
Variant calling in Exome-seq data using Varscan

Elodie Girard

U900 INSERM - Mines ParisTech - Institut Curie

ITMO Roscoff – 11/19/2013

Dataset

- Public data: exome sequenced by the International HapMap Project
- Single-end reads of 100bp, Illumina Genome Analyzer Ix
- RNA-seq data of this exome available (Pickrell *et al.*, Nature, 2010)

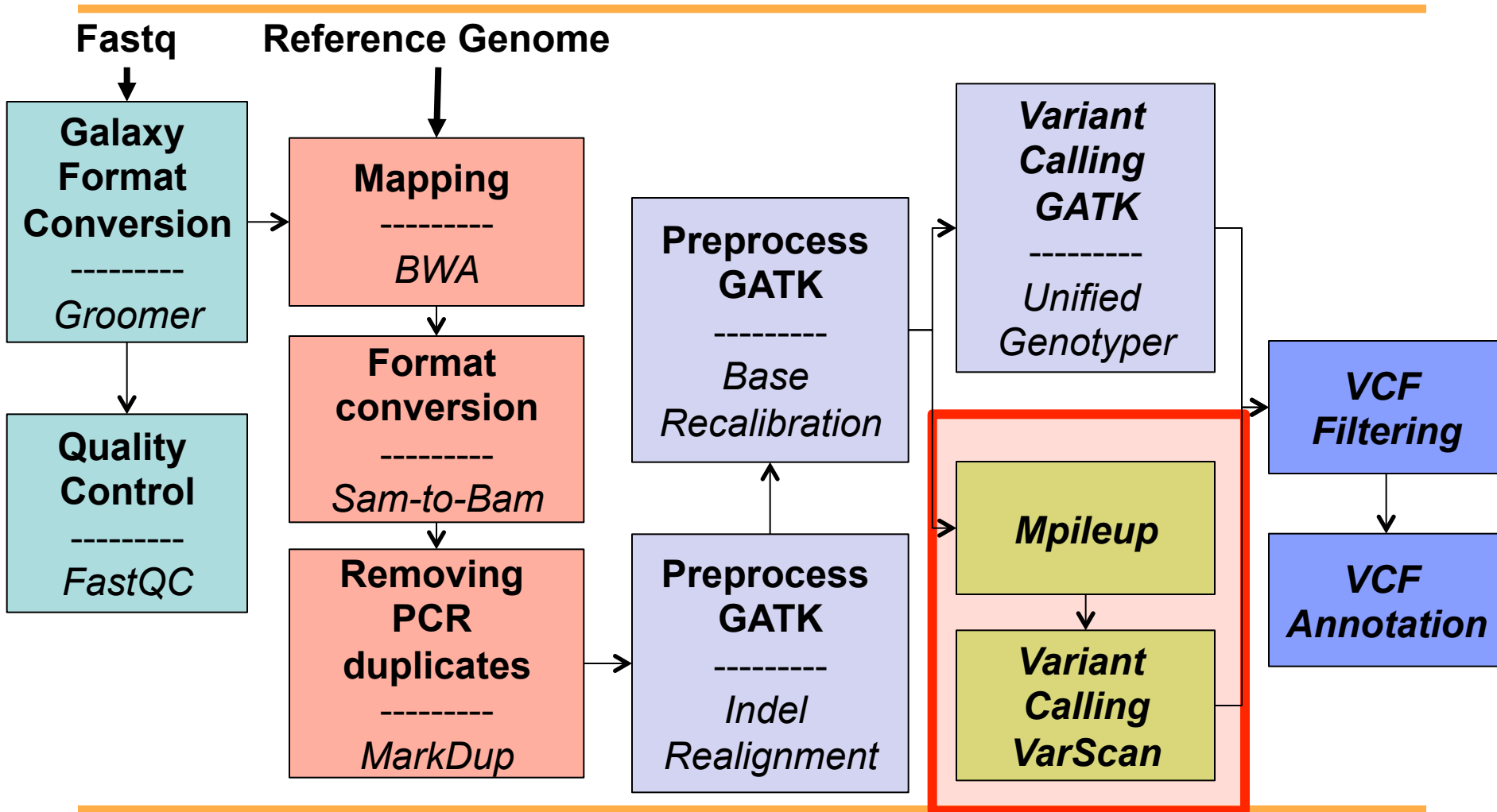
Objectives of the workshop:

- Variant calling, filtering and annotation in exome-seq data
- Observing the potential impact of these variants by looking at the corresponding RNA-seq data

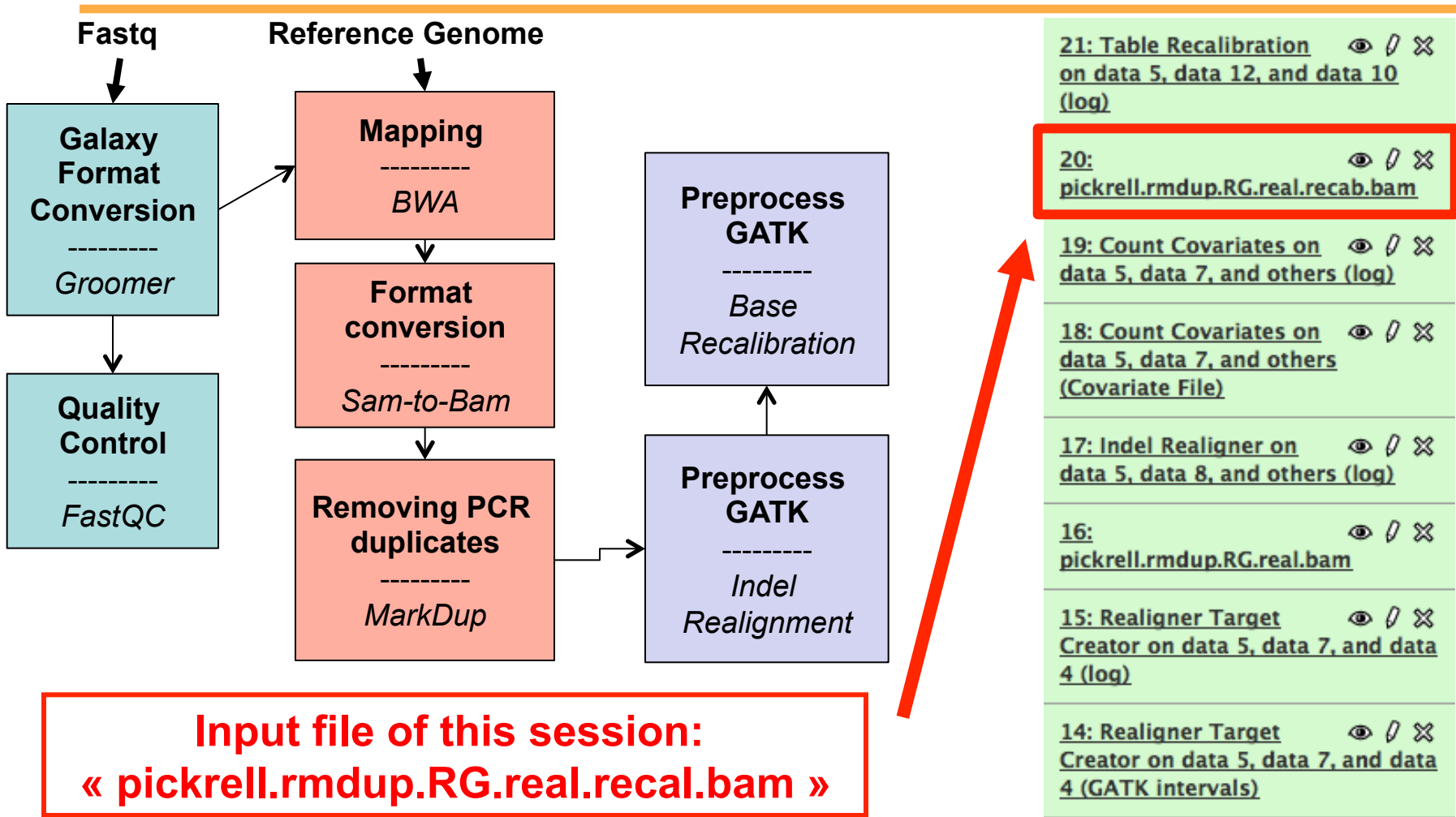
Objectif of this session :

