# Detection of mosaic and population-level structural variants with Sniffles2

In the format provided by the authors and unedited

# Supplementary Section 1: GIAB and HGSC benchmarking

Across the default coverage (30x for HiFi, 50x for ONT), Sniffles2 shows the best performance with respect to correctly identified and genotyped insertions (HiFi: GT F-score 0.909, ONT: GT F-score 0.915) and deletions (HiFi: GT F-score 0.934, ONT: GT F-score 0.944) (see **Supplementary Table 2** for details). Sniffles2 achieves a better result in a fraction of the time across data sets compared to Sniffles (v1.12), being over 17 times (HiFi) and 11 times (ONT) faster in processing a 30x coverage data set, respectively. **Figure 2A+B** shows the results for PacBio HiFi and ONT, respectively. In addition, Sniffles2 is also the fastest method overall, requiring 33.33 CPU minutes for processing a 30x coverage HiFi dataset , which was twice as fast as SVIM, the 2nd fastest method. For a 30x coverage ONT dataset, Sniffles2 was also close to twice as fast (1.92x) as the second fastest caller (SVIM), while also having an over 8.4% higher GT F-score. Taking into account Sniffles2 multi-processing capability (not supported by SVIM), the speedup increases even further, to more than 5.4-fold and 7.5-fold for HiFi 30x, ONT 30x data sets, respectively. When reducing the coverage from 30x to 10x we observe only a slight reduction in F-score for Sniffles2 (HiFi: reduction GT F-score 0.041, ONT: reduction F-score 0.058). This is in stark contrast to other programs such as cuteSV, where using default parameters, F-score dropped by an average of over 60% (HiFi: reduction F-score 0.56, ONT: reduction F-score 0.58). Even when using only 5x for Sniffles2, we still observe a high accuracy ONT (F-score: 0.74) and HiFi (F-score: 0.77). This is achieved as Sniffles2 includes an automated parameter selection for filtering of SV candidates based on the available coverage. In contrast, other SV callers rely on manual adjustment of these parameters to retrieve acceptable results across coverages and sequencing technologies. **Figure 2A-D** shows this clearly as all other SV callers show a decreased performance across lower coverage. Even

when tuning the parameters for other SV callers (**Figure 2A-D,** Tuned parameters), Sniffles2 remains the highest accuracy (see **Supplementary Table 1** for details).

Additionally, we tested Sniffles2 with an additional mapper (LRA) to showcase its versatility. When mapping HiFi data with LRA we observe a slight increase in performance when comparing the genotype F-scores (HiFi GT F-Score increase by 0.006), moreover when using LRA to map ONT data we observe a slight decrease in performance of Sniffles2 (ONT GT F-Score decrease 0.0136). Interestingly, wWhen testing Sniffles with CLR data the error rates of this outdated data impacts our variant calling and thus we suggest other SV callers (see **Supplementary Table 1**).

Next, we performed an evaluation with respect to the Tier2 GIAB dataset, which is a more challenging region of the GIAB benchmark set as it includes repeats and GIAB cannot guarantee the accuracy of the variants in these regions (see **Supplementary Table 3** for details). Again, Sniffles2 even increases the performance difference compared to other SV caller. Lastly, we benchmarked Sniffles2 across a more challenging SV data set across 386 medically relevant, but highly polymorphic/challenging genes[33]. GIAB has recently released this call set of ~200 SV covering around 70% of these genes[33]. **Figure 2E-H** shows the results. Again Sniffles2 outperforms the other SV callers in terms of accuracy and speed using default parameters. The next best performing SV caller (pbsv for HiFi, cuteSV for ONT) both achieved 3.6% and 5.5% lower genotyping accuracy (GT F1-Score) even at 30x coverage.

**Supplementary Table 4** contains the detailed results across all SV callers. Overall, Sniffles2 outperforms other state-of-the-art SV caller across the entire genome including the most challenging regions/genes.

Next, Sniffles2 improves insertion identification through two additional methods: First, the consensus module corrects sequencing-related errors in the recovered insertion sequences using a fast pseudo-alignment based approach. This allows Sniffles2 to attain the second highest mean sequence identity of (HiFi: 0.948, ONT: 0.939), after pbsv (HiFi: 0.953, ONT:

0.949), while Sniffles2 is over 14x (HiFi) and 36x (ONT) faster (see **Supplementary Figure 1, Supplementary Table 5** containing insertion sequence accuracy across all callers) at 30x coverage. Second, Sniffles2 increases the sensitivity for the detection of large insertions by recording additional supporting alignment signals in the affected regions (see **Supplementary Figure 2, Supplementary Table 6**) at much higher speed than pbsv, the only SV caller with a comparable accuracy for long insertions.

Lastly, GIAB only represents one individual benchmarked across most studies (HG002). Thus next, we used Dipcall[53] together with three T2T assemblies [54] (HG01243, HG02055, HG02080) to further assess the performance of Sniffles2. Clearly, we give Dipcall the benefit of the doubt, knowing that the accuracy will be lower than the GIAB vetted benchmark set. Overall, Sniffles2 performs the best across all samples having on average a F-Score of 0.79 at 30x coverage ONT and HiFi compared to 0.77 F-score for the next best SVcaller (cuteSV), at nearly 3.5 times the speed (74.03 CPU minutes versus 256.65 CPU minutes on average) for default parameters. **Supplementary Table 7** contains the detailed results for each benchmarked program across these three samples. Besides the here benchmarked insertions and deletions, we also benchmarked Sniffles2 on duplications, inversions and translocations using simulated data as no benchmark exists. Overall, Sniffles2 again outperformed all other methods in speed and accuracy (see **Supplementary Figure 3 and Supplementary Table 8**) (see methods for details).

# Supplementary Section 2: spike in experiments for mosaic SV benchmarking

To assess the performance of Sniffles2 across mosaic SVs, we first synthetically merged HG002 (at low concentrations: 5x (7%), 7x (10%), 10x (14%), 15x (21%) or 20x (28%)) with the 50-63x coverage read data from the sample of an unrelated individual (HG00733). This yielded multiple, synthetic samples with constant total coverage of ~70x (68x in one case), but varying concentrations of HG002 in them (7-28%). We only assessed this for ONT data as this technology often is sequenced at higher sequence depth than Pacbio HiFi. The latter can still be used, but is not benchmarked here. **Figure 4 A&B** highlights the results across this synthetic data set across different concentrations of HG002 (x-axis, **Supplementary Table 12**). **Figure 4A and 4B** show the precision and recall of SVs across the different concentrations. In blue, we highlight the performance of Sniffles2 default germline-mode, which shows low precision and increased recall with increase in concentration of HG002 in the mix. In red we show cuteSV which has lower precision and recall in both cases, when compared to sniffles2 germline-mode. In yellow we highlight the performance of Sniffles2 mosaic mode. Sniffles2 achieves high precision that increases with higher proportions of HG002 (**Figure 4A**) with the highest being at 21% HG002 proportion (84.12% precision). Recall also increases with higher proportions of HG002, moreover it only reaches a maximum of 54.05%. This is because Sniffles2 mosaic-mode only calls SV where the variant allele frequency (VAF) is between 5 and 20%.. We computed the adjusted recall by only using the SVs where the VAF matches in the range that mosaic-mode uses (**Figure 4B, Supplementary Table 12**). This adjusted recall is shown in green, which averages 94.47% across all comparisons (highest 96.39% at 7% mix of HG002). Thus following the default recall for germline variants.

# Supplementary Section 3: MSA analysis

We identified 3,049 Sniffles2 non-germline SVs that have read support unique to the cingulate cortex region from which 2,856 of these SV overlapped with 1,176 genes including at least 80 that are related to the brain and its development (**Supplementary Table 15**). Some examples are *GABRB3* which is implicated in many human neurodevelopmental disorders and syndromes, *NRXN3* which is involved in synaptic plasticity and Netrin receptor DCC which is involved in neuron migration and axonogenesis. *GRIN2A*, which encodes for NMDA receptor subunit 2a. Of these 29 SVs, 27 overlapped introns and 2 affected at least one exon. Furthermore, 4 SVs disturb regulatory regions associated with at least one gene, including *PEX26*, *DLL1* and *ABCA2*. This further highlights Sniffles2's ability to not only detect SVs that cannot be identified using Illumina data alone but also the likely unique presence of a subset of these calls within a brain region.

As highlighted above the role of repeat elements (Alu/L1) are important for the formation of germline and non-germline SV (**Supplementary Table 16**, **Figure 5C**). In summary, we found a similar fraction of germline and non-germline SVs was associated with repeat elements (60.61% and 59.70% respectively). When comparing specific repeat elements, we found Alu being more abundant in germline SV (4.82% difference) and simple repeats more abundant in non-germline/mosaic SVs (6.54% difference). The patterns of repeat association only shifted in duplication where LTR, L2 and other repeats differed from the norm and with inversions where no SV was detected in the mosaic SVs (**Figure 5C** and **Supplementary Figure 9**). For duplications, the fraction of non-germline SVs associated with LINE1 and simple repeats showed the highest difference over germline SVs. This highlights the different ways in which repeat elements are associated with somatic structural variants. As we could show above Alu insertion mediated mosaic deletions. **Extended Figures 3** and **4** shows two examples of read

alignments for a non-germline duplication SV that was solely called by Sniffles2, alongside with its relation to nearby and overlapping repeat elements. For deletions and insertions, we observe similar to slightly lower fractions of non-germline SVs associated with most repeat types, with the exception of simple repeats. For this repeat family, especially non-germline duplication and insertion SVs had a higher fraction associated. Repetitive elements may be associated with neurodegenerative disorders, through increased expression and / or de novo somatic genomic integration[41]. The observance of a higher fraction of non-germline insertion and deletion SVs being associated with simple repeat elements could suggest a further correlation for this on the level of an individual brain region. Overall, this also highlights the differences between repeat families in their effects on somatic SV generation[63].

## Supplement Section 4: Technology comparison for Brain MSA sample

From the total of 22,619 merged SV calls, only 214 are shared across all methods and 19,254 (85.12%%) were detected by Sniffles2 together with either Manta, OGM or both. Sniffles2 uniquely detected 10,720 (47.39%) SVs, while Manta uniquely detected 2,599 (11.49%) and OGM 749 (3.31%). We next identified the most abundant SV type per method for the unique calls. We found that Sniffles2 has a higher number of insertions and deletions (94.87%, 70.98% respectively%), while for Manta are duplications and inversion (98.35%, 92.12% respectively, **Supplementary table 17.2**,). At the same time, Illumina showed a higher overlap with Sniffles2 (8,022 SVs) compared to OGM (298 SVs), likely also because of the minimum length for OGM. The overlap between Sniffles2 and the other technologies was also much higher for deletions (4,801) than for insertions (3,607). These differences are likely explained by the individual difficulty in detecting larger (i.e., Illumina) and smaller (i.e., OGM, recommended threshold approximately 5kbp or larger) insertion events, respectively. Next, we took a closer look at putatively mosaic SVs detected by Sniffles2. For this, we separated Sniffles2 detected SVs by their reported variant allele frequencies (VAF) into germline (VAF > 0.3) and likely mosaic (VAF

≤ 0.2) calls. A total of 58 insertions were shared uniquely by the germline and mosaic calls of Sniffles2. No other SV types were observed overlapping between mosaic and germline. Interestingly wWe do observe size differences for these 58 SV but share the same start location.. Next, we compared the mosaic Illumina and OGM calls. Only 702 mosaic SVs reported by Sniffles2 could be detected by either Illumina or OGM or both (24 OGM, 666 illumina and 12 both), highlighting the difficulty in identifying rare SVs. For deletions, only 473 of SVs were also found by either OGMand/or Illumina. Only 218 mosaic insertions reported by Sniffles2 were detected by the other methods. For duplications and inversions, only the Illumina data showed overlap (one and ten, respectively) with Sniffles2 mosaic SVs.

