

667 Supplementary Material

668

669 S1 OMKar Report

670 For ease of use, OMKar optionally compiles a report that includes 1) the observed chromosomal
671 segment lists with reconstructed karyotype, 2) a list of interpreted events using ISCN notation, 3)
672 a visualization of each corresponding chromosome with cytoband and event labels, and 4) a table
673 of important genes that are near the breakpoints of an event or have copy number (CN) alteration.
674 An html version of the report is prepared for easy viewing.

675 By default, OMKar outputs a molecular karyotype in a text format to unambiguously describe
676 the karyotype as follows: It first lists all defined segments across the reference genome. For each
677 segment, the following information is provided: segment number, chromosome, start and end coordi-
678 nates, and the graph nodes representing the segment. Each segment is represented by two nodes,
679 connected by a segment edge. All segments are forward-oriented (i.e., the end coordinate is greater
680 than or equal to the start coordinate) and are sorted by chromosome groups (from 1 to Y) and
681 increasing coordinates. The segments are non-overlapping, with no gaps between them, although
682 telomeric regions of the reference genome may be excluded.

683 Following the segment definitions, OMKar reports the reconstructed paths that represent a
684 karyotype. Each path consists of a list of segments in the format “Path number = segment number
685 followed by direction.” The segments are traversed either in the forward (+) or reverse (-)
686 direction, where ‘+’ indicates traversal from the start to the end of the segment, and ‘-’ indicates
687 traversal from the end to the start of the segment. For example, “Path1 = 1+ 2+ 3-” means that
688 the path traverses segment 1 in the forward direction, segment 2 in the forward direction, and
689 segment 3 in the reverse direction. Additionally, the number of centromeres present in each path
690 is reported, which should ideally be one, indicating a valid chromosome structure.

691 S2 Terminal Event Simulation and Validation

692 We simulated a total of 117 terminal Structural Variations. During analyses, it was found that
693 only 50 of the 117 SVs (42.7%) were reconstructed under the 0.2 Mbp matching distance (same

694 distance used for non-terminal SVs; Supplementary Table S1). However, 78 of the 117 SVs (66.7%)
695 were reconstructed under the 5 Mbp matching distance. Because one of the endpoint was in the
696 terminal masking region, most terminal SV edge did not have an SV call. For unbalanced SVs,
697 OMKar reconstructed these variations using the information from CNV calls, which had much
698 greater error in the event boundaries. For inversion, as a balanced SV, missing the SV call meant
699 OMKar had no other information on the SV, thus, resulting in a much lower recall rate at all
7 distances. In addition, terminal events such as arm deletion may result in a false loss of centromere
7 1 call, which result in a relatively higher false positive chromosomal-loss aneuploidy reconstruction
7 2 from OMKar.

7 3 From the real data we received, we observed no terminal structural variation, hinting that ter-
7 4 minal SVs were far less frequent than our simulation. For future reference, OMKar may incorporate
7 5 additional information such as alignment contigs to improve the boundary accuracy and recall rates.

7 6 **S3 OMKar details**

7 7 **S3.1 Filtering SV and CNV calls**

7 8 OMKar utilizes the CNV and SV calls from the Bionano pipeline. We applied several filters to
7 9 improve data quality. First, we filtered out low-confidence CNV calls, defined as those with a
7 10 confidence level of 0.95 or below, as well as those located in masked regions that could interfere
7 11 with the analysis. Next, we excluded CNVs smaller than 200 kbp. However, if an excluded CNV
7 12 was supported by a corresponding SV, it was retrieved during later processing, ensuring no relevant
7 13 variations were missed.

7 14 Following these steps, SV calls were filtered based on variation-specific confidence thresholds
7 15 established by BioNano pipeline: translocations ($T_{trans} = 0.05$), inversions ($T_{inv} = 0.7$), and indels
7 16 ($T_{indel} = 0$). Additionally, SV calls in masked regions were discarded. Breakpoints from SVs were
7 17 processed by sorting them based on chromosomal and genomic coordinates and merging adjacent
7 18 breakpoints within a 50 kbp window to simplify the breakpoint graph. To ensure accurate repre-
7 19 sentation, CNV segments were split if breakpoints occurred within their boundaries, guaranteeing
7 20 that all breakpoints exclusively connect the terminal coordinates of segments.

721 **S3.2 Converting to linear constraints**

722 To tackle the inherent non-linearity of the absolute value and sign functions within the optimization
 723 problem, we introduce new variables and constraints to linearize these functions.