- Variant detection using Varscan

Workflow

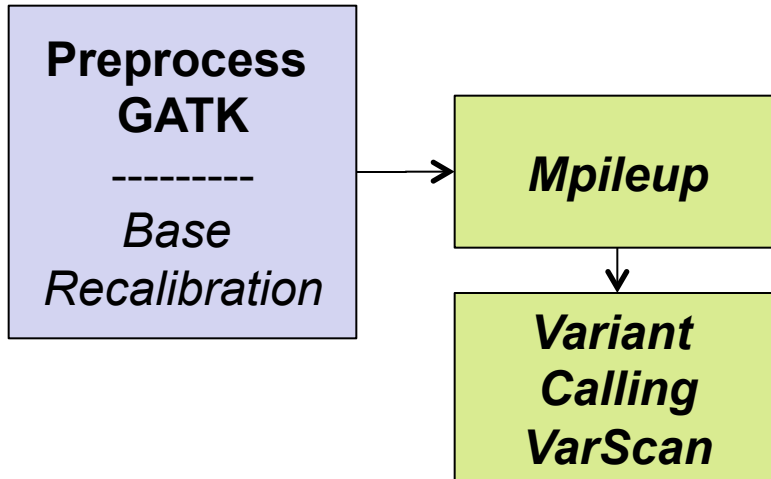


Galaxy: summary of the previous steps



Tools

All the tools are in the left panel: **Mardi 19 Détection de variants VarScan**



Mardi 19 Detection de variants VarScan

VarScan VarScan analysis.

POSTPROCESS TOOLS

Tag and merge multiple VarScan analysis

VarScan compare Compare two varscan results files (intersect / merge / unique).

VarScan Filter To filter a varscan input file.

Filter SNPs on same ref position

SAM TOOLS

MPileup SNP and indel caller

Mpileup

MPileup (version 0.0.1)

Choose the source for the reference list:

History

← 1 : History

BAM files

BAM file 1

BAM file:

14: pickrell.rmdup.RG.real.recab.bam

← 2 : Select the bam

Add new BAM file

Using reference file:

5: chr12.fa

← 3 : chr12.fa

Genotype Likelihood Computation:

Do not perform genotype likelihood computation

Set advanced options: ← 4 : Set Advanced options

Advanced

Minimum mapping quality for an alignment to be used:

20

← 5 : MapQ = 20

Minimum base quality for a base to be considered:

13

Only generate pileup in region:

chr12:1128!

← 6 : Region

chr12:112850000-113395000

Output per-sample Phred-scaled strand bias P-value:

Execute

← 7 : Execute

VarScan

- Mutation caller written in **Java** (no installation required) working with **Pileup files** of Targeted, **Exome**, and Whole-Genome sequencing data
- **Multi-platforms**: Illumina, SOLiD, Life/PGM, Roche/454
- Detection of different kinds of variants (SNVs/Indels) :
 - Germline variants in individual samples
 - Multi-sample variants **shared or private** in multi-sample datasets
- VarScan specificity is to be able to work with **Tumor/Normal pairs**:
 - Somatic and germline mutation, LOH events in tumor-normal pairs
 - Somatic copy number alterations (CNAs) in tumor-normal exome data

