

## Article

# Copy number variant detection with low-coverage whole-genome sequencing is a viable replacement for the traditional array-CGH

Marcel Kucharík<sup>1,2,\*</sup>, Jaroslav Budiš<sup>1,2,3</sup>, Michaela Hýblová<sup>4</sup>, Gabriel Minárik<sup>4</sup>, and Tomáš Szemes<sup>1,2,5</sup>

<sup>1</sup> Geneton Ltd., Bratislava, Slovakia; [geneton@geneton.sk](mailto:geneton@geneton.sk)

<sup>2</sup> Comenius University Science Park, Bratislava, Slovakia; [info@cusp.uniba.sk](mailto:info@cusp.uniba.sk)

<sup>3</sup> Slovak Centre of Scientific and Technical Information, Bratislava, Slovakia; [cvti@cvtisr.sk](mailto:cvti@cvtisr.sk)

<sup>4</sup> Trisomy Test Ltd., Bratislava, Slovakia; [info@trisomytest.sk](mailto:info@trisomytest.sk)

<sup>5</sup> Department of Molecular Biology, Faculty of Natural Sciences, Comenius University, Bratislava, Slovakia; [kmb@fns.uniba.sk](mailto:kmb@fns.uniba.sk)

\* Correspondence: [marcel.kucharik@geneton.sk](mailto:marcel.kucharik@geneton.sk); Tel.: +421-904-855-365

Received: date; Accepted: date; Published: date

**Abstract:** Copy number variations (CNVs) are a type of structural variant involving alterations in the number of copies of specific regions of DNA, which can either be deleted or duplicated. CNVs contribute substantially to normal population variability; however, abnormal CNVs cause numerous genetic disorders. Nowadays, several methods for CNV detection are used, from the conventional cytogenetic analysis through microarray-based methods (aCGH) to next-generation sequencing (NGS). We present GenomeScreen - NGS based CNV detection method based on a previously described CNV detection algorithm used for non-invasive prenatal testing (NIPT). We determined theoretical limits of its accuracy and confirmed it with extensive *in-silico* study and already genotyped samples. Theoretically, at least 6M uniquely mapped reads are required to detect CNV with a length of 100 kilobases (kb) or more with high confidence (Z-score > 7). In practice, the *in-silico* analysis showed the requirement at least 8M to obtain >99% accuracy (for 100 kb deviations). We compared GenomeScreen with one of the currently used aCGH methods in diagnostic laboratories, which has a 200 kb mean resolution. GenomeScreen and aCGH both detected 59 deviations, GenomeScreen furthermore detected 134 other (usually) smaller variations. Furthermore, the overall cost per sample is about 2-3x lower in the case of GenomeScreen.

**Keywords:** CNV detection, low-coverage WGS, CNV detection comparison, aCGH replacement

## 1. Introduction

Copy number variations (CNVs) are a phenomenon in which sections of the genome are repeated, and the number of repeats in the genome varies between individuals. CNVs contribute substantially to normal population variability. However, abnormal CNVs are a cause of numerous genetic disorders. Several methods for CNV analysis are used,

from the conventional cytogenetic analysis through microarray-based approaches to next-generation sequencing (NGS) [1].

Array-based comparative genomic hybridization (aCGH) provides genome-wide coverage at a great resolution, even in the scale of tens of kilobases (10–25 kb) [2]. This fact promoted aCGH for a golden standard in CNVs detection for several years [3] despite some limitations in resolution and accuracy [4].

In contrast, NGS provides a sensitive and accurate approach for the detection of the major types of genomic variations, including CNVs [5,6]. There are three basic strategies for NGS-based CNV analysis, including whole-genome, whole-exome, and targeted sequencing. Whole-genome sequencing allows even base-pair resolution of breakpoints with a sophisticated bioinformatics pipeline [7]. On the other hand, whole-exome sequencing and targeted sequencing aim to reduce the sequencing cost. Still, they are limited to certain regions (protein-coding, or custom), where most known disease-causing mutations occur [8].

We present GenomeScreen - low-coverage whole-genome NGS based CNV detection method (based on the previously published NIPT CNV detection method [9,10]) and estimate its accuracy in theoretical and *in-silico* settings. Furthermore, we compare its sensitivity to the more conventional aCGH method on 106 laboratory prepared clinical samples.

