



Webinars On:
Basic and Advance Next Generation Sequencing (NGS) data analysis
(July 24th, 2020)



Variant calling, annotation and Visualization

*By: RAHILA SARDAR
ICMR- Senior Research Fellow
Translational Bioinformatics
Group, ICGEB*

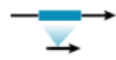


Genomic Variation

Single Nucleotide Variant



Deletion



Insertion



Tandem Duplication



Interspersed Duplication



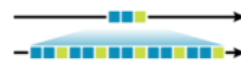
Inversion



Translocation



Copy Number Variant



Reference	1	2	3			
Insertion	1	2	5	3		
Deletion	1	3				
Inversion	1	3	2			
Copy Number Variation	1	1	1	1	2	3
Tandem Duplication	1	1	2	3		
Dispersed Duplication	1	2	1	3		
Mobile Element Insertion	1	2	Mobile Element	3		
Translocation	1					
	10	11	12	2	3	



Types of mutations

Somatic mutation



- Synonymous mutation
- Non-synonymous mutation

Germline mutation



- Synonymous mutation
- Non-synonymous mutation

- Synonymous mutations: which do not alter amino acid sequences and are (sometimes) silent mutations.
- Non-synonymous mutations: nucleotide mutation that alters the amino acid sequence of a protein.
- Germline mutations: occur in gametes and can be passed onto offspring (every cell in the entire organism will be affected)
- Somatic mutations: occur in a single body cell and cannot be inherited (only tissues derived from mutated cell are affected)

Variant calling workflow



Stage	Examples/explanation	File formats
Laboratory work	Experimental design Library preparation Enrichment (capture)	
Next-generation sequencing	Platforms include Illumina, SOLiD, Pacific Biosciences, other	Output: FASTQ-Sanger, FASTQ-Illumina SRA, EBI,ENA etc
Analysis pipeline	Trimming, filtering Software: FastQC	FASTQ TrimGalore, Trimmomatic, PrinSeq
	Alignment to reference genome Software: BWA, Bowtie2	Reference: FASTA Output: SAM/BAM
	Variant identification Single nucleotide variants (SNVs), structural variants (e.g. indels) Software: GATK, SAMTools Realignment, recalibration	Variant Call Format (VCF/BCF)
	Annotation Comparison to public database (dbSNP, 1000 Genomes); functional consequence scores	
Visualization	Variant visualization; read depth; comparison to other samples Software: IGV, BEDTools, BigBED	
Prioritization	Discovery of relevant variants Software: PolyPhen-2, VEP, VAAST SIFT, Panther	VCF
Storage	Deposit data in ENA, SRA, dbGaP	BAM, VCF



Reference databases

Reference genome database

- **NCBIRefSeq**: <https://www.ncbi.nlm.nih.gov/refseq/>
- **UCSC** : <https://genome.ucsc.edu>
- **Ensembl**: www.ensembl.org

