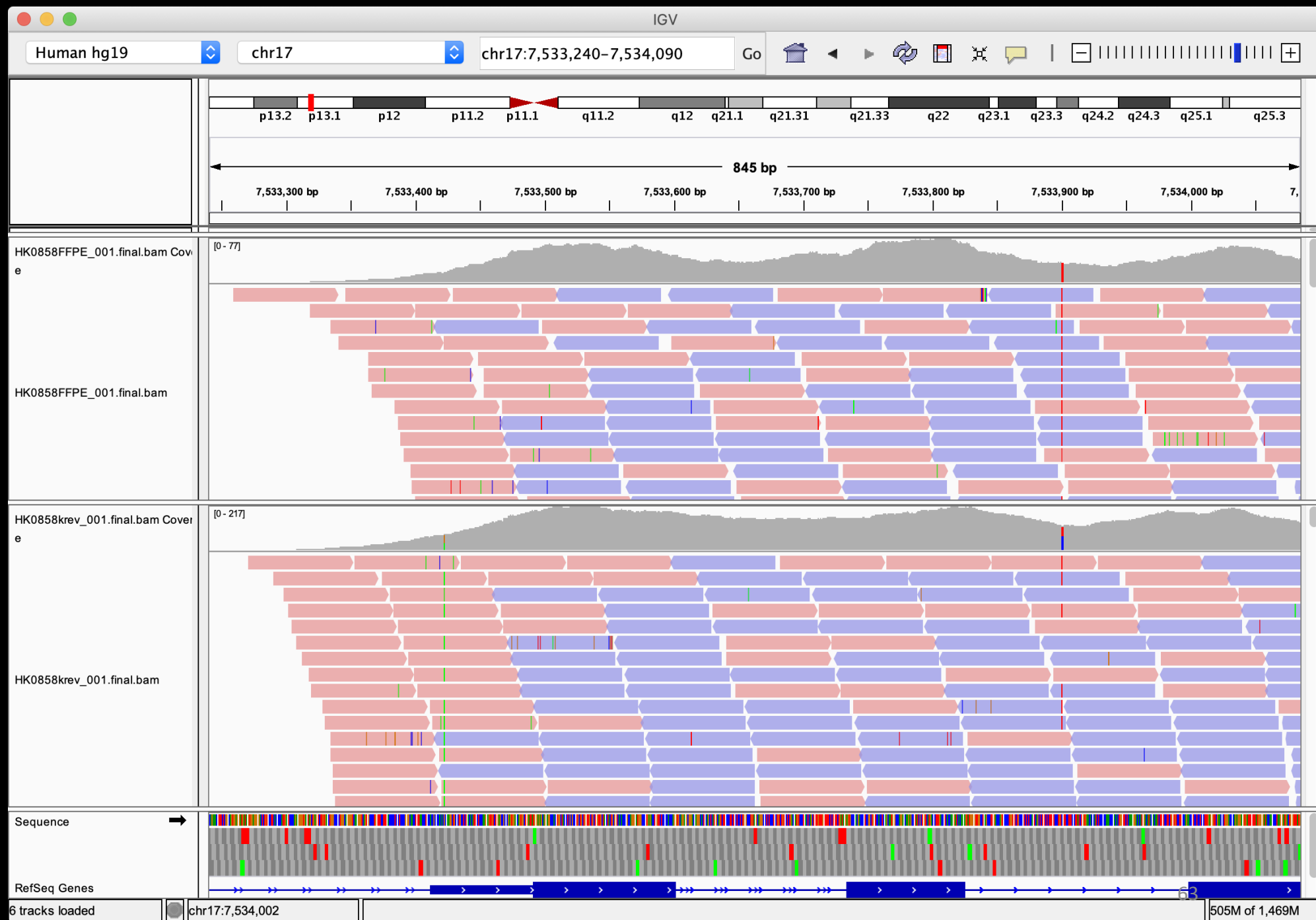
chr21:46,320,261

303M of 339M

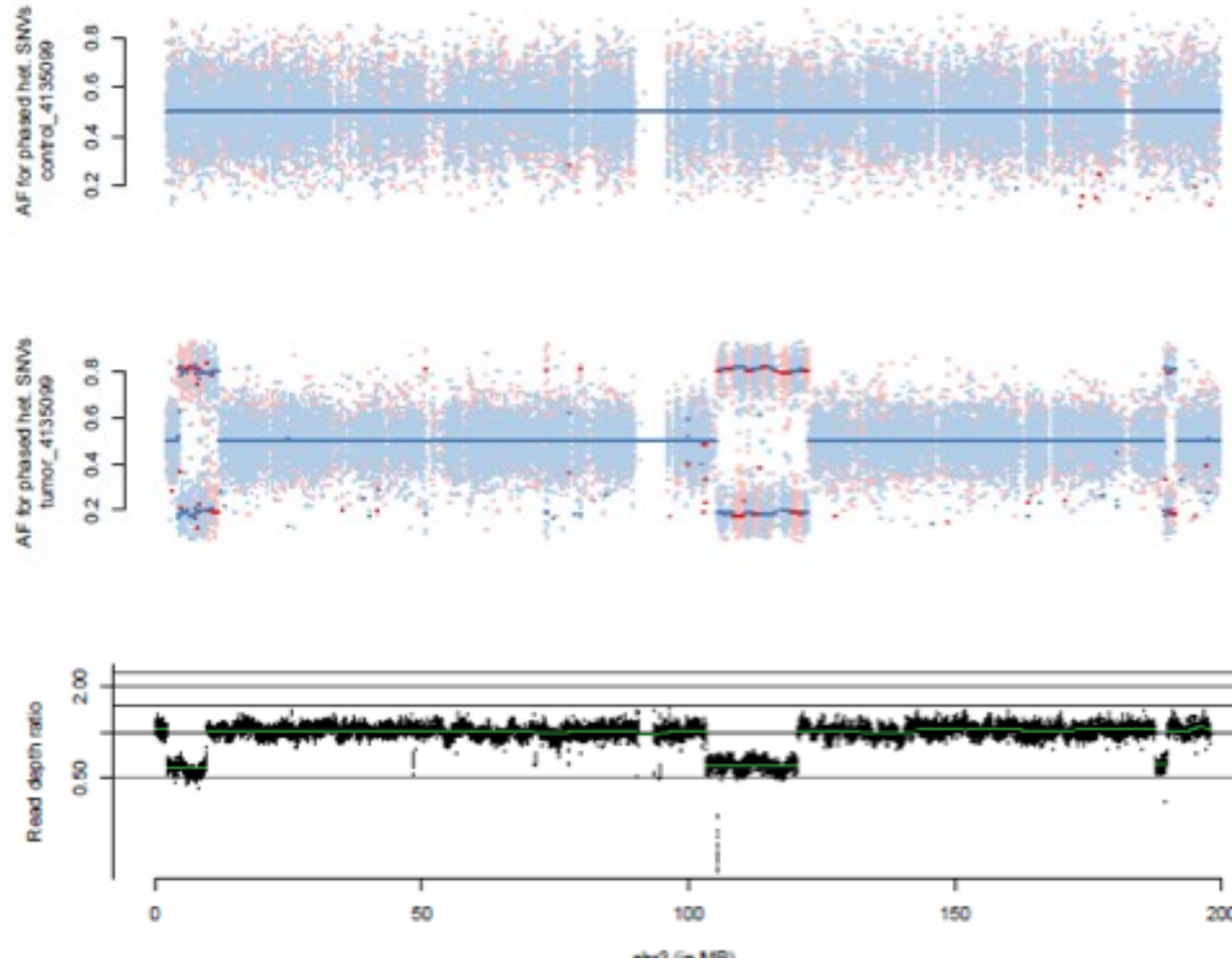
Somatic Variants

- Paired analysis
- Look for differences from the reference genome that do not occur in the control