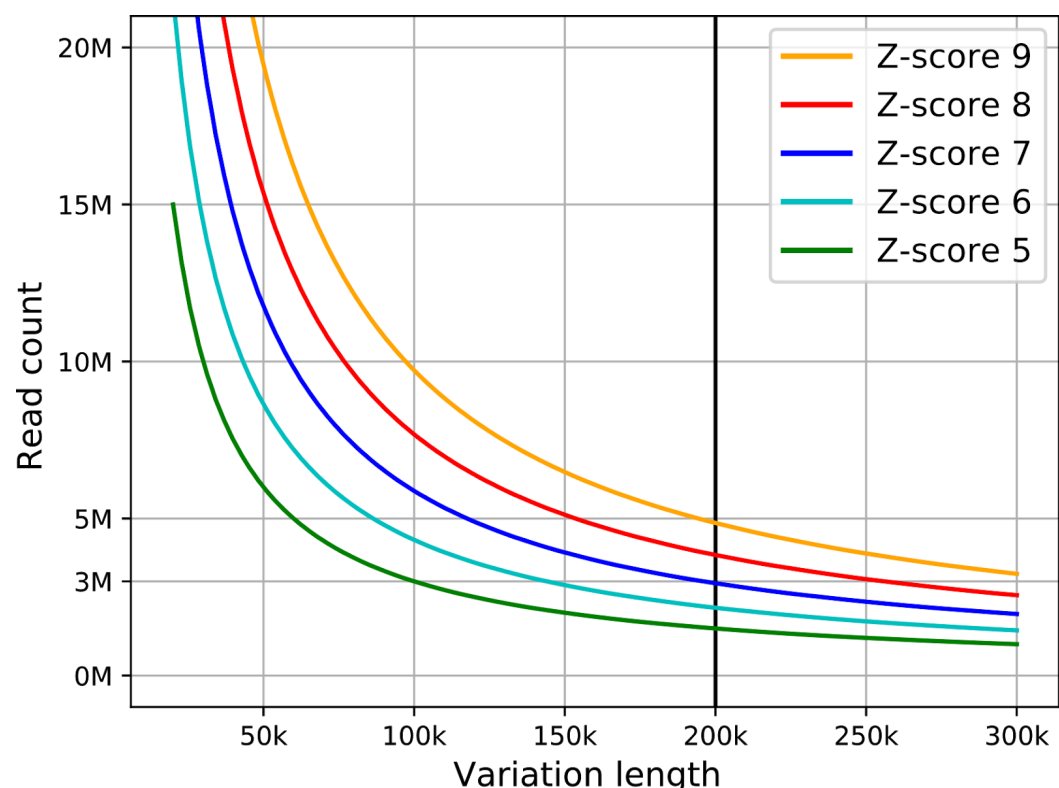
## 2. Results

### 2.1 Theoretical minimal read count

The theoretical minimum of reads for predicting a variation with length  $l_c$  with the desired Z-score ( $Z$ ) is estimated as (see Section 4.2):

$$n \geq \frac{4Z^2(l_g - l_c)}{l_c}$$

Standardly, Z-score of 4 is used in the detection of whole chromosomal aneuploidies [11,12]; however, there are inherently more possible CNVs as whole chromosomal aneuploidies. Thus, the desired Z-score should be much higher in this instance to decrease the number of false positives. Moreover, in practice, the number of reads needed would be even larger due to the uncertainty of sequencing, mapping, and inherent biological biases [13,14]. The theoretical minimal read count estimation for different Z-scores can be seen in Figure 1.

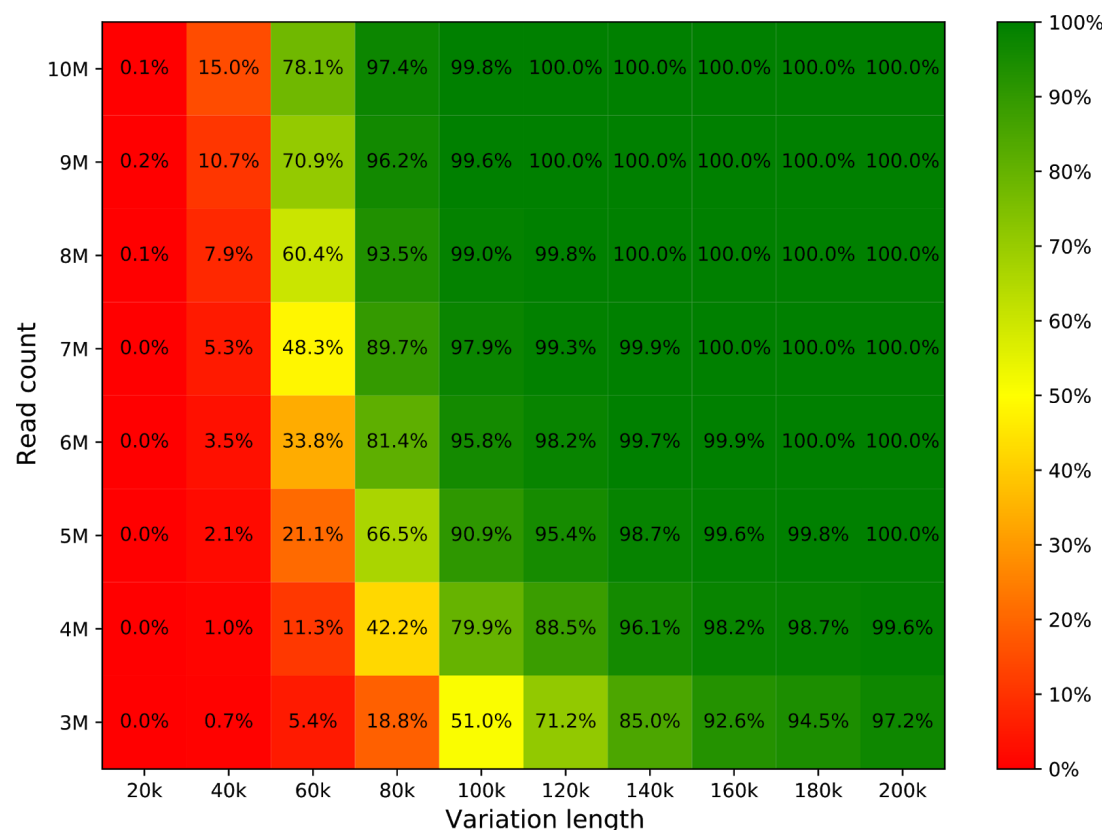


**Figure 1.** Theoretical minimal read count for successful estimation of CNV with specified variation length. Different lines represent different Z-score confidence levels.

## 2.2 Detection accuracy for variable CNV lengths and read count (*in-silico*)

To verify the theoretically estimated limitations, we first conducted a simulated *in-silico* experiment. Artificial samples with simulated CNV were created from healthy samples by multiplication of bins corresponding to the simulated region randomly on the genome. Only regions that did not span into filtered positions were kept for further analysis (about 85% of the genome). The details can be found in Section 4.3.

