# eurofins | **Genomics**

GATC Biotech AG, Jakob-Stadler-Platz 7, 78467 Konstanz

## Data Analysis Report: ONCOPANEL ALL-IN-ONE v1.4

Project / Study: FE-0242

Date: February 28, 2018

# Table of Contents

# 1 Results

## 1.1 Variant discovery

Single nucleotide variants (SNVs), Insertions and deletions (InDel) are detected in each sample using LoFreq[1], and are filtered based on mutation allele frequency (>1%) and coverage ($\geq$ 10% of average coverage excluding duplicated fragments; coverage metrics can be found in chapter 2.5). Variants that pass these thresholds are summarised in the following table(s).

Table 1: Variant metrics for HD753_3, HD753_4.

|  | HD753_3 | HD753_4 |
|---|---|---|
| Total SNV | 72893 | 72170 |
| Known SNV | 9817 | 9750 |
| Unknown SNV | 63076 | 62420 |
| Total InDel | 26616 | 28059 |
| Known InDel | 1145 | 1191 |
| Unknown InDel | 25471 | 26868 |

Known SNV / InDel: in reference variant databases (dbSNP, COSMIC[2] and / or ClinVar[3]).

Unknown SNV / InDel: currently not listed in reference variant database (as aforementioned).

## 1.2 Sample-wise known clinical significant variants

Variants detected are screened for known clinical significance in ClinVar (released 02. Oct 2017) [3] database. The ClinVar database aggregates information about genomic variation and its relationship to human health. It is hosted by the National Center for Biotechnology Information (NCBI). Detailed explanation of clinical significance in ClinVar database can be found at https://www.ncbi.nlm.nih.gov/clinvar/docs/clinsig/.

Variants which have clinical significance state as "Likely pathogenic", "Pathogenic" and "Drug response" are filtered from the complete list of variants and are reported in following table(s). For more detailed information navigate to the Clinvar database and type in the dbIDs of your variant of interest. Variant effects for multiple transcripts for the same variant are listed as separate entries. In case of multiple transcripts, transcripts which have missense, splice junction, UTR, frameshift, disruptive frameshift insertion / deletion variant types are listed.

### 1.2.1 HD753_3 Results

Table 2: Variants (SNV and InDels) in sample - **HD753_3.** Entries are sorted by gene.

| Location | Gene | AA Change | Codon Change | Mutation Freq. | Depth | ClinVar ID | ClinVar Significance |
|---|---|---|---|---|---|---|---|
| chr7:87160618 | ABCB1 | p.S893A<br>p.S829A | c.2677T>G<br>c.2485T>G | 64.9 % | 1394x | 166622 | drug response |
| chr5:131931451 | AC116366.3 | .<br>.<br>.<br>. | c.*2341delA<br>c.*1321delA<br>c.*2025delA<br>c.*2155delA | 15.0 % | 1372x | 408407 | pathogenic |

| Location | Gene | AA Change | Codon Change | Mutation Freq. | Depth | ClinVar ID | ClinVar Significance |
|---|---|---|---|---|---|---|---|
| chr14:105246551 | AKT1 | . <br> p.E17K | n.80G>A <br> c.49G>A | 5.4 % | 948x | 13983 | pathogenic |
| chr10:96540410 | AL583836.1 | . | c.*394G>A | 5.0 % | 1191x | 16899 | drug response |
| chr10:96541616 | AL583836.1 | . | c.*439G>A | 19.8 % | 1024x | 16897 | drug response |
| chr20:31022441 | ASXL1 | p.G641fs <br> p.G646fs | c.1919insG <br> c.1934insG | 4.5 % | 661x | 426927 | pathogenic |
| chr11:108205832 | ATM | p.V2716A | c.8147T>C | 4.1 % | 1861x | 142700 | pathogenic |
| chr17:63532584 | AXIN2 | p.G600fs <br> p.G665fs | c.1799delG <br> c.1994delG | 5.3 % | 452x | 5880 | pathogenic |
| chr19:49458970 | BAX | p.R24fs <br> p.E24fs <br> p.E41fs | c.69delG <br> c.70delG <br> c.121delG | 6.3 % | 867x | 9512 | pathogenic |
| chr7:140453136 | BRAF | . <br> p.V207E <br> p.V600E <br> p.V28E | c.*1249T>A <br> c.620T>A <br> c.1799T>A <br> c.83T>A | 18.5 % | 1685x | 13961 | pathogenic |
| chr17:41234451 | BRCA1 | . | c.*4110C>T | 5.0 % | 1824x | 17675 | pathogenic |
| chr13:32937354 | BRCA2 | p.I2675fs | c.8021insA | 5.2 % | 1044x | 267050 | pathogenic |
| chr9:21971186 | CDKN2A | . <br> p.P72L | c.*95C>T <br> c.215C>T | 8.1 % | 211x | 376310 | likely pathogenic |
| chr15:93545433 | CHD2 | p.Q1392fs <br> . <br> . | c.4173insA <br> c.*406insA <br> c.*344insA | 7.6 % | 262x | 218395 | pathogenic |
| chr3:41266101 | CTNNB1 | p.S33Y <br> p.S26Y | c.98C>A <br> c.77C>A | 5.6 % | 1473x | 17577 | pathogenic |
| chr10:17113456 | CUBN | p.S865N | c.2594G>A | 4.6 % | 798x | 265086 | pathogenic |
| chr19:41512841 | CYP2B6 | p.Q172H | c.516G>T | 24.3 % | 1945x | 29671 | drug response |
| chr19:41515263 | CYP2B6 | p.K262R | c.785A>G | 25.2 % | 547x | 120171 | drug response |
| chr10:96702047 | CYP2C9 | p.R144C | c.430C>T | 9.0 % | 1718x | 8409 | drug response |
| chr10:96741053 | CYP2C9 | p.I359L | c.1075A>C | 5.7 % | 1307x | 8408 | drug response |
| chr22:42524947 | CYP2D6 | . <br> . <br> . <br> . <br> . | c.353-1G>A <br> c.440-1G>A <br> c.506-1G>A <br> n.1198-1G>A <br> c.173-1G>A | 32.4 % | 550x | 16889 | drug response |
| chr22:42526694 | CYP2D6 | p.P34S <br> p.P12S | c.100C>T <br> c.34C>T | 52.0 % | 252x | 16893 | drug response |
| chr7:55242464 | EGFR | p.E746_A750del <br> p.E701_A705del <br> . | c.2235_2249delGGAATTAAGAGAAGC <br> c.2100_2114delGGAATTAAGAGAAGC <br> c.*225_*239delGGAATTAAGAGAAGC | 2.8 % | 1960x | 163343 | drug response |
| chr7:55248998 | EGFR | p.A767_V769ins <br> . <br> p.A722_V724ins | c.2300_2308insCCAGCGTGG <br> c.*290_*298insCCAGCGTGG <br> c.2165_2173insCCAGCGTGG | 3.4 % | 1442x | 177678 | drug response |
| chr22:41565529 | EP300 | p.D1399N | c.4195G>A | 11.0 % | 1138x | 376401 | likely pathogenic |