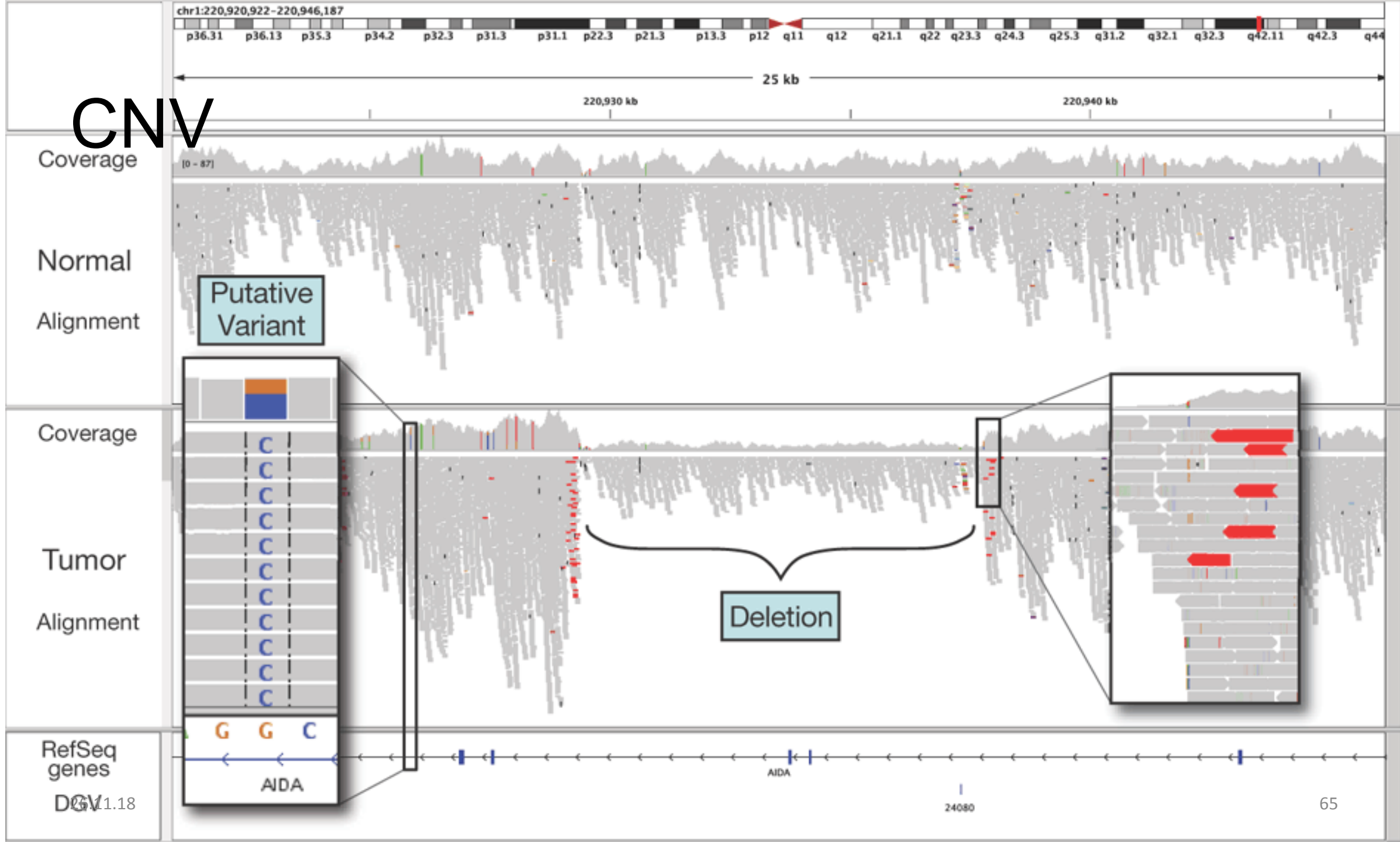
IGV inspector variant Soma



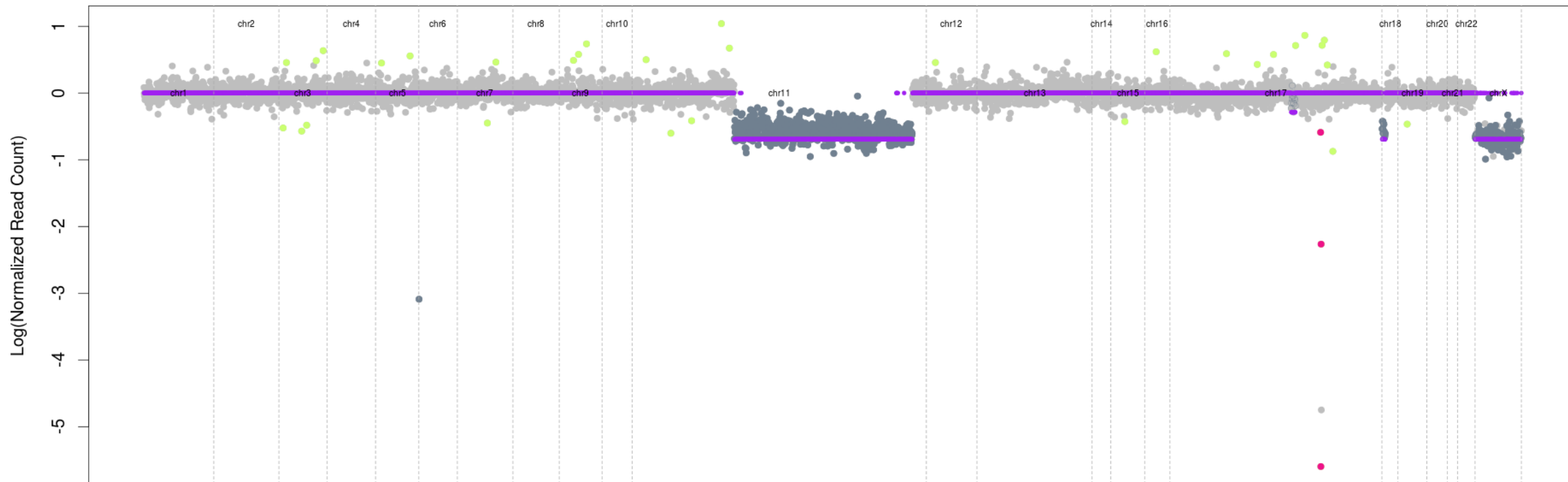
Phased het. SNPs germline vs tumor



CNV



CNV



IGV soft clips

Human hg19

General Tracks Mutations Charts **Alignments** Probes Proxy IonTorrent Advanced

Visibility range threshold (kb): Range at which alignments become visible

Downsample reads Max read count: per window size (bases):

Filter and shading options

Coverage allele-fraction threshold: Show coverage track

Filter duplicate reads Flag unmapped pairs