The *in-silico* analysis shows the influence of read count and CNV length for prediction accuracy (Figure 2). Based on the findings, we recommend using a read count of at least 8M to achieve >99% prediction accuracy for variations with 100 kb and more. Thus we recommend following the line for Z-score of 8 (red on Figure 1.) for estimation for different CNV lengths.

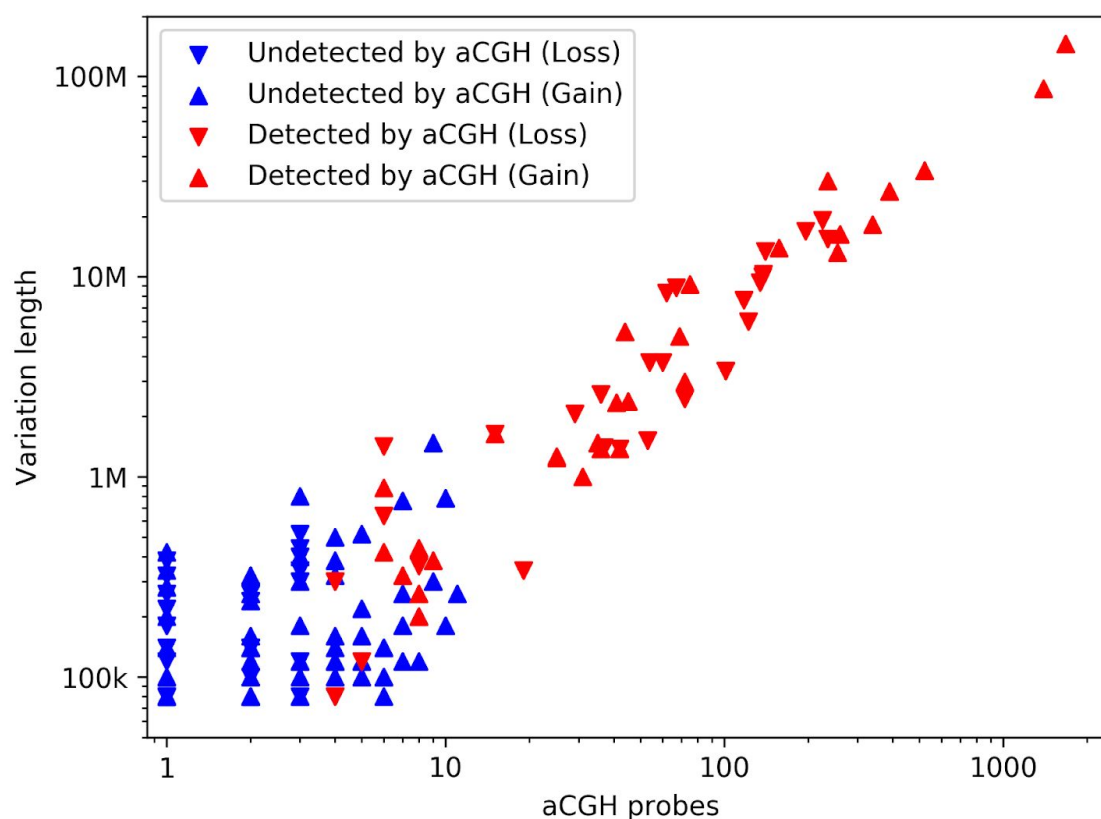


**Figure 2.** Prediction accuracy computed with *in-silico* analysis based on the length of variation and read count. Each cell number is generated from 8,300 simulations (100 randomly generated aberrations, 83 samples).

### 2.3 Evaluation of clinical samples

Finally, we ran an evaluation of samples analyzed previously in diagnostic settings using aCGH method (Human Genome CGH Microarray 4x44K Agilent [15]) and GenomeScreen. The chosen aCGH method has 42,494 probes, which result in mean accuracy of detection approximately 200 kb; however, the probes are focused mainly in gene regions and very sparsely in intergenomic regions, therefore the accuracy will be better within and worse outside of genes.

From the 106 tested samples, 58 did not show any detection on aCGH, and the rest contained 59 detections together (lengths from 39 kb to 146 Mb), from which GenomeScreen also detected all. The detections on GenomeScreen and on aCGH have excellent concordance - median 94.37% overlap (more data in Supplementary material Table 1.). GenomeScreen furthermore detected 134 additional variations with ranges from 80 kb up to 1.48 Mb, mainly in regions with a low number of aCGH probes and protein genes, where aCGH has low coverage (Figure 3. and Supplementary material Table 1. and Figure 1.).



**Figure 3.** Detection of GenomeScreen (all) and aCGH (red) based on the variation length and number of aCGH probes in the detected interval (by GenomeScreen). Deletions and duplications are visualized by lower and upper triangles, respectively.

### 3. Discussion

According to the presented results, our method is currently able to detect all tested and almost all simulated variations longer than 100 kb in mappable regions on the human genome; thus, it can replace the more conventional aCGH method. The detection accuracy of GenomeScreen heavily depends on variation length, read count, and genomic position of the variation. The variation length and position cannot be influenced but can be partly compensated with an increase of the read count, which comes inevitably with a higher operational cost.

The disadvantage of the binning approach used here is that the variation location is always reported as a multiplier of the bin-size (20 kb). On the other hand, the aCGH method uses probes, which can be seen as variable size bins, where the resolution is equal to the probe distance (which is sometimes larger than the 20 kb bin-size). Moreover, the bin-size of GenomeScreen can be lowered to increase the precision of the variation location at the cost of deeper sequencing.

The false-positive rate of GenomeScreen is not studied in this work and should be adequately addressed in the future. However, the loss or gain of the (non-mosaic) deviation with a length of at least 100 kb is so substantial, that we do not expect to see any false positive detections.

