# Structural Variation Detection using NextGENe® software

February 2021        Kevin LeVan, Jonathan Liu, Megan McCluskey, Ni Shouyong

## Introduction

Over 68,000 structural variants (SVs) have been reported in the human genome, including deletions, duplications/insertions, and translocations such as gene fusions (1). Alone, there are more than 2,200 catalogued gene fusions (2) making genomic disorders an important area of research (3).

There are several techniques available for detecting SVs. Microarrays, for one, can be useful for detecting copy number variation but are limited to detecting events relative to specific design probes (4). Next Generation Sequencing can produce higher resolution results and enable the identification of novel events as well as the simultaneous detection of SNPs and small indels.
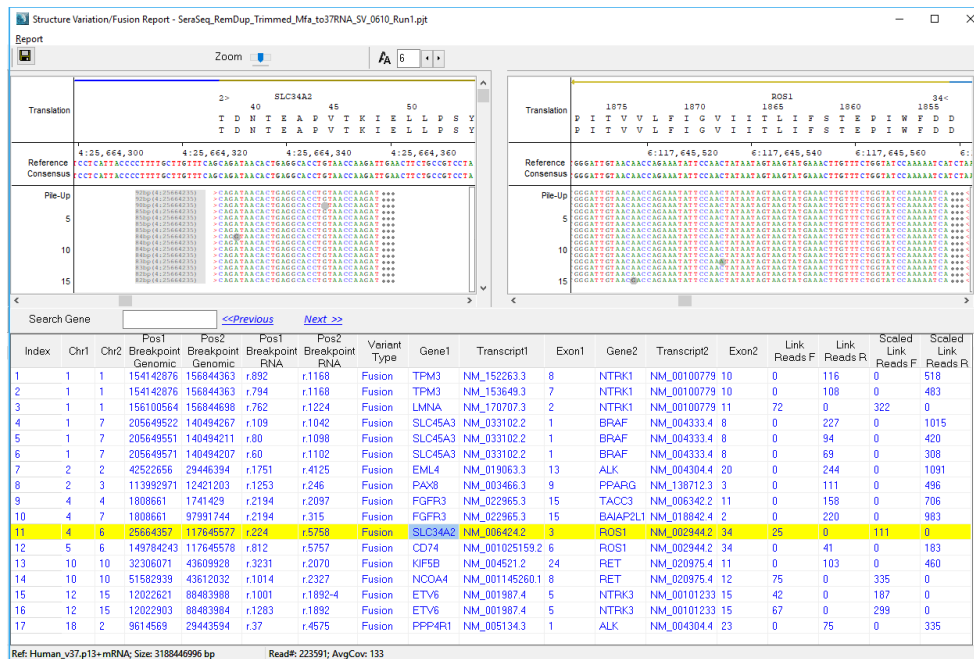


*Figure 1: NextGENe Viewer's Structural Variation Report is an interactive tool for displaying both breakpoints of each detected SVs.*

NextGENe software can align reads simultaneously to genomic and transcriptomic references. Reads from RNA samples can align across exon junctions more easily, enabling higher sensitivity for the detection of breakpoints. Next generation sequencing reads from instrument manufactures such as Illumina® can be analyzed by NextGENe software for the detection of SVs. Reads with a high level of mismatch at an end are split into two reads and remapped to the DNA/RNA reference. These Link Reads identify the breakpoints of structural variations including gene fusions.

NextGENe software includes templates to easily analyze gene fusion samples, streamlining the analysis of many SVs.

# Procedure

NextGENe software has many integrated features, including tools to import references, to merge paired reads into a single long read, to batch process samples, and to view the results.

### Import Reference

For detection of structural variations with samples derived from RNA, NextGENe can align to a reference that contains both DNA and RNA. Open the NextGENe Tools menu and select the Reference and Track Manager to import references such as Human_v37.p13+mRNA and Human_v38.p12+mRNA.

### Overlap Merger Tool

The Overlap Merger tool can be used for Illumina Paired End files to generate a single read from the two pairs (see Figure 2).

- Open the Overlap Merger tool in the NextGENe "Tools" menu
- Add the R1.fastq and R2.fastq files into the Overlap Merger's Input field
- Select the Merge overlapping paired reads option (Overlap min bases 10)
- Select Keep 5' end
- Select Merge length of 30-300bp
- Click OK