Filter vendor failed reads Show soft-clipped bases

Show center line Filter secondary alignments

Filter supplementary alignments Quality weight allele fraction

Mapping quality threshold:

Shade mismatched bases by quality: to

Filter alignments by read group

Flag insertions larger than: bases

Splice Junction Track Options

Show junction track Min flanking width: Min junction coverage:

Show flanking regions

Insert Size Options

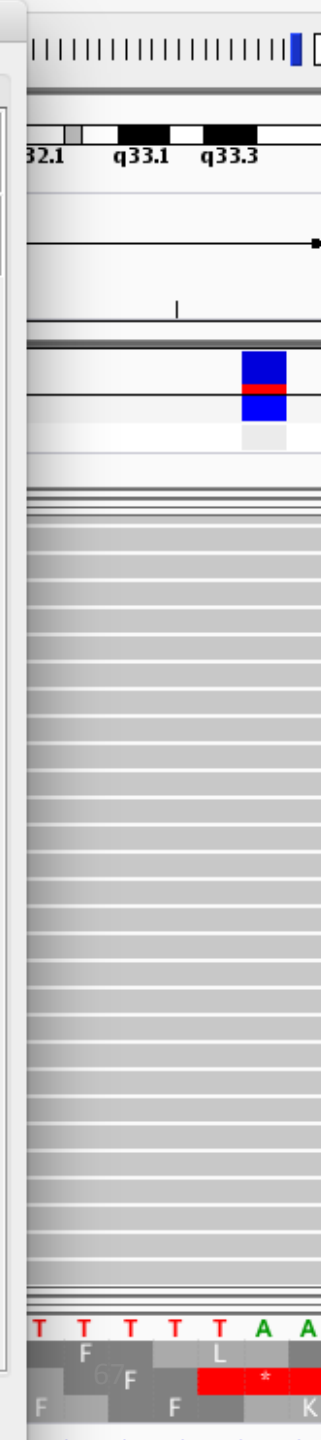
Defaults Minimum (bp): Compute Minimum (percentile):

Maximum (bp): Maximum (percentile):