We further noted a shift in the allele frequencies across the Manta calls compared to the Sniffles2 calls **(Figure 5B)**. As expected, for Sniffles2, we observe a multimodal distribution with three peaks, representing homozygous, heterozygous and non-germline SVs, respectively. In contrast, Manta shows two main peaks in their allele frequency distribution. A homozygous (~0.9-1 AF) and a broad peak around 0.3 AF, which would be below the typical expected heterozygous AF peak. For Sniffles2, we furthermore observe a ~140% increase in the area under the curve in the putative non-germline range of allele frequencies (0.0-0.2) when compared to the illumina data, thus showing the potential of Sniffles2 to detect low frequency SV.

# Supplementary Section 5: Identification of cancer-specific somatic SVs by Sniffles2

We have shown that Sniffles2 can accurately identify SVs across the whole range of allele frequencies. Furthermore, we showed that Sniffles2 enables fast and accurate population SV calling. Next, we made use of the population merge strategy of Sniffles2 to investigate its ability to identify tumor-specific somatic SVs using paired tumor/normal samples. For this endeavor we used the well-known and highly studied COLO829 cancer cell line (with blood control COLO829BL). By merging the SV calls from both tumor and normal we are able to identify tumor-specific somatic SVs based on the genotype, alternative read count and and the support vector. **Figure 5D** shows shared and unique SV between the COLO829/COLO829BL cell-lines for two long-read technologies: ONT MinION and Pacbio Revio.

First, we tested both germline and mosaic SV calls against a COLO829 benchmark dataset[66]. This benchmark was done by analyzing the COLO928 sample using multiple technologies. We filtered out all SV smaller than 50 bases, and divided the analysis in two parts: translocations (13 BNDs in the benchmark) and the rest of the SVs (49 SVs in the benchmark). For the case of the translocations (BND) Sniffles2 reported five BND, from which four are present in the benchmark with all within five bases of either breakpoint. Nine BNDs from the benchmark were not detected by Sniffles2, and one BND was called but not present in the benchmark (**Supplementary Table 18**). For the rest of the 49 SVs (INS, DEL, INV, DUP), Sniffles2 was able to identify 36 in either the ONT, Revio datasets or both with breakpoints within 300 bases (max. allowed distance between breakpoints = 1kb, **Supplementary Figure 10**), moreover in some cases the technologies didn't agree on the genotype (het vs. homozygous alt.). This highlights potential changes across the cell lines as the ONT and Pacbio data were released a few years apart. From the remaining unidentified SV, four had no reads for the alternative allele (0% AF), three were identified by Sniffles as mosaic (8.2-21% AF), four had heterozygous genotype

(31-41% AF) and two had homozygous alternative genotype (81-100% AF, **Supplementary table 18, Supplementary figure 12).** For Sniffle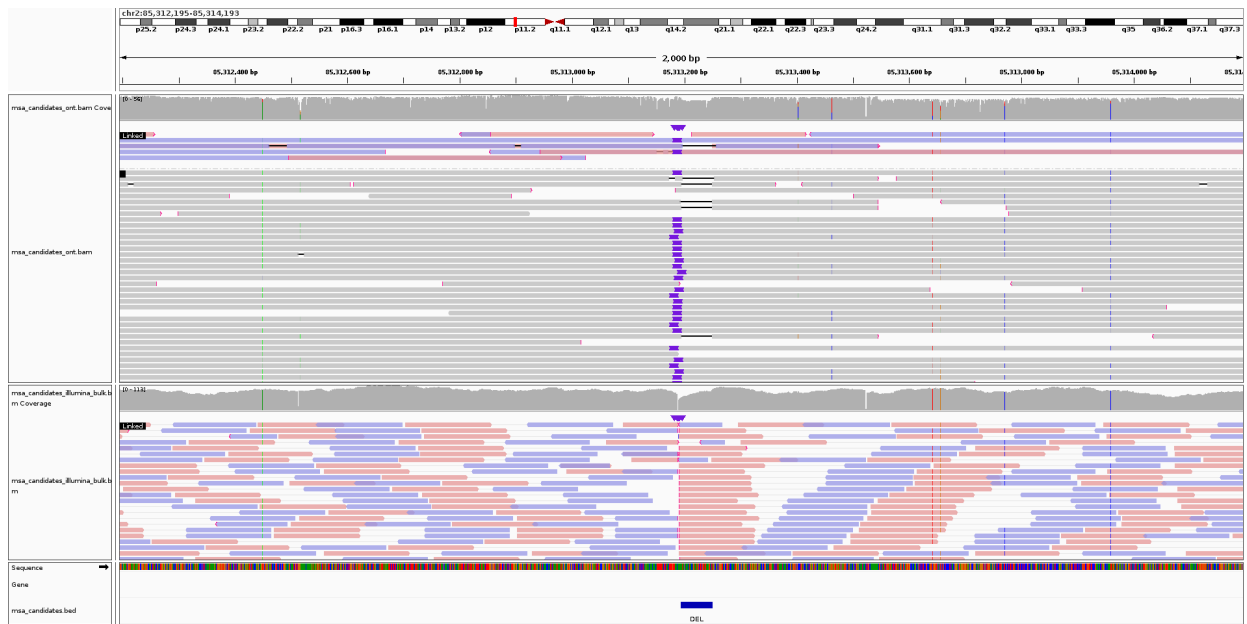s2 it shows that we missed only six germline and three mosaic SVs from the benchmark (81.63% recall, **Supplementary table 18**). These results show the difficulty of defining a benchmark dataset for an ever-evolving cancer cell-line such as COLO829 or any other.

We further annotated 79 tumor specific somatic SVs identified by Sniffles (see Methods). These SVs overlap with cancer related genes such as *PTEN, PMS2, ARHGEF5, PAK2, WWOX* to mention some. Additionally, some SVs overlap with olfactory receptor pseudogenes. **Figure 5E** and **Supplementary Figure 11** shows examples of cancer-specific germline SV. Manual inspection in IGV identified homozygous, heterozygous and LoH events.

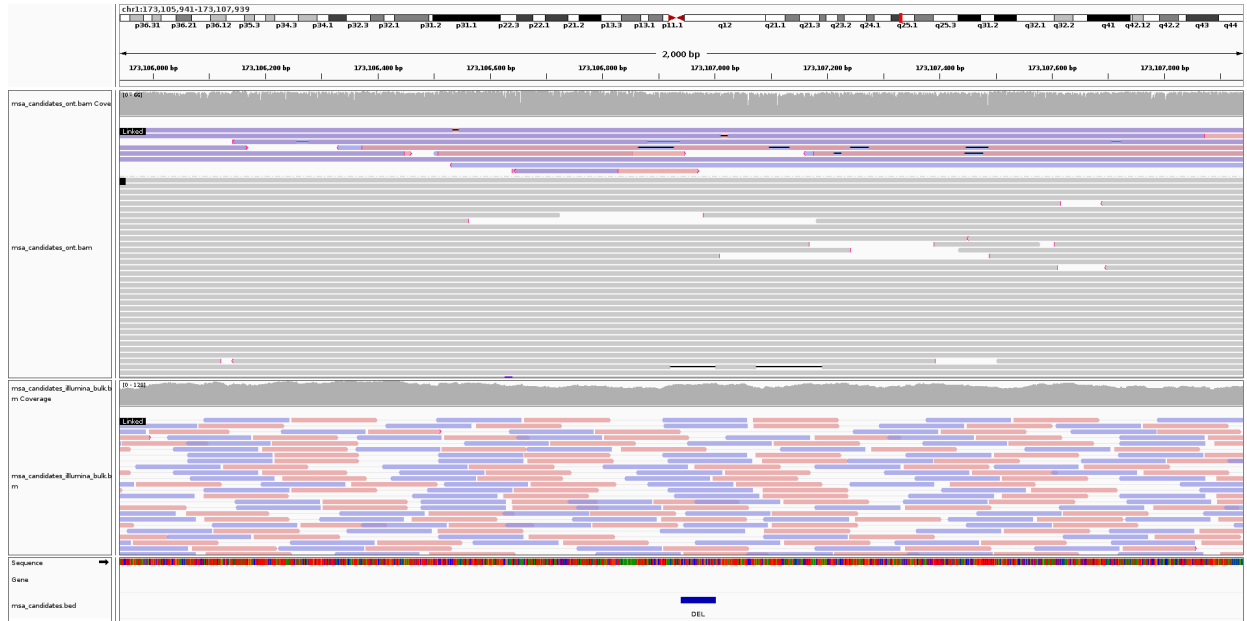# Supplementary figures

A)



B)

**Supplementary Figure 1**

Accuracy of recovered insertion sequences. All callers were benchmarked using default parameters as described in the methods section on benchmarking methodology. For ONT the number of insertions (INS) analyzed is 3,048 for every tool, which is the intersection of all TP from all callers. For the case of PacBio HiFi, the number of INs is 3440 for every tool, following the same strategy as in ONT.

**Supplementary Figure 2**

Genotype F1-measure for the detection of insertions by SV length. Measured by using the Truvari summary output files for the evaluation of each caller on the GIAB Pacific Biosciences HiFi and Oxford Nanopore 30x coverage benchmark. All callers were benchmarked using

default parameters. For those cases that no data is shown means that the result was not defined (NA). **Supplementary Table 6** also shows the number of SVs tested by SV length (Insertions only)

Benchmark for Duplications, Inversions and Translocation-type SVs (ONT 30x, SURVIVOR Simulation)

**Supplementary Figure 3**

Benchmark for duplications, inversions and translocations (simulation generated and evaluated using SURVIVOR). Precision, recall and F-measure were calculated from the SURVIVOR eval command output. All callers were benchmarked using the default parameters described in the methods section on benchmarking methodology.

**Supplementary Figure 4**

Performance of Sniffles2 population / multi-sample merge as a function of the number of samples combined. Performance is near-linear for the shown range of 3 to 777 samples. As input data, multiples of the HG002, HG003 and HG004 trio were merged (representing 1 to 259 families).

**Supplementary Figure 5**

Mendelian consistency, inconsistency and missing genotypes (NA) from a fully genotyped population VCF of 31 genomes. We show the information of 7 complete families (father, mother, proband, 21 genomes). We measured an average of 5.82% Mendelian inconsistency rate and 0.56% of missingness across the seven families (See **supplementary table 11**).

**Supplementary Figure 6**

Complete PCR gel confirmation of the two mosaic SVs shown in **Figure 4D** and **4E**.



**Supplementary Figure 7**

IGV screenshot of a 1,224 bases overlapping a mosaic DEL shown in **Figure 4E** in chr17. Here are shown five and four read support of the mosaic SV in Illumina data for two brain regions: cingulate cortex and adjacent white matter for the same brain.

chr9:12,558,472-12,559,470

p24.2  p23  p22.3  p22.1  p21.3  p21.2  p21.1  p13.3  p13.1  p11.1  q12  q13  q21.11  q21.13  q21.31  q21.33  q22.2  q22.32  q31.1  q31.2  q32  q33.1  q33.3  q34.11  q34.2

1,000 bp

12,558,500 bp   12,558,600 bp   12,558,700 bp   12,558,800 bp   12,558,900 bp   12,559,000 bp   12,559,100 bp   12,559,200 bp   12,559,300 bp   12,559,400 bp

msa_candidates_ont.bam

msa_candidates_illumina_bulk.b
m Coverage

[0 - 140]

Linked

msa_candidates_illumina_bulk.b
m

Sequence

Gene

msa_candidates.bed

DEL

chr9:111,043,806-111,044,804

p24.2  p23  p22.3  p22.1  p21.3  p21.2  p21.1  p13.3  p13.1  p11.1  q12  q13  q21.11  q21.13  q21.31  q21.33  q22.2  q22.32  q31.1  q31.2  q32  q33.1  q33.3  q34.11  q34.2

1,000 bp

111,043,900 bp   111,044,000 bp   111,044,100 bp   111,044,200 bp   111,044,300 bp   111,044,400 bp   111,044,500 bp   111,044,600 bp   111,044,700 bp   111,044,

msa_candidates_ont.bam

msa_candidates_illumina_bulk.b
m Coverage

[0 - 132]

