

# BCGES short courses, session 4: Multi-sample calling, VCFtools and variant annotation

## Contents

<b>1</b>	<b>Multi-sample variant calling (about 1h30)</b>	<b>2</b>
1.1	Installing the reference fasta file (5 minutes) . . . . .	2
1.2	Samtools pileup and mpileup (30 minutes) . . . . .	2
1.3	GATK v2.0 Unified Genotyper (30 minutes) . . . . .	3
1.4	GATK v3.0 HaplotypeCaller version of multi-sample calling (30 minutes) . . . . .	3
<b>2</b>	<b>VCFtools (about 45 minutes)</b>	<b>5</b>
<b>3</b>	<b>Variant annotation (about 45 minutes)</b>	<b>6</b>
3.1	Running the variant effect predictor (VEP) . . . . .	6
3.2	Interpreting the output of the VEP . . . . .	6
3.3	ANNOVAR as an alternative . . . . .	6

# 1 Multi-sample variant calling (about 1h30)

Multi-sample calling has several advantages:

- Borrow information from well covered samples to identify indels and other small variants, in turn improving the calling for other samples.
- Uniform calls across samples, which improves sample comparability.
- Computational improvements associated with the joint management of potentially large collection of samples.

Both **GATK** and **samtools** provide options for multi-sample calling and we will explore how these scripts are setup today.

## 1.1 Installing the reference fasta file (5 minutes)

As far as I can tell, and unlike **samtools**, **GATK** requires that the fasta file matches the sequence dictionary as specified in the header of the BAM file.

Hence a required element for variant calling with human data is the reference sequence in fasta format. It is also necessary to index this sequence in order to provide a “dictionary” of what sequences are present, and how long they are. This file is large and is therefore not on the **github** of the class. It can be downloaded using a **ftp** site and the **wget** utility

```
wget ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/human_g1k_v37.fasta.gz
gunzip human_g1k_v37.fasta.gz
samtools faidx human_g1k_v37.fasta
```

These three steps download the compressed reference, unzip it and generate an index for that file. Please save it in your **data** folder (at the same level as all the “session” directories), because my scripts will assume that this is where this file can be found.

## 1.2 Samtools pileup and mpileup (30 minutes)

**samtools** is less feature-rich than **GATK** but it remains a dependable, solid piece of software that can provide very reliable callsets. It has been used for large scale projects like UK10K, and there is no reason to doubt the quality of these calls. Calling samples requires the presence of the reference sequence, in fasta format. The compressed file **chr21.fasta** is available on the **github** site to make it easy. It is however necessary to uncompress this file using **gunzip**.

Once this is done, the script **multi\_sample\_calling\_samtools.sh** describes the steps that generate a VCF file using **samtools**.

**Exercise:** Have a look at the scripts and make sure you can run them (in the simplest case just calling **sh script\_name.sh** might do). They may need minor tweaking to point to the appropriate locations. It is important to understand all the options that are used.

Now we can look at the output and make sure we understand what is there. For example, using a simple **awk** script, let us look at position 21:43847096:

```
awk '{if ( ($2 == "POS") || ($2 == 43847096)) print $1,$2,$4,$5,$9, $10, $15}' \
    results/UBASH31A_samtools.vcf

## #CHROM POS REF ALT FORMAT HG00250 HG00255
## 21 43847096 G A GT:PL:DP:GQ 0/1:119,0,54:9:70 0/0:0,27,169:9:37
```

Make sure that you understand the meaning of all the fields in the output vcf file. Of note, you can see the total depth for each sample in the output VCF file, but I do not think that an option exists in **samtools** to show the depth per allele, unlike **GATK** (see below).

Another thing that is important is to understand what **samtools mpileup** actually does, and what format the calls are based on. I use the command below routinely to make sure that my calls are not buggy, that I have the right sample... that sort of things. It is usually referred to as MAQ pileup.

```
samtools mpileup -f ../data/reference/chr21.fasta -r 21:43847096-43847096 \
    ../data/BAM_files/HG00250.mapped.ILLUMINA.bwa.GBR.exome.20121211.bam

## 21 43847096 G 9 ,aaa,Aa,a B@BD<@CA2
```

In class we will discuss this format and how to interpret this output.

