



Bionano Solve Theory of Operation: Structural Variant Calling

Document Number: 30110

Document Revision: K

Table of Contents

Introduction	6
Structural variant calling with Bionano Solve	6
<i>De novo assembly.....</i>	<i>6</i>
<i>SV calling</i>	<i>7</i>
Insertions and deletions	8
Translocation breakpoints	8
Inversion breakpoints	9
Duplications.....	9
Large SV Types	10
<i>Zygoty classification.....</i>	<i>11</i>
<i>AOH/LOH Detection</i>	<i>12</i>
Input	12
Output	12
Theory	13
<i>Confidence modeling.....</i>	<i>14</i>
<i>Variant Allele Fraction (VAF) calculation</i>	<i>15</i>
<i>Copy number variant (CNV) calling</i>	<i>16</i>
Introduction	16
Input	17
Output	17
Environment setup	18
Usage.....	18
Theory	21
<i>SVMerge</i>	<i>37</i>
<i>Rare Variant Pipeline (RVP).....</i>	<i>37</i>
<i>Variant Allele Fraction (VAF) calculation</i>	<i>38</i>
Structural variant calling performance: Rare Variant Pipeline.....	39
<i>Structural variant calling performance using simulated data</i>	<i>39</i>
Summary	39
Performance for simulated insertions and deletions.....	39
Performance for simulated duplications, translocation breakpoints, and inversion breakpoints	40
<i>Structural variant calling performance using real data.....</i>	<i>42</i>
Structural variant calling performance: de novo assembly pipeline.....	43
<i>Structural variant calling performance using simulated data</i>	<i>43</i>
Performance for simulated insertions and deletions.....	43
Performance for other simulated events.....	44
<i>Structural variant calling performance using real data.....</i>	<i>45</i>
Performance for insertions and deletions	45

Performance for translocation breakpoints ²	46
Performance for inversion breakpoints ²	47
AOH/LOH detection performance	48
Copy number variant calling performance: fractional CN analysis	50
<i>Copy number variant calling performance using simulated data</i>	<i>50</i>
<i>Copy number variant calling performance using real data</i>	<i>51</i>
Copy number variant calling performance: integer CN analysis	53
<i>Copy number variant calling performance using simulated data</i>	<i>53</i>
<i>Copy number variant calling performance using real data</i>	<i>54</i>
Copy number variant calling performance: chromosomal aneuploidy detection	55
<i>Chromosomal aneuploidy detection performance using simulated data</i>	<i>55</i>
<i>Chromosomal aneuploidy detection performance using real data</i>	<i>56</i>
Bionano Access integration	58
VCF conversion	58
Additional Considerations	58
FAQs	60
Appendix A: SV Calling Performance Evaluation	66
<i>RVP SV calling performance at lower coverage levels</i>	<i>66</i>
<i>Methods for assessing SV calling performance</i>	<i>68</i>
SV calling performance with simulated data	68
<i>Simulation of large structural variants</i>	<i>69</i>
Data simulation	69
Analytical Performance	69
<i>Terminal deletions</i>	<i>70</i>
Data simulation	70
Analytical Performance	70
Artifacts of the simulator	71
<i>Analytical Repeatability</i>	<i>71</i>
RVP	71
De novo assembly	72
<i>SV calling performance with real data</i>	<i>72</i>
CHM1/13 analysis	72
CEPH trio analysis	72
<i>Method for assessing CNV calling performance</i>	<i>72</i>
CNV calling performance with simulated data	72
CNV calling performance with real data	73
Appendix B: Mask generation	74
<i>Comparing Solve 3.6.1 and 3.7 masks</i>	<i>74</i>

<i>Compiling common translocation breakpoint list</i>	79
<i>Compiling annotated segmental duplication regions</i>	79
<i>Tools for generating mask BED files</i>	79
Appendix C: Descriptions of output files	81
<i>Output files from the de novo assembly pipeline</i>	81
<i>Output files from RVP</i>	82
Appendix D: In-depth description of the <i>de novo</i> assembly pipeline	85
<i>Setting up a de novo assembly run</i>	85
<i>Key stages of the de novo assembly pipeline</i>	85
<i>Runtime performance of the de novo assembly pipeline</i>	90
<i>Evaluation of de novo assembly results</i>	90
Appendix E: Interpretation of inversion breakpoint calls	91
<i>Introduction</i>	91
<i>SMAP output for inversion breakpoints</i>	92
<i>Inversion breakpoints report file</i>	93
<i>Inversion confidence score</i>	94
<i>FAQs</i>	94
Appendix F: Confidence modeling for inversion and translocation breakpoints	98
<i>Introduction</i>	98
<i>Datasets</i>	98
<i>Models for the scores</i>	98
Model for the translocation breakpoint scores	98
Model for inversion breakpoint scores	101
<i>Validation</i>	103
<i>Simulated datasets</i>	103
<i>Cytogenetics datasets</i>	104
<i>References</i>	105
Appendix G: VCF conversion	106
Appendix H: Variant Allele Fraction calculation	110
<i>Method</i>	110
<i>Performance on simulated data</i>	111
<i>Reproducibility</i>	113
<i>Performance on experimental data</i>	114
<i>Segmentation of the whole genome VAF plot</i>	114
<i>References</i>	115
Appendix I – Custom CNV control data	116
Technical Assistance	122

Legal Notice

For Research Use Only. Not for use in diagnostic procedures.

This material is protected by United States Copyright Law and International Treaties. Unauthorized use of this material is prohibited. No part of the publication may be copied, reproduced, distributed, translated, reverse-engineered or transmitted in any form or by any media, or by any means, whether now known or unknown, without the express prior permission in writing from Bionano Genomics. Copying, under the law, includes translating into another language or format. The technical data contained herein is intended for ultimate destinations permitted by U.S. law. Diversion contrary to U. S. law prohibited. This publication represents the latest information available at the time of release. Due to continuous efforts to improve the product, technical changes may occur that are not reflected in this document. Bionano Genomics reserves the right to make changes in specifications and other information contained in this publication at any time and without prior notice. Please contact Bionano Genomics Customer Support for the latest information.

BIONANO GENOMICS DISCLAIMS ALL WARRANTIES WITH RESPECT TO THIS DOCUMENT, EXPRESSED OR IMPLIED, INCLUDING BUT NOT LIMITED TO THOSE OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. TO THE FULLEST EXTENT ALLOWED BY LAW, IN NO EVENT SHALL BIONANO GENOMICS BE LIABLE, WHETHER IN CONTRACT, TORT, WARRANTY, OR UNDER ANY STATUTE OR ON ANY OTHER BASIS FOR SPECIAL, INCIDENTAL, INDIRECT, PUNITIVE, MULTIPLE OR CONSEQUENTIAL DAMAGES IN CONNECTION WITH OR ARISING FROM THIS DOCUMENT, INCLUDING BUT NOT LIMITED TO THE USE THEREOF, WHETHER OR NOT FORESEEABLE AND WHETHER OR NOT BIONANO GENOMICS IS ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Patents

Products of Bionano Genomics® may be covered by one or more U.S. or foreign patents.

Trademarks

The Bionano Genomics logo and names of Bionano Genomics products or services are registered trademarks or trademarks owned by Bionano Genomics in the United States and certain other countries.

Bionano Genomics®, Irys®, IrysView®, IrysChip®, IrysPrep®, IrysSolve®, Saphyr®, Saphyr Chip®, Bionano Access®, and Bionano EnFocus™ are trademarks of Bionano Genomics, Inc. All other trademarks are the sole property of their respective owners.

No license to use any trademarks of Bionano Genomics is given or implied. Users are not permitted to use these trademarks without the prior written consent of Bionano Genomics. The use of these trademarks or any other materials, except as permitted herein, is expressly prohibited and may be in violation of federal or other applicable laws.

© Copyright 2021 Bionano Genomics, Inc. All rights reserved.

Revision History

Revision	Notes
K	<ul style="list-style-type: none">Added timing for analysis on Bionano Access serverAdded text on limitations of mosaicism simulation

Introduction

Structural variation (SV) is a common source of sequence variation, involving deletion, insertion, and rearrangement of genomic material. It has the potential to impact large stretches of sequence, disrupting genes and regulatory elements. Structural variants are associated with genetic disorders and are used as disease markers in clinical diagnosis of diseases such as DiGeorge syndrome and cancer.

Bionano Solve contains a suite of tools that analyze raw single-molecule data and generate fully annotated SVs. Bionano Solve features a haplotype-aware assembler specifically designed to detect and differentiate allelic differences, thus enabling effective and comprehensive detection of homozygous and heterozygous SVs. It also includes Rare Variant Pipeline (RVP), which has high sensitivity for low variant allele frequency (VAF) SVs. All major SV types are supported, and extensive validation based on simulated and experimental data showed high detection performance.

Large copy number variants are detected based on coverage depth information using a copy number analysis pipeline embedded in the *de novo* assembly pipeline and RVP. The copy number analysis pipeline can detect fractional copy number changes and chromosomal aneuploidy events. Previously released separately, the Introduction to Copy Number Analysis document (PN 30210) has now been consolidated into this document.

Variant Annotation Pipeline (VAP) annotates the SV calls and provides information key to understanding relevance to a biological question or phenotype of interest (see Bionano Solve Theory of Operation: Variant Annotation, PN 30190). It supports analysis of paired datasets (tumor-control, for example) and trio datasets (proband and two parents, for example).

Bionano Solve is fully integrated with Bionano Access®, which provides a user-friendly interface for streamlining analysis. Bionano Access is useful for managing projects, analyzing run results, and visualizing data.

Structural variant calling with Bionano Solve

De novo assembly

Bionano's *de novo* assembly algorithm is built on the overlap-layout-consensus strategy with a maximum likelihood model for scoring alignments. For more information, please refer to Appendix D. Following pairwise

alignment of the input single-molecule maps, an overlap graph is constructed. Spurious edges are removed and redundant edges collapsed. The assembler outputs the longest paths in the graph and constructs a set of draft consensus maps. The consensus maps are further refined, extended, and merged.

In order to optimize assembly of non-homozygous variants, during the extension stages of the assembly, the pipeline analyzes molecule-to-genome map alignments, identifies clusters of molecules with coordinated disrupted alignment, and assembles these clusters separately. This is critical for assembling haplotype maps with large differences and for detecting a wide range of variants. Additionally, we implemented haplotype-aware components to optimize assembly of heterozygous variants. In the last refinement stage, molecules are aligned to a given genome map and clustered into two alleles. The allele-specific molecules are used to generate the final set of allele-differentiated consensus genome maps.

During assembly, large non-unique regions (which we call complex multi-path regions, or CMPRs) in the genome maps are recognized and marked. These regions, often associated with large segmental duplications in the genome, create ambiguity in the assembly graph and are prone to mis-assembly. They are detected in a *de novo* fashion - assembled maps are aligned with each other, and maps that share significant stretches of sequence but are otherwise divergent are identified. If the CMPR, or the shared sequence, is at least 140 kbp, and if the option to split CMPRs is enabled, the maps are split (Figure 1) in order to avoid mis-assembly. The labels encompassing the CMPR would be marked in the Mask column in the CMAP output (see CMAP File Format Specification Sheet; PN 30039) and highlighted in Bionano Access. Users are also provided the option to not split these maps; the options are presented in Access when an assembly run is set up.

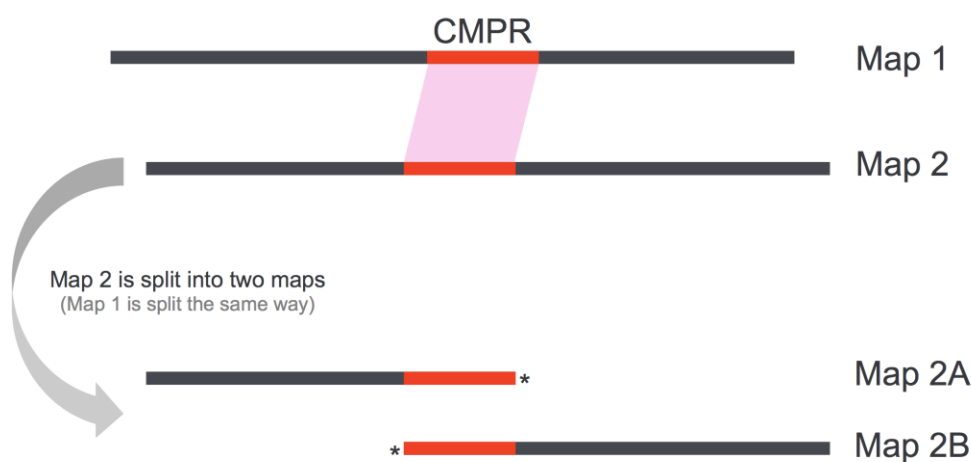


Figure 1. Schematic of how maps that share a CMPR (highlighted in red) are split. The sequences flanking the CMPR are divergent between Map 1 and Map 2. The resulting Map 2A and Map 2B each contains a copy of the CMPR. CMPR ends are marked (*) in figure) such that the maps would not be merged in subsequent steps.

SV calling

SV calls are obtained by aligning consensus genome maps to a reference using a Multiple Local Alignment algorithm and analyzing the alignments for SV signatures. Pairs of alignments within a map are analyzed and inconsistencies representing possible SV events between the genome maps and the reference are identified.

Insertions and deletions

An alignment outlier is defined by two well-aligned regions that flank a poorly aligned or unaligned region. An outlier is identified as a deletion if the reference range in the outlier region is larger than the corresponding range on the map, and an insertion if the converse is true. Sequence substitutions may also create alignment outliers without significant net gain or loss of sequence. In those cases, insertion and deletion calls (often small) may be called, even though they may not be simple insertions or deletions.

Translocation breakpoints

A fusion point between distant regions of the genome is identified as a translocation breakpoint. Intrachromosomal fusion breakpoints involve regions typically (but not always) at least 5 Mbp away from each other *on the same chromosome*. Interchromosomal translocation breakpoints involve regions *on different chromosomes*. Translocation breakpoint orientation is denoted in the “orientation” column in the SMAP output (Figure 2). Reciprocal translocations are expected to be detected as separate translocation breakpoints with opposite orientations. The pipelines do not explicitly pair translocation breakpoints. Also, transpositions of small pieces of sequence (typically less than 50 kbp) may be detected as insertions, if the transposed sequence cannot be confidently aligned to the reference.



Figure 2. Schematic of how translocation breakpoint orientation is defined in SMAP. There are four possible scenarios. Each translocation breakpoint is annotated in one of the four ways: “+/+”, “+/-”, “-/+”, and “-/-”.

Inversion breakpoints

Inversion breakpoint calls involve neighboring alignments with opposite orientations. Please refer to Appendix E for more information interpreting inversion breakpoint calls. Small inversions (whose inverted regions contained fewer than five labels) are identified by searching in a limited space for potential inverted alignments. They may be spanned by single genome maps and represented in sets of two linked, paired inversion breakpoint entries (“inversion_paired”), specifying eight coordinates of interest. For other inversion breakpoint calls, there are two linked SMAP entries (“inversion” and “inversion_partial”), specifying six coordinates of interest (Figure 3). If the inverted alignment overlaps partially with the neighboring non-inverted alignment, the inverted alignment is trimmed. Inversions with a reference gap larger than 5 Mbp are called as intra-chromosomal fusion breakpoints, provided there was no query alignment overlap of 140kb or more. Starting with Solve3.7, insertions or deletions located inside the inversion breakpoint are located at the correction breakpoint interval : in previous release the deletion or insertion would appear to span the entire inversion region including both breakpoints.

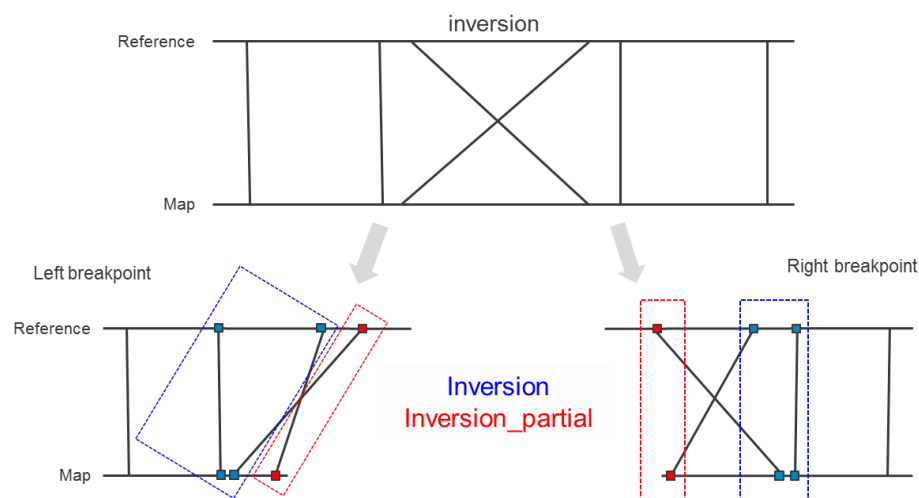


Figure 3. Schematic of how inversion breakpoint calls are output in SMAP. “Inversion” (highlighted in blue) and “inversion_partial” (highlighted in red) entries are linked, and together, they encode six coordinates of interest.

Duplications

Duplications are detected based on direct or indirect evidence of duplication within single maps. There is direct evidence when two stretches of sequence on the (consensus) map align to the same stretch of sequence on the reference, implying that there are (at least) two copies of the same reference sequence on the map (Figure 4). They are output as “duplication” events. However, sufficiently large duplications need to be detected differently. In those cases, the map may not span whole copies of the duplicated sequence. Based on the alignment signature, one could still infer duplication events, even though two full duplicated copies may not be observed on the map, and the alignments may not overlap on the reference (Figure 5). They are output as “duplication_split” events. The duplication detection algorithm supports tandem duplications and non-tandem duplications whose copies are sufficiently close to each other on a given map OR both copies are fully spanned by the same map (NOT a

duplication_split). The previous descriptions are relevant for duplications where both copies of the duplicated sequence have the same orientation; they are named differently to highlight the distinction between the detection methods. Inverted duplications are called if the two alignments on the map are in opposite orientations. Size cannot be determined for inverted duplications unless the map fully spans both copies.

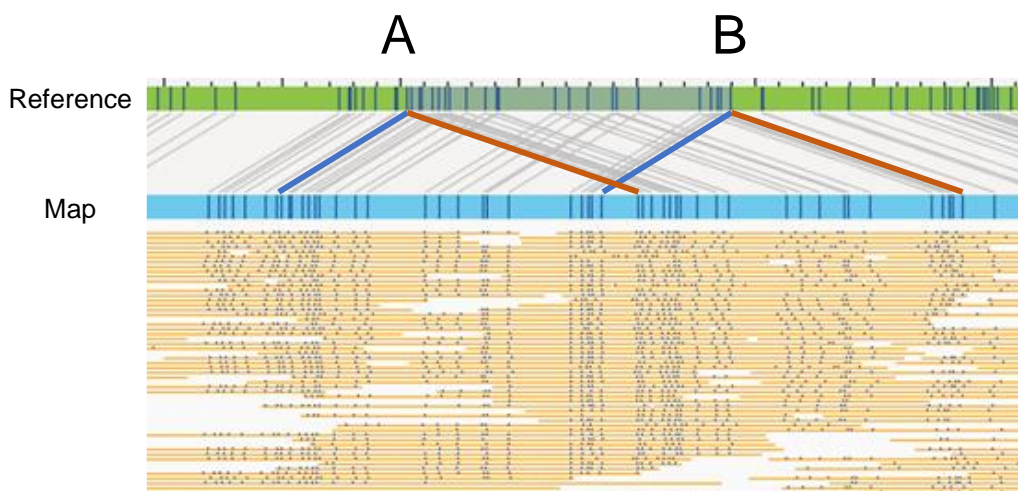


Figure 4. Example of a “duplication” event. Two regions in the map (one shaded in blue, and one in red) align to the same region on the reference (the region between points A and B). This suggests that there are two copies of the sequence in the map, and that they were next to each other in the same orientation.

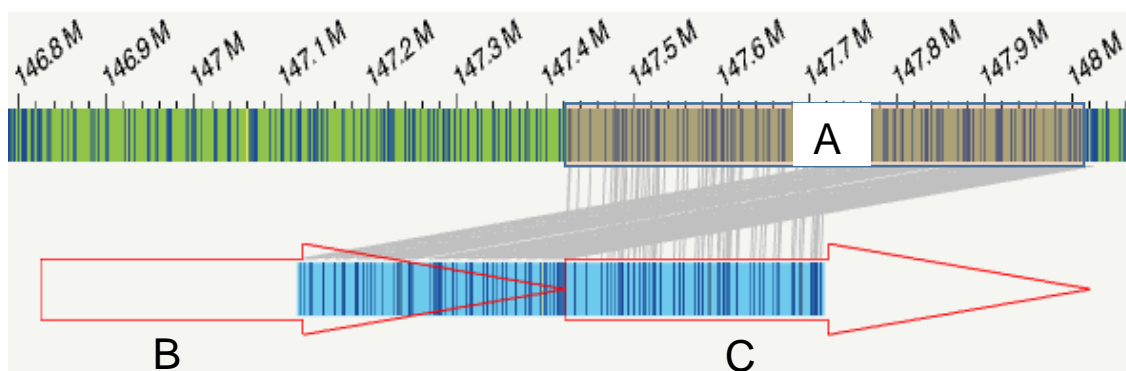


Figure 5. Example of a “duplication_split” event. In this case there is no overlap in the alignments on the map. However, based on the alignment, segment B on the map, represented by the red arrow on the left, is fused with segment C represented by the arrow on the right. The parsimonious interpretation would be that there is a direct tandem duplication of the sequence in region A. The direction of the arrows show the orientation of the segments. The outline areas of the arrows show the portion of segments B and C that are inferred and do not appear in the assembled map. While the evidence is indirect and the event inferred, one may confirm such signal using the copy number analysis data (for large duplications).

Large SV Types

Solve 3.7 uses updated sizing criteria to assign more probable types to large SVs > 200 kbp. Please see Table 1 for details of the size conditions and the alternative assigned types.

Type	Criteria	Type Called pre-Solve 3.7	Type Called in Solve 3.7
Duplication	Reference distance between fusion points plus the size of duplication > 5Mbp	Intra-chromosomal fusion	Intra-chromosomal fusion
Duplication	Reference distance between fusion points is greater than 300 Kbp and less than 5Mbp	Insertion	Duplication
Deletion	Deletion size > 5Mbp	Deletion	Intra-chromosomal fusion
Inversion	Reference distance between fusion points > 5Mb	Inversion Intra-chromosomal fusion	Inversion
Intra-chromosomal fusion	Reference distance between fusion points < 1Mb	Insertion/Deletion	Intra-chromosomal fusion

Table 1. Summary of type assignment for large SVs

Zygosity classification

Zygosity is a classification of a SV call as homozygous, heterozygous, or unknown. It is currently assigned to only insertion, deletion, translocation breakpoint, and inversion breakpoint calls. Conceptually, we try to determine whether there is presence of a reference allele or additional SV alleles. To do that, we look at two criteria: 1) whether the call in question overlaps another call (which might be from a different SV allele), and 2) whether the call in question overlaps alignment of another genome map that shows no SV (which might be from the reference allele).

In the case of overlap with another call, there is a determination of whether the two calls are likely to represent the same allele. An SV call is categorized as homozygous if there is no overlapping alignment or the same SV is called on another genome map. An SV call is categorized as heterozygous if there is overlapping alignment or a different SV is called on another genome map. If both another alignment and a different SV or multiple different SVs are present at the same location, zygosity is set as heterozygous. For the purpose of determining if SV calls are the same, different sets of criteria are applied for different SV types. For insertion and deletion calls, they must have overlapping ranges and at least 80% overlap and size similarity. For translocation breakpoint calls, the breakpoints must be within 50 kbp and on the same strand and chromosomes. For inversions breakpoints, the breakpoints must be within 50 kbp and on the same strand and chromosomes. The pipeline does not require exact breakpoint and size matches, in order to account for slight differences.

Zygosity is not currently assigned to duplication calls. Also, zygosity is not an output for calls from the Rare Variant Pipeline.

AOH/LOH Detection

Detecting regions of absence/loss of heterozygosity (AOH/LOH) was introduced as part of Bionano Solve 3.7 for human samples. Absence of heterozygosity refers to a specific type of genetic mutation that can be a result of uniparental disomy or consanguinity, in which there is an absence of one normal copy of a gene, and which may result in increased susceptibility to recessive disease. Regions of homozygosity are identified by a consistent decrease in heterozygous SV calls across a genomic region in the case sample compared to the level observed genome-wide in controls. SV zygosity is determined as part of the *de novo* assembly pipeline; therefore, AOH detection is only available for the *de novo* assembly pipeline.

Input

Required input to the AOH/LOH detection pipeline is an SMAP file with SV calling results that includes zygosity calls. Additionally, an output file from the CNV calling pipeline, `cnv_chr_stats.txt`, is used to determine sex; this information is used to suppress outputting AOH calls on chr X for males. Other necessary inputs pertaining to human datasets are built into the package.

Output

As part of the standard *de novo* assembly pipeline output, the results from the AOH detection pipeline are stored in `contigs/exp_refineFinal1_sv/merged_smaps/loh/`. There are two expected output files in the directory. Selected results files are imported by Bionano Access; the data are plotted in the AOH/LOH track below the CN track in the map-to-reference alignment view and in the whole-genome CN view. In these views and in the Circos view, AOH/LOH calls are denoted as yellow highlighted regions in the SV track. Details about the calls are listed in the AOH/LOH Regions tab, where users can sort and filter the calls.

loh_calls.txt contains information about the AOH calls; specifically the start and end positions and a confidence score. The confidence score is based on the size of the call, and it represents the model's precision for calling AOH events in the given size range bin. The confidence score is directly correlated with the size of the call. Note that AOH/LOH calls with a size under 25 Mbp are not output, as the model has low precision in this size range, unless the call is near another call (within 25 Mbp), since a true AOH/LOH region may be called as multiple smaller regions with a small gap between. Further size filtering can be applied in Bionano Access.

loh_per_sv_info.txt contains details about the subset of SVs from the input SMAP file that were used in AOH/LOH detection, with a few additional columns containing information about the SVs that were obtained from the HMM, such as the probability that the SV is in a AOH/LOH region, as calculated by a Hidden Markov Model (details in the following section). The SVs listed in this file are those that remained after variant filtering and were the SVs used in AOH/LOH calling. Only these SVs will have a corresponding orange dot in the Bionano Access AOH/LOH track denoting the probability that the SV is in an AOH/LOH region.

Please refer to Absence/Loss of Heterozygosity Pipeline File Format Specification Sheet (PN 30460) for additional details about these two files.

Theory

AOH regions are determined using a Hidden Markov Model (HMM), in which the hidden state is whether an SV belongs to an AOH or background region, and the observable states are whether the SV was called as homozygous or heterozygous. Model parameters were estimated by fitting the model to a simulated dataset that was generated by splicing together SMAP files from Bionano control samples and two haploid samples. Haploid samples are expected to have only homozygous SVs, with the exception of false positive heterozygous SV calls, and are the closest approximation to how an AOH region would appear in a real sample. The fact that males have AOH regions across chromosome X, with the exception of the pseudoautosomal region, was accounted for and also used in simulating AOH events and training the HMM.

Additional variant filtering was applied in order to enrich the data for informative sites. The following SVs were filtered out: confidence score below 0.95; unknown zygosity; interchromosomal SVs; duplicate SVs where the start and end position were identical. Finally, SVs were filtered out if they were known commonly homozygous sites, where at least 95% of Bionano controls had homozygous SV calls in the same 10 kbp region. If a bin contained a known common homozygous SV, all SVs in the same bin were filtered out in the query sample.

Performance was evaluated using samples with known AOH/LOH events, and using additional simulated data. It was observed that large AOH regions may be called as multiple smaller AOH calls with small gaps between, if there are false positive heterozygous SV calls in the region. On the other hand, many small AOH calls appear to be false positive, and these tend to stand alone and are often the only AOH call on a chromosome. Therefore, the following size-based filtering is applied to calls that are output: AOH calls under 25 Mbps are filtered out if they are at least 25 Mbps away from another AOH call. This was found to filter out approximately 90% of false positive calls while retaining true AOH calls.

Known issues:

- In the whole genome plot, we have observed cases where the AOH algorithm results and the segments appearing in the VAF plot do not agree. If this occurs, the AOH call should be prioritized. Please see the Variant Allele Fraction (VAF) calculation section for details.
- Homozygous SVs that were filtered are displayed in the AOH/LOH track in Bionano Access as green dots without a corresponding orange dot representing the confidence score. AOH/LOH calls that have a large number of homozygous calls without accompanying confidence scores may be indicative of a false positive.

Confidence modeling

The raw SV calls are made based solely on interpretation of the map alignments; however, other factors also impact whether the SVs are truly present in a sample. For example, SV calls may be less reliable in complex regions, where there may be conflicting alignments. The SV confidence scores are intended to provide a measure of “confidence” of the SV calls and a way for users to sort, filter, or prioritize the calls.

The SV confidence scores are computed based on models that were trained using simulated, isolated, non-complex SVs and detected SVs from real samples where orthogonal data were available. Real SVs in a genome tend to cluster and are found in certain structurally complex regions. They are difficult to simulate, and the simulation data do not capture the full range of complexities observed in complex regions and in real samples. Thus, the confidence scores may not be as informative and reliable for SV calls in those regions. As discussed in the following sections, the default confidence cutoffs in Bionano Access serve to provide a starting point. When selecting the recommended thresholds, we tried to balance sensitivity and the number of false positive calls. We looked at each SV type using simulated data and real data (for which we had orthogonal SV data) for this process, and we also considered feedback from internal teams and the field. Users may experiment with other cutoffs and vary the stringency based on project needs. For example, adjusting the cutoff upward should increase the fraction of true calls, but other lower confidence true calls may be filtered out as a result.

Users are encouraged to use confidence scores as one of the possible tools for filtering SV calls. The confidence scores may be used in conjunction with the SV masks and variant annotations.

Confidence for insertions and deletions

The confidence of an insertion or deletion call, which ranges from 0 to 1, reflects an estimate of the probability of the call being a true positive, or positive predictive value (PPV). These estimates are derived based on simulation studies, where we simulated insertions and deletions in the genome and assessed the resulting SV calls. It takes into account the SV size, the non-normalized p-value (\log_{10}) of the two well-aligned regions, and the non-normalized log-likelihood ratio of the poorly aligned or unaligned region. The genomic context is not explicitly taken into account; confidence scores of calls around complex regions may be less reliable. The confidence cutoff in Bionano Access is defaulted to 0 for insertions and deletion calls. Calls with lower confidence may be of interest in a discovery setting; however, users may increase the stringency based on project needs.

Confidence for insertion and deletion calls <500 bp in size is undefined and set to -1.

Confidence for inversion and translocation breakpoints

New confidence models for inversion and translocation breakpoints were introduced in the Bionano Solve 3.6 release. They were trained using a machine learning approach based on a combination of simulated and real human data. A different set of models is applied for non-human datasets. For Bionano Solve 3.7, the

translocations model has been re-trained and a new minimum confidence threshold set. Users are encouraged to develop custom score cutoffs for the non-human models.

The recommended score cutoffs for the human models are:

- Intrachromosomal fusion breakpoints: 0.05
- Interchromosomal translocations breakpoints: 0.05
- Inversions breakpoints: 0.7

Further discussion on the methods and the validation results is included in Appendix F.

Confidence for duplications

Currently, confidence for duplications called by the SV calling pipeline is undefined and set to -1. Inspection of duplication calls is encouraged.

Variant Allele Fraction (VAF) calculation

In Bionano Solve 3.7 the variant allele fraction (VAF) is provided for each SV detected by the *de novo* assembly pipeline. Its value ranges from greater than 0 to 1 and records the percentage of molecules calling the structural variant out of the total coverage in the region. In a normal diploid sample, a homozygous SV will have a VAF approaching 1 while a heterozygous variant will have a VAF of approximately 0.5. A value of -1 is provided when the algorithm could not find a value or there was a processing error. The values are present by default in the output SMAP under the VAF column. A parameter is available in the command line tool to disable the computation.

The VAF calculation uses the coverage of the consensus map segments and consists of two stages. The first stage clusters the SV calls in the SMAP file to find equivalent ones (those that are the same SV event detected by two or more different consensus maps). In the second stage the coverages for the maps calling the same SV are aggregated using the alignments of the consensus maps to the reference as a guide. The aggregated coverage is used as input to a Bayesian inference that determines the probability of having a VAF value α_k for an allele k given the coverage for a set of labels D in the SV region, and the genotype G : $P(\alpha_k | D, G)$. The value that maximizes the probability is chosen as the VAF. For duplications, it is not possible to apply Bayesian inference and the VAF is simply the quotient between the SV coverage and the total coverage at the genomic locus.

We also provide a segmentation algorithm to help distinguish changes in the VAF pattern across the genome, which is useful to detect aneuploidies. Further details are provided in Appendix H.

Known issues:

- In the whole genome plot, we have observed cases where the AOH algorithm results and the segments

appearing in the VAF plot do not agree (e.g. AOH is called and VAF show variants with allele fractions <1). This may occur because the analyses are performed independently using different rules for which variants to include in the calculation. For AOH detection only variants with very high confidence scores are considered, while for the VAF segmentation all variants are employed. This is done as a trade-off to improve genome-wide segmentation accuracy in samples at the potential cost of discordant results in AOH regions. When observing a possible conflict between a AOH region and a VAF segment line created by some heterozygous variants, we recommend to guide your analysis based on the AOH results.

Copy number variant (CNV) calling

Introduction

Copy number (CN) analysis is performed as part of the *de novo* assembly pipeline and Rare Variant Pipeline (RVP). The CN analysis tool analyzes an input molecule-to-reference alignment, normalizes this raw molecule coverage profile using control data, and segments the genome (inferred local CN states across the genome) based on detected changes in the underlying copy number state. CNV calls are output and annotated with confidence scores. Chromosomal aneuploidy events are detected after post-processing the initial CNV calls; however, whole-genome aneuploidy events cannot be detected.

The CN analysis tool has two main components – a fractional CN analysis module and an integer CN analysis module. Released as part of Bionano Solve 3.4, the fractional CN module is optimized for detecting events in genomes with multiple CN state changes and events at lower allelic frequencies (AF)¹. These events are frequently observed in highly heterogeneous genomes such as cancer. The fractional CN pipeline is intended for DLE-1 datasets. The integer CN module was released as part of Bionano Solve 3.2.1 and is intended for human Nt.BspQI and Nb.BssSI datasets. It is intended for homogenous genomes. The fractional and integer modules are integrated into the pipeline; there is not an explicit switch to use one or the other. The pipeline automatically determines which one to use based on the input data and available auxiliary data. Generally, DLE-1 datasets are expected to go through the fractional module.

Necessary auxiliary input for the standard human reference builds (hg19 and hg38) and the latest mouse build (mm10) is packaged with the tool. With user-provided data, the pipelines can also support additional reference and enzyme combinations.

The tool was implemented in the R programming language. It can be run on the command line as a standalone tool, or as part of a *de novo* assembly or RVP run. The resulting compressed output is compatible with and can be visualized in Bionano Access after import². The normalized coverage profile is plotted in Bionano Access, and

¹ Allele frequency (AF) is proportional to the fraction of chromosomes in a sample that carry a given events.

² Results from a standalone, command-line CN analysis run may not be directly imported into Access at this time.

the CNV calls can be sorted and filtered in the Copy Number tab.

Input

The *de novo* assembly pipeline and RVP automatically generate and use the necessary input for the CN analysis for supported input datasets. A reference needs to be specified at the start of the pipelines such that the molecule-to-reference alignment is generated.

Instead of using default data packaged with the tool, users could provide custom control data. For command-line usage, the minimum required input in the default mode includes the r.cmap and .xmap output from the molecule-to-reference alignment (alignmolvref) stage of the assembly or RVP, a name for the sample, and an output path for the results.

Output

As part of the standard pipeline output, the results from the CN analysis tool are stored in alignmolvref/copynumber/. There are several expected output files in the directory. Selected results files are imported by Bionano Access; the CN data are plotted in the CN track above the reference map in the map-to-reference alignment view, in the Circos view, and in the whole-genome CN view. Details about the CNV calls are listed in the Copy Number tab, where users can sort and filter the calls. Aneuploidy calls are in the Aneuploidy tab.

1. **cnv_rcmap_exp.txt** contains per-label coverage information. The format of this file is similar to the standard CMAP format, but with several additional columns, the definitions of which depend on which pipeline is run (see table below).