The Overlap Merger opens a message when the process is complete. The resultant files with a PairMerge.fasta suffix are ready to be aligned.
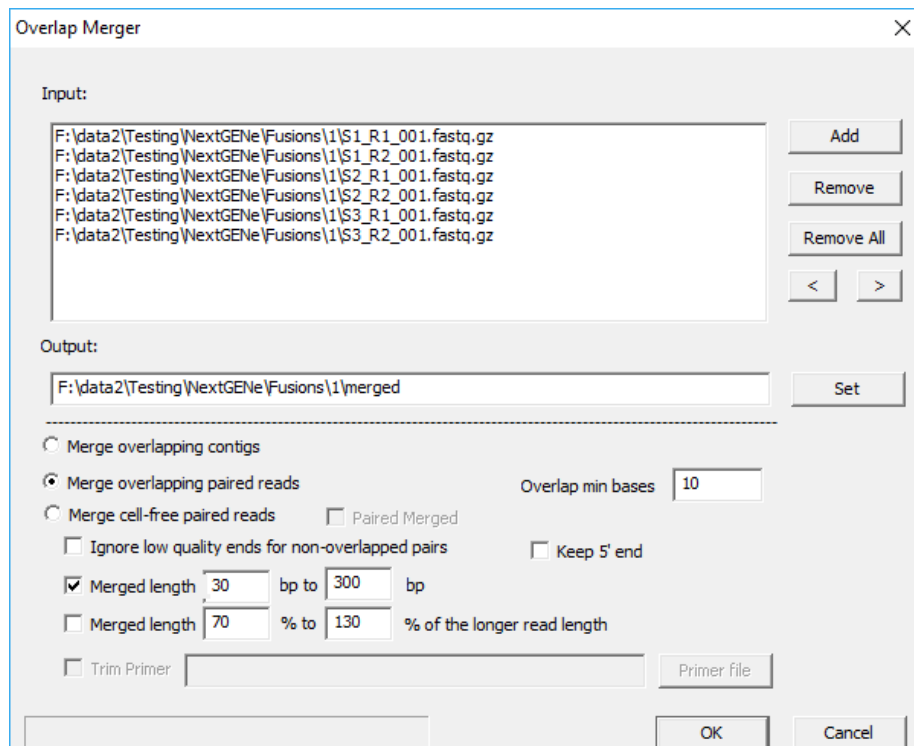


*Figure 2: Overlap Merger can be used to combine paired end reads with these settings.*

SOFTGENETICS®
Software PowerTools for Genetic Analysis

NextGENe®
2nd Generation Sequence Analysis Software

### NextGENe AutoRun Tool

A preinstalled template is included in NextGENe software for structural variation detection. Panel templates are used to specify all settings for the analysis when creating ngjob files (see Figure 3). Any setting(s) can be adjusted to create a new template. These ngjob files are detected and processed by the AutoRun tool.

- Open the AutoRun tool from the NextGENe "Tools" menu.

- Open the "Job File Editor" from the AutoRun "Tool" menu.

- Select a template from the drop-down list. A template for the structural variation detection is included with the software. The "Save As" button allows the templates to be modified into new templates or for completely custom templates to be created.

- Load the PairMerge.fasta output from the Overlap Merger Tool.

- Select one of the Human+mRNA preloaded references to use for alignment.

- Set an output location.

- Click the "Group Jobs" button to open the grouping tool. This makes it easy to split the list of raw data files into separate sample projects.

- Click OK to save the ngjob file.

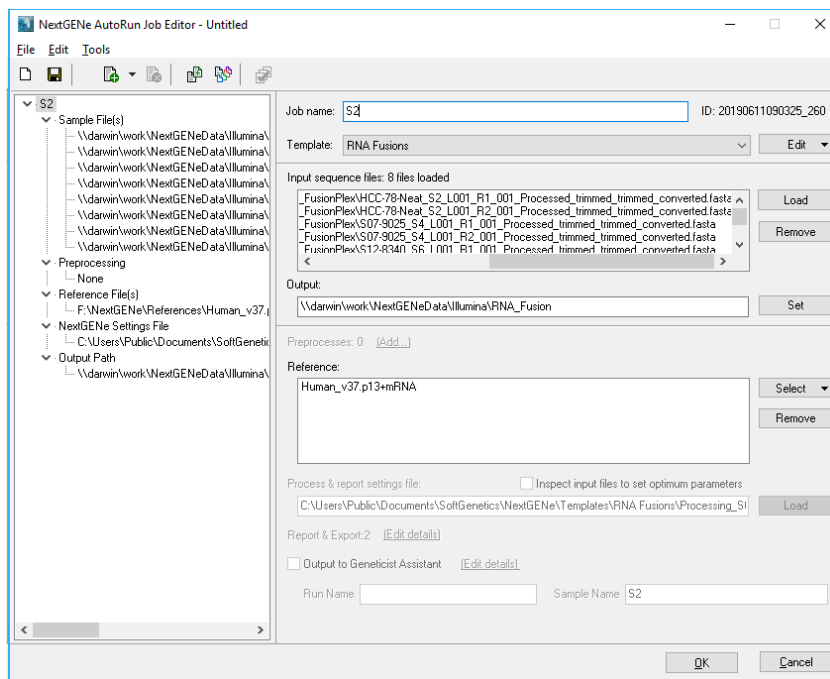- Start the AutoRun tool to begin processing.