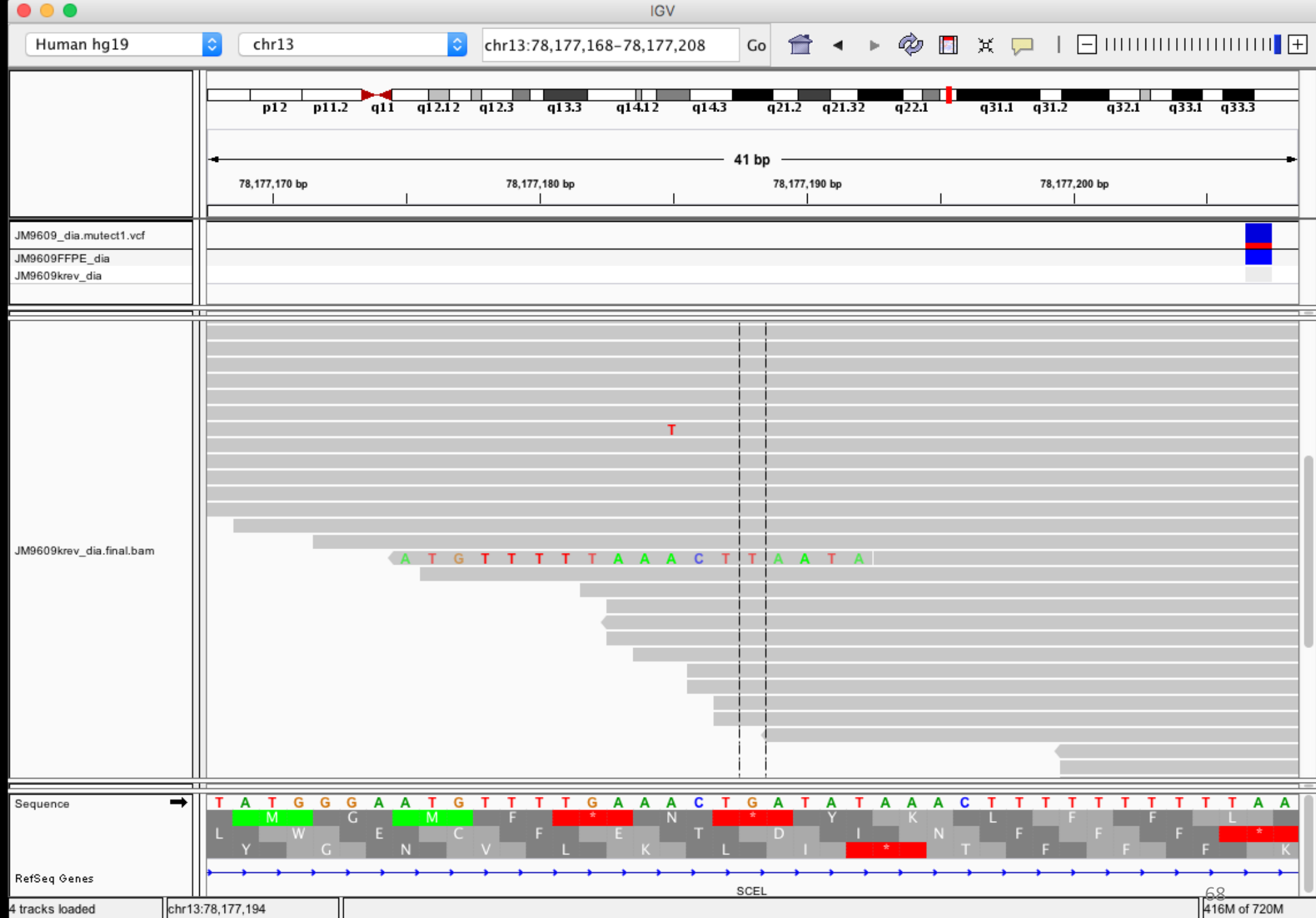
JM9609_dia.mutect1.vcf
JM9609FFPE_dia
JM9609krev_dia