VarScan

- Most published variant callers use **Bayesian statistics** (a probabilistic framework) to detect variants and assess confidence in them (*e.g.*: GATK)
- VarScan uses a robust **heuristic/statistic** approach to call variants that meet desired thresholds for read depth, base quality, variant allele frequency, and statistical significance
- In Stead *et al.* (2013), they compared 3 different **somatic callers** : MuTect, Strelka, VarScan2
 - **VarScan2 performed best** overall with sequencing depths of 100x, 250x, 500x and 1000x required to accurately identify variants present at 10%, 5%, 2.5% and 1% respectively

VarScan

VarScan (version 2.0)

VarScan version:

VarScan V.2.2.8

Pileup File:

22: MPileup on data 6 and data 20

Type of analysis:

Pileup with Cns

← 2 : Pileup with Cns (calls SNVs + Indels)

Ignore variants with >90% support on one strand [Yes]:

Yes, I use this option

Output Format.:

VarScan format [tabular]

← 3 : Choose VarScan Tabulated format

Execute

← 4 : Execute

VarScan VCF Format

- 2 types: **VCF** and **Tabulated** Formats - **VarScan Specific**
 - VCF output might not work with some filtering or annotating tools
- **VarScan VCF format**: classic VCF header (#) but specific variant lines

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	GENO
chr12	250239	.	A	G	20	PASS	ADP=104;WT=0; HET=1 ;HOM=0; NC=0	GT:GQ:SDP:DP:RD:AD: FREQ: PVAL:RBQ:ABQ: RDF:RDR:ADF:ADR	0/1:153:111:104:61:43: 41,35%: 4,5644E-16:38:32: 48:13:36:7

ADP = Average per-sample depth of bases with Phred score = 20
WT = Number of samples called reference (wild-type)
HET = Number of samples called heterozygous-variant
HOM = Number of samples called homozygous-variant
NC = Number of samples not called

Useful in
multi-sample
studies

VarScan VCF Format

- 2 types: **VCF** and **Tabulated** Formats - **VarScan Specific**
 - VCF output might not work with some filtering or annotating tools
- VarScan VCF format:** classic VCF header (#) but specific variant lines

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	GENO
chr12	250239	.	A	G	20	PASS	ADP=104;WT=0; HET=1;HOM=0; NC=0	GT:GQ:SDP:DP:RD:AD: FREQ: PVAL:RBQ:ABQ: RDF:RDR:ADF:ADR	0/1:153:111: 104:61:43: 41,35%: 4,5644E-16:38:32: 48:13:36:7

GT=Genotype (1/1: Homozygous ; 0/1 : Heterozygous) / GQ= Genotype Quality
SDP= Raw Read Depth as reported by SAMtools
DP= Quality Read Depth of bases with Phred score >= 20
RD= Depth of reference-supporting bases
AD= Depth of variant-supporting bases
FREQ= Variant allele frequency

VarScan VCF Format

- 2 types: **VCF** and **Tabulated** Formats - **VarScan Specific**
 - VCF output might not work with some filtering or annotating tools
- **VarScan VCF format**: classic VCF header (#) but specific variant lines

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	GENO
chr12	250239	.	A	G	20	PASS	ADP=104;WT=0; HET=1;HOM=0; NC=0	GT:GQ:SDP:DP:RD:AD: FREQ: PVAL:RBQ:ABQ: RDF:RDR:ADF:ADR	0/1:153:111:104:61:43: 41.35%: 4,5644E-16:38:32: 48:13:36:7

PVAL= P-value from Fisher's Exact Test (**not computed here : default value**)
RBQ= Average quality of reference-supporting bases
ABQ= Average quality of variant-supporting bases
RDF / RDR= Depth of reference-supporting bases on forward/reverse strand
ADF / ADR = Depth of variant-supporting bases on forward/reverse strand

VarScan VCF Format

- 2 types: **VCF** and **Tabulated** Formats - **VarScan Specific**
 - VCF output might not work with some filtering or annotating tools
- **VarScan VCF format**: classic VCF header (#) but specific variant lines

