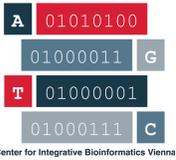


# NextGenMap-LR: Highly accurate read mapping of third generation sequencing reads for improved structural variation analysis

PHILIPP RESCHENEDER<sup>1</sup>, FRITZ J. SEDLAZECK<sup>2</sup>, ARNDT VON HAESLER<sup>1,3</sup>, MICHAEL C. SCHATZ<sup>2</sup>

<sup>1</sup> CENTER FOR INTEGRATIVE BIOINFORMATICS VIENNA, MAX F. PERUTZ LABORATORIES, DR.-BOHR-GASSE 9, A-1030 VIENNA, AUSTRIA, <sup>2</sup> DEPARTMENT OF COMPUTER SCIENCE, JOHNS HOPKINS UNIVERSITY, JOHNS HOPKINS UNIVERSITY, BALTIMORE, MD, USA, <sup>3</sup> BIOINFORMATICS AND COMPUTATIONAL BIOLOGY, FACULTY OF COMPUTER SCIENCE, UNIVERSITY OF VIENNA, VIENNA, AUSTRIA

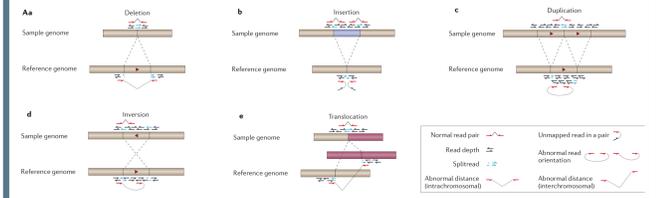


## 1. INTRODUCTION

Characterizing genomic **structural variations (SV)** is vital for understanding the mechanisms behind a wide range of diseases including cancer, autism, and schizophrenia. Nevertheless, due to their complexity they remain harder to detect and less understood than SNPs. Recently, **third-generation sequencing** has proven to be an invaluable tool for detecting SVs. The higher read length allows single reads to span a SV and reliable mapping to repetitive regions of the genome. However, current technologies show a raw read error rate of 10% or more consisting mostly of insertions and deletions. Due to this high error rate

**current mapping methods fail to find exact borders for SVs**, split up large SVs into several small ones, or completely fail to detect certain SVs. Here we present *NextGenMap-LR* for long single molecule PacBio and Nanopore reads which addresses these issues. Using simulated data, verified SVs from healthy human samples and human cancer samples we show how the combination of highly accurate NextGenMap-LR alignments and *Sniffles*, a structural variation caller specifically developed for noisy long-read data, **enables the full characterization of complex SVs** even at low coverage.

## 2. STRUCTURAL VARIATIONS

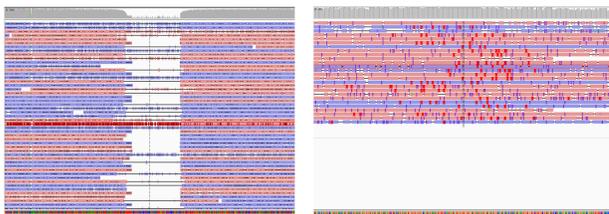


Different types of structural variations (SVs) <sup>a</sup>

<sup>a</sup>Weischenfeldt, J., Symmons, O., Spitz, F., Korbel, J.O. Phenotypic impact of genomic structural variation: Insights from and for human disease (2013) Nature Reviews Genetics, 14 (2), pp. 125-138.

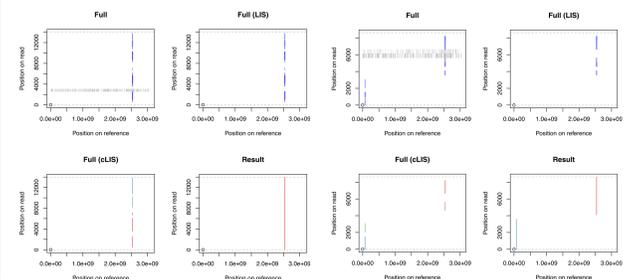
## 3. EXAMPLE ALIGNMENTS

Example regions containing a 300bp deletion (left) and a 200bp insertion (right).



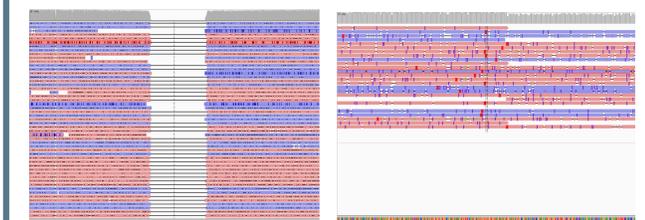
SKBR3 breast cancer reads aligned using BWA-mem with "-x pacbio". Although, BWA-mem in general produces very accurate alignments, larger SVs often cause misalignments.

