

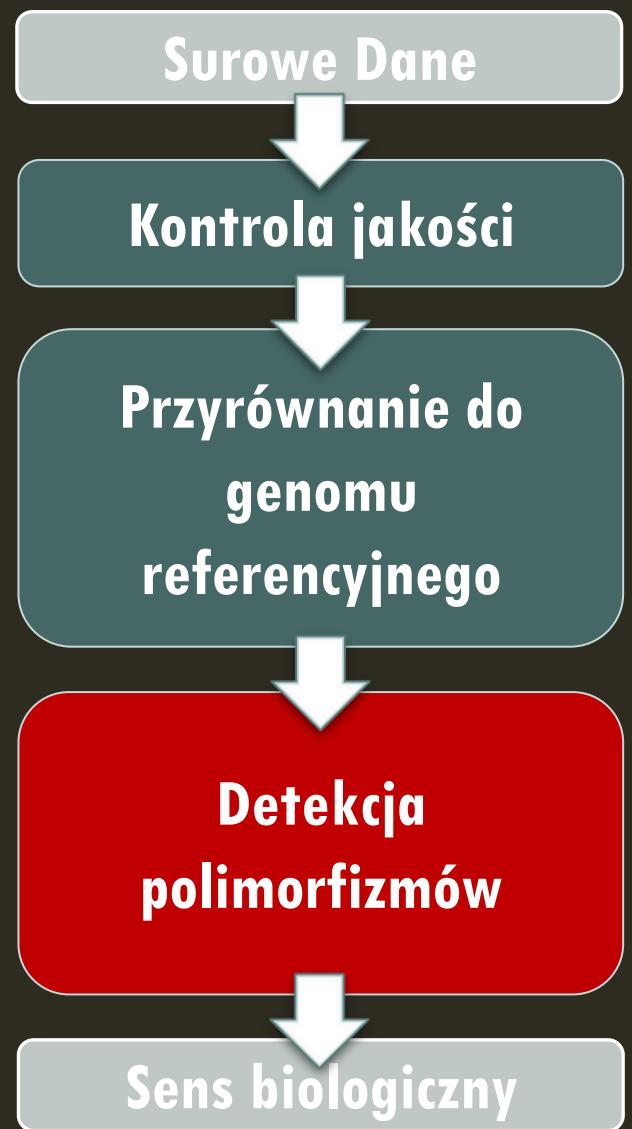
Niniejsze opracowanie zostało stworzone przez dr Magdę Mielczarek, pracownika Uniwersytetu Przyrodniczego we Wrocławiu w ramach wykonywania obowiązków związanych z kształceniem studentów i jest przeznaczone dla studentów Bioinformatyki (Wydział Biologii i Hodowli Zwierząt) na potrzeby dydaktyczne bez prawa do dalszego rozpowszechniania.

DETEKCJA POLIMORFIZMÓW

Analiza danych NGS
Wykład 6

PIPELINE

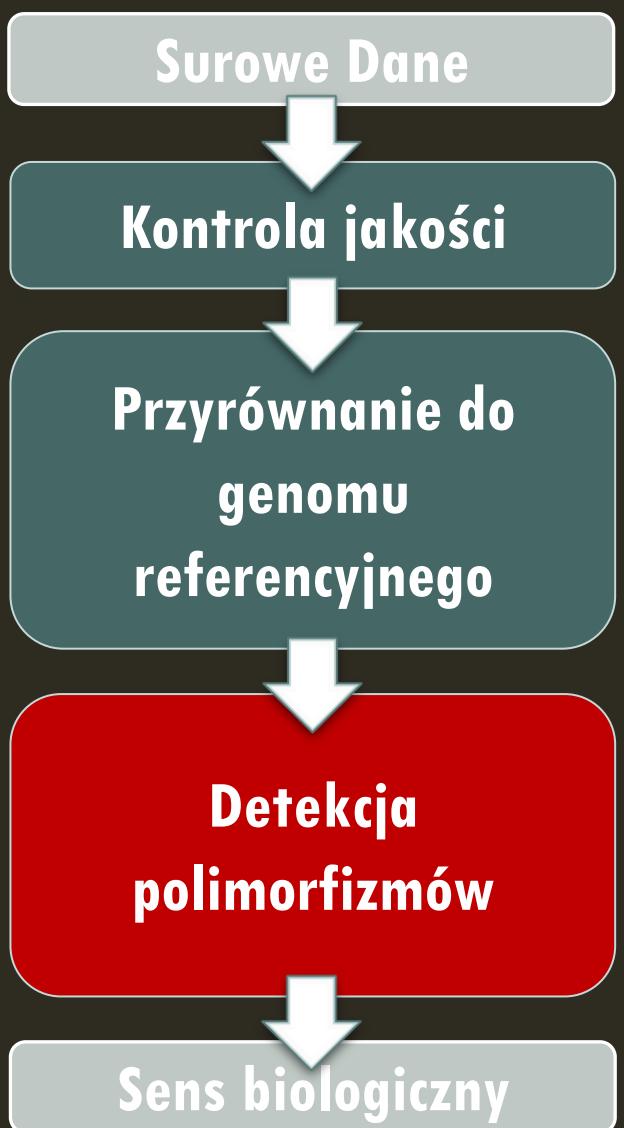
Pipeline = łańcuch przetwarzania danych



DETEKCJA POLIMORFIZMÓW

- Single Nucleotide Polymorphisms
- Insertions/Deletions
- Copy Number Variations
- Loss of Heterozygosity
- Inversions
- Translocations

SNP
INDEL
CNV
LOH
INV
TRANS



DETEKCJA POLIMORFIZMÓW



Genom referencyjny

Surowe Dane

Kontrola jakości

Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

DETEKCJA POLIMORFIZMÓW – PAKIET GATK

The screenshot shows the official GATK website at <https://gatk.broadinstitute.org/hc/en-us>. The top navigation bar includes links for User Guide, Tool Index, Blog, Forum, DRAGEN-GATK, Events, Download GATK4, and Sign in. The main header features the "gatk" logo and the text "Genome Analysis Toolkit" and "Variant Discovery in High-Throughput Sequencing Data". Below the header is a diagram illustrating the analysis pipeline: "Sequencing" leads to "READS", which then flows through the "gatk best practices™" processing engine to produce "VARIANTS". A descriptive text block below the diagram states: "Developed in the Data Sciences Platform at the Broad Institute, the toolkit offers a wide variety of tools with a primary focus on variant discovery and genotyping. Its powerful processing engine and high-performance computing features make it capable of taking on projects of any size. [Learn more](#)". At the bottom, there are three call-to-action boxes: "Find answers to your questions. Stay up to date on the latest news." (with a yellow arrow pointing to the "Getting Started" section), "Ask questions and help others." (with a yellow arrow pointing to the "Technical Documentation" section), and "Announcements" (with a yellow arrow pointing to the "Blog and events" section). Each box contains a small icon and a brief description.

User Guide Tool Index Blog Forum DRAGEN-GATK Events Download GATK4 Sign in

Genome Analysis Toolkit

Variant Discovery in High-Throughput Sequencing Data

gatk best practices™

Developed in the Data Sciences Platform at the Broad Institute, the toolkit offers a wide variety of tools with a primary focus on variant discovery and genotyping. Its powerful processing engine and high-performance computing features make it capable of taking on projects of any size. [Learn more](#)

Find answers to your questions. Stay up to date on the latest news.

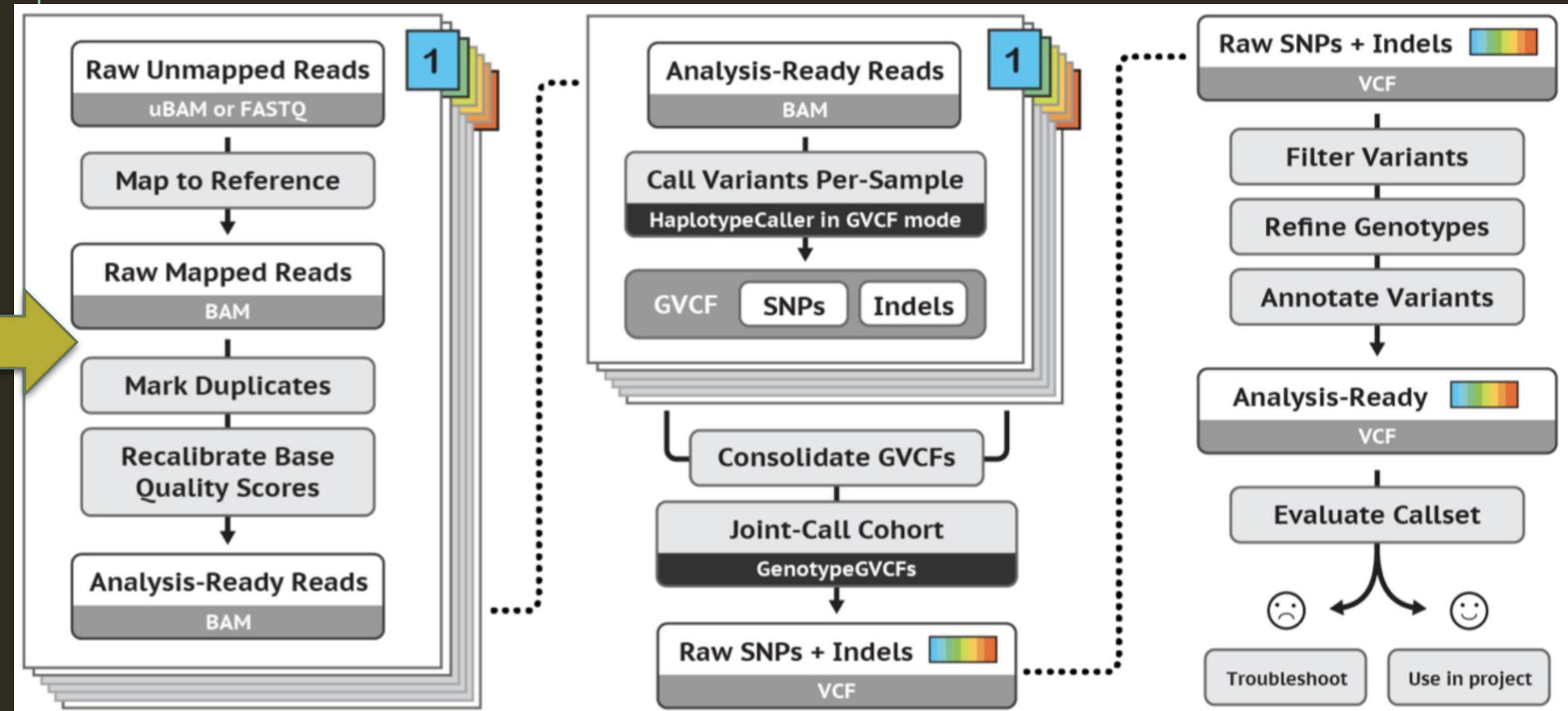
Ask questions and help others.

Getting Started
Best practices, tutorials, and other info to get you started

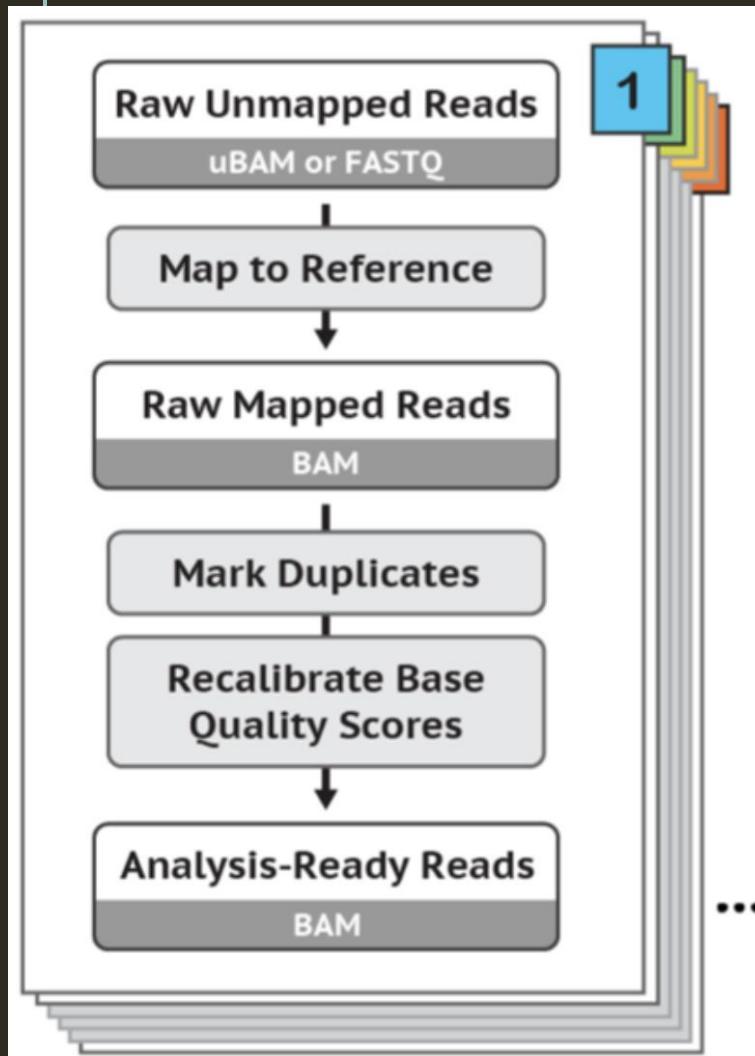
Technical Documentation
Algorithms, glossary, and other detailed resources

