

Supplementary Information

ILLUMINA SYNTHETIC LONG READ SEQUENCING ALLOWS RECOVERY OF MISSING SEQUENCES EVEN IN THE “FINISHED”

C. elegans Genome

Runsheng Li¹, Chia-Ling Hsieh², Amanda Young², Zhihong Zhang², Xiaoliang Ren¹, Zhongying Zhao^{1,3,*}

¹Department of Biology, Hong Kong Baptist University, 224 Waterloo Road, Kowloon Tong, Hong Kong, China, ²Illumina Inc., 5200 Illumina Way, San Diego, 92122, USA, ³State Key Laboratory of Environmental and Biological Analysis, Hong Kong Baptist University, Hong Kong, China

* To whom correspondence should be addressed: Zhongying Zhao

Email: zyzhao@hkbu.edu.hk

Phone: 852-3411-7058

This file includes:

Supplementary Figures S1 to S6 and the Figure legend

Supplementary Figure Legends

Supplementary Figure S1. Flowchart of data analysis.

Supplementary Figure S2. Differential coverage of genome or repetitive sequences by read of various lengths. (A) Histogram of read coverage on *C. elegans* chromosome III. Read length is scaled in grey as indicated. SNP density is shown on the top. Track of mapping gap is shown in black while those of repetitive sequences are differentially color coded and shown at the bottom. (B) A magnified view of the biggest gap in panel A which occupies a 47 Kb region consisting exclusively of satellite repeat as highlighted at the bottom. The aligned reads of both orientations are shown on each side of the gap. (C) Output of “dotplot” using the gap sequence in panel (B) as an input. Scale of sequence similarity is indicated on the top.

Supplementary Figure S3. Read coverage of ribosomal genes. (A) Read coverage of the ribosomal gene cluster consisting of 18S/5.8/28S genes currently placed at the end of chromosome I. Annotation of the gene cluster is shown at the bottom. Note, only a single copy of the ribosomal gene cluster was placed in the *C. elegans* reference genome. (B) Read coverage of the ribosomal gene 5S genomic region located on chromosome V. The red and blue bars indicate the unique mapping of forward and reverse reads, while the yellow bar indicates ambiguous mapping, meaning that a read was mapped to different genomic regions with equal efficiency.

Supplementary Figure S4. Density of deletion (red), insertion (green) and SNV (blue) events identified using the long reads over *C. elegans* chromosomes.

Supplementary Figure S5. Validation of a revised gene model based on an insertion that is located within an intron. (A) A revised gene model of *sul-2* with an insertion that adds a new exon. Pink, existing exon, yellow, newly added exon. (B) A magnified view of the affected exon that is confirmed by RNA-seq data with its mapping reads shown at the bottom. (C) Partial multiple alignment using the *C.elegans* SUL-2 protein sequence and that of its orthologs in *C.briggsae* (CBR), *C. remanei* (CRE) and *C. japonica* (CJP). Discrepancies in alignment are highlighted in grey.

Supplementary Figure S6. An example of translocation event found in the contig assembled by Celera. (A) Shown is the entire region of the contig assembled with two smaller contigs that are nevertheless derived from chromosome V (red thick bar) and I (blue thick bar) respectively. Chromosomal coordinates for the two fragments are indicated on the top. Note, the two are merged together by the Celera due to the single read which was preassembled incorrectly likely because of the overlapping sequences of around 1.2 Kb in size that contain nearly identical DNA repeats (thick yellow bar) between the two fragments. Reads used to assemble the “wrong” contig shown in panel A are drawn in thin bars. (B and C) Mapping of the incorrectly preassembled read (indicated by arrow) and its cognate genomic regions on chromosome V and I respectively with all other reads except for the wrong one. Note that the left and right arm of the read can align properly with its cognate genomic sequences respectively while the other arms remain unaligned (shown as wavy lines).