## 5. CANDIDATE SEARCH



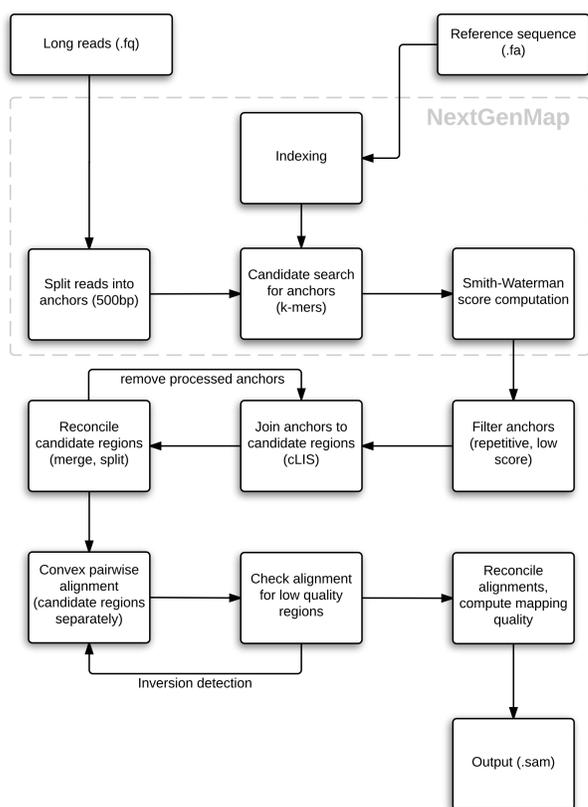
The high-quality anchors retrieved from the initial *k*-mer search are used to determine whether a read spans a large or none linear SV and has to be **split** (right) or can be **aligned contiguously** (left).

## 6. ALIGNMENT STEP



To compute the final alignment(s) we use a Smith-Waterman algorithm. To account for both the **sequencing error** (short and randomly distributed indels) and real **genomic variations** (typically, longer indels), we employ a heuristic non-affine gap model (gap decay) that penalizes gap extensions for longer gaps less than for shorter ones and does not increase the time complexity of the algorithm.

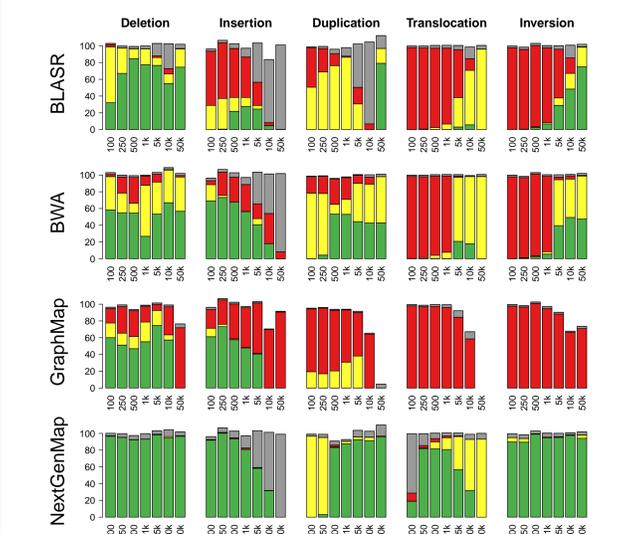
## 4. WORKFLOW



*NextGenMap-LR* comprises four main steps:

1. Identify initial anchors
2. Verify anchors with vectorized Smith-Waterman algorithm (scores only)
3. Filter anchors and find candidate regions for the alignments
4. Compute the full alignment between the read and the respective candidate reference regions

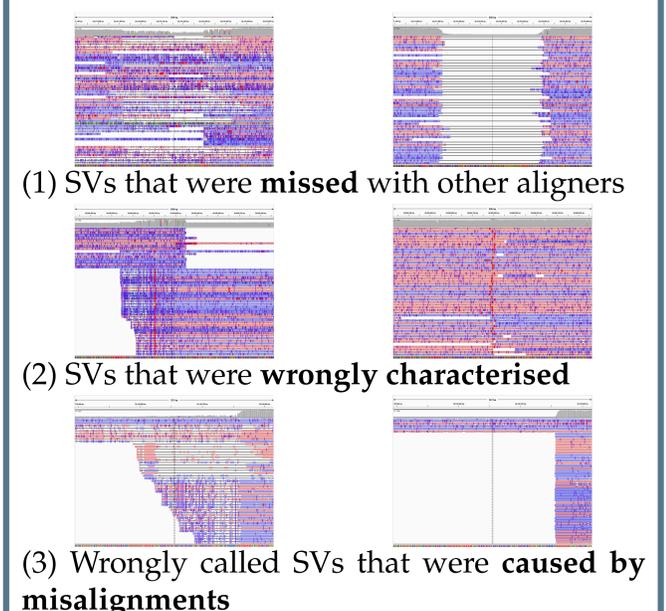
## 7. SIMULATED DATA



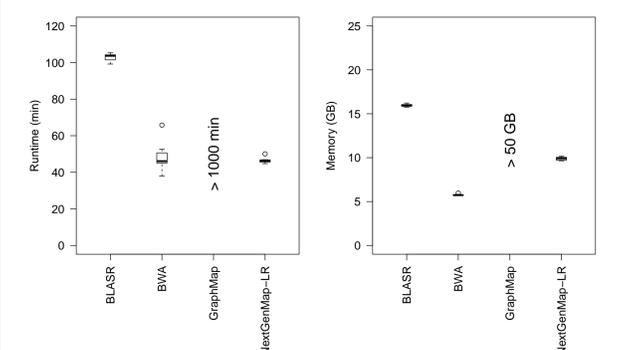
Green: correct SV indicated and correct break points (+/- 10bp) found, Yellow: correct SV indicated, Red: wrong SV indicated, Grey: no SV indicated

## 8. SKBR3 DATA

To evaluate *NextGenMap-LR* we mapped a full SKBR3 and the genome in a bottle PacBio data set and compared the results to BLASR and BWA-MEM alignments:



## 9. RUNTIME & MEMORY



Ten fastq files containing 1 Gbp of read data were aligned using 10 CPU threads.

## 10. SNIFFLES

To accurately detect structural variations in PacBio and Oxford Nanopore data, *NextGenMap-LR* is tuned to work hand in hand with the structural variation caller *Sniffles*. Using *NextGenMap-LR* alignments, *Sniffles* outperforms all other SV analysis approaches in both sensitivity and specificity.