Announcements
Blog and events

DETEKCJA POLIMORFIZMÓW – PAKIET GATK



DETEKCJA POLIMORFIZMÓW – PAKIET GATK

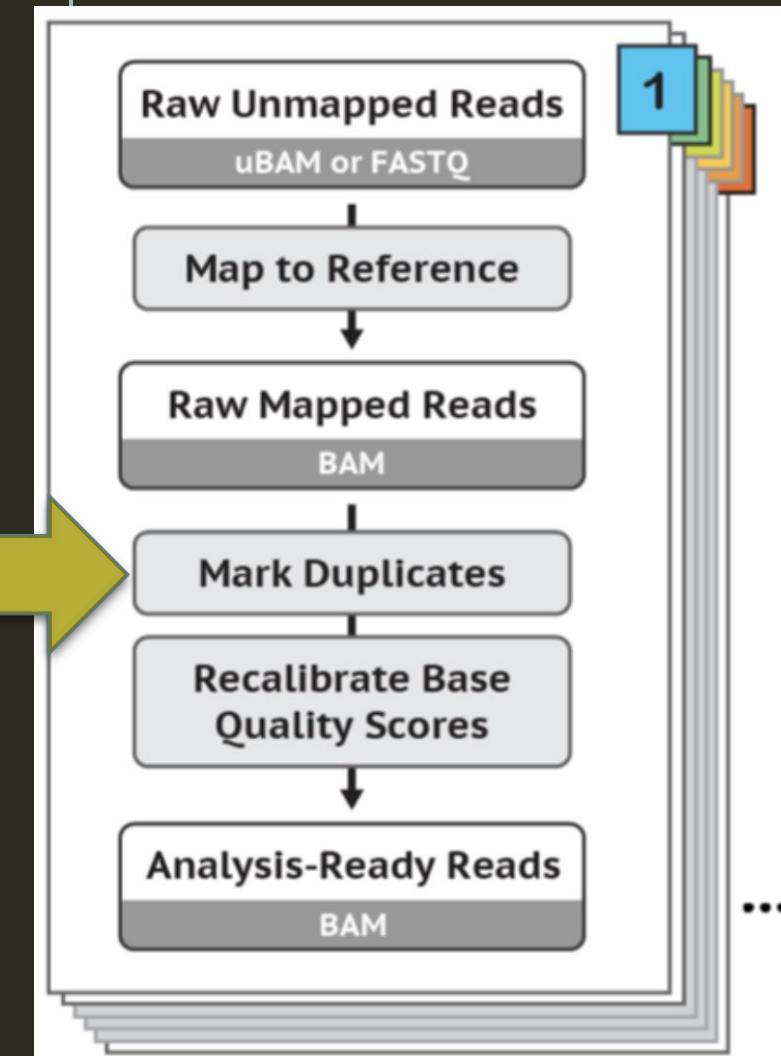


Mark Duplicates

„This tool locates and tags duplicate reads in a BAM or SAM file, where duplicate reads are defined as originating from a single fragment of DNA. Duplicates can arise during sample preparation e.g. library construction using PCR.”



DETEKCJA POLIMORFIZMÓW – PAKIET GATK

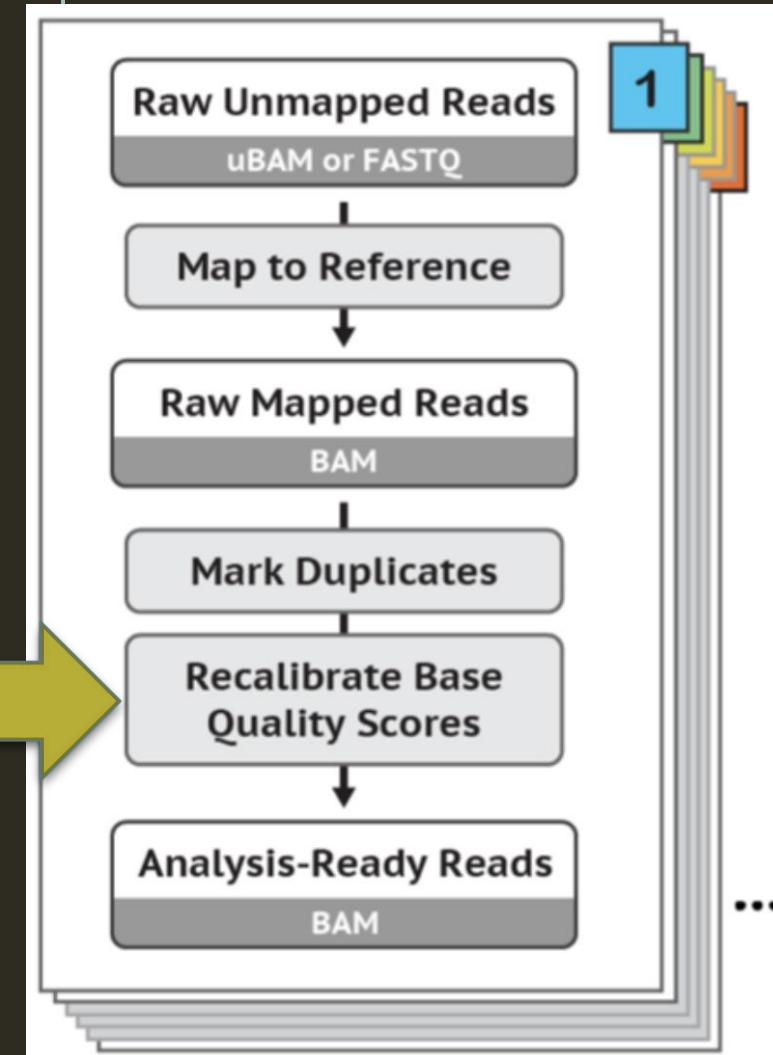


Mark Duplicates

„This tool locates and tags duplicate reads in a BAM or SAM file, where duplicate reads are defined as originating from a single fragment of DNA. Duplicates can arise during sample preparation e.g. library construction using PCR.”

```
java -jar picard.jar MarkDuplicates \
    I=input.bam \
    O=marked_duplicates.bam \
    M=marked_dup_metrics.txt
```

DETEKCJA POLIMORFIZMÓW – PAKIET GATK



Base Quality Score Recalibration (BQSR)

„Base quality scores are per-base estimates of error emitted by the sequencing machines (...). For example, let's say the machine reads an A nucleotide, and assigns a quality score of Q20 in Phred-scale, that means it's 99% sure it identified the base correctly. This may seem high, but it does mean that we can expect it to be wrong in one case out of 100; so if we have several billion base calls (we get ~90 billion in a 30x genome), at that rate the machine would make the wrong call in 900 million bases -- which is a lot of bad bases. The quality score each base call gets is determined through some dark magic jealously guarded by the manufacturer of the sequencing machines.”

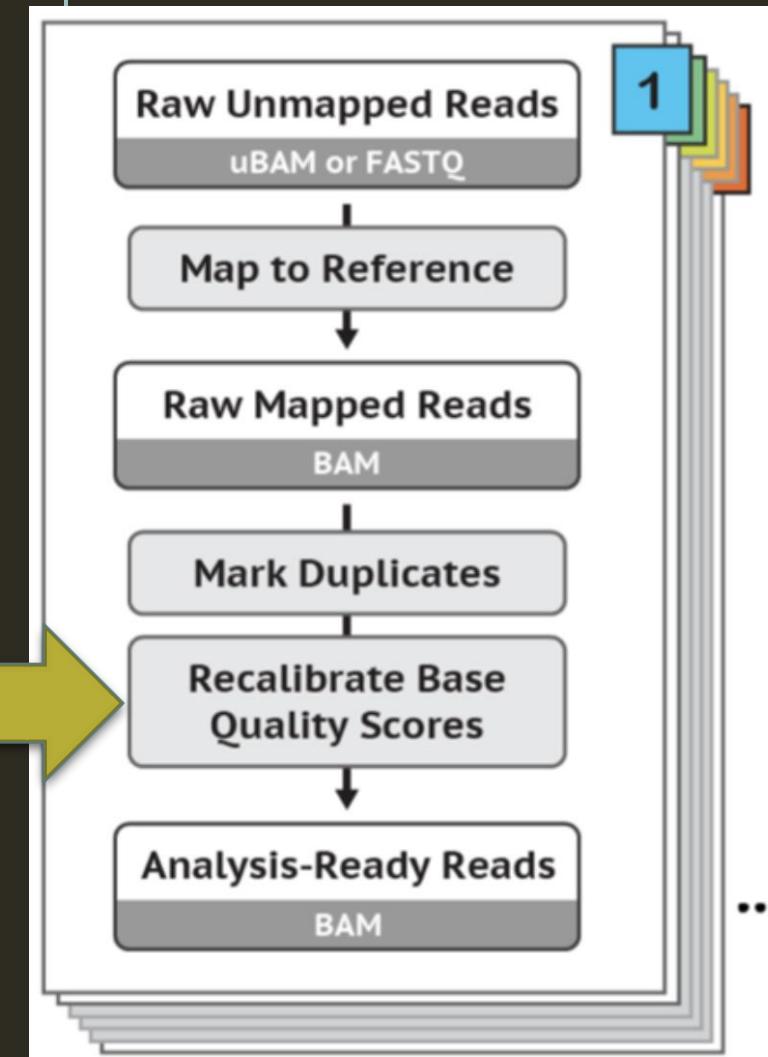
DETEKCJA POLIMORFIZMÓW – PAKIET GATK

Base Quality Score Recalibration (BQSR)

„(...) the scores are subject to various sources of systematic (nonrandom) technical error, leading to over- or under-estimated base quality scores in the data. Some of these errors are due to the physics or the chemistry of how the sequencing reaction works, (...) o manufacturing flaws in the equipment. (...). We apply machine learning to model these errors empirically and adjust the quality scores accordingly. For example (...) whenever we called two A nucleotides in a row, the next base we called had a 1% higher rate of error. So any base call that comes after AA in a read should have its quality score reduced by 1%. We do that over several different covariates (mainly sequence context and position in read, or cycle) in a way that is additive. So the same base may have its quality score increased for one reason and decreased for another.”

Szczegóły → gatk.broadinstitute.org/hc/en-us/articles/360035890531

DETEKCJA POLIMORFIZMÓW – PAKIET GATK

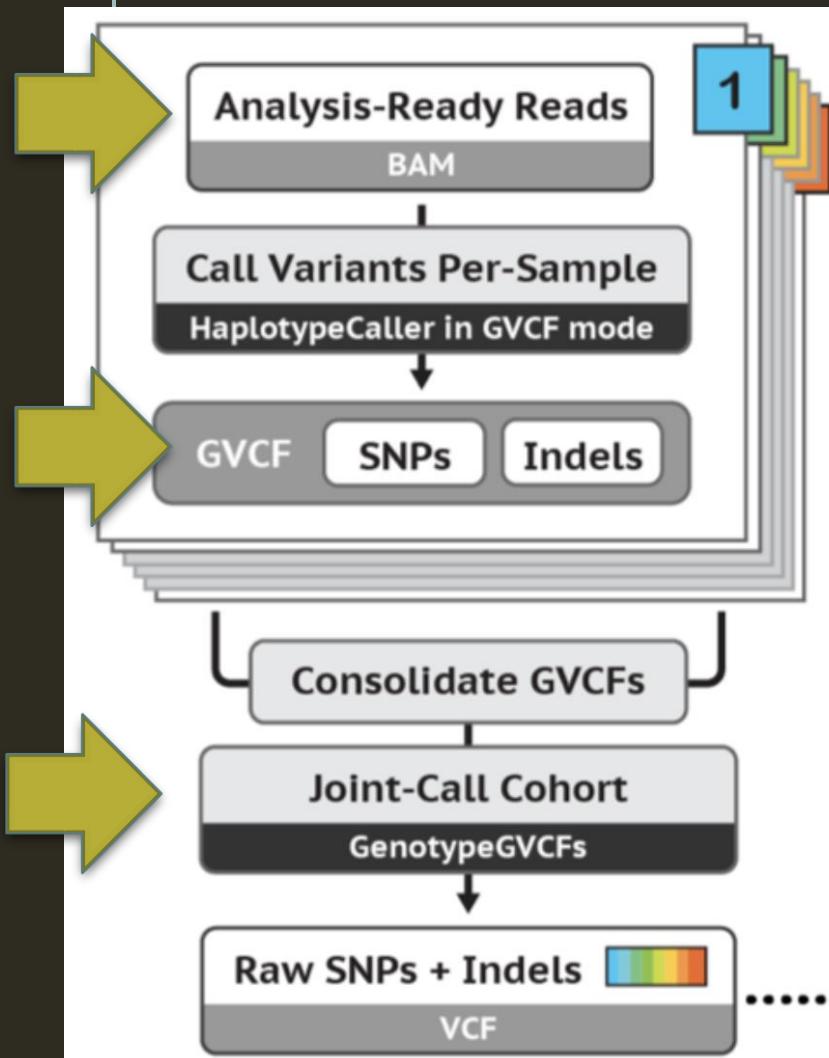


Base Quality Score Recalibration (BQSR)

```
gatk BaseRecalibrator \  
-I my_reads.bam \  
-R reference.fasta \  
--known-sites sites_of_variation.vcf \  
--known-sites another/optional/setOfSitesToMask.vcf \  
-O recal_data.table
```

```
gatk ApplyBQSR \  
-R reference.fasta \  
-I input.bam \  
--bqsr-recal-file recalibration.table \  
-O output.bam
```