| Location | Gene | AA Change | Codon Change | Mutation Freq. | Depth | ClinVar ID | ClinVar Significance |
|---|---|---|---|---|---|---|---|
| chr8:118849384 | EXT1 | p.R129H<br>p.R340H | c.386G>A<br>c.1019G>A | 11.4 % | 1252x | 265129 | pathogenic |
| chr5:176520243 | FGFR4 | p.G23R<br>p.G388R | c.67G>A<br>c.1162G>A | 35.0 % | 443x | 16326 | pathogenic |
| chr19:3118942 | GNA11 | p.Q57L<br>p.Q209L | c.170A>T<br>c.626A>T | 5.8 % | 950x | 376002 | pathogenic |
| chr11:67352689 | GSTP1 | .<br>p.I105V | c.*137A>G<br>c.313A>G | 61.6 % | 1461x | 37340 | drug response |
| chr12:121432114 | HNF1A | p.P291fs<br>p.G226fs<br>. | c.864delG<br>c.677delG<br>c.*304delG | 6.3 % | 748x | 435424 | pathogenic |
| chr7:142640113 | KEL | p.L597P | c.1790T>C | 16.0 % | 1504x | 31082 | pathogenic |
| chr12:25398281 | KRAS | p.G13D | c.38G>A | 4.3 % | 1689x | 12580 | pathogenic |
| chr15:66727451 | MAP2K1 | p.Q56P | c.167A>C | 3.7 % | 1341x | 375978 | pathogenic |
| chr15:66729147 | MAP2K1 | p.H119Y | c.355C>T | 4.5 % | 1772x | 40741 | pathogenic |
| chr5:79970914 | MSH3 | p.K383fs | c.1148delA | 31.5 % | 1272x | 8738 | pathogenic |
| chr2:48030639 | MSH6 | p.F56fs<br>p.F1088fs<br>.<br>p.F958fs<br>p.F786fs | c.165delC<br>c.3261delC<br>c.*2608delC<br>c.2871delC<br>c.2355delC | 1.1 % | 1595x | 89363 | pathogenic |
| chr2:48030639 | MSH6 | p.F786fs<br>p.F1088fs<br>.<br>p.F958fs<br>p.F56fs | c.2355insC<br>c.3261insC<br>c.*2608insC<br>c.2871insC<br>c.165insC | 5.1 % | 1595x | 89364 | pathogenic |
| chr17:29553477 | NF1 | p.P678fs<br>p.P344fs<br>.<br>p.P712fs | c.2033delC<br>c.1031delC<br>c.*1434delC<br>c.2135delC | 5.6 % | 1194x | 428991 | pathogenic |
| chr17:29553477 | NF1 | p.I345fs<br>p.I679fs<br>.<br>p.I713fs | c.1031insC<br>c.2033insC<br>c.*1434insC<br>c.2135insC | 4.2 % | 1194x | 141513 | pathogenic |
| chr3:178936091 | PIK3CA | p.E545K | c.1633G>A | 5.0 % | 1291x | 13655 | pathogenic |
| chr3:178947865 | PIK3CA | p.G914R | c.2740G>A | 5.6 % | 2226x | 39703 | pathogenic |
| chr3:178952085 | PIK3CA | p.H1047R | c.3140A>G | 14.4 % | 1957x | 13652 | pathogenic |
| chr5:131931451 | RAD50 | p.K722fs<br>.<br>.<br>p.?661fs | c.2165delA<br>c.*1791delA<br>c.*351delA<br>c.1982delA | 15.0 % | 1372x | 408407 | pathogenic |
| chr12:21331549 | SLCO1B1 | p.V174A | c.521T>C | 19.6 % | 1539x | 37346 | drug response |
| chr7:141672604 | TAS2R38 | p.I296V | c.886A>G | 58.6 % | 2304x | 2906 | drug response |
| chr7:141673345 | TAS2R38 | p.A49P | c.145G>C | 55.0 % | 2233x | 2904 | drug response |
| chr17:7577559 | TP53 | p.S82F<br>p.S230F<br>p.S202F<br>p.S109F<br>p.S148F<br>p.S241F | c.245C>T<br>c.689C>T<br>c.605C>T<br>c.326C>T<br>c.443C>T<br>c.722C>T | 5.9 % | 752x | 12359 | pathogenic |

| Location | Gene | AA Change | Codon Change | Mutation Freq. | Depth | ClinVar ID | ClinVar Significance |
|---|---|---|---|---|---|---|---|
| chr17:7577559 | TP53 | p.S148C<br>p.S82C<br>p.S241C<br>p.S230C<br>p.S109C<br>p.S202C | c.443C>G<br>c.245C>G<br>c.722C>G<br>c.689C>G<br>c.326C>G<br>c.605C>G | 5.7 % | 752x | 177791 | likely pathogenic |
| chr17:7579472 | TP53 | p.P33R<br>p.P72R | c.98C>G<br>c.215C>G | 80.5 % | 931x | 12351 | drug response |
| chr21:44524456 | U2AF1 | p.S34F<br>.<br>. | c.101C>T<br>c.-186C>T<br>c.-119C>T | 8.1 % | 677x | 376025 | likely pathogenic |
| chr21:44524456 | U2AF1L5 | p.S34F<br>.<br>. | c.101C>T<br>c.-186C>T<br>c.-119C>T | 8.1 % | 677x | 376025 | likely pathogenic |
| chr3:14187449 | XPC | p.Q939K<br>. | c.2815C>A<br>c.*2268C>A | 40.5 % | 570x | 190215 | drug response |

## 1.2.2 HD753_4 Results

Table 3: Variants (SNV and InDels) in sample - **HD753_4.** Entries are sorted by gene.

| Location | Gene | AA Change | Codon Change | Mutation Freq. | Depth | ClinVar ID | ClinVar Significance |
|---|---|---|---|---|---|---|---|
| chr7:87160618 | ABCB1 | p.S893A<br>p.S829A | c.2677T>G<br>c.2485T>G | 67.4 % | 1281x | 166622 | drug response |
| chr5:131931451 | AC116366.3 | .<br>.<br>.<br>. | c.*2341delA<br>c.*1321delA<br>c.*2025delA<br>c.*2155delA | 15.4 % | 1267x | 408407 | pathogenic |
| chr14:105246551 | AKT1 | .<br>p.E17K | n.80G>A<br>c.49G>A | 5.4 % | 905x | 13983 | pathogenic |
| chr10:96540410 | AL583836.1 | . | c.*394G>A | 4.1 % | 1128x | 16899 | drug response |
| chr10:96541616 | AL583836.1 | . | c.*439G>A | 20.8 % | 942x | 16897 | drug response |
| chr20:31022441 | ASXL1 | p.G641fs<br>p.G646fs | c.1919insG<br>c.1934insG | 4.1 % | 580x | 426927 | pathogenic |
| chr11:108205832 | ATM | p.V2716A | c.8147T>C | 4.1 % | 1653x | 142700 | pathogenic |
| chr17:63532584 | AXIN2 | p.G600fs<br>p.G665fs | c.1799delG<br>c.1994delG | 5.6 % | 393x | 5880 | pathogenic |
| chr19:49458970 | BAX | p.E41fs<br>p.R24fs<br>p.E24fs | c.121insG<br>c.69insG<br>c.70insG | 0.8 % | 860x | 9511 | pathogenic |
| chr7:140453136 | BRAF | .<br>p.V207E<br>p.V600E<br>p.V28E | c.*1249T>A<br>c.620T>A<br>c.1799T>A<br>c.83T>A | 18.3 % | 1597x | 13961 | pathogenic |
| chr17:41234451 | BRCA1 | . | c.*4110C>T | 5.1 % | 1712x | 17675 | pathogenic |
| chr13:32937354 | BRCA2 | p.I2675fs | c.8021insA | 6.1 % | 906x | 267050 | pathogenic |
| chr9:21971186 | CDKN2A | .<br>p.P72L | c.*95C>T<br>c.215C>T | 11.8 % | 211x | 376310 | likely pathogenic |
| chr15:93545433 | CHD2 | p.Q1392fs<br>.<br>. | c.4173insA<br>c.*406insA<br>c.*344insA | 11.5 % | 252x | 218395 | pathogenic |
| chr3:41266101 | CTNNB1 | p.S33Y<br>p.S26Y | c.98C>A<br>c.77C>A | 6.1 % | 1382x | 17577 | pathogenic |
| chr10:17113456 | CUBN | p.S865N | c.2594G>A | 4.7 % | 654x | 265086 | pathogenic |
| chr19:41512841 | CYP2B6 | p.Q172H | c.516G>T | 25.8 % | 1905x | 29671 | drug response |
| chr19:41515263 | CYP2B6 | p.K262R | c.785A>G | 24.5 % | 477x | 120171 | drug response |
| chr10:96702047 | CYP2C9 | p.R144C | c.430C>T | 10.4 % | 1648x | 8409 | drug response |
| chr10:96741053 | CYP2C9 | p.I359L | c.1075A>C | 5.2 % | 1253x | 8408 | drug response |
| chr22:42524947 | CYP2D6 | .<br>.<br>.<br>.<br>. | c.353-1G>A<br>c.440-1G>A<br>c.506-1G>A<br>n.1198-1G>A<br>c.173-1G>A | 34.1 % | 554x | 16889 | drug response |
| chr22:42526694 | CYP2D6 | p.P34S<br>p.P12S | c.100C>T<br>c.34C>T | 53.3 % | 242x | 16893 | drug response |