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	GENO
chr12	250239	.	A	G	20	PASS	ADP=104;WT=0; HET=1;HOM=0; NC=0	GT:GQ:SDP:DP:RD:AD: FREQ: PVAL:RBQ:ABQ: RDF:RDR:ADF:ADR	0/1:153:111:104:61:43: 41.35%: 4,5644E-16:38:32: 48:13:36:7

PVAL= P-value from Fisher's Exact Test (**not computed here : default value**)
RBQ= Average quality of reference-supporting bases
ABQ= Average quality of variant-supporting bases
RDF / RDR= Depth of reference-supporting bases on forward/reverse strand
ADF / ADR = Depth of variant-supporting bases on forward/reverse strand

VarScan Tabulated Format

- 2 types: **VCF** and **Tabulated** Formats - **VarScan Specific**
 - Tabulated output works with other VarScan Tools
 - By default on Galaxy, VarScan outputs a line for each base covered by the selected minimum coverage **even if there is no alternative variant**
- **Pre-process**: use « **VarScan Filter** » to keep only variants

VarScan Filter (version 1.0.0)

VarScan File:

25: [VarScan] Results

Minimum Coverage:

8

Minimum read depth at a position to make a call [8].

Execute

← **2 : Execute**

1 : Select the VarScan Results and leave the default parameters

VarScan Tabulated Format

```

25: [VarScan] Results
26,055 lines
format: tabular, database: hg19
1 2 3 4 5 6
Chrom Position Ref Cons Reads1 Reads2
chr12 112883809 T T 8 0
    
```

VarScan Filter



```

29: [VarScan] VarScan File
31 lines
format: tabular, database: hg19
1 2 3 4 5 6
Chrom Position Ref Cons Reads1 Reads2
chr12 112887918 T W 13 2
    
```

Chrom	Position	Ref	Cons	Reads1	Reads2	VarFreq	Strands 1	Strands 2	Qual1	Qual2	Pvalue	Map Qual1	Map Qual2	R1 +	R1 -	R2 +	R2 -	Rs2 +	Rs2 -	Alt
chr12	113348849	C	Y	31	30	49.18%	2	2	27	27	0.98	1	1	19	12	25	5			T
chr12	113354329	G	R	72	2	2.70%	2	2	31	26	0.98	1	1	48	24	1	1			A
chr12	113357193	G	A	2	72	97.30%	1	2	28	24	0.98	1	1	2	0	45	27			A
chr12	113357209	G	A	0	77	100%	0	2	0	29	0.98	0	1	0	0	51	26			A

Cons : Consensus Genotype of Variant Called (IUPAC code):

M -> A or C	Y -> C or T	D -> A or G or T	W -> A or T	V -> A or C or G
R -> A or G	K -> G or T	B -> C or G or T	S -> C or G	H -> A or C or T

Variants visualization with IGV

20: pickrell.rmdup.RG.real.recab.bam
505.5 MB

format: bam, database: hg19

display at Ensembl Current
display with IGV web current local
display in IGB Local web