To convert the sign function into linear constraints, we employed the following approach: The ILP formulation for $s'_e = \text{sgn}(s_e)$ when $s_e \geq 0$ is:

$$s_e \geq s'_e,$$

$$s_e \leq 1000 \cdot s'_e$$

724 where s'_e is a binary variable. This formulation ensures that $s'_e = 1$ when $s_e > 0$ and $s'_e = 0$ when
 725 $s_e = 0$, while maintaining linearity in the constraints. The constant 1000 is used as a sufficiently
 726 large value to approximate an upper bound on s_e .

727 The absolute value function can be linearized $y_v = |x_v|$ as follows:

$$y_v \geq x_v,$$

$$y_v \geq -x_v,$$

$$y_v \geq 0,$$

728 where y_v represents the absolute value of x_v . These constraints ensure that y_v is equal to $|x_v|$ by
 729 considering both the positive and negative cases for x_v . Since our objective is to minimize y_v , it
 73 naturally corresponds to $|x_v|$. Thus:

$$\begin{aligned} \text{Objective Function} &= \gamma \sum_{v \in V} (c_v + x_v - r_v - \sum_{e \in SV} s_e) + \sum_{v \in V} \beta_v |x_v| - \sum_{s_e \in SV} e \text{sgn}(s_e) + \lambda \sum_{v \in V} o_v \\ &= \gamma \sum_{v \in V} (c_v + x_v - r_v - \sum_{e \in SV} s_e) + \sum_{v \in V} \beta_v y_v - \sum_{s_e \in SV} e s'_e + \lambda \sum_{v \in V} o_v \end{aligned}$$

731 **S3.3 Parameter determination**

732 In determining e , we consider the characteristics of structural variation e . If e signifies a high-
733 confidence call and connects two distinct chromosomes, we set $e = 72$ otherwise, $e = 9$. A
734 higher value is given to intra-chromosomal translocations to prioritize their preservation in the
735 final karyotype. For the calculation of β_v , we introduce a penalty parameter that becomes more
736 pronounced as segments grow in length. To achieve this, we define $\beta_v = 4 \times \lceil \frac{\text{length } v}{5 \cdot 10^6} \rceil$. This
737 parameter β_v ensures that the penalty for modifying longer segments is appropriately weighted.

738 The parameter γ is empirically determined and set to a value of 5, while λ is fixed at 1 to
739 ensure that if two candidate decompositions are nearly identical, our preference is to select the one
740 with fewer odd vertices. This empirical choice addresses the unique requirements of our approach,
741 providing a balanced framework for the optimization process (Supplementary Table S4).

742 **S3.4 OMKar Eulerian Path algorithm**

743 To reconstruct chromosomal structures from breakpoint graphs, OMKar computes an Eulerian
744 path, ensuring each edge is traversed exactly once while maintaining biologically meaningful con-
745 straints. It begins with a breakpoint graph, where vertices represent segment boundaries, and edges
746 denote segment continuity, reference adjacencies, or breakpoint rearrangements. To enforce chromo-
747 somal structure, the algorithm prioritizes reference edges (minimizing deviation from the reference
748 genome), followed by breakpoint edges (capturing structural variations), and lastly, segment edges.
749 Using a recursive depth-first traversal, it starts from a telomeric vertex—representing a natural chro-
750 mosomal endpoint—and evaluates edges with a validity function that prevents breaking segment
751 integrity or disconnecting essential graph components. The pseudo-code is given in Algorithm 1.

Algorithm 1: Find Eulerian Tour Starting from Vertex v in Graph G

Data: Graph G

Result: Eulerian Tour

Function `ValidEdge`(edge $e(v, w)$, edge $e(u, v)$):

```
  if  $\deg v == 1$  then
     $\quad$  return True
  else if  $e(v, w)$  or  $e(u, v)$  is segment edge then
    if  $e(v, w)$  is not bridge edge then
       $\quad$  return True
     $\quad$  return False
```

Function `FindEulerianTour`(current vertex v , previous vertex u):

```
752   Add  $v$  to the Eulerian tour list  $E$ ;
    Initialize an empty list  $W$ ;
    for each edge  $e(v, w)$  do
      if ValidEdge( $e(v, w)$ ,  $e(u, v)$ ) then
         $\quad$  Add  $w$  to the valid edges list  $W$ 
    Sort  $W$  based on the edge types ( $E_r, E, E_s$ );
    Pop  $w$  from list  $W$ ;
    remove  $e(v, w)$  from  $G$ 
     $\quad$  FindEulerianTour(vertex  $w$ , vertex  $v$ )
```