Column name	Integer CN pipeline	Fractional CN pipeline
ScaledCoverage	Sample coverage divided mean coverage for each chromosome	Sample coverage divided by average control coverage
NormalizedCoverage	ScaledCoverage normalized by SVD	Same as above (to maintain consistent format)
CopyNumber	Rounded copy number states	Rounded copy number states
fractionalCopyNumber	Smooth copy number states	Fractional copy number states
MeanCov	NA	Average coverage

The fractionalCopyNumber and CopyNumber columns are the same as those in cnv_calls_exp.txt.

ScaledCoverage and NormalizedCoverage are considered intermediate results (see Theory section below for detail). The CN-related columns reflect local changes in the CN state. For human, the normal diploid CN state

is 2. CN states of 0 and 1 typically correspond to homozygous and heterozygous sequence loss, respectively. CN states of 3 and higher correspond to sequence gain.

2. **cnv_calls_exp.txt** contains the start and end positions of copy number variant (CNV) calls (those whose CN states differ from baseline). The fractionalCopyNumber column contains the scaled and smoothed CN state for a given CNV call, and the CopyNumber column is rounded from fractionalCopyNumber. Confidence is defined differently for the two pipelines. For the integer CN pipeline, confidence refers to the probability of a CNV call as being a true call. For the fractional CN pipeline, confidence is 1 minus the probability of observing the mean coverage of a segment assuming that the segment had a baseline CN state.
3. **cnv_calls_exp_full.txt** is similar to cnv_calls_exp.txt. This file contains the combined results from the integer and fractional CN modules. An additional column named “Algorithm” denotes which module a call is from. Calls from the fractional CN module are denoted as “Region-based” while calls from the integer CN module are denoted as “Label-based”.
4. **cnv_chrAneuploidy.txt** contains per-chromosome aneuploidy calls. The fractChrLen column indicates the fraction of a given chromosome that was consistent with an aneuploidy event. The score column is the weighted probability of observing the segments in the aneuploidy event assuming the segments had baseline CN states. The fractCN column indicates the likely aneuploidy state of the chromosome.
5. **cnv_chr_stats.txt** contains more detailed per-chromosome statistics. The noise statistics for a sample and the inferred sex of the sample are noted in the header section.
6. **cnv.log** contains information about the CN analysis run. Warning and error messages may be useful for troubleshooting purposes.

Environment setup

The pipeline was tested with R version 3.6.1. The pipeline requires publicly available R packages and two Bionano-specific packages (BionanoR and FractCNV). For installation help, see Bionano Solve Installation Guide (PN 30172).

Usage

The *de novo* assembly pipeline and RVP, call the main R wrapper script “CNV.R” to perform CN analysis. The CN analysis is typically run using default Bionano-supplied control data. Additional parameters to the *de novo* assembly and RVP may be added, as shown in the table below. The parameters for the fractional CN analysis module (cr, cm, and ce) are new and can be supplied in the Bionano Access interface.

Module	Parameter	Notes (corresponding parameters in CNV.R in the first line of each entry)
--------	-----------	---

Integer	cd	<p>--controlsDir [dir]</p> <p>User-provided control data directory that contains r.cmap output from alignmolvref. We recommend using at least 10 control datasets in total and at least 5 of each gender. This option automatically triggers calculation of new parameters and new confidence tables for hg19 and hg38. If the control data is from an unrecognized source (for example, non-human, generated from a different enzyme), the CN pipeline would output CNV calls without confidence scores. Please refer to the “Pre-process custom control data” section on how to generate custom confidence tables.</p>
	pd	<p>--parametersDir [dir]</p> <p>User-provided path to either store processed control data and confidence tables when --controlsDir is defined, or as a different parameter input directory from the default parameter directory.</p>
	op	<p>--outliersProbability [float]</p> <p>“1 - Type I error rate” for outlier detection.</p>
Fractional	cr	<p>--controlRef [file]</p> <p>User-provided control data for the fractional copy number pipeline. This should be the output file generated from running the script generate_controlFractCNV.R (see “Pre-process custom control data”).</p>
	cm	<p>--cnvMask [file]</p> <p>User-provided CNV mask BED file for the fractional copy number pipeline. This should be the output file generated from running the script generate_cnv_mask_FractCNV.R (see “Pre-process custom control data”).</p>
	ce	<p>--chrExpectedCopyNumbersFile [file]</p> <p>Tab-delimited text file describing the expected copy numbers for each chromosome, for different sexes. The required column headers are cmapid, female, and male. Any chromosome ID missing from the table is assumed to be diploid for all sexes. See examples of this file in /Pipeline/Analysis/FractCNV/data/. An example of a sufficient chrExpectedCopyNumbersFile for human is shown below:</p> <pre> cmapid female male 23 2 1 </pre>

		24 0 1
		Note: the parameters <i>sexChrType</i> , <i>sexChr1</i> and <i>sexChr2</i> have been deprecated and are no longer used. The parameter <i>chrExpectedCopyNumbersFile</i> replaces these.

To run CNV.R as a standalone operation, please refer to the run_default.sh as an example in the main CN analysis code directory in the Bionano Solve build. The following input parameters may be used for additional flexibility when calling the CNV.R wrapper script.

Module	Parameter	Notes
General	--resultsDir [dir]	Path to output from the CN analysis pipeline. The default path is alignmolvref/copynumber in the <i>de novo</i> assembly pipeline and RVP.
	--sampleRcmap [rcmap]	Path to input r.cmap from the molecule-to-reference alignment. The default path is alignmolvref/merge/alignmolvref_merge_r.cmap.
	--sampleXmap [xmap or dir]	Path to input .xmap file(s) from the molecule-to-reference alignment. This can be the path to single xmap file that contains all the molecules-to-reference alignments, or a directory path to the "alignmolvref/merge" folder in the <i>de novo</i> assembly pipeline or RVP. The pipeline would retrieve the relevant xmap files from the directory path.
	--sampleName [string]	Sample name. The default is "exp".
	--chrRemove [string]	Chromosome ID in the CMapId column of r.cmap to exclude from the analysis. The format is the same as componentsRemoved.
Integer	--testingDir [dir]	Path to simulated data used for recalculating the confidence tables for the integer CN analysis module. The default path is /testing_data in the CN analysis pipeline package.
	--componentsRemoved [string]	Principal components derived from control coverage profiles to remove during normalization. The format is a:b where principle components a to b are removed, or a,b,c where principle components a, b, and c are removed.
	--controlSdCutoff [float]	The upper threshold of scaled coverage standard deviation for

		selecting control datasets.
--	--	-----------------------------

Theory

Fractional CN pipeline

The main steps of the fractional CN pipeline are:

1. Normalization
2. Segmentation
3. Baseline estimation
4. CNV detection and confidence score computation
5. Chromosomal aneuploidy detection

1. Normalization

This pipeline takes advantage of control data to model and reduce variation unrelated to true CNV events in raw coverage data. The set of control samples that are used are assumed to have very few or no true large CNV events. First, the coverage data is pooled from the control samples by summing the raw coverage at each label across the samples. The summed coverage at each label is then divided by the median coverage of all labels. This relative coverage accounts for the systematic biases in raw coverage data due to factors such as local label density patterns and repeat regions that may impact coverage, but are not directly related to true CNV events. For each query sample of interest, the raw coverage is normalized at each label by dividing it by the relative coverage estimated from the control samples. The normalized coverage data are used as input for segmentation in Step 2. An example illustrating the normalization step is shown in Figure 7.

2. Segmentation

Conceptually, the copy number (CN) state of a particular genomic region is proportional to the mean coverage of labels in that region. In order to detect complex CNV events with possibly multiple CN changes in the genome, a segmentation algorithm called Wild Binary Segmentation³ (WBS) is used to partition the genome into distinct segments. Each segment is inferred to have a statistically different mean coverage, and hence, a CN state that is different from its immediate neighboring segments in the genome. WBS is an extension to the Binary Segmentation and Circular Binary Segmentation algorithms that have been used extensively to detect CNV events in array-CGH or next generation sequencing (NGS) data. WBS has been shown to have better sensitivity in detecting smaller changes when compared to conventional methods. After segmentation, a set of genomic

³ Piotr Fryzlewicz. Wild binary segmentation for multiple change-point detection. Ann. Statist. Volume 42, Number 6 (2014), 2243-2281.

intervals is returned, where each interval is inferred to have a different CN state compared to its immediate neighboring regions. An example of the segmentation procedure is illustrated in Figure 8.

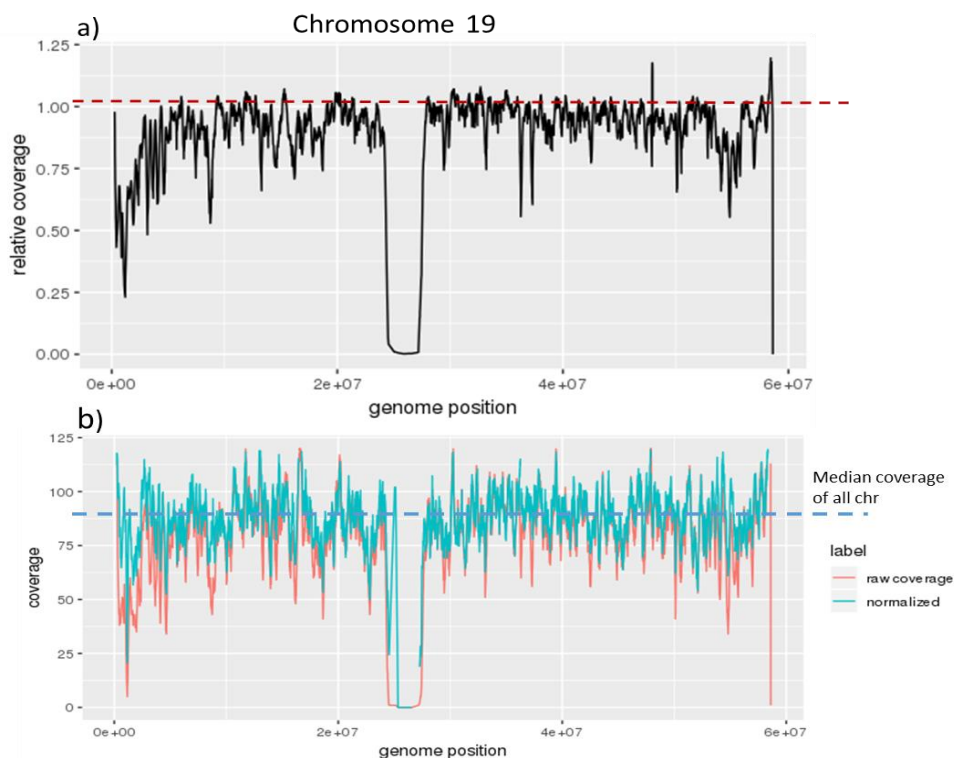


Figure 7. Normalization of raw coverage data. a) The relative coverage of chr19 from the pooled control data is shown. The relative coverage in chr19 was consistently below one. This means that on average, fewer molecules mapped to chr19 compared to other chromosomes. There is also local fluctuation (e.g. beginning of the chromosome) in the coverage due to factors such as label density. The normalization procedure is designed to correct for such systematic biases in raw coverage data. b) Comparison of coverage data from an example query sample before (red) and after (cyan) normalization. After normalization, the overall coverage became closer to the median coverage of the whole genome and local variation in the data was reduced.

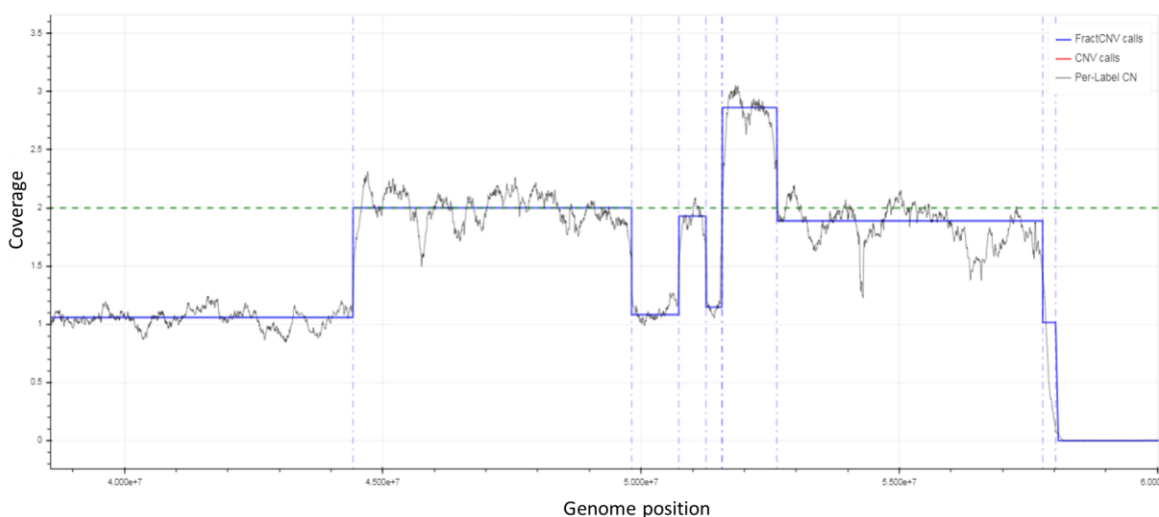


Figure 8. Segmentation of coverage profile. Genomic regions with the same CN states are assumed to have the same mean coverage with random local fluctuation around the mean. “Jumps” are expected when neighboring regions have different CN states. The segmentation procedure looks for such signature jumps and partitions the genome into non-overlapping intervals that are deemed to have different CN states. The figure shows example results from segmentation. The segment boundaries are marked by the vertical dashed lines. The solid blue line indicates the average coverage for a segment, and the black line shows the normalized coverage at each label along the genome.

3. Baseline estimation

One is interested in detecting regions of the genome where the CN states deviate from the normal or baseline CN state. In human, the normal CN state is diploid (copy number of 2). However, unless external data such as those from karyotyping are available, the part of the genome that is diploid is not known *a priori* and could only be inferred from the data. Also, there is not enough information to infer *absolute* copy number states of a genome solely from coverage data. Only *relative* copy number states can be inferred from the data. If the true ploidy is known based on external sources and is not diploid, users could scale the results accordingly. If the true ploidy is 3, one would multiply the inferred copy number states by 3/2 (or 1.5) for all chromosomes.

The pipeline makes a simplifying assumption that for a genome of interest, the fraction of the genome with a normal CN state is larger than the fraction of the genome with abnormal CN states. It then estimates the mean coverage of the normal diploid state by computing the mode of the coverage of all labels. This is defined to be the baseline coverage for a genome, and it is used as the basis for calling genomic regions with abnormal CN states. However, baseline estimation can be challenging if there are multiple overlapping underlying distributions of similar proportions, resulting in close peaks in the overall coverage distribution.

Two examples of baseline estimation is illustrated in Figure 9. The sex chromosomes are treated differently than autosomes. For human datasets, the pipeline first checks whether chrY is presented by determining if there is

non-trivial molecule coverage along chrY. If chrY is deemed present, the baseline coverage for the sex chromosomes is adjusted to be half of that of the autosomes. Otherwise, the baseline coverage for chrX is assigned to be the same as autosomes. The expected sex chromosome coverage values can be specified in the input for other species as discussed in the Usage section.

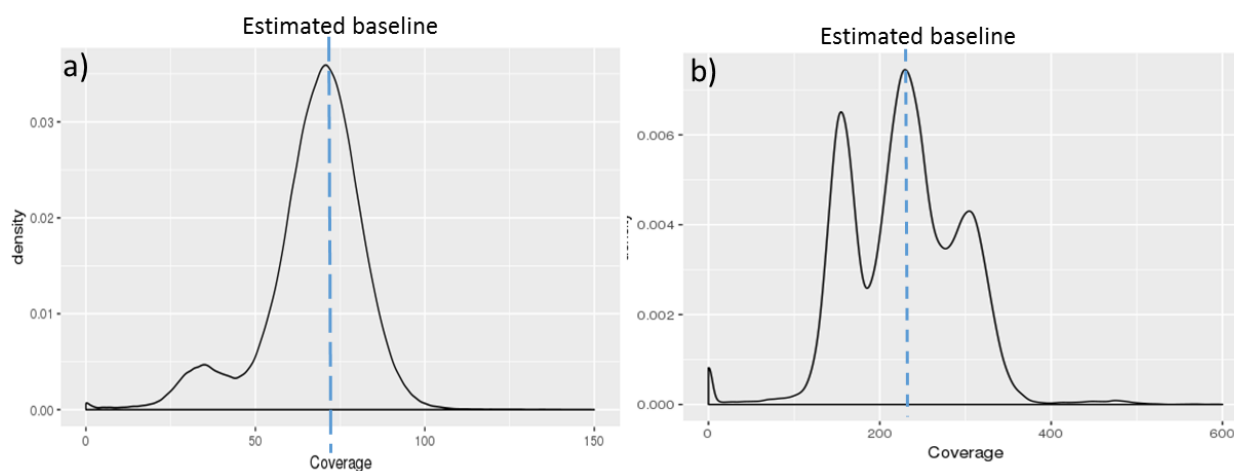


Figure 9. Baseline estimation. To estimate the baseline coverage, the distribution of the coverage data is analyzed, and the mode of the distribution is taken as the baseline coverage. a) The coverage distribution of a normal sample is shown. This sample is expected to have very few or no true CNV events; thus, the coverage data shows one single dominant distribution. The smaller distribution with the mode approximately half of that of the dominant distribution is the coverage distribution for sex chromosomes, since this is a male sample. b) The coverage distribution for a cancer sample with complex CNV events is shown. The various peaks in the distribution correspond to known chromosome aneuploidy events in the sample. The mode of the overall distribution (i.e. the tallest peak) is defined to be the baseline coverage of the genome.

4. CNV detection and confidence score computation

Once the baseline is established, the pipeline determines which genomic intervals from the segmentation step have abnormal CN states by checking if a particular genomic region has a mean coverage that is (statistically) significantly higher or lower than the baseline. It expects that genomic regions with same CN states have the same mean coverage equal to the baseline coverage plus some small variation due to random and systematic fluctuation in the data. Conceptually, the expected coverage (and hence, a CN state) under the null hypothesis should be consistent with the baseline. Intuitively, the smaller this probability is (i.e. the chances of observing the empirical mean coverage of a region's being much lower or higher than the baseline coverage or CN state), there is higher confidence that this particular region has a true CNV event or multiple events whose aggregated coverage deviates from the baseline. The probability is adjusted by accounting for multiple hypothesis testing. The confidence scores are computed as one minus the adjusted probability. Similar to Step 3, for human, the pipeline first determines if chrY is present. If it is, the confidence scores for sex chromosomes are computed by assuming the baseline coverage is half of those of autosomes. For non-human datasets, the pipeline handles the

sex chromosomes as specified in the chromosome configuration input. If the chromosome configuration is absent, the pipeline would assume that all chromosomes have the same baseline coverage.

Region-specific scoring

Before Bionano Solve 3.6, a single statistical model was used to model the coverage profile throughout the whole genome. However, as shown in the Figure below, the assumptions in the model are not always correct. Many regions in the genome in fact behave differently than a “typical” region (See Figure 9). These are often complex genomic regions such as centromeres or telomeres, or regions that contain segmental duplications or repetitive sequences. The coverage in these regions tends to be systematically higher or lower than the baseline and displays unusually high variability. Therefore, they are easily mistaken to be contain true CNV events, based on the whole-genome model. With the increasing availability of Bionano data from control samples (i.e. samples from healthy individuals where very few or no CNV events are known), we introduced a separate statistical model for each genomic region in Solve 3.6. Data from 99 control samples were used for building the model. Rather than using a single whole-genome statistical model to evaluate whether the coverage profile of a query genomic region is significantly different, now the coverage profile of the query region is evaluated against the background distribution of coverage profiles from a set of control samples in a region-specific manner (See Figure 10). A kernel density method is used to model the empirical distribution and compute region-specific p-values and confidence scores.

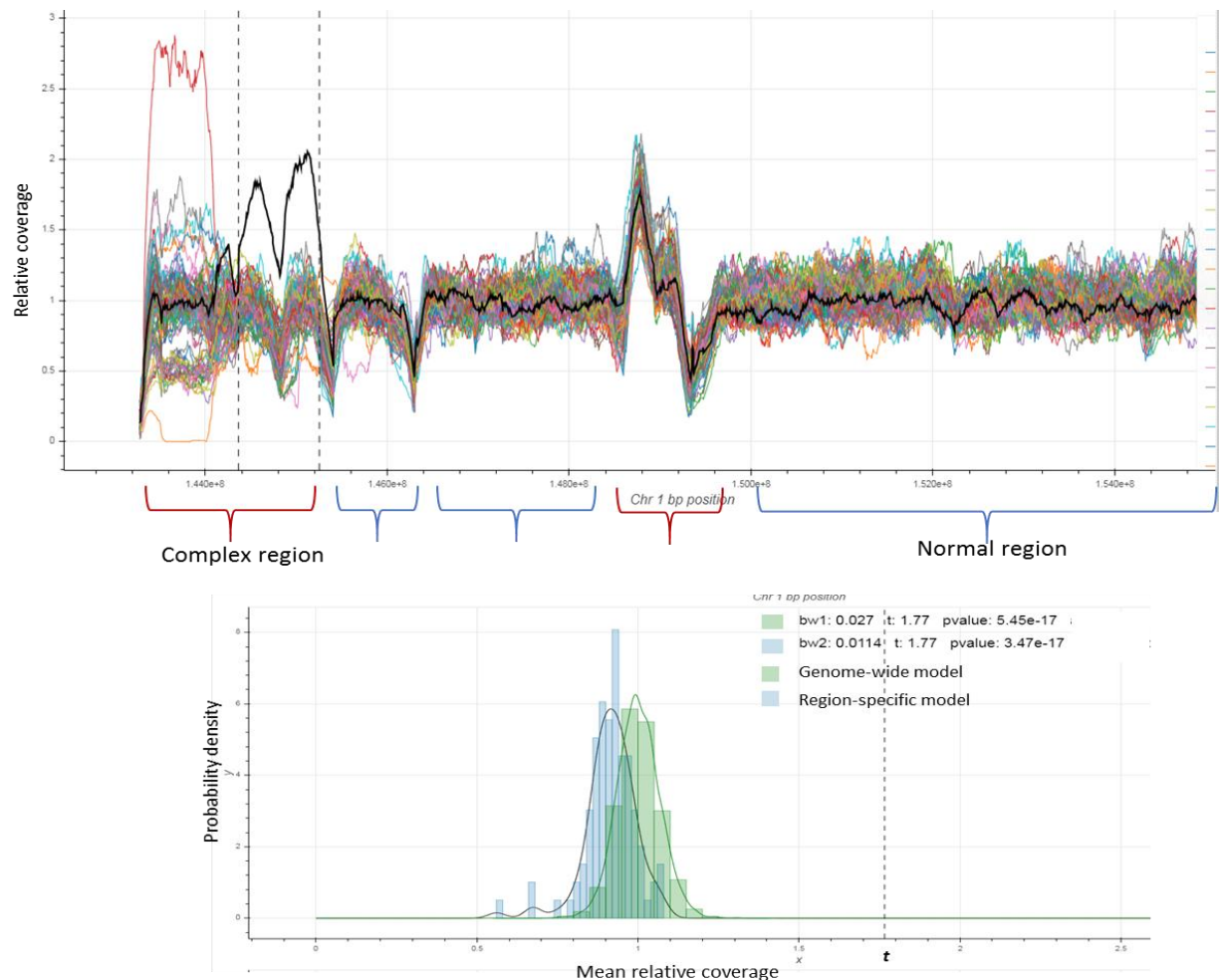


Figure 10. Top panel: Coverage profiles of an example region from approximately one hundred control samples are shown as colored lines. The coverage profile of a query sample is shown as bold black line. Genomic regions can be broadly divided into normal and complex categories. Normal regions are regions where the coverage profile follows an expected pattern – the coverage fluctuates randomly and uniformly around the baseline coverage (regions indicated by blue curly brackets). Complex regions are regions where the coverage profile has atypical patterns in a region/location-specific manner (indicated by red curly brackets). Starting from Solve 3.6, the CNV pipeline utilizes data from a large set of control samples to build a region-specific statistical model for each of the complex regions and determine whether the coverage profile of a query sample is statistically different from control samples. To illustrate this, a region of interest is highlighted between two vertical dashed line in the top panel. Bottom panel: in order to determine if the region has a CNV event, the mean coverage of all control samples is used to build an empirical distribution of coverage of the specific region (blue histogram). This is contrasted with the old whole-genome model (green histogram). The coverage profile of the query sample (vertical dash line in bottom figure and solid black line in top panel) is evaluated against the background distribution of the control samples (blue histogram), and the statistical significance is evaluated using a kernel density estimator.

Multi-allelic regions

Some regions of the genome are multi-allelic with respect to their copy number states in the control population. The empirical distribution of coverage profile in these regions is multi-modal (see Figure 10 below). In the current pipeline, we treated this multi-modal distribution as a single background distribution. As a result, some common CNV events in the control population will not be considered statistically significant events. However, if a rare or pathogenic variant that occupies a CN state that is substantially different than those in the control population, it would still be considered significant (see Figure 11, bottom panel).

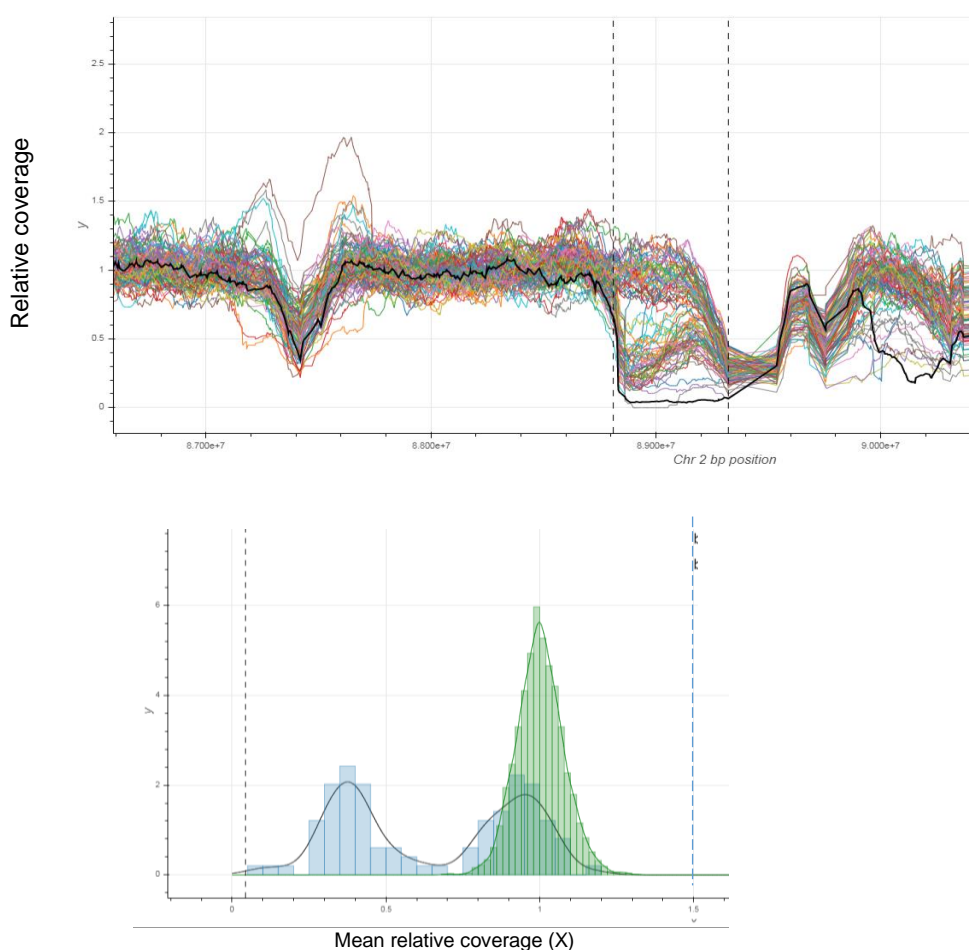


Figure 11. Example of a multi-allelic region. Top panel: Color lines are coverage profile from control samples. The region between two vertical dashed line is a region of interest with a multi-allelic CNV. Bottom panel: Blue histogram shows the distribution of mean relative coverage of all labels in the region of interest and green histogram shows the distribution of mean relative coverage across the genome with same size as the region of interest. The coverage of region of interest clearly has two distinct states around relative coverage of 1 and 0.5. As a result, samples with coverage near 1 or 0.5 in this region would not have significant CNV calls, but if a sample with relative coverage around 0, or 1.5 in this region (vertical dash line) would have statistically significant CNV calls.

CNV events in the Y chromosome

When analyzing the data from control samples, we observed that the variability of coverage profile in the Y chromosome is substantially higher compared to other chromosomes in the genome (see Figure 12). This means that the sensitivity of calling fractional CNV events will be relatively lower. In fact, a large portion of Y chromosome is in masked regions (see “Masking of high variance regions” for definition of mask regions). With the regions-specific scoring model introduced in Solve 3.6, CNV events could still be called in these variable regions, but users will have to ignore the mask and use the confidence score as a guide to look for true CNV events.

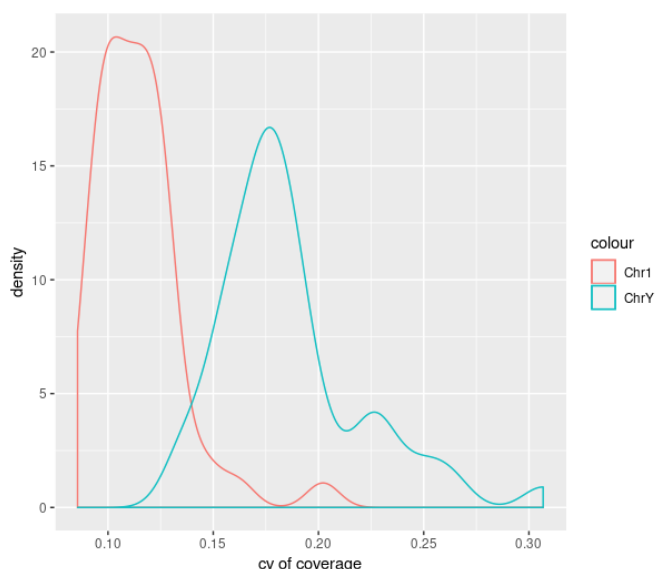


Figure 12. Coefficient of variation of coverage values across control samples for Chromosome 1 and Y.

CNV events near SV breakpoints

Traditionally, CNV pipelines detect gains and losses events by analyzing the coverage profile and detect any statistically significant increase and decrease in coverage depth. However, structural variation near the breakpoint of other types of SV such as inversion or translocation can also cause drops in coverage since molecules near the SV breakpoints cannot be fully aligned to the reference. In other technologies such as microarray or next generation sequencing, because the probes or the reads are very short compared to the size of typical CNV events, the impact on coverage around SV breakpoints tends to be minimal. However, due to the use of ultra-long molecules in the Bionano platform, the coverage drop can be significant and can lead to deletion calls. An example illustrating the coverage drop near the breakpoints of a large inversion is shown below (Figure 13). While there is likely not a genuine deletion, this drop in coverage is indicative of some rearrangement happening at the particular location. This needs to be taken into account when interpreting deletion events from the CNV pipeline.

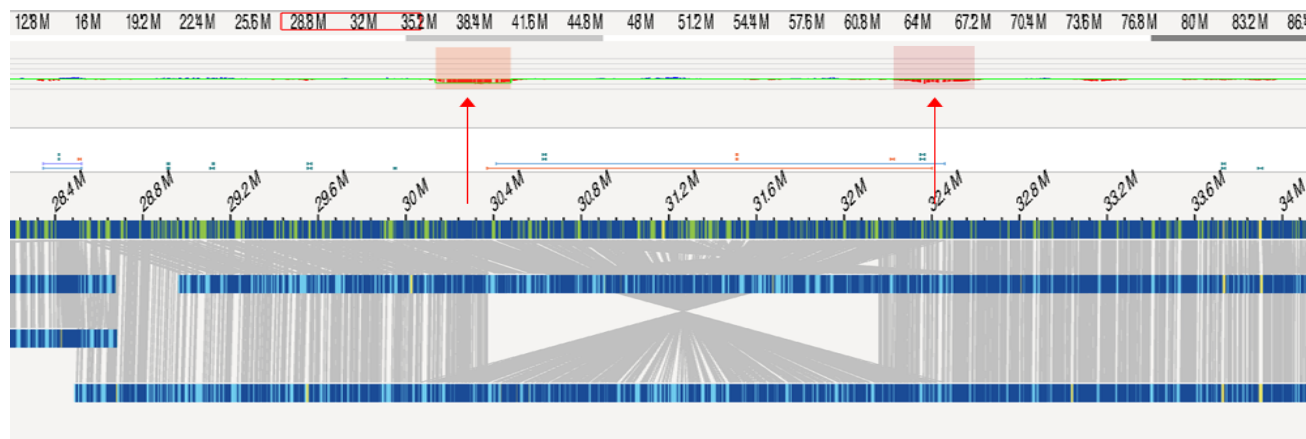


Figure 13. Example of deletion calls near SV breakpoints. The CNV pipeline detects statistically significant increase or decrease in the coverage profile data. In addition to detecting the classical duplication and deletion events, it sometimes detects drop in coverage near breakpoint around other SV type such as duplication. This is because molecules can only partially align to the reference near SV breakpoints and thus contribute to the lower coverage depth in such regions.

Systematic bias in coverage profile

In a small fraction of samples, it has been observed that certain systematic biases in the coverage profile exist. These biases cause systematic increase or decrease in coverage in a location-specific manner and do not correlated with true CNV events. The magnitude of these biases may differ across samples. Figure 14 shows examples of coverage profile with severe systematic bias (a), mild bias (b), and no bias (c). As illustrated in (c), the coverage profile looks relatively flat and uniform when there is no systematic bias. However, with bias (in a and b), the coverage profile looks like there are many CNV events. In order to assess the amount of bias in the input sample, the pipeline reports several statistics. The key observation is that when a sample contains significant systematic bias, the variation measure of the coverage profile would be elevated across the whole genome. This is often measured by the global coefficient of variation (cv) in the coverage profile. This value should usually be less than 15%.

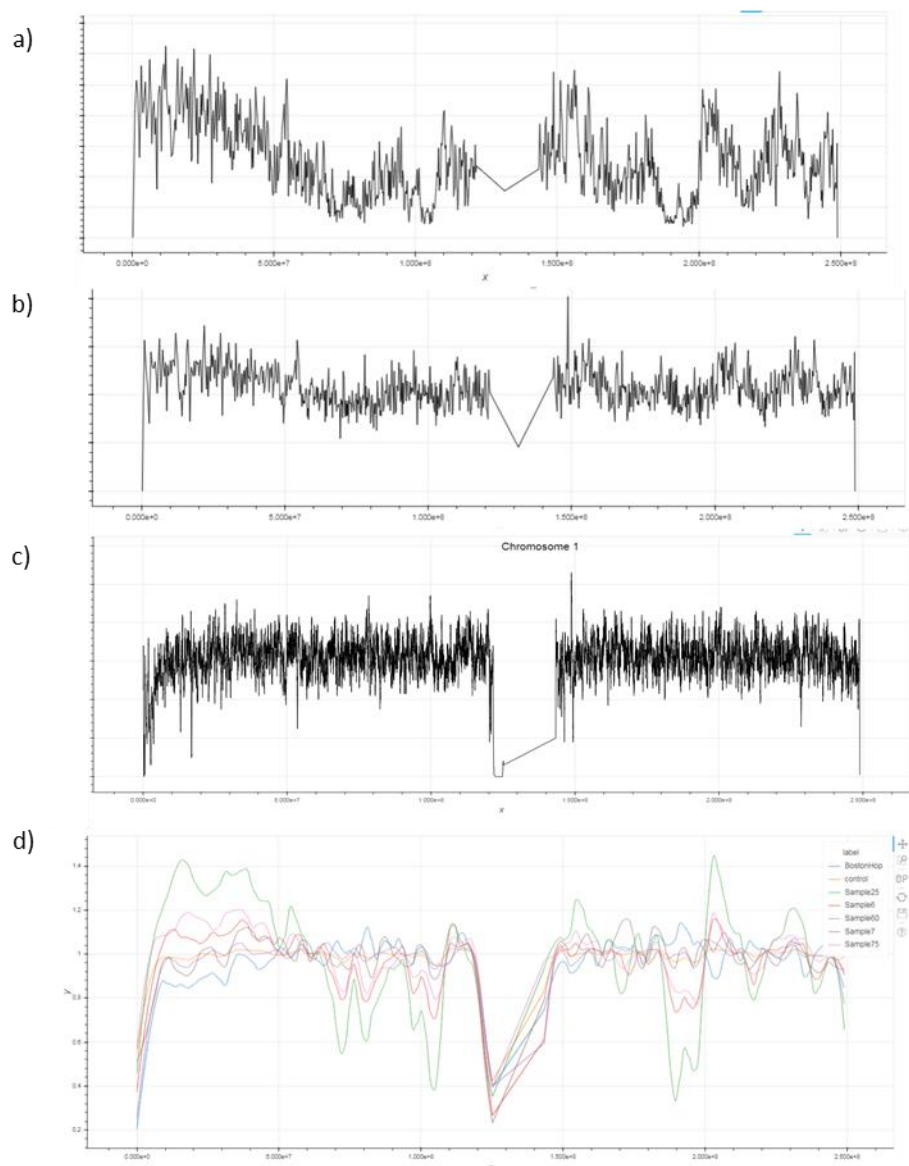


Figure 14. Example coverage profiles of samples with severe systematic bias (a), mild systematic bias (b), and no significant bias (c). The systematic bias shares a similar pattern across different samples. The colored lines in d) are the smoothed trend lines of the coverage profiles from multiple samples that are detected to contain different level of bias.