Binary bam alignments file

1: Click on the name to open the options

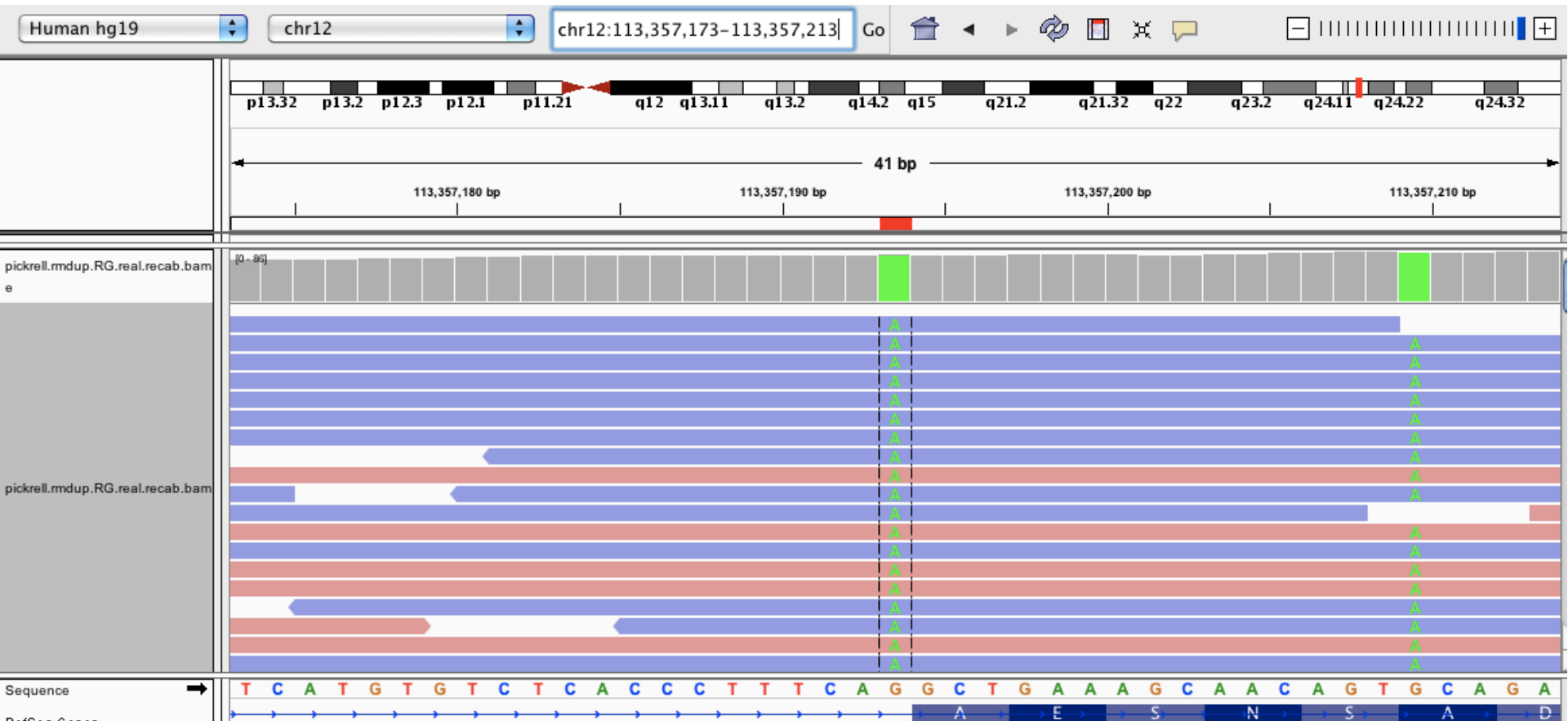
2: Click on « Local » to open it on IGV (if already open, otherwise « Web current »)

- Look at those two variants:

Chrom	Position	Ref	Cons	Reads 1	Reads 2	VarFreq	Strands 1	Strands 2	Qual 1	Qual2	Pval	Map Qual1	Map Qual2	R1 +	R1 -	R2 +	Rs2 -	Alt
chr12	113357193	G	A	2	72	97.30%	1	2	28	24	0.98	1	1	2	0	45	27	A
chr12	112888239	C	Y	52	56	51.85%	2	2	24	28	0.98	1	1	20	32	24	32	T