Initialize an empty list E ;

`FindEulerianTour`(Telomeric vertex v , -1);

753 **S3.5 OMKar Report: Event Interpretation**

754 After segregating the chromosomes, OMKar interprets the structural variations in each chromoso-
755 mal cluster using ISCN notation(Hastings et al., 2024) as follows:

756 1. For each reconstructed chromosome in the chromosomal cluster, its chromosomal identity is
757 assigned based on the highest represented centromere, and if it is acentric, by the overall
758 highest-represented chromosome in the remaining segments.

759 2. In a pre-processing step, the Wild Type (WT) chromosome corresponding to the reconstructed
760 chromosome is segmented prior to alignment, using the segments in the reconstruction.

761 3. Next, OMKar performs an alignment between each reconstructed chromosome and the cor-
762 responding Wild Type (WT) using an alignment that maximizes the longest common subse-

763 quence, with a linear penalty for indels.

- 764 4. After alignment, blocks of aligned segments are separated into three types: 1) concordant
765 block represents matching between the reconstruction and the WT, 2) insertion block rep-
766 resents segments inserted in the reconstruction, and 3) deletion block represents segments
767 deleted in the reconstruction. Adjacent blocks of the same type, and representing contiguous
768 genomic coordinates are merged to minimize the total number of blocks.
- 769 5. The final interpretation assigns a representative SV name (using ISCN nomenclature) to each
770 insertion or deletion block, reporting a potential cause of the deviation from the WT. Each
771 SV has its unique signature in the combination of block types (Supplementary Table S5). For
772 example, an inversion is always an insertion block of inverted segments next to a deletion block
773 of non-inverted segments. In this example, it is more likely that a single inversion resulted
774 in both the deletion and the insertion, instead of two separate events (inverted insertion
775 and deletion). Similarly during the interpretation step, a compound SV is preferred over a
776 sequence of simpler SVs (preference is given to reducing the total number of events).

777 To minimize the number of events, event types are resolved from the most complex to the least,
778 using the following preference order: (1) inter-chromosomal balanced translocation and transposi-
779 tion; (2) intra-chromosomal balanced translocation and transposition; (3) other intra-chromosomal
780 variations. Finally, if any deletions or insertions remain unaccounted for, they are marked as simple
781 deletions or duplicated-insertions. During each variation type's resolution, each un-resolved block
782 is iterated over, with the goal of being associated with the signature of that variation type. Bal-
783 anced translocations attempt to associate anywhere in the cluster or chromosome, for inter- and
784 intrachromosomal search, respectively. All other variations associate adjacent blocks.

785 For balanced translocations, the following step is performed after all insertion and deletion
786 blocks are resolved. All balanced translocations are initially denoted as *transposition*, associating a
787 deletion block and a non-adjacent insertion block of an overlapping set of segments, with allowance
788 for small indels. When transpositions form a cycle, they are interpreted as *balanced translocations*.
789 This is implemented recursively, by jumping between each associated deletion-insertion pair, and
790 then looking for nearby transposition blocks of the opposite type. When exhausted, if a cycle
791 is formed, an n -break balanced translocation is interpreted, and otherwise, all transpositions are

792 interpreted as individual transpositions.

793 When associating between different blocks, a procedure called “seed-matching” is applied. For
794 each block that is currently being resolved, it searches for an associated block, such that the common
795 subsequence between the blocks (in indivisible-unit of each segment) is sufficiently long (10 kbp by
796 default). The sizes of the flanking segments not matched may be limited by an “indel allowance”.
797 For example, an insertion block of (B^- , A^-) next to a concordant block of ($B + C +$) will be
798 interpreted as a left-duplication-inversion only if A^- is less than 50 kbp. On the other hand, the
799 size of $C +$ is irrelevant for associating a duplication-inversion, as it is not between the two blocks.
8 These allowances are determined empirically (Supplementary Section S4).

8.1 S3.6 OMKar Report: Visualization

8.2 The visualization is achieved by intersecting the coordinates in the observed chromosomal segment
8.3 list with a given cytoband coordinate table. Each band pattern is then stacked and displayed using
8.4 Matplotlib (v3.7.5). In an alternative plot, the bands are segregated by the segment boundaries in
8.5 the OMKar output, instead of the cytoband. The event label is then applied to the corresponding
8.6 segment location on the band stack.