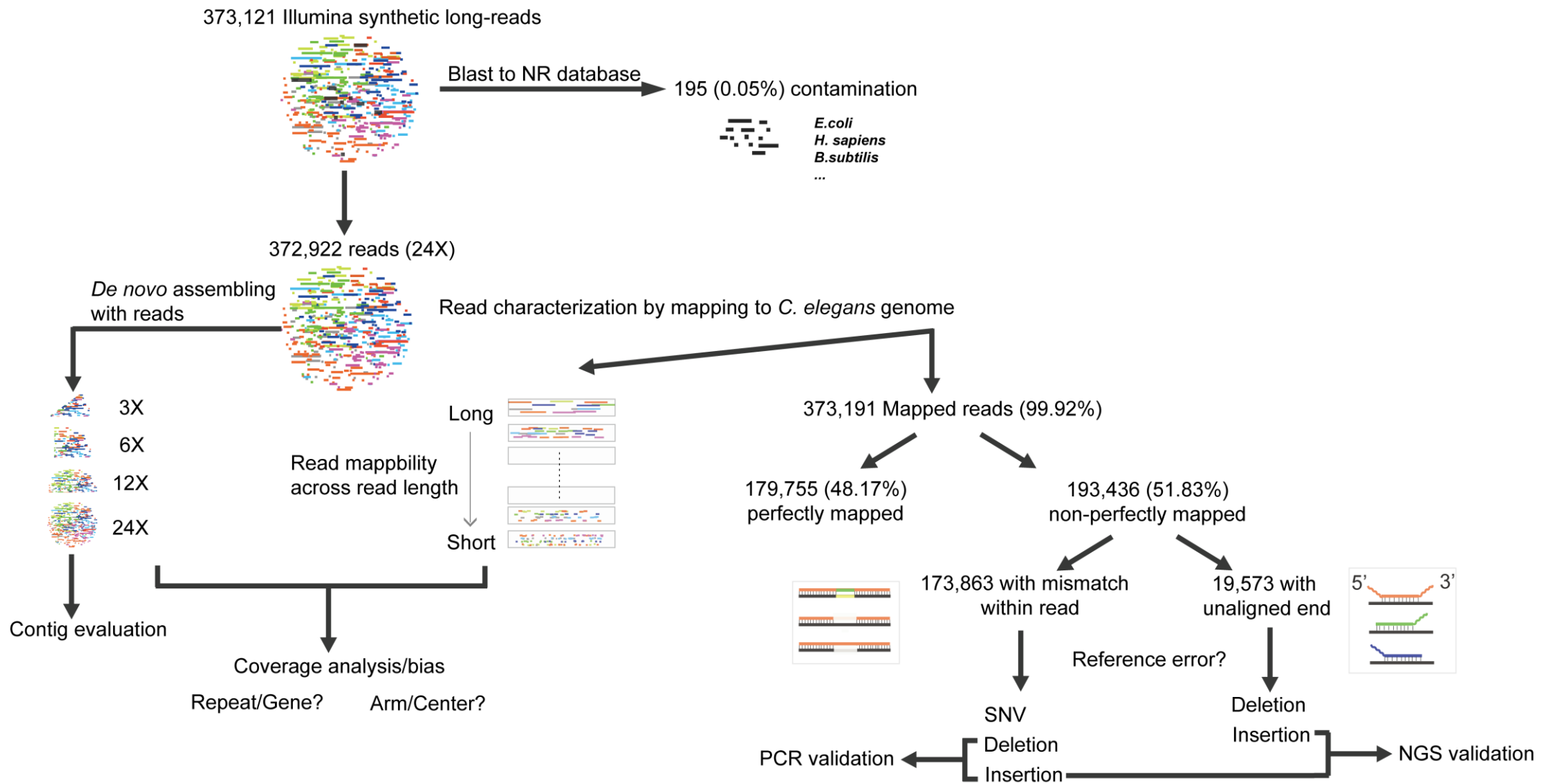


Figure S1