However, true CNV events can also lead to increase in variation in the coverage profile. In order to address this, the pipeline also computes the variation of the coverage in local windows. The idea is that nearby genomic locations are likely to be in the same CN states, so their variation should reflect the background variation of the data, despite the presence of true CNV events. We measure local variation in two window size: 2-Mbp windows and 6-Mbp windows. It provides variation estimates in different size scale. As mentioned in Section 4 above, the variation of coverage profile depends on the total coverage depth and size of the region being considered. The

higher the coverage or region size, we expect the variation to be smaller. Thus, based on empirical data, we built a model of what the variability of coverage data should be if a sample has no systematic bias. Then, we computed and reported the percent difference between the expected cv and the observed cv in the sample. The higher this difference is, the more severe the bias is in the sample. As a general guideline, if the observed percent different is greater than 20%, the sample should be considered to contain systematic bias. The noise estimates are in `cnv_chr_stats.txt`.

In Solve 3.7 we have further identified two types of systematic bias (Figure 15) in the coverage profiles and thus two metrics were added in the '`cnv_chr_stats.txt`' file that aim to quantify these systemic biases. These are 1) wave template correlation and 2) correlation with label density. The first metric measures the correlation between the coverage of the query sample and the coverage profiles from known samples with large systematic bias. The idea is to recognize known pattern of systematic bias in new samples. The second metric measures the correlation between coverage of a particular genomic region and the label density in that regions. It was found that this can be another source of systematic bias in coverage data. In both cases the higher the correlation the more serious the systemic bias. By analyzing the noise metrics in a set of control samples with or without known systematic bias, it was found for a sample to have reasonable low level of bias, the wave template correlation should be smaller than 0.4 and the absolute value of the label-density correlation should be less than 0.25.

There may be more false positive CNV calls in datasets with high systematic biases. The coverage profiles tend to fluctuate significantly even in the absence of true CNV events. Users noticing more than expected CNV calls are encouraged to contact Bionano Support for any assistance with such datasets.

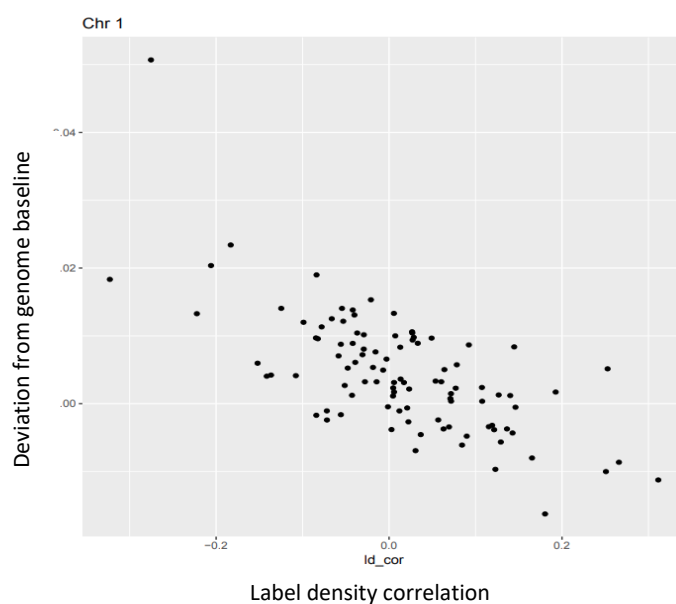
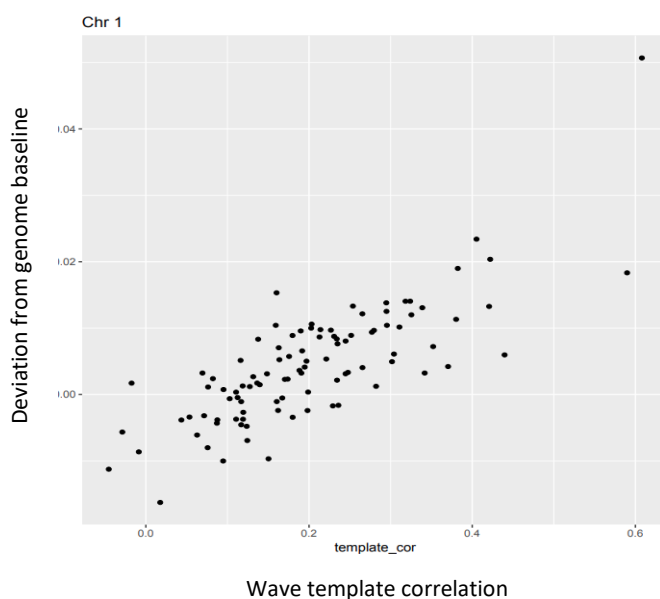


Figure 15. Quantifying systematic bias in coverage profile in control samples. In order to measure systematic bias in a set of control samples we computed the deviation of the coverage of a particular chromosome from that of the whole genome. Control samples are NOT expected to have chromosome aneuploidy events thus we expected the chromosome coverage to be equal to or very similar to the genome-wide baseline. Any deviation from it can be potential systematic bias in the coverage profile. Two possible sources contribute to this systematic bias. In the left figure we computed wave template correlation which is the correlation between the coverage of a sample versus samples with known systematic bias. In the right panel we computed the correlation between label density and coverage profiles. As shown in the figure, the higher these correlation the higher is the deviation of coverage of a chromosome from the genome-wide baseline, thus indicating that these two correlation can be used as a metric to quantify the systematic bias in new samples.

After calling potential CNV events, the pipeline performs an additional summary operation by grouping all the deletion or duplication events for each chromosome. If more than 80% of a chromosome by length has either lower or higher CN states than baseline, a whole-chromosome aneuploidy event would be reported in the output file with prefix "cnv_chrAneuploidy".

In Bionano Solve 3.7 a new algorithm was introduced to detect whole chromosome aneuploidy events with very low allele fraction (i.e. down to 5%). This new module considers the coverage data of each chromosome as a whole and compare that against the rest of the genome. It then detects whether the chromosome coverage data deviates from the whole genome baseline coverage. By considering the coverage data from the whole chromosome, the algorithm takes advantage of the large number of data points available and is able to detect very small changes in coverage. However this also means the algorithm is sensitive to local copy number changes in the chromosome. It can potentially mistake a significant local CNV change as a low allelic fraction whole genome event. Therefore, genomic regions with significant local CNVs that were detected in the previous steps were removed from consideration before analyzing each chromosome for aneuploidy event. This means that such analysis is not really suitable to analyze samples with extensive local rearrangement or CNV events such as cancer samples. In such cases, we fall back to the default aneuploidy detection algorithm as described above.

In order to improve the sensitivity of aneuploidy detection, we have identified two types of systematic bias that can cause fluctuations in the coverage data that does not correspond to biological events (see section Systematic bias in coverage profile above). These biases are relatively small and were ignored in previous version of the CNV pipeline. However, they become significant when one is trying to detect very small changes in coverage data as in low allelic fraction cases. A bias correction/normalization procedure was implemented to remove such systematic bias from coverage data. As shown in Figure 16 below, by removing such systematic bias we are able to reduce the variation of chromosome coverage significantly and which allow us to detect small fractional changes in coverage data.

After removing systematic bias from the coverage data, the empirical distribution of the relative coverage for each chromosome (relative to genome-wide baseline) was constructed based on a set of control samples. When a new sample was being analyzed, the statistical significance of the deviation of the chromosome coverage from its

genome-wide baseline was evaluated base on this empirical model and any significant deviation are considered as chromosome aneuploidy events. The confidence score reported is 1-pvalue base on the empirical statistical model.

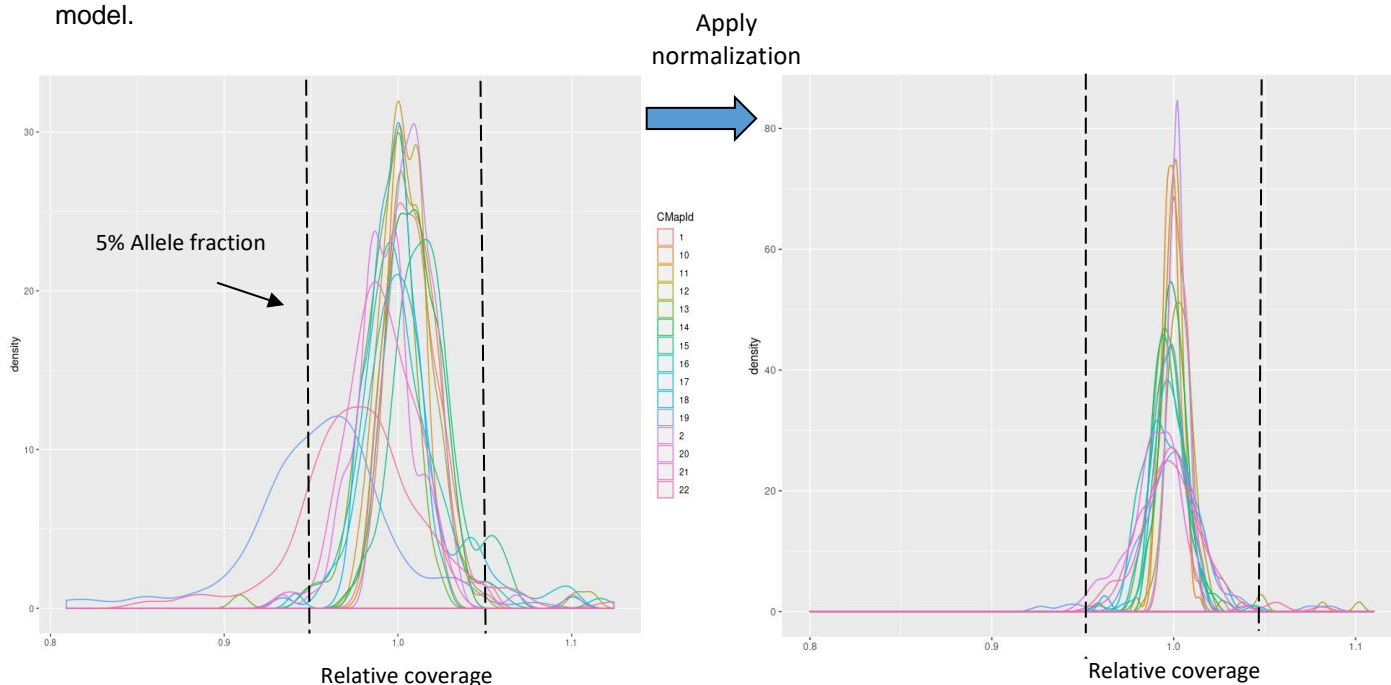


Figure 16. Systematic bias normalization reduces variation in chromosome coverage. The distribution of the relative coverage of each chromosome in a set of controls samples is shown. Relative coverage is the coverage of a particular chromosome relative to the whole genome coverage baseline. Different chromosomes have different variation in chromosome relative coverage. After the application of bias correction/normalization procedure the variation across all chromosome is significantly reduced.

Integer CN pipeline

The main steps of the integer CN pipeline are:

1. Scaling
2. Normalizing
3. Outlier detection
4. Smoothing
5. Confidence score calculation

1. Scaling

In the scaling step, the per-label coverage is standardized: coverage at each label is divided by the average coverage of all autosomal labels and then multiplied by 2. The pipeline currently assumes that the input comes from a diploid genome; therefore, the scaled per-label coverage should be non-negative and should average at 2.

2. Normalizing

The objective of the normalizing step is to reduce the non-CNV related variance of the per-label-coverage, minimizing noise in the coverage profile. To partition that variance out of the total variance, a dimension reduction method is implemented. Features that contribute to the variance in the control coverage profiles (from non-diseased samples that have no known large CNV events) are extracted using singular vector decomposition (SVD). The variance attributed to noise (e.g. rare CNV events, coverage fluctuation around fragile sites (for nickase data), and segmental duplication and N-base regions, and random sampling error) accounts for a large portion of the total variance. Theoretically, they can be explained by the first k principle axes, which are removed during normalization by subtracting them out from the scaled coverage profile such that only the variance related to CNV events is retained.

Processed control data are stored in the “parameters” directory. Normalization for a given chromosome may be skipped if the query profile is deemed incompatible with the control data. This is sometimes observed when analyzing the sex chromosomes and for highly rearranged genomes.

3. Outlier detection

The pipeline looks for labels with outlier coverage values; presumably, these labels are from CNV regions. Within the outlier detection step, there are two main sub-steps. First, a Hartigan’s dip statistics test for unimodality is applied to each chromosome to check if there is a large CN event up to 25% of the chromosome length. If so, the chromosome is split into segments, each with a separate baseline CN state. This allows the pipeline to identify large CNV changes. Then, a Chi-square test is performed for each label to determine if the coverage of a label is significantly different from the mean of all labels within the same segment. The default type I error rate for the Chi-square test is 0.05, corresponding to parameter [outlierProbability] 0.95.

4. Smoothing

In the smoothing step, several methods are used together to “stitch” together the per-label data from Step 3, smooth the normalized coverage profile, and generate the fractionalCopyNumber and CopyNumber columns.

The fractionalCopyNumber column is generated by a “gentle” smoothing method. First, the normalized coverage is rounded up to integer CN. Next, all CN values within median CN ± 0.25 inside a segment are smoothed to median CN. Then, the coverage profile is smoothed by applying running median of window size of 11 labels. Lastly, within each segment, neighboring outlier labels of the same direction (for example, either elevation or depletion relative to the segment median) are clustered, and their values are replaced by the extreme value within that cluster.

An “aggressive” smoothing method is applied to obtain CopyNumber. Neighboring segments with similar fractionalCopyNumber levels are merged together if 1) there are at least 55 labels in them, and 2) the lengths of the segments that have CN larger (or smaller) than the median CN are longer than 80% of the total size of all

close by fractionalCopyNumber. Their values are replaced by the median of the cluster. CNV events are sometimes broken up into small segments; this step ensures that these small segments are merged.

5. Confidence score calculation

The Confidence column in CNV calls in `cnv_calls_exp.txt` comes from the confidence tables in `/parameters/confidence_enzyme_alignmentType_reference.RData`. Each confidence table contains positive predictive value (PPV) by CN state and by size, starting from 500 kbp based on simulations. CNV calls falling into the same size bin and the same CN state in the confidence table would be assigned the same confidence score of that bin. Typically, smaller and heterozygous CNV calls have lower confidence scores.

The simulation was conducted for each enzyme (Nt.BspQI, Nb.BssSI, or DLE-1), platform (Irys or Saphyr) and reference (hg19 or hg38) combination, with two base control datasets being chosen from which CNV events were simulated. Duplication and deletion events with CN states 0, 1, 2, 3, and 4 were randomly simulated across the genome. The molecules overlapping with simulated CNV events were sampled according to the simulated CN states. Then, the sampled molecule set was aligned to the reference, and the alignment output was used as input to the CNV pipeline. CN calling performance was assessed, and PPV data were stored in the confidence tables.

If the input `r.cmap` is from an unrecognized source, the user needs to provide custom control datasets as input. Confidence tables need to be manually generated as well; otherwise, there would be no confidence output in the `cnv_calls_exp.txt` file.

Masking of high variance regions

Based on analysis of DLE-1 control datasets, we defined regions of the genome that had unusually high variance in their relative coverage compared to typical loci across control datasets. Relative coverage for each sample is defined as the raw coverage divided by the median coverage across the genome. We compute an average and standard deviation for each locus. A position is deemed variable if its variance is above three standard deviations away from the average.

These high variance regions are often concentrated around centromere and telomere regions, where molecule alignment can be unreliable. CNV calls overlapping with these regions (by at least 45% by length) are annotated as “_masked”; they are more likely to be false positive calls. The mask does not influence the initial CNV calling step. Calls are made with no knowledge of the mask. The mask is applied only after the calls have been made.

The masks are available for hg19, hg38, and mm10, and are automatically applied during the pipeline run. They are also available as BED files for visualization in Bionano Access. It is possible that common, highly polymorphic regions are included in the masks. The masked calls are shown in the CN table in Access and can be manually inspected.

Notes

- Performance for sex chromosomes is not as well-understood; less control data was available for sex chromosomes. Also, while supported, analysis of non-human samples has not been validated. Manual curation is encouraged.
- For female samples, no CN calls on chromosome Y are output.
- If there are no CNV calls, the CN table (cnv_calls_exp.txt) and the aneuploidy table (cnv_chrAneuploidy.txt) may be empty.
- The integer CN pipeline does not run on non-human datasets without the necessary auxiliary data, and corresponding columns in the output will be zeros.
- The pipeline assumes that the query dataset is of sufficient quality relative to the control datasets that were used for normalization.

Pre-process custom control data

Scripts for pre-processing custom control data are for advanced users to prepare their own control data for the CN pipeline. See Appendix I for detailed instructions.

SV masking

The *de novo* assembly pipeline and RVP include BED files for annotating insertion and deletion calls overlapping N-base gaps in the reference and putative FP translocation breakpoint calls so that they could be filtered in Access. The former would be annotated with a suffix “_nbase” in the SV type. Insertion and deletion calls in N-base regions may simply be due to mis-sizing of the N-base gaps in the reference and not genuine SVs. The latter would be annotated with a suffix “_common” or “_segdupe” in the SV type, depending on whether they overlap with common FP calls in control samples or annotated segmental duplication regions, respectively. For example, it includes selected sub-centromeric and sub-telomeric regions that are prone to generating putative FP translocation breakpoint calls. We provide BED files for the *de novo* assembly pipeline and RVP, separately. Details on custom BED generation are described in the Appendix B and FAQ sections.

The putative FP translocation breakpoint calls are derived based on from translocation breakpoint calls in control samples, and Bionano has continued to increase the number of control samples.

Here is the amount of sequence (in bp) for each of the categories in the **SV mask** BED files for Solve 3.6.1 and Solve 3.7 (the same “segdupe” and “gap” regions are shared between BED files):

De novo assembly	RVP
------------------	-----

		hg38	hg19	hg38	hg19
	segdup	88,295,249	85,579,474	88,295,249	85,579,474
	gap	85,960,835	234,345,079	85,960,835	234,345,079
Solve 3.6.1	common	15,054,403	16,416,297	24,390,440	24,158,950
	<i>Total</i>	<i>189,310,487</i>	<i>336,340,850</i>	<i>198,646,524</i>	<i>344,083,503</i>
Solve 3.7	common	26,158,264	25,709,013	10,241,633	12,456,961
	<i>Total</i>	<i>200,414,348</i>	<i>345,633,566</i>	<i>184,497,717</i>	<i>332,381,514</i>

Solve 3.7 uses the same **CNV mask** BED files as Solve 3.6.1 for both the *de novo* assembly pipeline and RVP. Below is the amount of sequence in the file (bp):

	hg38	hg19
Solve 3.7	278,537,677	277,270,477

SVMerge

SVMerge is no longer supported since the Bionano Solve 3.6 release. It was initially designed for merging SV calls from BspQI and BssSI datasets. For human SV analysis, we recommend the single-enzyme DLE-1 workflow.

Rare Variant Pipeline (RVP)

The ability to detect constitutional and somatic SVs is important for studies of cancer genomes. Bionano's *de novo* assembly-based SV detection approach – first constructing a *de novo* assembly and then comparing the resulting assembly with a reference – is highly sensitive in detecting homozygous and heterozygous insertions, deletions, duplications, inversions, and translocations in a diploid genome as small as 500 bp. However, this approach may miss SVs at low variant allele frequencies. Introduced in Bionano Solve 3.4, the Rare Variant Pipeline (RVP) is designed specifically to identify variants at low variant allele frequencies in the single-molecule data, while still finding homozygous and heterozygous variants. Low variant allele frequency variants are prevalent in heterogeneous samples such as cancers or samples with allelic mosaicism.

RVP consists of two major components: 1) the “split-read” analysis, and 2) the copy number analysis. Both analyses use information from a sample's molecule alignment against a reference assembly, bypassing the *de novo* assembly, and calling SVs based on these alignments. Neither analysis makes assumptions about the ploidy of the sample when detecting low allele-frequency variation. The “split-read” analysis calls SVs by examining the molecule alignments and searches for clusters of molecules with internal alignment gaps, and

multiple alignments. The copy number analysis detects copy number variation (CNV) by identifying regions of the genome with significant coverage elevation or depression.

There are three major steps in detecting SVs by the “split-read” analysis: initial molecule alignment and clustering of SVs, consensus generation by molecule extension refinement, and final SV calling.

Initial molecule alignment and clustering by SV

Molecules in an input BNX file are aligned directly to the reference, for example, human reference build GRCh38. Significant internal alignment gaps (termed outliers) and end alignment gaps (endoutliers) may indicate the presence of SVs in the sample. The pipeline requires only a default minimum of three molecules calling the same SV. Molecules are determined to confirm the same insertion, deletion, duplication or inversion if the inferred variant positions overlap and their inferred variant sizes are similar (default of 20% size similarity). To confirm the same translocation, the inferred translocations must be in the same orientation and their breakpoints in close proximity (within a default distance of 35 kbp).

Consensus generation by molecule extension refinement

A consensus is built using clusters of molecules that identified the same SV. The purpose of the consensus step is to verify that those SV-calling molecules truly agree and can form a consensus that represents the variant allele. In this step, the loci on the reference assembly flanking the inferred SV are extracted, and the molecules that called the SVs are aligned to each of the two SV-flanking reference fragments. The pipeline attempts to reconstruct the SV allele by using the molecules to extend into the SV region. If the molecules come from the same variant, they would have similar label patterns and form a consensus map that represents the variant allele. Note that for each SV, the same extension procedure is performed twice: one extension from the fragment left of the SV and one extension from the right.

Final SV calling

Finally, the new local consensus maps are realigned to the reference to check if the same initial SV calls are made. The pipeline will only report SVs that are confirmed in the final SV calling step. Note that potentially two consensus maps could form for each SV, but only one is kept in the end.

The SVs identified by RVP are listed in Bionano SMAP format. The SMAP file can then be used as input into downstream analysis workflows such as Variant Annotation Pipeline.

Variant Allele Fraction (VAF) calculation

Since Bionano Solve 3.7, the VAF values are provided by default in the output SMAP file of RVP, with the same interpretation as for the *de novo* assembly pipeline. Similarly, it is also possible to disable the computation from

the command line tool. The calculation uses consensus map coverage and the same Bayesian inference as for *de novo*, but requires some extra steps to obtain the coverage for the reference allele as it is not one of the RVP outputs. The extra steps are as follows: 1) Extract the fragments of the reference map spanning the SVs detected by the pipeline. 2) Competitively realign the input molecules to both the consensus maps and the extracted reference fragments; this way the input molecules aligned to the reference fragments provide the coverage for the reference allele 3) Run SVs calling on the consensus maps newly created during the competitive realignment. These maps and the newly called SVs are expected to minimally differ from those obtained after RVP but they will contain more accurate coverage values, and 4) Once the VAF is calculate for these new SVs, the value is transferred to their equivalent in the original RVP output. In the case that an original variant cannot receive a transferred VAF, the value will appear as -1.

Refer to appendix H for further information.

Structural variant calling performance: Rare Variant Pipeline

For performance validation, we took two complementary approaches. We used simulated data and real data to assess SV detection performance at different variant allele frequencies. Detailed methods are included in the Appendix A section.

Structural variant calling performance using simulated data

Summary

We observed at least 90% sensitivity at 300X effective coverage at 5% variant allele frequency:

- Insertions between 5 – 50 kbp
- Deletions > 7 kbp
- Translocations (or transpositions where the sizes of the translocated fragments are > 70 kbp)
- Inversions > 70 kbp
- Duplications > 150 kbp

Performance for simulated insertions and deletions

DLE-1 molecules were simulated from unedited and edited versions of the hg19 reference genome (with 1600 insertions and 1600 deletions from 200 bp to 1 Mbp), 300X effective coverage was mixed at different proportions, and used for SV calling with Rare Variant Pipeline (Figure 17). Sensitivity increases with variant allele frequency. Insertion sensitivity drops when insertions are larger than 50 kbp, and when using molecules with an average size is 250 kbp. Complementarily, as discussed in the next section, duplication performance is high above 50 kbp, and

most larger insertions are, in fact, duplications.

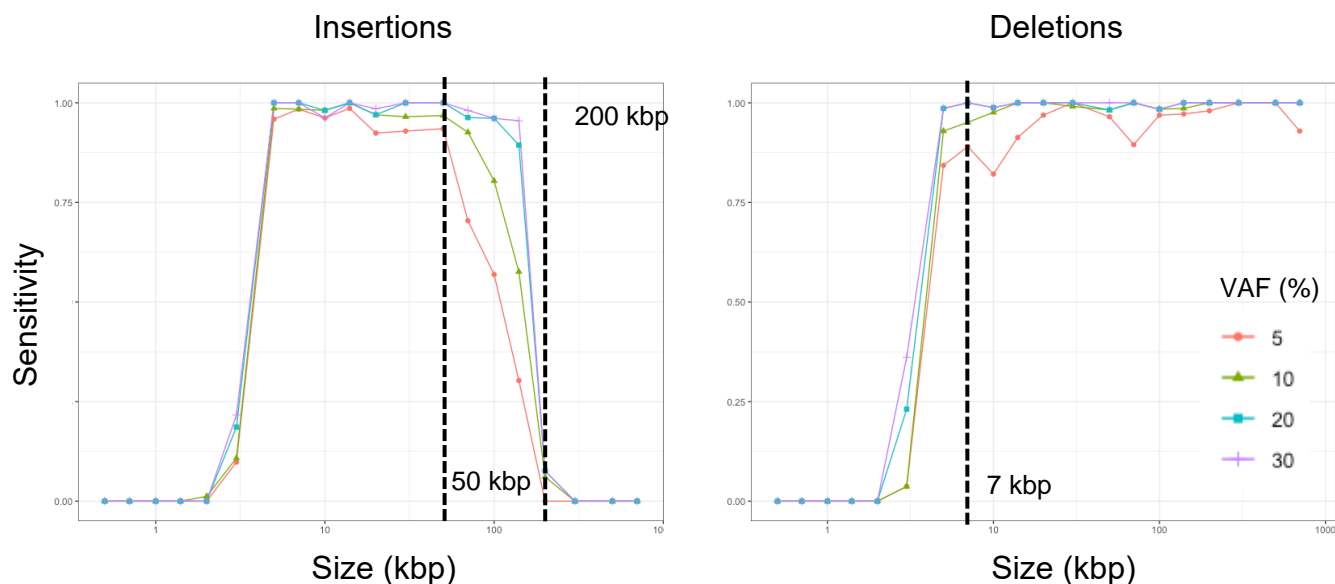


Figure 17. Insertion and deletion calling performance with simulated 300X DLE-1 data at different variant allele frequencies.

Performance for simulated duplications, translocation breakpoints, and inversion breakpoints

Performance data are shown in Figure 18. Duplications, translocations, and inversions need to be at least a certain size for the pipeline to recognize that a piece of sequence has been duplicated, transposed, or inverted. Small duplications and translocations will be classified as insertions, if the involved sequence is not large enough to be uniquely aligned. Small inversions with gain or loss of sequence relative to the reference may be classified as insertions or deletions.

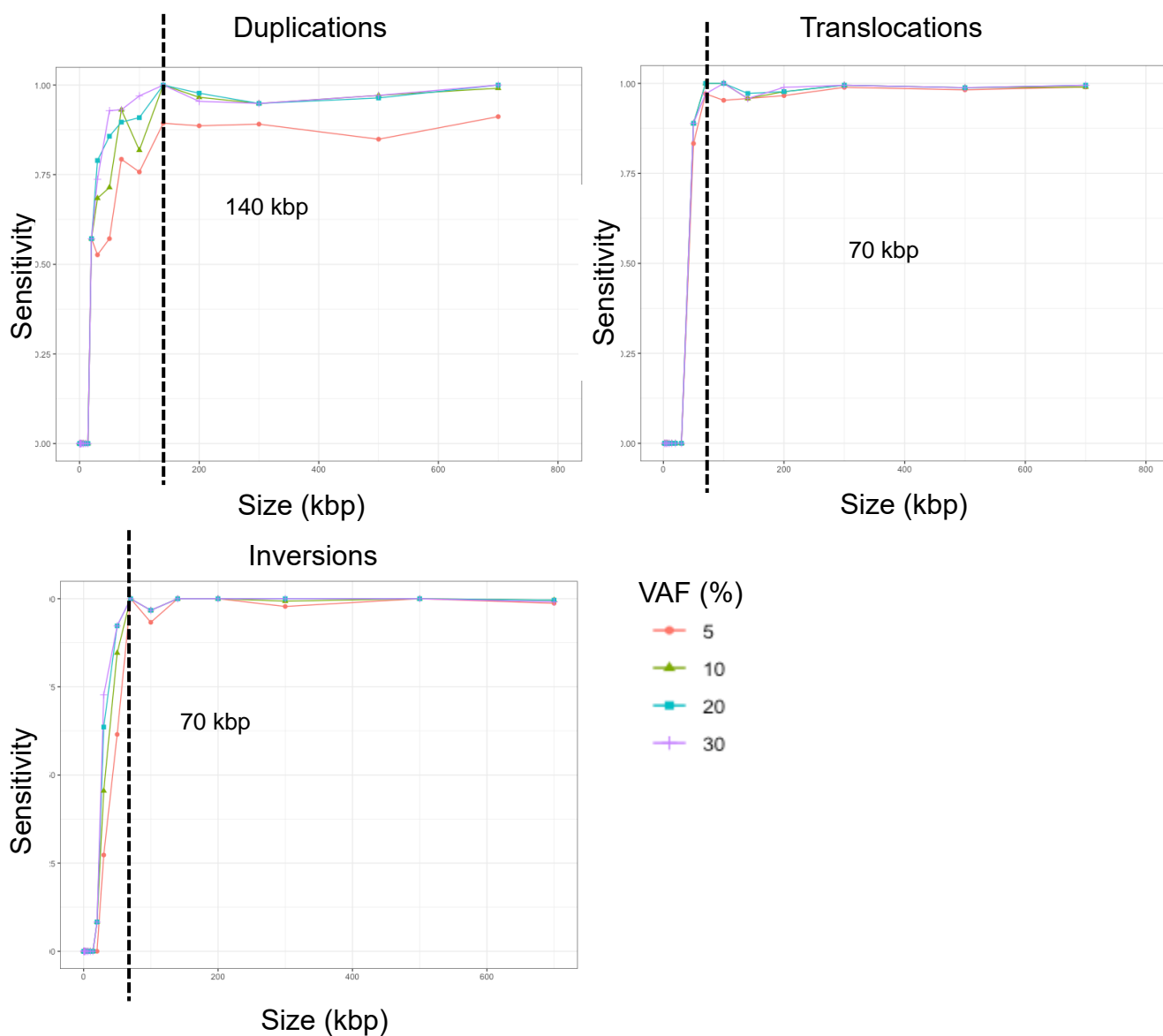


Figure 18. Duplication, translocation breakpoint, and inversion breakpoint calling performance with simulated 300X DLE-1 data at different variant allele frequencies.

Structural variant calling performance using real data

We made *in silico* 300X effective coverage mixtures of a cancer dataset and a NA12878 dataset at different proportions. SV calls were compared between the full coverage cancer dataset and the mixtures (Figure 19).

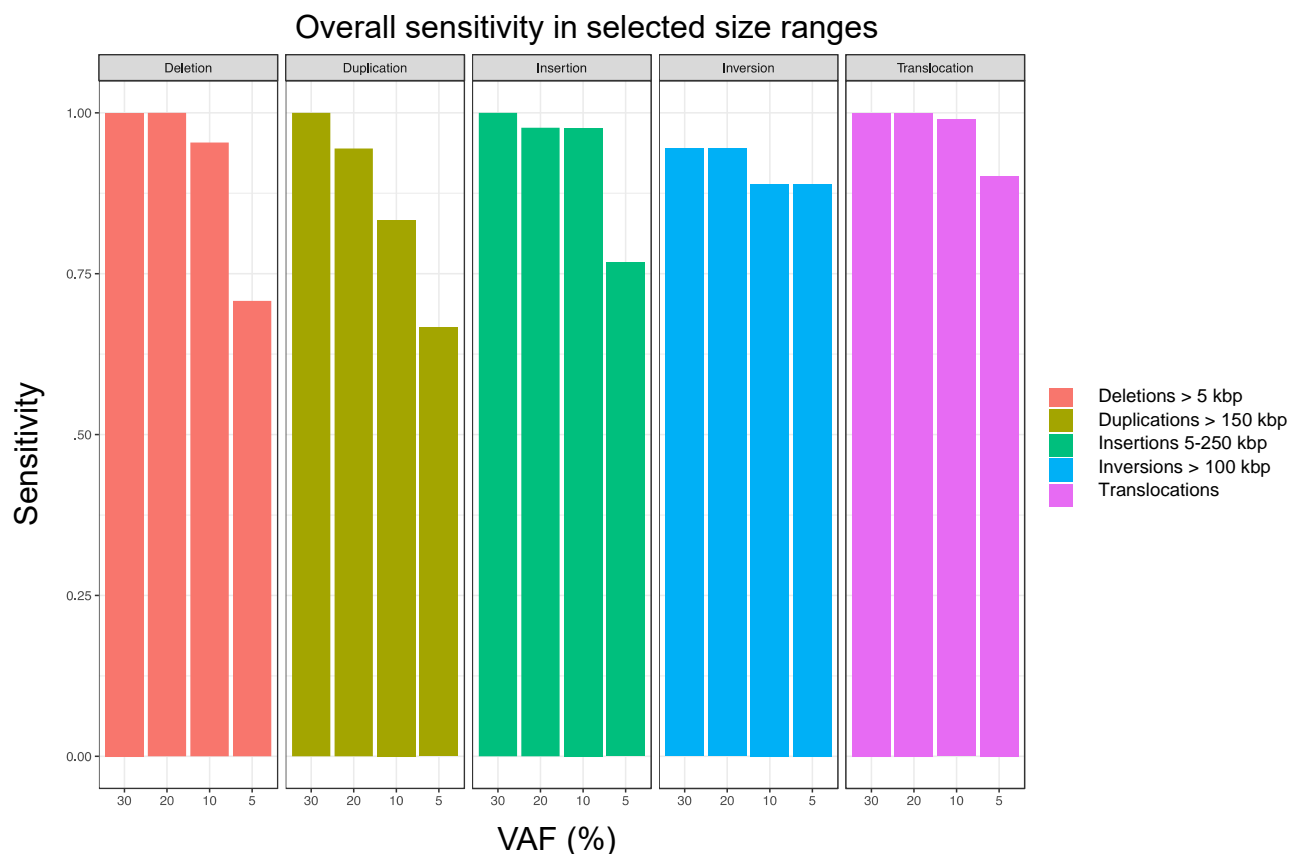


Figure 19. SV calling performance within selected size ranges with 300X mixtures of a cancer dataset and a NA12878 dataset at different proportions. There were 43 insertions, 65 deletions, 18 inversions, 18 duplications, and 193 translocations in the selected size ranges from the full coverage cancer dataset.

