

# SV evaluation report sample FR05812606

February 10, 2019

$\mu$  mean in control,  $\sigma$  stdev in control,  $z$  z-score. OK if  $|z| \leq 2$ , else !!!

## 1 QC from bam alignment file

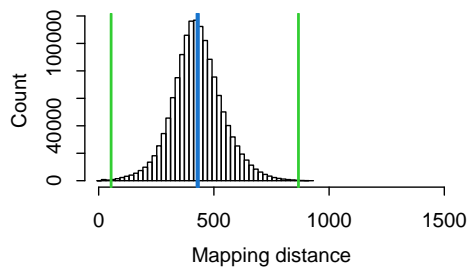
Input bam file: FR05812606.bam

Percent reads pairs not mapping as proper pair: 1.3 [ $\mu$  2.1,  $\sigma$  1.5,  $z$  -0.5] OK

Percent reads pairs mapping on different chromosomes: 0.68 [ $\mu$  1.2,  $\sigma$  1,  $z$  -0.5] OK

Stdev of coverage: 8.1 [ $\mu$  8.1,  $\sigma$  0.6,  $z$  0] OK

### Mapping distance distribution:



Concordant size min (green): 54 [ $\mu$  54,  $\sigma$  6.7,  $z$  0] OK

Concordant size max (green): 867 [ $\mu$  895,  $\sigma$  60,  $z$  -0.5] OK

Mean mapping distance size (blue): 430 [ $\mu$  444,  $\sigma$  26,  $z$  -0.5] OK

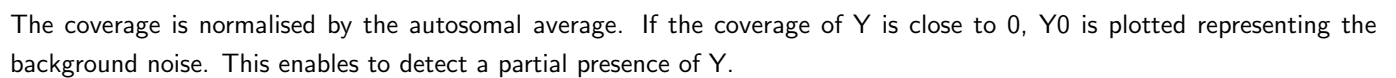
Stdev mapping distance size: 110 [ $\mu$  109,  $\sigma$  7.4,  $z$  0.1] OK

## 2 Input discordant pairs and split reads

Input discordant pairs: 10181760 [ $\mu$  12025111,  $\sigma$  7866150,  $z$  -0.2] OK

Input split reads: 1032913 [ $\mu$  1122920,  $\sigma$  174515,  $z$  -0.5] OK

Average sequence coverage: autosomes 38.4x, X 38.4x, Y 0x  
Inferred gender: XX (XX female, XY male, XXY Klinefelter's Syndrome)



## 4 Number called SVs

### By call confidence:

High confidence variant: 2044 [ $\mu$  1844,  $\sigma$  106,  $z$  1.9] OK<sup>b</sup>

Pass and High confidence variants: 5724 [ $\mu$  5841,  $\sigma$  323,  $z$  -0.4] OK

Low confidence variants: 4528 [ $\mu$  3983,  $\sigma$  666,  $z$  0.8] OK

### By SV type <sup>a</sup>:

CNVs: 4634 [ $\mu$  4730,  $\sigma$  187,  $z$  -0.5] OK

Loss: 3959 [ $\mu$  4025,  $\sigma$  173,  $z$  -0.4] OK

Gains: 675 [ $\mu$  705,  $\sigma$  54,  $z$  -0.6] OK

Balanced events including inversions and translocations: 1090 [ $\mu$  1110,  $\sigma$  147,  $z$  -0.1] OK

### By caller <sup>a</sup>:

Lumpy calls: 4154 [ $\mu$  3492,  $\sigma$  255,  $z$  2.6] OK<sup>b</sup>

CNVnator calls: 848 [ $\mu$  865,  $\sigma$  37,  $z$  -0.5] OK<sup>b</sup>

In both: 722 [ $\mu$  1484,  $\sigma$  87,  $z$  -8.8] !!!<sup>b</sup>

### Gene affecting / rare <sup>a</sup>:

Affecting genes: 2763 [ $\mu$  2832,  $\sigma$  170,  $z$  -0.4] OK

Rare Pass and High confidence: 87 [ $\mu$  41,  $\sigma$  23,  $z$  2] OK<sup>b</sup>

Rare Pass, High confidence, affecting genes: 49 [ $\mu$  22,  $\sigma$  12,  $z$  2.2] OK<sup>b</sup>

Rare Pass, High confidence CNV, affecting genes: 41 [ $\mu$  17,  $\sigma$  8,  $z$  3] OK<sup>b</sup>

