

## Bioinformatic Note



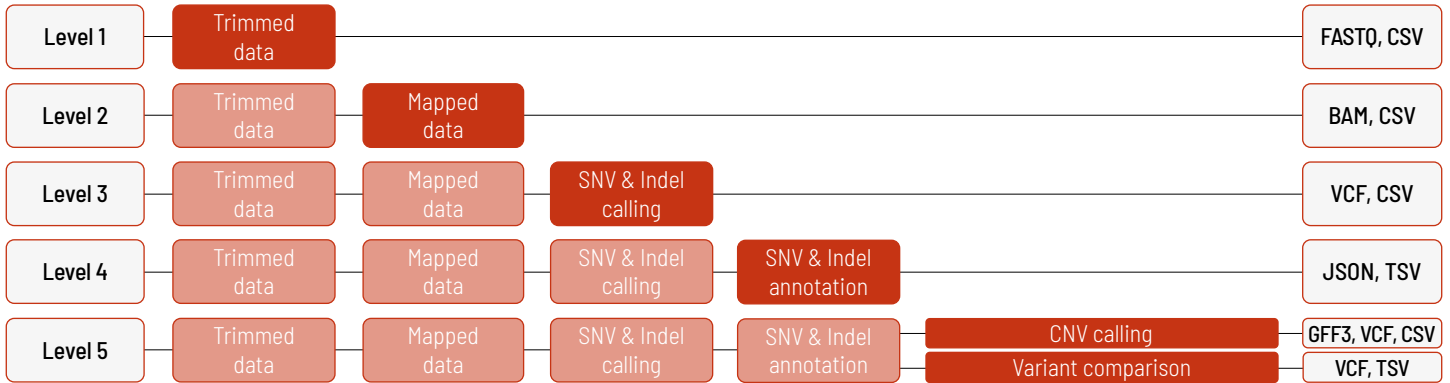
# Exome Sequencing

The exome represents the entirety of all known coding exons of the human genome. Although exons only comprise 1%-2% of the genome, 89% of all known disease-causing mutations are estimated to be located in these regions. Therefore, it is often reasonable to perform a targeted exome analysis.

Whole exome sequencing can be used for various application areas and goals, such as:

- ✕ population genetics
- ✕ genetic disorders
- ✕ rare diseases
- ✕ tumor research

Different levels of bioinformatic data analysis are available:



With increasing bioinformatic level, more data are delivered. All higher levels include the data from the lower levels, e.g., in Level 2, trimmed data and mapped data are provided. In addition to the data, and independent of the analysis level, a project report is generated.

## Level 1

If you wish to analyze your data yourself, we recommend the Levels 1 or 2. The default level for raw data is Level 1, where trimmed reads in FASTQ format are delivered. In this level, the sequencing data are demultiplexed and trimmed. This level is provided for every project, regardless of additionally purchased bioinformatic analyses.

Quality control of the samples is performed, resulting in a metrics file in CSV format. This file contains different sections for each metric type, such as read mean quality, positional base mean quality, positional base content, read GC content and quality, sequence positions, or positional quality. Each of these sections includes separate rows for length, position, or other relevant category variables.

The generated project report provides information for every sample about the laboratory protocol, including data about quality control of the starting material, library preparation, sequencing parameters, and the Q30 value of the sequencing. For the trimmed data, the number of sequenced fragments and bases is reported, and the sequence length, quality of the reads, and the GC content are illustrated in bar plots for all samples.

Additionally, a multiQC report in HTML format is generated. It covers the results from the DRAGEN FastQC module and – if performed – further DRAGEN analyses. In contrast to the generated project report, this multiQC report facilitates the interactive exploration of the analysis results: Samples can be highlighted, renamed, or hidden. Additionally, figures can be customized, edited, and saved.

## Level 2

If you wish to receive Level 2, the trimmed reads are aligned to the reference genome, and duplicates are marked. In addition to the trimmed reads in FASTQ format, you will receive the mapped reads as BAM files.

Together with the mapped reads, you will receive a mapping metrics file in CSV format that includes mapping and aligning metrics. The metrics are available over all input data as well as on a per-read-group level. Examples of the mapping and aligning metrics are the number of total input reads, the number of duplicate marked reads, the number of unique reads, reads with mate sequenced, QC-failed reads, and mapped reads. Table 1 shows an excerpt of the mapping metrics file.

Table 1 | Excerpt of the mapping metrics file in CSV format (header line added for the sake of clarity).