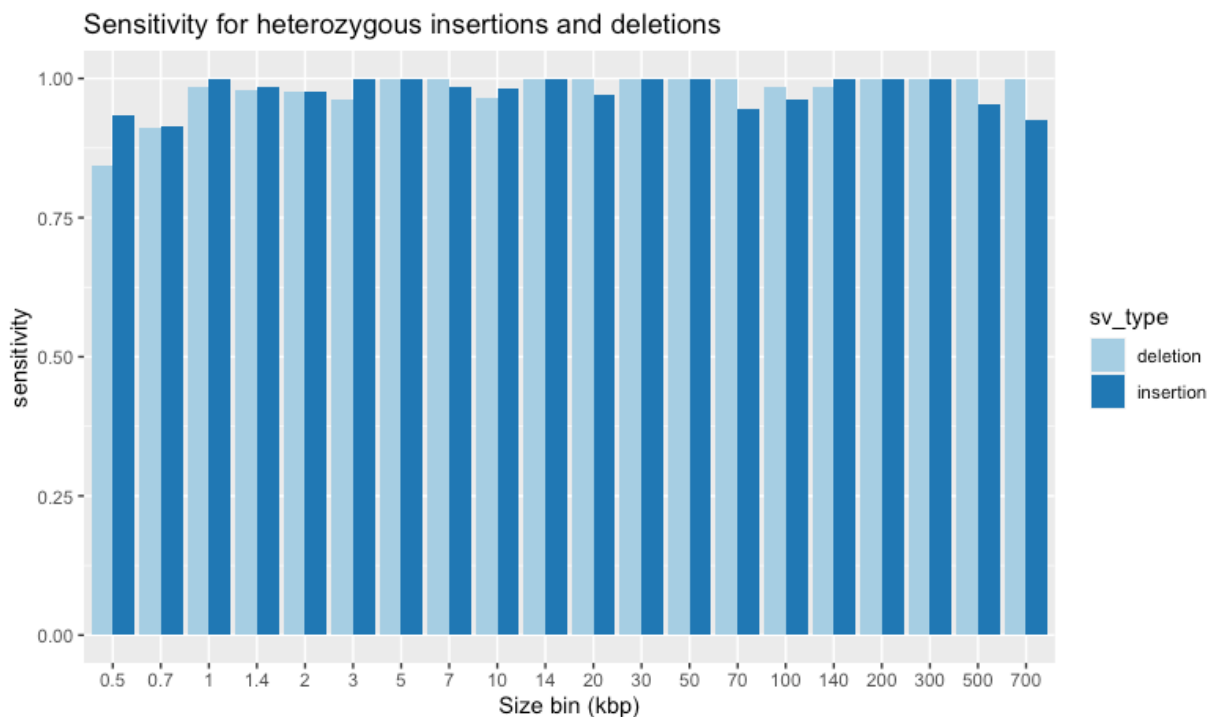
Structural variant calling performance: *de novo* assembly pipeline

We also used simulated data and real data to assess SV detection performance of the *de novo* assembly pipeline. Detailed methods are included in the Appendix A section.

Structural variant calling performance using simulated data

Performance for simulated insertions and deletions

DLE-1 molecules were simulated at 80X effective coverage from unedited and edited versions of the hg19 reference genome (with 1600 insertions and 1600 deletions from 200 bp to 1 Mbp) and used for assembly and SV calling with the Bionano Solve pipeline (Figure 20).



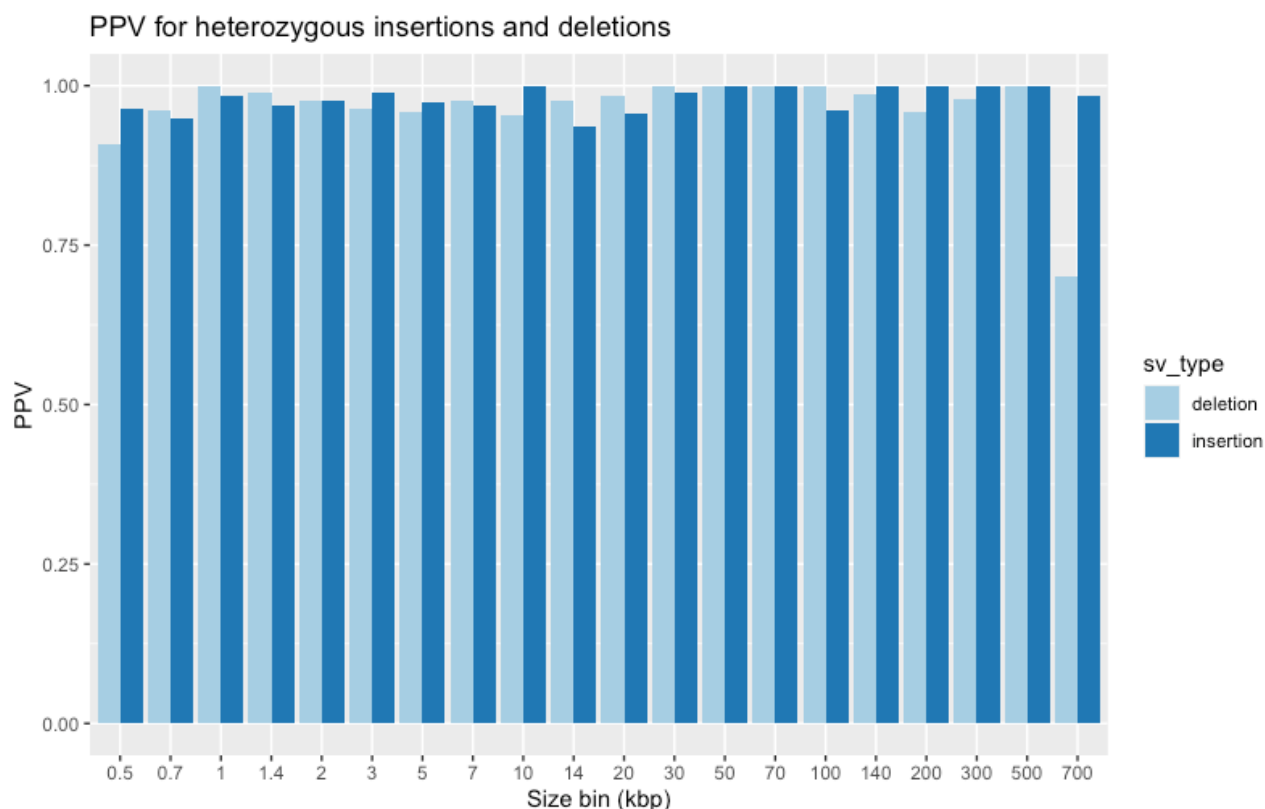


Figure 20. Heterozygous insertion and deletion calling performance with simulated 80X DLE-1 data. Top: sensitivity data; Bottom: positive predictive value (PPV) data.

Performance for other simulated events

Similarly, DLE-1 molecules were simulated from unedited and edited versions of the hg19 reference genome (with translocated fragments, or transposition events, inversions, and duplications of different sizes) and used for assembly and SV calling with the Bionano Solve pipeline. The performance data are shown in Figure 21.

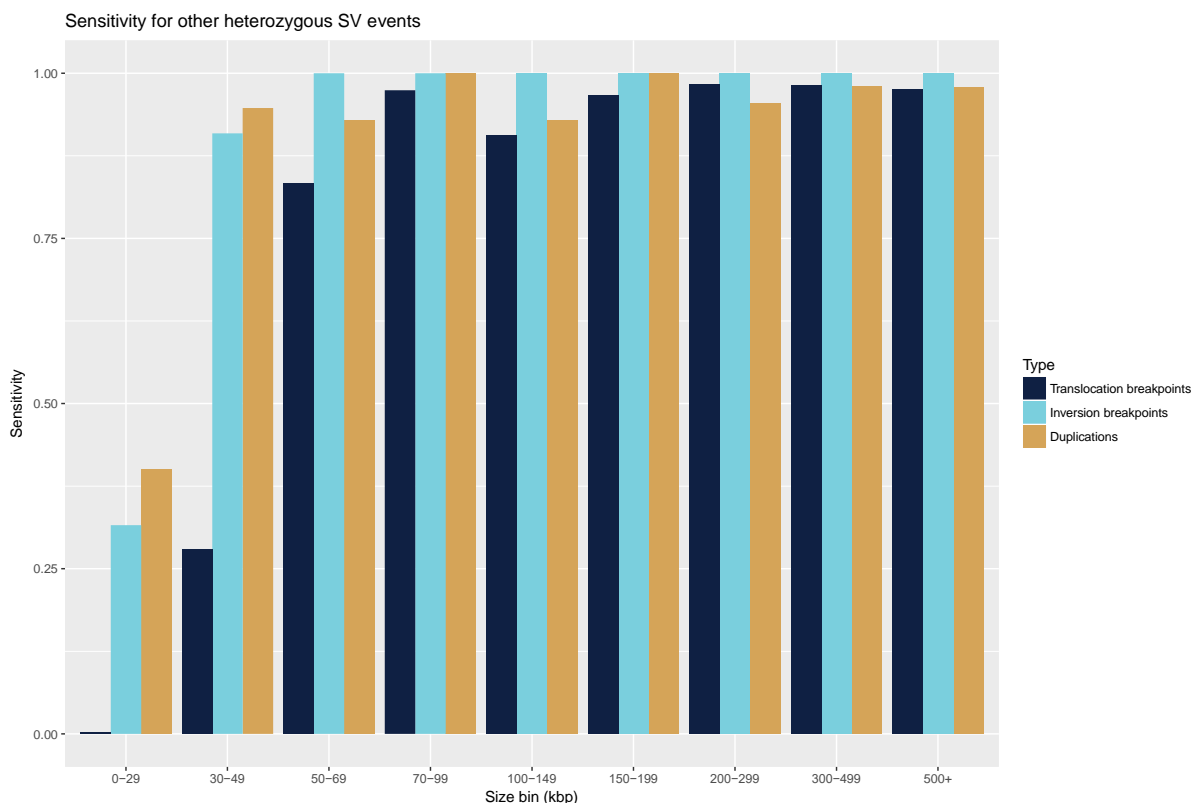


Figure 21. Heterozygous sensitivity of various SV types with simulated 80X DLE-1 data.

Structural variant calling performance using real data

Performance for insertions and deletions

CHM1/13 analysis was performed on an *in silico* mixture of CHM1 and CHM13 DLE-1 datasets. SV calls were compared from pure CHM1 and CHM13 assemblies against SV calls from the mixture. Overall performance is presented for SVs larger than 700 bp (Table 2).

Table 2: Insertion and deletion (> 700 bp) calling performance from CHM1/13 datasets using DLS.

Size cut-off	Type	Expected zygosity in mixture	<i>In silico</i> mixture	Individual assemblies	Fraction captured (%)	PPV (%)
700 bp	Insertions	Homozygous	1,975	1,963	99.4	97.5
		Heterozygous	1,737	1,480	85.2	
	Deletions	Homozygous	695	687	98.9	97.1
		Heterozygous	1,189	1,085	91.3	

CEPH trio Nt.BspQI datasets used in Mak *et al*⁴. were re-analyzed, and the resulting SV lists were compared with the published SV lists (Table 2⁵). Overall, we detected most of the published calls and made a large number of new SV calls. For example, we detected 93% of the published deletions in NA12878 and made an additional 764 deletion calls.

Performance for translocation breakpoints²

Nine samples with annotated single translocations were tested for translocation calling performance (Table 3). We found the expected translocation breakpoints in eight out of nine samples. For the remaining sample, only 10% of the cells contained the expected translocation based on karyotyping results. Reciprocal breakpoints were found in six out of eight samples for which expected translocations were found. In addition, putative FP calls were masked and/or filtered (Table 4).

Table 3. Translocation calling performance from real datasets (at close to or less than 80X effective coverage).

Sample	Diagnosis	Annotation	Sample prep	Found expected breakpoint?	Found reciprocal breakpoints?
Sample01	CML	FISH: t(9;22)	Plug lysis BspQI	Yes	Yes
Sample02	CLL	FISH: t(11;14)	Plug lysis BspQI	Yes	No
Sample03	CML	FISH: t(9;22)	Plug lysis BspQI	Yes	Yes
Sample04	CML	Kary.: t(9;22)	Plug lysis BspQI and BssSI	Yes	Yes
Sample05	AML	Kary.: t(7;11)	Plug lysis BspQI	Yes	No
Sample06	AML	Kary.: t(8;21)	Plug lysis BspQI	Yes	Yes
Sample07	AML	Kary.: t(4;5)	Plug lysis BspQI and BssSI	No	N/A

Cell lines GM16736 and GM21891 obtained from the Coriell Institute for Medical Research each contained one known translocation t(9;22) and t(4;15), respectively. Based on 80X effective coverage assemblies, we found the expected translocation breakpoints for both samples and additional translocation breakpoints. The additional breakpoint in GM21891 was filtered out after applying a confidence threshold of 0.1. An additional breakpoint with confidence score of 0.19 remained in the GM16736 BssSI assembly, which may be an unannotated translocation call in the sample. The BssSI masks were generated based on a small control database; thus, they may be less

⁴ Mak AC et al. Genome-Wide Structural Variation Detection by Genome Mapping on Nanochannel Arrays. Genetics. 2016 Jan;202(1):351-62.

⁵ The data were presented in a previous version of this document and are included here for completeness.

effective than the BspQI masks.

Table 4. Translocation breakpoint call masking and filtering for Coriell samples GM16736 and GM21891.

		Coriell (BspQI)		Coriell (BssSI)	
#Translocation calls		GM16736	GM21891	GM16736	GM21891
hg19	Known Translocation	1	1	1	1
	Remaining calls after masking	1	2	2	2
	As expected?	Yes	1 add. call	1 add. call	1 add. call
	As expected after applying confidence threshold ≥ 0.1	Yes	Yes	1 add. call	Yes
hg38	Remaining calls after masking	1	2	5	2
	As expected?	Yes	1 add. call	4 add. call	1 add. call
	As expected after applying confidence threshold ≥ 0.1	Yes	Yes	1 add. call	Yes

Performance for inversion breakpoints²

Analysis of samples with known events

Three Coriell cell lines with known inversion events were used for assessing inversion calling performance. Two-enzyme data were generated for two of the three samples, and expected inversion events in all three cell lines were observed at the annotated locations (Table 5).

Table 5. Inversion calling performance from real datasets (at close to or less than 80X effective coverage). *Called as intra-chromosomal translocation breakpoints, since those inversions were larger than 5 Mbp.

Sample	Diagnosis	Annotation	Sample prep	Found expected inversion?
GM19238	Phenotypically normal	Inv(15q13.3)	Plug lysis BspQI	Yes
GM14266	Micrognathia	Inv(4q34.2-35.2)	SVMerge	Yes*
GM21074A	Developmental delay	Inv(2p23-q31)	SVMerge	Yes*

Comparison to Valor calls

For the NA12878 sample, we compared Bionano calls to inversion calls predicted by a new NGS based algorithm.

Recently, Rasekh et. Al.⁶ described a novel algorithm, Valor, to detect large inversions based on “split clone” signatures. Out of the 18 inversion calls predicted by Valor, 15 overlapped with Bionano calls. Bionano detected an additional 82 putative inversion regions.

AOH/LOH detection performance

We evaluated the performance of AOH detection in two ways. First, we simulated datasets by splicing together data from healthy controls and haploid genomes. We found that we were able to detect AOH regions in the 40-50 Mbp size range with 92% sensitivity and 97% precision, where true positives were defined as an 80% overlap between simulated and predicted regions (Figure 22).

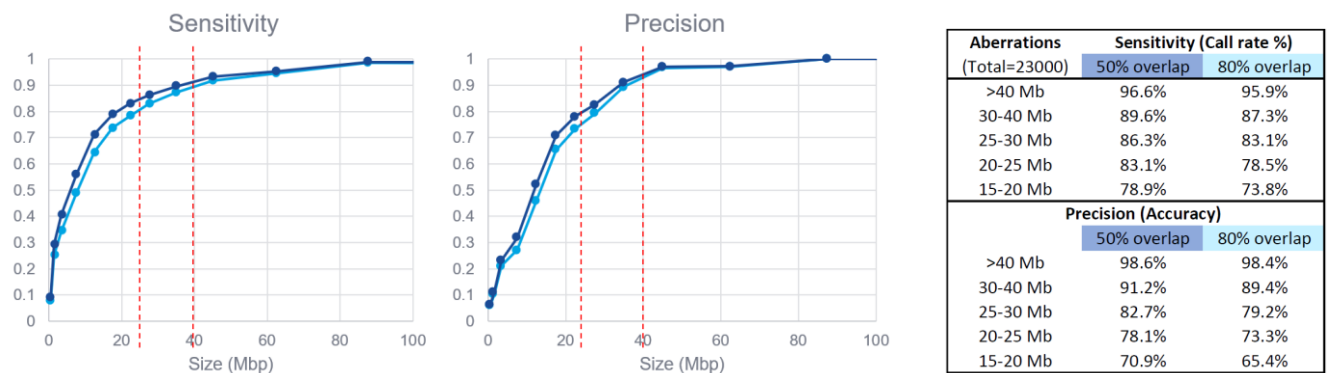


Figure 22. Performance of AOH/LOH detection with simulated human data, for simulated events of a given size.

Based on this evaluation, calls less than 25Mbp are not reported unless they are within 25 Mbp of another call. Second, we ran AOH/LOH detection on samples with known events that were previously identified using microarrays. In three samples with 4 events greater than 25 Mbp, 4/4 events were called correctly, some with 94% overlap, and some were called as multiple smaller regions with a total of 60-70% overlap between the known and predicted regions (Table 6). A fourth sample had 4 smaller known events, 3 of these were detected and one was filtered out from the output due to its small size. The detected AOH regions that were smaller than 25 Mbp were not filtered out from the output because there were other small calls nearby.

⁶ M Eslami Rasekh et al. Discovery of Large Genomic Inversions Using Long Range Information. BMC Genomics 18 (1), 65. 2017 Jan 10.

Table 6. Performance of AOH/LOH detection with constitutional human samples. Large events (>25 Mbp) are in bold, all large events were detected.

Sample	CNLOH identified using microarray	AOH/LOH detection results from OGM	Putative FP calls
1	7q11.22qter(71748536_159119707)x2 hmz, [CNLOH 87 Mb; EZH2]	7: 76,496,409 - 143,709,638 (67.2 Mbp) 7: 144,329,449 - 159,345,973 (15.0 Mbp) Total 82.2 Mb (94.5% overlap)	1 = 14.9 Mbp 6 < 6.5 Mbp
2	7q22.1qter(98411980_159119221)x2 hmz[0.7], [CNLOH 60.7 Mb; CUX1 and EZH2]	7:101,355,164 - 108,191,148 (6.8 Mbp) 7:124,805,963 - 143,709,638 (18.8 Mbp) 7:148,325,410 - 159,345,973 (11.0 Mbp) Total 36.8 Mb (60.6% overlap)	1 = 17.8 Mbp 5 < 6.2 Mbp
3	11q12.1qter(57994955_134942626)x2 hmz, [CNLOH 77 Mb; CBL] 13q13.1qter(32553841_115107733)x2 hmz [CNLOH 83 Mb]	11:56,370,292 – 99,821,083 (43.5 Mbp) 11:21,178,985 – 29,596,627 (8.4 Mbp) 11:33,078,807 - 35,086,622 (2.0 Mbp + 1.9 Mbp) Total 53.7 Mb (69.7% overlap) 13:27,998,292 – 114,364,328 (86.4 Mbp) 98.6% overlap, reciprocal overlap is 94.7%	13 < 8.5 Mbp
4	1q25.3q31.1(181,349,929-189,805,990)x2 hmz 5q23.1q31.3(119,409,742-139,707,439) x2 hmz 14q24.3q32.11(78,367,380-91,128,309)x2 hmz 22q12.1q12.3(27,696,986-34,194,745)x2 hmz	1:182,268,699 – 207,724,864 (25.4 Mbp) 14:73,748,304 - 95,896,776 (22.1 Mbp) 22:27,136,865 – 39,043,969 (11.9 Mbp)	16 < 7.2 Mbp

In Solve 3.7, AOH/LOH calling has been validated only for constitutional disorders. Additional development and validation are needed to enable AOH/LOH detection in cancer samples.

Copy number variant calling performance: fractional CN analysis

For performance validation, we took two complementary approaches. We used simulated data and real data to assess CNV detection performance. Detailed methods are included in the Appendix A section.

Copy number variant calling performance using simulated data

CNV detection performance for events ranging from 175 kbp to 10 Mbp was assessed using *in silico* mixtures of simulated datasets. The new fractional copy number pipeline provides improved sensitivity to smaller events and to lower variant allele frequency events (Figure 23). The PPV for detection was close to 100% for all size bins and VAF levels tested.

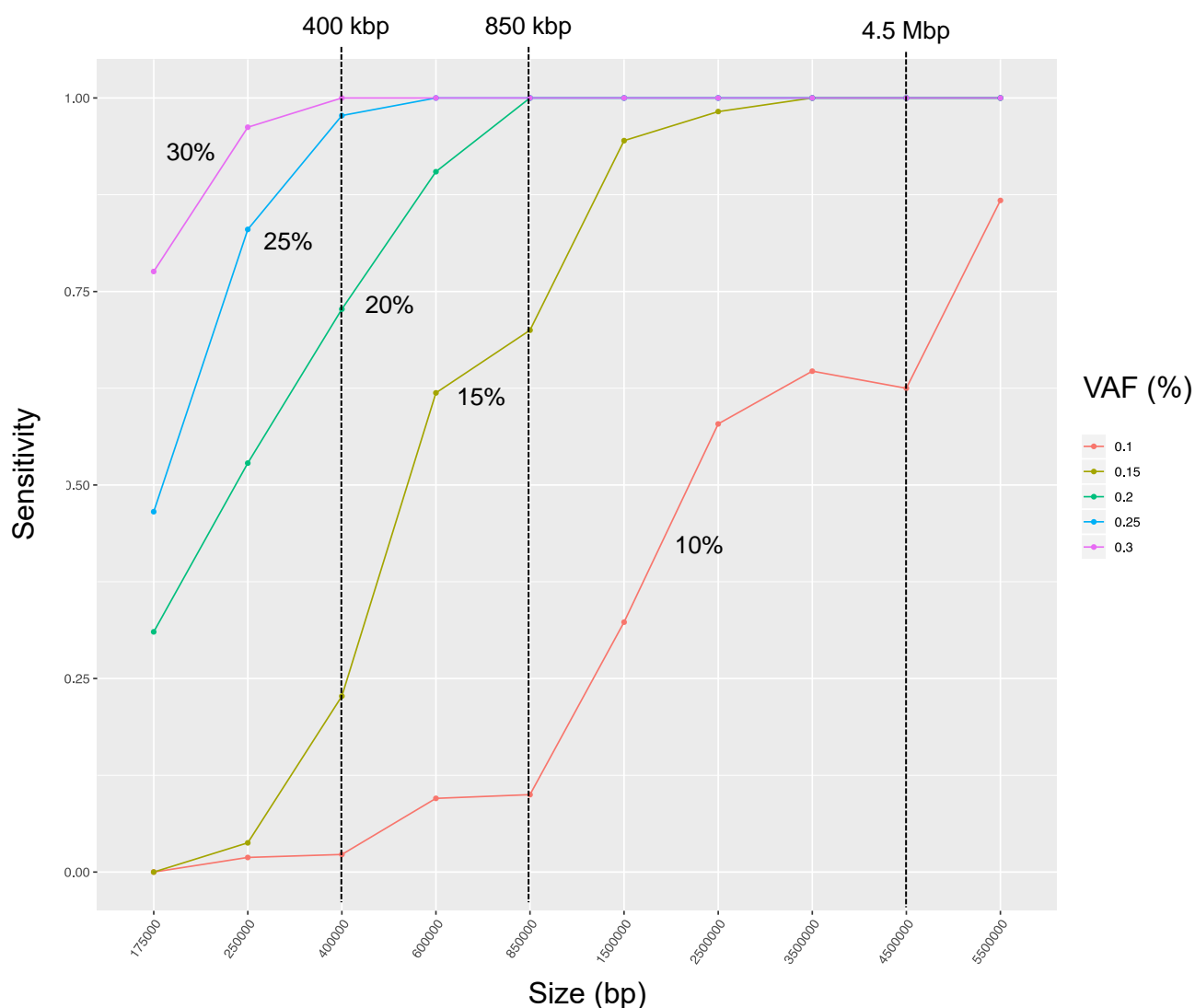


Figure 23. CNV detection sensitivity with simulated 300X DLE-1 data.

Breakpoint accuracy and VAF estimation accuracy are shown in Figure 24. The accuracy improves with variant allele frequency. Intuitively, changes in the copy number profiles are more obvious when VAF is higher, making it easier to infer breakpoint locations and VAF.

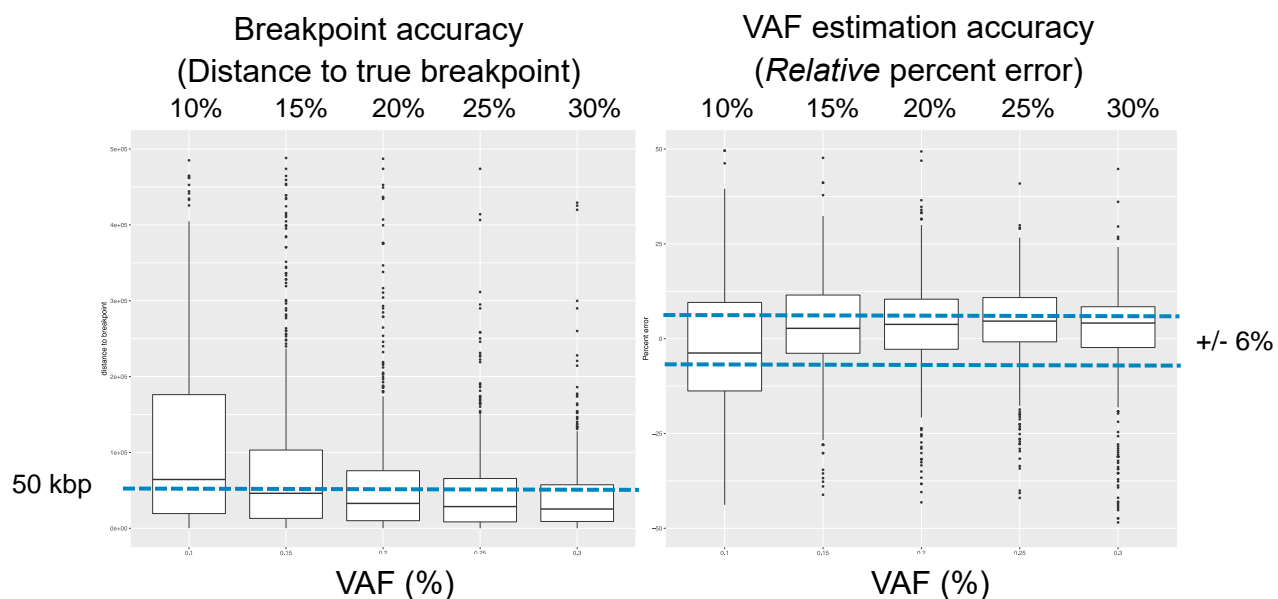


Figure 24. Breakpoint and VAF estimation accuracy with simulated 300X DLE-1 data.

Copy number variant calling performance using real data

Performance data are shown in Table 7 based on nine cancer samples (ranging from 200X to 300X coverage) showed high sensitivity to events at > 0.15 VAF. Detection of 0.1 VAF events was impacted by local noise in and complexity of CN profile.

Table 7a. Fraction CNV detection sensitivity with real DLE-1 datasets (top: duplications; bottom: deletions).

Duplications				
Size/AF	0.1	0.15	0.2	0.3
0 – 500 kbp	0.5	0.68	0.79	0.91
500 kbp – 1 Mbp	0.5	0.9	0.9	1
1 – 3 Mbp	0.69	0.92	0.98	1
3+ Mbp	0.78	0.97	0.97	0.97
Deletions				
Size/AF	0.1	0.15	0.2	0.3

0 – 500 kbp	0.85	0.85	0.87	0.95
500 kbp – 1 Mbp	0.97	0.97	0.97	1
1 – 3 Mbp	0.95	0.98	0.98	0.98
3+ Mbp	0.53	1	0.84	0.95

Table 7b. Fraction CNV detection PPV with real DLE-1 datasets.

VAF	PPV
0.1	0.98
0.15	0.98
0.2	0.97
0.3	0.97
1	1.00

Copy number variant calling performance: integer CN analysis

Copy number variant calling performance using simulated data

The integer CN pipeline is automatically run when the pipeline detects that the input is either a BspQI or BssSI dataset. High sensitivity and PPV were observed for events > 500 kbp. Generally, increasing coverage increases the signal-to-noise ratio; however, sensitivity for simulated events did not significantly increase beyond 80X (Figure 25). Higher PPV was observed at lower coverage levels, but sensitivity was significantly lower at those coverage levels.

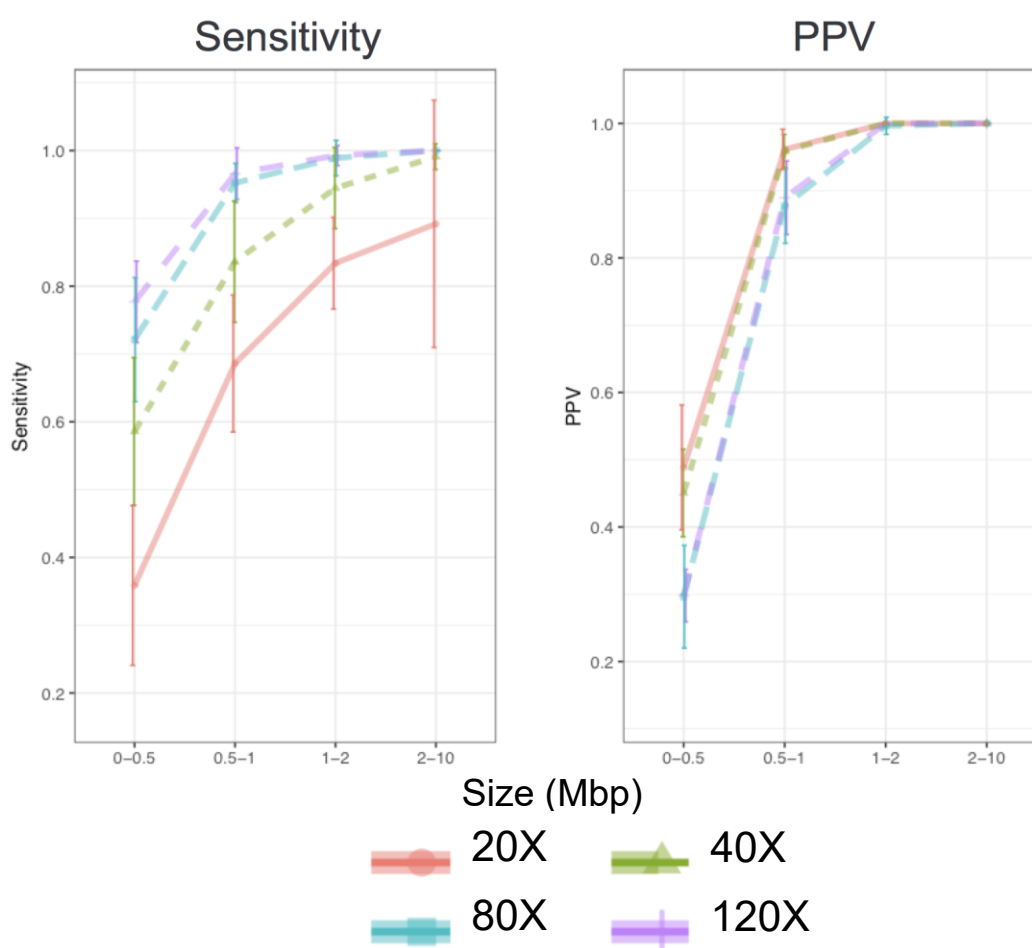


Figure 25. Copy number variant calling performance at different coverage levels with simulated DLE-1 data.

Copy number variant calling performance using real data

Copy number events of different sizes and zygosity states were expected in five well-characterized samples (Table 8). The smallest expected event (160 kbp) was not detected in either BspQI or BssSI data; all other events were detected.

Table 8. CNV calling performance from real datasets (at close to or less than 80X effective coverage).

	Chr	Size	Type	Zygosity	Detected in Nt.BspQI?	Detected in Nb.BssSI?
Sample 1	chr2	160 kbp	Duplication	Heterozygous	No	No
	chr2	380 kbp	Duplication	Heterozygous	Yes	Yes
Sample 2	chr3	540 kbp	Duplication	Heterozygous	Yes	Yes
Sample 3	chr6	1 Mbp	Deletion	Heterozygous	Yes	N/A
Sample 4	chr17	18 Mbp	Duplication	Homozygous	Yes	Yes
	chr8	75 Mbp	Duplication	Homozygous	Yes	Yes
Sample 5	chr6	46 Mbp	Deletion	Heterozygous	Yes	N/A
	chr18	65 Mbp	Duplication	Heterozygous	Yes	N/A
	chr11	67 Mbp	Duplication	Heterozygous	Yes	N/A
	chr1	106 Mbp	Duplication	Heterozygous	Yes	N/A

Copy number variant calling performance: chromosomal aneuploidy detection

Chromosomal aneuploidy detection performance using simulated data

The aneuploidy analysis module has been updated in Solve 3.7. Performance is shown in Table 9. Sensitivity and PPV for detecting chromosomal aneuploidy events were close 100% for events > 10% VAF.

Table 9. Chromosomal aneuploidy performance with simulated 300X DLE-1 data for events at 10-30% allele frequency.

	VAF	Sensitivity	PPV
Duplications	0.1	0.96	1
	0.2	1	1
	0.3	1	1
Deletions	0.1	0.94	1
	0.2	1	1
	0.3	1	1

Chromosomal aneuploidy events with low allele frequency were simulated in 48 samples with 65-80x coverage. This was done by adjusting the coverage up or down by a given percentage for each chromosome. 5% allele frequency aneuploidy was detectable in most chromosomes; however, sensitivity was lower for smaller chromosomes: sensitivity was below 90% for chromosomes 16, 17, 19, 20, 21, 22, X, Y (Figure 26). At 9% AF, all chromosomes except chrY had sensitivity >95%. There were false positive aneuploidy calls on chromosome 19 for 3 out of 48 samples; all other chromosomes had 100% precision at all allele frequency levels tested. Results were similar gains and losses.

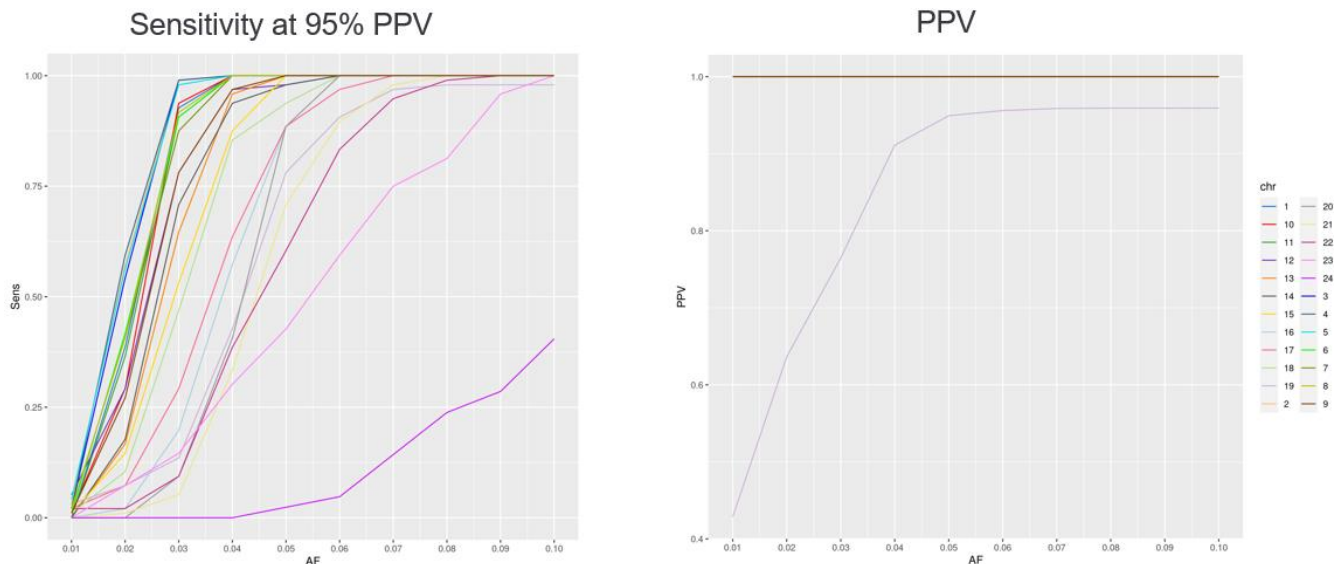


Figure 26. Chromosomal aneuploidy detection performance at different allele frequency levels for simulated aneuploidy events in samples with 65-80x coverage.