### 1.3 GATK v2.0 Unified Genotyper (30 minutes)

The GATK Unified Genotyper (UG) has been the standard approach for GATK for a while, and still delivers high quality and reliable calls. For reasons that will be outlined in the next subsection, GATK has now transitioned toward a more powerful system called HaplotypeCaller. But there is nothing wrong with still using the UG module of GATK, and this subsection shows how to run such tools in the context of multi-sample calling.

The Unified Genotyper (UG) version of the GATK pipeline is implemented in the script called `multi_sample_calling_GATK-UG`. The key output of this script is the file `results/GATK_msamplcalling.vcf`. Note that GATK sorts the samples by alphabetical order, so the order may not match what was specified in the input BAM list, and it may also not match the VCF output of `samtools`. However, one important details about these scripts: for these to run on my computer, I had to specify the path to the GATK java file. You will need to do the same thing, and replace my path with `/usr/local/gatk/GenomeAnalysisTK.jar`. Similarly, the path to the fasta file of the human genome may need to be updated. My scripts assume that you save this file in the `data` folder.

**Exercise:** As for the `samtools` subsection, make sure you can run these scripts, tweak them a bit and understand the options that I used. Don't hesitate to look at the [manual page](#) for the UG tool to help you out.

As before for `samtools`, we can look at the output VCF file and make sure we understand what is shown:

```
awk '{if ( ($2 == "POS") || ($2 == 43847096)) print $1,$2,$4,$5, $9, $11, $16}' \
    results/GATK_msamplcalling.vcf

## #CHROM POS REF ALT FORMAT HG00250 HG00255
## 21 43847096 G A GT:AD:DP:GQ:PL 0/1:3,6:9:78:178,0,78 0/0:10,0:10:27:0,27,358
```

**Exercise:** Make sure you fully understand the output from GATK. It differs from GATK in terms of what is shown, but there are also differences in terms of content (look for example at the sequencing depth for HG00255 at this locus). In general, can you make sense of why the read depth would differ? Can you find a GATK update related to this point in the [release notes](#) of the version 3.2? Observe also that the genotype likelihoods also differ significantly, even with the same input data.

### 1.4 GATK v3.0 HaplotypeCaller version of multi-sample calling (30 minutes)

The issue with the UG approach to large cohort is, to a large extent, computational. It is very impractical to call very large cohorts of samples together because of the  $(n+1)$  problem: if you add one sample, you ideally would like to not recall the 2,000 other samples at the same time. How do we go around that? This is the point of the updated GATK pipeline which is well suited for the analysis of very large cohorts. The principle is simple:

- Step1: each BAM file is processed individually by the `HaplotypeCaller` module. The output is a format called gVCF, which extends the concept of the VCF. It essentially keeps information at all positions, not only the polymorphic ones. It is essential to separate regions without coverage from the ones with coverage but no polymorphism.
- Step 2: these gVCF files are combined together using the `GenotypeGVCF` module of GATK to create a single gVCF file, which can then be processed as usual.

An example of the gVCF output is shown below:

```
awk '{if ( (NR == 199) || (NR == 200) || (NR == 201) ) print $1,$2,$8,$9,$10}' \
    results/HG00257.mapped.ILLUMINA.bwa.GBR.exome.20121211.gvcf

## 21 43829243 END=43829243 GT:DP:GQ:MIN_DP:PL 0/0:14:12:14:0,12,434
## 21 43829244 END=43829244 GT:DP:GQ:MIN_DP:PL 0/0:14:33:14:0,33,495
## 21 43829245 END=43829246 GT:DP:GQ:MIN_DP:PL 0/0:13:9:13:0,0,388
```

**Exercise:** These two steps are described in two independent scripts called `multi_sample_calling_GATK_HC_step1.sh` and `multi_sample_calling_GATK_HC_step2.sh`. These scripts will need light editing to run properly, essentially modifying the path to the fasta reference genome file and the path to GATK. Make sure you can run the scripts, and that you understand the options that I used.