| Location | Gene | AA Change | Codon Change | Mutation Freq. | Depth | ClinVar ID | ClinVar Significance |
|---|---|---|---|---|---|---|---|
| chr7:55242464 | EGFR | p.E746_ A750del p.E701_ A705del . | c.2235_ 2249delGGAATTAAGAGAAGC c.2100_ 2114delGGAATTAAGAGAAGC c.*225_ *239delGGAATTAAGAGAAGC | 2.6 % | 1911x | 163343 | drug response |
| chr7:55248998 | EGFR | p.A767_ V769ins . p.A722_ V724ins | c.2300_ 2308insCCAGCGTGG c.*290_ *298insCCAGCGTGG c.2165_ 2173insCCAGCGTGG | 4.5 % | 1325x | 177678 | drug response |
| chr22:41565529 | EP300 | p.D1399N | c.4195G>A | 9.8 % | 1018x | 376401 | likely pathogenic |
| chr8:118849384 | EXT1 | p.R129H p.R340H | c.386G>A c.1019G>A | 11.1 % | 1140x | 265129 | pathogenic |
| chr5:176520243 | FGFR4 | p.G23R p.G388R | c.67G>A c.1162G>A | 28.7 % | 401x | 16326 | pathogenic |
| chr19:3118942 | GNA11 | p.Q57L p.Q209L | c.170A>T c.626A>T | 4.3 % | 857x | 376002 | pathogenic |
| chr11:67352689 | GSTP1 | . p.I105V | c.*137A>G c.313A>G | 58.7 % | 1398x | 37340 | drug response |
| chr12:121432114 | HNF1A | p.P291fs p.G226fs . | c.864delG c.677delG c.*304delG | 5.0 % | 723x | 435424 | pathogenic |
| chr7:142640113 | KEL | p.L597P | c.1790T>C | 17.4 % | 1431x | 31082 | pathogenic |
| chr12:25398281 | KRAS | p.G13D | c.38G>A | 4.6 % | 1534x | 12580 | pathogenic |
| chr15:66727451 | MAP2K1 | p.Q56P | c.167A>C | 5.1 % | 1213x | 375978 | pathogenic |
| chr15:66729147 | MAP2K1 | p.H119Y | c.355C>T | 5.0 % | 1676x | 40741 | pathogenic |
| chr5:79970914 | MSH3 | p.K383fs | c.1148delA | 30.6 % | 1058x | 8738 | pathogenic |
| chr2:48030639 | MSH6 | p.F56fs p.F1088fs . p.F958fs p.F786fs | c.165delC c.3261delC c.*2608delC c.2871delC c.2355delC | 0.9 % | 1662x | 89363 | pathogenic |
| chr2:48030639 | MSH6 | p.F786fs p.F1088fs . p.F958fs p.F56fs | c.2355insC c.3261insC c.*2608insC c.2871insC c.165insC | 5.2 % | 1662x | 89364 | pathogenic |
| chr17:29553477 | NF1 | p.I345fs p.I679fs . p.I713fs | c.1031insC c.2033insC c.*1434insC c.2135insC | 4.8 % | 1140x | 141513 | pathogenic |
| chr17:29553477 | NF1 | p.P678fs p.P344fs . p.P712fs | c.2033delC c.1031delC c.*1434delC c.2135delC | 5.8 % | 1140x | 428991 | pathogenic |
| chr16:23646980 | PALB2 | p.M1fs p.M296fs | c.1delA c.886delA | 0.3 % | 1369x | 143979 | pathogenic |
| chr3:178936091 | PIK3CA | p.E545K | c.1633G>A | 4.8 % | 1238x | 13655 | pathogenic |
| chr3:178947865 | PIK3CA | p.G914R | c.2740G>A | 4.6 % | 2110x | 39703 | pathogenic |
| chr3:178952085 | PIK3CA | p.H1047R | c.3140A>G | 15.2 % | 1846x | 13652 | pathogenic |

| Location | Gene | AA Change | Codon Change | Mutation Freq. | Depth | ClinVar ID | ClinVar Significance |
|---|---|---|---|---|---|---|---|
| chr5:131931451 | RAD50 | p.K722fs<br>.<br>.<br>p.?661fs | c.2165delA<br>c.*1791delA<br>c.*351delA<br>c.1982delA | 15.4 % | 1267x | 408407 | pathogenic |
| chr12:21331549 | SLCO1B1 | p.V174A | c.521T>C | 18.1 % | 1494x | 37346 | drug response |
| chr7:141672604 | TAS2R38 | p.I296V | c.886A>G | 59.1 % | 2209x | 2906 | drug response |
| chr7:141673345 | TAS2R38 | p.A49P | c.145G>C | 54.4 % | 2126x | 2904 | drug response |
| chr17:7577559 | TP53 | p.S148C<br>p.S82C<br>p.S241C<br>p.S230C<br>p.S109C<br>p.S202C | c.443C>G<br>c.245C>G<br>c.722C>G<br>c.689C>G<br>c.326C>G<br>c.605C>G | 4.2 % | 782x | 177791 | likely pathogenic |
| chr17:7577559 | TP53 | p.S82F<br>p.S230F<br>p.S202F<br>p.S109F<br>p.S148F<br>p.S241F | c.245C>T<br>c.689C>T<br>c.605C>T<br>c.326C>T<br>c.443C>T<br>c.722C>T | 6.8 % | 782x | 12359 | pathogenic |
| chr17:7579472 | TP53 | p.P33R<br>p.P72R | c.98C>G<br>c.215C>G | 81.2 % | 807x | 12351 | drug response |
| chr21:44524456 | U2AF1 | p.S34F<br>.<br>. | c.101C>T<br>c.-186C>T<br>c.-119C>T | 8.2 % | 622x | 376025 | likely pathogenic |
| chr21:44524456 | U2AF1L5 | p.S34F<br>.<br>. | c.101C>T<br>c.-186C>T<br>c.-119C>T | 8.2 % | 622x | 376025 | likely pathogenic |
| chr3:14187449 | XPC | p.Q939K<br>. | c.2815C>A<br>c.*2268C>A | 34.0 % | 580x | 190215 | drug response |

## 1.3   Copy number analysis

Copy number variations (CNV) are detected using the software package CNVkit[4] which uses normalized read depths to infer copy number evenly across the exome/genome. CNVkit uses both the on-target reads and the nonspecifically captured off-target reads to calculate log2 copy ratios across the genome for each sample. Briefly, off-target bins are assigned from the genomic positions between targeted regions, with the average off-target bin size being much larger than the average on-target bin to match their read counts. Both the on and off target locations are then separately used to calculate the mean read depth within each interval. The on and off target read depths are then combined, normalized to a reference derived from control samples, corrected for several systematic biases (GC content, sequence complexity and targets) to result in a final table of log2 copy ratios. Then, the segmentation algorithm uses log2 ratio values to infer discrete copy number events.Copy number events with minimum 100 x coverage are reported.

Table 4: Case vs Control setup.

| Case | Control(s) |
|------|------------|
| HD753_3 | NA12878_Pv2_052017_CNV_R1, NA12878_Pv2_052017_CNV_R2 |
| HD753_4 | NA12878_Pv2_052017_CNV_R1, NA12878_Pv2_052017_CNV_R2 |

Table 5: Summary of CNV events detected in each sample.