Chromosomal aneuploidy detection performance using real data

Entire-chromosome aneuploidy events were detected in 7/7 constitutional samples with known events (Table 10). Single-arm events are not detected. Cancer samples had lower sensitivity and some had high number of false positive calls.

Table 10. Chromosomal aneuploidy performance with constitutional and cancer 300X DLE-1 data.

SampleType	Annotation	TP	FN	Comments
Constitutional	trisomy 21	1	0	Called
Constitutional	trisomy 13	1	0	Called
Constitutional	trisomy 13	1	0	Called
Constitutional	mosaic Y loss	1	0	Called
Constitutional	XYY	1	0	Called
Constitutional	XXX	1	0	Called
Constitutional	trisomy 13	1	0	Called
Cancer	chr 6 loss chr 14 gain	2	0	Called
Cancer	chr12 gain	1	0	Called

Cancer	chr 19 loss chr 20 gain	1	1	chr20 gain called; chr19 loss not called. Cancer sample shows large coverage variation that complicates calls
Cancer	chr 3, chr 4 and chr 6 duplicated	3	0	Called
Cancer	whole-chr 4 deletion	1	0	Called most chromosomes

Bionano Access integration

With Bionano Access, users are able to select from the dropdown menus the set of assembly parameters for de novo assembly, depending on the application. Bionano Access enables import, filtering, and visualization of SV and copy number variant calls. Users can import a mask file to annotate insertion and deletion calls and filter out likely false positive translocation calls. We provide different versions of the mask to support the (hg19 and hg38) human reference assemblies and the enzymes recommended for human analysis (Nt.BspQI, Nb.BssSI, and DLE-1). We also provide masks for the mouse reference assembly mm10 and the DLE-1 enzyme. The sets of masks for the de novo assembly pipeline and the RVP are separately maintained. Users can also import custom gap files.

VCF conversion

The python-based VCF converter supports conversion of both SV and CNV variant calls. The resulting VCF output is compliant with VCF format version 4.2. For more information, please refer to Appendix G and Bionano VCF File Format Specification Sheet (PN 30459). For insertion and deletion calls, the SV breakpoints are output as the midpoint between the last aligned label and the next label. The uncertain breakpoint is indicated in the CIEND and CIPOS fields. It can take SMAP files and CNV output from the *de novo* assembly pipeline, Rare Variant Pipeline, and Variant Annotation Pipeline. As of Solve 3.7, variant annotations are included in the VCF output.

Consistent with the VCF v4.2 format, the VCF file output passed testing with vcf-tools/vcf-validator. The information that is added to the output includes the reference accession, orientation information for translocations, and additional annotation.

Additional Considerations

SV sensitivity generally improves with depth of coverage. The minimum recommended effective coverage (defined as the product of raw coverage and molecule alignment rate) for a diploid sample is 80X, assuming that the quality and average length of the input molecules are good. More coverage will yield small additional sensitivity for homozygous and heterozygous SVs with diminishing returns beyond 120X.

For mosaic/heterogeneous samples, RVP should be used for detection of low allelic fraction SVs. Data provided in this document generally uses the output of a single Saphyr 2x1300 or 3x1300 flow cell, i.e. 1.3 Tbp. As of the Bionano Solve 3.6 release, we have updated our recommendation to collect 1.5 Tbp of data to ensure that the RVP provides reliable 5% allele fraction performance even with sub-optimal samples. Additional coverage can further improve sensitivity, we currently do not recommend inputting more than 5 Tbp.

The Bionano Solve pipeline was validated on human diploid and heterogeneous samples.

FAQs

1. How does coverage affect heterozygous and homozygous SV calling performance using the *de novo* assembly pipeline?

Based on consideration of SV calling performance, runtime, and system throughput, we recommend a minimum of 80X effective coverage for haplotype-sensitive assembly for all SV types. Sensitivity does increase with coverage, but PPV is comparable across coverage levels (the lowest coverage level tested was 50X effective coverage).

2. What is the sensitivity to small insertions and deletions?

We recommend that users focus on insertions and deletions larger than 500 bp.

3. Why do alignment boundaries appear to be off sometimes?

Based on visual inspection, we noticed cases where there was room for improvement in the accuracy of the alignment boundaries. Even though we correctly detected that there was an event, there was an impact on the SV boundaries for these calls. We have optimized parameters to improve the accuracy of SV boundaries and maintain overall performance.

4. Do we classify translocations?

A translocation can be balanced or unbalanced, and reciprocal or non-reciprocal. We currently detect single translocation breakpoints and do not attempt to further classify them. Also, orientation information is currently encoded in SMAP. Users may infer the orientation when visualizing the calls in Access. Additional custom secondary analysis may aid pairing and classification of translocation breakpoints.

5. What are the different versions of cluster parameters and assembly parameters files?

Different versions of cluster parameters and assembly parameters files were created for supporting different platforms and different applications. Cluster parameter files are distributed with the python assembly pipeline; assembly parameter files are distributed with RefAligner. To facilitate selection, dropdown menus are available in Access when an assembly is set up.

6. How are reciprocal translocation breakpoints detected?

Each translocation breakpoint of a reciprocal translocation is independently detected. We currently do not pair potential reciprocal translocation breakpoints. However, translocation orientation information is provided; translocation breakpoints in the same region and with opposite orientations may be part of a

reciprocal translocation.

7. How does the assembler handle ambiguity associated with segmental duplication regions?

Large segmental duplication regions, a type of CMPRs, appear at least twice in the genome and are connected to different sequences (potentially on different chromosomes). If a given segmental duplication region is larger than an average molecule's length, the assembler likely does not have enough information to resolve the inherent ambiguity in the connectivity. Maps containing CPMRs larger than 140 kbp are split by default; there is also the option to not split those maps.

8. How are masked translocation breakpoint calls annotated?

A suffix (“_common” or “segdupe”) would be appended to any masked translocation breakpoint call in the SMAP SV detection output. A suffix of “_common” indicates that the call overlaps with translocation breakpoints detected in genomes not known to contain translocations. A suffix of “_segdupe” indicates that the call overlaps annotated segmental duplication regions. If a breakpoint overlaps with both types of masked regions, “_common” would be appended.

9. How do I create a custom mask?

Users interested to generate custom masks are encouraged to model the custom masks based on the masks that Bionano provides. The Bionano masks are in BED format (<https://genome.ucsc.edu/FAQ/FAQformat#format1>). They can be opened and edited using a text editor. They contain three types of BED entries. The three sources of data were pre-processed separately and concatenated together to generate the final mask BED files. The data include 1) N-base gaps, 2) annotated segmental duplication regions, and 3) translocation breakpoint regions detected in control genomes not known to contain translocations. The chromosomal locations of N-base gaps can be obtained from analyzing a given reference fasta file and outputting locations of the N bases (this is output by the *in silico* digestion tool included in Bionano Access). The list of annotated segmental duplication regions (> 50 kbp) for human can be obtained from the UCSC Genome Browser database. When we pre-processed the regions, we further verified map-level similarity for these regions. Using a database of control samples, we defined a list of translocation breakpoints that we detected in more than a specified number of genomes.

10. How are translocation and inversion confidence scores different from insertion and deletion confidence scores?

Confidence scores for insertions and deletions are computed as PPV estimates; the confidence scores for inversion and translocation breakpoints are computed as probabilities from the machine learning models. Because they were trained using different methods and training data, they require different

thresholds.

11. What is the expected runtime performance?

Below are representative runtime data for a Saphyr Compute running the pipelines using example human datasets.

	Enzyme	Collected data	Effective coverage	Saphyr compute	Pipeline
Sample 1	DLE-1	400 Gbp	80X	12-16 hr	<i>De novo</i> assembly
Sample 2	DLE-1	1.5 Tbp	300X	8-12 hr	RVP
Sample 3	DLE-1	1.5 Tbp	300X	19 hr	<i>De novo</i> assembly

12. Why do I see fewer maps in regions where I've seen many more previously? Am I losing any information?

We have optimized parameters for allele separation. We expect that homozygous regions are now more likely to be represented by a single allele map. Heterozygous regions are expected to be represented by allele maps containing allele-specific SVs.

13. Could I take advantage of the new Solve release with my existing assemblies?

The SV detection step may be re-run on the command line using the new Solve tools in order to take advantage of improved SV calling performance. After SV detection, one would replace the existing SV output directory with the new SV output directory, compress the entire output, and import into Access for SV visualization and downstream analysis.

The computation of the inversion and translocation breakpoint confidence relies on data only generated in more recent versions of the Solve pipeline. If expected data were not available in the input consensus maps for SV detection, “-1” would be output for inversion and translocation breakpoint confidence.

14. How do I take advantage of having both NLRS and DLS data?

Data generated with a second enzyme can provide a valuable, independent validation of SVs of interest. Although Bionano has not developed and tested a particular workflow, one common workflow is to first discover relevant SVs using DLS and the Variant Annotation Pipeline, as well as, potentially, other data and tools. Then, generate a second dataset and *de novo* assembly for any samples that require validation. It should not be necessary to pass the data through the discovery pipeline. Simply query the SV list from the confirmation sample to find a match. Key considerations are that the SV coordinates may not match exactly because the coordinates are based on the enzyme recognition sites and have a

median accuracy of ~3-5 kbp. A second key consideration is that because of coverage dropout from NLRs-produced “fragile sites,” confirmation could fail, i.e. be negative. That result should not be considered a refutation of the SV being tested. Users are recommended to use Bionano Access to inspect the region of interest to determine whether there is sufficient coverage and whether there is any SV “signal” that was not called but could partially confirm the SV (i.e. the map diverging from the reference at the SV location).

15. How close are the detected breakpoint coordinates to the actual coordinates?

The detected breakpoints are typically around 3 kbp from the actual breakpoints (with the 90th percentile at around 11 kbp).

16. How are molecule alignment confidence scores calculated?

Confidence scores are useful for evaluating alignments. We estimate the probability that the labels on a map match the labels on the reference purely by chance and that the maps are unrelated. The scores are calculated as $-\log_{10}$ of the probability. Generally, alignments with higher scores are of higher confidence.

17. How do I interpret SV counts in the informatics reports?

The informatics report displayed in Access contains counts of clustered SV calls. Sometimes, the same SV may be present in more than one map. For example, a homozygous SV may be present in both allelic maps. We cluster similar SV calls in order to more accurately estimate the number of SVs in the assembly. In the full informatics report (part of the compressed assembly output), SV counts for both pre- and post-clustering are available.

18. How do I run the VCF converter?

The VCF converter is packaged with Bionano Solve. It can be run on the command-line; it is typically run as part of the analysis pipelines.

19. How does the *de novo* assembly pipeline handle high-coverage datasets?

To efficiently assemble high-coverage datasets and to make use of the best-quality data, an adaptive downsampling procedure is applied by default (could be turned off via RefAligner options specified in the xml files). Briefly, the “-minlen” parameter is increased up to 240 kbp to reduce coverage, assuming that the longer molecules are more informative. If coverage is still above 250X, molecules are further randomly subsampled until 250X is reached. In refineB and extension stages, molecules with lowest confidence alignments are not used if local effective coverage exceeds 80X. In the refineFinal1 stage, molecules with lowest confidence alignments are not used if local effective coverage exceeds 100X.

20. What is the maximum input coverage recommended for Rare Variant Pipeline?

300X effective coverage is recommended for typical low variant allele frequency applications, where events are expected at as low as 5% allele fraction. To achieve that performance with most samples, we recommend 1.5 Tbp to be collected on Saphyr. RVP has been validated with up to 1000X effective coverage (equivalent of 5 Tbp of data collected on Saphyr). We do not recommend running this pipeline with more than 1000X.

21. Why am I seeing small deletions that seemed suspicious (weak molecule support, for example) in the RVP?

We observed that with RVP, PPV was slightly lower (~80%) for deletions under 25 kbp. This was not observed with the *de novo* assembly pipeline.

22. Where can I get information about variant allele fraction?

The allele fraction for each SV call is part of the Bionano Access UI output and can also be found in the output SMAP file, under the column VAF.

23. What is new in the copy number analysis pipeline?

As of Solve 3.6, the fractional copy number analysis module includes a region-based scoring method that handles complex regions better. Taking advantage of the availability of control data, the pipeline is able to avoid making false positive calls in regions where the coverage variation is high even in the absence of true events.

24. When should I use the *de novo* assembly pipeline and when should I use Rare Variant Pipeline? Should I run both?

This depends on the sample type and availability of compute resources. We recommend running just the *de novo* assembly pipeline (80X) for germline or “typical” diploid samples. For cancer or somatic samples, there may be three options (ordered by increasing sensitivity and computational costs):

1. RVP only (1.5 Tbp / 300X effective): 8-12 hours
2. Assembly (400 Gb / 80X effective) + RVP (1.5 Tbp / 300X effective) for cancer or somatic samples: 12-16 hours + 8-12 hours
3. Assembly (1.5 Tbp / 300X effective) + RVP (1.5 Tbp / 300X effective) for cancer or somatic samples: 19 hours + 8-12 hours

Note: For most applications, we do not recommend running an assembly at >80X effective coverage as the number of new SVs will not increase substantially.

25. Why are effective coverage estimates different between MQR and the assembly report?

Slightly different alignment parameters are used for MQR and during the assembly pipeline, so the effective coverage estimates may differ. The coverage values used in this document refer to the MQR estimates. Please refer to 30173 Data Collection Guideline for more details.

26. Why do I get different results when analyzing the same sample?

There is some non-determinism in the analysis software so that repeated analyses of identical input may result in slightly different outputs. This is due to randomness in the de novo assembly and haplotype assignment processes. With Solve 3.7 we have tried to quantify the expected amount of variability (see the Analytical Repeatability section in Appendix A). We observed 7 differences in clustered structural variant calls between repeated runs of de novo assembly. RVP replicates had identical structural variant calls with small differences in the SVFreq calculation. Note that technical replicates of samples that include repeating lab procedures and data collection will yield additional differences as each step in the process introduces its own level of variation. We are working on quantifying the expected variance of these type of replicates.

Appendix A: SV Calling Performance Evaluation

RVP SV calling performance at lower coverage levels

We also tested RVP SV calling performance at lower coverage levels and different variant allele frequencies. Figures 27 and 28 show sensitivity performance at 80X and 160X effective coverage, respectively. The same coverage dependency was observed.

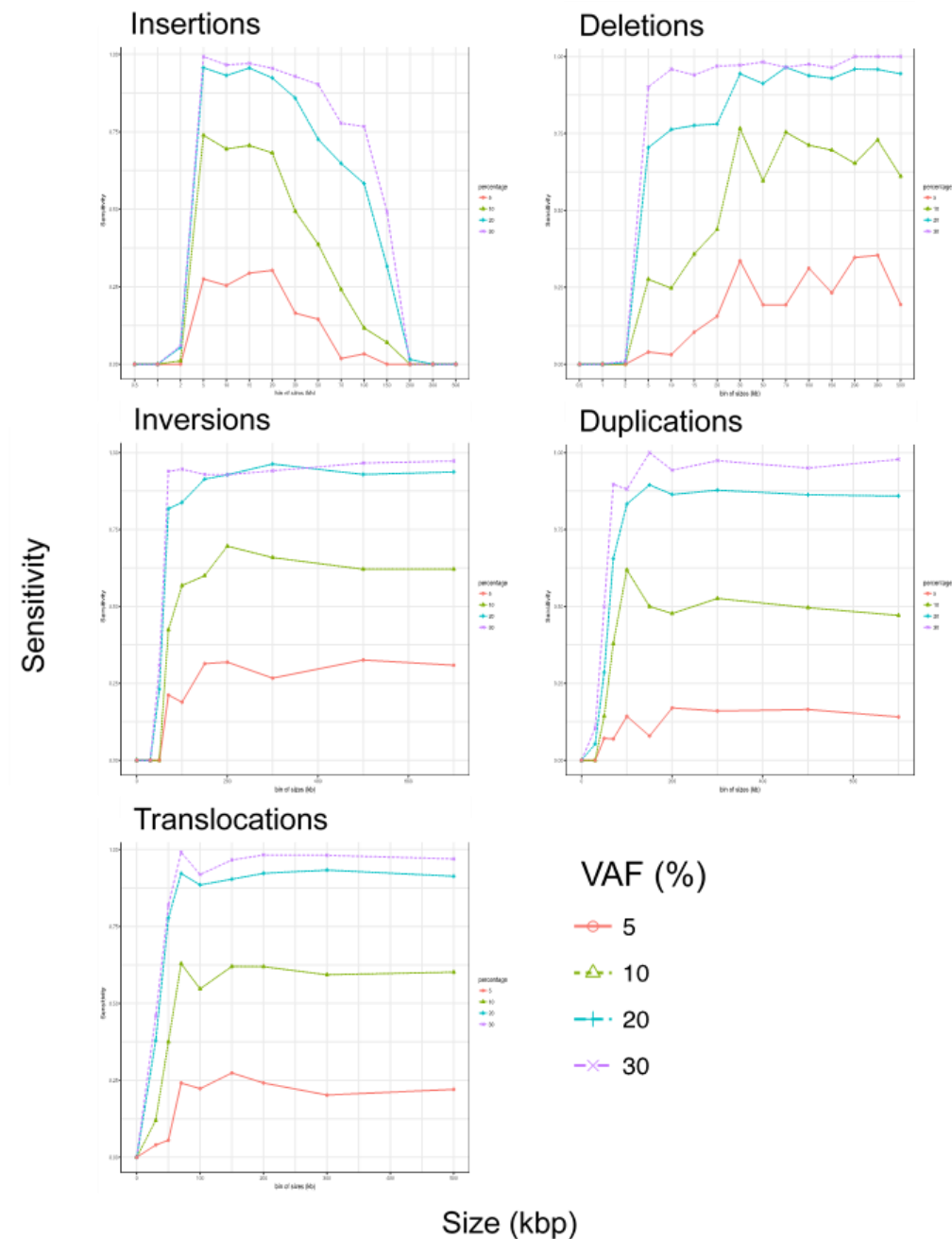


Figure 27. SV calling performance with simulated 80X DLE-1 data at different variant allele frequencies.

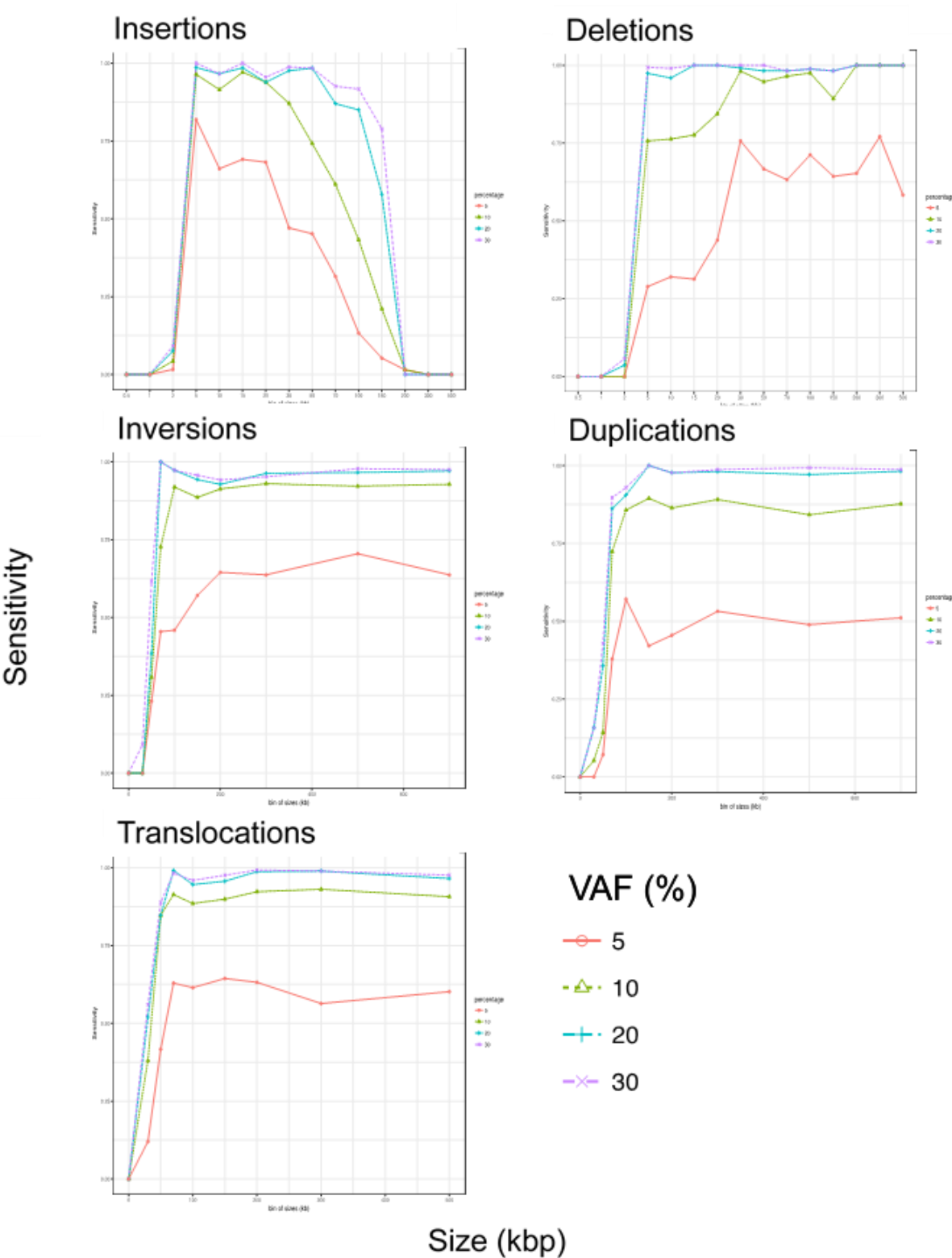


Figure 28. SV calling performance with simulated 160X DLE-1 data at different variant allele frequencies

Methods for assessing SV calling performance

SV calling performance with simulated data

We simulated random SV events so that we could estimate our genome-wide SV calling performance accurately. The human reference assembly hg19 was used as an “SV-free” base genome (in our performance analyses, SVs were called against hg19).

Simulation of insertions and deletions

We randomly introduced 1600 insertions and 1600 deletions into an *in-silico* map of hg19. The insertions involved addition of new simulated material (random sequences of defined sizes) and deletions involved removing material; we did not simulate replacement of sequence (substitutions). The simulated events were at least 500 kbp from each other or N-base gaps. They ranged from 200 bp to 1 Mbp, with smaller SVs more frequent than larger ones.

Based on the edited hg19, molecules were simulated to resemble actual molecules collected on a Bionano system. This entailed adding sizing error in accordance with the model used in RefAligner, outliers created by stitching and DNA knots and folds, and fragile sites. In addition, we simulated molecules from the unedited hg19 and mixed with molecules from the edited hg19 such that all events would be heterozygous. Four such datasets with effective coverage levels from 87X-110X were generated. The simulated molecules were used as input for the haplotype-aware assembly pipeline, and SV calls were compared to the ground truth.

Simulation of translocations

Random translocation events were simulated to form an edited genome (Figure 29). 918 segments were randomly selected across hg19 and randomly inserted elsewhere into the genome. We avoided N-base gaps, and the size of the translocation fragments ranged from 50 kbp to 1 Mbp. Breakpoints were at least 500 kbp away from each other. For simulated intrachromosomal fusions, the breakpoints were at least 5 Mbp away. Each translocated fragment is expected to generate two translocation breakpoints. These can also be considered transpositions.

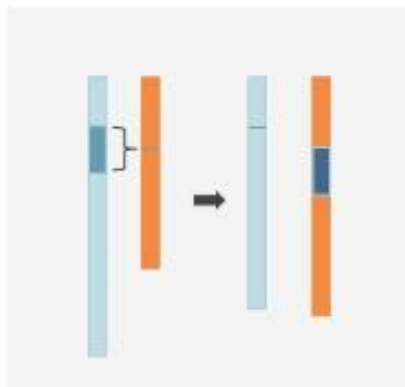


Figure 29. Simulation of translocated fragments. Random fragments are removed from a donor chromosome and inserted into an acceptor chromosome.

Simulation of inversions

About 900 intervals of 5 kbp to 1 Mbp in size were randomly sampled across hg19 and inverted to create an edited genome with simulated inversions. The inversion events were at least 500 kbp away from each other and were required to not overlap with N-base gaps.

Simulation of duplications

About 900 intervals of 5 kbp to 1 Mbp in size were randomly sampled across hg19. For each sampled interval, an extra copy of the sequence was inserted in tandem next to the original segment. The new copy could either be in the same or opposite orientation such that performance for detecting tandem and inverted duplications could be assessed. The duplication events were at least 500 kbp away from each other and were required to not overlap with N-base gaps.

Simulation of molecules

Molecules were simulated from edited genomes according to empirically derived error and size characteristics. Error-free molecules were simulated; then, errors (such as sizing errors and FP and FN labels) were added. Simulated molecules were used as input for de novo assembly and SV detection.

Simulation of large structural variants

Data simulation

We validated the large SV calling (>200kb) by simulating 726 duplication, deletion, inversion, and intra-chromosomal fusion SVs with varied sizes in hg19 human reference (Table 11).

Type	SVs Simulated	Size	Distance between copies	Distance between SVs
Duplication	190	200 kb - 1Mb	300 kb - 7 Mb	500 kb-1Mb
Deletion	175	5 Mb - 10Mb	NA	500 kb-1Mb
Inversion	175	5 Mb - 10Mb	NA	500 kb-1Mb
Intra-chromosomal fusion	186	30 kb - 7Mb	< 1Mb	500 kb-1Mb

Table 11. Simulation Summary

Analytical Performance

We measured the sensitivity and precision of large SV calling with 20% size tolerance and 20 kb breakpoint error and showed that the sensitivity and precision across all types are > 95% and > 85% respectively (Table 12).

Type	Sensitivity (%)	Precision (%)
Deletions > 200 kbp	98.5	87.4
Inversions > 200 kbp	98.3	97.6
Duplications > 200 kbp	98.1	86.1
Intra-chromosomal fusion > 200 kbp	95.5	85.1

Table 12. Large SV calling sensitivity & precision per variant type

Terminal deletions

Solve 3.7 adds preliminary support for detection of deletions at the terminal end of chromosomes for human samples only.

Data simulation

We validated this module by simulating terminal deletions using the hg38 human reference. This simulation involves 771 synthetic homozygous terminal deletions of sizes > 50Kb in the genomic regions corresponding to 13 terminal deletion syndromes.

Analytical Performance

We measured the sensitivity of our method per deletion size and per syndrome and showed that the average sensitivity across the syndromes is 95.6% with average sensitivity of > 90% on deletions of size >100Kb. The average precision across all sizes and syndromes is 86.4%. Please see Tables 13, 14 for detailed results.

syndrome	simulated	called	sensitivity
1p36 microdeletion syndrome	63	58	92
2q37 deletion syndrome	63	63	100
3q29 deletion syndrome	54	54	100
4p deletion syndrome - Wolf-Hirschhorn	63	57	90
5p deletion - Cri-du-Chat	50	46	92
6p25 deletion syndrome	54	49	91
9q34.3 sub telomere deletion syndrome	55	55	100
15q26 overgrowth syndrome	63	63	100
22q13 deletion syndrome - Phelan	63	59	94
Cat-eye syndrome	63	61	97
Xq28 microduplication syndrome	63	62	98
11q deletion - Jacobsen	63	62	98
1q terminal deletion	54	49	91

Average	<u>95.6%</u>
---------	---------------------

Table 13. Sensitivity/Syndrome

	50 kb	100 kb	250 kb	500 kb	1 Mb	5 Mb	> 8 Mb
simulated	46	135	117	117	117	144	63
called	36	123	115	117	117	144	63
sensitivity (%)	78	91	98	100	100	100	100

Table 14. Sensitivity per deletion size

Artifacts of the simulator

In Table 3, we explain some artifacts of our molecule simulator affecting the analytical performance of the terminal deletion caller module. Our simulator currently generates low coverage data towards the end of the chromosomes, which limits our analysis in two ways: 1- we could only simulate and analyze the terminal deletions of size > 50kb and 2- the unrealistic low coverage in terminal deletions of size 50Kb causes low sensitivity of 78%. The other artifact is that the simulated molecules near the unlabeled regions of size >100Kb are not spanning through the region and hence forming an artificial group of lined up molecules which mimics an actual terminal deletion. This would affect the precision of the method; hence we carefully identified and removed these artifacts from the analysis.

Artifact	Effects on the analysis	Affects
Extreme low coverage towards the end of chromosome.	Insufficient coverage to call smaller deletion sizes Analysis restricted to > 50Kb	Sensitivity
Artificial pileup of molecules near unlabeled regions from > 100 kb	RefAligner incorrectly calls large terminal deletions in these regions	Precision

Artifacts of molecule simulator

Analytical Repeatability

We repeatedly ran *de novo* Assembly and RVP pipelines with the same input to evaluate the determinism in results.

RVP

We ran RVP pipeline five times with the same input and compared the outputs of the SV calling, CNV and VAP analysis in a pairwise manner between 5 runs (resulting in 10 different combinations). Based on our results the number of clustered variants called by RVP pipelines is the same in all runs. However, comparing the raw smaps showed that while entries in 4 of 10 pairwise comparisons were identical, in the other 6 the value of column

SVfreq differed by almost 0.001 in 0.24% of the entries. CNV and VAP outputs were identical in all runs.

De novo assembly

Similar to the RVP pipeline evaluation, we ran the *de novo* assembly pipeline five times with the same input and compared the outputs of SV calling, CNV and VAP analysis in a pairwise manner between 5 runs (resulting in 10 different combinations). Based on our results the number of clustered variants called by *de novo* Assembly pipelines differ by at most 7 clusters across all runs. Similarly, comparing the raw smaps showed that the number of called SVs was 7717 in 2 of 5 runs and 7734 in the other 3 runs, differing by 17 entries. We also compared the contents of the entries in the runs with same number of called SVs and found out that the values of columns Zygosity and VAF differ in at most 12 and 10 entries respectively. Moreover, while CNV outputs were identical across all runs, VAP outputs differ according to the differences between the smaps.

SV calling performance with real data

CHM1/13 analysis

Data was generated from homozygous CHM1 and CHM13 cell lines initially derived from hydatidiform moles. Single-molecule maps from CHM1 and CHM13 were evenly sampled and combined *in silico* to simulate a diploid genome at 80X effective coverage. The CHM1, CHM13, and the CHM1/13 mixture molecule sets were assembled separately.

A three-way SV comparison analysis was performed for each trio set of CHM1, CHM13, and CHM1/13 mixture SV calls. The SV calls from the CHM1 and CHM13 pure assemblies were considered as the (conditional) ground truth. The sensitivity and positive predicted value (PPV) at different coverage levels were analyzed. Sensitivity was defined as fraction of SV calls in pure assemblies that were called in the mixture assembly, and PPV was defined as fraction of calls in mixture assembly that were called in the pure assemblies.

CEPH trio analysis

Using genome mapping, Mak *et al.* analyzed a Caucasian trio from the 1000 Genomes Project (the parents NA12891 and NA12892, and the daughter NA12878) and published an expert-curated SV list. We reanalyzed the same starting data and compared the resulting SV list from the automated pipeline against the curated list.

Method for assessing CNV calling performance

CNV calling performance with simulated data

The simulations are used for performance validation and for constructing confident tables. The methods are described in Introduction to Copy Number Analysis (PN 30210). Briefly, we randomly introduced copy number events across the genome. Starting from a molecule-to-reference alignment, molecules that overlap with

simulated copy number events are sampled accordingly based on their simulated copy number states. The sampled molecules were used as input to the copy number analysis pipeline, and the detected events were compared with the simulated ground truth to derive sensitivity and PPV data.

CNV calling performance with real data

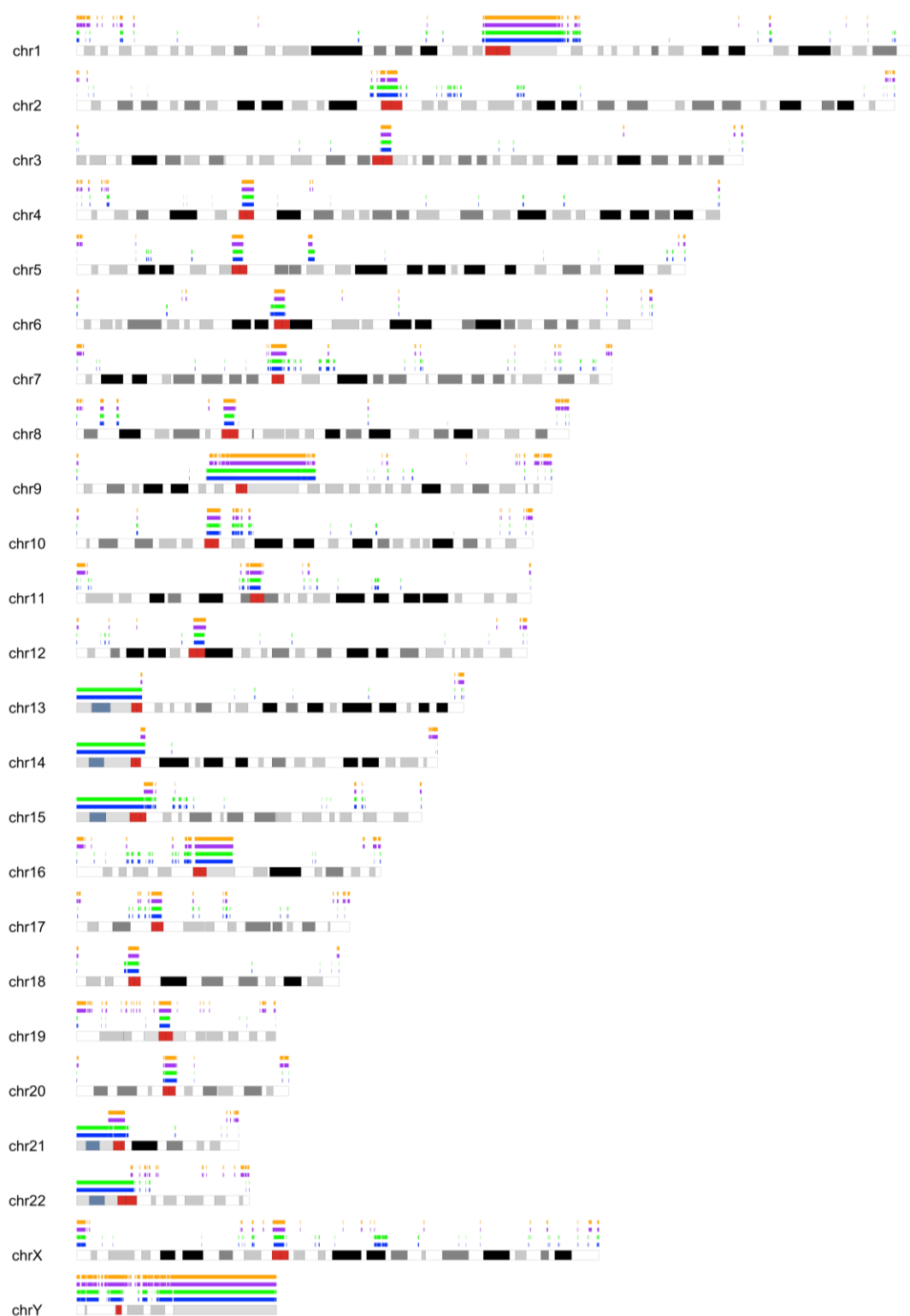
Samples with known copy number events were analyzed with BspQI, BssSI, and/or DLE-1. The datasets are analyzed using the Solve pipeline, and the copy number analysis output were assessed based on coordinates where CNV events were expected.