**Exercise:** In the gVCF output format highlighted above, I show 3 lines. Can you understand why GATK decides to go to the next line in that manner? A key tip comes from the command below. Can you see where these options come from in the gVCF file, what they mean, and how this relates to the new lines in the gVCF format?

```
grep GVCFBBlock results/HG00254.mapped.ILLUMINA.bwa.GBR.exome.20121211.gvcf
```

```
## ##GVCFBBlock=minGQ=0(inclusive),maxGQ=10(exclusive)
## ##GVCFBBlock=minGQ=10(inclusive),maxGQ=20(exclusive)
## ##GVCFBBlock=minGQ=20(inclusive),maxGQ=60(exclusive)
## ##GVCFBBlock=minGQ=60(inclusive),maxGQ=2147483647(exclusive)
```

*# Note the GQ value that changes from one block to the next at each line of the gVCF file. This is  
# what defines the presence of new lines, i.e. the level of compression essentially.*

**Exercise:** An badly intentioned colleague, jealous of your rapid progress with bioinformatics tools, has sabotaged your GATK script by adding two subtle errors that prevent the script from running. The script name is `multi_sample_calling_GATK_HC_step2_buggy.sh`. Can you understand what these two errors are and debug the script?

*# Error 1: an additional space has been added after a new line sign, which breaks the GATL scripts.  
# Error 2: the file containing the gVCF files has been modified and does not have a .list extension  
# anymore, which confuses GATK.*

## 2 VCFtools (about 45 minutes)

VCFtools is a swiss army type tool very handy for VCF manipulations. It contains a suite of Perl scripts, which we won't use today, as well as a binary executable, which will be the focus of this section. The documentation for the binary executable is located [here](#) and [here](#).

The first thing I propose to do is to compare the output of the GATK v2.0 output (Unified Genotyper, or UG) with the output from `samtools`.

```
vcftools --vcf results/UBASH31A_samtools.vcf --freq --out results/samtools
vcftools --vcf results/GATK_msamplercalling.vcf --freq --out results/GATK_UG
```

Note that from this output it is easy to compare the number of calls for both algorithms.

```
wc -l results/samtools.frq results/GATK_UG.frq

## 135 results/samtools.frq
## 112 results/GATK_UG.frq
## 247 total
```

So where does the difference come from? A first thing is the fact that the call to GATK used a minimum quality filter of 30, which was not applied to the `samtools` output. `VCFtools` allows you to recode VCF files after having applied relevant filters. Try for example:

```
vcftools --vcf results/UBASH31A_samtools.vcf --minQ 30 --out results/samtools_filtered --recode
```

Note the use of the `recode` argument. Without this, no new VCF file will be created. I suggest you try, and you will see nothing but a log file.

**Exercise:** Using the VCF output from `VCFtools`, create a new file with the frequency values for each site. Is the result more consistent with the output from GATK UG in terms of number of polymorphic sites? See how you could have created the same frequency file without the intermediate VCF file `results/samtools_filtered.recode.vcf`. Make sure that this “direct” frequency file is identical to the one that used the intermediate VCF file.