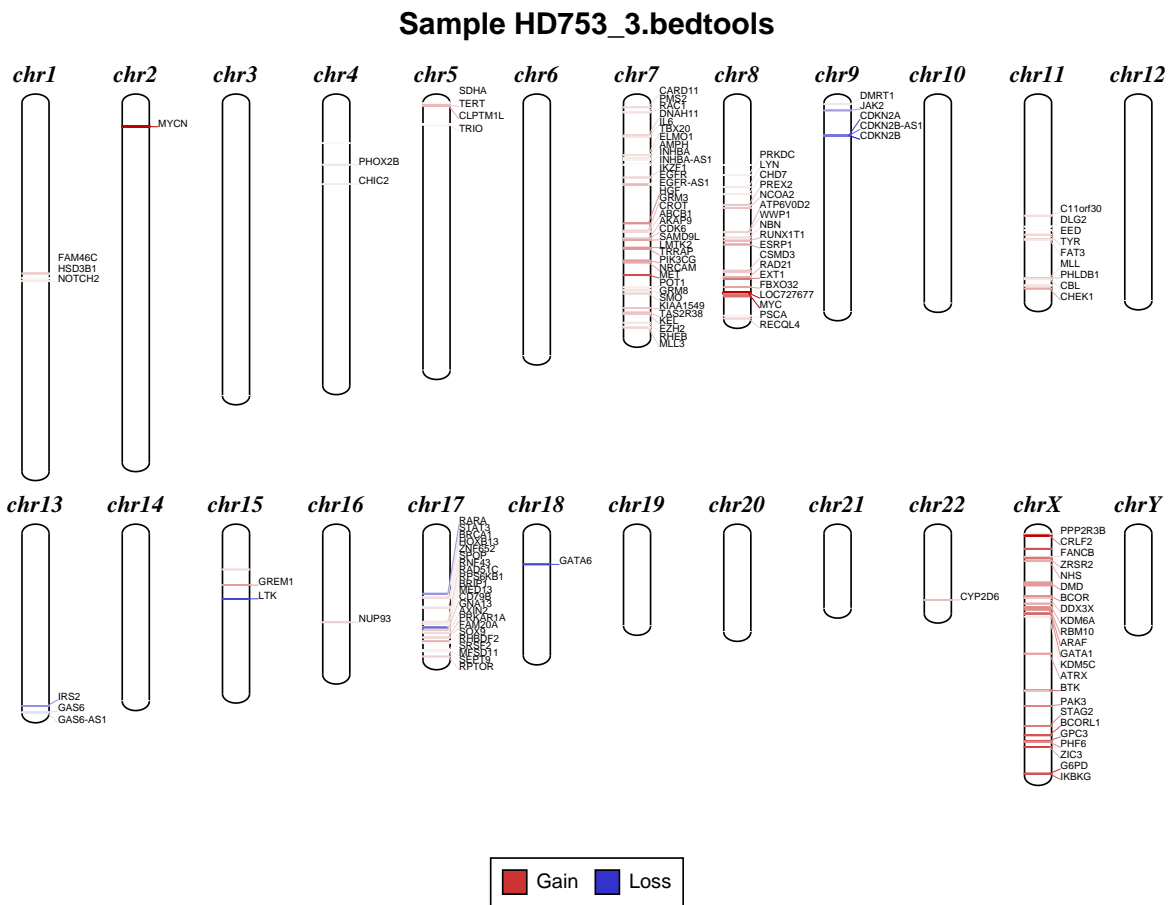| Sample | Duplication Events | Deletion Events |
|--------|-------------------:|----------------:|
| HD753_3 | 24 | 10 |
| HD753_4 | 23 | 14 |

### 1.3.1    HD753_3 Results



Figure 1: Ideogram representing chromosome wise copy number events observed in sample HD753_3. Copy gain events are drawn in red and copy loss events are drawn in blue.

Table 6: Duplication events detected in sample HD753_3. Gene column lists the name of genes (HGNC convention), CN column contains copy number observed and Depth column displays the coverage depth at the location (Loci column).

| Gene | CN | Depth | Loci |
|------|----|-------|------|
| MYCNOS, MYCN | 9 | 3407.48 | chr2:15729629-16086340 |
| LOC727677, MYC | 9 | 4076.32 | chr8:128093106-129192460 |
| MET | 4 | 5155.89 | chr7:116312530-116339669 |
| MET | 4 | 2222.36 | chr7:116339669-116436309 |
| TAS2R38, KEL, EZH2, RHEB, MLL3 | 3 | 1453.71 | chr7:141672348-151879671 |
| STAT3, BRCA1 | 3 | 1440.28 | chr17:40469026-41248090 |
| RNF43, RAD51C, RPS6KB1, BRIP1, MED13, CD79B, GNA13, AXIN2, PRKAR1A, FAM20A, SOX9, RHBDF2, SRSF2, SRSF2, MIR636, SRSF2, MIR636, MFSD11, SRSF2, MFSD11, SEPT9, RPTOR | 3 | 1046.06 | chr17:56432135-78938285 |
| PSCA, RECQL4 | 3 | 699.78 | chr8:129543758-145743186 |
| PREX2, NCOA2, ATP6V0D2, WWP1, NBN, RUNX1T1, ESRP1 | 3 | 1219.44 | chr8:68864809-95718352 |
| POT1, GRM8, SMO, KIAA1549 | 3 | 1291.67 | chr7:124503268-138604337 |

| Gene | CN | Depth | Loci |
|------|-----|-------|------|
| NUP93 | 3 | 1229.2 | chr16:56782028-56878711 |
| NRAS, FAM46C, HSD3B1, NOTCH2 | 3 | 1236.54 | chr1:115258835-120572569 |
| NIPA2, GREM1 | 3 | 1016.07 | chr15:23021376-33023568 |
| MLL, PHLDB1, CBL, CHEK1 | 3 | 1552.92 | chr11:118318180-125497553 |
| KAT6A, PRKDC, LYN, CHD7, PREX2 | 3 | 1130.77 | chr8:41906529-68864809 |
| HOXB13, ZNF652, SPOP | 3 | 1150.96 | chr17:46804028-47700323 |
| HGF, GRM3, CROT, ABCB1, AKAP9 | 3 | 1282.82 | chr7:81399053-91719132 |
| EXT1, FBXO32 | 3 | 1330.89 | chr8:118811823-124553306 |
| CYP2D6 | 3 | 702.1 | chr22:42522943-42525265 |
| CSMD3, RAD21 | 3 | 1136.85 | chr8:113236859-117879116 |
| CCDC127, SDHA, SDHA, TERT, CLPTM1L, TRIO | 3 | 941.97 | chr5:218285-14532057 |
| CARD11, PMS2, RAC1, DNAH11, IL6, TBX20, ELMO1, AMPH, INHBA, INHBA, INHBA-AS1, IKZF1, EGFR, EGFR, EGFR-AS1 | 3 | 1072.63 | chr7:2946153-55273441 |
| C11orf30, DLG2, EED, TYR, FAT3 | 3 | 1487.38 | chr11:76157827-92624397 |
| AKAP9, CDK6, SAMD9L, LMTK2, TRRAP, PIK3CG, NRCAM | 3 | 1483.3 | chr7:91722300-107880675 |

Table 7: Deletion events detected in sample HD753_3. Gene column lists the name of genes (HGNC convention), CN column contains copy number observed and Depth column displays the coverage depth at the location (Loci column).

| Gene | CN | Depth | Loci |
|------|-----|-------|------|
| RARA | 1 | 2956.98 | chr17:38497282-38499676 |
| PPP2R3B, CRLF2, FANCB, ZRSR2, NHS, DMD, BCOR, DDX3X, KDM6A, RBM10, ARAF, GATA1, KDM5C | 1 | 665.76 | chrX:320389-53222482 |
| PHOX2B, CHIC2 | 1 | 598.86 | chr4:27162572-54930627 |
| LTK | 1 | 113.2 | chr15:41803307-41804311 |
| IRS2, GAS6-AS1, GAS6, GAS6 | 1 | 322.92 | chr13:110434353-114566778 |
| GATA6 | 1 | 109.41 | chr18:19751090-19757211 |
| DMRT1, JAK2 | 1 | 648.78 | chr9:845380-5126957 |
| CDKN2A, CDKN2B-AS1, CDKN2B, CDKN2B-AS1 | 1 | 357.65 | chr9:21968129-22062228 |
| BTK, PAK3, STAG2, BCORL1, GPC3, PHF6, ZIC3, G6PD, G6PD, IKBKG | 1 | 733.32 | chrX:100604718-153775209 |
| ATRX | 1 | 764.09 | chrX:76776109-77041667 |

### 1.3.2 HD753_4 Results



Figure 2: Ideogram representing chromosome wise copy number events observed in sample HD753_4. Copy gain events are drawn in red and copy loss events are drawn in blue.