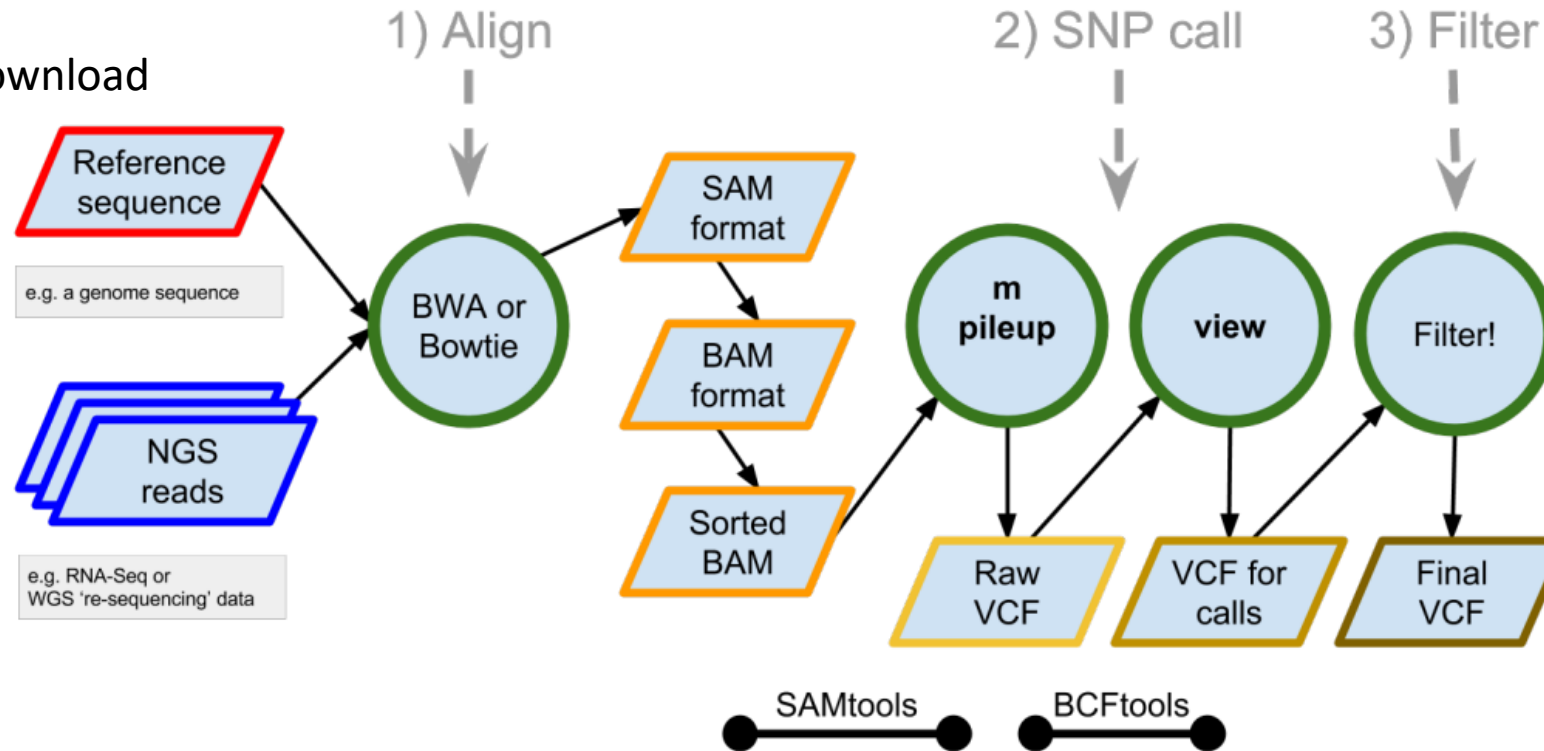
SNP database

- **dbSNP** : Database of short Genetic Variants
- **dbVar** : Database of genomic structural variants
- **ClinVAR**: Genomic variations and their relationship to human health and disease
- **dbGap**: Database of genotypes and phenotypes



Variant calling by samtools pipeline

wget link to download



Source: EBI

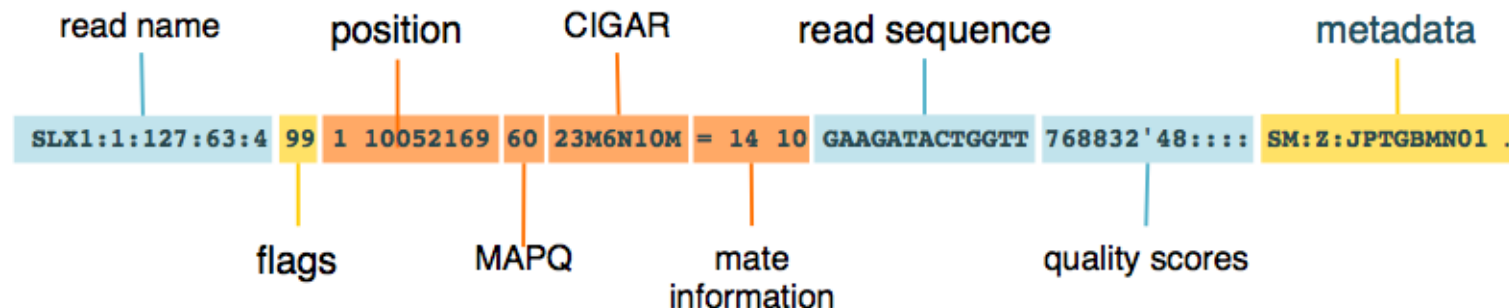


Sequence alignment/map format (SAM) and BAM

- SAM is a common format having sequence reads and their alignment to a reference genome.
- BAM is the binary form of a SAM file.
- Aligned BAM files are available at repositories (Sequence Read Archive at NCBI, ENA at Ensembl)
- SAMTools is a software package commonly used to analyze SAM/BAM files.
- Visit <http://samtools.sourceforge.net/>

HEADER containing metadata (sequence dictionary, read group definitions etc)

RECORDS containing structured read information (1 line per read record)





Align reads to reference (using BWA)