DETEKCJA POLIMORFIZMÓW – PAKIET GATK

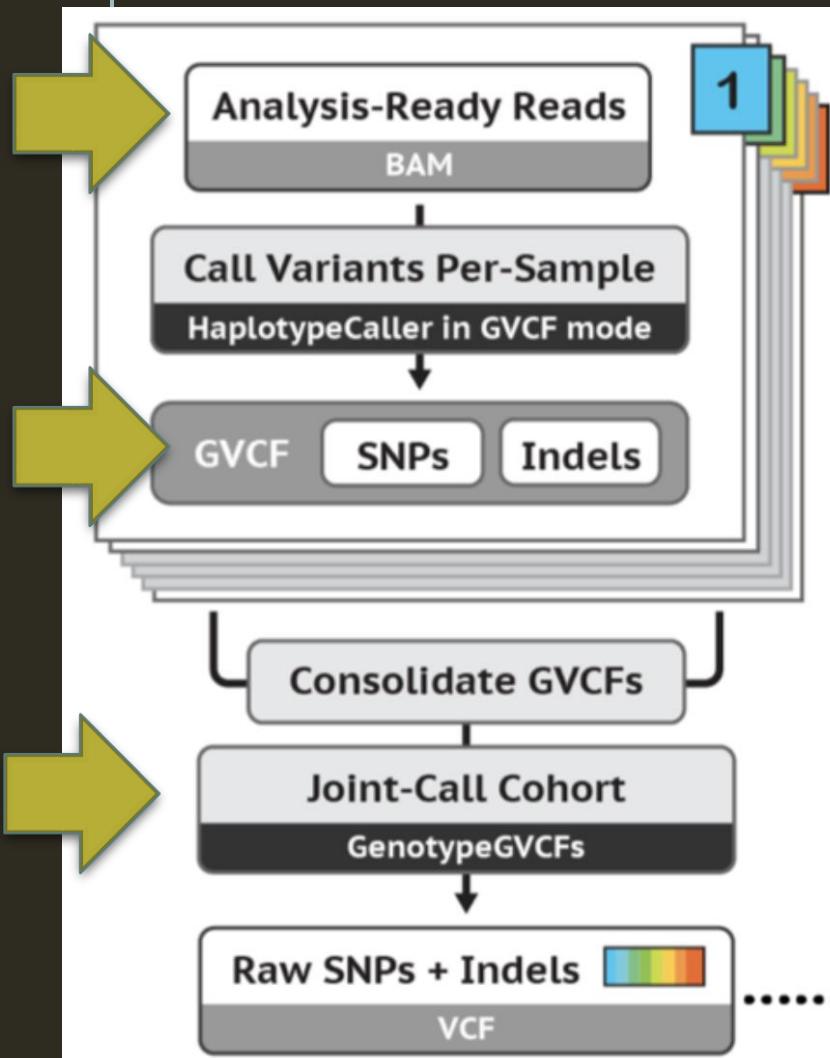


SNP and InDel calling

HaplotypeCaller „is capable of calling SNPs and indels simultaneously via local de-novo assembly of haplotypes in an active region. In other words, whenever the program encounters a region showing signs of variation, it discards the existing mapping information and completely reassembles the reads in that region. This allows the HaplotypeCaller to be more accurate when calling regions that are traditionally difficult to call, for example when they contain different types of variants close to each other.”

CombineGVCFs „combine per-sample gVCF files produced by HaplotypeCaller into a multi-sample gVCF file.”

DETEKCJA POLIMORFIZMÓW – PAKIET GATK

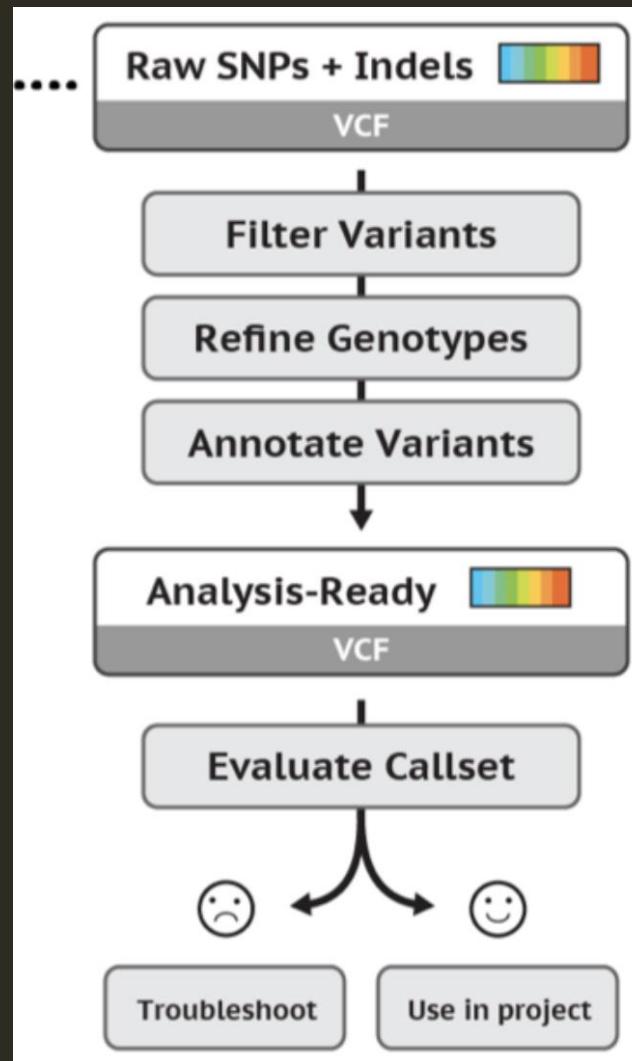


SNP and InDel calling

```
gatk --java-options "-Xmx4g" HaplotypeCaller \
-R Homo_sapiens_assembly38.fasta \
-I input.bam \
-O output.g.vcf.gz \
-ERC GVCF
```

```
gatk CombineGVCFs \
-R reference.fasta \
--variant sample1.g.vcf.gz \
--variant sample2.g.vcf.gz \
-O cohort.g.vcf.gz
```

DETEKCJA POLIMORFIZMÓW – PAKIET GATK



Filtrowanie zbioru polimorfizmów
Adnotacja
Sens biologiczny

DETEKCJA POLIMORFIZMÓW – SAMTOOLS

Samtools

Samtools is a suite of programs for interacting with high-throughput sequencing data. It consists of three separate repositories:

Samtools Reading/writing/editing/indexing/viewing SAM/BAM/CRAM format

BCFtools Reading/writing BCF2/VCF/gVCF files and calling/filtering/summarising SNP and short indel sequence variants

HTSlib A C library for reading/writing high-throughput sequencing data

Samtools and BCFtools both use HTSlib internally, but these three packages contain their own copies of htseq, so they can be built independently.

Download

Source code releases can be downloaded from [GitHub](#) or [Sourceforge](#):

 Source release details

Workflows

We have described some standard workflows using Samtools:

- FASTQ to BAM / CRAM
- WGS/WES Mapping to Variant Calls
- Filtering of VCF Files
- Using CRAM within Samtools

Documentation

- Manuals
- HowTos
- Specifications
- Duplicate Marking
- Zlib Benchmarks
- CRAM Benchmarks
- Publications

Support

- Mailing Lists
- HTSlib issues
- BCFtools issues
- Samtools issues

DETEKCJA POLIMORFIZMÓW – SAMTOOLS

WGS/WES Mapping to Variant Calls

The standard workflow for working with DNA sequence data consists of three major steps:

- Mapping BWA
- Improvement GATK
- Variant Calling bcftools
- <http://www.htslib.org/workflow/wgs-call.html>

The screenshot shows a portion of the Samtools website. At the top, there is a navigation bar with the Samtools logo on the left and links for Home, Download, Workflows, Documentation, and Support on the right. Below the navigation bar, the main content area has a dark background. The title "WGS/WES Mapping to Variant Calls" is displayed in large, bold, white font. Below the title, a sub-section title "The standard workflow for working with DNA sequence data consists of three major steps:" is shown in white. Underneath this, a bulleted list of four items is provided, each with a white bullet point and white text. The fourth item in the list is a link that is highlighted in yellow.

Samtools

Home Download ▾ Workflows ▾ Documentation ▾ Support ▾

WGS/WES Mapping to Variant Calls

The standard workflow for working with DNA sequence data consists of three major steps:

- Mapping
- Improvement
- Variant Calling
- <http://www.htslib.org/workflow/wgs-call.html>

DETEKCJA POLIMORFIZMÓW – SAMTOOLS

WGS/WES Mapping to Variant Calls

The standard workflow for working with DNA sequence data consists of three major steps:

- Mapping BWA
- Improvement GATK
- Variant Calling bcftools

Variant Calling

To convert your BAM file into genomic positions we first use mpileup to produce a BCF file that contains all of the locations in the genome. We use this information to call genotypes and reduce our list of sites to those found to be variant by passing this file into bcftools call.

You can do this using a pipe as shown here:

```
bcftools mpileup -Ou -f <ref.fa> <sample1.bam> <sample2.bam> <sample3.bam> | bcftools call -vm0 z -o <study.vcf.gz>
```

DETEKCJA POLIMORFIZMÓW – SAMTOOLS

WGS/WES Mapping to Variant Calls

The standard workflow for working with DNA sequences

- Mapping BWA
 - Improvement GATK
 - Variant Calling bcftools

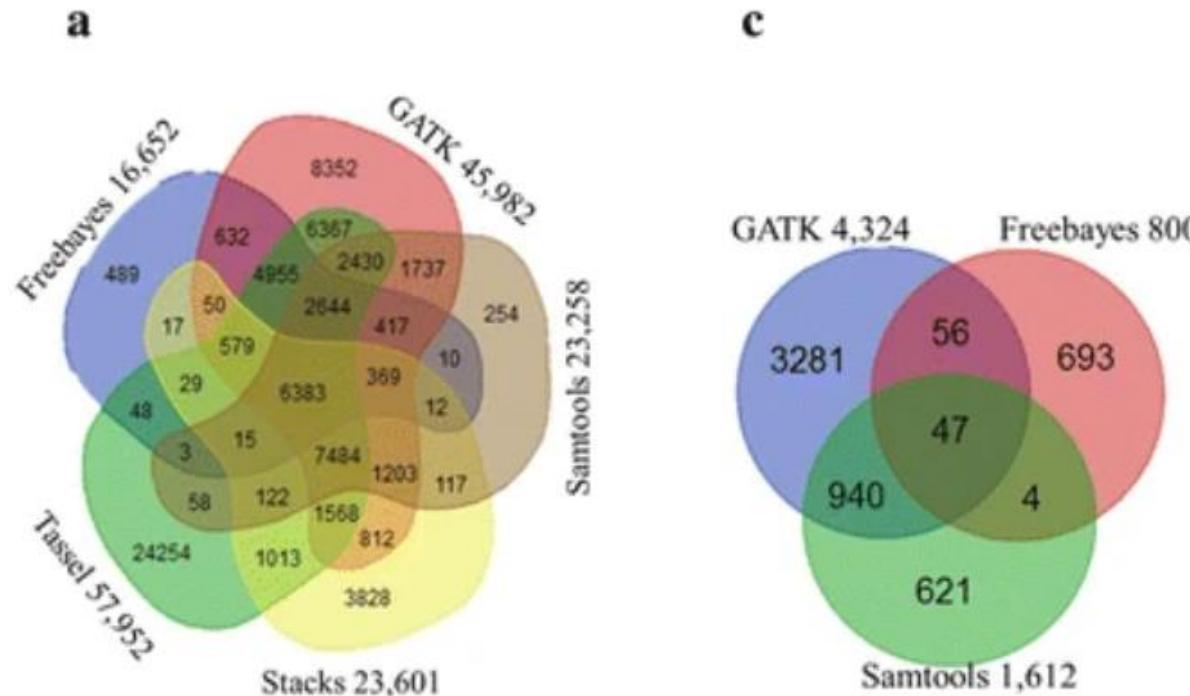
BCFtools

- Installation
 - Latest development version
 - Calling
 - CNV calling
 - Consequence calling
 - Consensus calling
 - ROH calling
 - Variant calling and filtering
 - Tips and Tricks
 - Converting formats
 - Extracting information
 - Filtering expressions
 - Plugins