<sup>a</sup> Without low confidence variants

<sup>b</sup> OK if  $|z| \leq 4$

## 5 Explanation QC and SV evaluation report

### 5.1 General QC

Low quality input data may lead to unexpected results. To guarantee the quality of the results several variables that have an impact on or are an indicator for the quality of the results are measured and compiled in this automated QC report. QC metrics measured for a particular sample are compared to the expected range obtained from analyzing 500 germline controls samples. The control samples represent previously analyzed healthy individuals (MGRB cohort) that passed QC. The expected range is generally defined as two times the standard deviation ( $|z| \leq 2$ ) from the mean of the control cohort, unless specified otherwise. If a measured metric is within expectations, it is marked with a green OK, else with three orange exclamation marks. ClinSV is robust to a few metrics being outside the expected range, but within 4 times the standard deviation.

#### Re 1. QC from bam alignment file

Read pairs from Illumina paired end sequencing do not always align to the reference with their expected distance (roughly 450 bp, depending fragmentation size and size selection), regardless of the presence of structural variation. The sequencing process produces a small percentage of chimeric read pairs. These pairs originate from distant genomic locations. Despite these chimeric reads being randomly distributed; elevated numbers will impact the SV calling. Indicators for the relative abundance of chimeric pairs is the percentage of reads not mapping as proper pairs and the percentage of pairs mapping on different chromosomes. An uneven read coverage can affect CNVnator resulting in an elevated number of false CNV calls. The un-evenness of the read coverage is reflected by an increased standard deviation of the read coverage. The number of discordantly mapping pairs the prediction program Lumpy can handle is finite. The threshold for when the read mapping distance is considered discordant for pairs mapping with the expected read orientation is automatically determined (see online methods section), and results are shown here. The insert size distribution and resulting thresholds for concordant mapping distances have an impact on the smallest detectable deletions. To save computing time, metrics in this section are estimated for a 10 mega base pair region on chromosome 1 (chr1:20,000,001-30,000,000).

#### Re 2. Input discordant pairs and split reads

Number of discordant read pairs and split reads used as input for Lumpy. Deviating numbers could indicate library preparation or sequencing issues, deeper coverage, or samples with high numbers of structural variation, as expected for cancer samples.

#### Re 3. Coverage by chromosome

The average sequence coverage was determined for all chromosomes (see methods). The number of sex chromosomes is inferred from the sequence coverage. Sex chromosome aneuploidy is visible here. This section also displays the chromosome wide coverage in intervals of 1 Mega bases. Grey dots below the black dots represent the average coverage in 500 control samples plus minus two times the standard deviation. The black dots indicate the coverage of the current sample. Truncated alignment files will not cover all grey dots. One Mega base segments greater than five times the standard deviation of the control are colored blue, highlighting regions that have a copy number gain, and segments less than five times the standard deviation are colored in red, highlighting regions of copy number loss. The standard deviation is used, because regions close to the centromere tend to show a greater variation that is still considered normal, thus will not get highlighted in blue or red. Large deletions or duplications, that are likely clinical significant will be visible in this representation. N-regions usually correspond to centromeric or telomeric regions of the chromosome. The sex chromosomes will be compared to the expected coverage of X, XX, Y and/or Y0 (Y-zero), depending the average coverage of X and Y. Y0 (Y-Zero) indicates unspecific background read coverage of the Y chromosome and is helpful to reveal a partial presence of Y.

#### Re 4. Number of called SVs

Number of called variants by call confidence, SV type, caller, and number of variants affecting genes and being rare. Some metrics in this section show greater variation and are allowed 4 times the standard deviation from the control average. For

instance the number of rare variants could be increased in an individual of a race underrepresented in the control cohort.

## **5.2 NA12878 SV evaluation**

The following two sections are to evaluate the SV recall rate of a NA12878 sample and allow assessing the fitness of the entire SV detection pipeline. Metrics are compared to average values of nine NA12878 control samples. Here z values greater or equal to -2 are acceptable, in order to not penalize a greater concordance than expected from the nine NA12878 control samples. This section appears when option -eval is set.

### **Re 5. Sensitivity**

This section shows the sensitivity of detecting gold standard deletion calls, as published by GIAB (Parikh et al. 2016), excluding 12 false positives > 500 bases (Minoche et al. 2017).

### **Re 6. Comparison to NA12878 sample**

Concordance of SV between NA12878 control sample (FR05812662) and current test sample, shown in percent of FR05812662 calls or in percent of test sample calls. High confidence calls generally have a higher reproducibility compared to all pass variants. CNVs and all SVs are tested separately.