1. Index the reference (genome) sequence

- `bwa index e.coli.fasta`
- # The various index files are output in the CWD

2. Perform the alignment

- `bwa aln [opts] e.coli.fasta input1.fastq > output.sai`

3. Output results in SAM format

- `bwa samse e.coli.fasta output.sai input1.fastq > output.sam {Single end}`
- `bwa sampe e.coli.fasta output.sai input1.fastq input2.fastq > output.sam {paired end}`
- `bwa mem e.coli.fasta output.sai input1.fastq input2.fastq > ouput.sam {paired end}`

Index file= `e.coli..fasta.ann`, `e.coli.fasta.bwt`, `e.coli.fasta.fai`,`e.coli.fasta.pac`, `e.coli.fasta.sa`

*`bwa sampe`: single end , `sampe` for paired end. BWA-MEM also has better performance than BWA-backtrack for 70-100bp Illumina reads.



1... convert alignments (using SAMtools)

1. Convert SAM to BAM for sorting

➤ `samtools view -S -b output.sam > output.bam`

2. Sort BAM for SNP calling

➤ `samtools sort output.bam output-sorted`

Alignments are both:

- compressed for long term storage and
- sorted for variant discovery. Sort alignments by leftmost coordinates, or by read name when `-n` is used.



2) Call SNPs (using SAMtools)

1. Index the genome assembly (again!)

➤ **samtools faidx e.coli.fasta**

All we did so far (roughly) is to perform a formatconversion from BAM to VCF!

faidx: Index reference sequence in the FASTA format or extract subsequence from indexed reference sequence. If no region is specified, faidx will index the file and create <ref.fasta>.fai on the disk. If regions are specified, the subsequences will be retrieved and printed to stdout in the FASTA format.

2. Run 'mpileup' to generate VCF format

➤ `samtools mpileup -f e.coli.fasta output-sorted1.bam output-sorted-2.bam output_sorted-n.bam > output_raw.bcf`

```
Usage: samtools mpileup [options] in1.bam [in2.bam [...]]

Input options:
  -G, --illumina1.3+      quality is in the Illumina-1.3+ encoding
  -A, --count-orphans     do not discard anomalous read pairs
  -b, --bam-list FILE      list of input BAM filenames, one per line
  -B, --no-BAQ            disable BAQ (per-Base Alignment Quality)
  -C, --adjust-MQ INT     adjust mapping quality; recommended:50, disable:0 [0]
  -d, --max-depth INT    max per-file depth; avoids excessive memory usage [8000]
  -E, --redo-BAQ         recalculate BAQ on the fly, ignore existing BQs
  -f, --fasta-ref FILE    faidx indexed reference sequence file
  -G, --exclude-RG FILE   exclude read groups listed in FILE
  -l, --positions FILE    skip unlisted positions (chr pos) or regions (BED)
  -q, --min-MQ INT        skip alignments with mapQ smaller than INT [0]
  -Q, --min-BQ INT        skip bases with baseQ/BAQ smaller than INT [13]
  -r, --region REG        region in which pileup is generated
  -R, --ignore-RG         ignore RG tags (one BAM = one sample)
  --rf, --incl-flags STR|INT required flags: skip reads with mask bits unset []
  --ff, --excl-flags STR|INT filter flags: skip reads with mask bits set
                               [UNMAP,SECONDARY,QCFAIL,DUP]
  -x, --ignore-overlaps  disable read-pair overlap detection
  -X, --customized-index use customized index files

Output options:
  -o, --output FILE      write output to FILE [standard output]
  -O, --output-BP        output base positions on reads
  -s, --output-MQ        output mapping quality
  --output-QNAME         output read names
  --output-extra STR     output extra read fields and read tag values
  --output-sep CHAR     set the separator character for tag lists [,]
  --output-empty CHAR   set the no value character for tag lists [*]
  --reverse-del         use '#' character for deletions on the reverse strand
  -a                    output all positions (including zero depth)
```

mpileup: Generate text pileup output for one or multiple BAM files.



3) Call SNPs (using bcftools)