Mining sequence variations in representative polyploid sugarcane germplasm accessions

SNP + INDELS

Xiping Yang ¹, Jian Song ¹, Qian You ¹, Dev R Paudel ¹, Jisen Zhang ², Jianping Wang ^{3 4 5}



The number of variants called in the 14 sugarcane germplasm accessions. **a.** Venn diagram showing overlapping SNPs among five genome reference-based callers, Tassel, Stacks, Samtools, GATK, and Freebayes. (...) **c.** Venn diagram showing overlapping short insertions or deletions (InDels) among three reference-based callers, Samtools, GATK, and Freebayes

```
##fileformat=VCFv4.2
##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

This example shows (in order): a good simple SNP, a possible SNP that has been filtered out because its quality is below 10, a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error), a site that is called monomorphic reference (i.e. with no alternate alleles), and a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T). Genotype data are given for three samples, two of which are phased and the third unphased, with per sample genotype quality, depth and haplotype qualities (the latter only for the phased samples) given as well as the genotypes. The microsatellite calls are unphased.

header section

polimorphisms (body)

VCF

#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE

Chr1 238 . C T 48 . DP=8;VDB=4.789490e-02;RPB=-1.551181e+00;AF1=0.5;AC1=1;DP4=4,1,3,0;MQ=44;FQ=51; PV4=1,1,0.16,0.41
GT:PL:GQ 0/1:78,0,135:81

#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

This example shows (in order): a good simple SNP, a possible SNP that has been filtered out because its quality is below 10, a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error), a site that is called monomorphic reference (i.e. with no alternate alleles), and a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T). Genotype data are given for three samples, two of which are phased and the third unphased, with per sample genotype quality, depth and haplotype qualities (the latter only for the phased samples) given as well as the genotypes. The microsatellite calls are unphased.

VCF

#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE

Chr1 238 . C T 48 . DP=8;VDB=4.789490e-02;RPB=-1.551181e+00;AF1=0.5;AC1=1;DP4=4,1,3,0;MQ=44;FQ=51; PV4=1,1,0.16,0.41
 GT:PL:GQ 0/1:78,0,135:81

#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

This example shows (in order): a good simple SNP, a possible SNP that has been filtered out because its quality is below 10, a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error), a site that is called monomorphic reference (i.e. with no alternate alleles), and a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T). Genotype data are given for three samples, two of which are phased and the third unphased, with per sample genotype quality, depth and haplotype qualities (the latter only for the phased samples) given as well as the genotypes. The microsatellite calls are unphased.

VCF

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	SAMPLE
Chr1	238	.	C	T	48	.	DP=8;VDB=4.789490e-02;RPB=-1.551181e+00;AF1=0.5;AC1=1;DP4=4,1,3,0;MQ=44;FQ=51; PV4=1,1,0.16,0.41	GT:PL:GQ	0/1:78,0,135:81

#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

This example shows (in order): a good simple SNP, a possible SNP that has been filtered out because its quality is below 10, a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error), a site that is called monomorphic reference (i.e. with no alternate alleles), and a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T). Genotype data are given for three samples, two of which are phased and the third unphased, with per sample genotype quality, depth and haplotype qualities (the latter only for the phased samples) given as well as the genotypes. The microsatellite calls are unphased.

VCF

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	SAMPLE
Chr1	238	.	C	T	48	.	DP=8;VDB=4.789490e-02;RPB=-1.551181e+00;AF1=0.5;AC1=1;DP4=4,1,3,0;MQ=44;FQ=51; PV4=1,1,0.16,0.41	GT:PL:GQ	0/1:78,0,135:81

#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

This example shows (in order): a good simple SNP, a possible SNP that has been filtered out because its quality is below 10, a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error), a site that is called monomorphic reference (i.e. with no alternate alleles), and a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T). Genotype data are given for three samples, two of which are phased and the third unphased, with per sample genotype quality, depth and haplotype qualities (the latter only for the phased samples) given as well as the genotypes. The microsatellite calls are unphased.

VCF

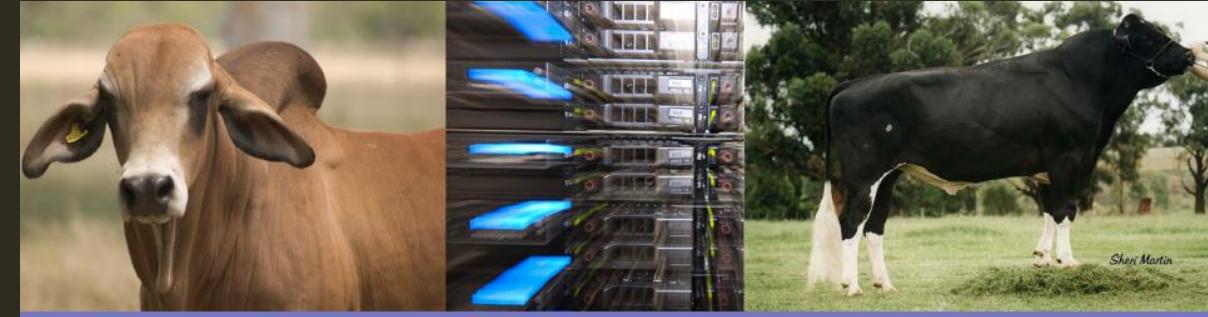
SINGLE- AND MULTI-SAMPLE

#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

This example shows (in order): a good simple SNP, a possible SNP that has been filtered out because its quality is below 10, a site at which two alternate alleles are called, with one of them (T) being ancestral (possibly a reference sequencing error), a site that is called monomorphic reference (i.e. with no alternate alleles), and a microsatellite with two alternative alleles, one a deletion of 2 bases (TC), and the other an insertion of one base (T). Genotype data are given for three samples, two of which are phased and the third unphased, with per sample genotype quality, depth and haplotype qualities (the latter only for the phased samples) given as well as the genotypes. The microsatellite calls are unphased.

VARIANTS FILTERING

EXAMPLES



The 1000 bull genomes project

- Removal of variants with 2 or more alternative alleles
- Overall quality (QUAL)
- Setting minimum and maximum read depths for filtering (DP)
 - Set minimum number as 10 across all animals
 - Set maximum as: median read depth + 3 * standard deviation read depth
- Remove variants with the same basepair position
 - If two variants have the same bp position both are removed
 - Resolves issue with SNP and INDEL calls at same position
- (...)

VARIANTS FILTERING

Toward better understanding of artifacts in variant calling from high-coverage samples

Heng Li¹

Several universal filters applicable to most callers:

- **Low-complexity (LC) filter:** filtering variants overlapping with low-complexity regions (LCRs)
- **Maximum depth (MD) filter:** filtering sites covered by excessive number of reads. It should be noted that different callers may define the depth differently. For example, Platypus apparently only counts reads with unambiguous realignment. The read depth reported in the Platypus VCF is noticeably smaller in comparison with other callers.
- **Quality filter (QU):** filtering sites with the reported variant quality below a threshold.
- others

VARIANTS FILTERING

Next-generation data filtering in the genomics era

William Hemstrom , Jared A. Grummer, Gordon Luikart & Mark R. Christie 

Nature Reviews Genetics 25, 750–767 (2024) | [Cite this article](#)