Linked

msa_candidates_illumina_bulk.b
m

Sequence

Gene

msa_candidates.bed

DEL

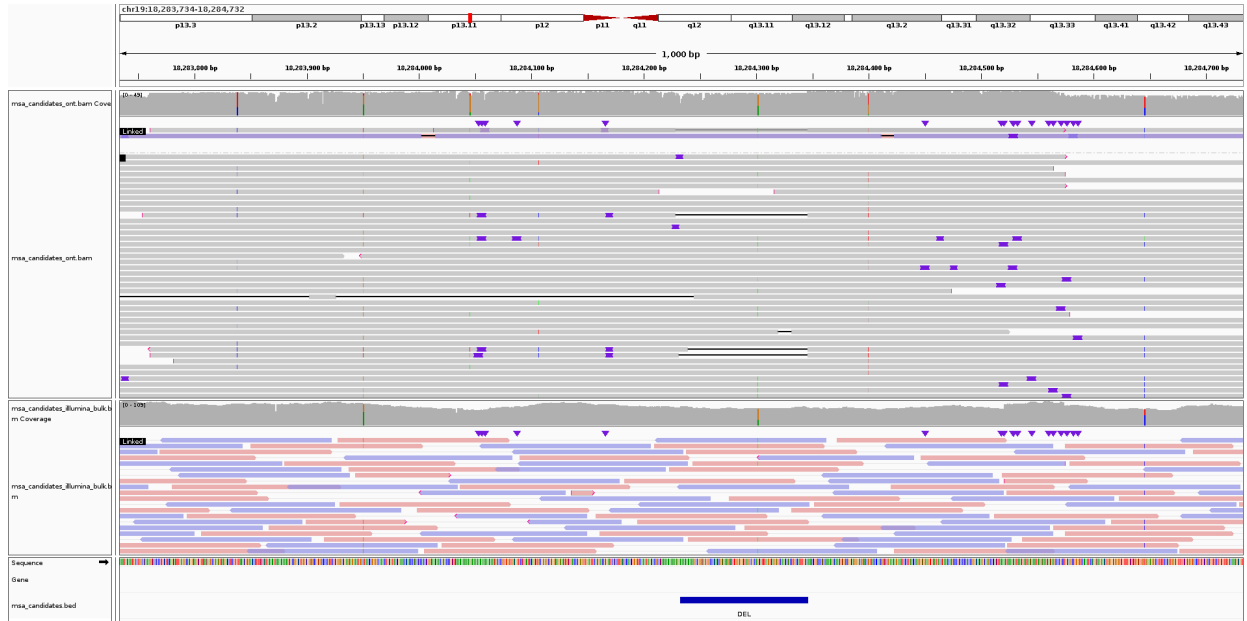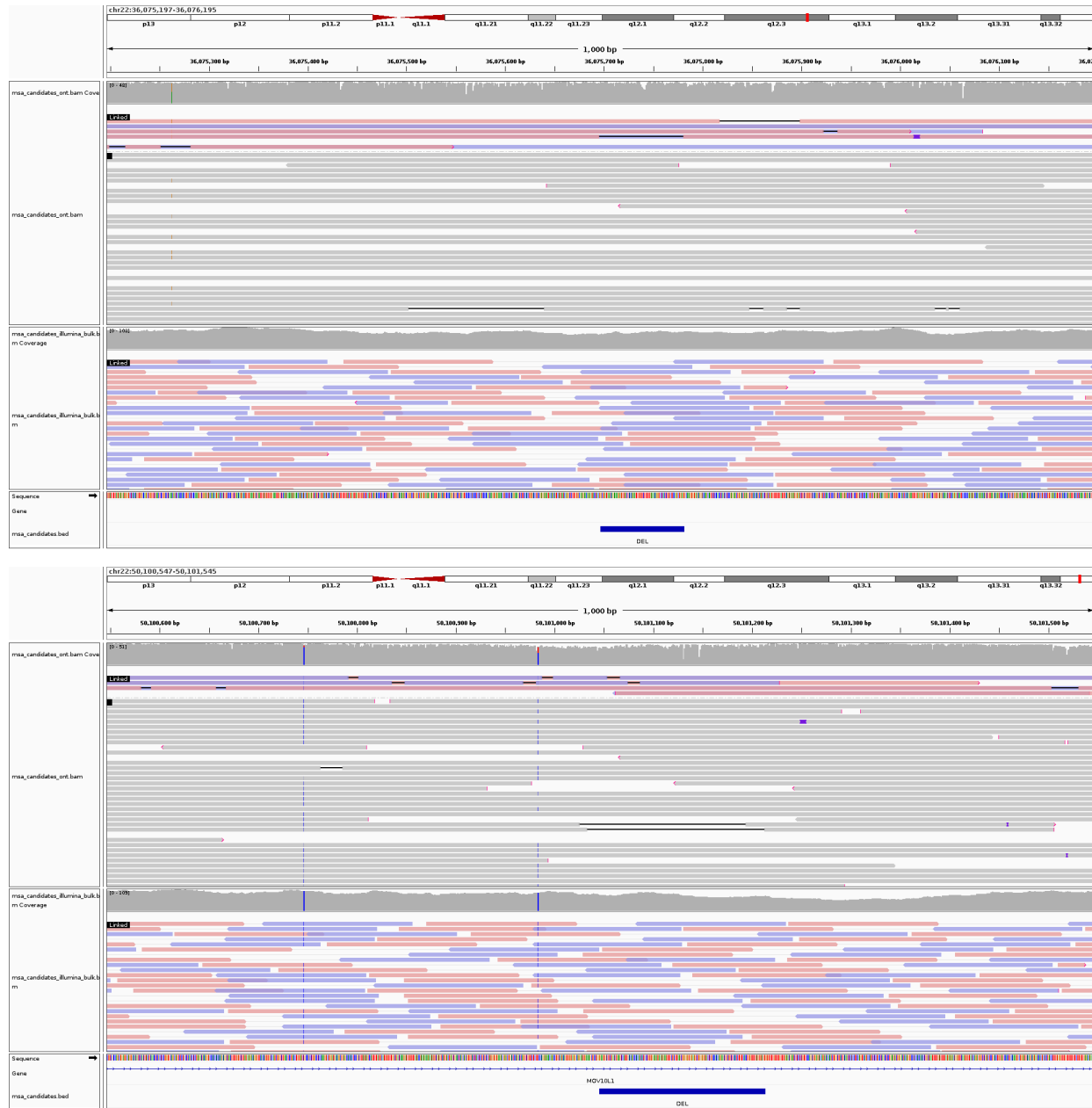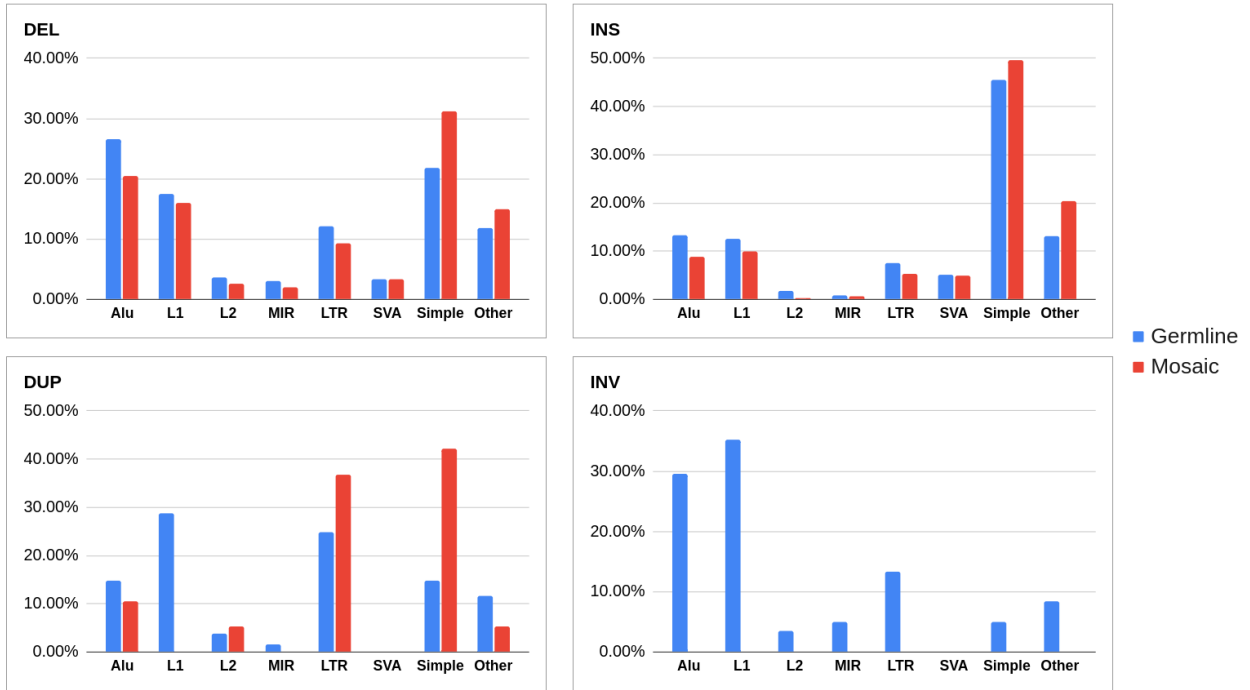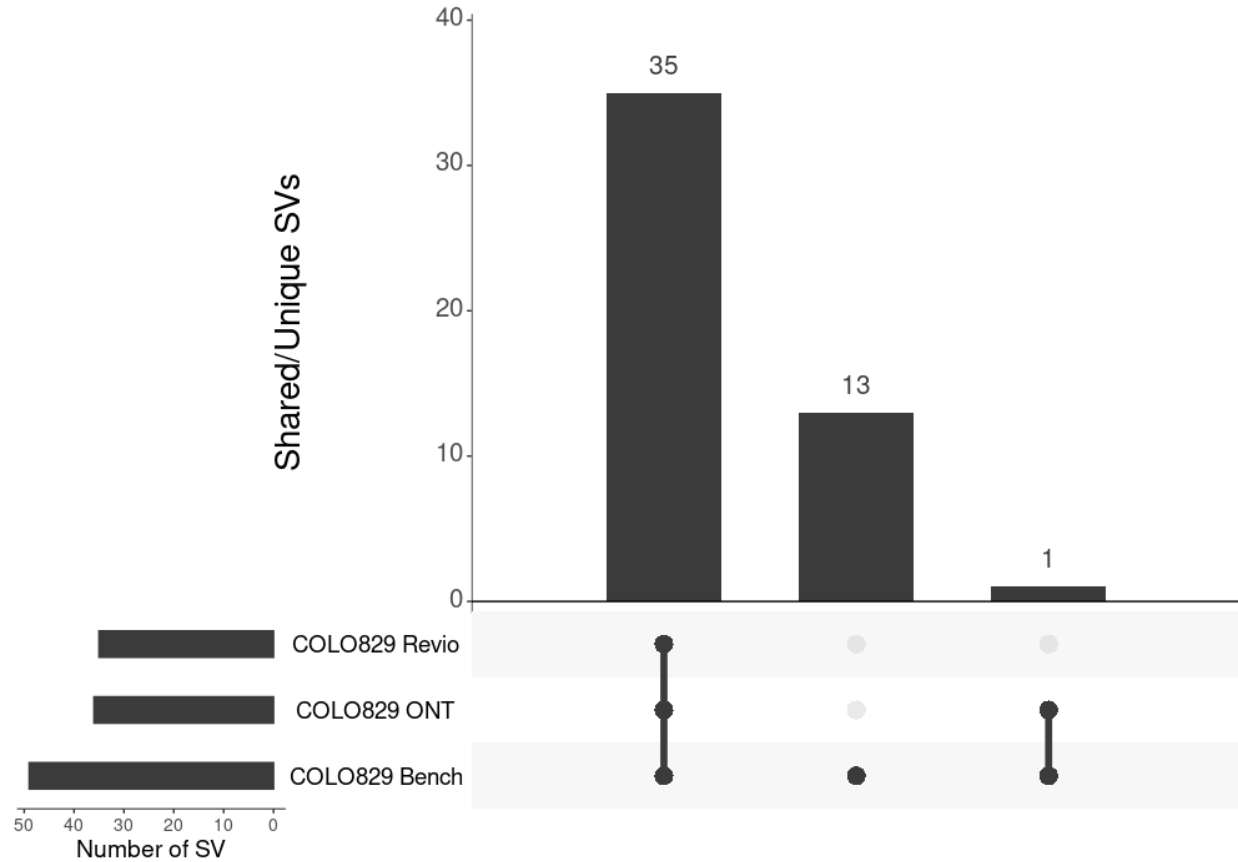chr15:98,576,830-98,577,828