In our setting, the GenomeScreen method is 2-3x cheaper than aCGH, while the detection accuracy on variations larger than 100 kb is comparable or even better for GenomeScreen [1]. The aCGH method has probes more frequently inside genes, thus the accuracy for the aCGH approach would be higher there (lower outside of genes). In practice, we observed more variations detected by the GenomeScreen method, mainly in the intergenic regions. However, further research is needed to correctly assess the prediction accuracy for both methods for variation lengths between 50 kb and 100 kb and different genomic regions.

## 4. Materials and Methods

### 4.1 Sample collection and processing

Samples of chorionic villi, amniotic fluid, placenta, tissue, or peripheral blood were obtained from 106 patients in the clinical sample group and 789 in the training group. All patients signed written consent for participation in the research. Peripheral blood was drawn in EDTA or STRECK tubes, inverted several times after collection, stored in a chilled environment (4–10 °C) for EDTA and at room temperature for STRECK tubes, and transported to the laboratory within 36 hours. DNA was extracted from 200 µl of whole blood or 700 µl of amniotic fluid using the QIAamp DNA Blood Mini Kit (Qiagen, Hilden, Germany) according to the manufacturer's protocol and stored at –20°C until further analysis.

Genomic DNA from clinical samples was fragmented using 1U/µl dsDNA Shearase™ Plus (Zymo Research, Irvine, CA, USA) and incubated 23 min at 42°C to generate 100-500bp fragments. For adapter-ligated DNA library construction, a TruSeq Nano kit (Illumina, San Diego, CA, USA) with an in-house optimized protocol was used. Low coverage sequencing (0.3×) was performed on Illumina NextSeq 500/550 platform (Illumina, San Diego, CA, USA) with paired-end setting 2×35 using High Output Sequencing Kit v2.5. Library quantity and quality were measured by fluorometric assay on Qubit 2.0 (ds DNA HS Assay Kit, Life Technologies, Eugene, Oregon, USA). Fragment analysis was performed on 2100 Bioanalyzer (High Sensitivity DNA Kit, Agilent Technologies, Waldbronn, Germany). We targeted 5M uniquely mapped reads per sample; however, none of the analyses were excluded due to lower (or higher) read counts (more info in Supplementary material Table 1.).

### 4.2 Theoretical minimal read count estimation

Suppose we model sequencing as a random choice of reads from the whole (mappable) genome. Then we can theoretically deduce the number of needed uniquely mapped reads for a certain accuracy criterion. The random choice for a target region is described by the binomial distribution with mean  $\mu = np$  and variance  $\sigma^2 = np(1 - p)$ . Here,  $p$  is the probability of choosing a read from the target region, and  $n$  is the number of reads sequenced. The probability  $p$  can be furthermore expressed as the ratio of the region length  $l_c$  to whole-genome length  $l_g$  ( $p = l_c/l_g$ ). When predicting a CNV, we need to have certain confidence traditionally determined by the Z-score ( $Z$ ), defined as:

$$Z = \frac{\delta - \mu}{\sigma}$$

Here  $\delta$  is the number of reads that we observe in the target region. We assume that the number of reads in the target region will be proportional to the number of present copies of gonosomes, i.e., either  $\delta = n(p + p/2)$  for duplication or  $\delta = n(p - p/2)$  for deletion of the region on a single chromosome. If we solve for  $Z^2$  and substitute:

$$Z^2 = \frac{(\delta - \mu)^2}{\sigma^2} = \frac{(n(p+p/2) - np)^2}{np(1-p)} = \frac{n^2 p^2}{4np(1-p)} = \frac{np}{4(1-p)} = \frac{nl_c}{4(l_g - l_c)}$$

Then we can estimate the minimal number of reads ( $n$ ) to be able to predict a variation with length  $l_c$  with the desired Z-score ( $Z$ ):

$$n \geq \frac{4Z^2(l_g - l_c)}{l_c}$$

### 4.3 Variant identification

To identify variations, we performed the following pipeline:

1. Mapping and binning
  - a. mapping reads using bowtie2 [16]
  - b. binning reads into same-size 20 kb bins
  - c. normalizing bin counts
2. Normalization (similar to one published previously by [17])
  - a. LOESS-based GC correction [18]
  - b. PCA normalization to remove higher-order population artifacts on autosomal chromosomes
  - c. subtracting per-bin mean bin count to obtain data normalized around zero
3. Filtration of unusable bins
  - a. unmappable or badly mappable regions (zero or low mean of bin count)
  - b. repetitive regions or areas with some systematically increased mappability (high mean of bin count)
  - c. highly variable regions (high variance of bin count)
4. Segment identification and reporting
  - a. circular binary segmentation algorithm [19] to identify consistent segments of similar coverage
  - b. assigning significance to segments based on the proportion of reads
  - c. visualization of findings (Figure 3.)

Scripts (Python 3.7) and data are available on the website <https://github.com/marcelTBI/GenomeScreen>.

#### 4.3.1 Mapping and binning

Firstly, the reads are mapped to a reference using Bowtie 2 [16] with *--very-sensitive* setting. We use hg19 reference in all applications, but other references can be used without changes to the algorithm. The reads are then filtered for map quality at least 40 and binned according to their starts to same-size 20 kb bins. All subsequent analyses are performed on the bin counts, and the algorithm does not use any other information about reads (for example, sequence). For training purposes, the bin counts corresponding to autosomal chromosomes for each sample are normalized to the same number of reads (i.e., each bin is divided so the sum of all bins on autosomal chromosomes would be the same for each sample). Furthermore, the same is done separately for chromosome X and chromosome

Y. A consequence of separate normalization of sex chromosomes is that the current approach can detect only small sex chromosomal variations and not the whole sex chromosomal aneuploidies.