3. Call SNPs Bcf to vcf

bcftools view output.var.bcf >output.var.vcf

```
Output options:
-G, --drop-genotypes          drop individual genotype information (after subsetting if -s option set)
-h/H, --header-only/--no-header  print the header only/suppress the header in VCF output
-l, --compression-level [0-9]   compression level: 0 uncompressed, 1 best speed, 9 best compression [-1]
                                --no-version                          do not append version and command line to the header
-o, --output-file <file>        output file name [stdout]
-O, --output-type <b|u|z|v>     b: compressed BCF, u: uncompressed BCF, z: compressed VCF, v: uncompressed VCF [v]
-r, --regions <region>         restrict to comma-separated list of regions
-R, --regions-file <file>       restrict to regions listed in a file
-t, --targets [^]<region>       similar to -r but streams rather than index-jumps. Exclude regions with "^" prefix
-T, --targets-file [^]<file>    similar to -R but streams rather than index-jumps. Exclude regions with "^" prefix
--threads <int>                use multithreading with <int> worker threads [0]

Subset options:
-a, --trim-alt-alleles          trim ALT alleles not seen in the genotype fields (or their subset with -s/--s)
-I, --no-update                 do not (re)calculate INFO fields for the subset (currently INFO/AC and INFO/AN)
-s, --samples [^]<list>         comma separated list of samples to include (or exclude with "^" prefix)
-S, --samples-file [^]<file>    file of samples to include (or exclude with "^" prefix)
--force-samples                 only warn about unknown subset samples

Filter options:
-c/C, --min-ac/--max-ac <int>[:<type>]  minimum/maximum count for non-reference (nref), 1st alternate (alt1), least frequent
                                         (minor), most frequent (major) or sum of all but most frequent (nonmajor) alleles [nref]
-f, --apply-filters <list>          require at least one of the listed FILTER strings (e.g. "PASS,")