Table 8: Duplication events detected in sample HD753_4. Gene column lists the name of genes (HGNC convention), CN column contains copy number observed and Depth column displays the coverage depth at the location (Loci column).

| Gene | CN | Depth | Loci |
|---|---|---|---|
| MYCNOS, MYCN | 9 | 3210.06 | chr2:15729629-16086340 |
| LOC727677, MYC | 9 | 3808.59 | chr8:128093106-129192460 |
| NRCAM, MET | 4 | 4735.58 | chr7:107880235-116335700 |
| MET | 4 | 3046.81 | chr7:116335700-116436309 |
| STAT3, BRCA1 | 3 | 1354.15 | chr17:40468930-41248090 |
| SAMD9L, LMTK2, TRRAP, PIK3CG, NRCAM | 3 | 1436.99 | chr7:92760906-107878400 |
| RNF43, RAD51C, RPS6KB1, BRIP1, MED13, CD79B, GNA13, AXIN2, PRKAR1A, FAM20A, SOX9, RHBDF2, SRSF2, SRSF2, MIR636, SRSF2, MIR636, MFSD11, SRSF2, MFSD11, SEPT9, RPTOR | 3 | 985.26 | chr17:56432135-78938285 |
| PSCA, RECQL4 | 3 | 683.55 | chr8:129543758-145743186 |
| PRKDC, LYN, CHD7, PREX2 | 3 | 1057.83 | chr8:48686615-68864809 |
| PREX2, NCOA2, ATP6V0D2, WWP1, NBN, RUNX1T1, ESRP1 | 3 | 1134.78 | chr8:68864809-95718352 |

| Gene | CN | Depth | Loci |
|------|----|-------|------|
| POT1, GRM8, SMO, KIAA1549, BRAF, TAS2R38, KEL, EZH2, RHEB, MLL3 | 3 | 1383.66 | chr7:124462349-151879671 |
| NUP93 | 3 | 1156.22 | chr16:56782028-56878711 |
| MLL, PHLDB1, CBL, CHEK1 | 3 | 1443.7 | chr11:118307674-125497553 |
| HOXB13, ZNF652, SPOP | 3 | 1078.12 | chr17:46804028-47700323 |
| HGF, GRM3, CROT, ABCB1, AKAP9, CDK6, SAMD9L | 3 | 1198.66 | chr7:81399053-92760906 |
| FBXO32 | 3 | 1051.59 | chr8:124526431-124553306 |
| EXT1, FBXO32 | 3 | 1301.2 | chr8:118811823-124526431 |
| EED, TYR, FAT3 | 3 | 1636.83 | chr11:85989283-92624397 |
| CYP2D6 | 3 | 680.21 | chr22:42522943-42525369 |
| CSMD3, RAD21 | 3 | 1062.6 | chr8:113236859-117879116 |
| CCDC127, SDHA, SDHA, TERT, CLPTM1L, TRIO | 3 | 896.26 | chr5:218285-14532057 |
| CARD11, PMS2, RAC1, DNAH11, IL6, TBX20, ELMO1, AMPH, INHBA, INHBA, INHBA-AS1, IKZF1, EGFR, EGFR, EGFR-AS1 | 3 | 1010.54 | chr7:2946153-55273441 |
| C11orf30 | 3 | 1191.31 | chr11:76157827-76261312 |

Table 9: Deletion events detected in sample HD753_4. Gene column lists the name of genes (HGNC convention), CN column contains copy number observed and Depth column displays the coverage depth at the location (Loci column).

| Gene | CN | Depth | Loci |
|------|----|-------|------|
| USP25 | 1 | 682.48 | chr21:17135058-17250927 |
| RARA | 1 | 2787.03 | chr17:38497482-38499177 |
| PPP2R3B, CRLF2, FANCB, ZRSR2, NHS, DMD, BCOR, DDX3X, KDM6A, RBM10, ARAF, GATA1, KDM5C | 1 | 620.9 | chrX:320389-53222482 |
| PHOX2B | 1 | 399.86 | chr4:27162287-41748357 |
| MIR548AN, FGF14, BIVM-ERCC5, ERCC5, IRS2, GAS6-AS1, GAS6, GAS6 | 1 | 634.49 | chr13:100041547-114566778 |
| FUBP1 | 1 | 704.69 | chr1:78414264-78435748 |
| DMRT1, JAK2 | 1 | 602.78 | chr9:845380-5126957 |
| CDKN2A, CDKN2B-AS1, CDKN2B, CDKN2B-AS1 | 1 | 347.12 | chr9:21968129-22062228 |
| CDA, ARID1A | 1 | 155.91 | chr1:20945002-27023944 |
| CD274, PDCD1LG2, PTPRD, MLLT3, MTAP, CDKN2A | 1 | 786.59 | chr9:5455949-21968129 |
| CCND3, VEGFA | 1 | 172.01 | chr6:41909068-43739196 |
| BTK, PAK3, STAG2, BCORL1, GPC3, PHF6, ZIC3, G6PD, G6PD, IKBKG | 1 | 693.44 | chrX:100604718-153775209 |
| ATRX | 1 | 705.82 | chrX:76763817-77041667 |
| TCF7L2 | 0 | 139.28 | chr10:114710993-114711498 |

## 1.4   Fusion gene discovery

Fusion events are detected using the software DELLY2[5]. From the genome alignments, DELLY discovers fusion events (translocations and inversions) by integrating insert distances determined by the paired-end reads and split-read alignments to accurately detect genomic rearrangements at single nucleotide resolution. Fusion events are tagged as "Known fusions" if they match the entry in ChimerDB[6] (collection of known fusion events). Known fusion events are reported in the following sections (if any).

Table 10: Summary of fusion events detected in each sample.

| Sample | Known events | Unknown events |
|--------|-------------:|---------------:|
| HD753_3 | 1 | 2 |
| HD753_4 | 1 | 2 |

## 1.4.1  HD753_3 Results



Figure 3: Circos plot displaying fusion events in relation to chromosome location for sample HD753_3. Fusion events observed on the same chromosome are drawn in red whereas fusion events that are on different chromosomes are drawn in blue. Gene annotations are drawn at the tip of the arcs.

Table 11: Fusion events detected in sample HD753_3. Associated disease and source of annotation are mentioned in Disease and Source column, respectively.

| Fusion genes | Fusion location | Supporting fusion reads | Supporting paired reads | Disease | Source |
|---|---|---|---|---|---|
| RET–CCDC6 | chr10:43609952-chr10:61638611 | 33 | 36 | adenocarcinoma | Mitel-man,OMIM,GenBank |

## 1.4.2  HD753_4 Results



Figure 4: Circos plot displaying fusion events in relation to chromosome location for sample HD753_4. Fusion events observed on the same chromosome are drawn in red whereas fusion events that are on different chromosomes are drawn in blue. Gene annotations are drawn at the tip of the arcs.

Table 12: Fusion events detected in sample HD753_4. Associated disease and source of annotation are mentioned in Disease and Source column, respectively.

| Fusion genes | Fusion location | Supporting fusion reads | Supporting paired reads | Disease | Source |
|---|---|---|---|---|---|
| RET–CCDC6 | chr10:43609952-chr10:61638611 | 22 | 21 | adenocarcinoma | Mitel-man,OMIM,GenBank |

## 2 Quality Metrics

### 2.1 Sequence Quality Metrics

The base quality of each sequence read is inspected. Low quality calls are removed before proceeding with further processing. Using a sliding window approach, bases with low quality are removed from the 3' and 5' ends. Bases are removed if the average phred quality is below 15. Finally only mate pairs (forward and reverse read) were used for the next analysis step. The total amount of raw sequence data and the results of the quality filtering is collected and reported in the following table.

Table 13: Sequence quality metrics per sample

| Sample | Total Reads | LQ Reads | Single Reads | HQ Reads |
|--------|-------------|----------|--------------|----------|
| HD753_3 | 83,483,850 | 1,900,592 (2.3%) | 1,791,324 (2.1%) | 79,791,934 (95.6%) |
| HD753_4 | 77,279,854 | 1,821,082 (2.4%) | 1,720,580 (2.2%) | 73,738,192 (95.4%) |

Total Reads: Total number of sequence reads analysed for each sample.
LQ Reads: Number of low quality reads.
Single Reads: Number of high quality reads without mates (2nd read).
HQ Reads: Number of high quality reads used for further analysis.

### 2.2 Alignment Metrics

Mapping to the reference sequence/database is done using BWA[7] with default parameters. The following table contains the number of reads mapped to the reference for each sample. Please note that the mapping efficiency depends on the accuracy of the reference and the quality of sequence reads.

Table 14: Mapped read metrics observed per sample

| Sample Name | HQ Reads | Mapped Reads |
|-------------|----------|--------------|
| HD753_3 | 79,791,934 | 79,692,755 (99.88%) |
| HD753_4 | 73,738,192 | 73,635,493 (99.86%) |

### 2.3 Alignment Classification

The alignment classification table includes the following read categories:

- Mapped: Reads mapped to reference.

- Unique: Reads mapped to exactly one site on the reference.

- Non-unique: Reads mapped to more than one site on the reference.

- Singletons: Mapped reads without mates (read not paired).