### 4.3.2 Normalization

Normalization consists of three steps: firstly, a sample-wise LOESS-based GC correction is employed on the bin counts [18]. Afterward, the principal component analysis (PCA) normalization is used to remove higher-order population artifacts on autosomal chromosomes [17]. For training of the PCA, LOESS-corrected bin counts of 789 NIPT samples with female fetuses were converted to principal component space, and the first 15 principal components were stored. The bin count vector of a new sample is then transformed into principal component space defined by these first 15 components and transformed back to the bin space to obtain residuals that are then removed from the bin counts. The first principal components represent noise commonly seen in euploid samples, and their removal helps to normalize the data. Currently, the PCA normalization is done only on autosomal chromosomes due to the unavailability of enough male samples for training. In the future, the training of PCA on both male and female samples is likely to increase the prediction precision for sex chromosomes. Lastly, we subtract per-bin mean bin counts to obtain data normalized around zero. This last step is trained already on the PCA normalized bin counts (where available) and helps compensate for the mapping inequality between various genomic regions.

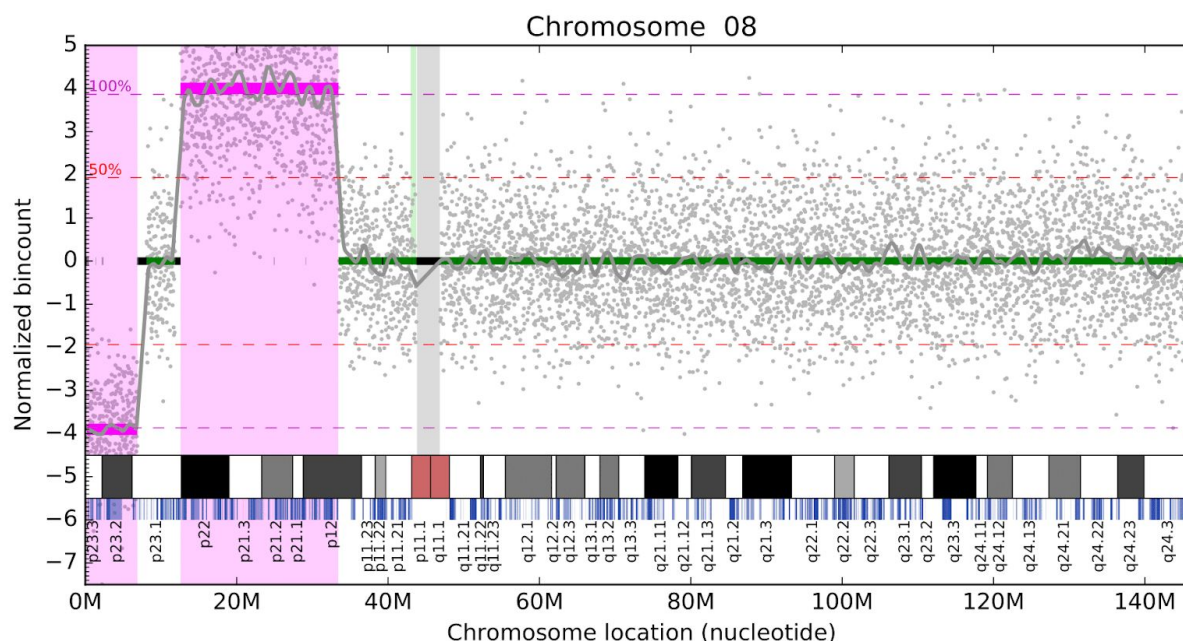
### 4.3.3 Filtration of unusable bins

To further improve accuracy, we filter bins that have an unusual signature - either low mean (this signals bad mappability of the region), high mean (repetitive regions or regions with some systematic bias), or high variance (highly variable regions). Furthermore, the filtered regions were manually curated to reduce the scatter of filtered regions, mainly around centromeres and sex chromosomes. The filtration leaves out around 15% of the genome, mainly due to the low mappability, especially in and around centromeres.

### 4.3.4 Segment identification and reporting

After normalization and filtering, we have a signal (grey dots in Figure 4.) that needs to be segmented into the same level parts to be evaluated. For this purpose, we use the circular binary segmentation (CBS) algorithm implemented in the R package DNACopy [19]. After segmentation, each segment is assigned a significance level based on its length and difference from zero. Since we know the mean bin counts, we can estimate the level for a complete deletion or duplication of one copy of a chromosome (magenta dashed lines in Figure 4.). We then differ between five color-coded levels of significance: magenta - at least 75%, at least 200kb, red - at least 25%, at least 200kb, orange - at least 25%, at least 40kb, yellow - at least 12,5%, at least 40kb, and green - all others (very short segments or segments around zero). The findings are then reported as a text file for further machine processing, and each chromosome is visualized (Figure 4.).





**Figure 4.** Visualization of detected deviations on chromosome 8. Chromosome location is on X-axis. Normalized bin count is on Y-axis. Green lines represent normal bin count segments (normalized around zero), magenta lines visualize aberrations (one deletion in the start of the chromosome, one duplication on p22-p12). Filtered bins are depicted as black bars on the zero line on Y-axis. The unmapped region around centromere is visualized with a grey bar. Grey dots represent the normalized individual bin counts for each bin.