Appendix B: Mask generation

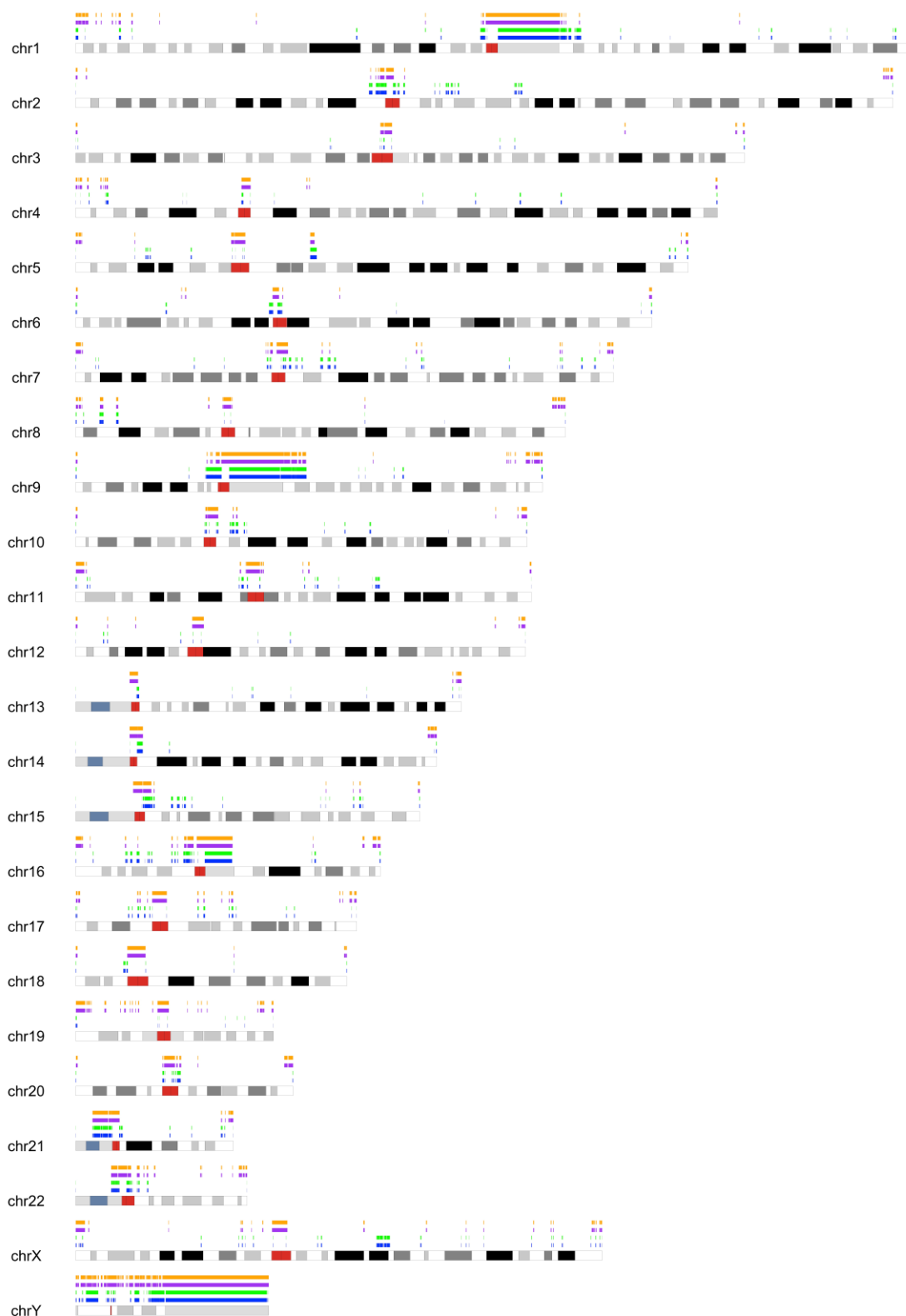
Comparing Solve 3.6.1 and 3.7 masks

The ideogram plots below were generated based on the SV and CNV mask BED files. The top two tracks are the CNV masks from Solve 3.6.1 (purple) and 3.7 (orange). The bottom tracks are the SV masks from Solve 3.6.1 (blue) and 3.7 (green). The same CNV masks are used for both the *de novo* assembly pipeline and RVP.

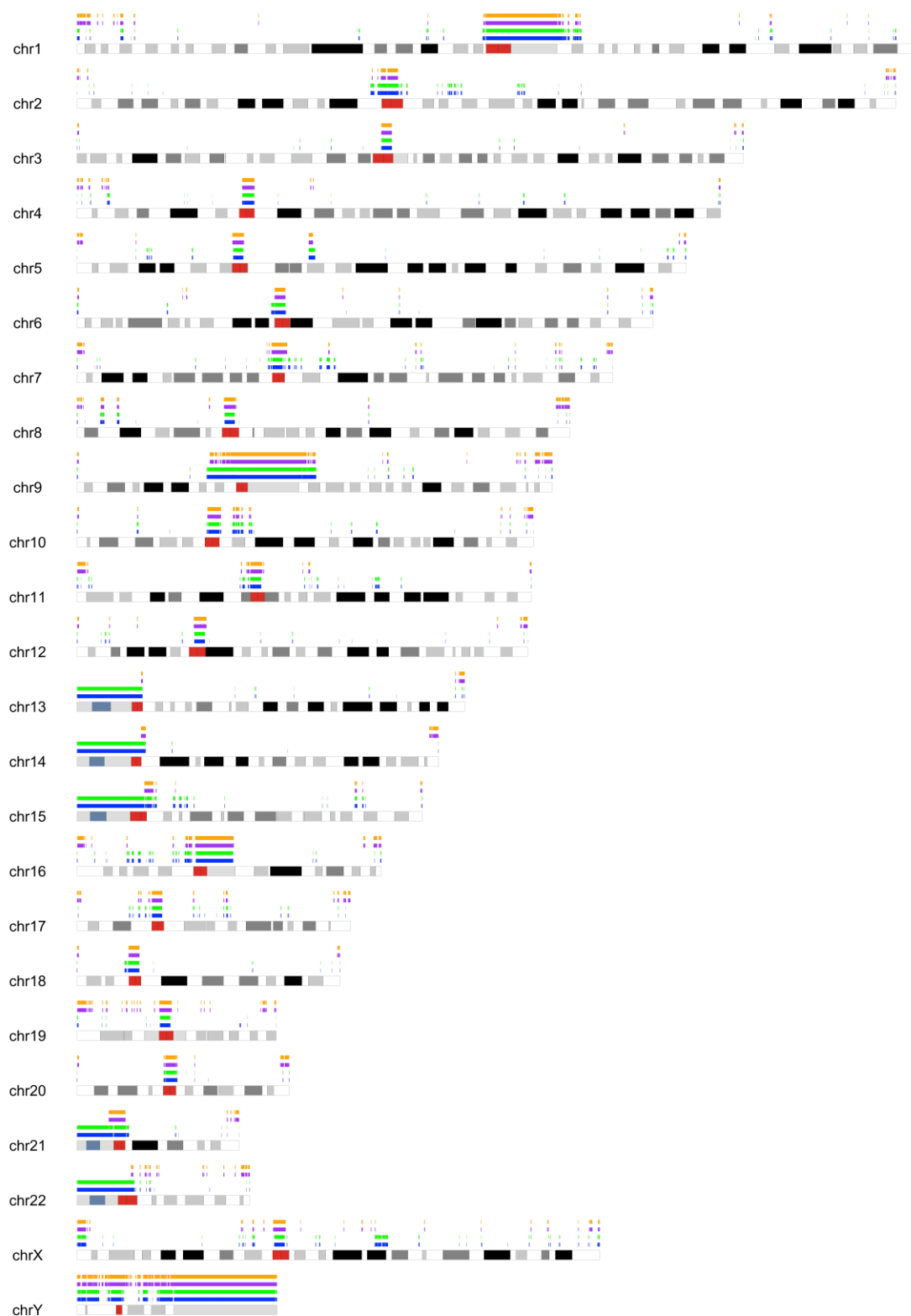
hg19/*de novo* assembly masks:



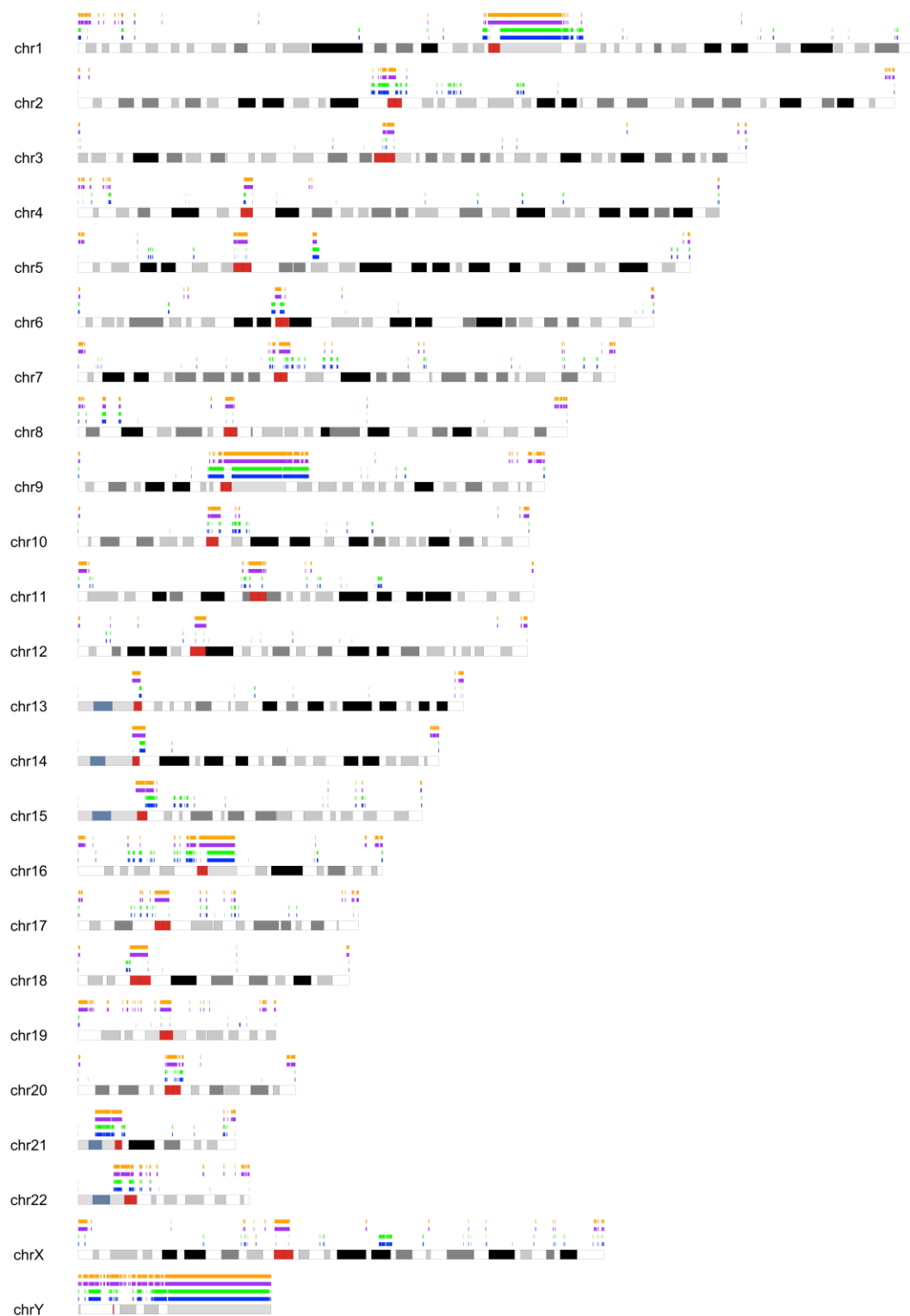
hg38/*de novo* assembly masks:



hg19/RVP masks:



hg38/RVP masks:



Compiling common translocation breakpoint list

Genome map assemblies for control human samples were used as input. Given that they were phenotypically normal samples, we assumed that the control samples did not contain translocations. We performed SV detection on each assembly, and detected translocation breakpoints were compiled.

Common translocation breakpoint calls were merged and included into the masks if they were observed in 10 or more samples in the control database. Two neighboring breakpoints were merged if they were within 25 kbp of each other. After merging, for each breakpoint or merged region, a 10-kbp buffer was added on both sides. These common breakpoint regions were labeled as “common” in the masks.

Compiling annotated segmental duplication regions

Annotated segmental duplication regions of at least 50 kbp were compiled from the UCSC Genome Browser database. Sequences from the annotated segmental duplication regions were extracted and converted into *in silico* maps. For each pair of segmental duplication regions, we checked for map-based similarity. If the *in silico* maps aligned with each other with a p-value of less than 1E-4, the segmental duplication region would be included in the mask and labeled as “segdupe”.

Tools for generating mask BED files

In Bionano Solve 3.6, we provide tools for generating custom mask BED files. They are contained in the “process_control_datasets” directory in the Solve build. There are three R scripts in the directory, containing functionalities for creating mask BED files, and checking and applying these masks to an SMAP file. These scripts require the R package “maskTranslocation_0.0.1.0.tar.gz” also included in the directory; they were tested with R version 3.6.3. Selected parameters are described below.

- 1) “makeTranslocationMasks.R” is used to generate mask BED files based on a set of control samples.

Example command: `Rscript --vanilla path-to-makeTranslocationMasks.R --inputFile inFile --outputDir outDir --reference mm10 --enzyme DLE1 --bedtoolPath path-to-bedtools --minCount 10`

“inputFile”: a CSV file with the first column containing the names of the samples and the second column the paths for their SMAP files.

“reference”: “hg19” or “hg38” for human samples and “mm10” for mouse samples.

“enzyme”: “DLE1”, “BSPQI” or “BSSSI” for human samples, and “DLE1” for mouse samples.

Two other parameters “segdupPath” and “gapPath” are not required for “hg19”, “hg38” or “mm10” since they were included in the R package. However, if the gap and segmental duplication annotations are

updated on the UCSC Genome Browser, they can be extracted and prepared for this input with the script "make_bed4masks_translocation.R".

- 2) "check_set_transMask2Samp.R" is used for either checking masked calls in an SMAP file overlapping with chromosomal locations in a mask BED file, or masking calls in an SMAP file with an input mask BED bed file.

Example command: Rscript --vanilla path-to-check_set_transMask2Smap.R --smapFile smapFile --maskFile maskFile --outputDir outDir --check[/set]

"check" and "set" are switches and cannot both be presented in the same command.

- 3) "make_bed4masks_translocation.R" is a utility script for preparing gap and segmental duplication annotations to be used by the script "makeTranslocationMasks.R".

Example command: Rscript --vanilla path-to-makeTranslocationMasks.R inputFile --outputDir outDir --gap -- refName referenceName [--segdup --refName referencename --refPath path-to-reference --RefAlignerPath path-to-RefAligner]

"inputFile": text file downloaded from the UCSC Genome Browser, with gap or segmental duplication annotations for a particular reference genome.

"gap" and "segdup" are switches and cannot both be presented in the same command. "refName" is a shortened name for the reference genome in the parameter "refPath". They could either be "hg19", "hg38", or "mm10". The two other parameters "segdupSize" and "confidence" are for the parameter "segdup". Their default values (50 kbp and 5) are optimized for the human references (hg19 and hg38).

Appendix C: Descriptions of output files

The following sections detail output files contained in the compressed *de novo* assembly pipeline and RVP output. A subset of the files is used for Bionano Access visualization. Some are considered intermediate files that may be helpful for troubleshooting.

Output files from the *de novo* assembly pipeline

Here is an example of the output directory structure after unzipping the *de novo* assembly results (*pipeline_results.zip):

```
/output/  
  
/output/contigs/alignmolvref/copynumber  
  
/output/contigs/auto_noise  
  
/output/contigs/exp_refineFinal1/  
  
/output/contigs/exp_refineFinal1/alignref_final  
  
/output/contigs/exp_refineFinal1_sv/merged_smaps
```

In the *output* directory:

exp_informaticsReport.txt – summary statistics of the assembly stages. Final assembly statistics are in the section "Stage Summary: CharacterizeF_refineFinal1",.

exp_informaticsReportSimple.txt – a simplified version of exp_informaticsReport.txt. This is shown as the *de novo* assembly report in Bionano Access.

exp_informaticsReportSimple.json – contents of exp_informaticsReportSimple.txt in JSON format.

exp_optArguments.xml – parameters used in the assembly.

exp_pipelineReport.txt – assembly pipeline log file.

status.xml - assembly pipeline status log.

In the *output/contigs/alignmolvref/copynumber* subdirectory:

cnv_calls_exp.txt – unfiltered copy number calls. This is shown in the Copy Number tab in Bionano Access.

cnv_chrAneuploidy_exp.txt – per-chromosome aneuploidy calls. This is shown in the Aneuploidy tab in Bionano Access.

cnv_rcmap_exp.txt - coverage data for each label in the reference (values are not computed for map ends).

cnv_chr_stats.txt – per chromosome coverage statistics

cnv_mask.bed – BED file used to mask CNV calls

The */output/contigs/auto_noise* subdirectory contains error estimation of molecules:

auto_noise1.errbin - error estimation of molecules

autoNoise1_rescaled.bnx – molecule file generated after scaling and applying other filters (e.g. min len, min label density) defined in the configuration file (exp_optArguments.xml).

In */output/contigs/exp_refineFinal1* subdirectory:

EXP_REFINEFINAL1.cmap – the full set of assembled genome maps.

The */output/contigs/exp_refineFinal1/alignref_final* contains alignments between the assembled consensus genome map and the reference:

EXP_REFINEFINAL1.err – human readable error parameter estimates for each EM iteration

EXP_REFINEFINAL1.errbin – binary version of error parameter estimates

In */output/contigs/exp_refineFinal1_sv* subdirectory:

EXP_REFINEFINAL1_full.xmap – contains complete set of alignments used for SV calling, before filtering overlapping alignments on query maps. Used for debugging SV calling

The */output/contigs/exp_refineFinal1_sv/merged_smaps* contains alignments and SV calls between the assembled consensus genome map and the reference:

exp_refineFinal1_merged_filter_inversions.smap – structural variants detected against the reference. This is shown in the SV tab (no SV filtering) in Bionano Access.

exp_refineFinal1_merged_r.cmap – the reference maps used in SV calling.

exp_refineFinal1_merged_q.cmap – the genome maps aligning to the reference.

exp_refineFinal1_merged.xmap – detailed alignment information used for SV calling between the genome maps and the reference maps.

sv_mask.bed – BED file used to mask SV calls

Output files from RVP

Here is an example of the output directory structure after unzipping Rare Variant Pipeline results (*_pipeline_results.zip*):

/output/

/output/data

/output/data/alignmolref/copynumber

/output/data/auto_noise

/output/data/consensus_check/extension/refine1/merge

/output/data/consensus_check/extension/sv/merged_smaps

In the */output* directory:

exp_informaticsReportSimple.txt – summary statistics of the pipeline stages and SV results. This is shown as the Rare Variant Analysis Report in Bionano Access.

exp_informaticsReportSimple.json – contents of exp_informaticsReportSimple.txt in JSON format.

exp_optArguments.xml - parameters used in the analysis run

In the */output/data/alignmolref/copynumber* subdirectory:

cnv_calls_exp.txt – unfiltered copy number calls. This is shown in the Copy Number tab in Bionano Access.

cnv_rcmap_exp.txt – coverage data of each label in the reference (values are not computed for map ends).

cnv_chrAneuploidy_exp.txt – per-chromosome aneuploidy calls. This is shown in the Aneuploidy tab in Bionano Access.

cnv_chr_stats.txt – per chromosome coverage statistics

cnv_mask.bed – BED file used to mask CNV calls

The */output/data/autonoise* subdirectory contains error estimation of molecules:

autoNoise1.err – error estimation of molecules

autoNoise1_rescaled.bnx – molecule file generated after scaling and applying other filters (e.g. min len, min site density) defined in the configuration file (optArguments.xml).

The */output/data/consensus_check/extension/refine1/merge* contains molecules alignments to the assembled genome maps:

exp_refine1_contig*_r.cmap – the assembled genome map (the anchor of the alignment, by genome map IDs).

exp_refine1_contig*_q.cmap – the molecules aligning to the genome maps (the query of the alignment).

exp_refine1_contig*.xmap – detailed alignment information between the molecules and the genome maps.

In */output/data/consensus_check/extension/sv* subdirectory:

SV_LOCI_CLUSTER_full.xmap – contains complete set of alignments used for SV calling, before filtering overlapping alignment on query maps. Used for debugging SV calling.

The */output/data/consensus_check/extension/sv/merged_smaps* contains alignments and SV calls between the assembled consensus genome maps and the reference:

EXP_REFINEFINAL1.smap – structural variants detected against the reference. This is shown in the SV tab (no SV filtering) in Bionano Access.

EXP_REFINEFINAL1.cmap – the full set of maps used in SV calling.

EXP_REFINEFINAL1_r.cmap – the reference used in SV calling.

EXP_REFINEFINAL1_q.cmap – the genome maps aligning to the reference (may be a subset EXP_REFINEFINAL1.cmap).

EXP_REFINEFINAL1.xmap – detailed alignment information used in SV calling between the genome maps and the reference. This is shown in the Match tab in Bionano Access.

sv_mask.bed – BED file used to mask SV calls

Appendix D: In-depth description of the *de novo* assembly pipeline

This is advanced material for users interested in the details of the assembly process. We recommend reading the “*De novo* assembly” in the main text first. The following sections consider features available in the Bionano Solve 3.7 release. Some features mentioned may be missing in earlier releases.

Setting up a *de novo* assembly run

A *de novo* assembly run can be set up with a “molecules” object in Bionano Access or a BNX molecules file on the command line. For more detail, please refer to the Bionano Access User Guide (PN 30142) and the Command Line document (PN 30205) for the two ways of setting up a *de novo* assembly run. A reference *in silico* map (CMAP) can be provided but is optional. If a reference is provided, the molecules and the consensus maps would be aligned to the reference. The alignments could provide helpful information for quality assessment. Also, the reference is necessary for structural variant (SV) calling, and a mask BED file can be selected for annotating the resulting SV calls.

Several default assembly parameter sets are provided for different use cases. The parameter sets are specified in XML files, whose names suggest the target usage. Which parameter set to use depends on several factors. In the Access interface, the dropdown selection menus reflect the specific considerations that impact the selection of parameter sets. For example, certain parameter sets have been optimized for larger genomes (typically > 5 Gbp) that require more computation and memory.

For genome assembly or finishing projects, if a high-quality reference is not available, we recommend using the pre-assembly option. When this option is enabled, a preliminary assembly is constructed first, and this assembly is used as reference for estimating noise parameters in the molecules in the subsequent assembly steps.

For SV analysis projects for typical diploid samples (such as human), we recommend using the haplotype-aware options and cutting complex multiple-path regions (CMPRs) in the map. When haplotype-aware options are enabled, we try to separate possible alleles. Because of the presence of large segmental duplications in the genome, which can confuse the assembler, we recommend cutting maps that cover these CMPR regions to avoid assembly errors.

Key stages of the *de novo* assembly pipeline

Here are brief descriptions of the key steps in the basic *de novo* assembly pipeline (Figure 30). By default, the output data from the intermediate pipeline steps are not imported or visualized in Bionano Access. For Bionano Access runs, the key files in the final steps are compressed before being imported into Access. The intermediate files are only kept temporarily on the compute server to reduce disk space use. However, when needed and available, certain intermediate files may be useful for troubleshooting.

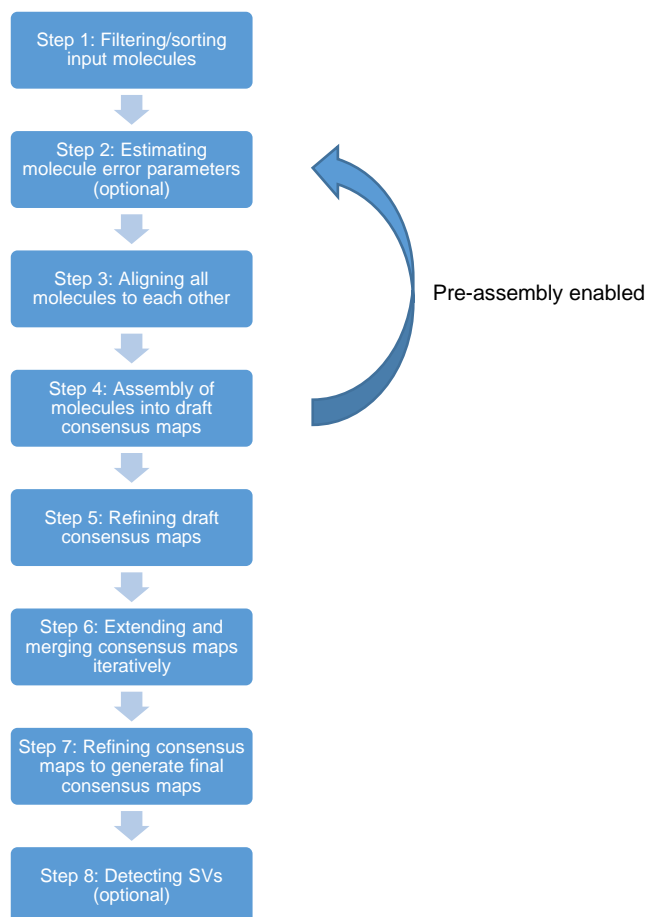


Figure 30. *De novo* assembly workflow.

The pipeline uses RefAligner and Assembler for alignment and assembly operations. It serves to chain RefAligner and Assembler operations together, handle job submissions, and generate intermediate data summaries. In addition to the pipeline report (typically named exp_pipelineReport.txt), RefAligner- and Assembler-generated output includes information in the headers about the command to generate the output and other useful information (such as the version of the RefAligner/Assembler). There are .stdout files that log information during RefAligner/Assembler operations. Although expected to be rare, assertion errors can occur, and error message can be found in the .stdout files. Please contact Bionano Support if assertion errors are observed.

If the pre-assembly option is enabled, steps 1, 3, 4 from below are performed first. Then, the consensus maps would be used as the input reference, and the assembly would then restart from step 2.

1. Filtering and sorting input molecules

The raw input to the assembly pipeline would be an all.bnx file typically. Depending on what filters have been applied previously, it might contain short molecules or molecules that do not have sufficient labels for alignment and assembly. Therefore, the pipeline filters out some of the molecules using a length (typically with -minlen 120)

and a site (typically with -minsites 9) filter. This filtering step helps filter out molecules that are likely not useful and limit the file size. The pipeline would generate an `all_sorted.bnx` file in the main output directory after the filters are applied. As the `.bnx` name suggests, the molecules are sorted. The sorting is done based on the numeric molecule IDs.

2. Estimating molecule error parameters (optional)

In the auto-noise stage, the pipeline attempts to align all molecules in the `all_sorted.bnx` file to the reference *in silico* map (`.cmap`) to estimate molecule noise parameters. The output is stored in the “auto_noise” directory under the main output directory. The parameters include for example, the amount of missing or extra labels compared to the reference, and the interval sizing differences. They are represented in the `.err` files. These parameters are helpful for scoring alignments in subsequent steps. Conceptually, they form an expectation of the errors and help quantify whether an alignment is good or bad. Using a good-quality reference allows the pipeline to set the expectation correctly. If the reference quality is questionable, we would expect the molecules to be different from the reference (even if the molecules were perfect), thus artificially inflating the error estimates.

The alignment is done in several iterations, and the noise estimates are updated in each round. After the iterations, the noise estimates should converge. The very first iteration starts with default noise parameters. In rare cases, if the default parameters are too different from the actual parameters, it's possible that very few alignments are found, making RefAligner unable to estimate the parameters. In such cases, one may adjust the default noise parameters in the `parameters.xml`. If the error rates are too high (because either the molecules or the reference is of poor quality), the auto-noise stage can take a long time, and the noise estimates may not converge. This might impact subsequent assembly stages.

For computational efficiency, the auto-noise stage is split up into two sub-stages: `auto-noise0` and `auto-noise1`. The output from `auto-noise1` is used for subsequent steps.

Another calculation done during the autonoise stage is the determination and correction of the per-cohort stretch and the per-scan stretch factors for the molecules, measured in basepairs per pixel (bpp). During image acquisition, each scan is divided into a number of subgroups (“cohorts”) for real-time analysis. All the molecules in the same cohort or scan (one scan contains multiple cohorts) may be stretched differently on average compared to molecules in a different cohort or scan due to a variety of reasons. Initially a BNX file is generated per cohort, but all BNX files are merged before assembly, so the calculation of the stretch factors aims to correct for these differences. Note that the correction refers to the overall or average stretch of molecules in a cohort or scan, which is different from local or per-interval stretch variation in the individual molecules. To convert the pixel-based data from the molecule images into basepairs, the pipeline uses information from the molecule-to-reference alignment. The bpp for each cohort and scan is normalized to account for the average cohort or scan stretch differences. After normalization, the molecules are expected to be more uniform.

The main output of the auto-noise stage are the error parameters (.err) and a rescaled BNX file with stretch-corrected molecules. These are used for subsequent steps. The BNX file is sometimes split into pieces (named with pattern all_X_of_N.bnx) for more efficient computation. If a reference is not provided the auto-noise stage is skipped and default error parameters are used instead. This can be sub-optimal if the default error parameters are not appropriate for the dataset of interest. We recommend enabling the pre-assembly option such that the preliminary assembly could be used for noise estimation.

3. Aligning all molecules to each other (pairwise alignment)

The input is the rescaled molecules, and the output is stored in the "align" directory. The pipeline pairwise-aligns the molecules to other molecules to form the basis of the overlap graph needed for the overlap-layout-consensus, OLC, assembly. The output files are binary (.align); they contain information about which molecules are aligned to which molecules, the alignment confidence, and the offsets. The fraction of molecules that have pairwise alignments should be comparable to the map rate. The statistics from the pairwise alignment step are available in the .stdout output.

4. Assembling molecules into draft consensus maps

An initial assembly is created from the pairwise alignments of the previous step. The output is a set of draft consensus maps (.cmap files) in the "exp_unrefined" directory. The assembler uses the pairwise alignments to form an overlap graph and analyzes it for potential paths. The graph can be quite complex due to both genuine genomic features (such as repeats) and false positive pairwise alignments, so we apply several clean up steps to output the longest paths (which consensus maps represent). The assembly size should be close to the genome size, even though maps continue to be extended and refined in subsequent stages. The N50 (which is an indication of consensus map lengths) does tend to increase in the later stages. If a reference is provided, the consensus maps are aligned to it for quality control purposes.

5. Refining draft consensus maps

The initial refinement of the draft consensus maps is done by analyzing molecule-to-map alignments during the stages refineA and refineB. The input is the draft consensus maps and the output is refined consensus maps. RefAligner tries to determine whether the labels on the draft consensus maps need to be adjusted (added, deleted, or moved) and whether the maps need to be split due to lack of molecule support in certain regions of the map. The difference between the stages refineA and refineB is that the former uses only molecules that were in the initial assembly. The latter tries to recruit additional molecules by aligning all molecules to the refineA consensus maps (assumed to be of higher quality) and refining the maps based on the larger set of molecules.

6. Extending and merging consensus maps iteratively

Extension and pairmerge are paired stages that are run iteratively in an effort to make the maps more contiguous. Input molecules are aligned to the ends of the consensus maps to "extend" them. After maps are extended, some of them are expected to overlap each other. The pipeline aligns the extended maps to each other and merges overlapping maps. The map N50 is expected to increase because of this process. This is similar in concept to the gap closing step that some sequence assembly tools perform. By default, the extension and merge process is repeated five times, as the consensus map N50 tends to plateau after 5 rounds.

During the extension step, if the haplotype-aware options are turned on, the pipeline also looks for alternative alleles. If molecules align but diverge from the maps, the pipeline forms new maps that represent alternative alleles. This tends to split off large SVs, allowing them to be detected. Rare alleles such as those found in cancer clones may be found, but the pipeline requires at least a certain number of molecules to support an allele and makes certain assumptions about the minimum fraction of molecules that support an allele. The *de novo* assembly is not optimized for detection of rare allele. Rare Variant Pipeline would be more suited for that purpose.

In the final merging stage, the CMPR-related option would become active if enabled. RefAligner would analyze the pairwise alignments of the consensus maps. If pairs of maps align to each other but do diverge, they likely contain CMPR sequences, and they may be cut.

7. Refining consensus maps to generate the final set of maps

This stage of final refinement of the consensus maps generates the final final product of the assembly pipeline, a set of refined maps. It tends to be the most time-consuming stage. The pipeline aligns again the molecules to the maps and analyzes the alignment to see if any maps need to be updated. If the haplotype-aware options are turned on, the pipeline also looks for alternative alleles. This tends to split off smaller SVs.

The pipeline makes the default assumption that there can be two alleles, so one often sees a pair of maps per genomic region. For genome assembly or finishing projects, the presence of these seemingly redundant maps may create issues during scaffolding. We recommend not using the haplotype-aware assembly options for those projects.

8. Detecting SVs using the final consensus maps (optional)

The final consensus maps are used for SV calling, if a reference is provided. The pipeline aligns the maps to the reference using a multiple alignment algorithm (where all possible alignments are found) and analyzes the differences. The SV calls are recorded in the .smap files.

Runtime performance of the *de novo* assembly pipeline

The runtime performance data are included in earlier sections. Generally, runtime depends on several factors. The amount of input data, the size of the genome, the data quality, and the load of the compute server can contribute to runtime differences. The *de novo* assembly automatically down-samples the dataset if the amount of data is greater than 220X. Please refer to Data Collection Guidelines (PN 30173) for more detail. The downsampling step helps limit the runtime for high coverage datasets, and it should not impact the quality of the assembly. The pre-assembly option requires additional steps in the assembly process and might take longer. The haplotype-aware options also lengthen the assembly time because of the extra analyses needed to detect alternative alleles.

Evaluation of *de novo* assembly results

Please refer to the Bionano *De Novo* Assembly Report (PN 30255) for more detail.

Appendix E: Interpretation of inversion breakpoint calls

Introduction

All the structural variant (SV) calls are obtained by aligning consensus genome maps (query maps) to a reference map using a Multiple Local Alignment algorithm and analyzing the alignments for SV signatures. Pairs of alignments within a map are analyzed and inconsistencies representing possible SV events between the genome maps and the reference are identified. The signature for inversion breakpoints is a pair of neighboring alignments with opposite orientations (each alignment has an orientation relative to the reference), indicating that two locations (breakpoints) on the reference map are fused in the query map (Figure 31). The orientation information is in the XMAP file; it can either be “+” or “-”, relative to the reference.

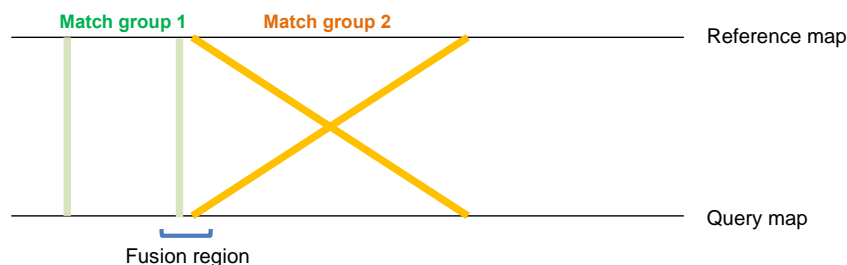


Figure 31. Signature of an inversion event. Two alignments from the XMAP file (yellow and green solid lines) are close to each other with opposite orientations.

There are two modules for detecting inversion breakpoints: one for detecting breakpoints of relatively “large” inversions, and one for detecting breakpoints of relatively “small” inversions (typically involving fewer than five labels in the inverted region). Both modules are automatically run during SV detection; there is not an explicit switch or parameter to turn them on. The definitions of large and small inversions are informal but are relevant in the context of inversion breakpoint detection using Bionano data. Whether an inversion is considered small and large depends on whether the map spans both breakpoints of the inversion. A map could span both breakpoints of a small inversion, and both breakpoints are called together. If a map does not span both breakpoints, but there is evidence of an inversion, a single breakpoint is output. The other breakpoint of the inversion may be detected in another map (which likely also does not span both breakpoints of the inversion).

We generally use the term “breakpoints” to describe the two sides of an inversion. Conceptually, an inversion involves two breaks in the sequence (relative to the reference) and a subsequent inversion of that sequence. This could happen in different ways biologically. Alternatively, the concept of “fusions” may be helpful. There would be two fusion points where the inverted sequence is joined with the surrounding non-inverted sequence. By analyzing the alignment data, one could infer breakpoints and which sequences are fused together.