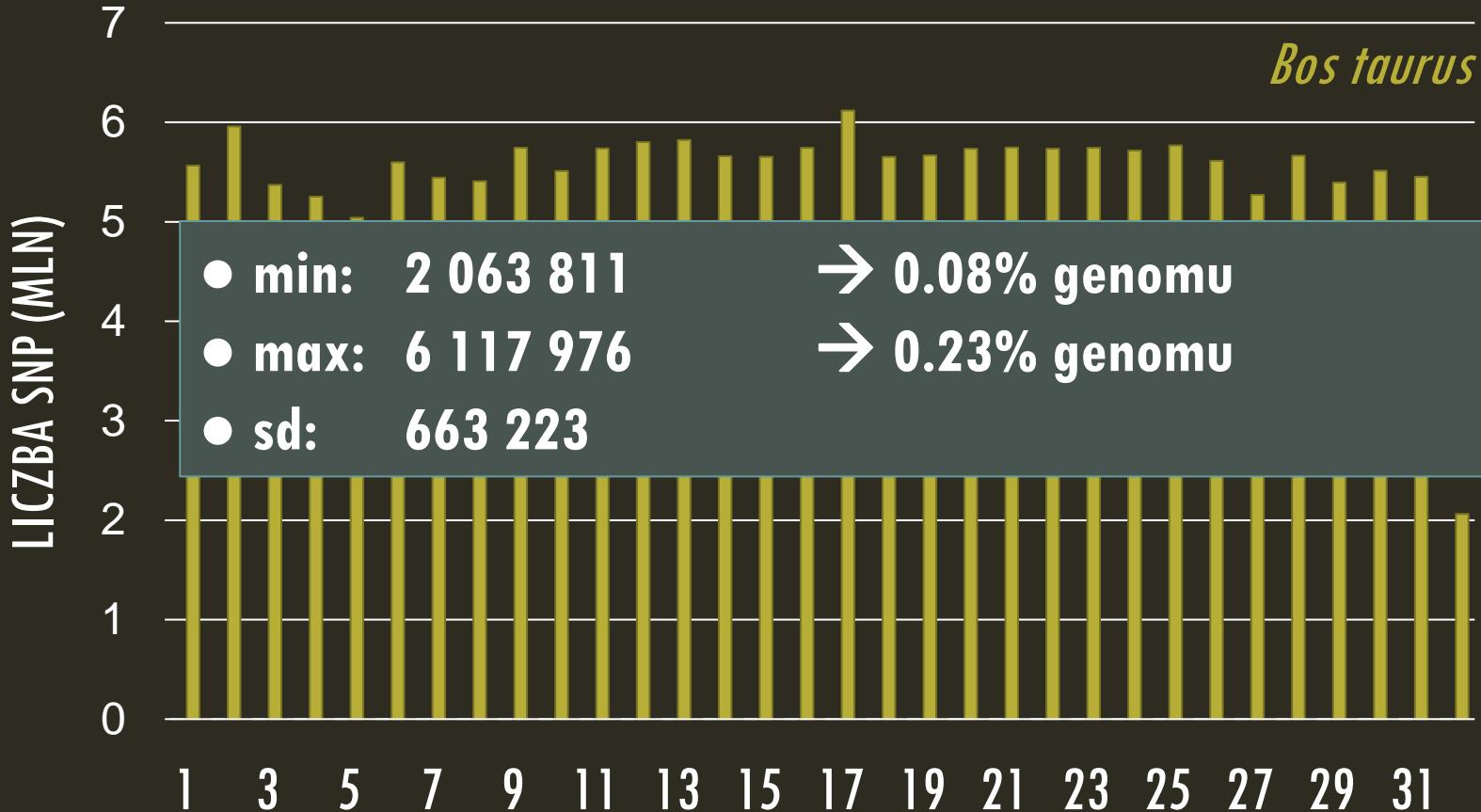
12k Accesses | 62 Altmetric | [Metrics](#)

Abstract

Genomic data are ubiquitous across disciplines, from agriculture to biodiversity, ecology, evolution and human health. However, these datasets often contain noise or errors and are missing information that can affect the accuracy and reliability of subsequent computational analyses and conclusions. A key step in genomic data analysis is filtering – removing sequencing bases, reads, genetic variants and/or individuals from a dataset – to improve data quality for downstream analyses. Researchers are confronted with a multitude of choices when filtering genomic data; they must choose which filters to apply and select appropriate thresholds. To help usher in the next generation of genomic data filtering, we review and suggest best practices to improve the implementation, reproducibility and reporting standards for filter types and thresholds commonly applied to genomic datasets. We focus mainly on filters for minor allele frequency, missing data per individual or per locus, linkage disequilibrium and Hardy–Weinberg deviations. Using simulated and empirical datasets, we illustrate the large effects of different filtering thresholds on common population genetics statistics, such as Tajima’s D value, population differentiation (F_{ST}), nucleotide diversity (π) and effective population size (N_e).



DETEKCJA POLIMORFIZMÓW → SNP



Surowe Dane

Kontrola jakości

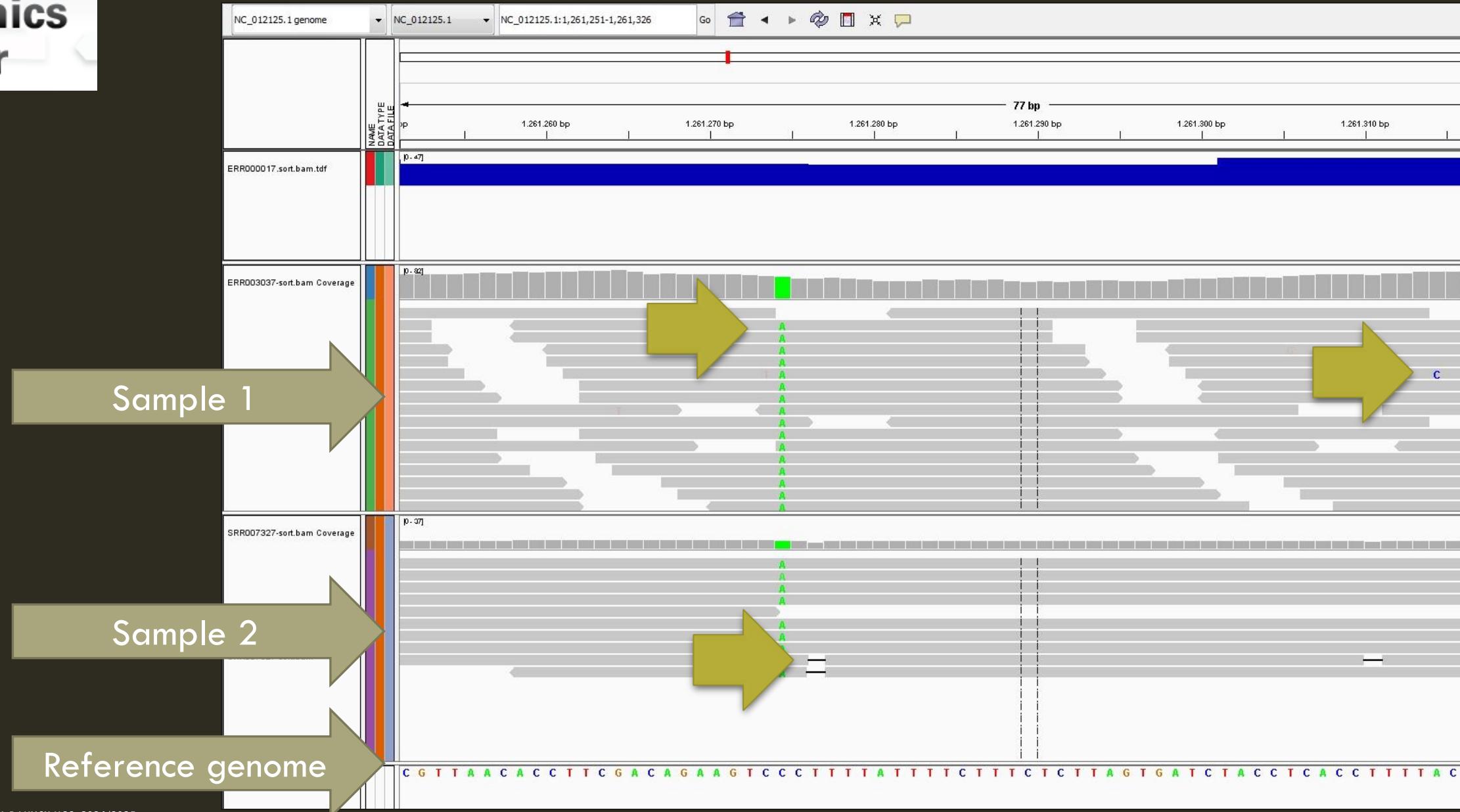
Przyrównanie do
genomu
referencyjnego

Detekcja
polimorfizmów

Sens biologiczny

Integrative Genomics Viewer

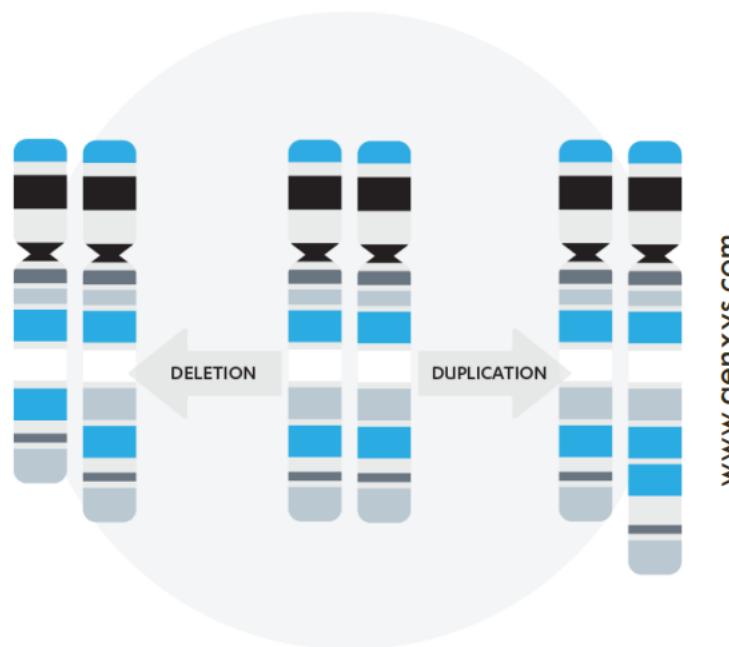
<http://software.broadinstitute.org/software/igv/>



DETEKCJA CNV



DETEKCJA POLIMORFIZMÓW → CNV

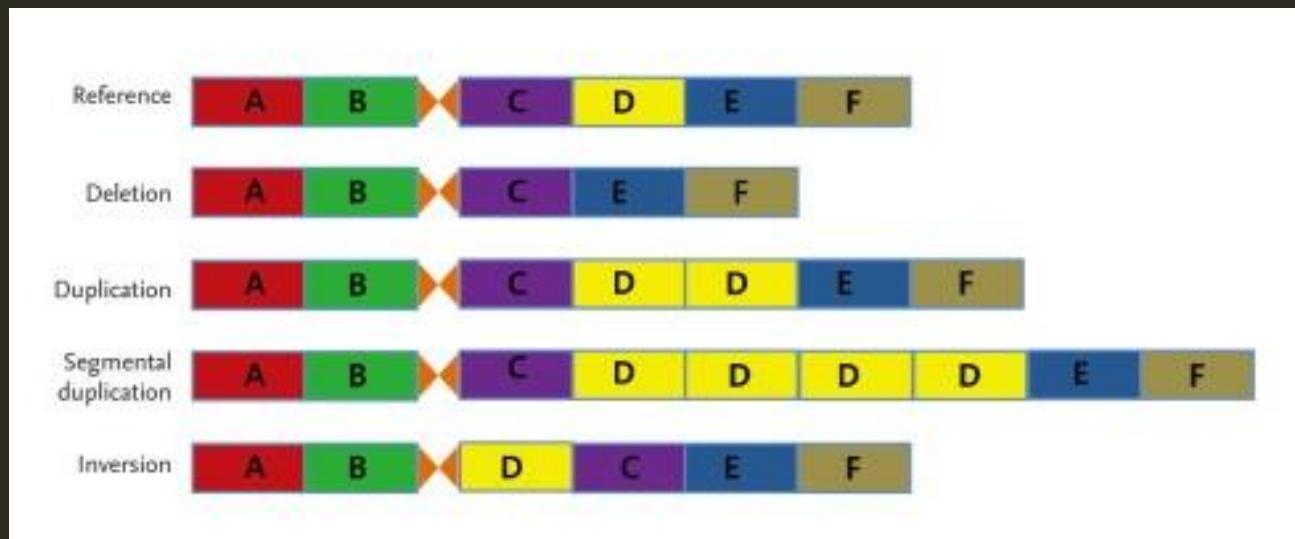


www.genixys.com

Rozmiar: bp - Mbp

COPY NUMBER VARIATION (CNV)

Zmienność liczby kopii obejmuje **duplikacje i delekcje dłuższe niż 1 000 pz** (różne definicje).



„Different types of CNVs and an example of genome-wide detection of CNVs. The plot illustrates deletion, duplication, and multiple segmental duplication of the "D" locus compared with the reference genome. Inversion of "C" and "D" loci is also illustrated.”

COPY NUMBER VARIATION (CNV)

CNVs → **gains** and **losses** (insertions and deletions) of genomic sequence greater than **50 bp** between two individuals of a species (Mills et al. 2011)

SNPs more frequent than CNVs

CNVs higher percentage of genomic sequence, greater effects, including the changing of gene structure and expression

Interrogation of the genome for both common and rare variations, including SNPs and CNVs, was proposed as an effective way to elucidate the causes of complex disease and traits (Manolio et al. 2009)

CNV → FENOTYP

CNVs:

- zmieniają strukturę i ekspresję genów
- są związane z wieloma chorobami
- mogą wpływać pozytywnie na proces adaptacji do zmieniającego się środowiska

Klopocki and Mundlos 2011. Copy-Number Variations, Noncoding Sequences, and Human Phenotypes. Annual Review of Genomics and Human Genetics. 12:1, 53-72.

Locus	CNV	Target gene	Disorder
2q31-q32	dup	<i>HOXD</i>	Mesomelic dysplasia Kantaputra type
2q35	dup	<i>IHH</i>	Syndactyly type 1, craniosynostosis Philadelphia type
7q36	dup dup dup dup	<i>SHH</i>	Triphalangeal thumb-polysyndactyly syndrome Preaxial polydactyly type 2 Werner mesomelic syndrome Haas polysyndactyly Laurin-Sandrow syndrome
17q24.2-q24.3	dup		Congenital generalized hypertrichosis terminalis with or without gingival hyperplasia
17q24.3	dup del dup	<i>SOX9</i>	Cooks syndrome Pierre Robin syndrome Female-to-male sex reversal
20p12	dup	<i>BMP2</i>	Brachydactyly type A2
Xp	del	<i>SHOX</i>	Idiopathic short stature, Léri-Weill dyschondrosteosis
Xq22	del dup	<i>PLP1</i>	Pelizaeus-Merzbacher disease Spastic paraplegia type 2
Xp21.2	del	<i>DAX1</i>	Male-to-female sex reversal

CNV → FENOTYP



Table 2 Functional impacts of cattle CNV

References	Phenotype	Mechanism
Ohba et al. (2000) and Hirano et al. (2000)	Renal tubular dysplasia	Deletions of <i>CLDN-16</i> gene: Type 1—Deletion of 37 kb (exons 1 to 4); Type 2—56 kb (exons 1 to 4 and 821 bp of exon 5)
Drogemuller et al. (2001)	Anhidrotic ectodermal dysplasia	A deletion including exon 3 of the <i>EDA</i> gene
Sugimoto et al. (2003)	Myopathy of diaphragmatic muscles	A deletion of the <i>HSPA1B</i> gene
Dreger and Schmutz (2010)	Variant Red	Partially or completely duplicated <i>DEFB103</i> (β Defensin 3)
Meyers et al. (2010)	Osteopetrosis	A 2.8-kb deletion including exons 2 and 3 of the <i>SLC4A2</i> gene
Flisikowski et al. (2010)	Abortions and stillbirths	A 110-kb deletion of exons 3 and 4 of the <i>MIMT1</i> gene
Hou et al. (2011a, b)	Parasite resistance	A case-control study of CNV events associated with parasite resistance
Durkin et al. (2012)	Color sidedness	Serial translocation by means of circular intermediates encompassing <i>KIT</i>

METODY DETEKCJI CNV

PEM/RP

Paired-end mapping

SR

Split-read

DOC/RD/RC

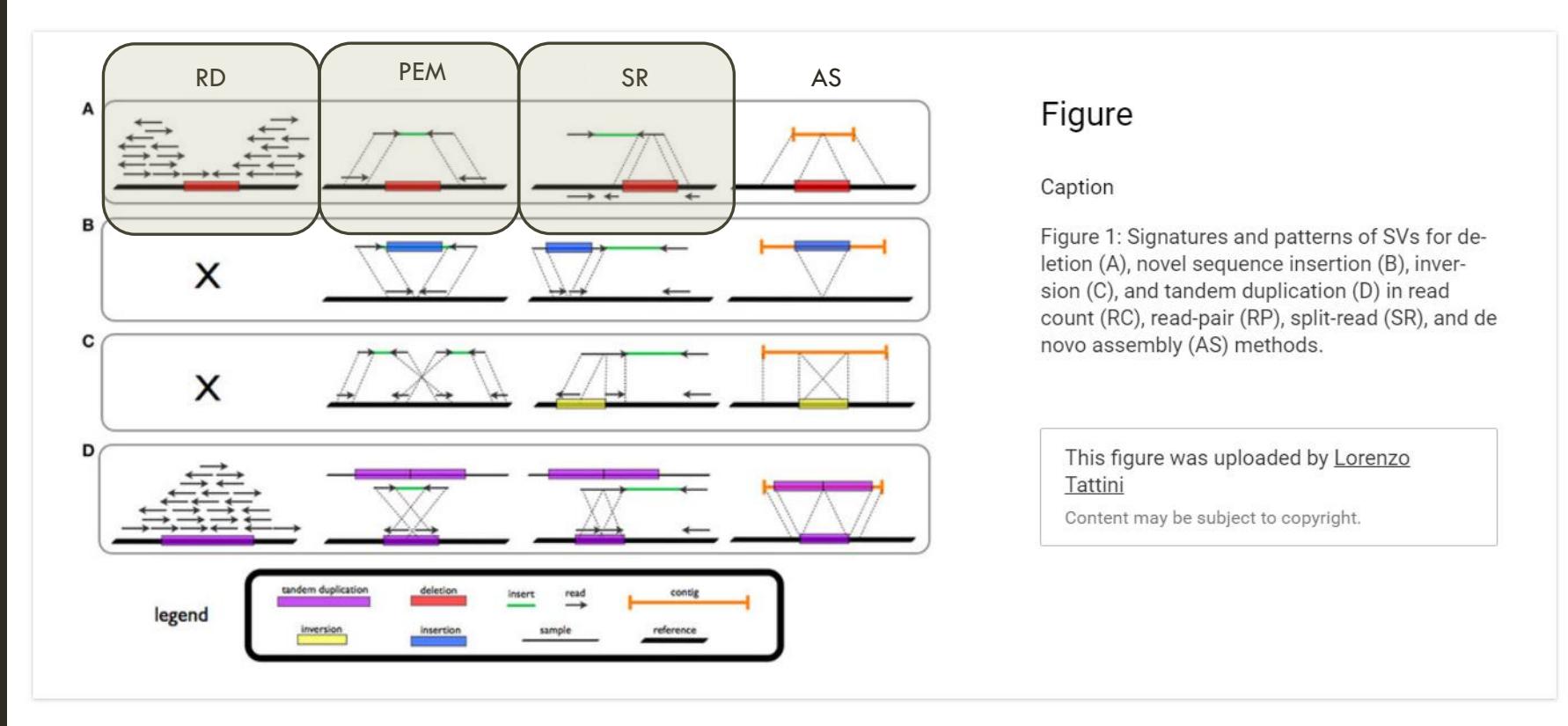
Depth of coverage/Read Depth/
Read Count

AS

Assembly- based

CB

Combinations of Methods



Figure

Caption

Figure 1: Signatures and patterns of SVs for deletion (A), novel sequence insertion (B), inversion (C), and tandem duplication (D) in read count (RC), read-pair (RP), split-read (SR), and de novo assembly (AS) methods.

This figure was uploaded by [Lorenzo Tattini](#)

Content may be subject to copyright.

METODY DETEKCJI CNV

PEM/RP

- limited by the insert size
- detects short CNVs

SR

- applicable to the unique regions in the reference genome

DOC/RD/RC

- detect larger CNVs in complex genomic region classes, which are difficult to detect using PEM and SR

AS

- overwhelming demand on computational resources
- perform poorly for repeated regions

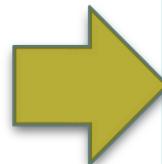
CB

- none of tools is able to detect the full spectrum of all types of CNVs with high sensitivity and specificity



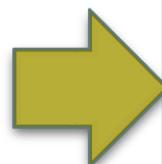
An evaluation of copy number variation detection tools for cancer using whole exome sequencing data

Fatima Zare¹, Michelle Dow², Nicholas Monteleone¹, Abdelrahman



Abstract

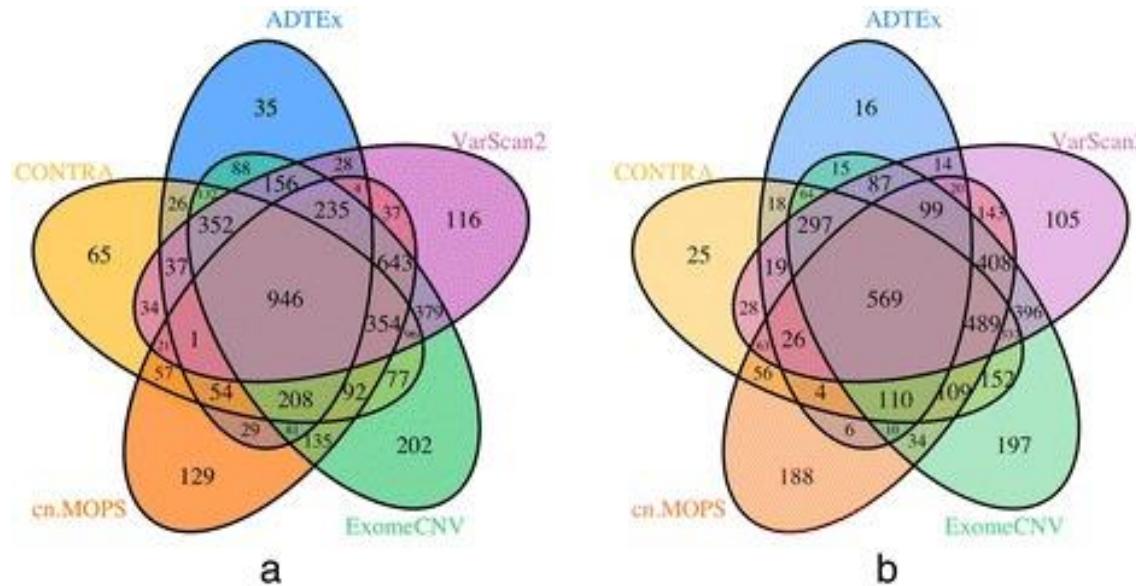
Background: Recently copy number variation (CNV) has gained considerable interest as a type of genomic/genetic variation that plays an important role in disease susceptibility. Advances in sequencing technology have created an opportunity for detecting CNVs more accurately. Recently whole exome sequencing (WES) has become primary strategy for sequencing patient samples and study their genomics aberrations. However, compared to whole genome sequencing, WES introduces more biases and noise that make CNV detection very challenging. Additionally, tumors' complexity makes the detection of cancer specific CNVs even more difficult. Although many CNV detection tools have been developed since introducing NGS data, there are few tools for somatic CNV detection for WES data in cancer.