#### 4.4 *In-silico* analysis

For *in-silico* analysis, we chose 83 samples without any aberration and with read count at least 10M. Firstly, the samples were downsampled to the studied read count (3M - 10M with the step of 1M). Then, for each of the tested variation lengths (20 kb - 200 kb with the step of 20 kb), 100 random variations on autosomal chromosomes were generated that do not overlap with the filtered regions (see Section 4.3.3). To create a sample with an artificial aberration, the bins corresponding to the generated random variation were multiplied accordingly (thus, the most time-consuming mapping step was performed only once per sample). Afterward, variant identification was performed without changes.

In total, we gradually created 664,000 artificial samples (100 variations \* 83 samples \* 10 variation lengths \* 8 read counts) and performed variant identification on them to analyze the impact of read count and variant length. The results are displayed in Figure 2.

**Supplementary Materials:** Supplementary materials can be found at a Google Drive folder [https://drive.google.com/drive/folders/1VQU\\_ljq5m0CdNzantB-tgG-3ImS79Id9](https://drive.google.com/drive/folders/1VQU_ljq5m0CdNzantB-tgG-3ImS79Id9).

**Author Contributions:** Conceptualization, T.S., G.M., and J.B.; methodology, M.K., M.H., and G.M.; software, M.K.; validation, M.K., M.H., and G.M.; investigation, M.K.; resources, M.H.; data curation, M.K., M.H.; writing—original draft preparation, M.K.; writing—review and editing, M.K., M.H., G.M., and J.B.; visualization, M.K.; supervision, J.B.; funding acquisition, G.M., and T.S. All authors have read and agreed to the published version of the manuscript.

**Funding:** The presented work was supported by projects “REVOGENE - Research centre for molecular genetics” (ITMS 26240220067) and by the “Long term strategic research and development focused on the occurrence of Lynch syndrome in the Slovak population and possibilities of prevention of tumors associated with this syndrome” (ITMS 313011V578) co-financed by the European Regional Development Fund.

**Conflicts of Interest:** We declare a potential competing financial interest in the form of employee contracts (see affiliations for each author) with Geneton Ltd. and TrisomyTest Ltd. Geneton Ltd. participated in the development of a commercial NIPT test in Slovakia; however, it is not a provider of this commercial test, but continues to do basic and applied research in the field of NIPT. On the other hand, TrisomyTest Ltd. is the commercial provider of NIPT testing in Slovakia. Its participation in the study was limited to the routine NIPT testing that generated the genomic results reused in our research. Related to this work, there are no patents, products in development, or marketed products to declare. The authors declare no other conflict of interest.