*Figure 3: AutoRun can be used to save panel templates and batch process samples.*

## Results

The Seraseq Fusion RNA Mix v2 from SeraCare Life Sciences, Inc. includes many gene fusions, like TPM3-NTRK1, LMNA-NTRK1, SLC45A3-BRAF, EML4-ALK, PAX8-PPARG, FGFR3-TACC3, FGFR3-BAIAP2L1, SLC34A2-ROS1, CD74-ROS1, KIF5B-RET, NCOA4-RET, and ETV6-NTRK3.

A sample was run on an Illumina MiSeq instrument and produced about 255K 2x150bps reads. Trimming of adapters yielded reads averaging 112bp. After utilizing the Overlap Merger tool, 90.5% of the pairs were merged into single reads with an average length of 140bp, extending the average length by about 30bp, (see Figure 4).

NextGENe software was run on a standard Windows Desktop computer. 223K (96.4%) of these reads are aligned to the DNA/RNA reference genome in 3 minutes and Structural Variation detection is completed in 8 additional minutes by NextGENe software. Reports in TXT, VCF and PDF formats are automatically generated.

The final project can be opened in the NextGENe Viewer software, which offers an interactive tool for viewing the SVs and producing additional reports (see Figure 1). Reads that are split and mapped across breakpoints are called Link Reads, and about 2K of these reads are mapping to all 12 expected fusion events in the SeraCare fusion controls.
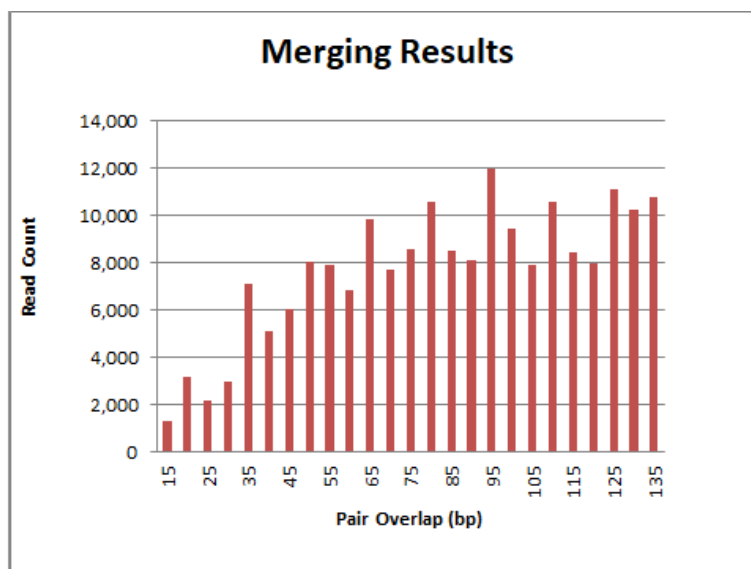


*Figure 4: Overlap Merger converted over 200K paired end reads into single reads averaging 140bp. The original paired reads averaged 112bp, and about 2000 pairs had only 15 bases of overlap with each other, roughly doubling the read length.:*

## Discussion

One feature of the NextGENe Overlap Merger tool can create a single read from paired end reads. Depending on the template size, this could increase 2x150bp reads to a single 290bp read. This step utilizes the quality information within the FASTQ files to assemble the pair with the highest base quality. Longer reads can span more exons, improving the accuracy structural variation breakpoint detection. NextGENe software contains a powerful tool for the analysis of structural variations from high throughput sequencers. Reads are aligned to a reference including both DNA and RNA to enhance sensitivity, reads spanning SVs are split into pseudo paired reads to better identify breakpoints.

NextGENe software is widely used for the analysis of samples from Illumina and Ion Torrent instruments and can also be used for applications such as SNP/Indel discovery, including somatic detection and rare disease research, Copy Number Variation detection, RNA-Seq, de novo assembly and more.

## References

1. An integrated map of structural variation in 2,504 human genomes. **Sudmant, Peter H., Eichler, Evan E. and Korbel, Jan O**. s.l. : Nature, 2015, Nature, Vol. 526, pp. 75-81. 75-81.

2. A comprehensive transcriptional portrait of human cancer cell lines. **Klijn, Christiaan and Zhang, Zemin**. s.l. : Nature Biotechnology, 2014, Vol. 33. 306-312.

3. Implications of human genome architecture for rearrangement-based disorders: the genomic basis of disease. **Shaw, Christine J. and Lupski, James R.** s.l. : Human Molecular Genetics, 2004, Vol. 13. 57–64.

4. Genome structural variation discovery and genotyping. **Alkan, Can, Coe, Bradley P. and Eichler, Evan E.** s.l. : Nature Reviews Genetics, 2011, Vol. 12. 363-376.

Trademarks are property of their respective owners.

**SOFTGENETICS**®
Software PowerTools for Genetic Analysis

**NextGENe**®
2nd Generation Sequence Analysis Software