- Cross-Contig: Read pairs with the mate mapped to a different contig.

- On target: Reads mapped to target +/- 100 bp extension.

Percentage of reads in categories **Non-unique, Unique, Singletons, Cross-Contig** are calculated based on the number of reads mapping to entire reference.

Percentage of reads in category **On target** is calculated based on the number of reads mapped uniquely (excluding **Singletons** and **Cross-Contig** - if any).

Table 15: Read metrics for HD753_3, HD753_4.

| Read category | HD753_3 | HD753_4 |
| --- | --- | --- |
| Mapped | 79,692,755 | 73,635,493 |
| Unique | 78,148,831 (98.06%) | 72,185,992 (98.03%) |
| Non-unique | 1,543,924 (1.94%) | 1,449,501 (1.97%) |
| Singletons | 16,467 (0.02%) | 15,256 (0.02%) |
| Cross-Contig | 185,612 (0.23%) | 167,679 (0.23%) |
| On target | 60,223,567 (77.26%) | 53,068,587 (73.70%) |

Reads in categorie(s) **Non-unique** , **Singletons** and **Cross-Contig** are excluded from analysis.

## 2.4   Alignment Refinement Metrics

The removal of PCR duplicates is done using Picard[8] in order to remove the artificial coverage brought on by the PCR amplification step during the library preparation. If a read maps to the same genomic location and has same orientation as the read already mapped it is considered as duplicated. For paired-end, both reads should fulfill the criteria in order to designate as PCR duplicate. One copy of the duplicate read pair is kept in the alignment.

Local realignment serves to transform regions with misalignments due to indels into clean reads containing a consensus indel suitable for standard variant discovery approaches. GATK is used for this purpose.

The goal of Base Quality Recalibration is to improve the base quality score of reads for downstream processing and also correct for error covariates like machine cycle and dinucleotide context. A base quality score represents the probability of a particular base mismatching the reference genome. After recalibration quality scores are more accurate in that they are closer to the true probability of mismatch. This process is achieved by analyzing the covariation among several different features of a base. The reported quality score, sequencing cycle, and sequencing context are considered for this step. GATK modules are used for achieving this.

The following table contains the number of high-quality reads after read mapping, alignment and refinement.

Table 16: HQAligned reads per sample

| Sample Name | Input Reads | Duplicate Reads | HQ Reads |
| --- | --- | --- | --- |
| HD753_3 | 60,223,567 | 25,157,050 (41.77%) | 35,066,517 (58.23%) |
| HD753_4 | 53,068,587 | 20,017,660 (37.72%) | 33,050,927 (62.28%) |

## 2.5   Coverage Report

The coverage plot showing the base coverage distribution from the HQ aligned data. Depth of coverage is plotted on X-axis and the percentage of the respective reference covered is plotted on Y-axis. The coverage plot is restricted to the target region without extension. The shape of the curve defines the uniformity of the reference coverage in the samples analysed. Samples with high uniformity usually have >90% covered at 0.2x average coverage (e.g. 100x for 500x average coverage)
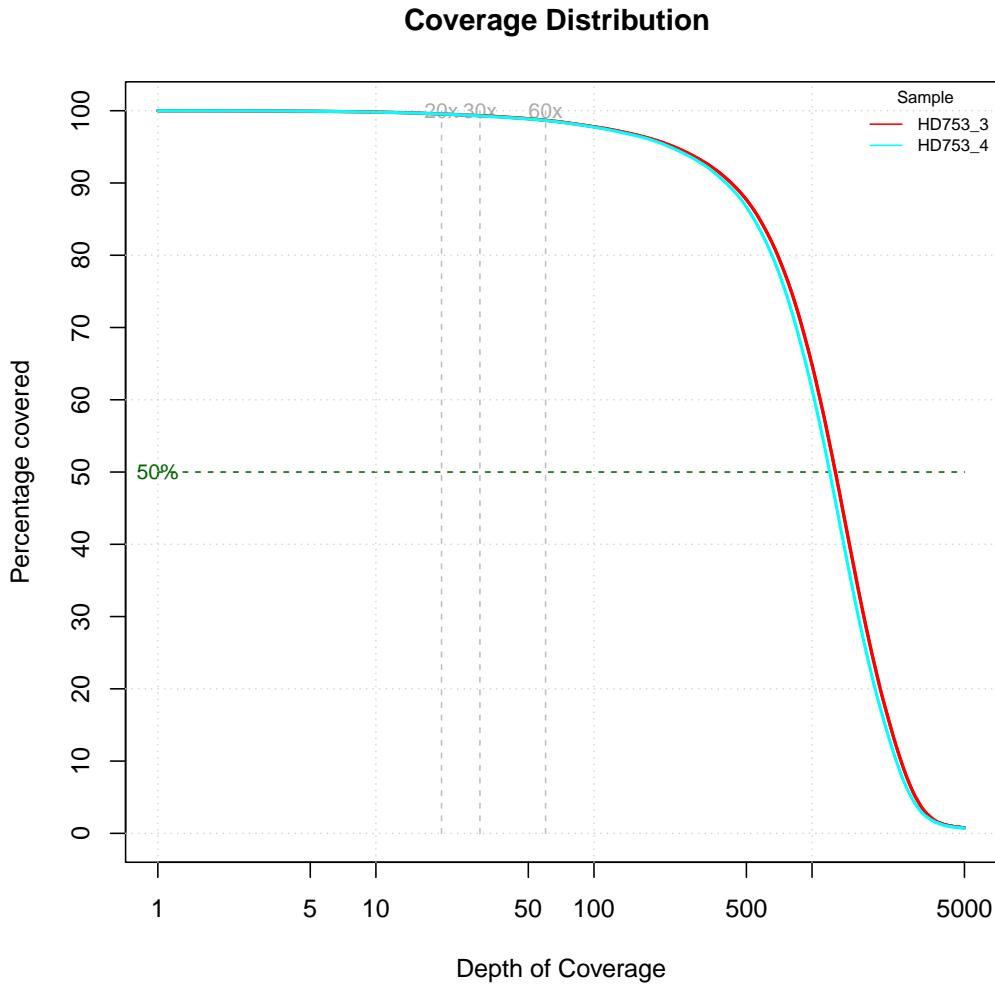
### Coverage Distribution



Figure 5: Coverage plot (excluding duplicated fragments).

Table 17: Depth of coverage summary (excluding duplicated fragments).

| | target coverage | | % of target covered with at least | | | | |
|---|---|---|---|---|---|---|---|
| sample | total bases | average (x) | 2x | 50x | 100x | 300x | 500x |
| HD753_3 | 4.14 GB | 1425.33 | 100 | 98.9 | 97.8 | 93.4 | 87.7 |
| HD753_4 | 3.92 GB | 1350.25 | 100 | 98.9 | 97.7 | 93.0 | 86.7 |

## Coverage Distribution



Figure 6: Coverage plot (including duplicated fragments).

Table 18: Depth of coverage summary (including duplicated fragments).

| | target coverage | | % of target covered with at least | | | | |
|---|---|---|---|---|---|---|---|
| sample | total bases | average (x) | 2x | 50x | 100x | 300x | 500x |
| HD753_3 | 7.23 GB | 2486.13 | 100 | 99.1 | 98.3 | 95.3 | 91.9 |
| HD753_4 | 6.40 GB | 2200.72 | 100 | 99.0 | 98.2 | 94.8 | 91.0 |

## 2.6   Library Report

Fragment insert size histogram of the paired-end library observed from all the samples analysed. The insert size is determined by mapping individual read pairs on the reference sequence. The distance between 5'prime ends of both sequenced reads in a pair that are mapped to the reference is the observed length of the sequenced fragment. By performing this operation for all mapped reads the distribution can be generated. X-axis shows the insert size in bp and Y-axis shows the number of fragments with the observed fragment insert sizes.



Figure 7: HD753_3 .

Table 19: Sample wise insert size metrics for HQ aligned reads. The mean insert size (Mean) and its standard deviation (Stddev) is given in base pairs.

| Sample | Pair orientation | Mean | Stddev | # Read pairs |
|--------|-----------------|------|--------|--------------|
| HD753_3 | FR | 232 | 81 | 17,495,031 |
| HD753_4 | FR | 230 | 80 | 16,489,444 |

Figure 8: HD753_4 .