JM9609krev_dia.final.bam

Sequence →

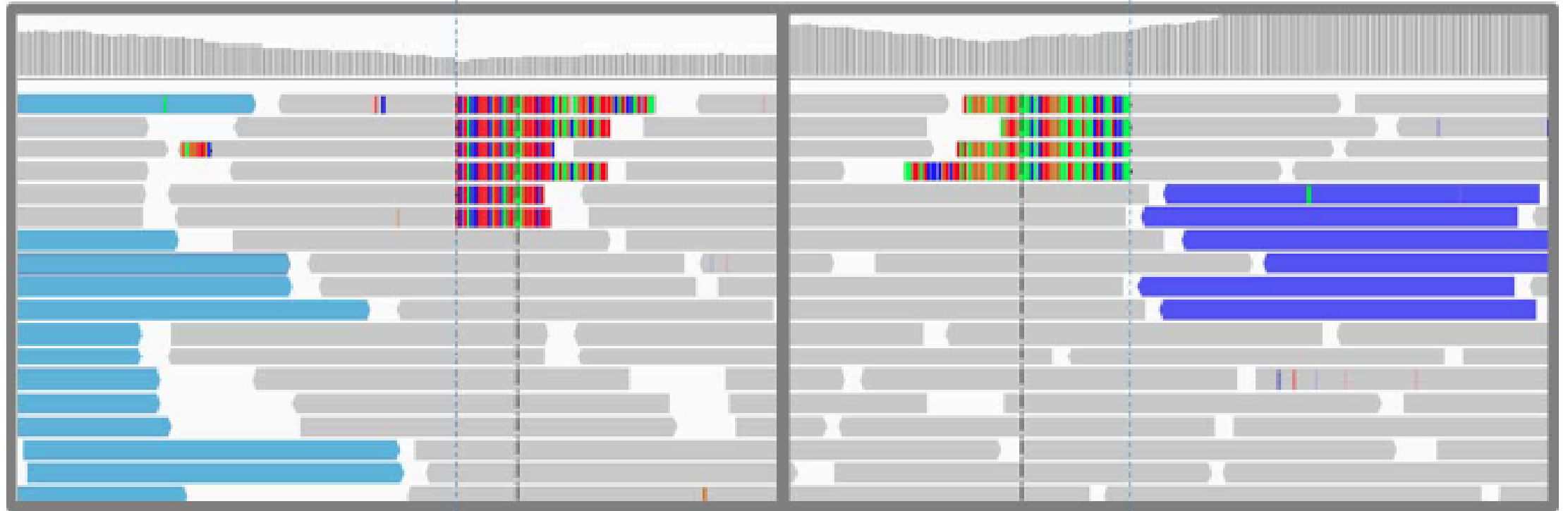


IGV soft clips



Breakpoint from Chr1: 176733607

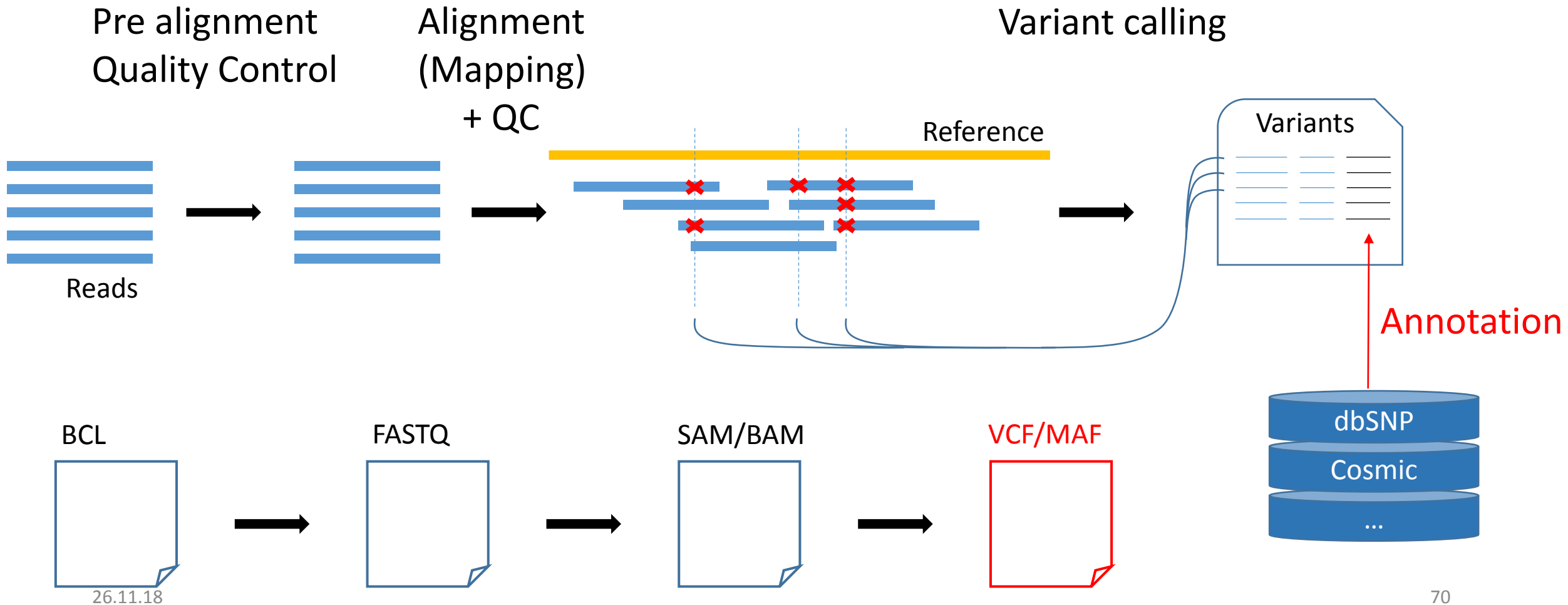
Breakpoint to Chr7: 127139813



A
Evidence from
discordant mapped
pair reads

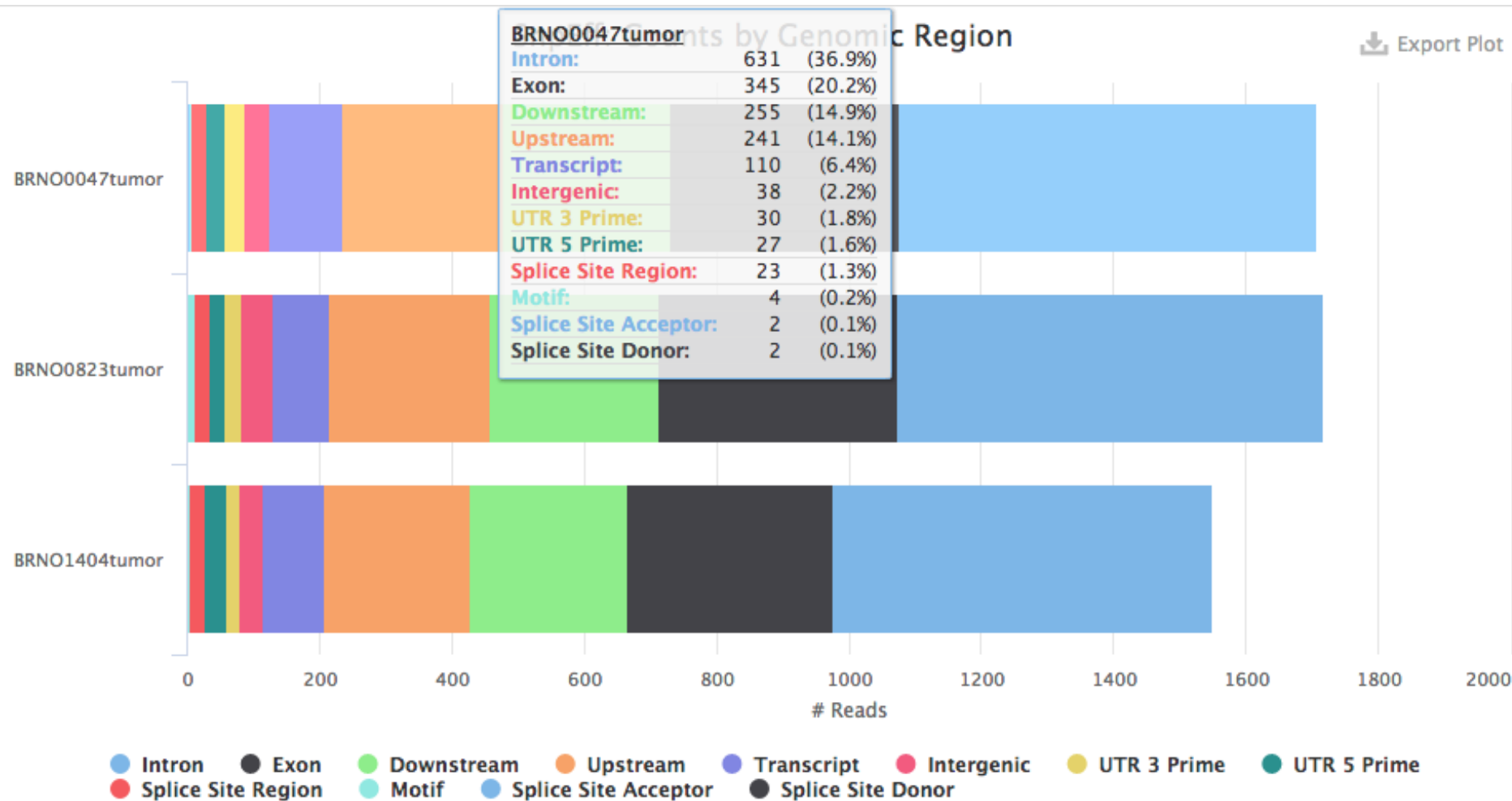
B
Evidence from soft clipped
bases

Recap – Data analysis pipeline



Variant annotation

MultiQC output



Variant annotation + report

BRCA
2 genes

REPORTED 1 BRCA1
BRCA MASTR™ Dx germline

Variant List - sorted by: PRED_CAT > PATHOGENICITY_CLASS > GENE

P...	Pat.	id	type	cod. cons.	gene	refSeqId	c.DNA	Protein	VF%	refSeq	altSeq	depth	SC
C	5	29	SNP	intronic	BRCA1	NM_007294	c.4485-63C>G		49.87			770	
C	5	27	SNP	intronic	BRCA1	NM_007294	c.4987-68A>G		53.63			716	
C	4	2	SNP	missense	BRCA1	NM_007294	c.2077G>A	p.Asp693Asn	51.43	GAC	AAC	525	
C	2	33	SNP	5'UTR	BRCA2	NM_000059	c.-26G>A		50.0			1020	
C	1	25	SNP	intronic	BRCA1	NM_007294	c.5152+66G>A		51.94			258	
C	1	18	INDEL	intronic	BRCA2	NM_000059	c.6841+80...		51.14			1095	
C		36	SNP	5'UTR	BRCA1	NM_007294	c.-134T>C		55.74			540	
C		26	SNP	intronic	BRCA1	NM_007294	c.5075-53C>T		55.41			231	
C		32	INDEL	intronic	BRCA1	NM_007294	c.548-58delT		51.32			793	
C		31	SNP	intronic	BRCA1	NM_007294	c.4097-141A...		46.97			264	
C		28	SNP	intronic	BRCA1	NM_007294	c.4987-92A>G		53.63			716	

OVERVIEW DETAILS COMMENTS VIEWER SIMILAR PATIENTS WARNINGS

reads: 716 DEPTH: 231 min, 6491 max
frequencies: 4/11 RUN: 24% ACCOUNT: 53.4% COMMUNITY: 53.4%