All input data/per read group level	Metrics	Value	Percentage
MAPPING/ ALIGNING SUMMARY	Total input reads	120000000	100
MAPPING/ ALIGNING SUMMARY	Number of duplicate marked reads	7637844	6,36
MAPPING/ ALIGNING SUMMARY	Number of duplicate marked and mate reads removed	NA	
MAPPING/ ALIGNING SUMMARY	Number of unique reads (excl. duplicate marked reads)	112362156	93,64

Another file in CSV format reports the coverage metrics. It provides metrics over the target region, such as the aligned bases in the region, average alignment coverage over the region, or the uniformity of the coverage. Similar to the mapping metrics file, the coverage metrics file provides the metrics and respective values per row.

GC biases can arise from library prep, capture kits, sequencing system differences, and mapping. We can perform a GC bias correction. In a Gnu Zipped file (GZ), the GC-corrected target counts are stored. The file contains columns with the contig identifier, start position, end position, target interval name, (GC-corrected) count of alignments in this interval, and (GC-corrected) count of improperly paired alignments in this interval. An excerpt of this file can be found in table 2.

A ploidy estimator calculates the sequencing depth of the coverage for each autosome and allosome to subsequently estimate the sex karyotype of the sample. The result is provided in a ploidy estimation metrics file in CSV format. In this file, the median coverage of the autosomal chromosomes, the X chromosome, and the Y chromosome is provided, as well as the ratio of each chromosome and the autosomal median. Based on these ratios, a ploidy estimation is provided eventually.

For Level 2, the project report also includes a table with statistics of the mapped reads, including the number of mapped reads, the proportion of sequenced reads, the proportion of PCR duplicates, the median insert size, and the average coverage.

### Level 3

In Level 3, single nucleotide variants (SNVs) and small insertions and deletions (indels) are called. Calling of the SNVs and indels is performed with default germline parameters. The resulting VCF file contains the quality-filtered calls of the SNVs and small indels. It includes information for every small variant about its location (chromosome and position), the identifier, the reference base(s) and the alternate base(s), the quality and filter status, and additional information in several columns. In the header of the VCF file, additional information about the FILTER, INFO, and FORMAT columns is provided, and the abbreviations are explained. An excerpt of the small variant VCF file is shown in table 3. For additional information on the VCF format, we refer the reader to the [VCF specification](#).

An SNV and indels metrics file in CSV format includes information about the variant calling statistics. Amongst others, this metrics file contains information about the number of processed reads, the number of insertions, deletions, and SNVs before and after the quality filtering. Every entry consists of a section description, the sample, the metric, its value as count or ratio, and, where applicable, the percentage.

Table 2 | Excerpt of the GC-corrected target counts file.

Contig	Start	Stop	Name	Sample 1	Improper_pairs
chr1	65564	65569	target-wes-chr1-65564:65569/1	212.3316056	0
chr1	65569	65573	target-wes-chr1-65569:65573/2	131.4680721	0
chr1	69036	69522	target-wes-chr1-69036:69522/1	1338.344186	6

Table 3 | Excerpt of the SNVs and indels VCF file. The abbreviations in the INFO and FORMAT columns are explained in the respective VCF file.

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	SAMPLE1
chr1	601436	.	C	T	5.16	PASS	AC=1;AF=0.500;AN=2;D-P=8;FS=0.000;MQ=24.84;MQRankSum=0.421;QD=3.38;ReadPosRankSum=0.922;SOR=0.693;FractionInformativeReads=1.000	GT:AD:AF:DP:-F1R2:F2R1:GQ:-PL:GP:PRI:SB:MB	0/1:4,4:0.5000:8:1,3:3,1:5:38,0,18:5.1642e+00,1.6103e+00,2.2745e+01:0.00,34.77,37.77:0,4,0,4:3,1,1,3
chr1	611317	.	A	G	15.80	PASS	AC=1;AF=0.500;AN=2;D-P=72;FS=3.135;MQ=28.25;MQRankSum=0.489;QD=0.42;ReadPosRankSum=0.535;SOR=0.274;FractionInformativeReads=1.000	GT:AD:AF:DP:-F1R2:F2R1:GQ:-PL:GP:PRI:SB:MB	0/1:45,27:0.3750:72:23,12:22,15:14:50,0,17:1.5800e+01,1.6288e-01,1.9784e+01:0.00,34.77,37.77:5,4,0,5,22:25,20,16,11

## Level 4

In Level 4, the SNVs and indels are annotated. The annotations are available as compressed JSON files and as TSV files. The JSON files are usually very large and not optimized for human readability. However, it is useful for automated processing steps. We provide an additional annotation file in a tabular format that contains selected information from the JSON annotation file in a tabular format.