# 3 Deliverables

Table 20: List of delivered files, format and recommended programs to access the data.

| File | Format | Program To Open File |
| --- | --- | --- |
| PROJECT_supplementary_tables.tar.gz | GZ | Unzip tool |
| SAMPLE.CNV_deletion.tsv | TSV | Spreadsheet Editor |
| SAMPLE.CNV_duplication.tsv | TSV | Spreadsheet Editor |
| SAMPLE.fusion_events.tsv | TSV | Spreadsheet Editor |
| SAMPLE.hg19.HQ.alignment.bam | BAM | IGV, Tablet |
| SAMPLE.hg19.HQ.alignment.bam.bai | BAI | None |
| SAMPLE.hg19.alignment.bam | BAM | IGV, Tablet |
| SAMPLE.hg19.alignment.bam.bai | BAI | None |
| SAMPLE.indels.tsv | TSV | Spreadsheet Editor |
| SAMPLE.indels.vcf | VCF | Text Editor |
| SAMPLE.snps.tsv | TSV | Spreadsheet Editor |
| SAMPLE.snps.vcf | VCF | Text Editor |

SAMPLE.hg19.alignment.bam was used for Fusion Gene discovery (see chapter 1.4)

SAMPLE.hg19.HQ.alignment.bam was used for Variant discovery (see chapter 1.1) and for Copy number analysis (see chapter 1.3)

PROJECT_supplementary_tables.tar.gz contains the variant calls (SNVs and InDels) that were observed in the sample(s) but filtered out due to QC checks.

# 4 Formats

Table 21: References and descriptions of file format.

| Format | Description |
| --- | --- |
| TSV | Tab separated table style text file. This can be imported into spreadsheet processing software like MS OFFICE Excel. |
| FASTQ[9] | Text-based format for storing both a biological sequence (usually nucleotide sequence) and its corresponding quality scores. Both the sequence letter and quality score are encoded with a single ASCII character for brevity. |
| BAM[10] | Compressed binary version of the Sequence Alignment/Mapping (SAM) format, a compact and index-able representation of nucleotide sequence alignments. |
| VCF[11] | Variant Call Format (VCF) is a format to describe and report the variants. |

# 5 FAQ

Q: How can I open a TSV file in Excel?
A: Start Excel and click File -> Open and select the TSV file you want to open. Next an assistant dialog should show up. Make sure that you select tab as separator. Set the format of all rows without numbers to text. The TSV files use the dot as decimal separator and comma as thousands separator. Make sure that you set both correctly.

# 6 Bibliography

[1] Andreas Wilm, Pauline Poh Kim P. Aw, Denis Bertrand, Grace Hui Ting H. Yeo, Swee Hoe H. Ong, Chang Hua H. Wong, Chiea Chuen C. Khor, Rosemary Petric, Martin Lloyd L. Hibberd, and Niranjan Nagarajan. LoFreq: a sequence-quality aware, ultra-sensitive variant caller for uncovering cell-population heterogeneity from high-throughput sequencing datasets. *Nucleic acids research*, 40(22):11189–11201, December 2012.

[2] Simon A. Forbes, David Beare, Prasad Gunasekaran, Kenric Leung, Nidhi Bindal, Harry Boutselakis, Minjie Ding, Sally Bamford, Charlotte Cole, Sari Ward, Chai Y. Kok, Mingming Jia, Tisham De, Jon W. Teague, Michael R. Stratton, Ultan McDermott, and Peter J. Campbell. COSMIC: exploring the world's knowledge of somatic mutations in human cancer. *Nucleic Acids Research*, 43(D1):gku1075–D811, October 2014.

[3] Melissa J. Landrum, Jennifer M. Lee, Mark Benson, Garth Brown, Chen Chao, Shanmuga Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Jeffrey Hoover, Wonhee Jang, Kenneth Katz, Michael Ovetsky, George Riley, Amanjeev Sethi, Ray Tully, Ricardo Villamarin-Salomon, Wendy Rubinstein, and Donna R. Maglott. ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic Acids Research*, 44(D1):D862–D868, January 2016.

[4] Eric Talevich, A. Hunter Shain, Thomas Botton, and Boris C. Bastian. CNVkit: Genome-Wide Copy Number Detection and Visualization from Targeted DNA Sequencing. *PLoS Comput Biol*, 12(4):e1004873+, April 2016.

[5] Tobias Rausch, Thomas Zichner, Andreas Schlattl, Adrian M. Stütz, Vladimir Benes, and Jan O. Korbel. DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics*, 28(18):i333–i339, September 2012.

[6] Pora Kim, Suhyeon Yoon, Namshin Kim, Sanghyun Lee, Minjeong Ko, Haeseung Lee, Hyunjung Kang, Jaesang Kim, and Sanghyuk Lee. ChimerDB 2.0 - a knowledgebase for fusion genes updated. *Nucleic acids research*, 38(suppl 1):D81–D85, 2010.

[7] Heng Li and Richard Durbin. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics (Oxford, England)*, 25(14):1754–1760, July 2009.

[8] Picard. http://picard.sourceforge.net.

[9] Peter J. A. Cock, Christopher J. Fields, Naohisa Goto, Michael L. Heuer, and Peter M. Rice. The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Research*, 38(6):1767–1771, 2010.

[10] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16):2078–2079, 2009.

[11] Petr Danecek, Adam Auton, Goncalo Abecasis, Cornelis A. Albers, Eric Banks, Mark A. DePristo, Robert E. Handsaker, Gerton Lunter, Gabor T. Marth, Stephen T. Sherry, Gilean McVean, Richard Durbin, and 1000 Genomes Project Analysis Group. The variant call format and vcftools. *Bioinformatics*, 27(15):2156–2158, 2011.

[12] Derek Barnett, Erik Garrison, Aaron Quinlan, Michael Strömberg, and Gabor Marth. BamTools: a C++ API and toolkit for analyzing and managing BAM files. *Bioinformatics*, 27(12):btr174–1692, April 2011.

[13] Mary Kate Wing. "bamUtil is a repository that contains several programs that perform operations on SAM/BAM files.". http://genome.sph.umich.edu/wiki/BamUtil, 2015.

[14] Aaron R. Quinlan and Ira M. Hall. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, 26(6):841–842, March 2010.

[15] Aaron McKenna, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis, Andrew Kernytsky, Kiran Garimella, David Altshuler, Stacey Gabriel, Mark Daly, and Mark A. DePristo. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research*, 20(9):1297–1303, 2010.

[16] Mark A DePristo, Eric Banks, Ryan Poplin, Kiran V Garimella, Jared R Maguire, Christopher Hartl, Anthony A Philippakis, Guillermo del Angel, Manuel A Rivas, Matt Hanna, Aaron McKenna, Tim J Fennell, Andrew M Kernytsky, Andrey Y Sivachenko, Kristian Cibulskis, Stacey B Gabriel, David Altshuler, and Mark J Daly. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet*, 43:491–498, 2011.

[17] R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2010. ISBN 3-900051-07-0.

[18] Artem Tarasov, Albert J. Vilella, Edwin Cuppen, Isaac J. Nijman, and Pjotr Prins. Sambamba: fast processing of NGS alignment formats. *Bioinformatics*, February 2015.

[19] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16):2078–2079, 2009.

[20] Pablo Cingolani. "snpEff: Variant effect prediction". http://snpeff.sourceforge.net, 2012.

[21] Marc Lohse, Anthony M. Bolger, Axel Nagel, Alisdair R. Fernie, John E. Lunn, Mark Stitt, and Björn Usadel. RobiNA: a user-friendly, integrated software solution for RNA-Seq-based transcriptomics. *Nucleic Acids Research*, 40(W1):W622–W627, July 2012.