chr16:81,576,933-81,577,931

**Supplementary figure 8 A-AH**: IGV screenshots of 34 mosaic SVs from the MSA sample. See **Supplementary table 14** for more details.
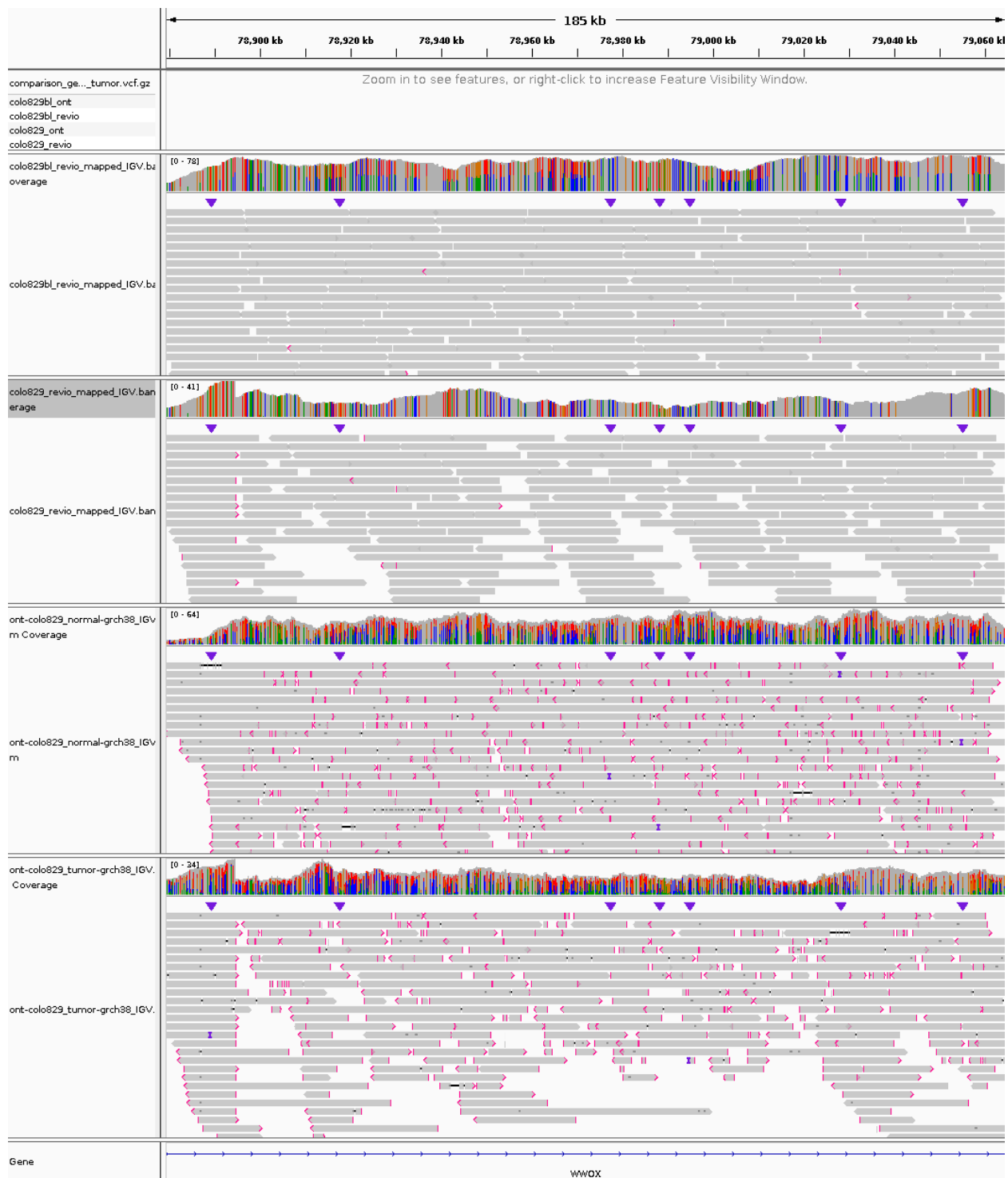
**Supplementary Figure 9**

Fraction of germline and non-germline/ mosaic SVs associated to different elements of repeat families in the cingulate cortex sample (ONT), split by SV type.
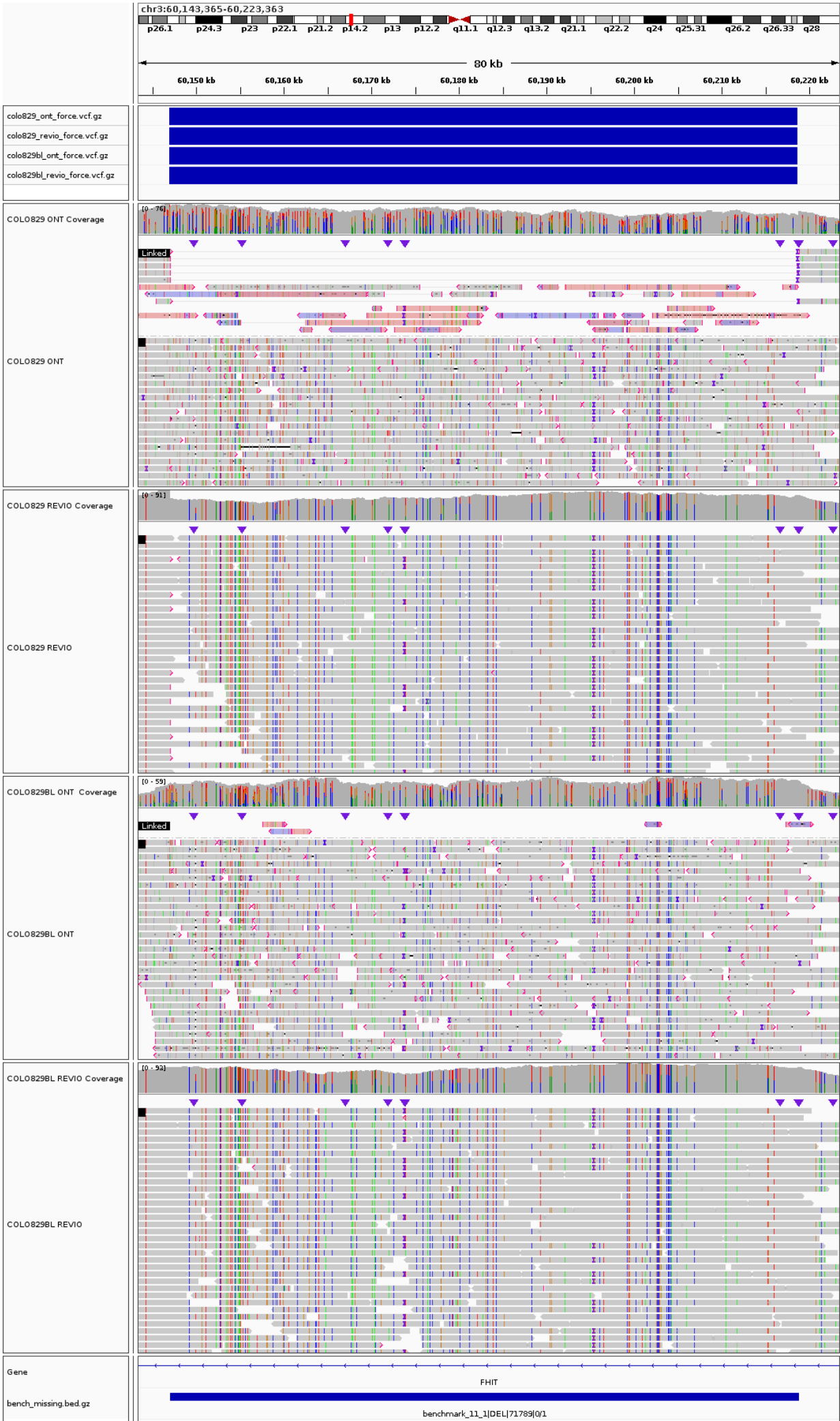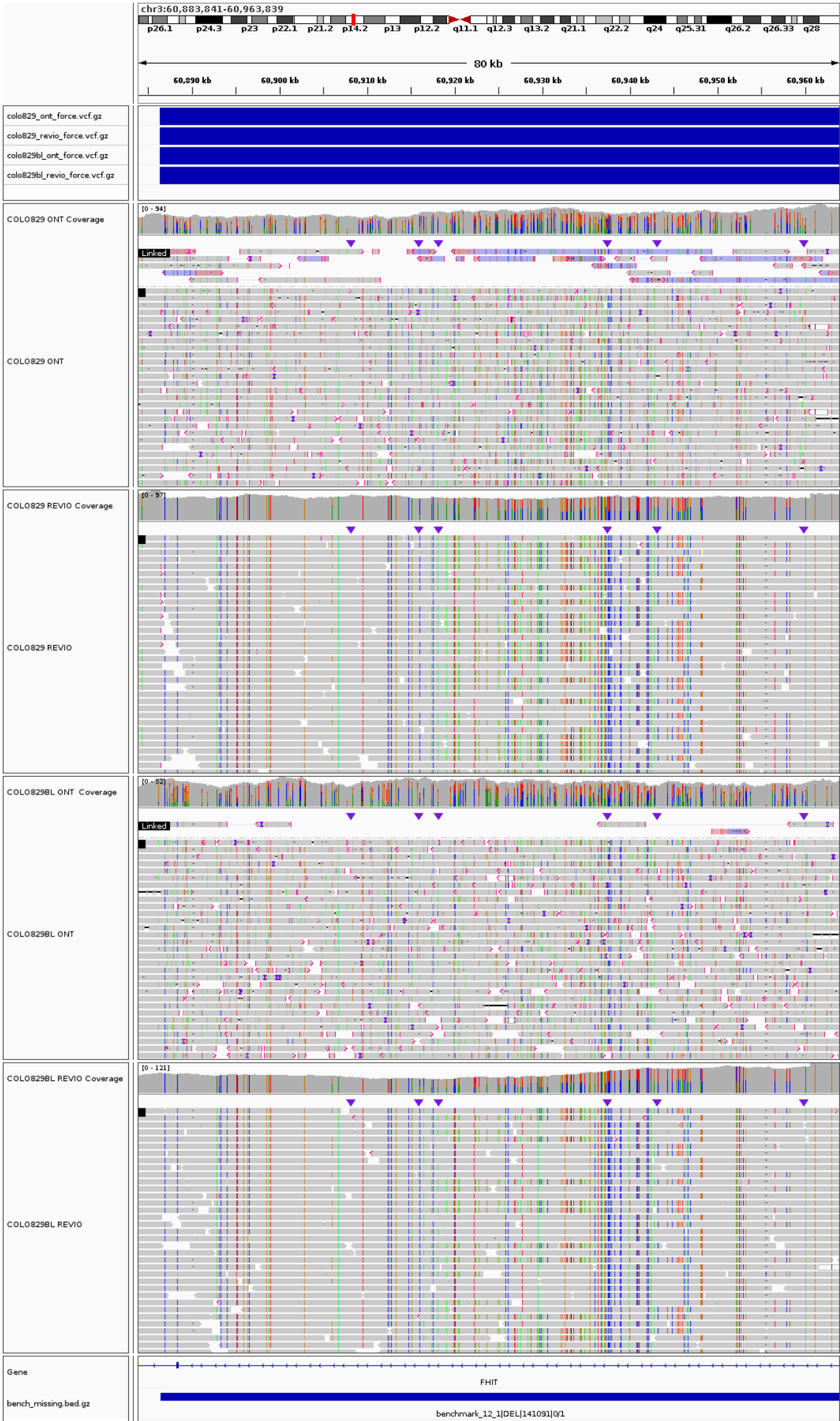
**Supplementary Figure 10**