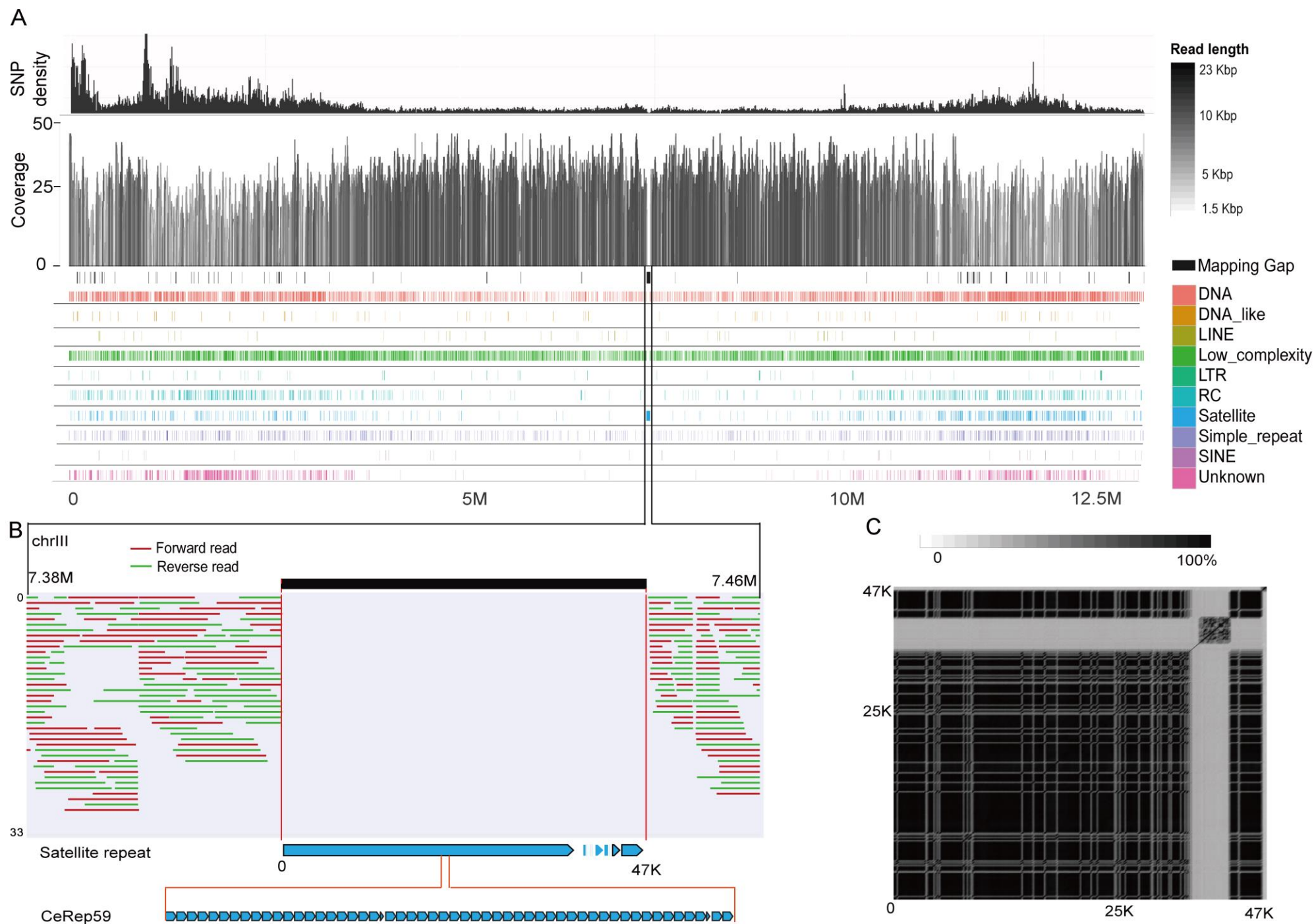


Figure S2

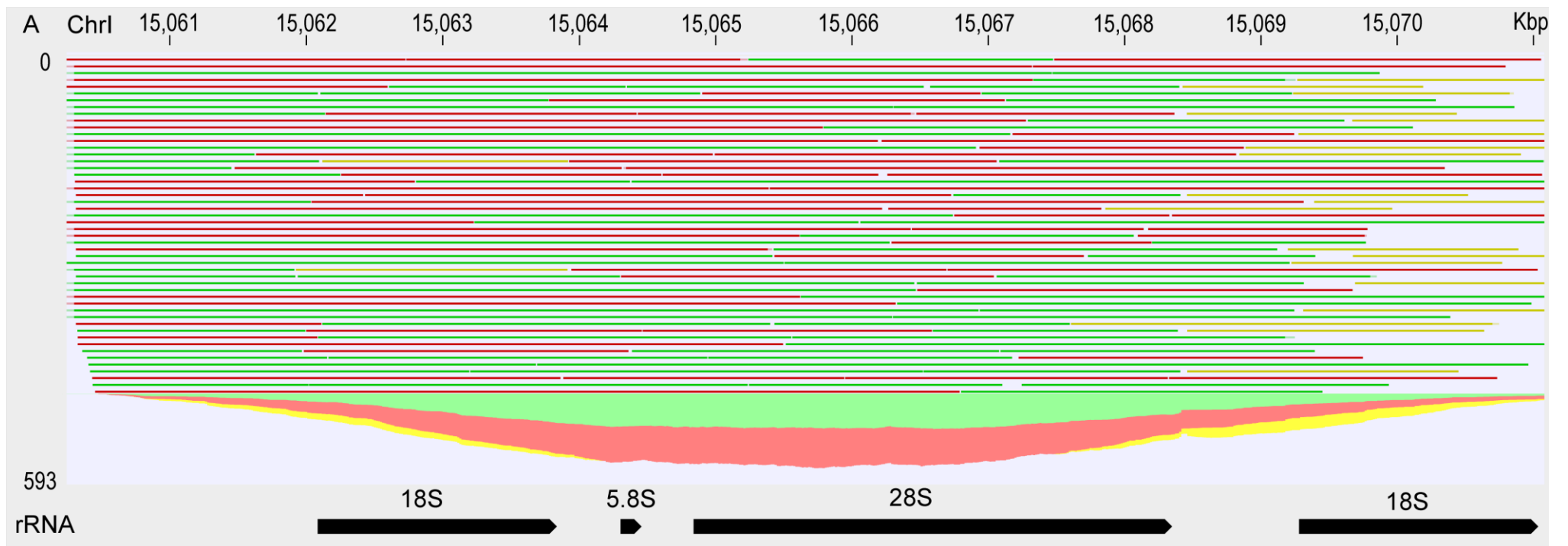