8.7 By default, the table of important genes comes from genes that have CN alteration or within
8.8 20 kbp of an event’s breakpoint. This list of genes are further filtered to only include those in
8.9 the Developmental Disorders panel in the Gene2Phenotype database (DDG2P)(Thormann et al.,
81 2019).

811 S4 Analyses of Event Distances

812 The resolution of OGM was determined by the matching distances between the Truth SV edge and
813 the reconstructed SV edge, using KarCheck. Observations from this analyses were used to estimate
814 the sensitivity of OMKar event interpretation (Supplementary Table S5).

815 When applying KarCheck between the simulated and the reconstructed karyotypes, history logs
816 from KarSim’s output was used to label SV edges on K_t . This step marks each SV edge with its
817 causal rearrangement. During KarCheck matching, the distance between a pair of matched edges
818 were recorded with the causal SV type. This distance is further separated as the two distances of

819 the endpoint.

82 Then, the matched edge from K_t was searched in OGM’s SMAP output (contains all SV calls),
821 with a proximity of 50 kbp, equal to the SV call merging distance of OMKar. From Fig. S7, we
822 observed duplication inversions tend to have much larger distances than all remaining structural
823 variations simulated. Most of the non-duplication-inversion were reconstructed with less than 5
824 kbp distance from the truth, much lower than the average gene size of 10-15 kbp. Thus, if an SV
825 contains or interrupts a gene, OGM and OMKar are likely to reconstruct the correct boundary to
826 perform genotype-to-phenotype inference. Additionally, for all distances greater than 50 kbp, the
827 true SV edge was not captured with a high proximity in the SMAP, therefore, OMKar either had
828 to reconstruct based on a distant SV call or had to infer the missing SV call based on CNV call.

829 **S5 KarSim**

83 **S5.1 KarSim module for simulating karyotypes**

831 The KarSim module outputs a molecular karyotype file, a corresponding FASTA file, and a history
832 log with event segments and edges noted are outputted for downstream usage. The molecular
833 karyotype and history log can be used for KarCheck comparison, while the FASTA file can be used
834 as input for simulating the sequencing technology of choice, given that many of such simulators
835 are already available. KarSim is publicly available at <https://github.com/MolecularKaryotype/KarSimulator>.

837 Three steps are taken in order, with step 2 being optional:

- 838 1. A template karyotype is created given the counts of autosomes and sex chromosomes.
- 839 2. (Optional) a series of SVs are applied to the karyotype
- 84 3. The FASTA formatted sequence of the rearranged chromosomes is generated.

841 All intermediate files in steps 1 and 2 are molecular karyotypes, which can be read-in for multiple
842 parallel edits or outputting the FASTA file.

843 There are two methods to introduce additional SVs to a karyotype: 1) manual addition of SVs
844 given the variation type and exact boundaries, and 2) using a parameter JSON file that contains
845 the number of SVs, each SV type’s likelihood, max/min size for each SV type, terminal occurrence

846 likelihood etc. Both processes can be applied multiple times and in an interleaved fashion. The
847 types of SVs supported can be found in Supplementary Table S6. In addition, the user has the
848 option to include a masking file such that SVs are not generated within proximity to any of these
849 regions. Another parameter can be passed in to prevent events that result in the formation of
850 segments smaller than a certain size, useful for sequencing technologies which have a resolution
851 limit.

852 **S5.2 Implementation**

853 KarSim is implemented using Python (version 3.9). A flowchart of the algorithm is in Supplemen-
854 tary Fig. S8A. Each chromosome in a karyotype is represented as a sequence of oriented segment
855 objects of the hg38 genome, denoted by the chromosome, start index, and end index. To simulate
856 an SV, left and right breakpoints are first introduced to the chromosome, breaking up the chromo-
857 some into segments. This results in each SV breakpoint being on the exact border of a segment.
858 Then, the corresponding rearrangement is applied to the segments to simulate the intended SV.

859 For the parameterized-random SV selection, SVs are selected one at a time, for the number of
860 SVs indicated on the parameter file as follows:

- 861 1. The SV type is randomly selected based on its likelihood.
- 862 2. The event size is selected from a uniform distribution between the maximum and minimum
863 size indicated for the SV type.
- 864 3. The left breakpoint of the event location is selected uniformly among all chromosomes.
- 865 4. The right breakpoint is calculated based on the size of the event selected earlier. For a
866 balanced reciprocal translocation, it is selected similar to the left breakpoint.
- 867 5. The breakpoints are applied to partition the segments, and if a masked region file or a smallest
868 segment allowance is applied, the resulting breakpoints and segments are checked for legality.
869 If the result is illegal under the parameters, steps three through five are recomputed.