-g, --genotype [^]<hom|het|miss>     require one or more hom/het/missing genotype or, if prefixed with "^", exclude sites with hom/het/missing genotypes
-i/e, --include/--exclude <expr>    select/exclude sites for which the expression is true (see man page for details)
-k/n, --known/--novel              select known/novel sites only (ID is not/is '.')
-m/M, --min-alleles/--max-alleles <int>  minimum/maximum number of alleles listed in REF and ALT (e.g. -m2 -M2 for biallelic sites)
-p/P, --phased/--exclude-phased     select/exclude sites where all samples are phased
-q/Q, --min-af/--max-af <float>[:<type>]  minimum/maximum frequency for non-reference (nref), 1st alternate (alt1), least frequent
                                         (minor), most frequent (major) or sum of all but most frequent (nonmajor) alleles [nref]
-u/U, --uncalled/--exclude-uncalled  select/exclude sites without a called genotype
-v/V, --types/--exclude-types <list>  select/exclude comma-separated list of variant types: snps,indels,mnps,ref,bnd,other [null]
```

bcftools view: Applies the prior and does the actual calling and convert bcf to vcf



4) Filter SNPs

1. Filter SNPs

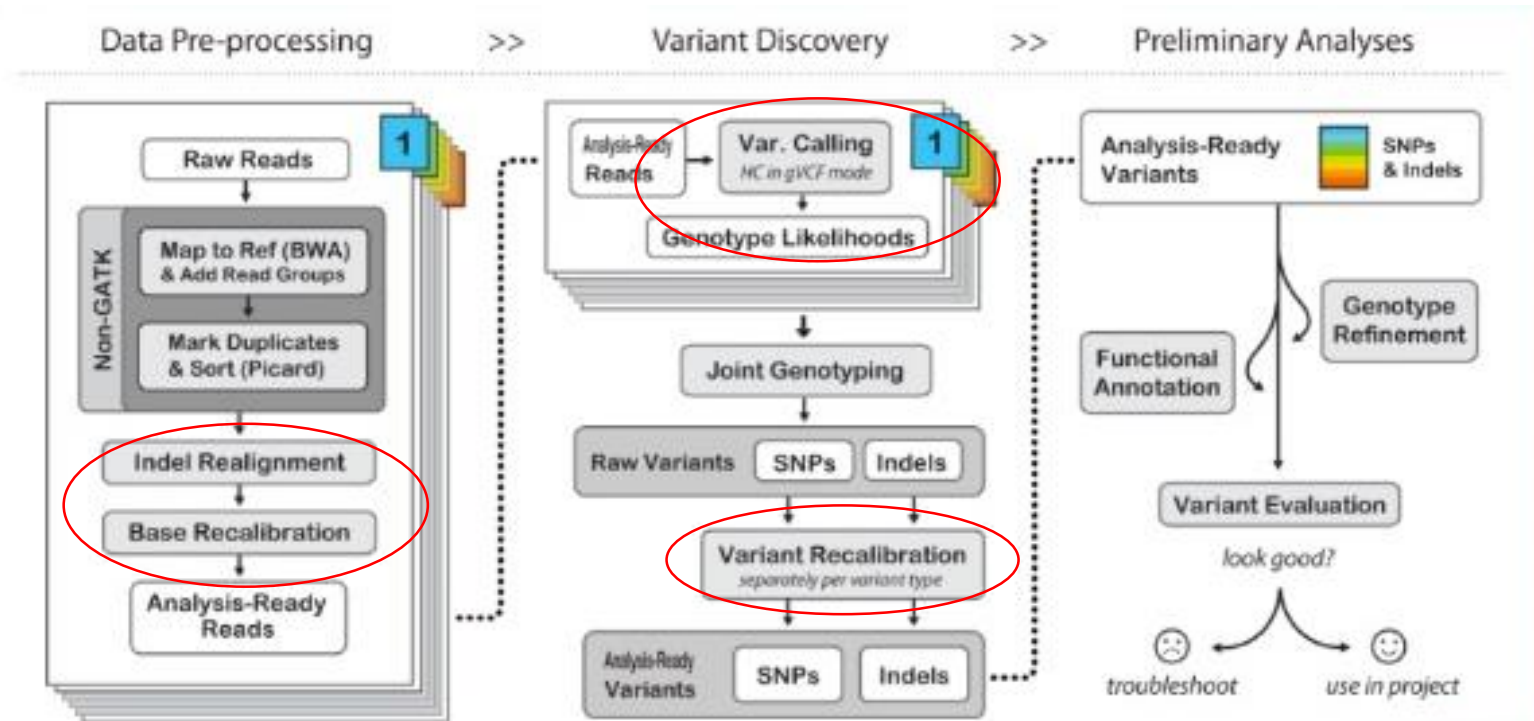
➤ `vcfutils.pl varFilter output.var.vcf > outputs.var-final.vcf`

```
Options: -Q INT    minimum RMS mapping quality for SNPs [10]
         -d INT    minimum read depth [2]
         -D INT    maximum read depth [100000000]
         -a INT    minimum number of alternate bases [2]
         -w INT    SNP within INT bp around a gap to be filtered [3]
         -W INT    window size for filtering adjacent gaps [10]
         -1 FLOAT  min P-value for strand bias (given PV4) [0.0001]
         -2 FLOAT  min P-value for baseQ bias [1e-100]
         -3 FLOAT  min P-value for mapQ bias [0]
         -4 FLOAT  min P-value for end distance bias [0.0001]
         -e FLOAT  min P-value for HWE (plus F<0) [0.0001]
         -p        print filtered variants
```