Germline and mosaic SV called by Sniffles2 in the COLO829 cell line compared to a benchmark containing 49 SV (no BND). In total 36 SVs from the benchmark are detected and 13 are missed. From those 13: four had no reads for the alternative allele (0% AF), three were identified by Sniffles as mosaic (8.2-21% AF), four had heterozygous genotype (31-41% AF) and two had homozygous alternative genotype (81-100% AF)
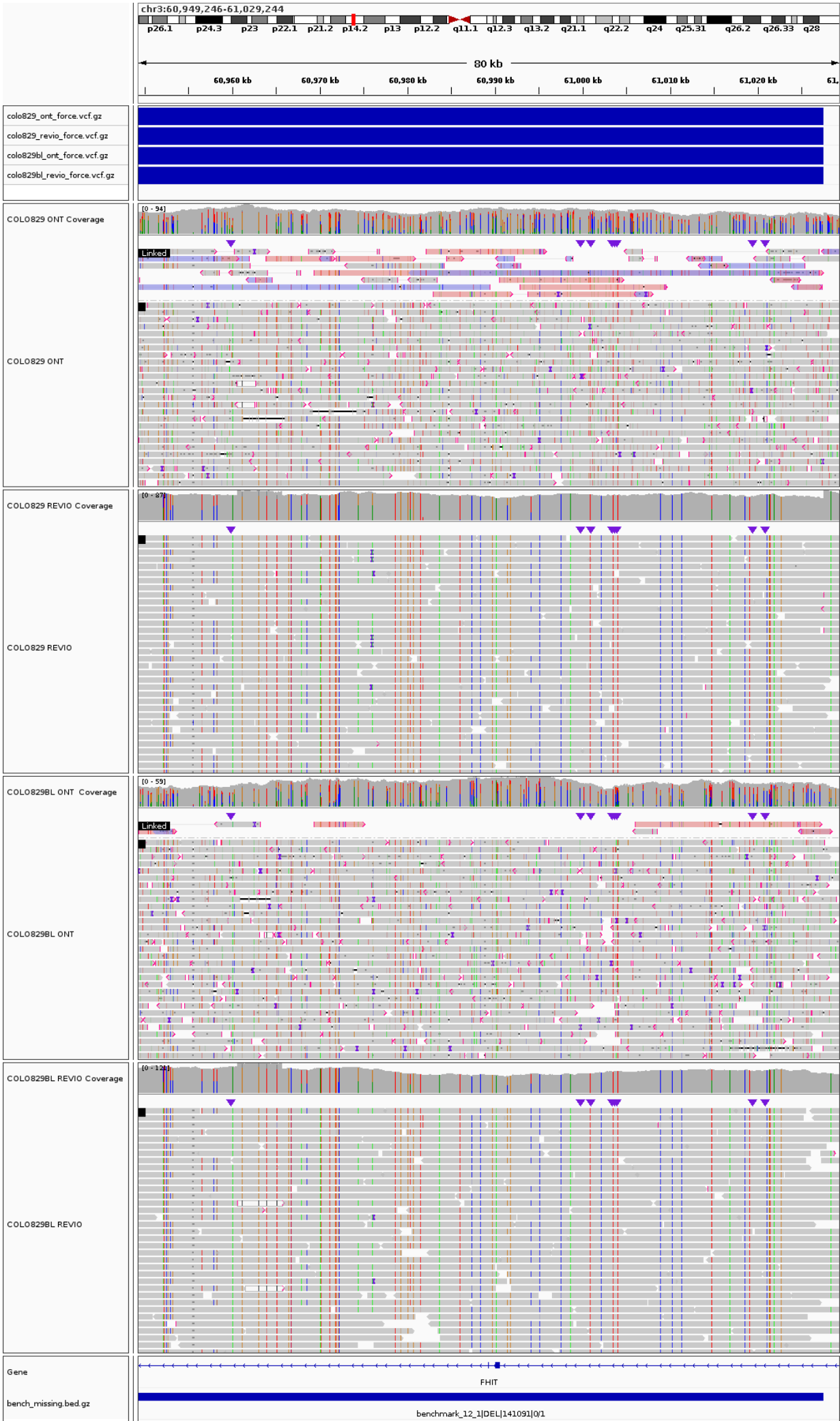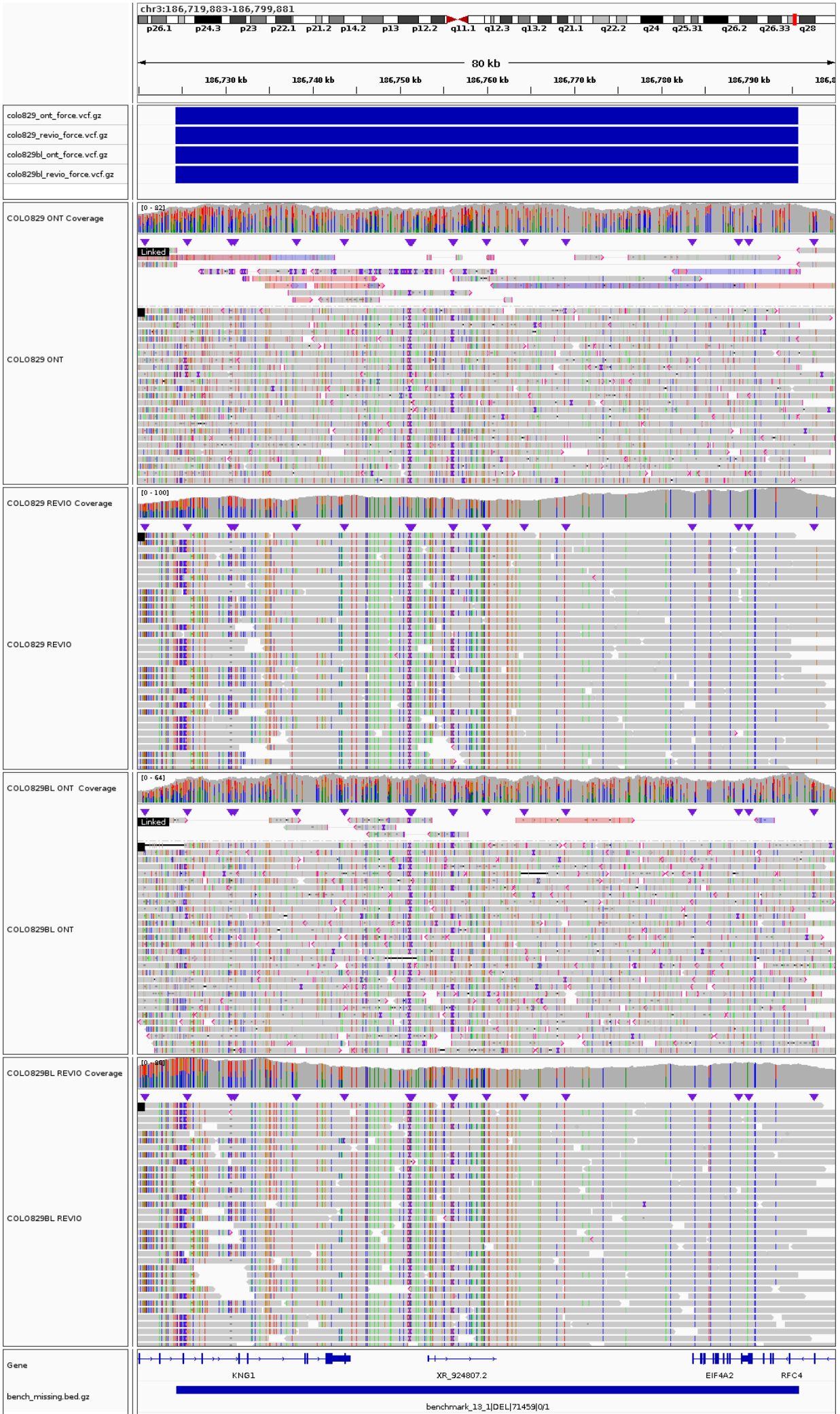
**Supplementary Figure 11**