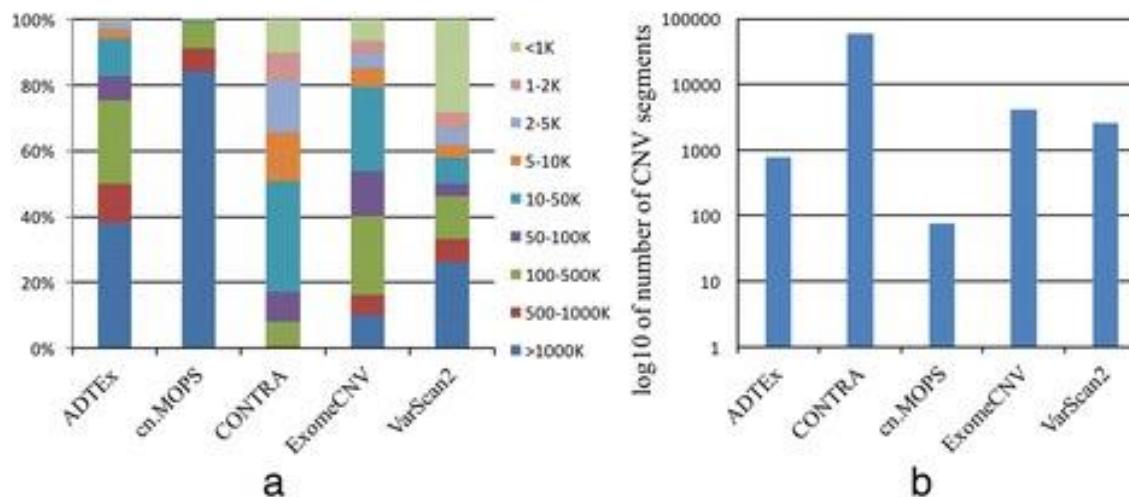
Results: In this study, we evaluated the performance of the most recent and commonly used CNV detection tools for WES data in cancer to address their limitations and provide guidelines for developing new ones. We focused on the tools that have been designed or have the ability to detect cancer somatic aberrations. We compared the performance of the tools in terms of sensitivity and false discovery rate (FDR) using real data and simulated data. Comparative analysis of the results of the tools showed that there is a low consensus among the tools in calling CNVs. Using real data, tools show moderate sensitivity (~50% - ~80%), fair specificity (~70% - ~94%) and poor FDRs (~27% - ~60%). Also, using simulated data we observed that increasing the coverage more than 10x in exonic regions does not improve the detection power of the tools significantly.

Conclusions: The limited performance of the current CNV detection tools for WES data in cancer indicates the need for developing more efficient and precise CNV detection methods. Due to the complexity of tumors and high level of noise and biases in WES data, employing advanced novel segmentation, normalization and de-noising techniques that are designed specifically for cancer data is necessary. Also, CNV detection development suffers from the lack of a gold standard for performance evaluation. Finally, developing tools with user-friendly user interfaces and visualization features can enhance CNV studies for a broader range of users.

Keywords: Copy number variation, Whole-exome sequencing, Somatic aberrations, Cancer

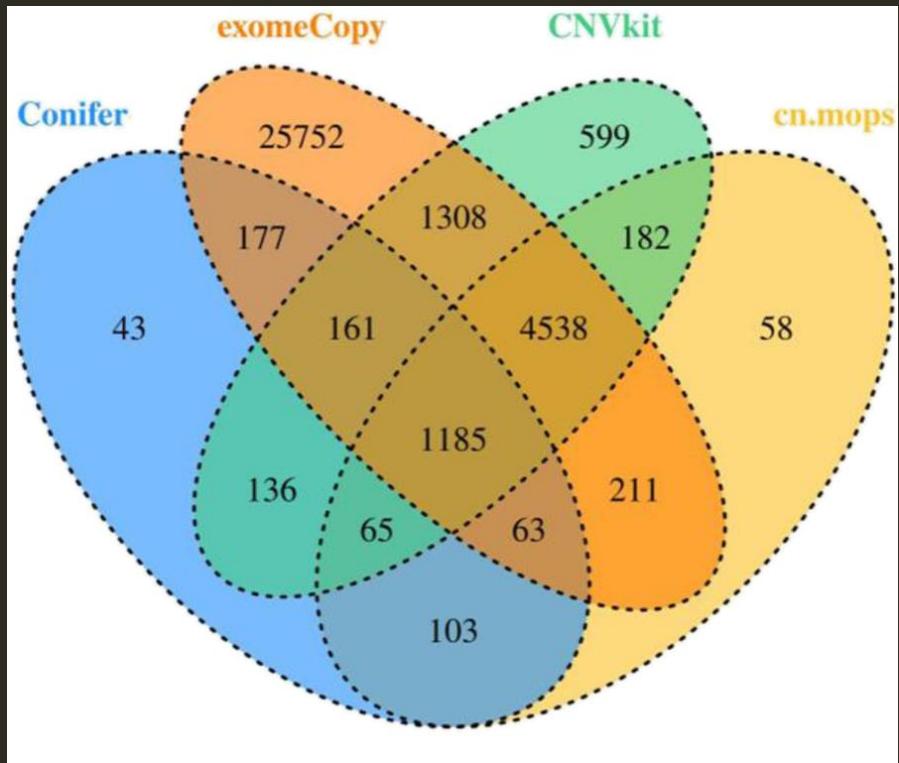


Venn diagrams of the average of the number of truly detected CNV genes from the 5 tools, **(a)** amplified genes, **(b)** deleted genes.



Characteristics of the detected CNV regions by the 5 tools. **(a)** Size distributions of CNV segments, **(b)** Number of detected CNV segments

ZGODNOŚĆ PROGRAMÓW DO DETEKCJI CNV



The overlapping consistency results. Fig shows the simulated data.

Zhao, L., Liu, H., Yuan, X. et al. Comparative study of whole exome sequencing-based copy number variation detection tools. BMC Bioinformatics 21, 97 (2020). <https://doi.org/10.1186/s12859-020-3421-1>



► Cancers (Basel). 2021 Dec 14;13(24):6283. doi: [10.3390/cancers13246283](https://doi.org/10.3390/cancers13246283)

A Comparison of Tools for Copy-Number Variation Detection in Whole Exome and Whole Genome Sequencing Data

Migle Gabrielaite ^{1,†}, Mathias Husted Torp ^{1,†}, Malthe Sebro Rasmussen ¹, Sergio Andre Vieira ¹, Christina Bligaard Pedersen ^{1,2}, Savvas Kinalis ¹, Majbritt Busk Madsen ¹, Miya Demircan ¹, Arman Simonyan ¹, Christina Westmose Yde ¹, Lars Rønn Olsen ^{1,2}, Rasmus Maria Rossing ^{1,3}, Finn Cilius Nielsen ¹, Ole Winther ^{1,4,5}, Frederik Otzen Bagger ^{1,6,7,*}

Editor: Dimitrios H Roukos

► Author information ► Article notes ► Copyright and License information

PMCID: PMC8699073 PMID: [34944901](https://pubmed.ncbi.nlm.nih.gov/34944901/)

5. Conclusions

In summary, by reviewing 50 tools for CNV calling, of which 11 were included for a benchmark (CLC Genomics Workbench (WGS and WES), cn.MOPS (WGS and WES), CNVkit (WES), CNVnator (WGS), CODEX2 (WGS), Control-FREEC (WGS), DELLY (WGS), ExomeDepth (WES), GATK gCNV (WES and WGS), Lumpy (WGS), and Manta (WES and WGS)), we conclude that CNV identification from NGS data remains challenging. For the best reliability of CNV calling from NGS data, we observed that even if the tools were developed for WES data or allowed it as input, they did not perform well. We suggest WGS as the only NGS-based option for broad calling of CNVs. Furthermore, low precision in all tools leads us to recommend a hypothesis-based approach for finding causative CNVs by NGS in the clinic, and further validation of these candidates by manual inspection, MLPA or array-based approaches. If multiple samples are available from the same protocol, we suggest using these to filter by commonly called CNVs. If only the WGS data is available for the sample, for a higher precision of CNV calls, multiple CNV calling tools should be used. We suggest combining tools which have the best recall (GATK gCNV, Lumpy, DELLY, and cn.MOPS) using consensus callers (e.g., [53]) and prioritize the CNV calls made by such tools.

Combining callers improves the detection of copy number variants from whole-genome sequencing

[Marie Coutelier](#), [Manuel Holtgrewe](#), [Marten Jäger](#), [Ricarda Flöttman](#), [Martin A. Mensah](#), [Malte Spielmann](#)

[Peter Krawitz](#), [Denise Horn](#), [Dieter Beule](#) & [Stefan Mundlos](#)✉

European Journal of Human Genetics **30**, 178–186 (2022) | [Cite this article](#)

15k Accesses | 33 Citations | 7 Altmetric | [Metrics](#)

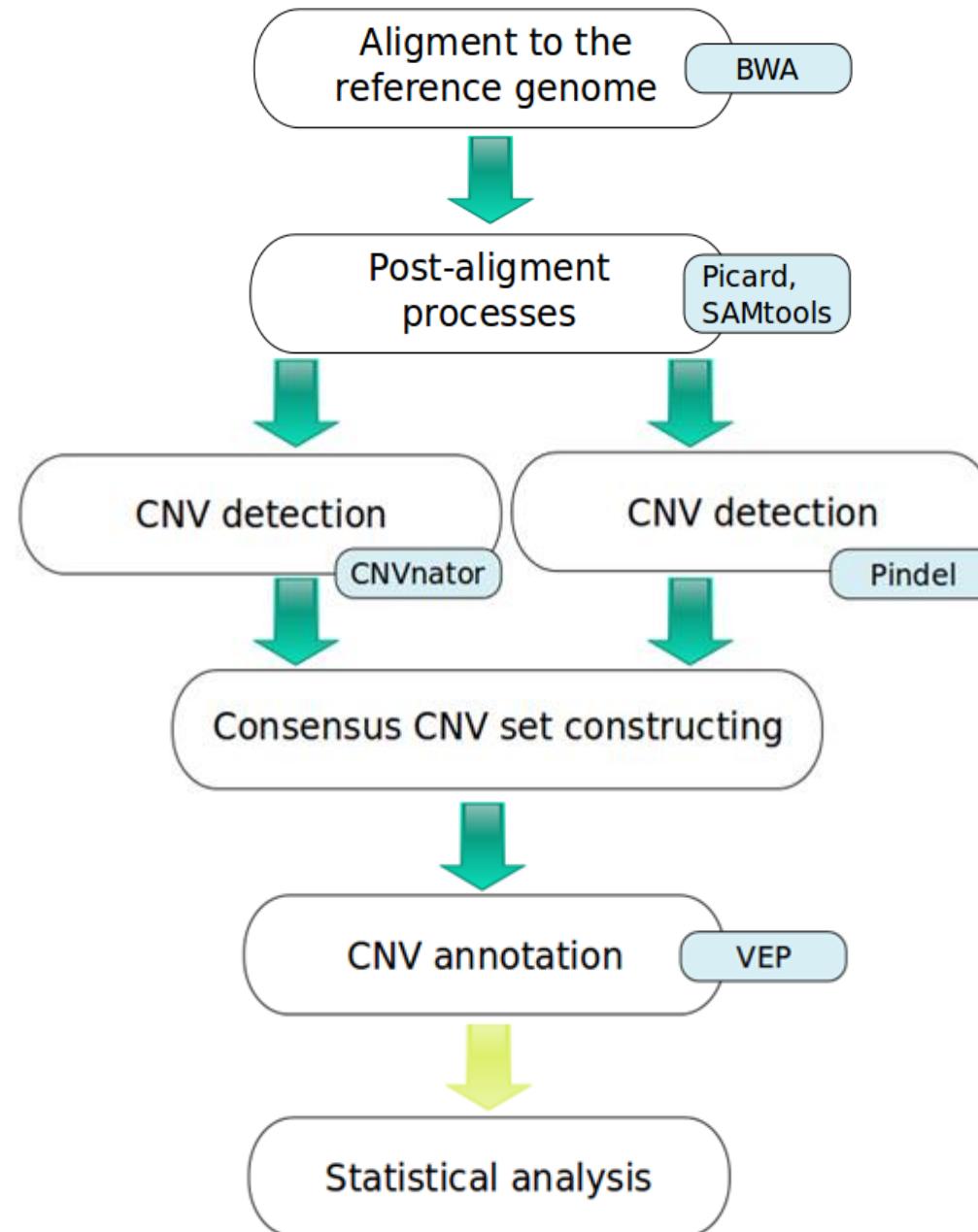
Abstract

Copy Number Variants (CNVs) are deletions, duplications or insertions larger than 50 base pairs. They account for a large percentage of the normal genome variation and play major roles in human pathology. While array-based approaches have long been used to detect them in clinical practice, whole-genome sequencing (WGS) bears the promise to allow concomitant exploration of CNVs and smaller variants. However, accurately calling CNVs from WGS remains

a difficult computational task, for which a consensus is still lacking. In this paper, we explore practical calling options to reach the best compromise between sensitivity and sensibility. We show that callers based on different signal (paired-end reads, split reads, coverage depth) yield

complementary results. We suggest approaches combining four selected callers (Manta, Delly, ERDS, CNVnator) and a regenotyping tool (SV2), and show that this is applicable in everyday practice in terms of computation time and further interpretation. We demonstrate the

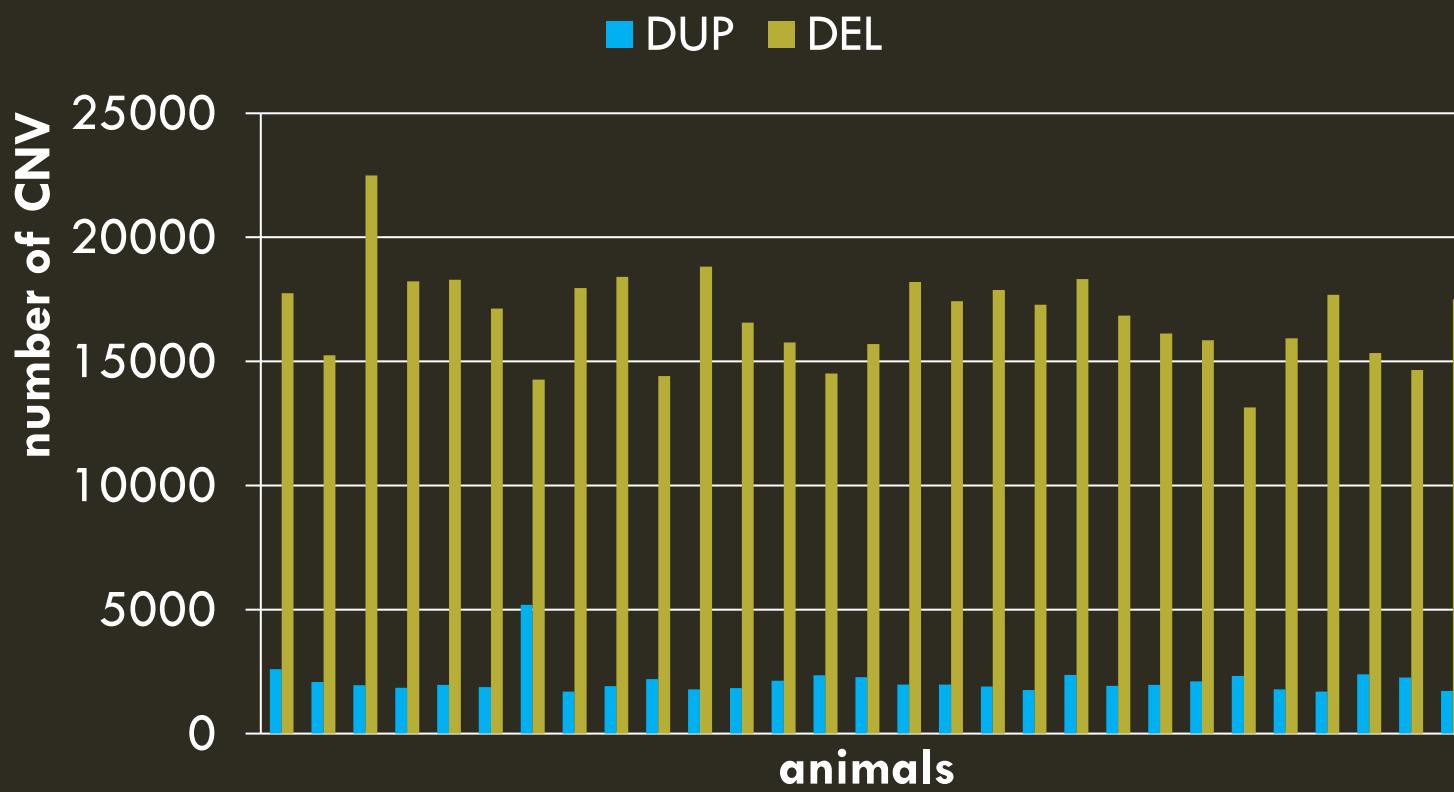
superiority of these approaches over array-based Comparative Genomic Hybridization (aCGH), specifically regarding the lack of resolution in breakpoint definition and the detection of potentially relevant CNVs. Finally, we confirm our results on the NA12878 benchmark genome, as well as one clinically validated sample. In conclusion, we suggest that WGS constitutes a timely and economically valid alternative to the combination of aCGH and whole-exome sequencing.



CNV PIPELINE

→ PRZYKŁAD OD THETA

CNV



The min and max number

Duplications: 1 694 - 5 187

Deletions: 13 149 - 22 496

The occupancy (%)

Duplications:

min 0.1 (BTA20, BTA22)

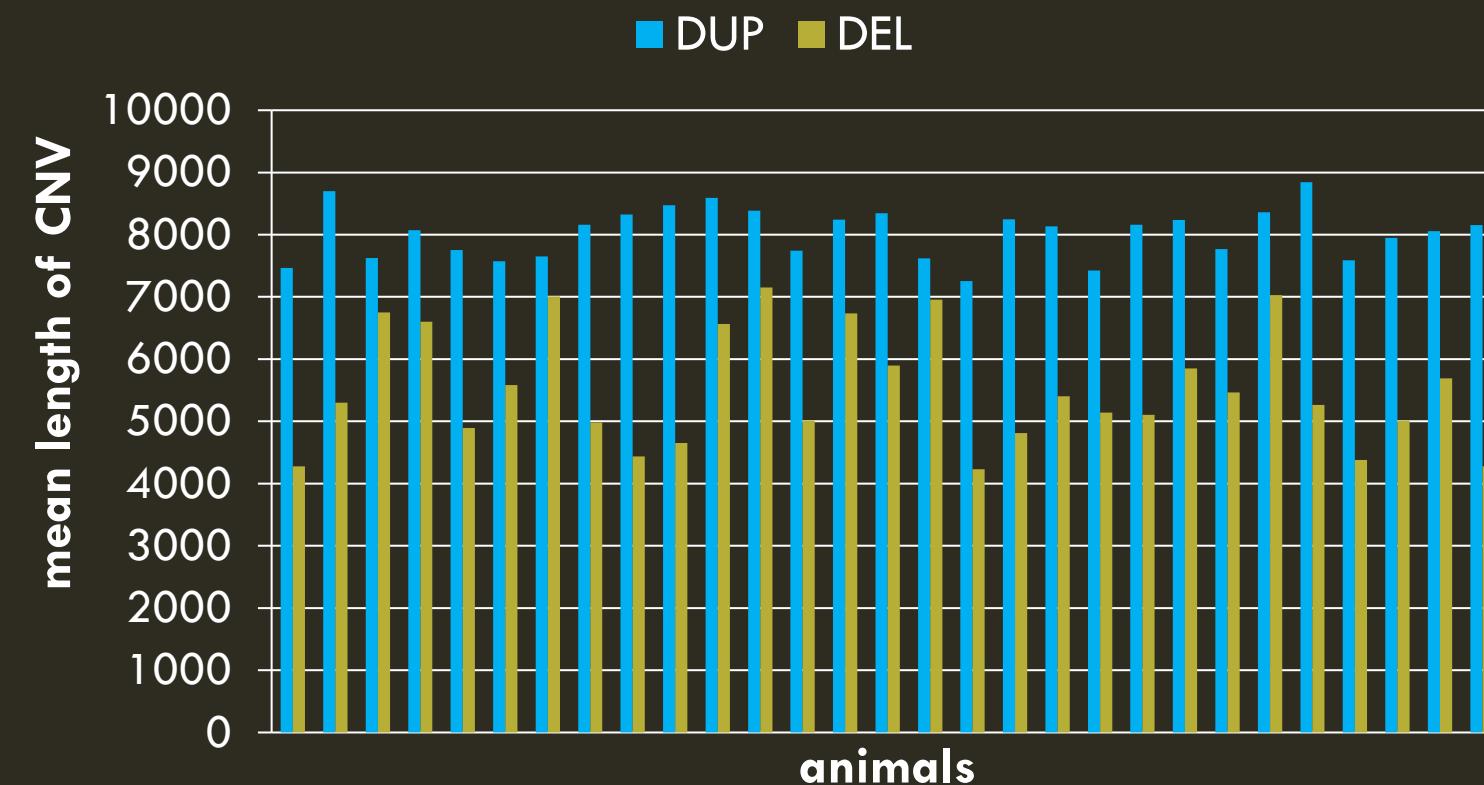
max 2.9 (BTA27)

Deletions:

min 1.2 (BTA24)

max 15.8 (BTA25)

CNV LENGTH



The average length (bp)
Duplications: 7 254 - 8 843
Deletions: 4 233 - 7 154

The min and max lengths (bp)
Duplications: 200 - 439 300
Deletions: 200 - 724 000

FORMATY DANYCH WYJŚCIOWYCH



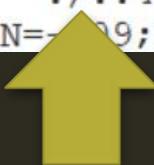
FORMATY DANYCH WYJŚCIOWYCH: PINDEL

