

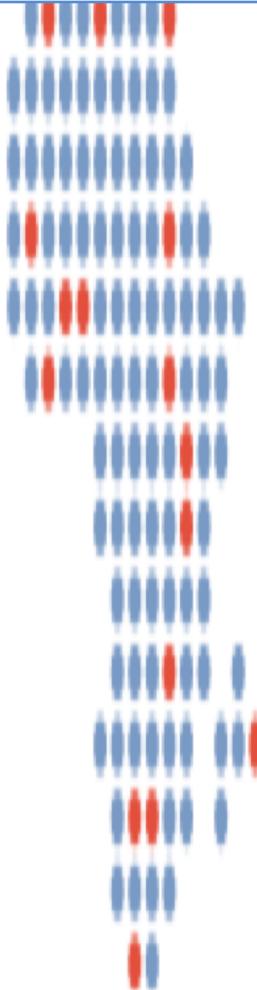


H3ABioNet

Pan African Bioinformatics Network for H3Africa

CONNECTING
SCIENCE

ADVANCED
COURSES +
SCIENTIFIC
CONFERENCES

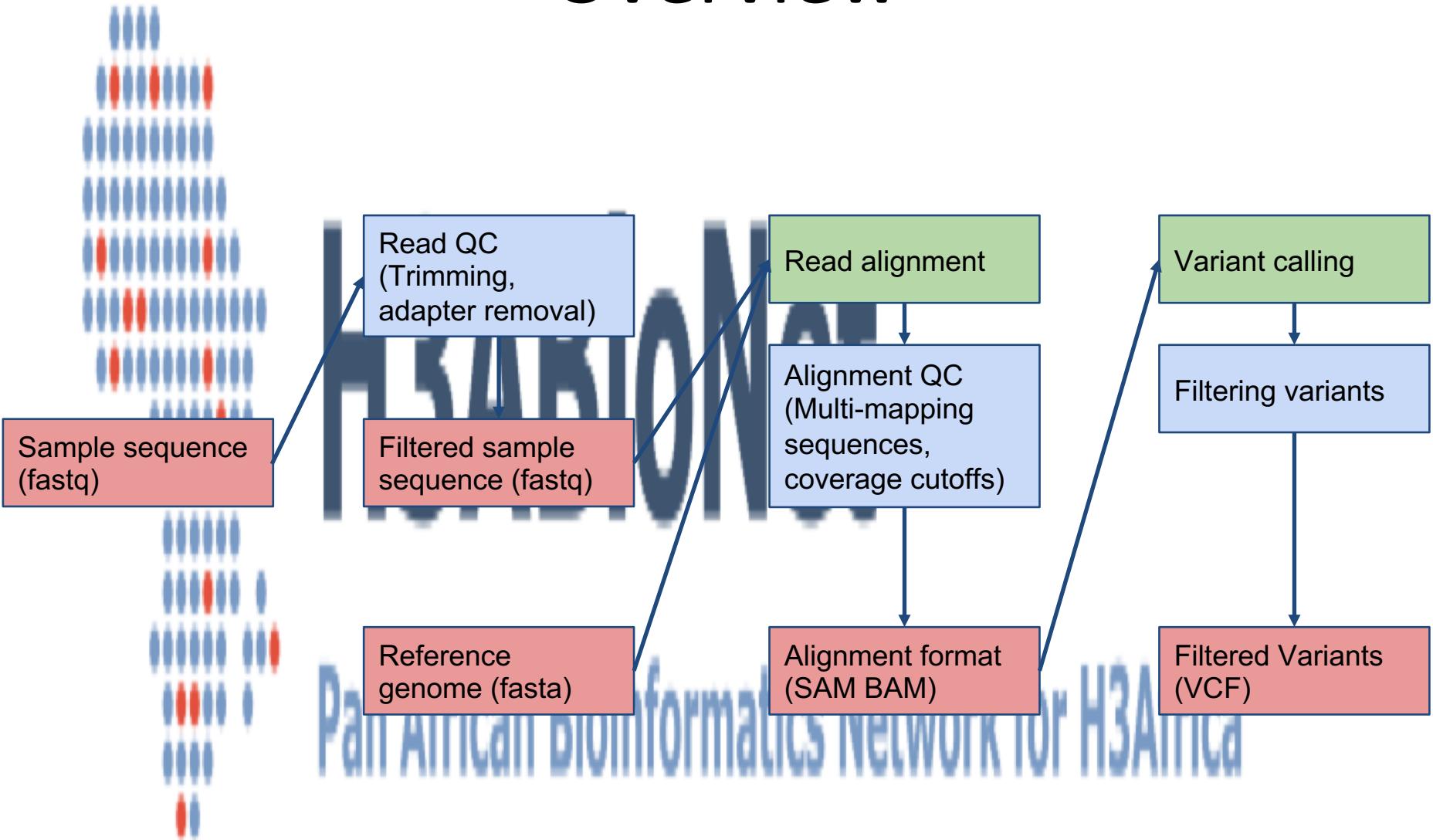


Next Generation Sequencing Bioinformatics Course 2021

Pathogen Variant Calling Variant calling

Pan African Bioinformatics Network for H3Africa

Overview



Overview

- Question: Investigating resistance in *M. tuberculosis*
- Next steps:
 - Where do errors come from?
 - How are variants called
 - Variant filtering
 - Investigate VCF files

Overview

What does variant calling do?

- Variant calling identifies small differences (variants) between a sample and a reference
- Resolution dependent on the read length & algorithm
- Most tools are limited to detecting events shorter than the length of a short-read sequencing alignment.