# A   Analysis Workflow

The schematic diagram of the data analysis steps that have been performed is shown in figure 1.
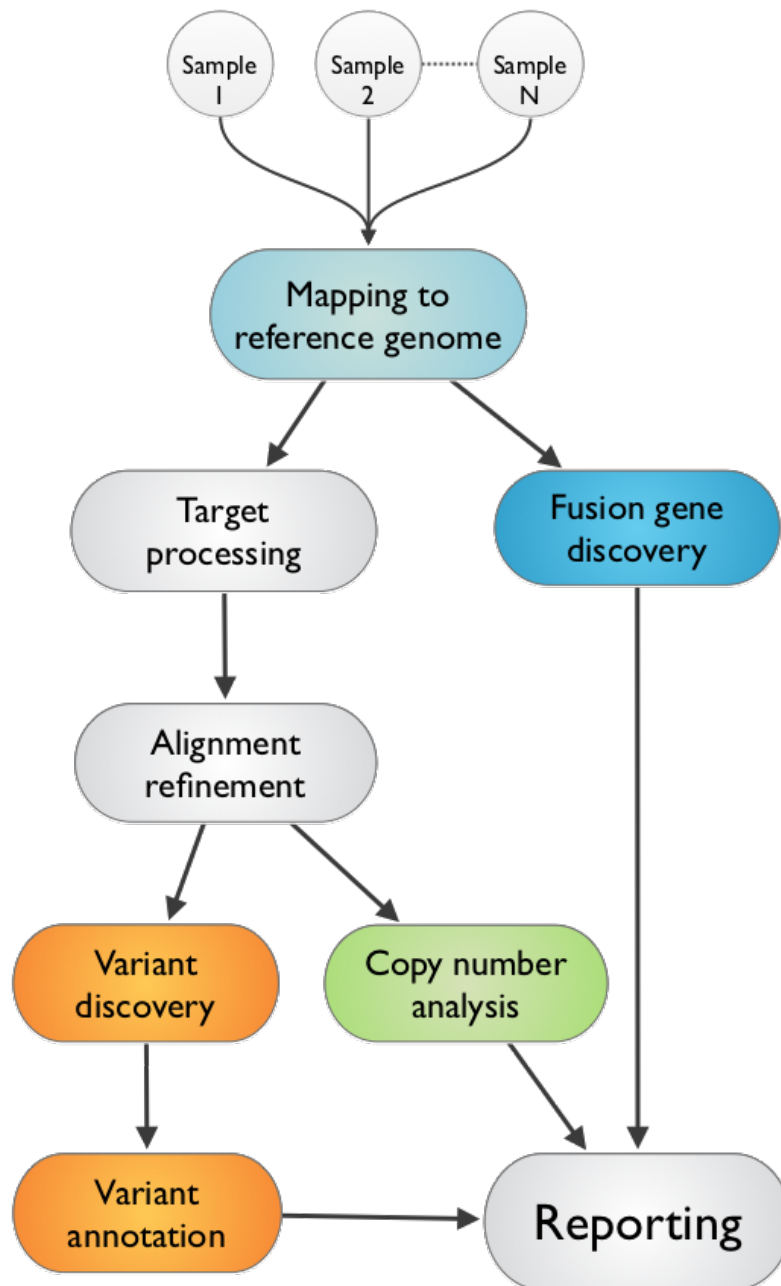


Figure 9: ONCOPANEL ALL-IN-ONE v1.4 Workflow

# B   Sequence Data Used

Table 22: Analysed samples (SE = single end, PE = paired end).

| Sample  | Read Type | File Name |
|---------|-----------|-----------|
| HD753_3 | PE | FE-0242_HD753_3_lib185570_5345_1_1.fastq |
|         |    | FE-0242_HD753_3_lib185570_5345_1_2.fastq |
| HD753_4 | PE | FE-0242_HD753_4_lib185571_5345_1_1.fastq |
|         |    | FE-0242_HD753_4_lib185571_5345_1_2.fastq |

# C   Reference Database

Table 23: Information about the Homo sapiens Reference Database.

| Tag | Description |
| --- | --- |
| Name | Homo sapiens |
| Version | hg19.chronly |
| Source | UCSC |
| Size (bp) | 3.095 GB |
| Sequences | 23 |

Table 24: Information about additional reference data used.

| Type | Version | Source |
| --- | --- | --- |
| Annotation | 19 | GENCODE |
| dbSNP | 150 | NCBI |
| ClinVar[3] | 02.10.17 | NCBI |
| COSMIC[2] | 71 | Sanger Institute |
| ChimerDB[6] | 2.0 | ERCSB |

Table 25: Information about the target region used.

| Tag | Description |
| --- | --- |
| Name | GATC All in One |
| Size (bp) | 2,908,369 |
| Source | GATC Biotech AG |

# D   Relevant Programs

Table 26: Name, version and description of relevant programs.

| Program | Version | Description |
|---|---|---|
| bamtools[12] | 2.3.0 | BamTools provides a small, but powerful suite of command-line utility programs for manipulating and querying BAM files for data. |
| BamUtil[13] | 1.0.10 | BamUtil is a repository that contains several programs that perform operations on SAM/BAM files |
| bedtools[14] | 2.26.0 | Bedtools allows one to intersect, merge, count, complement, and shuffle genomic intervals from multiple files in widely-usedgenomic file formats such as BAM, BED, GFF/GTF, VCF |
| BWA[7] | 0.7.15 | BWA is a software package for mapping low-divergent sequences against a large reference genome |
| CNVkit[4] | 0.9.1.dev0 | CNVkit is a Python library and command-line software toolkit to infer and visualize copy number from targeted DNA sequencing data |
| Delly2[5] | 0.7.6 | DELLY2: Structural variant discovery by integrated paired-end and split-read analysis |
| GATK[15, 16] | 3.7 | GATK is a java-based command-line toolkit that process SAM / BAM / VCF files. |
| LoFreq[1] | 2.1.2 | Lofreq is a fast and sensitive variant caller for inferring SNVs and indels from next-generation sequencing data. |
| Picard[8] | 1.131 | Picard is a java-based command-line utilities for processing SAM / BAM files. |
| R[17] | 3.2.4 | R is a programming language and environment for statistical computing. |
| sambamba[18] | 0.6.6 | Sambamba is a high performance modern robust and fast tool (and library), for working with SAM and BAM files. |
| SAMTools[19] | 0.1.18 | SAMtools provide various utilities for manipulating alignments in the SAM format. |
| snpEff[20] | 4.3 | SnpEff is a genetic variant annotation and effect prediction toolbox. |
| SnpSift[20] | 4.3 | SnpSift helps filtering and manipulating genomic annotated files . |
| Trimmomatic[21] | 0.33 | Trimmomatic performs a variety of useful trimming tasks for Illumina paired-end and single-end data. |

# E Tables

Table 27: Definition of fields of the tab delimited variant report (Sample.indels.tsv and Sample.snps.tsv).

| Name | Meaning |
| --- | --- |
| Ref ID | Name of chromosome or reference contig where the variant occurs. |
| Position | Position of reference contig or chromosome where the variant occurs. |
| Reference Base (s) | The reference base at the variant site. |
| Modified Base (s) | Alternative (observed) base in the samples in general [ VARIANT ]. |
| Mutation Frequency (%) | The mutation frequency with which a particular mutation occurs in a population. |
| Coverage Depth (x) | The total depth of the reads that passed the internal quality control metrics from all reads present at this site. |
| dbID | Known variant indentifier. |
| FILTER | Variants passing the filters will be tagged as "PASS" and the variants failing the filters will be tagged by the respective filter names. |
| AF | Allele (Mutation) frequency. |
| DP | Counts for ref-forward bases, ref-reverse, alt-forward and alt-reverse bases. |
| CLNDSDBID | Variant disease database ID. |
| CLNSIG | Variant Clinical Significance, 0 - unknown, 1 - untested, 2 - non-pathogenic, 3 - probable-non-pathogenic, 4 - probable-pathogenic, 5 - pathogenic, 6 - drug-response, 7 - histocompatibility, 255 - other. |

Table 28: Definition of genomic annotations as produced by snpEff (Sample.indels.tsv and Sample.snps.tsv).

| Name | Meaning |
| --- | --- |
| EFFECT | Variant's effect on protein. |
| IMPACT | Predicted impact from variant's protein effect. |
| HGVS_C | Variant's codon change (DNA level). |
| HGVS_P | Variant's codon change (Protein level). |
| GENE | The gene entry associated with the location of the variant call. |
| BIOTYPE | Variant's coding status. |
| TRID | Associated transcript IDs. |
| CDS_POS | Variant's codon change position. |
| AA_POS | Variant's amino acid position. |

Eurofins Genomics' products, services and applications reach the best quality and safety levels. They are carried out under strict QM and QA systems and comply with the following standards:

| | | | |
|---|---|---|---|
| ISO 9001 | Globally recognised as the standard quality management certification | GLP | The gold standard to conduct non-clinical safety studies |
| ISO 17025 | Accredited analytical excellence | GCP | Pharmacogenomic services for clinical studies |
| ISO 13485 | Oligonucleotides according to medical devices standard | cGMP | Products and testing according to pharma and biotech requirements |