Note: Some of the filters rely on annotations generated by SAMtools/BCFtools.

- Note: `vcfutils.pl` is a perl script that helps filtering variants according to a certain set of parameters

GATK

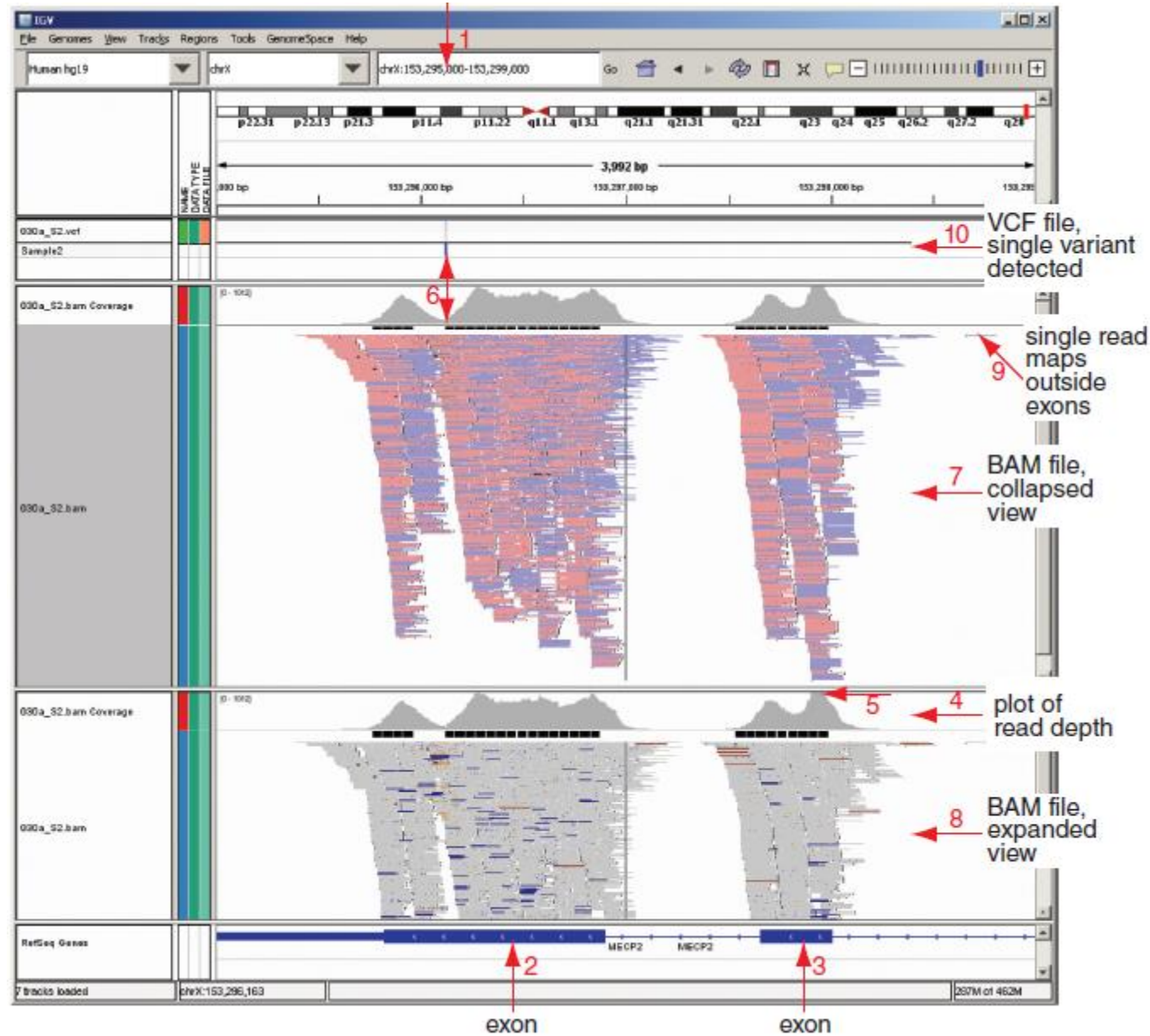


GATK vs Samtools



Variant callers	Samtools	GATK
Algorithm	bayesian approach to call the variants	bayesian approach to call the variants
Preprocess of alignment.	No	uses picard tools
GATK specific features	No other variant callers	haplotype caller, indel realignment, UnifiedGenotype among many others
Recalibration	No	Yes using ML algorithm
Organism specific	No	Yes
Utility	samtools is equally scalable to huge data sets in spite of its much simpler framewor	GATK is clearly the more sophisticated and the more complicated

Integrative Genomics Viewer (IGV)



Variant Call Format (VCF) file summarizes variation



Column	Mandatory	Description
CHROM	Yes	Chromosome
POS	Yes	1-based position of the start of the variant
ID	Yes	Unique identifier of the variant; the dbSNP entry rs1413368 is given in our example
REF	Yes	Reference allele
ALT	Yes	A comma-separated list of alternate nonreference alleles
QUAL	Yes	Phred-scaled quality score
FILTER	Yes	Site filtering information; in our example it is PASS
INFO	Yes	A semicolon-separated list of additional information. These fields include the gene identifier GI (here the gene is NEGR1); the transcript identifier TI (here NM_173808); and the functional consequence FC (here a synonymous change, T296T).
FORMAT	No	Defines information in subsequent genotype columns; colon separated. For example, GT:AD:DP:GQ:PL:VF:GQX in our example refers to genotype (GT), allelic depths for the ref and alt alleles in the order listed (AD), approximate read depth (reads with MQ=255 or with bad mates are filtered) (DP), genotype quality (GQ), normalized, Phred-scaled likelihoods for genotypes as defined in the VCF specification (PL), variant frequency, the ratio of the sum of the called variant depth to the total depth (VF), and minimum of {genotype quality assuming variant position, genotype quality assuming nonvariant position} (GXQ).
Sample	No	Sample identifiers define the samples included in the VCF file

SNP prioritization



PredictSNP 1.0

Consensus classifier for prediction of disease-related mutations

INPUT Load example

Insert protein sequence in FASTA format:

```
>HBA_HUMAN
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAGQQVGGHG
KRVADALTNVAHVDDMPALNSALSDLHAHLKRVDPVFNKLSHCLLVTLAAHLPAEFTF
AVHASLDKFLASVSTVLTISKYR
```

Load

MUTATIONS Manual input

Select positions:

1	M	V	L	S	P	A	D	K	T	N	V	K	A	A	W	K	V	G	A	H	A	G	E	Y	G	A	E	A	L	E	R	M	F	L	S	F	P	T	T		
41	K	T	Y	F	P	H	E	D	L	S	H	G	S	A	Q	V	K	G	H	G	K	K	V	A	D	A	L	T	N	A	V	A	H	V	D	D	M	P	N	A	
81	L	S	A	L	S	D	L	H	A	H	K	L	R	V	D	P	V	N	F	K	L	L	S	H	C	L	L	V	T	L	A	A	H	L	P	A	E	F	T	P	
121	A	V	H	A	S	L	D	K	F	L	A	S	V	S	T	V	L	T	S	K	Y	R																			