## Abbreviations

aCGH - array-based comparative genomic hybridization

CBS - circular binary segmentation

CNV - copy number variant

NGS - next-generation sequencing

NIPT - non-invasive prenatal testing

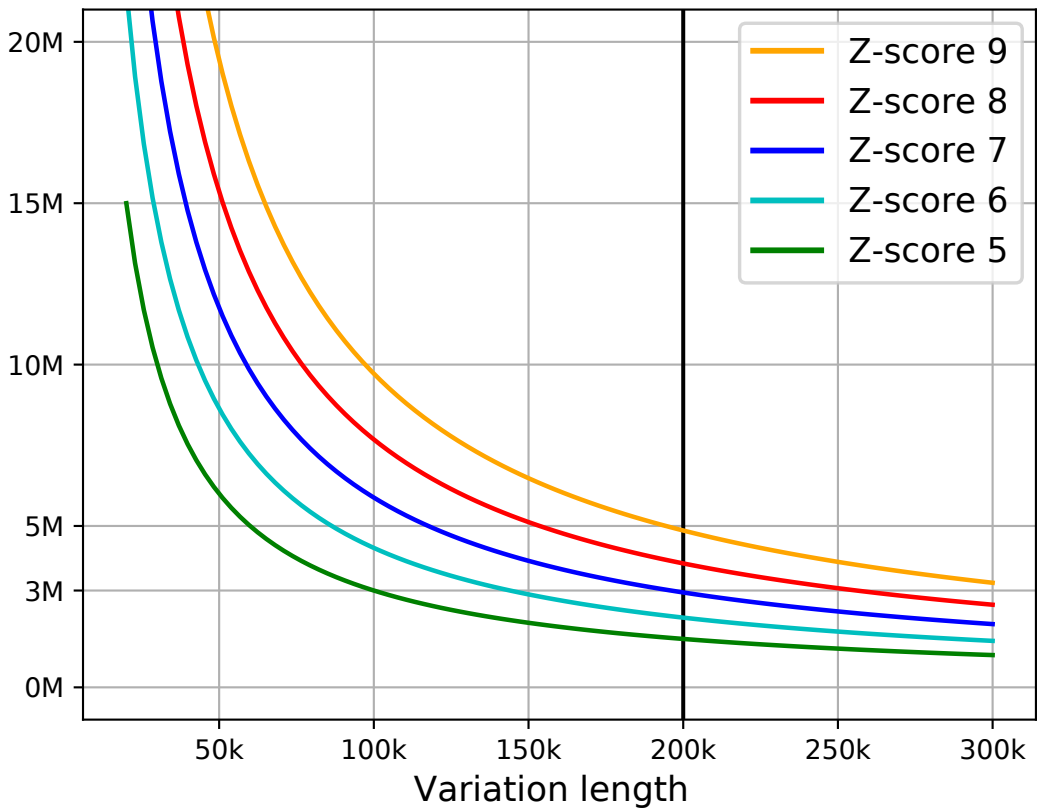
## References

1. Pös O, Budis J, Kubiritova Z, Kucharik M, Duris F, Radvanszky J, et al. Identification of Structural Variation from NGS-Based Non-Invasive Prenatal Testing. *Int J Mol Sci.* 2019;20. doi:10.3390/ijms20184403
2. Yoon S, Xuan Z, Makarov V, Ye K, Sebat J. Sensitive and accurate detection of copy number variants using read depth of coverage. *Genome Res.* 2009;19: 1586–1592.
3. Kearney HM, Thorland EC, Brown KK, Quintero-Rivera F, South ST, Working Group of the American College of Medical Genetics Laboratory Quality Assurance Committee. American College of Medical Genetics standards and guidelines for interpretation and reporting of postnatal constitutional copy number variants. *Genet Med.* 2011;13: 680–685.
4. Coughlin CR 2nd, Scharer GH, Shaikh TH. Clinical impact of copy number variation analysis using high-resolution microarray technologies: advantages, limitations and concerns. *Genome Med.* 2012;4: 80.
5. Russo CD, Di Giacomo G, Cignini P, Padula F, Mangiafico L, Mesoraca A, et al. Comparative study of aCGH and Next Generation Sequencing (NGS) for chromosomal microdeletion and microduplication screening. *J Prenat Med.* 2014;8: 57–69.
6. Wang H, Nettleton D, Ying K. Copy number variation detection using next generation sequencing read counts. *BMC Bioinformatics.* 2014;15: 109.
7. Magi A, Bolognini D, Bartalucci N, Mingrino A, Semeraro R, Giovannini L, et al. Nano-GLADIATOR: real-time detection of copy number alterations from nanopore sequencing data. *Bioinformatics.* 2019;35:

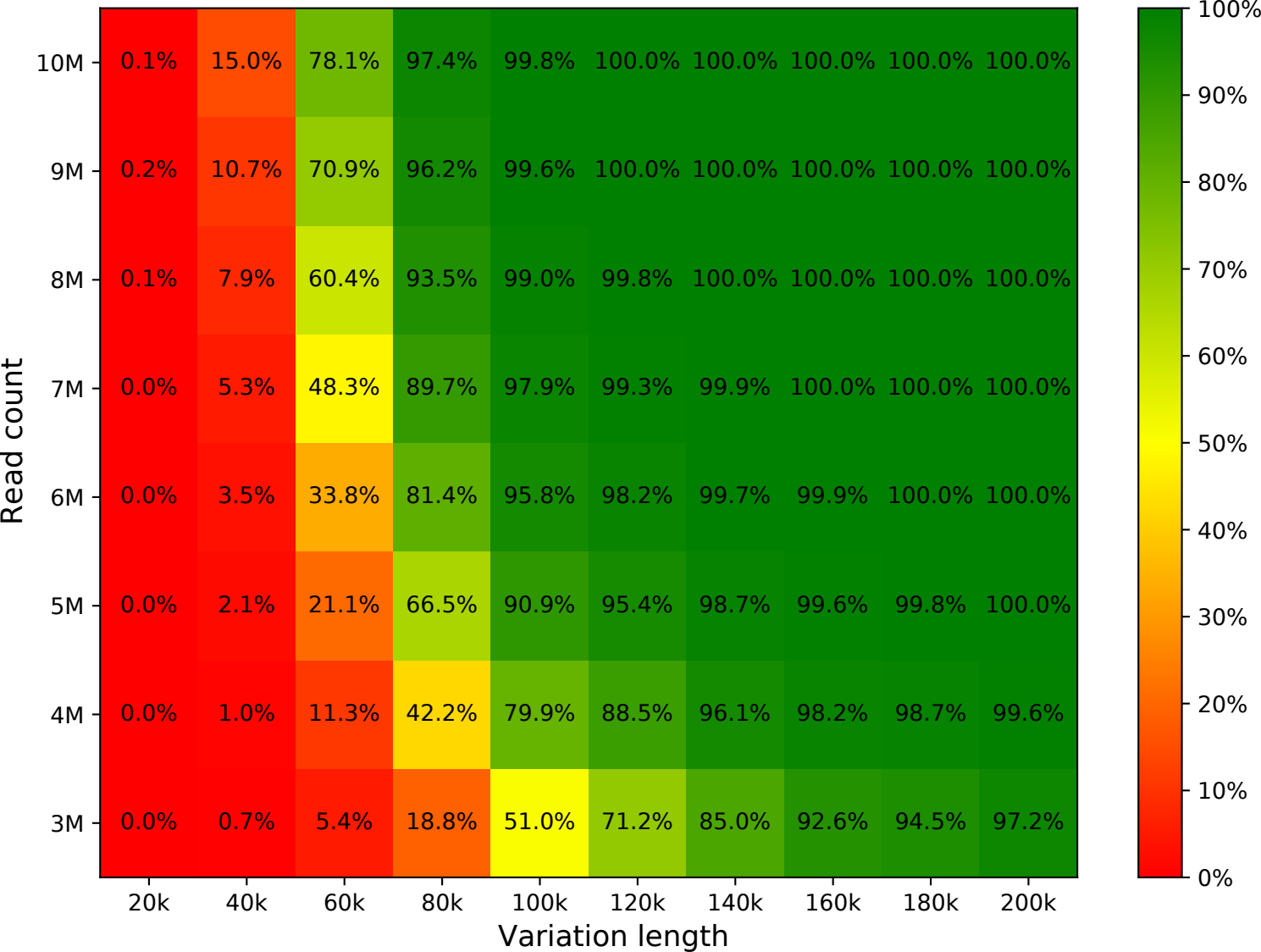
4213–4221.