SMAP output for inversion breakpoints

The SMAP file summarizes information about the labels involved in an inversion, both from the reference and the query maps. To fully interpret inversion events, we need to also consider the corresponding XMAP entries of the alignments used to make the call. Each line in the SMAP references four coordinates (two on the query map and two on the reference map). The coordinates of the labels in the query and reference maps are given in the QryStartPos, QryEndPos, RefStartPos, and RefEndPos columns, respectively. Please refer to SMAP File Format Specification Sheet (PN 30041) for the specific meaning of those columns.

SV types other than inversion breakpoints are represented by single SMAP lines. For inversion breakpoints, additional relevant coordinates are reported. To do so, all inversion breakpoints are represented by pairs of two lines in the SMAP file. The pairing can be identified using the link ID column. The two lines have link IDs that point to each other.

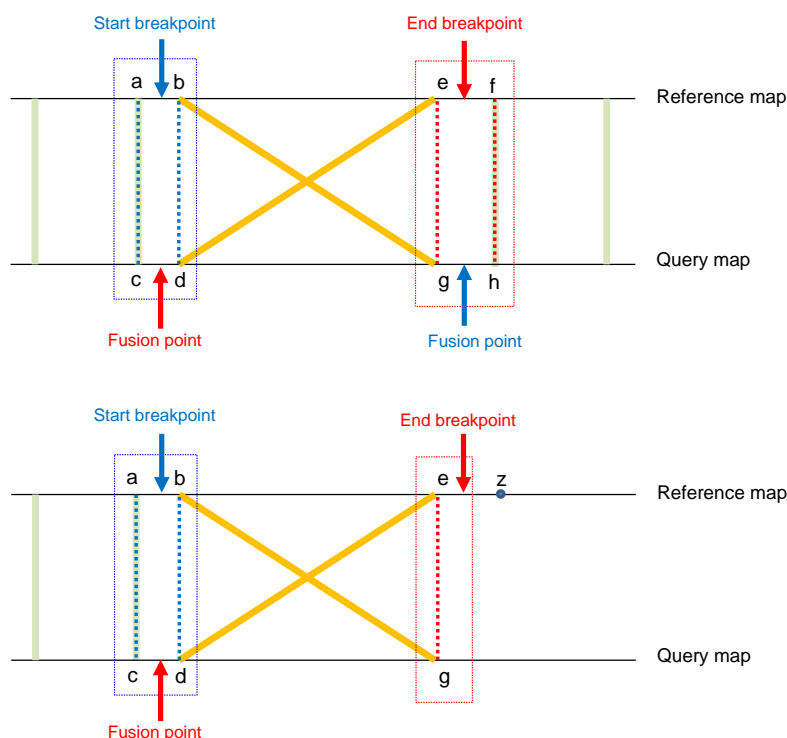


Figure 32. Signatures of inversions.

A) Signature of a fully described inversion: The blue box contains the four labels described in the first *inversion_paired* line of the SMAP, and the red box contains the four labels described in the second *inversion_paired* line. The dotted lines represent the pairing of coordinates in the columns of the Smap text line. E.g., in this diagram a=RefStartPos, b=RefEndPos, c=QryStartPos, and d=QryEndPos, and similarly, e=RefStartPos, f=RefEndPos, g=QryStartPos, and h=QryEndPos. The coordinates at basepair resolution of the start (start breakpoint) and end (end breakpoint) of the inverted material are not observed, but they are bounded by the intervals (a,b) and (e,f). In the query map, the start of the inversion is fused to the downstream material (right of the blue fusion point), and the end breakpoint is fused to the upstream material (left of the red fusion point). The coordinates of the fusion points are

bounded by (c,d) and (g,h), respectively.

B) Signature of a partially described inversion, where there is evidence of a single fusion event. The blue box contains the four labels appearing in the *inversion* line of the SMAP file, and the red box contains the two coordinates of the labels in the *inversion_partial* line. As described previously, the dotted lines represent the pairing of reference and query map coordinates. An end breakpoint in the interval (e,z) is fused with the material upstream of the red fusion point, bounded by (c,d). The red fusion point would align to a start breakpoint bounded by (a,b). The label z is the label of the reference map immediately after e. Note: The coordinate of label z is not provided in the SMAP file. Note: The situation where the coordinates of the *inversion* line are greater than those in the *inversion_partial* line is also possible, and the interpretation is the same as in 2B.

Small inversions (whose inverted regions contained fewer than five labels; typically less than 50 kbp) can be identified by searching in a limited space for two inversion signatures involving the same inverted alignment. They may be spanned by single genome maps, and there are eight coordinates (four on the query map and four on the reference map) involved in their descriptions (Figure 32). In the SMAP file, they are reported as two linked lines where the type is *inversion_paired*. They represent two inversion breakpoints that form the full inversion event. Paired inversions fully describe the material inverted and its bounds.

Neither the start nor the end breakpoints are observed at basepair resolution. However, the uncertainty in the breakpoint coordinates is bounded by the distances between the boundary labels. In Figure 32A, there is an inversion event where the genomic coordinate of the start is bounded by labels **a** and **b**, and the genomic coordinate of the end is bounded labels **e** and **f**. In the query map, the start of the inversion is fused to the material downstream the inversion, and the end of the inversion is fused to the upstream material.

In some cases (often for larger inversions), the signature for only one inversion breakpoint is found. These single inversions breakpoints are described by six label coordinates (Figure 32B). In the SMAP file, they are reported as two linked lines where the first one is of type *inversion* and the second one is of type *inversion_partial*. In the *inversion_partial* lines, some coordinate fields are not used (output as -1), hence the name “_partial”.

In the scenario of Figure 32B, there is a fusion in an inverted orientation between a point bounded by labels **a** and **b**, and a point bounded by labels **e** and **z**, but there is no evidence of a reciprocal fusion or some other type of fusion.

Currently, inversions larger than 5 Mbp are called as intra-chromosomal translocation breakpoints. The 5-Mbp cutoff can be considered arbitrary (and may be adjusted on the command line). These large inversions are typically only partially spanned, so there is more uncertainty for the SV type assignment.

Inversion breakpoints report file

Apart from the SMAP file, we also provide an intermediate file with a different format specifying just the bounds of the inversion breakpoints in the reference map and the key fusion points (InversionBreakpoint1 and InversionBreakpoint2). The file is generated during confidence score calculations (available in Solve 3.6). Each

inversion call is described with a single line and contains the following fields:

- RefcontigID1: The reference map (or chromosome) containing the inversion. It has the same value as the SMAP file.
- QryContigID: The query map used to call the inversion. It has the same value as the SMAP file.
- InversionBreakpoint1: The coordinate of the label in the reference map closest to the reference start breakpoint based on the alignment.
- Bp1LowerBound: Lower bound for InversionBreakpoint1
- Bp1UpperBound: Upper bound for InversionBreakpoint1
- InversionBreakpoint2: The coordinate of the label in the reference map closest to the reference end breakpoint based on the alignment.
- Bp2LowerBound: Lower bound for InversionBreakpoint2
- Bp2UpperBound: Upper bound for InversionBreakpoint2
- Bp1SmapId: The SmapEntryID of the SMAP line used to infer InversionBreakpoint1.
- Bp2SmapId: The SmapEntryID of the SMAP line used to infer InversionBreakpoint2.

The uncertainty for the bounds is based on a "closest label" criterion. The location of the true fusion point is within the label used to detect the inversion event and the next/previous label. Refer to Figure 27 for a graphical interpretation. For all entries in the file, it is guaranteed that:

- InversionBreakpoint1 < InversionBreakpoint2
- Bp1LowerBound <= InversionBreakpoint1 <= Bp1UpperBound
- Bp2LowerBound <= InversionBreakpoint2 <= Bp2UpperBound

There is no guarantee that Bp1SmapId < Bp2SmapId.

Inversion confidence score

The confidence score aims to give the user a measure of the correctness of an inversion call, assigning a high value to the score for calls that are true positives, i.e., events that are present in the sample, and low value of the score for false positives, i.e., calls that are made by a pipeline that are not present in the sample.

FAQs

1. Is there a way to get information about the orientation of the inversion breakpoint call? Could I tell whether I am looking at the left or right breakpoint?

The inversion_partial coordinates in the Smap refer to the inverted matchgroup. This determination is made by the rule that the matchgroup lines should not cross. See Figure 27B.

2. Why is the output in the SMAP file different from what I see in the SV table in Access?

Within Bionano Access, there is a script to "flatten" the two-line inversion entries into single-line entries for Access visualization. The entries for the other SV types are largely not modified. When an SMAP file is imported into Bionano Access, this script is run, and the resulting output is used for visualization. Any filtering applied in Access is then based on this flattened output instead of on the SMAP explicitly.

The "SuperType" and "Mask" columns are added to handle SV calls with additional annotations and for downstream filtering. Insertion and deletion calls may be output as "_nbase", if they overlap with N-base gaps; their SuperType would be "insertion" and "deletion", respectively. The SuperType is coded in numeric values:

SV type	SuperType value
Insertion	1
Deletion	2
Inversion breakpoint	3
Translocation breakpoint	4
Duplication	5
Catch all	-1

Depending on their types, the SV calls are classified as masked or not masked. For example, translocation breakpoint calls that overlap with annotated segmental duplication regions are output as "_segdupe" calls. These are considered masked calls. This allows Access to interactively filter out masked calls in the interface.

The script separates the SMAP entries based on the SV types. For all "inversion" types, it saves them in a hash table (base data structure in Perl) because they need to be flattened separately. After building the hash table, the script checks that the number of entries in the table is an even number (assuming that all inversion breakpoints are represented in paired entries). The flattened output has 6 columns for inversion visualization specifically: "BrkptEdges_Start_Ref_1", "BrkptEdges_Start_Ref_1_Pos", "BrkptEdges_Start_Ref_2", "BrkptEdges_Start_Ref_2_Pos", "BrkptEdges_End_Qry_1", and "BrkptEdges_End_Qry_2". The additional columns to be added contain information about where to draw the "hourglass", or the inverter alignment that is characteristic of an inversion. For example, in Figure 27B, **a** and **e** are reported to be the inversion boundaries.

The additional columns from the Variant Annotation Pipeline are carried over. However, for inversion entries, there is some loss in information during the flattening process. The flattening script only retains annotations from the first of the paired inversion entries.

3. How is inversion size estimated?

For inversion_paired calls:

In this case, both breakpoints are bounded, and the inversion size is estimated as the reference inverted alignment size plus 50% of the gaps on both sizes (Figure 33). For small inversions the gaps are often a significant fraction of the size. Taking into the gaps improves the accuracy of the size estimate.

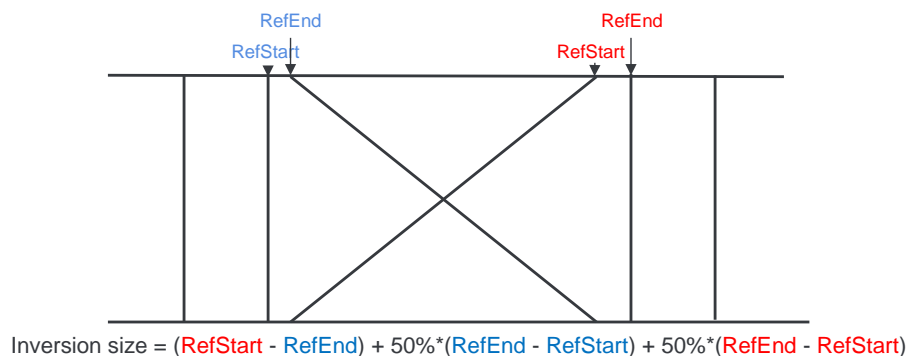


Figure 33. Estimation of inversion size on paired inversion calls.

For inversion/inversion_partial calls without overlap of the alignments on the reference:

The inversion size estimate is simply the size of the inverted alignment (Figure 34). This represents the lower bound estimate of the inversion size. Because the map does not fully span the inversion, we do not have full information about the entire inversion events. The inverted sequence may in fact be larger than the inverted alignment.

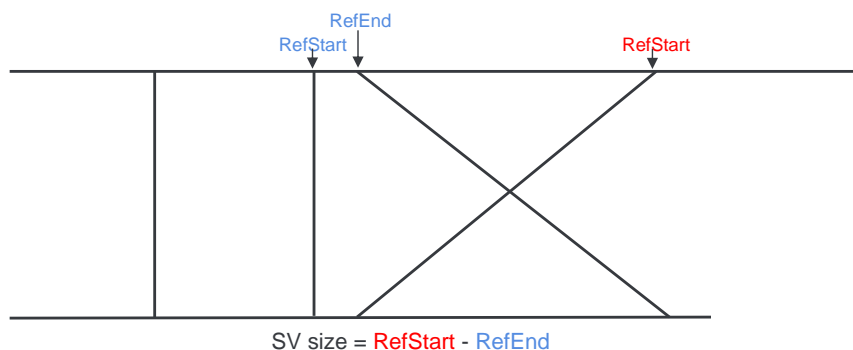


Figure 34. Estimation of inversion size on inversion/inversion_partial calls.

For inversion/inversion_partial calls with significant overlap of the alignments on the reference:

The overlap is taken into account, and the inversion size is estimated differently than describe in the previous section. Instead of taking the size of the full inverted alignment, it is “trimmed” to where the straight, non-hourglass alignment starts (Figure 35).

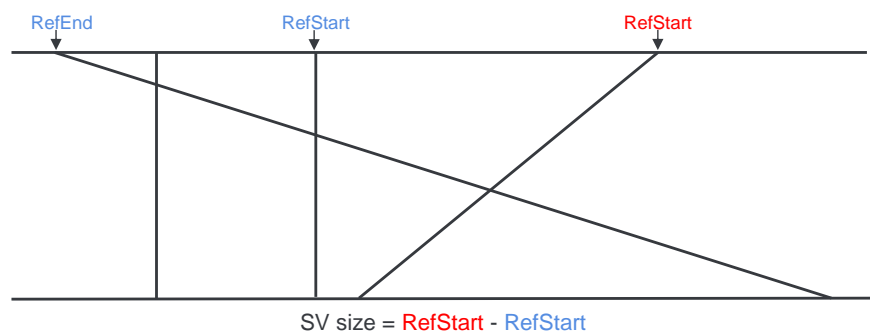


Figure 35. Estimation of inversion size on inversion/inversion_partial calls with alignment overlap.

Appendix F: Confidence modeling for inversion and translocation breakpoints

Introduction

The confidence scores aim to give the user a measure of how likely a structural variant (SV) call is correct, assigning high scores for calls that are likely true positives (TP), i.e. events that are present in the sample, and low scores for likely false positives (FP), i.e. calls that are made by a pipeline that are not present in the sample. Here we describe the algorithms used to produce the confidence scores for translocation and inversion breakpoints, together with their performance both in simulated and real datasets. The new confidence models were introduced in the Solve 3.6 release and updated for the Solve 3.7 release.

Datasets

For the purpose of developing the models, we considered as source of truth real samples processed by the Bionano pipelines, either the *de novo* assembly pipeline or the Rare Variant Pipeline (RVP), having orthogonal datasets of SV calls made with other technologies, as well as 10 simulated genomes containing SVs. For the real samples, we used the proband data from the Ashkenazi trio (NA24385) from the integrative benchmark from Genome in a Bottle (GIAB; Zook et al., 2019), five cancers samples (Dixon et al., 2018), seven samples sequenced with Pacific Biosciences (Audano et al., 2019), nine internal samples with manually curated translocations, three samples from the integrative analysis from 1000 Genomes (Chaisson et al., 2019), SK-BR-3 (a breast cancer cell line) with orthogonal data from Pacific Biosciences, and 2 colorectal cancer samples with orthogonal data from Oxford Nanopore.

The datasets used to validate the models consisted of 17 runs of both the assembly and RVP pipelines containing mixtures of simulated SVs together with the hg19 human reference genome and 21 samples with confirmed SV events detected by cytogenetic analysis.

Models for the scores

Model for the translocation breakpoint scores

The translocation scoring model has been retrained for Bionano Solve 3.7. The input dataset contained 7,914 translocation calls, of which 3,004 were confirmed as TP by the orthogonal datasets or manual curation, and 4,910 were FP calls. The features were derived from maps alignments, reference coverage, and genomic location. In particular:

- The SV type: inter- or intrachromosomal translocation breakpoints

- The confidence score assigned by RefAligner (Shelton et al., 2015) to each of the 2 alignments used for calling the translocation breakpoints.
- The values of label coverage and occurrence from the query map in the region used to call the translocation breakpoints (± 30 kbp buffer). These are proportional to the number molecules aligning to a query map. These values are divided by the global effective sequence coverage to make them relative and thus independent of the input coverage depth that may be different for different datasets.
- The average values of ChimQuality, SegDupL, SegDupR, FragileL, and FragileR in the consensus map region (± 30 kbp), as described in the Bionano CMAP File Format Specification Sheet (PN 30039).
- The genomic position of the translocation breakpoints bucketed into regions of 1 Mbp.
- Boolean variables indicating if the translocation breakpoints are in the centromeric or segmental duplication regions. The segmental duplication regions used during training have been updated.

The full dataset was randomly split into training (60%), validation (12%), and test (28%) datasets. The algorithm employed was Gradient Boosted Trees as implemented in LightGBM (Ke et al., 2017). After assessing model convergence and ability to improve its performance (Figure 36), it was optimized for the depth of the tree representation, the learning rate, and the number of node leaves using the package Optuna (Akiba et al., 2019) with the early-stopping convergence rule applied to the validation dataset. The area under the curve (AUC) measure on the receiver operator curve (ROC) was 0.991 (Figure 37a). Using the TPR/FPR plots (Figure 37b) we selected individual thresholds of 0.05 for interchromosomal translocations and 0.05 for intrachromosomal translocations (Note: These new thresholds for Bionano Solve 3.7 are different from those used in Bionano Solve 3.6). Users may select a different set of thresholds depending on the project goals.

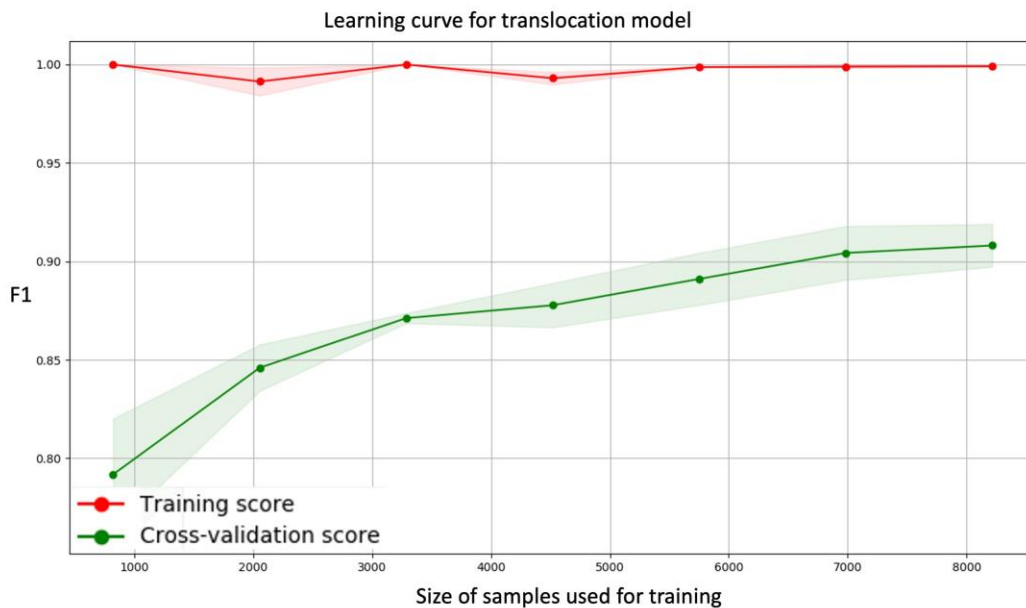


Figure 36. Learning curve of the model for the translocation score, built with the leave-one-out cross-validation schema (5 partitions). The x-axis shows the number of examples used for each of the models produced, and the y-axis is the f1-score (harmonic mean between sensitivity and PPV). The width of the shaded region shows the standard deviation of the f1-score for each of the cross-validation partitions.

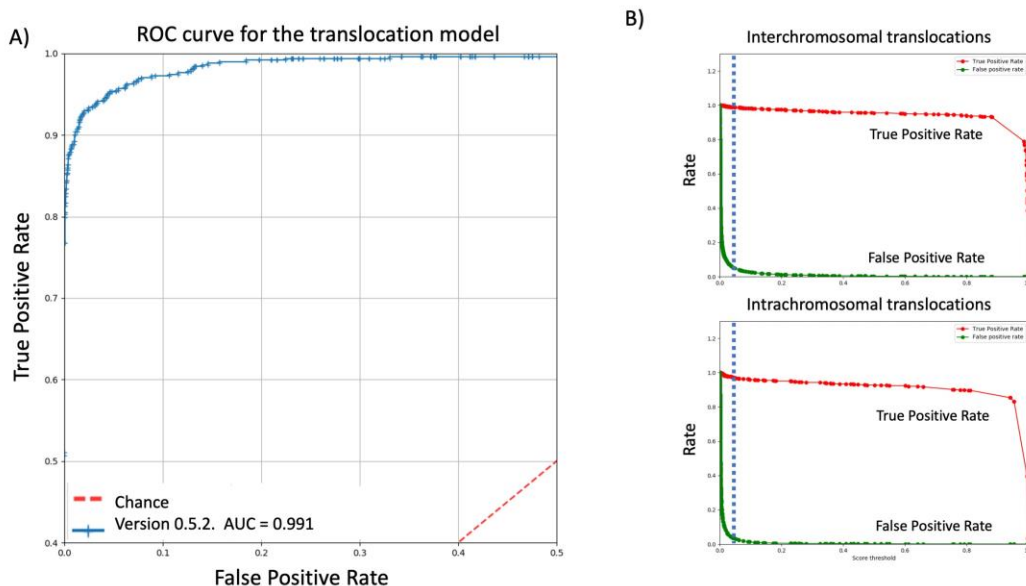


Figure 37 A) Receiver operator curve (ROC) for the translocation score model for its human version. The x-axis is the False Positive Rate (FPR) and the y-axis is the True Positive Rate (TPR). B) Plots of the TPR and FPR for the human model as functions of the

threshold used for classification. The dashed line indicates the selected threshold.

We also built a non-human version of the confidence score model that does not consider the human-specific features. Namely, the genomic position and the location in the centromeric/segmental duplication regions are not considered. This model remains the same as in Bionano Solve 3.6, with AUC=0.939.

Model for inversion breakpoint scores

To build the confidence score model for inversion breakpoints, each of the breakpoints of an inversion is considered separately, and thus each inversion call contributes with two data points to the model. We used 27,699 inversion breakpoints, 21,194 TP and 6,505 FP. The features considered for each breakpoint were:

- Minimum of the RefAligner confidences for the two alignments used to call the inversion.
- Length of the shortest alignment.
- Number of labels aligned in the inverted alignment.
- Distance between the lower and upper bounds for the position of the breakpoint.
- Number of labels in the query map of the inverted alignment that are in complex multi-path regions (CMPRs). RVP maps are not expected to have CMPR annotations.
- Number of upstream/downstream labels in the reference that are unaligned. If the breakpoint is the start of the inversion, this number is calculated counting unaligned labels upstream of the position. For the end breakpoint, the number of unaligned labels is counted downstream of the position.
- Mean of the outlier fraction (OutlierFrac column of the CMAP file) for the region of the query map contained in the inverted alignment (± 30 kbp buffer).
- Boolean variables indicating if the two alignments overlap in the reference map or overlap in the query map. An overlap in the alignments for the query map is a strong indication of a problematic call, as it implies that the query map aligns to the reference in both directions or that there is ambiguity in the alignment.
- The values of label coverage and occurrence from the query map in the inverted region, normalized by the effective sequence coverage. This feature is calculated the same way it is calculated for translocations.
- Boolean variables indicating if the breakpoint is located in the centromeric or segmental duplication regions (the list of annotated regions was downloaded from the UCSC Genome Browser).

We again split the full dataset into training, validation, and test datasets, with the same percentages, and used LightGBM and Optuna. Due to the imbalanced nature of the available dataset (77% TP, 23% FP), we added a pre-processing step of random oversampling with replacement for the minority class (FP) during each model training. Each time a model was optimized in search of the best hyperparameters, the training dataset contained the same data points for the TP calls, but different datasets of FP calls. The learning curve shows that the model cannot learn more with the data at hand (Figure 38). Compared to the model for translocations, the model for inversions obtains better performance metrics (Table 9) and increased area under the ROC curve (Figure 39a), which allows to set a classification threshold with very high TPR and low FPR. We selected as such threshold 0.7

(Figure 39b).

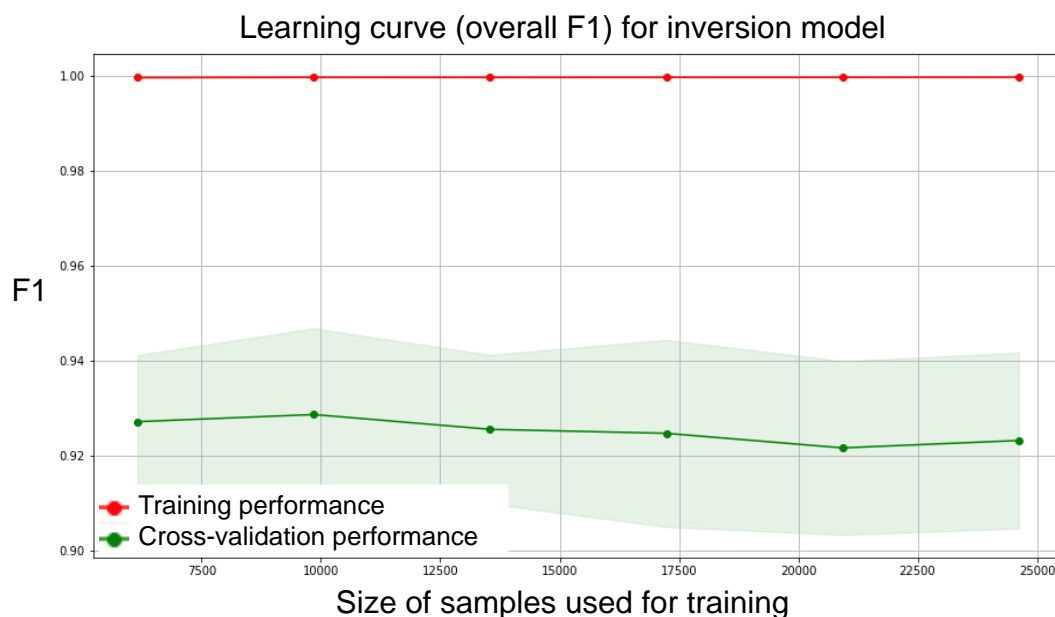


Figure 38. Learning curve of the model for the inversion score, built with the leave-one-out cross-validation schema (5 partitions). The x-axis shows the number of examples used for each of the models produced, and the y-axis is the f1-score. The width of the shaded region shows the standard deviation of the f1-score for each of the cross-validation partitions.

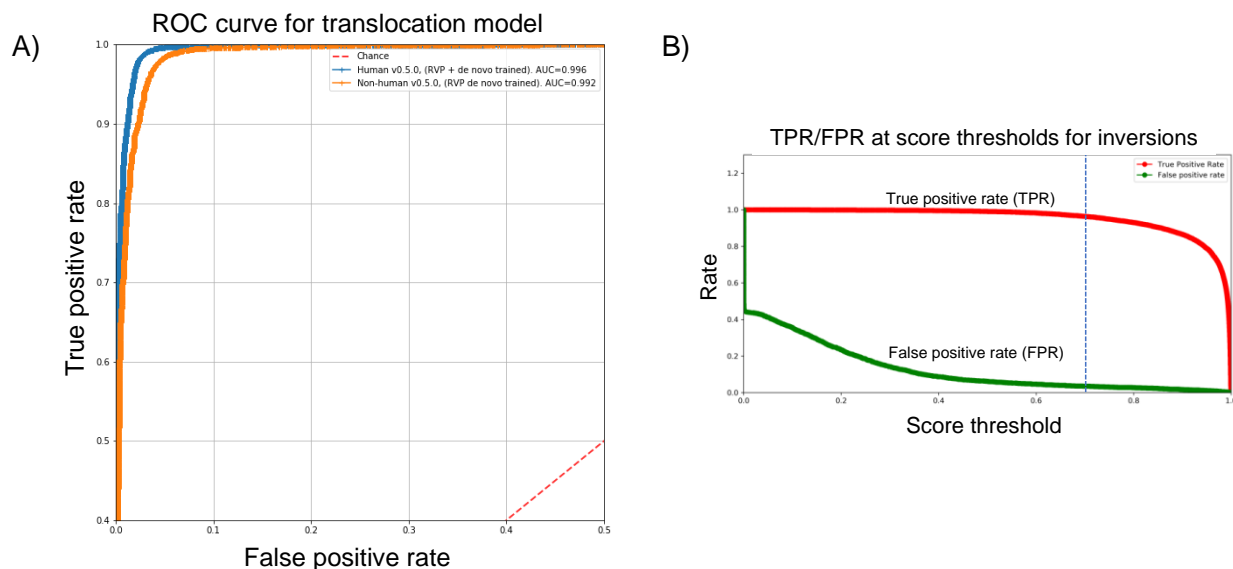


Figure 39. A) Receiver operator curves (ROC) for the inversion score model, both for its human and non-human versions. The x-axis shows the FPR and the y-axis shows TPR. B) Plot of the TPR and FPR for the human model as functions of the threshold used for classification. The dashed line indicates the selected threshold.

Dataset	Data points	Threshold	Percentage	Sensitivity	PPV	Accuracy
Training	16619	0.70	60%	100%	100%	100%
Test	7756	0.70	28%	95%	97%	94%

Table 15. Performance of the human version of the confidence score model for inversions.

A full inversion call is scored as the maximum of the individual scores of its two breakpoints.

Note: The model for inversions has not been updated for Bionano Solve 3.7.

Validation

We validated our models with two independent datasets not used during model building: A) Mixtures of simulated genomes based on introducing SVs into the hg19 human reference, and B) Real samples with translocations and inversions confirmed by cytogenetic analyses.

Simulated datasets

This dataset consisted of two simulated genomes, one containing only translocations and the other containing only inversions. We first ran each genome at 100% allele fraction through both the *de novo* and RVP pipelines, and subsequently mixed with the human reference hg19 at 5%, 10%, 20%, and 30% allele fraction. In total, we used the data from 17 pipeline runs. The overall performance for each of the SV types (Table 16) is similar to those obtained for the test split during model generation, with an improvement in the PPV of intrachromosomal fusions.

Dataset	Data points	Threshold	Sensitivity	PPV	Accuracy
Interchromosomal translocations	8,997	0.65	96%	98%	95%
Intrachromosomal fusions	1,128	0.70	99%	92%	91%
Inversions	22,062	0.70	94%	92%	91%

Table 16. Performance of the scoring models for 17 runs of the pipeline on simulated datasets.

Cytogenetics datasets

To validate the concordance between calls with high scores and those confirmed by cytogenetic analysis, we ran the scoring algorithms on a dataset composed of 21 samples containing 25 curated translocation (19) and inversion (6) variants and searched for Bionano calls consistent with those events. A Bionano call describing a certain event is considered consistent with the expected cytogenetic variant if it is of the same variant type, affects the same chromosomes, and it is located in the same cytoband described in the cytogenetics report (± 10 Mbp buffer). Multiple Bionano calls may be consistent with a given cytogenetic variant. When plotting the maximum confidence score of the calls consistent for each event, there is always at least one call with a score above the recommended thresholds (Figure 40). All the cytogenetic events would have been found by looking only at the Bionano calls of high confidence.

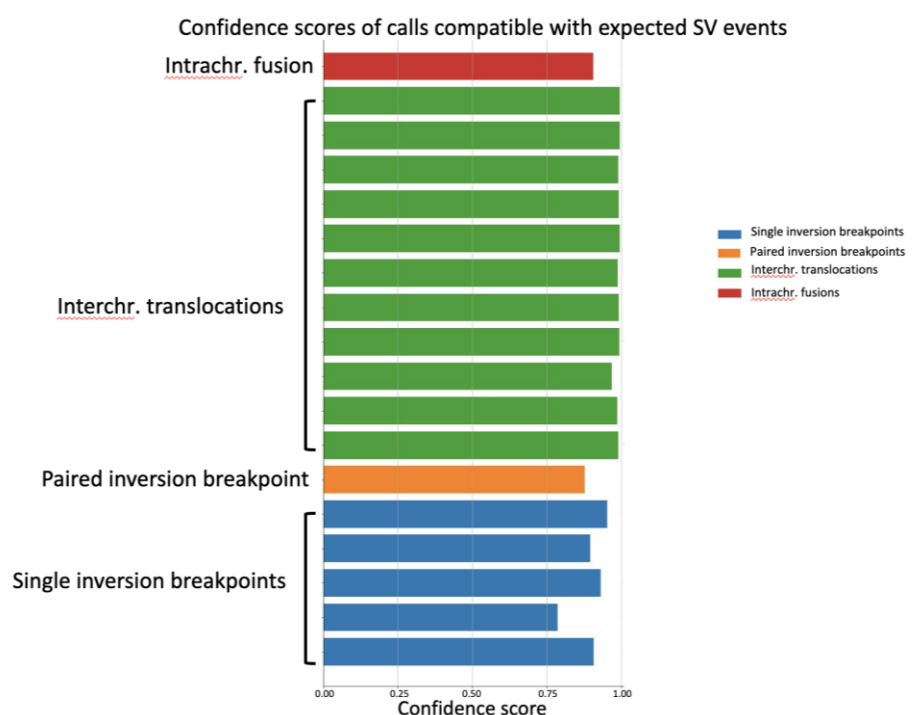


Figure 40. Translocation and inversion confidence scores for the Bionano calls consistent with curated cytogenetic variants. Left y-axis: Cytogenetic variant, one per row. x-axis: Maximum value of the confidence score for all consistent Bionano calls.

References

- Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M. (2019). Optuna: A Next-generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, (Anchorage AK USA: ACM), pp. 2623–2631.
- Audano, P.A., Sulovari, A., Graves-Lindsay, T.A., Cantsilieris, S., Sorensen, M., Welch, A.E., Dougherty, M.L., Nelson, B.J., Shah, A., Dutcher, S.K., et al. (2019). Characterizing the Major Structural Variant Alleles of the Human Genome. *Cell* 176, 663-675.e19.
- Bionano Genomics (2020). CMAP File Format Specification Sheet.
- Chaisson, M.J.P., Sanders, A.D., Zhao, X., Malhotra, A., Porubsky, D., Rausch, T., Gardner, E.J., Rodriguez, O.L., Guo, L., Collins, R.L., et al. (2019). Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat. Commun.* 10, 1784.
- Chen, T., and Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16*, (San Francisco, California, USA: ACM Press), pp. 785–794.
- Dixon, J.R., Xu, J., Dileep, V., Zhan, Y., Song, F., Le, V.T., Yardimci, G.G., Chakraborty, A., Bann, D.V., Wang, Y., et al. (2018). Integrative detection and analysis of structural variation in cancer genomes. *Nat. Genet.* 50, 1388–1398.
- Ke, G., Meng, Q., Finely, T., and Wang, T. (2017). LightGBM: A Highly Efficient Gradient Boosting Decision Tree. *Adv. Neural Inf. Process. Syst.*
- Shelton, J.M., Coleman, M.C., Herndon, N., Lu, N., Lam, E.T., Anantharaman, T., Sheth, P., and Brown, S.J. (2015). Tools and pipelines for BioNano data: molecule assembly pipeline and FASTA super scaffolding tool. *BMC Genomics* 16, 734.
- Zook, J.M., Hansen, N.F., Olson, N.D., Chapman, L.M., Mullikin, J.C., Xiao, C., Sherry, S., Koren, S., Phillippy, A.M., Boutros, P.C., et al. (2019). A robust benchmark for germline structural variant detection (*Genomics*).

Appendix G: VCF conversion

A standalone VCF conversion script is provided for converting Bionano structural variant (SV) and copy number variant (CNV) calls into VCF format (version 4.2). The Python script is typically run automatically during the *de novo* assembly pipeline and Rare Variant Pipeline, but it can be run as a standalone conversion tool on the command line. The script requires Python 3.7 and Python libraries including pandas and numpy.