87 **S5.3 Generation and Processing of Simulated OGM Data**

871 The FASTA files generated by the KarSim module were processed by OMSim(Miclotte et al., 2017)
872 to simulate a BNX file containing OGM molecules with added noise. Parameters for the enzyme,

873 OGM generation, and noise were sourced from publicly available repositories (RaeisiDehkordi, 2024).
874 The Bionano Solve pipeline (v3.7) (BionanoGenomics, 2018) was then used to compute CNVs, SV
875 calls, and contig alignments, and these outputs were used as input for OMKar to reconstruct the
876 final virtual karyotype.

877 When simulating structural variations, we used the masking files provided in Bionano Solve v3.7
878 (also included in repository (RaeisiDehkordi, 2024)). The masking file is on reference hg38, with
879 size of 423.1 Mbp (13.65%), including 128.1 Mbp (4.13%) in centromeres and telomeres, and 64.5
880 Mbp (2.15%) in acrocentric chromosomes' p-arm. None of the SVs was simulated with breakpoint
881 boundary within 200 kbp of any masking region. Events were simulated with size 50 kbp to 2 Mbp.

882 **S6 KarCheck**

883 KarCheck is designed to be a symmetric comparator between two unphased karyotypes, denoted
884 by K_r and K_t , by comparing their SV and CNV calls. In addition, to accommodate molecular
885 techniques that do not have a nucleotide-level resolution, KarCheck has an adjustable tolerance for
886 small breakpoint mis-matching for an SV that is present in both karyotypes. KarCheck is publicly
887 available at <https://github.com/MolecularKaryotype/KarComparator>.

888 **S6.1 KarCheck Preprocessing**

889 Preprocessing is first applied to K_t and K_r to partition chromosome groups into *chromosome*
890 *clusters*. Two chromosome groups are linked into the same cluster when there exists a breakpoint
891 connecting them (signaling an inter-chromosomal SV). A maximal connected component of linked
892 chromosome groups is denoted as a chromosome cluster.

893 Recall that we define each karyotype as an ordered list of segments. To make the two karyotypes
894 comparable, we further partition their segments such that the two karyotypes share an identical
895 set of segments (Supplementary Fig. S8B). To achieve this, all left/right endpoints of segments are
896 collected, and if a segment has an endpoint internal to it, it is split into two, ensuring the left/right
897 endpoint are on the boundaries (Supplementary Fig. S8B).

898 **S6.2 SV similarity computation.**

899 After pre-processing, the SVs between the two karyotypes are compared as a directed-multi-graph,
9 which offers ease of checking the orientation of each SV edge (Supplementary Fig. S8C). On this
9 1 graph, nodes represent either the start or the end of a segment (denoted as s and t), with the
9 2 addition of a source node (S) and a sink node (T) for transition at the start and end of each
9 3 chromosome. Edges represent the presence of a transition between two segments. For example, a
9 4 chromosome of $[k+, l-]$, where the WT is $[k+, l+]$, will be represented by edges $(S, k_s), (k_s, l_t)$, and
9 5 (l_s, T) .

9 6 Since the two karyotypes share the same set of segments, they also share the same set of
9 7 nodes. Therefore, the comparison between the two sets of edges on this graph is equivalent to the
9 8 comparison of the two sets of SVs. To allow tolerance in breakpoint matching, we define a linear
9 9 distance function. Denote two non-segment edges as (a, b, o) and (c, d, o') , where o (and o') describe
9 10 orientation. Define distance $D((a, b, o), (c, d, o'))$ as:

$$D((a, b, o), (c, d, o')) = \begin{cases} \infty & \text{Chr}(a) \neq \text{Chr}(c) \text{ or } \text{Chr}(b) \neq \text{Chr}(d) \\ \infty & o \neq o' \\ |\text{pos}(a) - \text{pos}(c)| + |\text{pos}(b) - \text{pos}(d)| & \text{otherwise} \end{cases}$$