- SNPs (single-nucleotide polymorphisms), Indels (insertions and deletions), MNPs (multi-nucleotide polymorphisms)

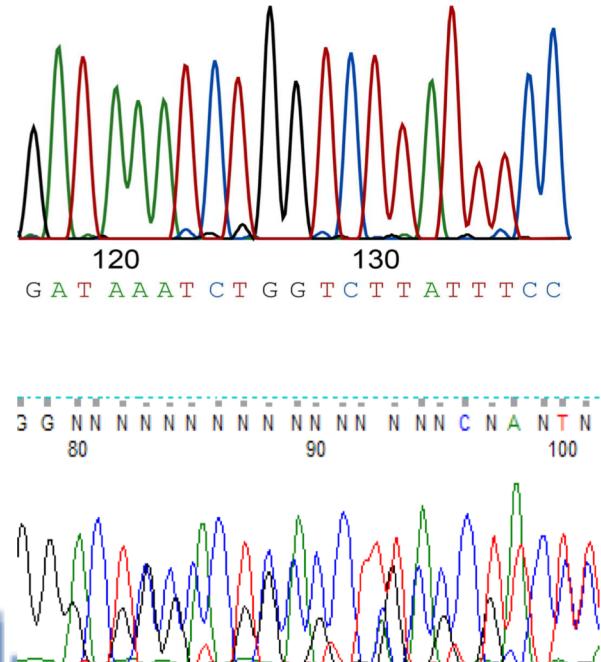
What does it not do?

- Identification of large structural events
 - Chromosomal rearrangements, large deletions, gene duplications
- Often variants are incorrectly called at structural breakpoints

QC - Base calling errors

Where do errors come from?

- Uncertainty from the base call
 - Normally quality falls off towards the end of the read
 - Low quality, wrong base
- PCR errors
 - Errors early in the PCR reaction
 - Bridge PCR (Illumina / Solexa solid phase amplification) or emulsion PCR
 - High quality, wrong base
- PCR vs PCR free sequencing
 - Nanopore sequencing = no PCR
- Platform specific error profiles



Variant calling

Goal:

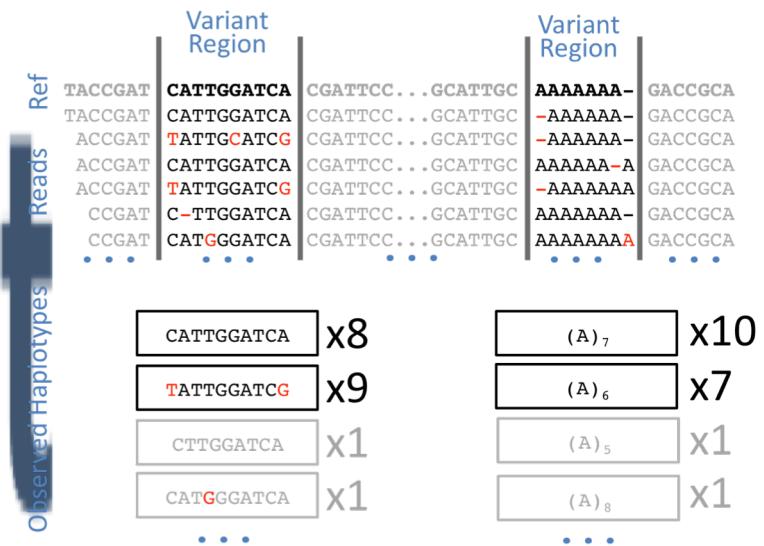
Given an alignment (bam), identify positions where reads support an alternate allele at a position compared to the reference sequence.

Tools:

- Freebayes
- SAMtools mpileup / bcftools
- GATK

Variant calling

- Freebayes
 - Uses a Bayesian approach to find variants
- GATK
 - Built for human genomes
 - Widely used, constantly improved
- SAMtools mpileup
 - bcftools mpileup generates genotype likelihoods
 - bcftools call makes the calls
- How to choose
 - Look at what other researchers are using for the organism
 - Use what you have support for / your team has experience with (but stay up to date)
 - Read papers to stay up to date



Source: <https://github.com/freebayes/freebayes>

Term definition: A haplotype is a group of genes within an organism that was inherited together from a single parent.

QC - Variant filtering

- What can you filter on?
 - Depth, ploidy, quality
- Depth
 - Why can high coverage be bad?
 - Could indicate a region with a lot of multi-mapping reads from a gene duplication
- Ploidy
 - In haploid organisms, things are simpler
 - Many tools built for diploid, still output genotype information (A/A or A/C or C/C)
- Quality
 - QUAL column: phred-scaled probability the ALT is wrong
 - GQ (Genotype quality) in the format column. Defined as $-10\log_{10}(\text{genotype call is wrong})$
 - Examples:
 - 1 / 1 000 chance the genotype call is wrong
 - $-10\log_{10}(0.001) = 30$
 - GQ = 12
 - $10^{-12/10} = \text{Probability of } 0.063095 \text{ that the call is incorrect}$

QC - Variant filtering

- Challenges

- Duplicate genes

- Sequenced genome had two copies of a region
 - Region appears only once in the reference genome

- Missing genes

- If a gene is missing in a sample but present in the reference, no variants will be called.

- Closely related reference genome NB

QC - Variant filtering

- Tools
 - gatk VariantFiltration
 - vcftools
 - bcftools filter / view
 - bcftools query
 - Useful to test filtering and report general quality
 - A useful cheat sheet for BCF:
 - <https://gist.github.com/elowy01/93922762e131d7abd3c7e8e166a74a0b>
 - By: Ernesto Lowy

QC - Anatomy of a VCF file



- Mapping of info in the Meta-information lines

```
##fileformat=VCFv4.0
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=1000GenomesPilot-NCBI36
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT N
A00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:
HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:.,.
20 17330 . T A 3 q10 PASS NS=3;DP=11;AF=0.017 GT:GQ:DP:
HQ 0|0:49:3:58,50 0|1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:
HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:
HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTCT G,GTACT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP
0/1:35:4 0/2:17:2 1/1:40:3
```