The minimum required input is the SMAP file (via the “-s” parameter), which contains SV calls that the pipelines generate. Based on the SMAP input, the script would look for the corresponding alignment files (.xmap and _r.cmap), if they are not explicitly provided as input. The VCF conversion script also optionally takes in CNV calls and converts them. If required input is not found, the script would exit with an error message.

There are also a number of optional parameters:

- -r: path to reference CMAP which is referred to in the SMAP supplied as “-s”. It defaults to _r.cmap in SMAP directory.
- -x: path to alignment xmap which is referred to in the SMAP supplied as “-s”. It defaults to .xmap in SMAP directory.
- -n: sample name for genotype data column. It defaults to “Sample1”.
- -o: prefix for the output VCF. It defaults to have the same prefix as the input SMAP.
- -a: RefSeq assembly accession version. It defaults to “GCA_000001405.1”, which corresponds to Genome Reference Consortium Human Build 37 (GRCh37).
- -b: whether the sample is human [True/False]. It defaults to true. If this parameter is set to “true”, the script would convert map ID (chromosome) 23 to “X” and 24 to “Y” for the sex chromosomes. Otherwise, the script makes no assumption about the input maps and simply outputs the map IDs in the VCF.
- -i: dbVar-required experiment ID. It defaults to 1.
- -m: whether to filter out masked SVs (typically with additional underscored annotations such as “_nbase”). Use 0 to output all calls and 1 to filter out masked calls. It defaults to 1, where masked calls are filtered out.
- -c: path to CNV output file (typically named “cnv_calls_exp.txt” from the copy number analysis pipeline).
- -u: whether to estimate breakpoint uncertainty [True/False]. Default = True.
- -cu: Uncertainty value used to correct breakpoint uncertainty for start and end positions. Default is 30kb

- `--species_reference`: Combined species and genome build information [human_hg38 / human_hg19 / non_human]. If provided with `--chr_cnv_stat`, zygosity on sex chromosomes are adjusted to be hemizygous as needed
- `--chr_cnv_stat`: Path to `cnv_chr_stats.txt` file produced by fractional CNV caller.
- `--skip_confidence_filter`: Do not apply recommended confidence filters and list all variants as PASS [True/False]. Default = False
- `--keep_duplicates`: Keep duplicate SMAP entries in file without filtering. [True/False]. Default = False

The VCF converter bases on the VCF output on the input SMAP/CNV data; however, several modifications are applied during the conversion, such that the values in the VCF may appear to be different than those in the SMAP. These modifications are discussed below.

How are coordinates in the VCF computed?

This depends on the SV type. *For insertions and deletions*, the coordinates in the VCF would appear different from the coordinates in the SMAP. In the VCF, we try to indicate the uncertainty in the breakpoints. Consensus maps are aligned to reference and “outlier” regions are output as structural variants.

Figure 36 shows an example of a deletion of reference sequence where the alignments involve multiple labels on the reference (which we consider to be the general case). In this example, the outlier region (where the SV is located) is bound by the two red match lines. The map has a deletion relative to the reference (reference range r larger than map range m). The “start” is assumed to be between the two labels I_1 and I_2 ; the midpoint between them is output (POS or END). The uncertainty is the distance between the midpoint and I_1 or I_2 (CIPOS or CIEND). On average, it would be closest to the true “start” coordinate. The “end” is output as the midpoint between the two labels I_3 and I_4 .

Figure 41 shows an example of a deletion of reference sequence where the alignments involve two adjacent labels on reference (which we consider to be a special case). In this example, the outlier region is bound by the two red match lines. The map has a deletion relative to the reference ($r > m$). Some sequence in the reference between the two labels I_1 and I_4 is deleted. This is a special case of the above general case. The “start” and “end” are both assumed to be between I_1 and I_4 . The “start” is output as the midpoint between I_1 and the next label I_2 ; the “end” is also output as the midpoint between I_3 and I_4 .

Figure 41. Example of a deletion involving multiple labels on the reference.

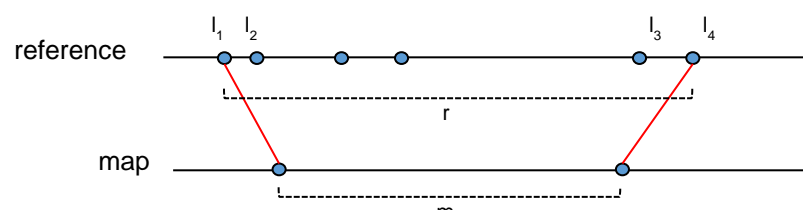
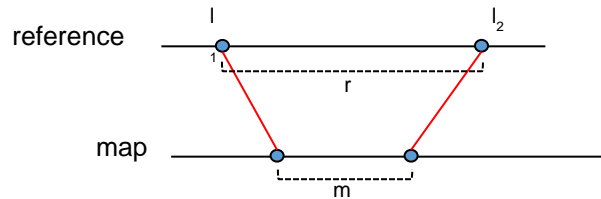


Figure 42. Example of a deletion involved two adjacent labels on the reference.



There is special handling of breakpoint labels with close neighboring labels (where a label on the map aligns to multiple labels on the reference). This could occur if the close labels are not resolved. The script attempts to find a neighboring label away from the cluster of close labels. For example in Figure 42, if l_1 and l_2 were not resolved hypothetically (if they aligned to the same map label), the script would continue to look for the next label on the right (one more label to the right of l_2) to define the possible uncertainty in the breakpoint of the start of the deletion.

For *translocation and inversions breakpoints, and duplications*, the procedure is different. Instead of always looking inwards in the outlier region for the next reference label to define the uncertainty, the script looks for the next reference label to the left and to the right for each breakpoint. The assumption is that because these SV regions tend to be complex, the alignment boundaries may be off by one label. It is equally likely that we over- or under-align by one label. The breakpoint uncertainty is thus the average between the left interval and the right interval. This tends to provide conservative (upper bound) estimates of the breakpoint uncertainty.

For *CNVs* the script simply uses the coordinates from the CNV output without modification, and breakpoint uncertainty is not computed. Generally, the breakpoint uncertainty for CNV calls is higher than for SV calls.

How is size computed (SVLEN)?

For insertion and deletion calls, the size is defined to be the difference between the reference and the query map intervals in the outlier region. For inversion breakpoints and duplications, the size is defined to be the difference between the reference start and end coordinates in the inverted or duplicated regions. For translocation breakpoints, the size is output as 0.

How is orientation for translocation breakpoints computed?

For more recent versions of the SMAP input, the script expects and would use the “orientation” column for annotating translocation breakpoints. If the column is not present, the VCF converter would compute the orientation accordingly. For each side of a translocation breakpoint, the script would compute the orientation based on the alignment orientation (either forward or reverse related to the reference). The orientation can be “+”

or “-“. The final orientation is the combination of the orientation of each side. There are four possible orientations: “+/+”, “+/-“, “-/+”, and “+/+”.

How is confidence computed?

Confidence in the VCF file is computed as the Phred score of the SMAP column “Confidence”: $-10 \cdot \log_{10}(1 - \text{confidence in the SMAP})$. The values are artificially capped at 20 (given that we have limited resolution for the confidence assignment beyond that). When the SMAP confidence score is -1 (e.g., for duplications) the value in the VCF file is “.”.

How is the genotyping (the GT field) computed?

GT is output depending on the zygosity annotation in the SMAP. GT would be “1/1” if the zygosity is “homozygous”, “0/1” if the zygosity is “heterozygous”, and “./.” if the zygosity is “unknown”. For previous versions of the SMAP file, the zygosity column may be absent; if the column is not found, GT would be output as “./.”.

Why SV counts in RVP or de novo assembly report is different from VCF?

The RVP or *de novo* assembly report includes SV counts after clustering and collapsing similar SV calls. The VCF conversion accounts for potentially redundant SV calls by recording only a single entry of a variant with identical type and position. Users may also notice that the number of entries is different. Inversion breakpoints are represented by two lines in the SMAP, but as a single line in the VCF while intra-chromosomal fusion and inter-chromosomal translocations will be represented as two related breakend entries in the VCF while being represented by a single line in the SMAP.

What filter is applied?

“Masked” is output for masked calls. “LowConfidence” is output for variants that do not meet the minimum recommended confidence score for that variant type. “PASS” is output for all other calls. Masking is performed by default during the SV and CNV calling steps using separate masks. The SV mask include regions where false positive translocation calls were made in control samples with no known translocations. These regions are often segmental duplication loci and cannot be aligned uniquely. The CNV mask include regions with elevated coverage noise, defined based on control samples with no known large CNV events. False positive CNV calls are more common in high coverage noise regions. SV and CNV calls overlapping with the masks are masked and of lower confidence.

Appendix H: Variant Allele Fraction calculation

Method

The VAF calculation uses as input the SV calls in the SMAP file and the alignments in the XMAP file between all the consensus maps and the reference genome created by the pipelines. From the alignments we obtain the coverage for each of the consensus map segments, where a segment is the stretch between 2 labels. Each segment of a consensus map is supported by a certain number of aligned molecules (Figure 42). As multiple SV calls may refer to the same underlying SV detected using different consensus maps, to avoid incorrectly considering multiple equivalent alleles at the same locus, we first cluster the SV calls based on the distance between their breakpoints and their SV type (Bionano Genomics proprietary clustering algorithm). The alignments for the consensus maps of the clustered calls are considered together to calculate the VAF of an SV, with the map of the reference genome providing the common reference frame for aggregation. For example, in Figure 43 there is a consensus map containing a deletion; if a second consensus map (not depicted in the figure) containing the same deletion would have been clustered together with it, the coverage of both maps for each for the reference labels l_1, l_2, \dots, l_8 would have been aggregated.

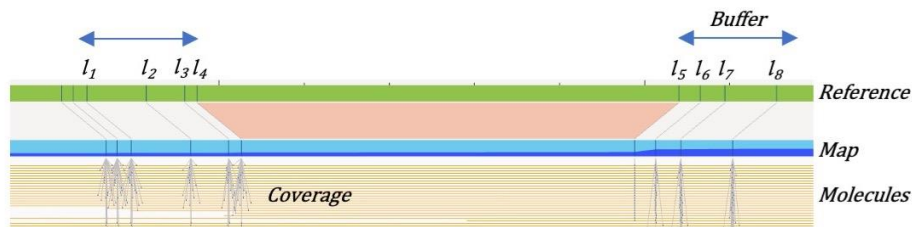


Figure 43. Concepts involved in the VAF calculation algorithm. A consensus map containing a deletion is aligned against the reference map, and the consensus map is created by multiple molecules. Each segment of the reference map may have different coverage on the consensus map, depending on the number of aligned molecules containing the segment. A buffer distance applied at both sides of the deletion breakpoints contains 8 labels, 4 before the start of the deletion, and 4 past the end. Those labels will be used to calculate the VAF.

After the aggregation step, the VAF α_k of a certain allele k for an SV can be inferred by the differences in coverage of a small set of reference labels $D = \{l_1, l_2, \dots, l_L\}$ in the surroundings of the SV breakpoints. In a Bayesian framework, the probability of having a VAF α_k for the allele k given the label set D and the genotype G can be expressed as:

$$P(\alpha_k | D, G) = \frac{P(D | \alpha_k, G)P(\alpha_k | G)}{P(D | G)} = \frac{P(D | \alpha_k, G)P(\alpha_k | G)}{\int_0^1 P(D | \alpha, G)P(\alpha | G)d\alpha}$$

Under the assumption that each label l_i behaves independently, the probability for the observed coverages of the labels in the set D is:

$$P(D | \alpha_k, G) = \prod_{i=1}^L P(l_i | \alpha_k, G)$$

We model the probability for each individual reference label l_i as a binomial distribution:

$$P(l_i | \alpha_k, G) = \binom{M}{M_k} (1 - p_k)^{M-M_k} p_k^{M_k}$$

Where M is the total coverage for the segment starting at label l_i , M_k is the aggregated coverage assigned to allele k after the clustering step, and p_k is the probability of observing a label from a molecule of the allele k aligned to the reference label l_i . If we consider p_k to be the same for all the labels in the set D , we have $p_k = \alpha_k$. We find the VAF for the SV by finding the value of α_k that maximizes $P(\alpha_k | D, G)$.

For duplications calls it is not possible to apply Bayesian inference because the consensus maps may not always contain both whole copies of the duplicated material, making the determination of the correspondence between the labels of the reference, first, and second copies impossible. Additionally, there is only one region in the consensus map where the coverage may be different from the reference map: Namely, around the label of the reference aligned to the beginning of the duplicated material (Figure 44).

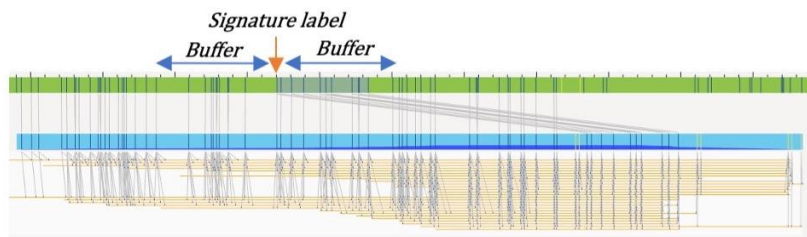


Figure 44. Alignment of a consensus map containing a duplication to a reference map. The reference label aligned to the beginning of the duplicated material (signature label) marks the center of the region used to calculate the VAF. The labels used are those contained inside a buffer distance

We average the coverages over the region defined by the signature label (plus buffers) and apply the naïve approach of:

$$\alpha_k = \frac{C_k}{C_R + C_k}$$

Where C_R is the mean coverage in the region the for the reference map, and C_k is the mean coverage around the critical breakpoint for the consensus maps from allele k .

Performance on simulated data

We tested our method on 10 samples created by mixing simulated molecules containing SVs inserted randomly in the hg38 reference genome together with non-modified simulated molecules from hg38. We ran the pipelines for values of 80, 180 and 300X effective coverage, and 5, 10, 20, 40, 80, and 100% simulated VAF. For RVP we had

159/160 successful runs (Table 11) and for *de novo* assembly 129/130 (Table 12). The median error (difference between the simulated and calculated VAF) values were in the range 0-0.11 for all combinations of VAF and coverage, for both pipelines. The third quartile of the error was 0-0.20 for *de novo* assembly and 0-0.50 for RVP, indicating that the VAF values for the RVP pipeline are expected to have wider variability and have poorer performance for the most incorrect SVs. Although minimal, we observed highly incorrect VAF values (Max column) in all runs, due to multiple causes: false positive calls, SVs where the first clustering stage failed to group all equivalent SV calls together, or genomic regions with shallow label coverage that strongly influenced the Bayesian inference procedure.

When adopting as a success criterion that the difference between the simulated and calculated VAFs should be under 10%, we have found as a trend that the method behaves differently in different contexts. For RVP (Table 17) the highest success rate is in the VAF range 5-20% (78-96%), with decreased performance afterwards (55-70%). For *de novo* (Table 18), the highest success range is at the 100% VAF range (94-96%), followed by the 5-40% VAF range (52-89%). The poorest performance is at 80% VAF, with success rate 48-52%.

Effective Coverage	VAF simulated	Samples	SV count	Median	3 rd quartile	Max	Variants with VAF error <= 10%
300	5%	10	12271	0.02	0.03	0.95	96%
300	10%	10	16849	0.02	0.05	0.90	91%
300	20%	10	18190	0.03	0.07	0.80	84%
300	40%	10	21135	0.05	0.14	0.60	67%
300	80%	10	24366	0.07	0.22	0.80	55%
300	100%	10	25742	0.07	0.29	1.00	56%
180	10%	10	13176	0.03	0.05	0.90	95%
180	20%	10	16143	0.04	0.08	0.80	83%
180	40%	10	18486	0.06	0.13	0.60	67%
180	80%	10	19345	0.07	0.20	0.56	56%
180	100%	10	22533	0.06	0.24	1.00	62%
80	10%	9	6972	0.03	0.50	0.90	97%
80	20%	10	12519	0.06	0.10	0.80	78%
80	40%	10	15223	0.08	0.15	0.60	59%
80	80%	10	17521	0.07	0.19	0.80	59%
80	100%	10	18374	0.04	0.15	1.00	70%

Table 17. Performance of the VAF calculation algorithm for RVP runs on 10 simulated samples at different VAF and effective coverage. Only variants of high confidence are considered for the SV counts. The columns for the Median, 3rd quartile, and Max describe the distribution of the difference between the simulated and calculated VAF. The last column shows the percentage of calls that had an error lower than 10% for the VAF.

Effective Coverage	VAF simulated	Samples	SV count	Median	3 rd quartile	Max	Variants with VAF error <= 10%
300	5%	5	705	0.08	0.14	0.95	52%
300	10%	5	2924	0.06	0.12	0.90	63%
300	20%	5	8415	0.09	0.15	0.80	54%
300	40%	5	14370	0.08	0.12	0.60	63%
300	80%	5	16473	0.09	0.18	0.80	48%
300	100%	5	23371	0.00	0.00	1.00	96%
180	10%	10	3557	0.05	0.08	0.90	76%
180	20%	10	15869	0.05	0.08	0.80	85%
180	40%	10	22424	0.05	0.10	0.60	72%
180	80%	10	36445	0.09	0.20	0.80	52%
180	100%	10	35573	0.00	0.00	0.99	94%
80	10%	10	1950	0.05	0.09	0.90	84%
80	20%	10	16104	0.04	0.07	0.80	89%
80	40%	10	25643	0.05	0.10	0.60	74%
80	80%	9	30398	0.11	0.20	0.77	49%

80	100%	10	44533	0.00	0.00	1.00	95%
----	------	----	-------	------	------	------	-----

Table 18. Performance of the VAF calculation algorithm for *de novo* assembly pipeline runs on 10 simulated samples at different VAF and effective coverage. Only variants of high confidence are considered for the SV counts. The columns for the Median, 3rd quartile, and Max describe the distribution of the difference between the simulated and calculated VAF. The last column shows the percentage of calls that had an error lower than 10% for the VAF.

Based on these simulated results, we expect the *de novo* assembly pipeline to produce generally more accurate and consistent results than RVP when considering all the variants. We nevertheless expect the VAF calculations for RVP to have good accuracy for rare variants (VAF 5-20%)

Reproducibility

We tested the reproducibility of our method by running the VAF calculations on a triplicated experiment of a tumor/normal pair of samples, for both RVP and *de novo* assembly pipelines. For the RVP runs, it is apparent that there is correspondence between the values of the VAF for each variant across the replicates (Figure 45). An analysis of variance (Table 3) confirmed that the null hypothesis of the sample not being a factor of variation is plausible, both for the tumor replicates (p-value 0.29) and normal replicates (p-value 0.20). The same results were obtained for the *de novo* assembly pipeline (data not shown).

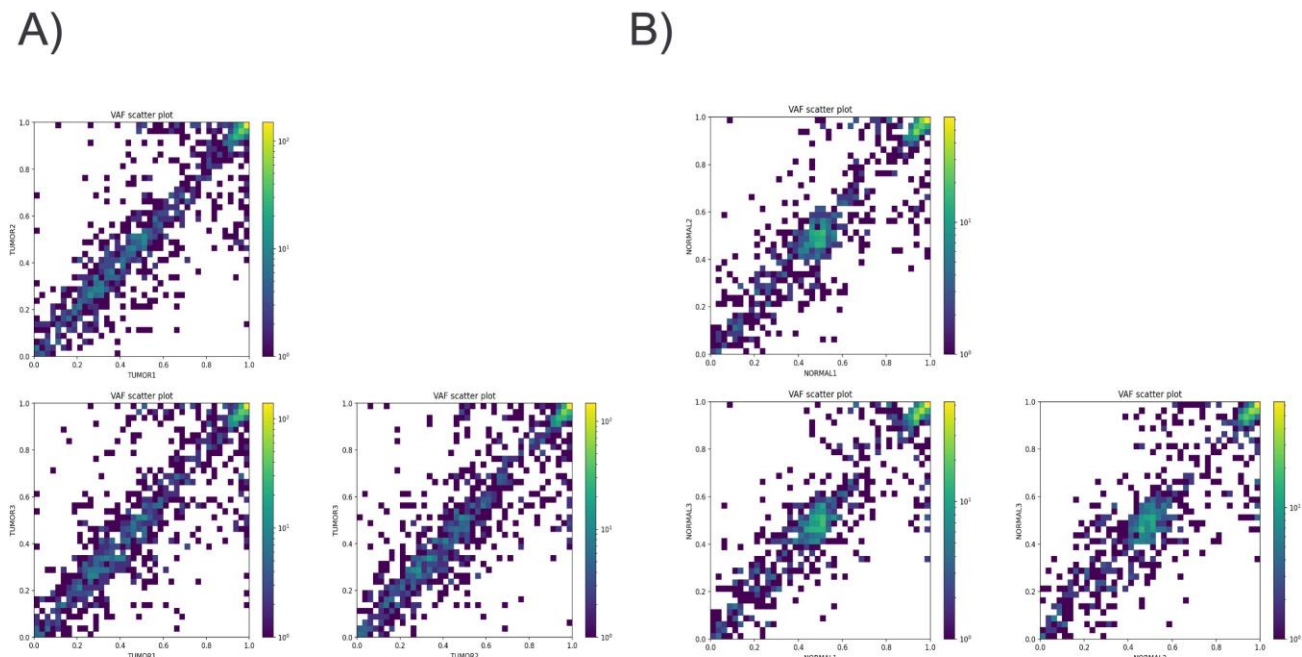


Figure 45. Pairwise visual comparison of the VAF values between replicate samples processed with RVP. For each plot, the x-axis contains the VAFs for the first sample replicate, the y-axis the VAFs for the second sample replicate, and each point represents a SV. The color scale indicates the number of SVs with a certain (x,y) value combination. The ideal perfect concordance would appear as a diagonal in the plots. A) Comparison between the tumor replicates. B) Comparison between the normal replicates.

Tumor samples					Normal samples				
Effect	Sum squares	Degrees of freedom	F statistic	p-value	Effect	Sum squares	Degrees of freedom	F statistic	p-value
Variant	824.80	1285	60.56	0.00	Variant	552.02	956	68.58	0.00
Sample	0.03	2	1.23	0.29	Sample	0.03	2	1.60	0.20
Residual	27.18	2565			Residual	16.08	1890		

Table 19. Blocked-ANOVA tests for the analysis of the repeatability of VAF results from RVP. Left, grey: Analysis for the 3 tumor replicates. Right, white: Analysis for the 3 normal replicates.

Performance on experimental data

It is difficult to obtain an experimental dataset where all the VAF values are known for all the SVs present in the sample. To overcome that difficulty and at the same offer results on the performance of the algorithm on experimental data, we created an in-silico dilution experiment where molecules from a tumor sample and its paired normal are mixed in different proportions, and we look at the VAF values for the somatic variants across the different mixes. We observed the expected progressive decrease in the VAF values (Figure 46), with values close to 100% when there are only tumor molecules in the mix and VAF distributions located around the simulated tumor purity.

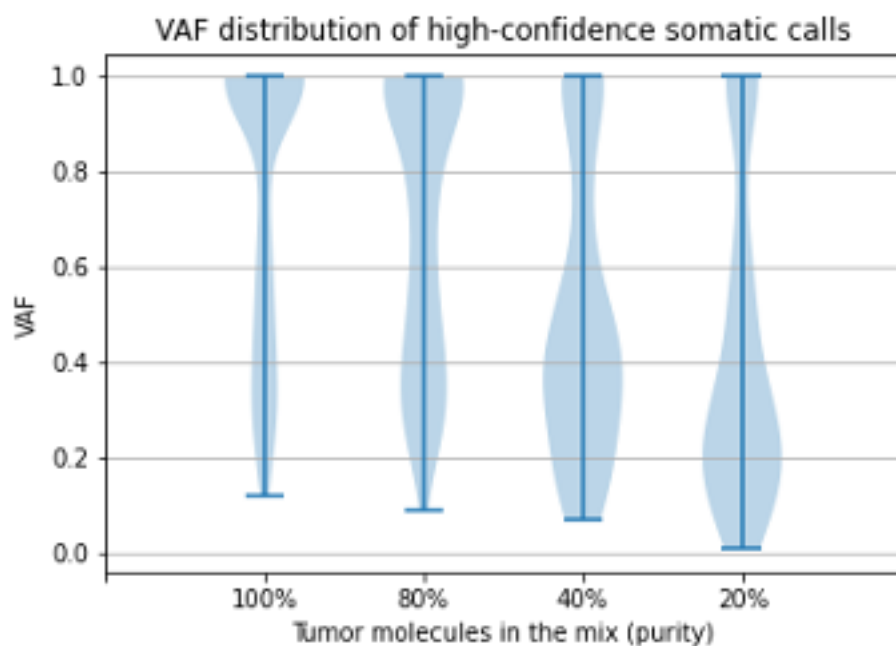


Figure 46. Violin plots with the histogram of VAF values for in-silico mixes of paired tumor/normal samples, at different tumor purity. Only somatic high-confidence SV calls are considered.

Segmentation of the whole genome VAF plot

The plot of VAF vs genomic coordinates shows patterns of variation of the VAF across the genome. For a normal healthy germline sample two main VAF bands are expected at VAFs 50% and 100% for heterozygous and homozygous variants, respectively. When there are changes in the chromosomal structure of the sample, such as aneuploidy events, it is possible to see stretches of VAF values (segments) with different VAF and changes in the number of bands. To help discovering and analyzing those events, we have developed an algorithm to calculate the boundaries and the median VAF of the segments of the autosomal chromosomes in the plot. The algorithm is as follows:

First, the non-informative homozygous SVs (default threshold VAF 97%) are not considered for segmentation. Next, the raw VAF values are transformed to the “major allele” VAF values by applying the transformation:

$$mVAF = \begin{cases} VAF & \text{if } VAF > 0.5 \\ 1 - VAF & \text{if } VAF \leq 0.5 \end{cases}$$

We filter the mVAF outliers using triplet filtering (*Staaf et al., 2008*) with default threshold 0.40 and segment the mVAF plot using circular binary segmentation (CBS) (*Olshen et al., 2004*) on the non-centromeric areas. As a measure more robust to outliers, we consider as the mVAF of a segment the median of the mVAFs from the points spanned by it. Next, we compare the median mVAFs of consecutive segments using the Man-Whitney test, and two segments are joined if the null hypothesis of having the same median is satisfied (p-value 0.0001). The procedure is repeated until no consecutive segments can be joined. The segment boundaries and median values obtained from the mVAF are transferred to the raw VAF plot. When the median VAF of a segment is below a certain threshold (default 0.56), we consider that the possible allelic imbalance present in the region has not been sufficiently resolved and recalculate the median VAF for the segment, this time using the VAF values and not the mVAF. The result is a more accurate median value for the segments in normal, heterozygous parts of the genome. Finally, to avoid the creation of many small, noisy segments, or segments in regions with loss of heterozygosity, segments shorter than 10 Mbp, supported by less than 10 SVs, or with variant density lesser than 0.3 SV/Mbp are discarded.

Note: Currently, the default thresholds values employed by the segmentation algorithm are not customizable.

References

- Olshen, A.B. *et al.* (2004) Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics*, **5**, 557–572.
- Staaf, J. *et al.* (2008) Segmentation-based detection of allelic imbalance and loss-of-heterozygosity in cancer cells using whole genome SNP arrays. *Genome Biol*, **9**, R136.

Appendix I – Custom CNV control data

Scripts for pre-processing custom control data are for advanced users to prepare their own control data for the CN pipeline. For example, custom control data are needed for non-human CNV detection, which the pipeline does not support by default. The scripts are located in the /preprocess directory under the main CN pipeline directory.

For DLE-1 data, the following two scripts can be used for generating the auxillary input for the fractional CN analysis module. There is no need to generate input for the integer CN analysis module.

Step	Objective	Script name	Input	Output
Generate control data for fractional CN module	Generate the necessary control reference for fractional CN module from control datasets	generate_control_FractCNV.R	A list of paths to alignmolvref output	Processed control r.cmap file
Generate CNV mask for fractional CN module	Generate files to mask high-variance regions for fractional CN module from control datasets	generate_cnv_mask_FractCNV.R	A list of paths to alignmolvref output, control reference file	BED file

generate_control_FractCNV.R

This wrapper script generates the necessary model parameter data from control samples for the fractional CN pipeline. It takes as input a file that contains a list of paths where each path is the output folder for the molecule-to-reference alignment for each control sample. The path can be the “alignmolvref/merge” folder from the *de novo* assembly or Rare Variant Pipeline. Alternatively, if the user independently generated the molecule-to-reference alignment for the control samples, each output path in the list needs to contain one .xmap file and one r.cmap file for each control sample. The script generates a r.cmap file that can be used as input to the CN pipeline. Default for sexChr1 is 23 (cmapID of human chrX), default for sexChr2 is 24 (cmapID of human chrY).

Usage

```
Rscript generate_control_fractCNV.R --path <input file> --output <output directory> --sexChr1 <cmapID of homogametic sex chromosome> --sexChr2 <cmapID of heterogametic sex chromosome>
```

generate_cnv_mask_FractCNV.R

This wrapper script generates CNV masks from control samples for the fractional CN pipeline. The script identifies regions with high variance in CN calls in controls; these often indicate centromeres and other regions difficult to call CNs in. Users may pass this file to the CNV pipeline in order to mask out these regions when calling CNVs in case samples.

Input

It takes as input a file that contains a list of paths where each path is the output folder for the molecule-to-reference alignment for each control sample. The paths can be the “alignmolvref/merge” folder from the *de novo* assembly or Rare Variant Pipeline. Alternatively, if the user independently generated the molecule-to-reference alignment for the control samples, each output path in the list needs to contain one .xmap file and one r.cmap file for each control sample. The input parameter filePostfix can be provided to indicate the correct r.cmap file to use. The input parameter refFile is the CNV control reference file output by generate_control_FractCNV.R. Defaults are: sexChrType “xy”, sexChr1 23 (cmapID of human chrX), sexChr2 24 (cmapID of human chrY).

Output

The script generates a BED file, masks_cnv.bed.

Usage

```
Rscript generate_cnv_mask_fractCNV.R --controlDirsList <input file> --refFile <cmap file> --outputDir
<output directory> --SVMask <optional SV mask BED file> --sexChr1 <cmapID of homogametic sex
chromosome> --sexChr2 <cmapID of heterogametic sex chromosome> --filePostfix <'merged_r.cmap'
or '_r.cmap'>
```

For nickase data, the follow scripts can be used for generating the auxillary input for the integer CN analysis module. The order of the rows in the table indicates the recommended sequence of use.

Step	Objective	Script name	Input	Output
Simulation	Simulate data for confidence table calculation and cross validation	001_sim_wrapper.R	autonoise and alignmolvref files from assembly pipeline	Simulated templates
Control samples preparation	Filter based on input quality and generate eigenvectors by SVD	002_control_wrapper.R	Set of control data r.cmap	Control RData file and eigenvector files
Select control samples	Select the best control datasets based on cross validation	003_crossValidation.R		Performance by selection of control datasets
(optional) Re-generate control samples	Re-process the control datasets based on results from previous step	002_control_wrapper.R	Set of control data r.cmap	Control RData file and eigenvector files
Determine principle components to use	Find the best number of principle components to use based on cross validation	004_componentsOptimization.R		Performance by number of PC
Generate confidence tables	Use the PCs, eigenvectors and templates from previous steps to generate confidence tables	002_control_wrapper.R		Confidence tables

After pre-processing, users can use the output to update the default_param_table.txt, the eigenvector files and confidence tables under /parameters and /testing_data in the main CNV pipeline.

Below details the input, output, and usage of each wrapper script. The usage of each script is stated in their top script lines. Here, we demonstrate general usage. Detailed explanations for each parameter are in the defaultParams section inside each program file.

001_simulation_wrapper.R

This wrapper uses assembly pipeline files as input. It simulates CN events, downsamples the molecules in the corresponding regions, and runs autonoise and alignmolvref steps to form a new r.cmap containing coverage

features that correspond to the simulated CN events. The templates and r.cmap are output to a testing_data/ directory.

Input

This wrapper requires an csv file pointing to the path of each input files. For example, inputs/simulation_dle1_saphyr_hg19_11092017.csv:

```
filePath sampleId inputType fileType
/saphyr_dle1_hg19/SAMPLE2/autoNoise1_rescaled.bnx sample18 simulation scaledbnx
/saphyr_dle1_hg19/SAMPLE2/autoNoise1.errbin sample18 simulation errbin
/saphyr_dle1_hg19/SAMPLE2/alignmolvref_merge.xmap sample18 simulation xmap
/saphyr_dle1_hg19/SAMPLE2/alignmolvref_merge_r.cmap sample18 simulation rcmap
```

If the first line is a raw bnx file, then the fileType should be replaced by “bnx”. Also, sampleId refers to the column name in the eigenVector.RData files generated by 002_control_wrapper.R. A wrong sampleId will lead to errors in subsequent step.

Output

1. testing_data/ containing raw_coverage.RData templates.RData templSignals.RData
2. New alignmolvref r.cmap
3. Intermediate bnx, moleculeID and copyNumberTemplate.RData based on the sampling strategy

Usage

The default assembly command line call is as below. For more parameters, please see the defaultParams section in 001_sim_wrapper.R.

```
Rscript 001_sim_wrapper.R --alignmentType [platform] --enzyme [enzyme] --reference [reference] --
simulationInput [.csv] --resultsDir [dir] --testingDir [dir] --outputFolderName [string] --nBaseFile [.bed] --
referenceCMAP [.cmap] --xml [optArg.xml] --refAligner [dir] --pipeline [dir]
```

002_control_wrapper.R

This wrapper filters the control datasets based on --controlSdCutoff value, which is the standard deviation of the scaled coverage across chromosomes. Once the sample selection is determined, it can also generate the confidence table using the testing_data files generated by 001_control_wrapper.R.

Input

This wrapper requires a directory containing alignmolvref r.cmap of the control samples and a testing_data directory containing simulation output.

Output

1. control.RData, a matrix storing scaled control sample coverages
2. eigenvectors.RData files
3. confidence.RData files if --confidence option is used

Usage

```
Rscript 002_control_wrapper.R --controlsDir [dir] --enzyme [enzyme] --parametersDir [dir] --outliersProbability [float] --componentsRemoved [1:max PC number] --selectControls --controlSdCutoff [float] --confidence --testingDir [dir]
```

003_crossValidation.R

This wrapper helps find the best control datasets to use. By default, it sorts the control datasets by the standard deviation of their scaled coverages across chromosomes and evaluates the performance from using the fewest samples to using all samples. The users can also use a random sampling strategy. In this cross validation, the number of components or principle axes (--componentsRemoved option) is fixed for different sampling. However, if the sample size was smaller than the specified number of components, then the number of components would be reduced to the sample size.

Input

This wrapper requires testing_data files from 001_sim_wrapper.R, and scaled coverage file control.RData from 002_control_wrapper.R. In addition, it requires a file specifying the column header (control sample names) in control.RData to be included in the control sample selection, and the column header in testing_data/raw_coverage.RData (simulated sample names) to be included in the simulation sample selection. An example file is 004_componentsOptimization_input.R

Output

1. Performance data by sample size in table and pdf plot format
2. Intermediate cnv_call and cnv_rcmap tables

Usage

```
Rscript 002_control_wrapper.R --controlsDir [dir] --enzyme [enzyme] --parametersDir [dir] --
outliersProbability [float] --componentsRemoved [1:max component number] --selectControls --
controlSdCutoff [float] --confidence --testingDir [dir]
```

004_componentsOptimization.R

This wrapper helps optimize the number of components to remove. It tests the performance starting from 0 to n components by default, where n is the control sample size. The input and output files are similar to those for 003_crossValidation.R, except that the performance table and plot are by number of components instead of sample size.

Usage

```
Rscript 004_componentsOptimization.R --testingDataDir [dir] --enzyme [enzyme] --alignmentType
[platform] --reference [reference] --parametersDir [dir] --controlCoverageRData [.RData] --resultsDir [dir]
--nGenomesPerSample [integer] --componentsN [min component number : max component number]
```

Technical Assistance

For technical assistance, contact Bionano Genomics Technical Support.

You can retrieve documentation on Bionano products, SDS's, certificates of analysis, frequently asked questions, and other related documents from the Support website or by request through e-mail and telephone.

Type	Contact
Email	support@bionanogenomics.com
Phone	Hours of Operation: Monday through Friday, 9:00 a.m. to 5:00 p.m., PST US: +1 (858) 888-7663
Website	www.bionanogenomics.com/support