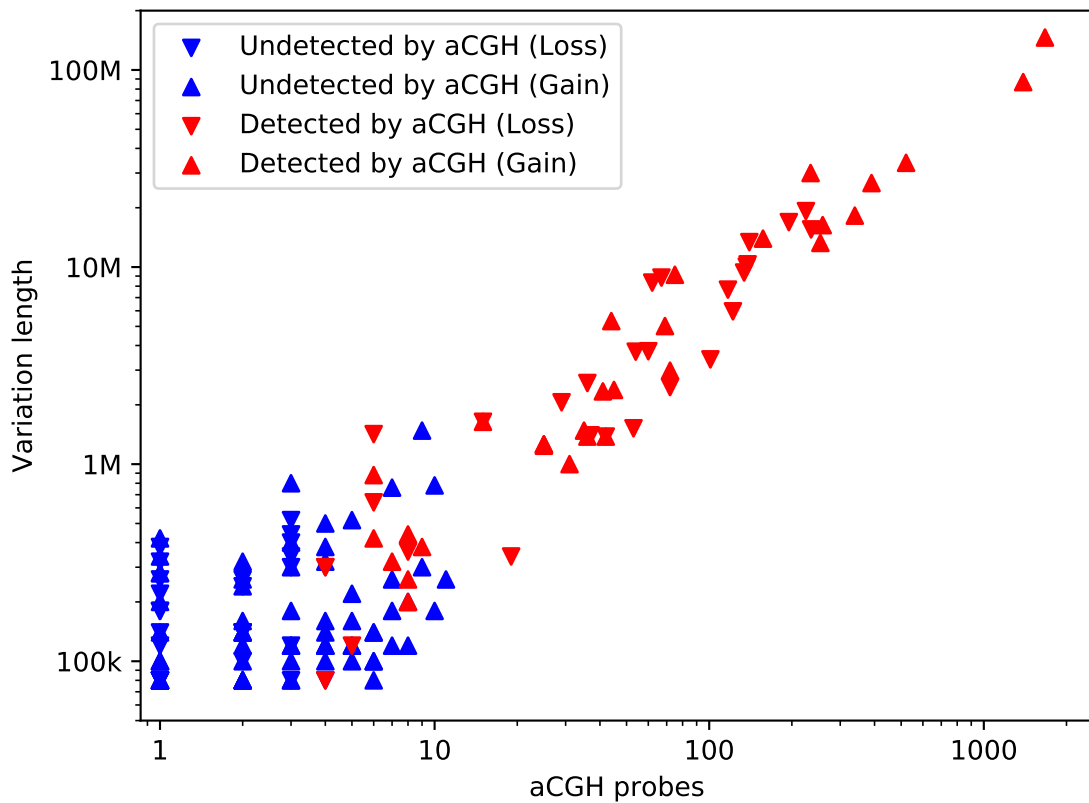
8. Bartha Á, Györfy B. Comprehensive Outline of Whole Exome Sequencing Data Analysis Tools Available in Clinical Oncology. *Cancers*. 2019;11. doi:10.3390/cancers11111725
9. Hyblova M, Harsanyova M, Nikulenkov-Grochova D, Kadlecova J, Kucharik M, Budis J, et al. Validation of Copy Number Variants Detection from Pregnant Plasma Using Low-Pass Whole-Genome Sequencing in Noninvasive Prenatal Testing-Like Settings. *Diagnostics (Basel)*. 2020;10. doi:10.3390/diagnostics10080569
10. Kucharik M, Gnip A, Hyblova M, Budis J, Strieskova L, Harsanyova M, et al. Non-invasive prenatal testing (NIPT) by low coverage genomic sequencing: Detection limits of screened chromosomal microdeletions. *PLOS ONE*. 2020. p. e0238245. doi:10.1371/journal.pone.0238245
11. Minarik G, Repiska G, Hyblova M, Nagyova E, Soltys K, Budis J, et al. Utilization of Benchtop Next Generation Sequencing Platforms Ion Torrent PGM and MiSeq in Noninvasive Prenatal Testing for Chromosome 21 Trisomy and Testing of Impact of In Silico and Physical Size Selection on Its Analytical Performance. *PLoS One*. 2015;10: e0144811.
12. Sekelska M, Izsakova A, Kubosova K, Tilandyova P, Csekes E, Kuchova Z, et al. Result of Prospective Validation of the Trisomy Test for the Detection of Chromosomal Trisomies. *Diagnostics (Basel)*. 2019;9. doi:10.3390/diagnostics9040138
13. Chandrananda D, Thorne NP, Ganesamoorthy D, Bruno DL, Benjamini Y, Speed TP, et al. Investigating and correcting plasma DNA sequencing coverage bias to enhance aneuploidy discovery. *PLoS One*. 2014;9: e86993.
14. Gazdarica J, Budis J, Duris F, Turna J, Szemes T. Adaptable Model Parameters in Non-Invasive Prenatal Testing Lead to More Stable Predictions. *Int J Mol Sci*. 2019;20. doi:10.3390/ijms20143414
15. Agilent Technologies, Inc. Human genome cgh microarray kit, 4x44k. [cited 25 Aug 2020]. Available: <https://www.agilent.com/en/product/cgh-cgh-snp-microarray-platform/cgh-cgh-snp-microarrays/human-microarrays/human-genome-cgh-microarray-kit-4x44k-228410>
16. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012;9: 357–359.
17. Zhao C, Tynan J, Ehrich M, Hannum G, McCullough R, Saldivar J-S, et al. Detection of fetal subchromosomal abnormalities by sequencing circulating cell-free DNA from maternal plasma. *Clin Chem*. 2015;61: 608–616.
18. Alkan C, Kidd JM, Marques-Bonet T, Aksay G, Antonacci F, Hormozdiari F, et al. Personalized copy number and segmental duplication maps using next-generation sequencing. *Nat Genet*. 2009;41: 1061–1067.
19. Seshan VE, Olshen A. DNACopy: DNA copy number data analysis.

Read count



Variation length





# Chromosome 08