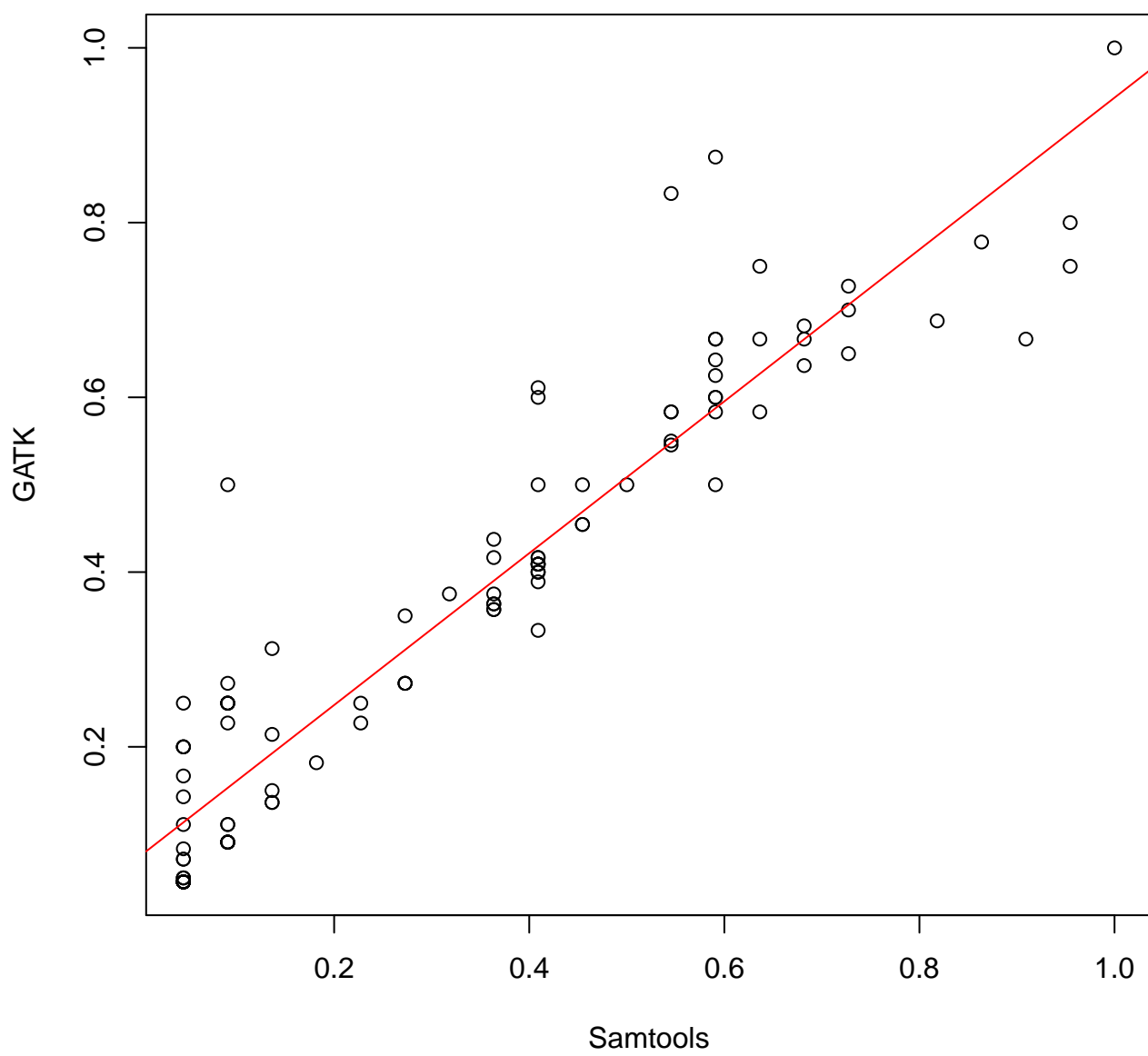
```
vcftools --vcf results/samtools_filtered.recode.vcf --freq --out results/freq_from_intermediate.vcf
vcftools --vcf results/UBASH31A_samtools.vcf --minQ 30 --freq --out results/freq_direct_from_samtools_outp
```

**Optional exercise:** For the subset of variants present in both `samtools` and GATK UG output, generate a graph to compare the estimated variant frequencies (using R most likely, but feel free to use something else if you prefer).

```
data.samtools <- read.table("results/samtools.frq", skip = 1, col.names = c("CHROM", "POS", "N_ALLELES",
  "N_CHR", "A111", "A112"))
data.gatk <- read.table("results/GATK_UG.frq", skip = 1, col.names = c("CHROM", "POS", "N_ALLELES", "N_CHR",
  "A111", "A112"))

### now I modify the string that contains the allele, to just keep the frequencies
data.samtools$frq.sam <- as.numeric(gsub(pattern = ".*:", data.samtools$A112, replacement = ""))
data.gatk$frq.gatk <- as.numeric(gsub(pattern = ".*:", data.gatk$A112, replacement = ""))

data <- merge(data.samtools, data.gatk, by = "POS")
plot(x = data$frq.sam, y = data$frq.gatk, xlab = "Samtools", ylab = "GATK")
abline(coef(lm(data$frq.gatk ~ data$frq.sam)), col = "red")
```