Source: <https://www.internationalgenome.org/wiki/Analysis/vcf4.0/>

QC - Anatomy of a VCF file



- Mapping of info in the Meta-information lines

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | N |
|--------|--|-----------|------|---------|------|--------|-----------------------------------|-----------|---|
| A00001 | | NA00002 | | NA00003 | | | | | |
| 20 | 14370 | rs6054257 | G | A | 29 | PASS | NS=3;DP=14;AF=0.5;DB;H2 | GT:GQ:DP: | |
| | HQ 0 0:48:1:51,51 1 0:48:8:51,51 1/1:43:5:,,,. | | | | | | | | |
| 20 | 17330 | . | T | A | 3 | q10 | NS=3;DP=11;AF=0.017 | GT:GQ:DP: | |
| | HQ 0 0:49:3:58,50 0 1:3:5:65,3 0/0:41:3 | | | | | | | | |
| 20 | 1110696 | rs6040355 | A | G,T | 67 | PASS | NS=2;DP=10;AF=0.333,0.667;AA=T;DB | GT:GQ:DP: | |
| | HQ 1 2:21:6:23,27 2 1:2:0:18,2 2/2:35:4 | | | | | | | | |
| 20 | 1230237 | . | T | . | 47 | PASS | NS=3;DP=13;AA=T | GT:GQ:DP: | |
| | HQ 0 0:54:7:56,60 0 0:48:4:51,51 0/0:61:2 | | | | | | | | |
| 20 | 1234567 | microsat1 | GTCT | G,GTACT | 50 | PASS | NS=3;DP=9;AA=G | GT:GQ:DP | |
| | 0/1:35:4 0/2:17:2 1/1:40:3 | | | | | | | | |

Source: <https://www.internationalgenome.org/wiki/Analysis/vcf4.0/>

OC - Anatomy of a VCF file

```
##fileformat=VCFv4.0
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=1000GenomesPilot-NCBI36
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT N
A00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:
HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:...
20 17330 . T A 3 q10 PASS NS=3;DP=11;AF=0.017 GT:GQ:DP:
HQ 0|0:49:3:58,50 0|1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:
HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:
HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTCT G,GTACT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP
0/1:35:4 0/2:17:2 1/1:40:3
```



Source: <https://www.internationalgenome.org/wiki/Analysis/vcf4.0/>

QC - Anatomy of a VCF file

- Mapping of info in the Meta-information lines

| #fileformat=VCFv4.0 ##fileDate=20090805 ##source=myImputationProgramV3.1 ##reference=1000GenomesPilot-NCBI36 ##phasing=partial ##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data"> ##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth"> | | | | | | | FORMAT | N |
|--|---|-----------|------|---------|------|--------|-----------------------------------|-----------|
| CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | |
| A00001 | | NA00002 | | NA00003 | | | | |
| 20 | 14370 | rs6054257 | G | A | 29 | PASS | NS=3;DP=14;AF=0.5;DB;H2 | GT:GQ:DP: |
| | HQ 0 0:48:1:51,51 1 0:48:8:51,51 1 1:43:5:,,. | | | | | | | |
| 20 | 17330 | . | T | A | 3 | q10 | NS=3;DP=11;AF=0.017 | GT:GQ:DP: |
| | HQ 0 0:49:3:58,50 0 1:3:5:65,3 0/0:41:3 | | | | | | | |
| 20 | 1110696 | rs6040355 | A | G,T | 67 | PASS | NS=2;DP=10;AF=0.333,0.667;AA=T;DB | GT:GQ:DP: |
| | HQ 1 2:21:6:23,27 2 1:2:0:18,2 2/2:35:4 | | | | | | | |
| 20 | 1230237 | . | T | . | 47 | PASS | NS=3;DP=13;AA=T | GT:GQ:DP: |
| | HQ 0 0:54:7:56,60 0 0:48:4:51,51 0/0:61:2 | | | | | | | |
| 20 | 1234567 | microsat1 | GTCT | G,GTACT | 50 | PASS | NS=3;DP=9;AA=G | GT:GQ:DP |
| | 0/1:35:4 0/2:17:2 1/1:40:3 | | | | | | | |

Source: <https://www.internationalgenome.org/wiki/Analysis/vcf4.0/>

QC - Anatomy of a VCF file



- Multiple samples
in the same VCF



```
##fileformat=VCFv4.3
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:HQ 0|0:48:1:51,51 1|0:48:8:51,51 1/1:43:5:..
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017 GT:GQ:DP:HQ 0|0:49:3:58,50 0|1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTC G,GTCT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP 0/1:35:4 0/2:17:2 1/1:40:3
```



Source: <https://www.internationalgenome.org/wiki/Analysis/vcf4.0/>

QC - Variant filtering

- bcftools filter -i 'type="snp" && QUAL>=50 && FORMAT/DP>5 && MQ>=30' -g10 -G10 MD001_variants.vcf -o MD001_SNPs_filtered_try1.vcf
- QUAL = Base quality score column
- DP = Read depth
- MQ = Mapping quality