Cancer specific germline SV call identified by sniffles2. We call germline SV for both the normal sample and cancer sample. After the merge we selected SV that were only detected in the cancer sample. Here is shown an example of an heterozygous deletion in an intron of the WWOX gene.
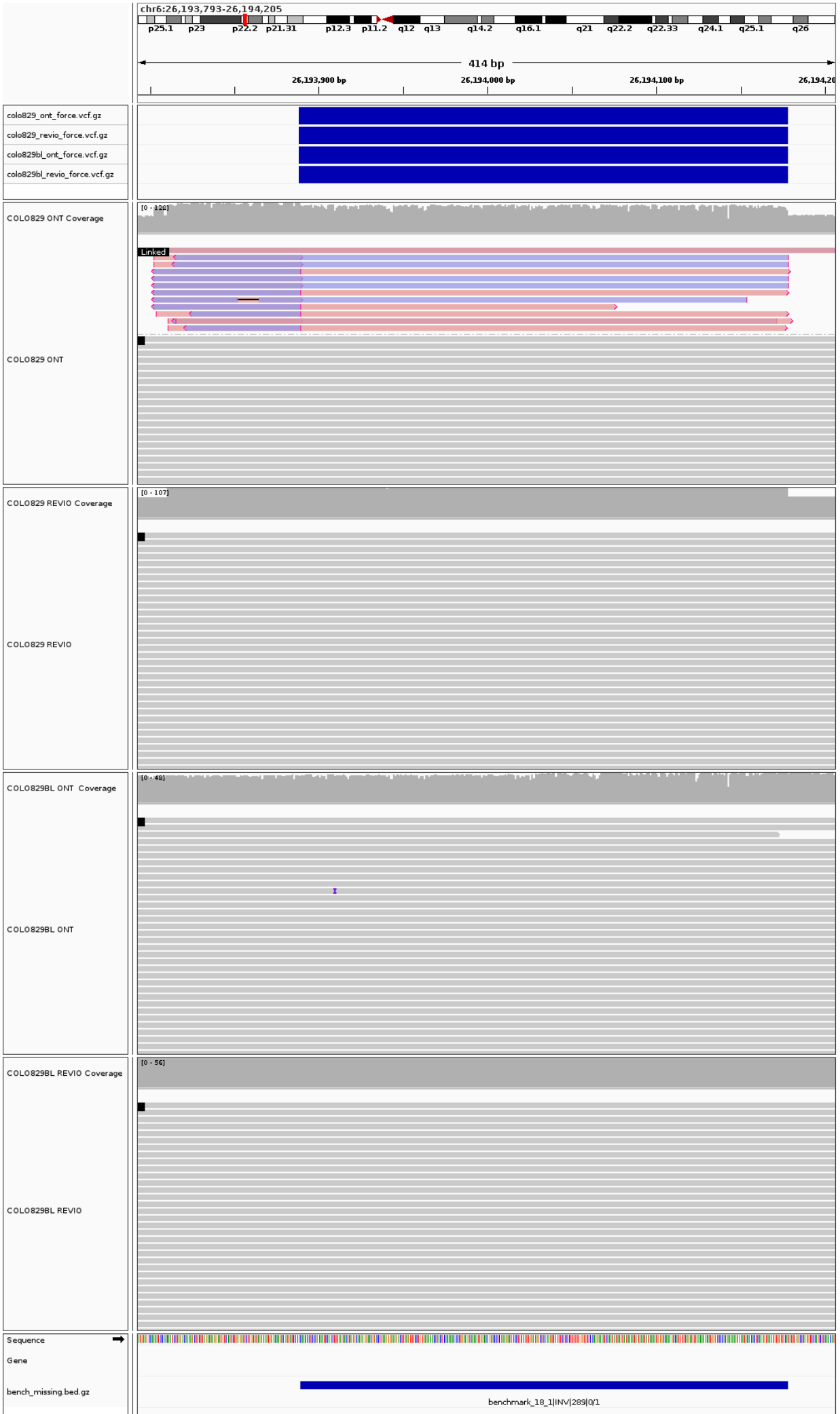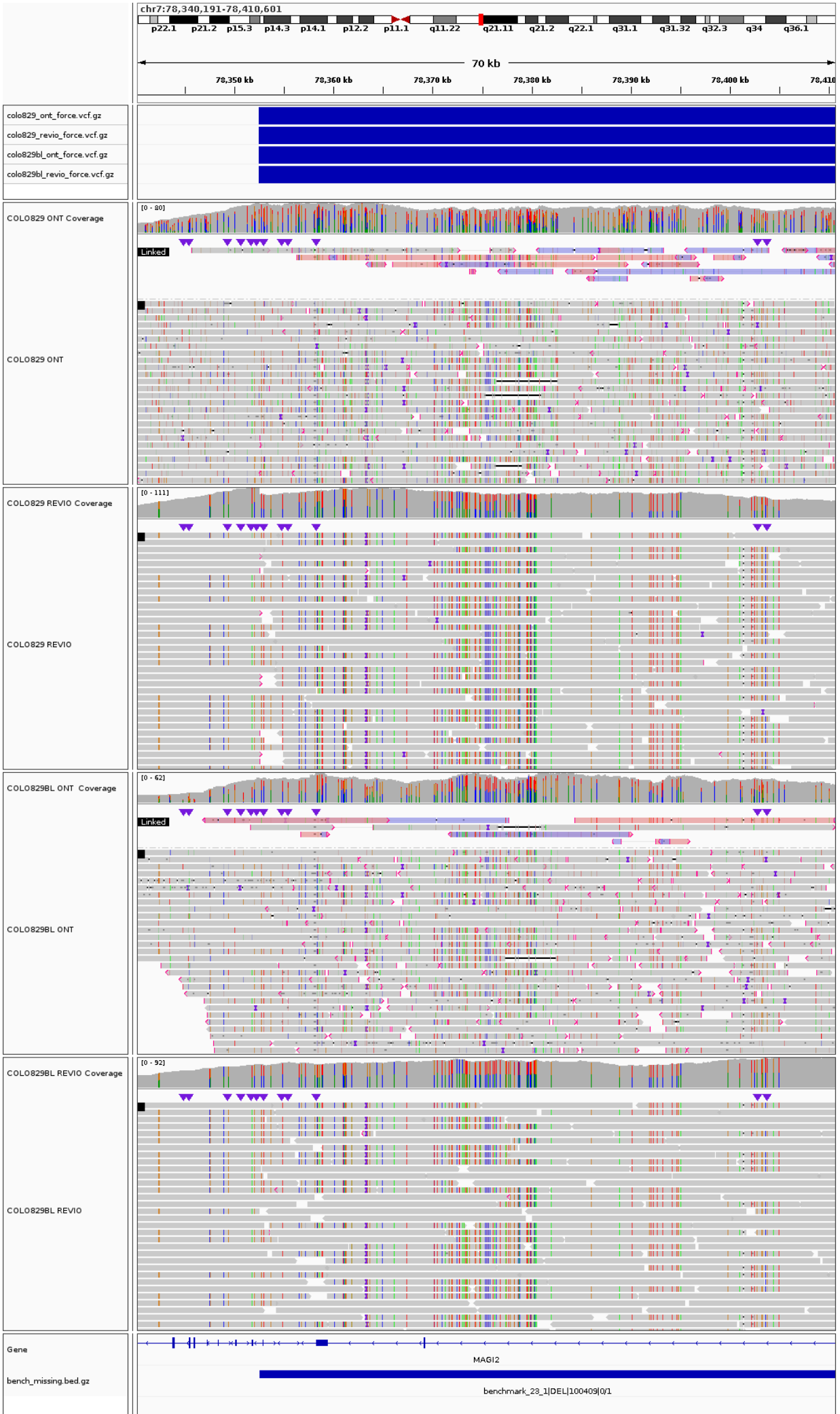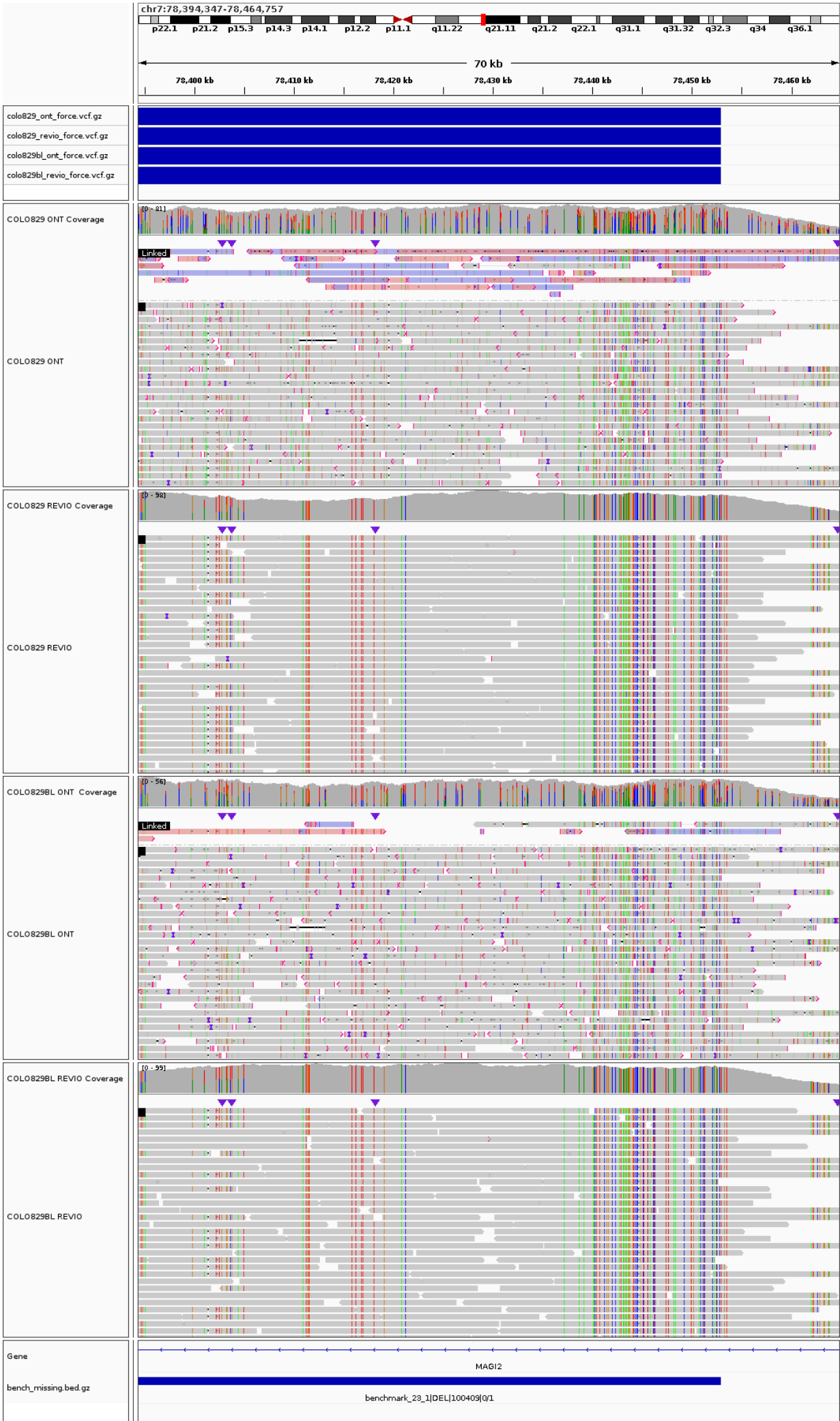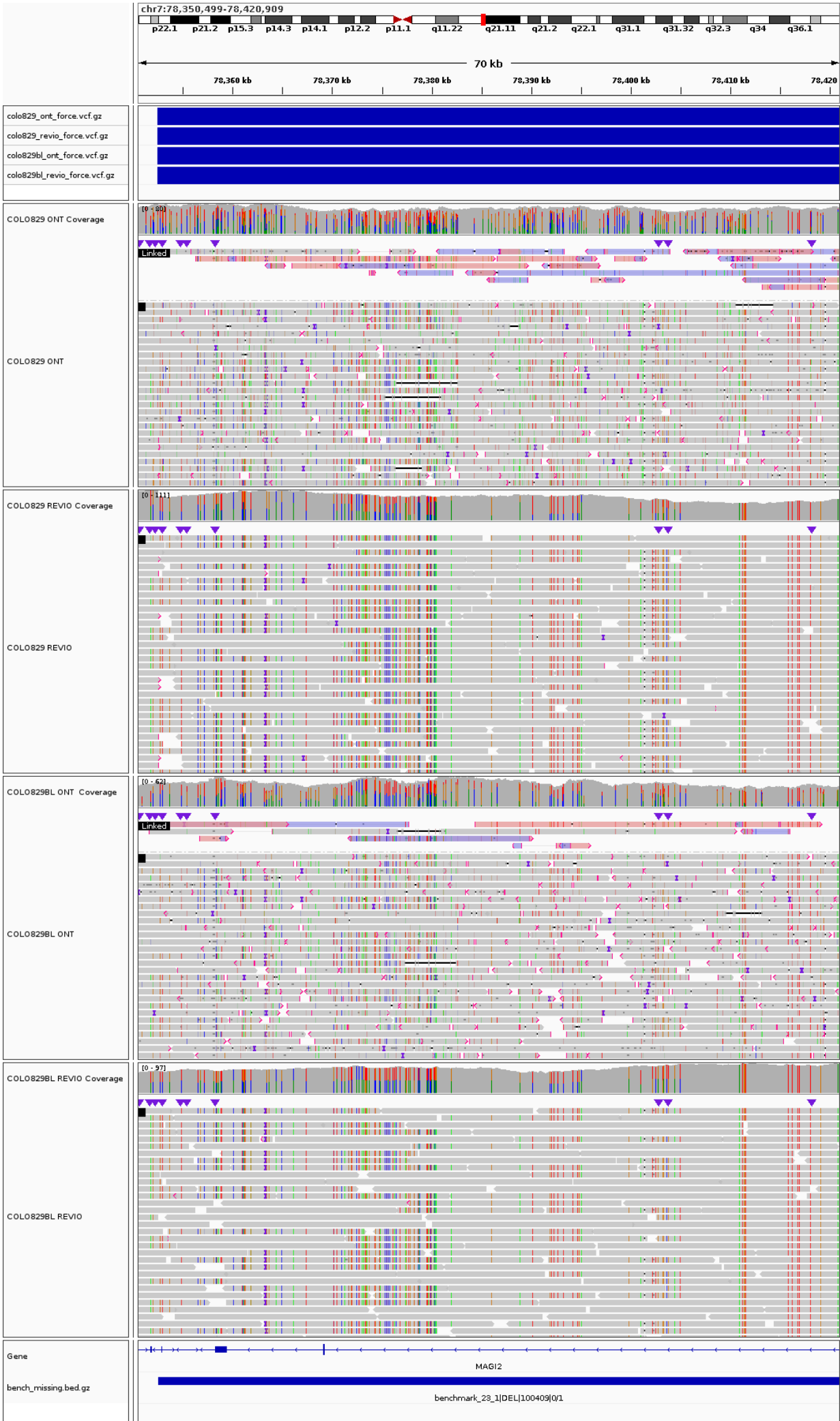
chr6:26,193,793-26,194,205

p25.1  p23  p22.2  p21.31  p12.3  p11.2  q12  q13  q14.2  q16.1  q21  q22.2  q22.33  q24.1  q25.1  q26

414 bp

26,193,900 bp        26,194,000 bp        26,194,100 bp        26,194,20

colo829_ont_force.vcf.gz

colo829_revio_force.vcf.gz

colo829bl_ont_force.vcf.gz

colo829bl_revio_force.vcf.gz

COLO829 ONT Coverage          [0 - 128]

COLO829 ONT

Linked

COLO829 REVIO Coverage        [0 - 107]

COLO829 REVIO

COLO829BL ONT  Coverage       [0 - 48]

COLO829BL ONT

COLO829BL REVIO Coverage      [0 - 56]