Figure S3

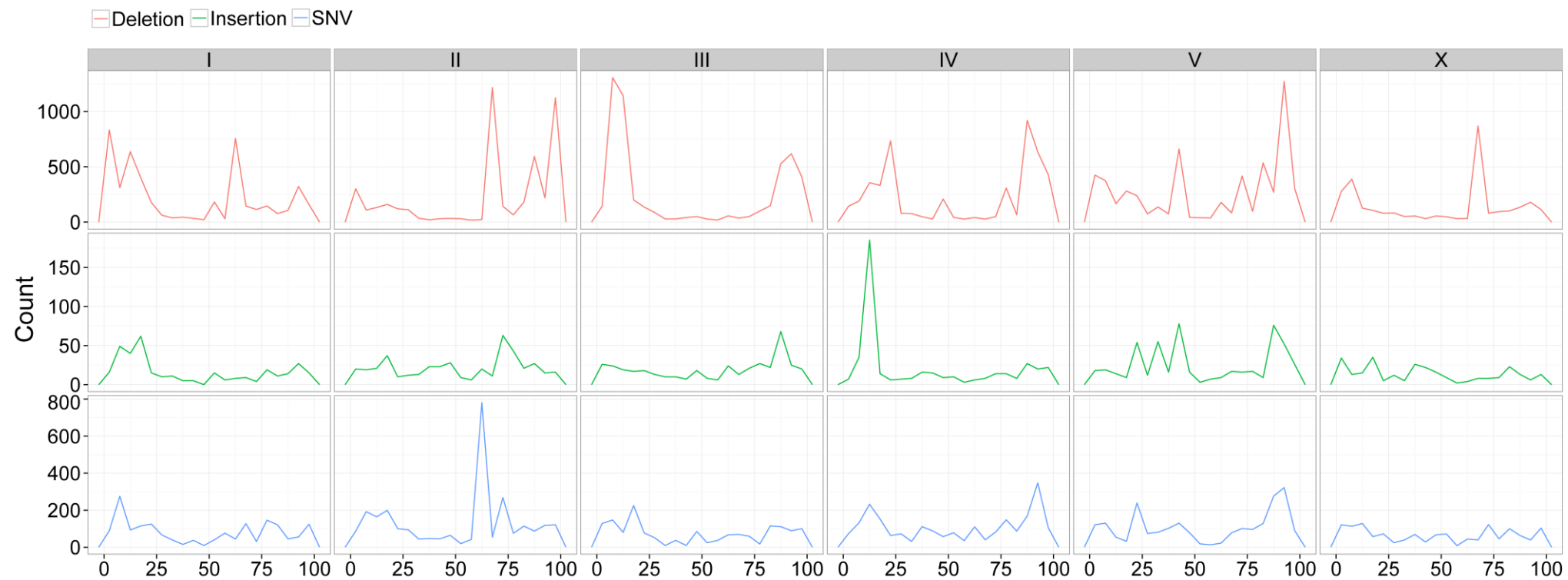