flagging: 13 (D, C, B, A) pred: In Report 2 Set To False + 0

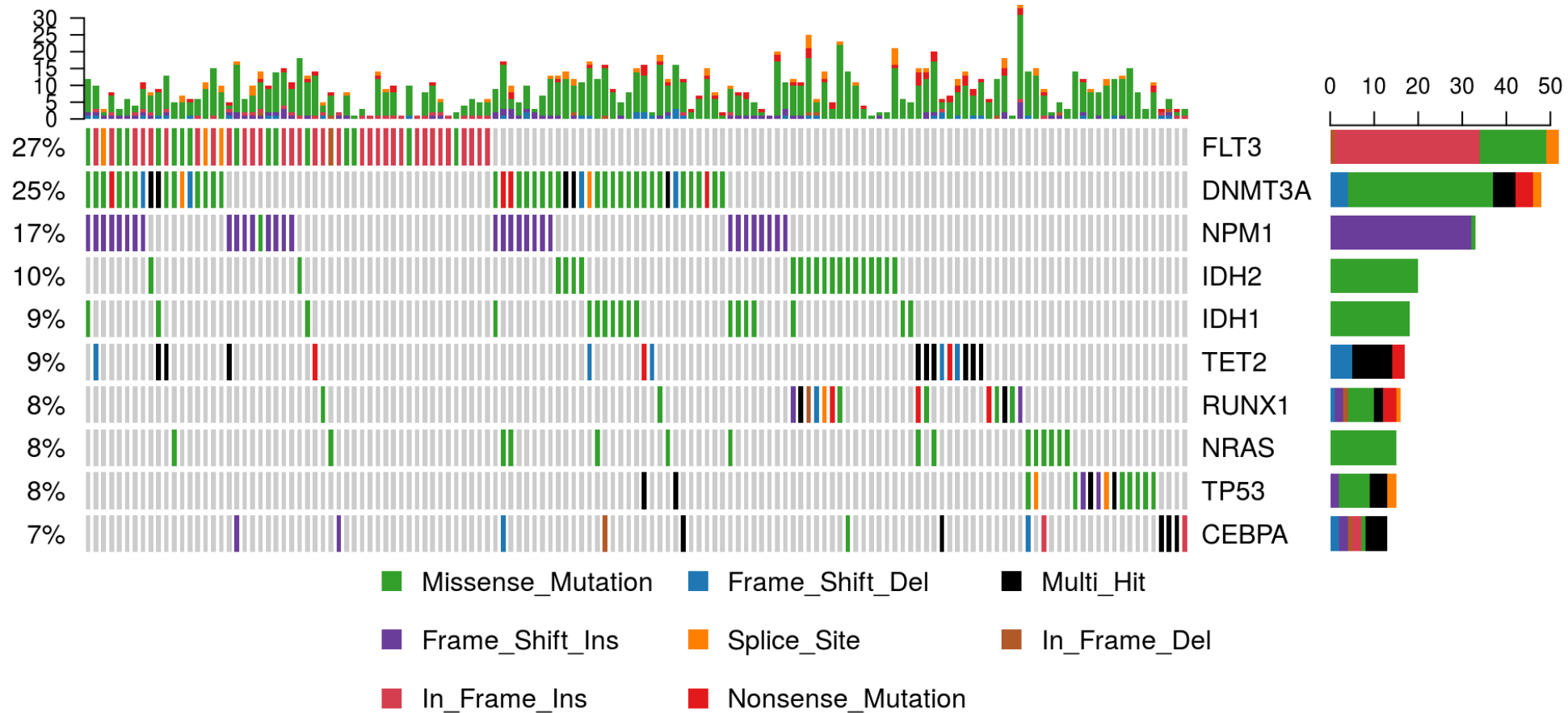
transcript: NM_007294
cDNA: c.4987-68A>G
rs number: rs8176234
SNP: 17-16
Intronic

SNP BRCA1
PolyPhen2 na
SIFT na
MutationTaster na
ESP5400 0.0
ExAC 0.0
cg99 0.26
G1000 0.35

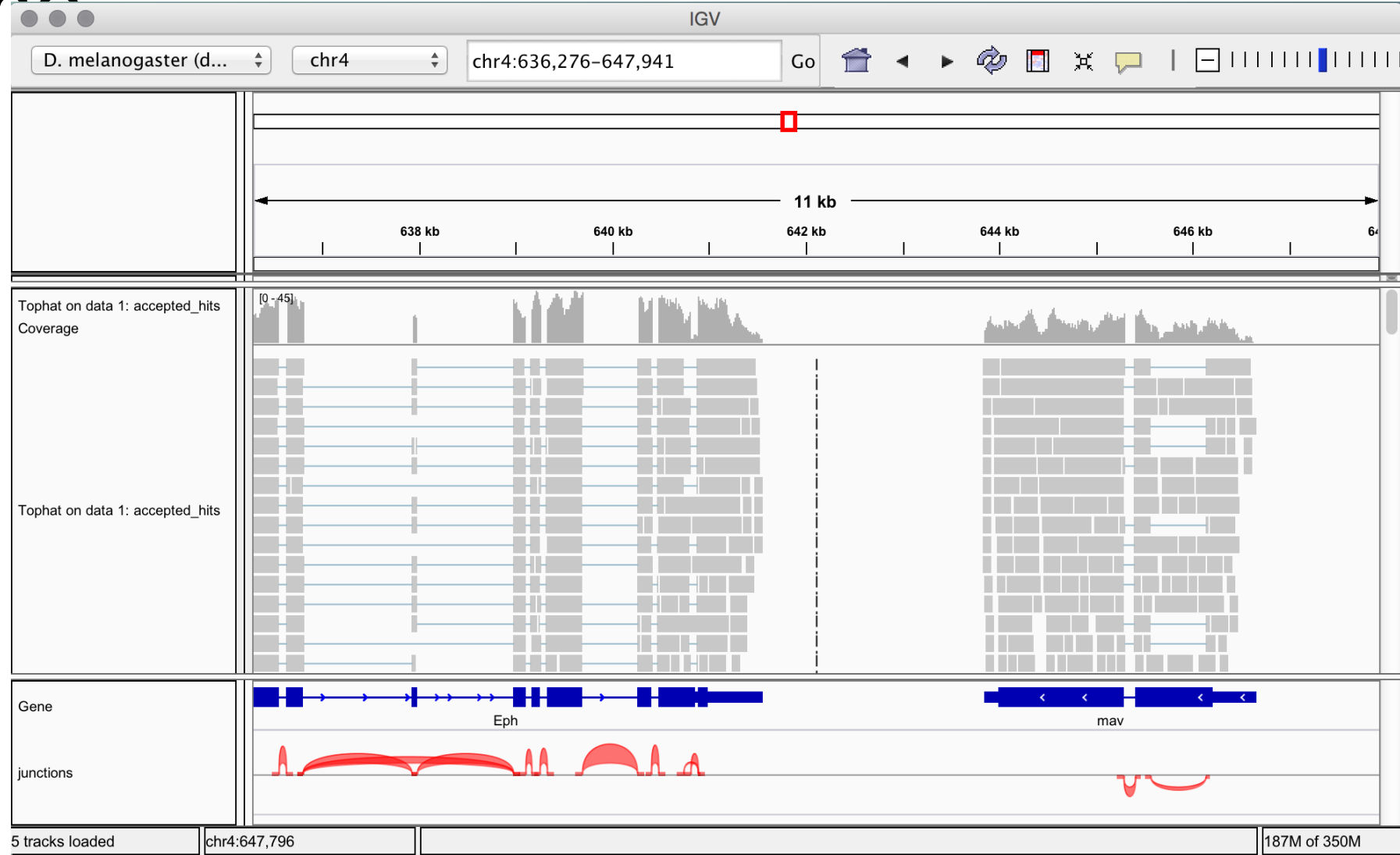
Values are scaled so that the most pathogenic scores are plotted towards the external circle.
ESP5400 & G1000 empty values are considered as 0.0