```
mielczarek@eagle: ~
#####
# 1       D 2      NT 0  "" ChrID Chr1      BP 510  513      BP_range 510    513      Supports 1      1      +
# 1       1      - 0      0      S1 2      SUM_MS 60      1      NumSupSamples 1 1      AVIESPM000000003134
# 74 74 1 1 0 0
#TGAGAGTGTATCCTCTACTCTAGATTTAATCTTGCTTTGGTATTGTTATCAATTGTACCTTAAGAACATCTGCAGTATCCATTCTCTTagGGG
#TGTGATTACTGGCTTTATTGCTCTCCTCTTTGACTCTCCCTTTCTCCCAGTCAGCTCTTCTCCTCCCTCCCCCTCTCTACTT
#                                ATCTTGCTTTGGTATTGTTATCAATTGTACCTTAAGAACATCTGCAGTATCCATTCTCTT      GGG
#TGTGATTACTGGCTTTATTGCTCTCCTC
#                                229      60      AVIESPM000000003134      @HWI-D00616:22:C536CACXX:1:1311:15774:58575/2
#####
# 2       D 1      NT 0  "" ChrID Chr1      BP 964  966      BP_range 964    969      Supports 1      1      +
# 1       1      - 0      0      S1 2      SUM_MS 53      1      NumSupSamples 1 1      AVIESPM000000003134
# 76 76 1 1 0 0
#TAGAATAAGACTGAAAGTTAGAGACAGGGAGGATTAAATCCAAAATTGAGAACACCAGGAAACTCCTGACTCCAGGAACATTAATCAACAAGAGCTCATCCaAAAG
#CCTCCATACCTACACGGAAACCAAGCTCCATCCAAGAGCCAACAAGTTCCAGATCAAGACATACCATGCTAATTCTCCAACACATAGGAACATAG
#                                AGGAGGATTAAATCCAAAATTGAGAACACCAGGAAACTCCTGACTCCAGGAACATTAATCAACAAGAGCTCATCC AAAG
#CCTCCATACCTACACGGAAAC
#                                752      53      AVIESPM000000003134      @HWI-D00616:22:C536CACXX:1:2316:1508:50443/1
#####
```

FORMATY DANYCH WYJŚCIOWYCH: CNVNATOR

FORMATY DANYCH WYJŚCIOWYCH: LUMPY (VCF)