COLO829BL REVIO

Sequence

Gene

bench_missing.bed.gz

benchmark_18_1|INV|289|0/1

chr9:28,031,100-28,101,785

p24.1  p22.3  p21.3  p21.1  p13.1  q12  q13  q21.12  q21.2  q21.33  q22.32  q31.1  q31.3  q33.1  q33.3  q34.2

70 kb

28,040 kb    28,050 kb    28,060 kb    28,070 kb    28,080 kb    28,090 kb    28,100 kb

colo829_ont_force.vcf.gz

colo829_revio_force.vcf.gz

COLO829 ONT Coverage    [0 - 133]

COLO829 ONT

COLO829 REVIO Coverage    [0 - 137]

COLO829 REVIO

COLO829BL ONT  Coverage    [0 - 164]

COLO829BL ONT

COLO829BL REVIO Coverage    [0 - 262]

COLO829BL REVIO

Gene

LINGO2                      LINGO2

bench_missing.bed.gz

benchmark_36_1|DEL|123393|0/1

chr14:105,865,689-105,869,832

p12   p11.2  q11.1      q12   q13.1   q21.1      q21.3  q22.2  q23.1   q23.3  q24.2      q31.1   q31.3  q32.12    q32.31

4,145 bp

105,866,000 bp          105,867,000 bp          105,868,000 bp          105,869,000 bp

colo829_ont_force.vcf.gz
colo829_revio_force.vcf.gz
colo829bl_ont_force.vcf.gz
colo829bl_revio_force.vcf.gz

COLO829 ONT Coverage          [0 - 64]

COLO829 ONT

COLO829 REVIO Coverage        [0 - 39]

COLO829 REVIO

COLO829BL ONT  Coverage       [0 - 10.00]

COLO829BL ONT

COLO829BL REVIO Coverage      [0 - 10.00]
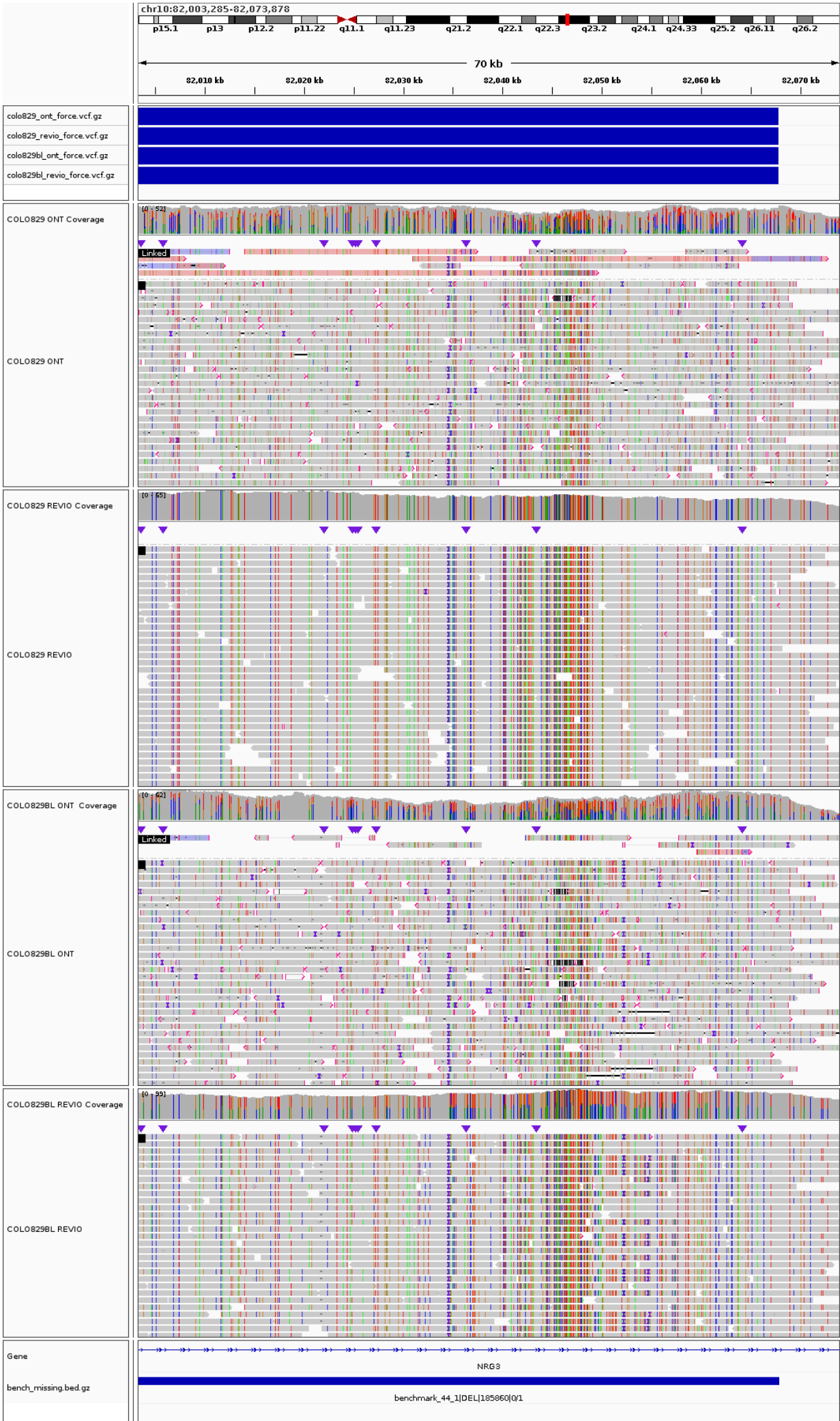
COLO829BL REVIO

Gene
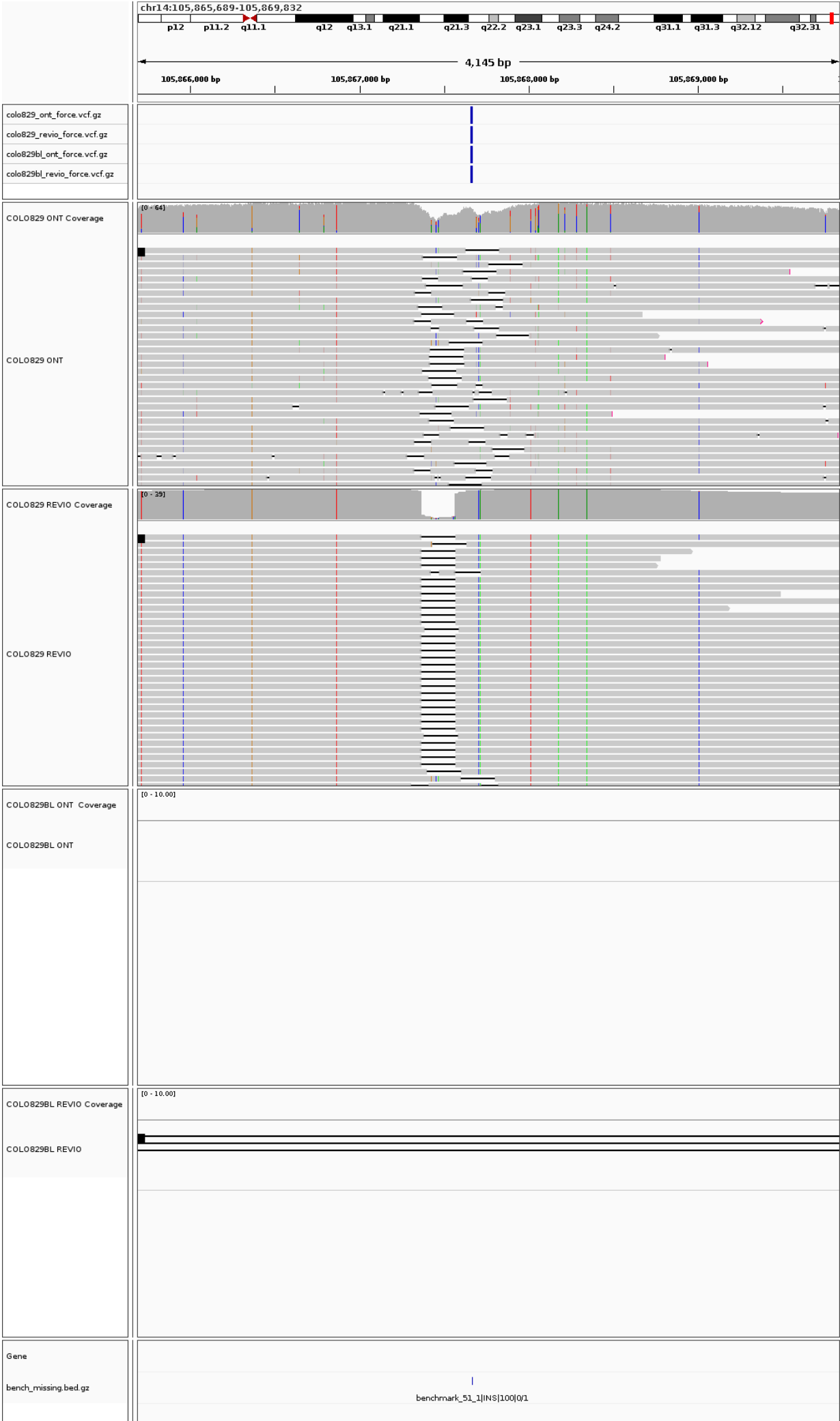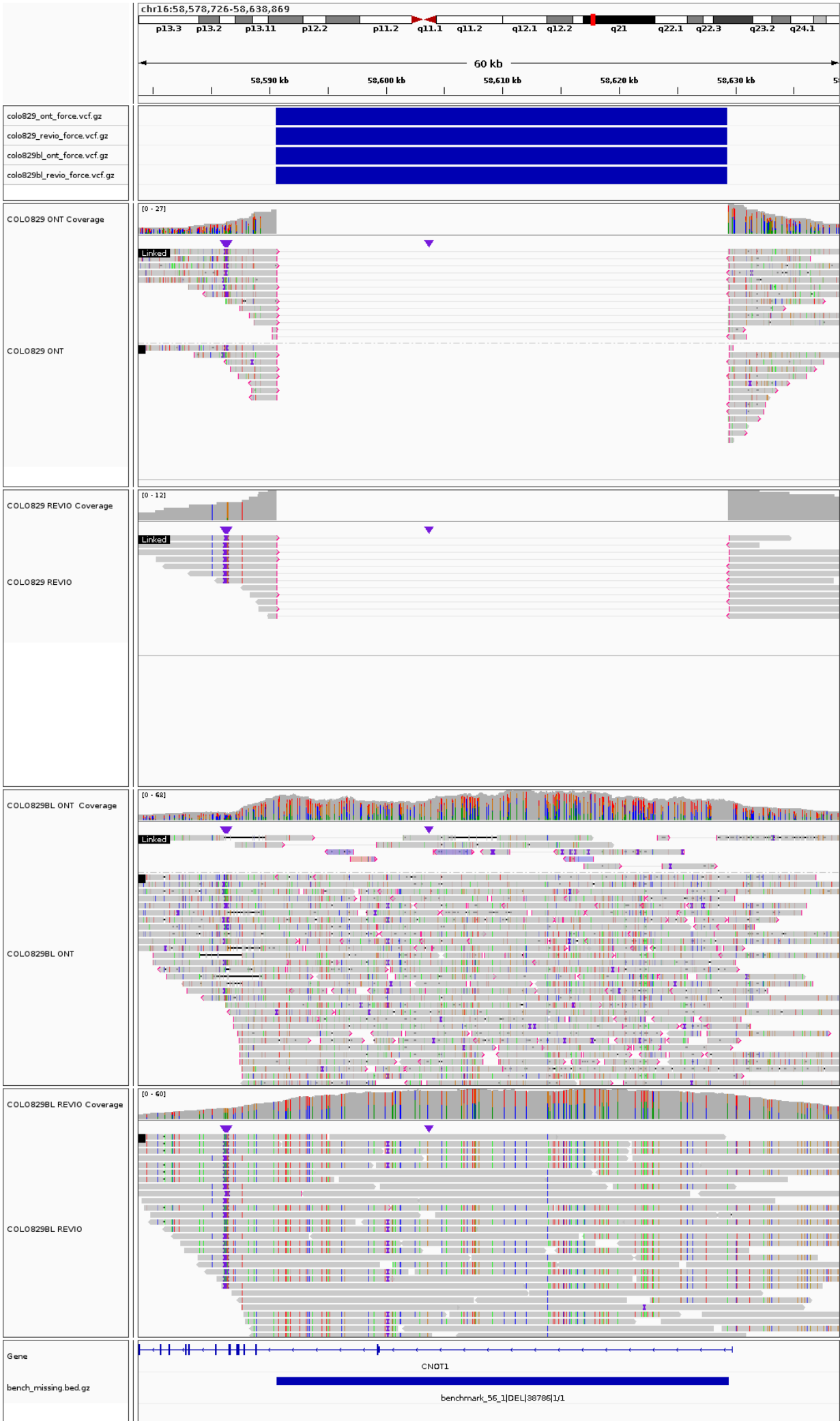
bench_missing.bed.gz
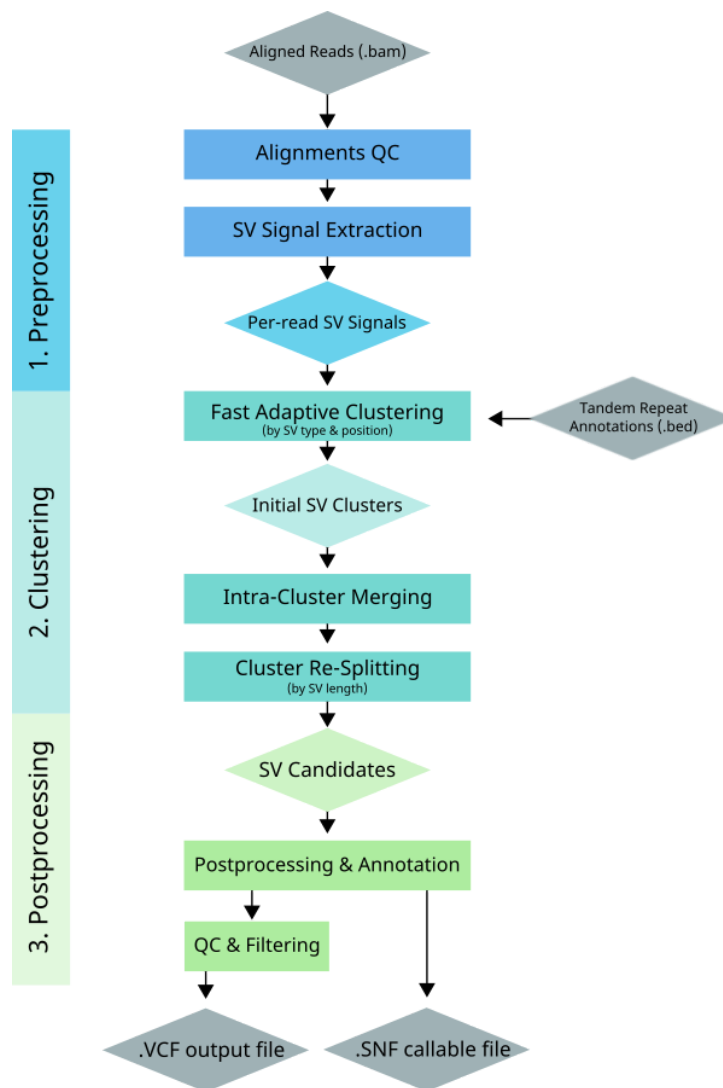
benchmark_51_1|INS|100|0/1