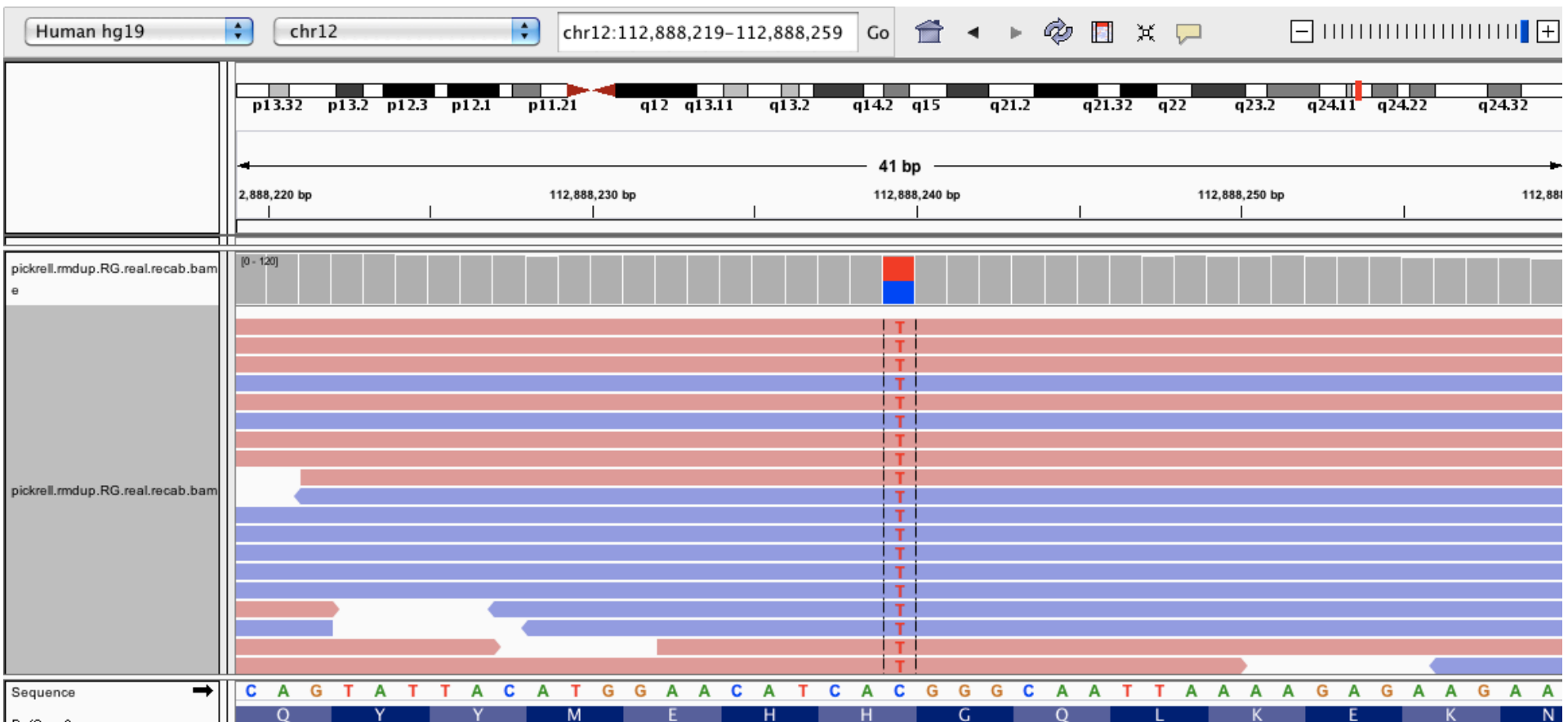
Variants visualization with IGV

Chrom	Position	Ref	Cons	Reads 1	Reads 2	VarFreq	Strands 1	Strands 2	Qual 1	Qual2	Pval	Map Qual1	Map Qual2	R1 +	R1 -	R2 +	Rs2 -	Alt
chr12	113357193	G	A	2	72	97.30%	1	2	28	24	0.98	1	1	2	0	45	27	A



Variants visualization with IGV

Chrom	Position	Ref	Cons	Reads 1	Reads 2	VarFreq	Strands 1	Strands 2	Qual 1	Qual2	Pval	Map Qual1	Map Qual2	R1 +	R1 -	R2 +	Rs2 -	Alt
chr12	112888239	C	Y	52	56	51.85%	2	2	24	28	0.98	1	1	20	32	24	32	T



Next Step