The annotations in tabular format include, amongst others,

- ✗ information about the chromosomal position and the observed variant,
- ✗ functional consequences of the variant in the context of a transcript,
- ✗ position and sequence changes in the context of the most affected transcripts,
- ✗ and information about the observed variant in the global population.

Together with the annotation list, you will receive a file that contains further information on the annotation of variant lists. Therefore, we will not go into detail about the provided TSV files. A database-versions file in TXT format provides information about the names, versions, and short descriptions of the used databases.

SNVs are commonly used markers in case-control association studies. Furthermore, some SNVs can have functional impact leading to disease susceptibilities and drug sensitivities. These functional impacts can concern the transcriptional machinery of a cell, alternative or aberrant splice isoforms when located at a splice site, or the translational machinery leading to protein folding, localization, stability, binding, or catalysis interference (Cline and Karchin, 2011). As SNVs, indels have an impact on certain diseases. With the annotated SNVs and indels files, many research questions regarding functional impacts on diseases, disease susceptibility, and drug sensitivity might be answered.

## Level 5

In Level 5, copy number variations (CNVs) can be called. The results are stored in three files: A VCF file, a GFF3 file, and a metrics file in CSV format. The VCF file looks like the small variant VCF file displayed in table 3. Due to its similarity and the additional explanations in the header of the VCF file, we do not show an excerpt of the file here again.

The information about the CNVs is also stored in a GFF3 file. Like the VCF file, the GFF3 file includes information about the chromosome, the position of the feature, and additional information. An excerpt of the GFF3 file is shown in table 4. The header line is added for explanatory reasons and is not provided in the delivered file.

Like the other metrics files, the metrics file for CNVs is in CSV format and contains information about the CNV calling statistics, including, e.g., the number of called amplifications, deletions, or segments. Every entry consists of a section description, the metric, its value as count or ratio, and, where applicable, the percentage.

In addition to the three CNV files, a PDF file includes the coverage deviation and variant frequency deviation for every chromosome as figures.

Similar to the SNVs and indels, the CNVs are annotated. The annotation results are available as JSON and TSV files. As these two files were already described for the SNV and indel annotation, we will not go into detail here again. However, the annotation files for CNVs do not include information about the observed variant in the global population.

## Resources

For further information, the DRAGEN Bio-IT Platform manual can be browsed [here](#).

Table 4 | Excerpt of the CNVs' GFF3 file (header line added for the sake of clarity).

Sequence ID	Source	Feature type	Feature start	Feature end	Score	Strand	Phase	Attributes
chr1	DRAGEN	CNV	183922	184158	9	.	.	Alt=DEL;LinearCopyRatio=0.652688;CopyNumber=1;Genotype=0/1;Qual=9;Filter=cnvQual;Start=183921;Stop=184158;Length=237;BinCount=2;ImproperPairsCount=4,5;color=#DDDDDD;
chr1	DRAGEN	CNV	685716	686654	9	.	.	Alt=DUP;LinearCopyRatio=1.29752;CopyNumber=3;Genotype=.;Qual=9;Filter=cnvQual;Start=685715;Stop=686654;Length=939;BinCount=2;ImproperPairsCount=5,16;color=#DDDDDD;
chr1	DRAGEN	CNV	13049808	13226106	32	.	.	Alt=DUP;LinearCopyRatio=1.30753;CopyNumber=3;Genotype=.;Qual=32;Filter=PASS;Start=13049807;Stop=13226106;Length=176299;BinCount=44;ImproperPairsCount=14,20;color=#FF0000;

## References

Cline, Melissa S; Karchin, Rachel (2011): Using bioinformatics to predict the functional impact of SNVs. In *Bioinformatics* 27(4), pp. 441 – 448.



## About Us

CeGaT was founded in 2009 in Tübingen, Germany. Our scientists are specialized in next-generation sequencing (NGS) for genetic diagnostics, and we also provide a variety of sequencing services for research purposes and pharma solutions. Our sequencing service portfolio is complemented by analyses suited for microbiome, immunology, and translational oncology studies.

Our dedicated project management team of scientists and bio-informaticians works closely with you to develop the best strategy to realize your project. Depending on its scope, we select the most suitable library preparation and conditions on our sequencing platforms.

We would be pleased to provide you with our excellent service.  
Contact us today to start planning your next project.



CeGaT GmbH  
Research & Pharma Solutions  
Paul-Ehrlich-Str. 23  
72076 Tübingen  
Germany



CLIA CERTIFIED ID: 99D2130225

Phone: +49 7071 56544-333  
Fax: +49 7071 56544-56  
Email: rps@cegat.com



Accredited by DAkks according to  
DIN EN ISO/IEC 17025:2018



For more details please visit  
[www.cegat.com/rps](http://www.cegat.com/rps)