Variant annotation + report

Altered in 141 (73.44%) of 192 samples.



RNASeq reads aligned to the reference genome

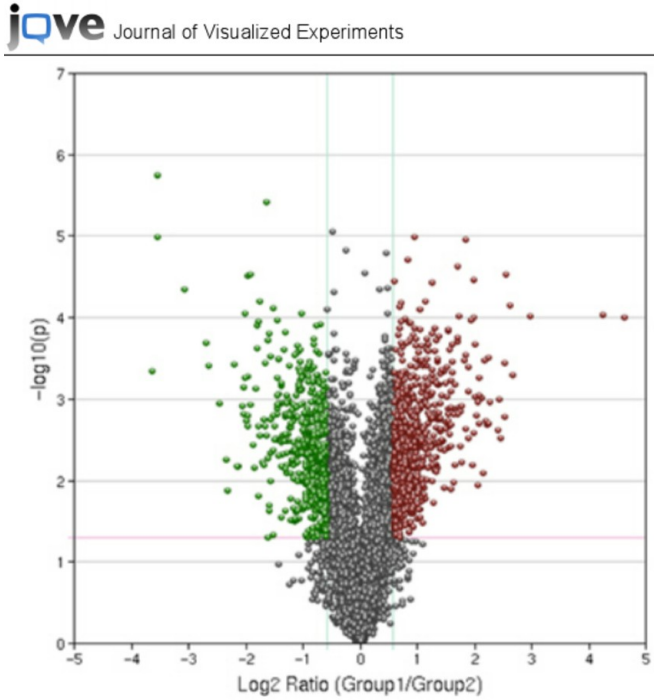


Annotation

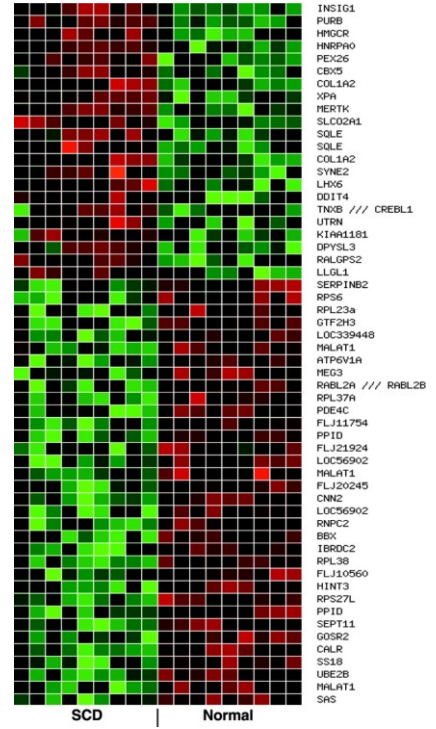
RNAseq pipeline

#	G1:MeOH	G1:MeOH	G1:MeOH	G2:R3G	G2:R3G	G2:R3G
#Feature	MeOH_Rep1	MeOH_Rep2	MeOH_Rep3	R3G_Rep1	R3G_Rep2	R3G_Rep3
LOC100288778	38	48	47	51	46	47
IQSEC3	0	0	0	0	0	0
CCDC77	51	51	51	40	40	39
B4GALNT3	4	4	3	6	6	11
WNK1	264	293	268	281	256	272
ERC1	55	55	68	83	57	49
LOC100292680	0	0	0	2	1	0
WNT5B	3	1	0	1	0	1
ADIPOR2	96	83	109	79	65	81
LRTM2	0	0	0	1	0	0
CACNA1C	5	1	2	7	3	4
CACNA1C-IT3	0	0	0	0	0	0
FKBP4	466	472	466	257	229	257
ITFG2	51	63	64	46	41	44
LOC100507424	5	1	2	0	1	4
RHNO1	73	82	74	61	58	66
TULP3	32	19	32	18	19	27
TEAD4	1	0	0	0	1	0
TSPAN9	0	0	1	1	1	0
PRMT8	1	0	1	0	0	0
CCND2	4440	4496	4694	2743	2739	2726

Feature counts
(+normalization)



Volcano Plot



Heat map

Combined approach



- Personalized medicine:**
1. The genetic changes in a person's cancer are discovered.
 2. Drugs that target these genetic changes are identified.
 3. The patient is treated and their response to therapy is monitored.

Taka away

- Terminology
- Interpreting Different QC metrics
- Interpreting NGS data visually
- Basic intuition (reads, alignments, references, variants)

Thank you for your attention