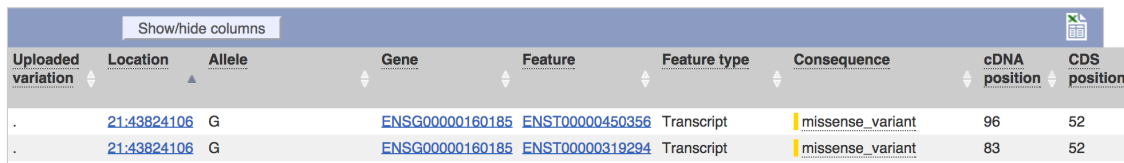
### 3 Variant annotation (about 45 minutes)

#### 3.1 Running the variant effect predictor (VEP)

After having called variants, the next step consists of annotating them. A widely used and reliable tool is the Variant Effect Predictor (VEP), which is part of the tools provided by Ensembl. The simplest way to interact with the VEP is the [web based interface](#). One possible input format, and probably the most widely used, is VCF. There is no need to copy all the fields, only the first 8 will be sufficient. For example you could use:

```
awk '{if ($1 == 21) print $1,$2,$3,$4,$5,$6,$7,$8}' \
    results/UBASH31A_samtools.vcf > results/samtools_for_VEP.vcf
```

**Exercise:** Note that there is a trap here. As the name indicates, these sequencing reads are supposed to cover the gene *UBASH31A*, however, none of the annotated variants is actually located in that gene. Can you figure out what went wrong? Fix it and download the output of the VEP in csv format (small excel like icon on the right side of the blue banner (see Figure 1)).



Uploaded variation	Location	Allele	Gene	Feature	Feature type	Consequence	cDNA position	CDS position
.	21:43824106	G	ENSG00000160185	ENST00000450356	Transcript	missense_variant	96	52
.	21:43824106	G	ENSG00000160185	ENST00000319294	Transcript	missense_variant	83	52

Figure 1: Small icon on the right side is what you need to download the output in csv format

Running the VEP directly from the web is a very practical thing to do but it is not effective for very large amount of data (tens of thousand of variants...). It is possible to run it directly from your computer. If you do that, the first time you run the VEP, it will download locally all the files that you need. The installation has been streamlined extensively by the Ensembl team and the `INSTALL.pl` will fetch and install the minimum set of tools required for the VEP to run.

```
wget -O VEP.zip https://github.com/Ensembl/ensembl-tools/archive/release/76.zip
unzip VEP.zip
cd ensembl-tools-release-76/scripts/variant_effect_predictor
perl INSTALL.pl
```

You can read more about this process at [this location](#).

Optional Exercise: Install the VEP locally, including the cached file for the NCBI build 37 of the human genome. Run the VEP against the VCF files that was generated by `samtools`. Note that this step can be time consuming, and you may hit small technical issues on these computers mostly set up for teaching, but it should be possible to make these scripts work on most linux installs.

#### 3.2 Interpreting the output of the VEP

You should have now downloaded the output of the VEP in csv format, but in case that is not the case my output is available (`samtools_output_VEP.csv`). R is probably the best tool to read a csv file like the one generated by the VEP.

**Exercise:** Using R, load the output of the VEP. Look at the number of lines and understand the logic of the output, in particular why the number of lines exceeds the number of variants that were provided as input.

#### 3.3 ANNOVAR as an alternative

Another software widely used is ANNOVAR. It is also written in perl, and you can find all the details [here](#). In class we will discuss the advantages and weaknesses of ANNOVAR compared to the VEP. They do deliver broadly comparable results, but the main difference comes from differences in the underlying databases (i.e. GENCODE, Ensembl, UCSC...).