```
magda@penguin:~/WGS_CNV_beans/Lumpy_files/Lumpy_out
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT S7_CATCAA
LG1 786100 1 N <DUP> . . SVTYPE=DUP;STRANDS=-+:6;SVLEN=330;END=786430;CIPOS
=-288,9;CIEND=-10,273;CIPOS95=-68,2;CIEND95=-2,89;IMPRECISE;SU=6;PE=6;SR=0 GT:SU:PE:SR ./.:6:6:0
LG1 357821 2 N <DEL> . . SVTYPE=DEL;STRANDS=-+:5;SVLEN=-36;END=357857;CIPOS
=0,7;CIEND=-9,0;CIPOS95=0,1;CIEND95=-9,0;IMPRECISE;SU=5;PE=0;SR=5 GT:SU:PE:SR ./.:5:0:5
LG1 411115 3 N <DUP> . . SVTYPE=DUP;STRANDS=-+:7;SVLEN=659;END=411774;CIPOS
=-190,9;CIEND=-10,36;CIPOS95=-64,2;CIEND95=-4,11;IMPRECISE;SU=7;PE=7;SR=0 GT:SU:PE:SR ./.:7:7:0
LG1 506922 4 N <DEL> . . SVTYPE=DEL;STRANDS=-+:5;SVLEN=-664;END=507586;CIPO
S=-10,8;CIEND=-10,9;CIPOS95=-1,1;CIEND95=-1,1;IMPRECISE;SU=5;PE=3;SR=2 GT:SU:PE:SR ./.:5:3:2
LG1 507175 5 N <DEL> . . SVTYPE=DEL;STRANDS=-+:5;SVLEN=-1220;END=508395;CIP
OS=-10,300;CIEND=-297,9;CIPOS95=-1,94;CIEND95=-97,1;IMPRECISE;SU=5;PE=5;SR=0 GT:SU:PE:SR ./.:5:5:0
LG1 584355 6 N <DUP> . . SVTYPE=DUP;STRANDS=-+:5;SVLEN=636;END=584991;CIPOS
=-131,9;CIEND=-10,107;CIPOS95=-50,2;CIEND95=-3,28;IMPRECISE;SU=5;PE=5;SR=0 GT:SU:PE:SR ./.:5:5:0
LG1 584536 7 N <DEL> . . SVTYPE=DEL;STRANDS=-+:4;SVLEN=-464;END=585000;CIPO
S=-10,311;CIEND=-301,9;CIPOS95=0,120;CIEND95=-103,1;IMPRECISE;SU=4;PE=4;SR=0 GT:SU:PE:SR ./.:4:4:0
LG1 796763 8 N <DEL> . . SVTYPE=DEL;STRANDS=-+:4;SVLEN=-23;END=796786;CIPOS
=-10,7;CIEND=-8,9;CIPOS95=-1,0;CIEND95=0,1;IMPRECISE;SU=4;PE=0;SR=4 GT:SU:PE:SR ./.:4:0:4
LG1 905653 9 N <DEL> . . SVTYPE=DEL;STRANDS=-+:6;SVLEN=-29;END=906162;CIPO
```





LONG READS



LONG READS

Evaluating nanopore sequencing data processing pipelines for structural variation identification

Anbo Zhou¹, Timothy Lin¹ and Jinchuan Xing^{1,2*} 

The majority of known SVs are **poorly assayed** using currently dominant **short-read sequencing technologies** → evidence supporting an SV event are indirect (e.g. read depth, mismatch read pairs).

SVs can be detected using **long-read sequencing technologies** from Pacific Biosciences (PacBio) and Oxford Nanopore Technologies (ONT) → **direct evidence**.

Long-reads:

- better mappability in repetitive regions, further extending the part of the genome in which variation can be called reliably (+)
- poor quality (-).

CONCLUSIONS

Evaluating nanopore sequencing data processing pipelines for structural variation identification

Anbo Zhou¹, Timothy Lin¹ and Jinchuan Xing^{1,2*} 

Nanopore sequencing is a rapidly developing technology in terms of both sequencing technology and data analysis.

SV callers' performance diverges between SV types → recommendations are tailored to the specific applications:

- For an initial analysis, we recommend **minimap2** + **Sniffles** → high speed and relatively balanced performance calling both insertions and deletions.
- For more detailed analysis, **multiple tools** and integrating their results for the best performance.
- When a high-quality true set can be defined, a **machine learning approach** can be used to further improve the call set.

BAZY DANYCH CNV PRZYKŁADY



BAZA DANYCH: DGV

*D*atabase of *G*enomic *V*ariants
A curated catalogue of human genomic structural variation

About the Project Downloads Statistics FAQ
Genome Browser Query Tool Sessions Contact Us Training Resources

Keyword, Landmark or Region Search: GRCh37/hg19 ▾

Examples: RP11-34P13; CFTR, 7q11.21; chr7:71890181-72690180

Find DGV Variants

[by Study](#) [by Sample](#)
[by Method](#) [by Variant](#)
[by Platform](#) [by Chromosome](#)

Summary Statistics

Stat	Merged-level	Sample-level
CNVs:	983845	7021692
Inversions:	4083	32044

[**Number of Studies:** 75](#)

BAZA DANYCH: DGVA



Services | Research | Training | About us



Overview | Data submission | Data download | Quick tour | Contact

Database of Genomic Variants archive

Phasing out support for the Database of Genomic Variants archive (DGVa).

The submission, archiving, and presentation of structural variation services offered by the DGVa is transitioning to the [European Variation Archive \(EVA\)](#). All of the data shown in the DGVa website is already searchable and browsable from the [EVA Study Browser](#).

Submission of structural variation data to EVA is done using the VCF format. The VCF specification allows representing multiple types of structural variants such as insertions, deletions, duplications and copy-number variants. Other features such as symbolic alleles, breakends, confidence intervals etc., support more complex events, such as translocations at an imprecise position.

We expect to cease accepting direct submissions to DGVa at the end of 2019, in the meantime we recommend submitters make SV submissions to the EVA. If there are specific difficulties with preparing SV submissions in VCF format, please contact the EVA helpdesk.

BAZA DANYCH: DGVA

Overview | Data submission | **Data download** | Quick tour | Contact

DGVa > Data download



Data download

Genomic structural variant study data can be downloaded via ftp by following the appropriate link.

DGVa studies

Study	Reference	Organism	Variants	Link
estd233	luo_et_al_2017b	Homo sapiens	1026	Download via FTP
estd232	Blanco-Kelly_et_al_-2017	Homo sapiens	8	Download via FTP
estd231	WONG_et_al_2016	Homo sapiens	7063	Download via FTP
estd229	Fakhro_et_al_2015	Homo sapiens	16676	Download via FTP
estd228	Ansari_et_al_2016	Homo sapiens	12	Download via FTP
estd226	Zlotina_et_al_2016	Homo sapiens	4	Download via FTP
estd225	Magnusson_et_al_2016	Homo sapiens	1917	Download via FTP
estd224	Suktitipat_et_al_2014	Homo sapiens	3576	Download via FTP
estd223	Boussaha_et_al_2015	Bos taurus	6426	Download via FTP

BAZA DANYCH: EVA

European Variation Archive

Home | Submit Data | Study Browser | Variant Browser | GA4GH | API | RS Release | Help | Feedback |

EVA / STUDY BROWSER

Study Brow

Search for studies archived at EVA using any combination of the filtering options on the left hand-side.

Individual studies can be further investigated using the in-depth study view page found by clicking the study ID in search results.

Filter

Studies found

ID	Name	Genome	Species	Type	Download
estd1	Redon_et_al_2006	Human	Homo sapiens	Control Set	FTP
estd118	Keane_et_al_2011	Mouse	Mus musculus, Mus musculus casta	Control Set	FTP
estd176	Banerjee_et_al_2011	Human	Homo sapiens	Control Set	FTP
estd180	Pang_et_al_2010	Human	Homo sapiens	Control Set	FTP
estd185	Yalcin_et_al_2012	Mouse	Mus musculus	Control Set	FTP
estd186	Thevenon_et_al_2012	Human	Homo sapiens	Case Set	FTP