**Supplementary Figure 12.** IGV screenshots of the 13 missing SVs from the COLO829 benchmark (see **Supplementary table 18**)

**Supplementary Figure 13**
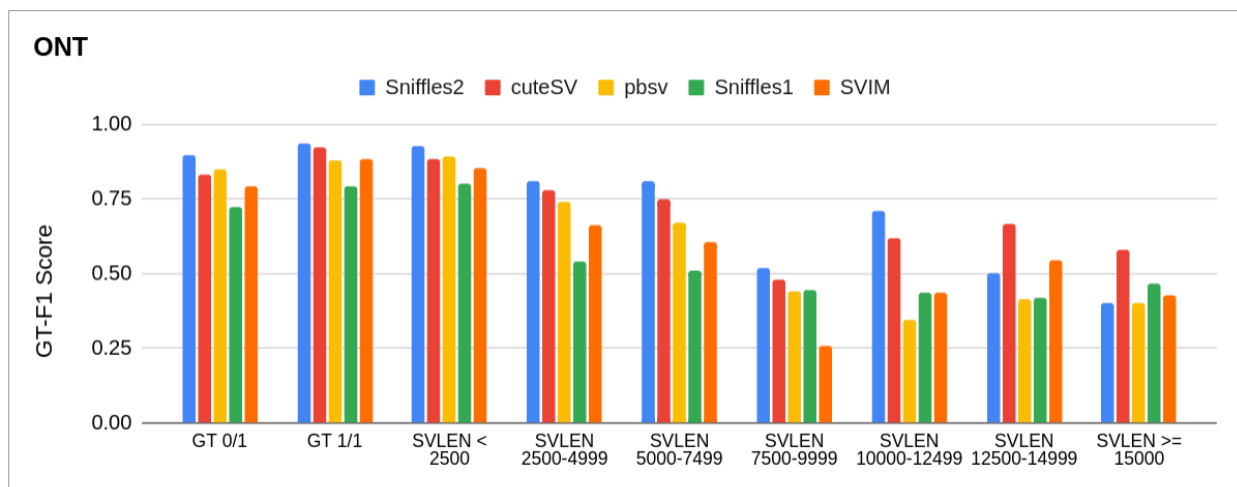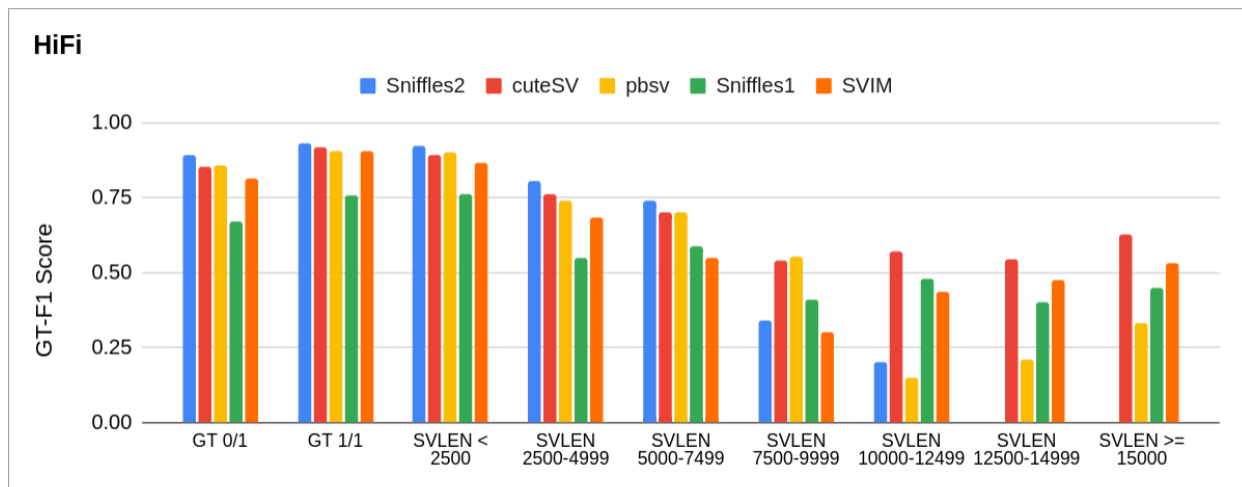
Steps in the Sniffles2 germline mode algorithm

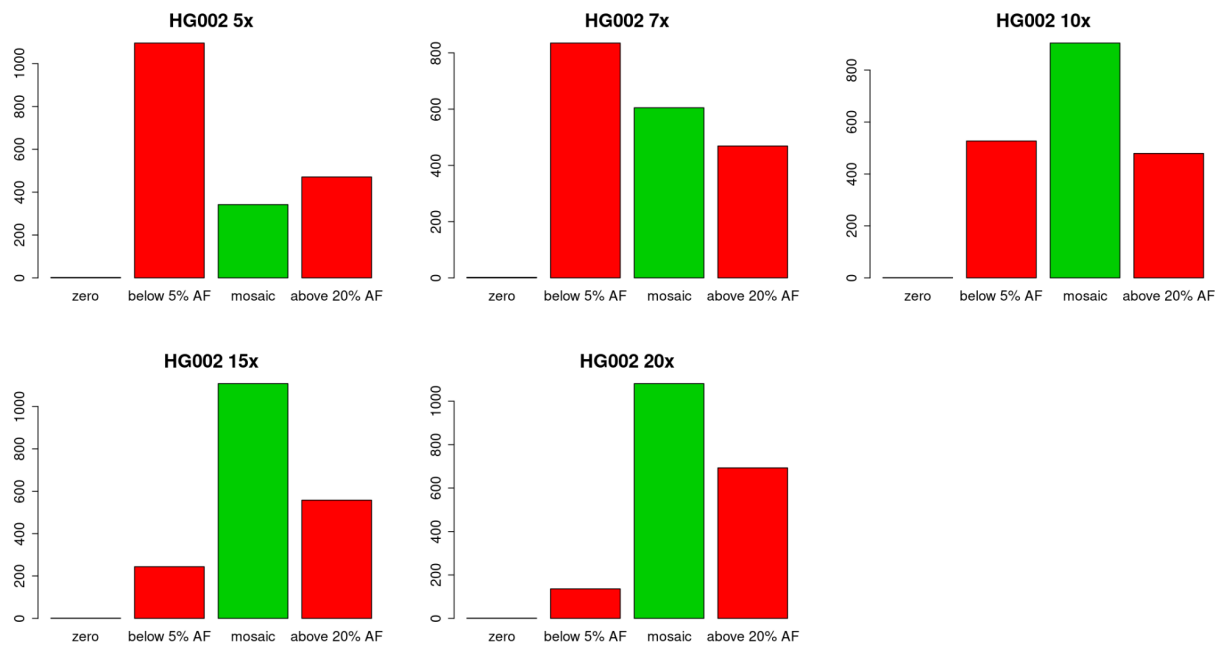| SAMPLE_1.snf | SAMPLE_2.snf | ... | SAMPLE_N.snf |
|---|---|---|---|
| INDEX | INDEX | ... | INDEX |
| Block 1: chr1 0-10000<br><br>SVCand1 (INS)<br>SVCand2 (DEL)<br>...<br>SVCandN (DUP)<br><br>Coverage Track | Block 1: chr1 0-10000<br><br>SVCand1 (DEL)<br>SVCand2 (INS)<br>...<br>SVCandN (INS)<br><br>Coverage Track | ... | Block 1: chr1 0-10000<br><br>SVCand1 (DUP)<br>SVCand2 (BND)<br>...<br>SVCandN (INS)<br><br>Coverage Track |
| Block 2: chr1 10000-20000 | Block 2: chr1 10000-20000 | ... | Block 2: chr1 10000-20000 |
| ... | ... | ... | ... |
| Block n | Block n | ... | Block n |

**Supplementary Figure 14**

Schematic for the SNF file format

**Supplementary Figure 15**

Genotype F1 measure for the detection of insertions and deletions by genotype and SV length. Measured by using the Truvari summary output files. For the evaluation of each caller we used the GIAB Oxford Nanopore and Pacbio HiFi 30x coverage titration benchmark. All callers were benchmarked using default parameters. The plot shows F1-genotype measure (Y-axis) per category (X-axis). The complete results are shown in **Supplementary Table 19.** Notice that over 94% of the SV are of smaller to 2500 bases (absolute size).

**Supplementary Figure 16**

Distribution of SV from HG002 by their AF in each of the five tested titrations (5x, 7x, 10x, 15x, 20x). In green are those that Sniffles2 mosaic mode can detect, the rest will not be assessed nor reported. Notice that even at 5x, some SVs do have AF greater than 20%.