Select mutations:

<input type="checkbox"/> ALL	<input type="checkbox"/> C Cys	<input type="checkbox"/> D Asp	<input type="checkbox"/> E Glu	<input type="checkbox"/> F Phe
<input checked="" type="checkbox"/> A Ala	<input type="checkbox"/> G Gly	<input checked="" type="checkbox"/> H His	<input type="checkbox"/> I Ile	<input type="checkbox"/> K Lys
<input type="checkbox"/> M Met	<input type="checkbox"/> N Asn	<input type="checkbox"/> P Pro	<input type="checkbox"/> Q Gln	<input type="checkbox"/> R Arg
<input type="checkbox"/> S Ser	<input type="checkbox"/> T Thr	<input type="checkbox"/> V Val	<input type="checkbox"/> W Trp	<input checked="" type="checkbox"/> Y Tyr

Pos	Wild-type	Mutations	Clear
59	H	Y - Tyr	-
60	G	D - Asp, V - Val	-
63	V	T - Thr	-
68	T	V - Val	-
72	A	E - Glu, V - Val	-

Clear all mutations

TOOLS FOR EVALUATION

Tool name	Time demands	Expected accuracy
<input checked="" type="checkbox"/> PredictSNP	38 min	73.4%
<input checked="" type="checkbox"/> MAPP	10 min	70.7%
<input checked="" type="checkbox"/> PHD-SNP	38 min	71.5%
<input checked="" type="checkbox"/> PolyPhen-1	15 min	68.1%
<input checked="" type="checkbox"/> PolyPhen-2	15 min	69.2%
<input checked="" type="checkbox"/> SIFT	15 min	70.3%
<input checked="" type="checkbox"/> SNAP	30 min	67.6%
<input type="checkbox"/> nsSNPAnalyzer	15 min	62.9%
<input type="checkbox"/> PANTHER	5 min	63.5%

E-mail (optional):

Evaluate!

PredictSNP 1.0

Consensus classifier for prediction of disease-related mutations

RESULTS

Annotation	Mutation	PredictSNP	MAPP	PHD-SNP	PolyPhen-1	PolyPhen-2	SIFT	SNAP	nsSNPAnalyzer	PANTHER
▶	A72E	74%	70%	58%	67%	87%	66%	77%	65%	71%
▶	A72V	60%	59%	73%	67%	76%	53%	71%	65%	63%
▶	N79H	74%	72%	55%	67%	87%	53%	67%	65%	65%
▶	V97W	52%	45%	45%	74%	81%	79%	58%	65%	76%
▶	L110R	87%	80%	88%	74%	81%	79%	62%	63%	68%
▶	A112T	83%	75%	58%	67%	74%	76%	83%	65%	70%
▶	P115S	63%	72%	59%	67%	75%	79%	77%	65%	68%
▶	E117A	83%	46%	55%	67%	87%	77%	67%	65%	67%
▶	L126P	79%	81%	82%	74%	81%	79%	50%	63%	72%
Natural variant: in Quong Siz; causes alpha-thalassemia mapped from position 126 in Uniprot P52806										
Disease: Hemoglobin-H disease mapped from position 125 in PMD 2502009										
▶	L126R	79%	80%	86%	74%	68%	79%	50%	63%	69%
▶	S132P	61%	51%	77%	59%	43%	79%	67%	63%	55%

DOWNLOAD

Summary table Raw results

JOB INFORMATION

Job ID: axlwi@0zbtvqj5wadshehntbkykx5he4kyw70vvtufuto
 The results will be mailed to you once the job is completed.

LOG RECORDS

2013-05-09 15:32:03	PMD search finished.
2013-05-09 15:32:03	PMD annotations are being searched.
2013-05-09 15:32:03	UniProt search finished.
2013-05-09 15:31:52	UniProt annotations are being searched.
2013-05-09 15:31:52	PredictSNP consensus calculation successfully finished.
2013-05-09 15:31:50	PredictSNP consensus calculation running.
2013-05-09 15:31:50	Phd-SNP calculation successfully finished.