Figure S4

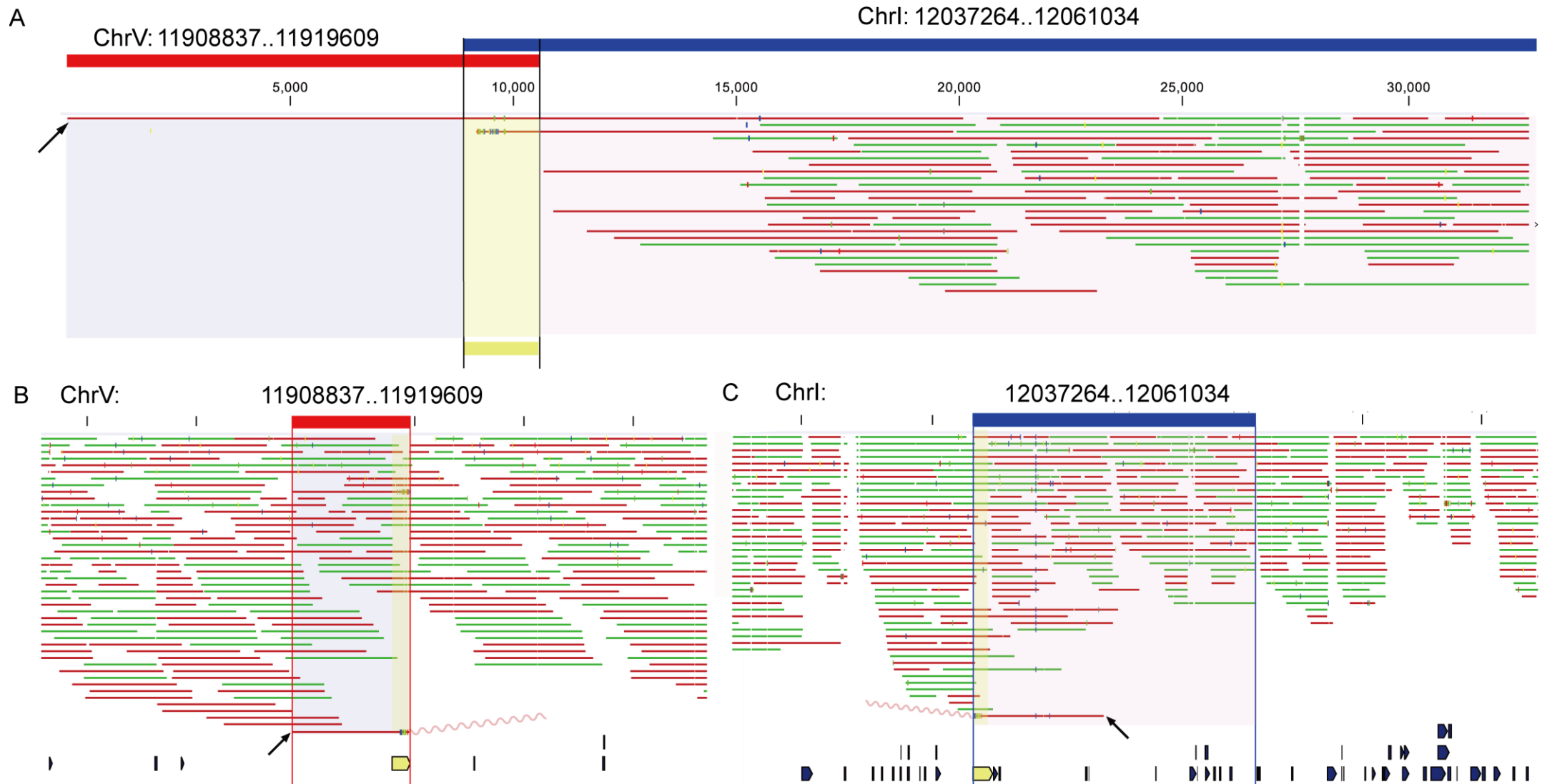


Figure S6