911 Prior to distance computation, identical edges between the two graphs are pruned. Second, if
912 a transition is s-to-t or t-to-s and is between two segments from the same chromosome with a
913 small distance (≤ 5 kbp), it is pruned. This is justified by that these transitions represent small
914 indels without change in orientation, which are not responsible if the technique has a minimum
915 resolution threshold. Finally, minimum weight bipartite matching is performed for the remaining
916 transition edges between the two karyotypes, with a maximum allowed matching distance 200 kbp.
917 Matching pairs are pruned, because they are considered similar SV edges. The matching distances
918 are recorded for downstream analyses upon resolution. Finally, all residual non-segment edges are
919 the differential SV edges between the two karyotypes. We use K_r to denote the reconstructed
920 chromosome, and K_t as the WT or true chromosome. Therefore, residual edges in K_t represent
921 false negatives, and residual edges in K_r represent false positives.

922 **SV Similarity Metrics.** The similarity of each individual SV edge is compared via the count
923 of K_t edges after the two initial pruning procedures, and the residual edges in K_t and K_r after
924 final matching. $TP = |\text{initial edge}| - |K_t \text{ residual}|$, $FN = |K_t \text{ residual}|$, and $FP = |K_r \text{ residual}|$. A
925 Jaccard Similarity is computed for each cluster of compared simulated data to penalize both false
926 positive and false negative events. For comparison against real data, only the recall is computed
927 for each cluster, as the cytogenetic methods employed for the reference calls do not necessarily have
928 small enough resolution or ability to catch de-novo balanced events.

929 **S6.3 Copy Number similarity comparison and metrics**

930 CN comparison is done by binning the whole genome (excluding prefix/suffix masked region) into
931 spanning, non-intersecting bins of 50 kbp +/- 100 bp (exact size chosen to maximize the size of
932 the last bin on the chromosome). Each bin is used to store the average CN within that region,
933 separately for the K_t and the K_r .

934 Then, for each cluster, the chromosome groups are determined by the union of all the chromo-
935 somal origins of the segments in the cluster. Each cluster's CN bins are the subset of the total CN
936 bins, to only include the chromosome groups within the cluster. A WT expected count is deter-
937 mined to each chromosome group by rounding the average bin CN from K_t (diploid for autosomes
938 and XX or XY for sex chromosomes).

939 The values of the CN bins' of K_t and K_r are computed with CNs from all corresponding
940 segments. If a bin has a gain or loss of more than 0.05 CN from the WT expected count, it is a
941 marked as "CN gain" or "CN loss". Otherwise, it is marked as "CN neutral". This forms a paired
942 CNV array where the Jaccard Similarity can be computed. The denominator of the similarity is
943 the count of bins marked as CN gain or loss in either K_t or K_r , and the numerator is the count of
944 bins where K_t and K_r agree on CN gain or loss.

945 **S6.4 Addtional Functionalities for Downstream Analyses**

946 Edges on both K_t and K_r can be labeled with additional input. For example, Table 2 was computed
947 where each edge in K_t from the simulation was labeled with the Structural Variation event type,
948 so a summary statistics was generated from the residual edge count of each event type.

949 When matching transition edges, the distances between the matched edges can be collected for

95 analysis. This can be further categorized by having K_t 's transition edge labeled with the causal
96 event type. A detailed analysis on OMKar's reconstructions' distances against the simulation can
97 be found in Supplementary Section S4.

98 **S6.5 Usage in validating real data reconstructions with previous cytogenetic
99 records**

100 A function was implemented in KarCheck to allow efficient input of a karyotype for the purpose
101 of validating a reconstruction against previous cytogenetic records on that karyotype. For each
102 karyotype, its aneuploidy (if any), each event's induced SV-breakpoints, and CN changes are taken
103 as input for the validation. This information is sufficient to populate a full molecular karyotype,
104 while KarCheck assumes the rest of the genome is WT.

105 For real data with cytogenetic records, the cytogenetic records were used as the “truth” karyo-
106 type (K_t) and compared against the reconstruction (K_r). CN gains called with microarray do not
107 specify the structural breakpoints of the amplification. For these, we first assumed each amplifica-
108 tion was a tandem duplication, and if the matching failed, we manually verified if the reconstruction
109 contained the amplification as a segmental duplication. Additionally, previous cytogenetic tech-
110 niques using microarray and staining did not fully capture every event on the genome, so the K_t
111 was treated as incomplete, with potentially missed TP events. Therefore, for our final statistics of
112 the comparison, we only computed the recall to verify if OMKar reconstruction included all the
113 previously identified TPs.