Ensembl “Variation Table” shows SIFT and PolyPhen scores for nonsynonymous variants



Missense variants ☰

[\[back\]](#)

Show All ▼ entries Show/hide columns Filter 📄									
ID	Chr: bp	Alleles	Class	Source	Type	AA	AA coord	SIFT	PolyPhen
rs121909815	11:5248247	A/G	SNP	dbSNP	Missense variant	V/A	2	0.01	0.119
rs121909830	11:5248247	A/C	SNP	dbSNP	Missense variant	V/G	2	0.07	0.007
rs121909815	11:5248247	A/G	SNP	dbSNP	Missense variant	V/A	2	0.01	0.119
rs121909830	11:5248247	A/C	SNP	dbSNP	Missense variant	V/G	2	0.01	0.007
rs33958358	11:5248248	C/T/A	SNP	dbSNP	Missense variant	V/L	2	0.01	0.001
rs33958358	11:5248248	C/T/A	SNP	dbSNP	Missense variant	V/M	2	0	0.271
rs33958358	11:5248248	C/T/A	SNP	dbSNP	Missense variant	V/L	2	0.02	0.001
rs33958358	11:5248248	C/T/A	SNP	dbSNP	Missense variant	V/M	2	0	0.271
rs35906307	11:5248245	G/A	SNP	dbSNP	Missense variant	H/Y	3	0.02	0.135

In Silico Analysis of SNPs in *PARK2* and *PINK1* Genes That Potentially Cause Autosomal Recessive Parkinson Disease

OPEN ACCESS PEER-REVIEWED
RESEARCH ARTICLE

In silico identification of genetic mutations conferring resistance to acetohydroxyacid synthase inhibitors: A case study of *Kochia scoparia*

Yan Li, Michael D. Netherland, Chaoyang Zhang, Huixiao Hong, Ping Gong

Published: May 7, 2019 • <https://doi.org/10.1371/journal.pone.0216116>

Yousuf Hasan, Yousuf Bakhit¹, Mohamed Osama Mirghani Ibrahim,² Mutaz Amin³

In silico analysis of a novel causative mutation in Cadherin23 gene identified in an Omani family with hearing loss

Show Mohammed Nasser Al-Kindi, Mazin Jawad Al-Khabouri, Khalsa Ahmad Al-Lamki, Flavia Palombo, Tommaso Pippucci, Giovanni Romeo & Nadia Mohammed Al-Wardy

Acad *Journal of Genetic Engineering and Biotechnology* 18, Article number: 8 (2020) | [Cite this article](#)
664 Accesses | [Metrics](#)

Abstract

Background

Hereditary hearing loss is a heterogeneous group of complex disorders with an overall incidence of one in every 500 newborns presented as syndromic and non-syndromic forms. Cadherin-related 23 (CDH23) is one of the listed deafness causative genes. It is found to be expressed in the stereocilia of hair cells and in the retina photoreceptor cells. Defective CDH23 have been associated mostly with prelingual severe-to-profound sensorineural hearing loss (SNHL) in either syndromic (USH1D) or non-syndromic SNHL (DFNB12) deafness. The purpose of this study was to identify causative mutations in an Omani family diagnosed with severe-profound sensorineural hearing loss by whole exome sequencing technique and analyzing the detected variant in silico for pathogenicity using several in silico mutation prediction software.

Article	Authors	Metrics	Comments	Media Coverage
---------	---------	---------	----------	----------------

Abstract

- Introduction
- Materials and methods
- Results
- Discussion
- Disclaimer
- Supporting information
- Acknowledgments
- References

Reader Comments (0)

Media Coverage (0)

Figures

Abstract

Mutations that confer herbicide resistance are a primary concern for herbicide-based chemical control of invasive plants and are often under-characterized structurally and functionally. As the outcome of selection pressure, resistance mutations usually result from repeated long-term applications of herbicides with the same mode of action and are discovered through extensive field trials. Here we used acetohydroxyacid synthase (AHAS) of *Kochia scoparia* (KSAHAS) as an example to demonstrate that, given the sequence of a target protein, the impact of genetic mutations on ligand binding could be evaluated and resistance mutations could be identified using a biophysics-based computational approach. Briefly, the 3D structures of wild-type (WT) and mutated KSAHAS-herbicide complexes were constructed by homology modeling, docking and molecular dynamics simulation. The resistance profile of two AHAS-inhibiting herbicides, tribenuron methyl and thifensulfuron methyl, was obtained by estimating their binding affinity with 29 KSAHAS (1 WT and 28 mutated) using 6 molecular mechanical (MM) and 18 hybrid quantum mechanical/molecular mechanical (QM/MM) methods in combination with three structure sampling strategies. By comparing predicted resistance with experimentally determined resistance in the 29 biotypes of *K. scoparia* field populations, we identified the best method (i.e., MM-PBSA with single structure) out of all tested methods for the herbicide-KSAHAS system, which exhibited the highest accuracy (up to 100%) in discerning mutations

Thank You and Questions

