

Human copy number variants are enriched in regions of low-mappability

Jean Monlong^{1,2}, Patrick Cossette³, Caroline Meloche³, Guy Rouleau⁴, Simon L. Girard^{1,5},
and Guillaume Bourque^{1,2,6}

¹Department of Human Genetics, McGill University, Montréal, H3A 1B1, Canada

²McGill University and Génome Québec Innovation Center, Montréal, H3A 1A4, Canada

³Centre de Recherche du Centre Hospitalier de l'Université de Montréal, Montreal, H2X 0A9, Québec, Canada.

⁴Montreal Neurological Institute, McGill University, Montréal, H3A 2B4, Québec, Canada.

⁵Département des sciences fondamentales, Université du Québec à Chicoutimi, Chicoutimi, G7H 2B1, Canada

⁶Correspondence: guil.bourque@mcgill.ca

April 27, 2016

Abstract

Germline copy number variants (CNVs) are known to affect a large portion of the human genome and have been implicated in many diseases. Although whole-genome sequencing can help identify CNVs, existing analytical methods suffer from limited sensitivity and specificity. Here we show that this is in large part due to the non-uniformity of read coverage, even after intra-sample normalization, and that this is exacerbated in regions of low-mappability. To improve on this, we propose **PopSV**, an analytical method that uses multiple samples to control for technical variation and enables the robust detection of CNVs. We show that **PopSV** is able to detect up to 2.7 times more variants compared to previous methods, with an accuracy of about 90%. Applying **PopSV** to 640 normal and cancer whole-genome datasets, we demonstrate that CNVs affect on average 7.4 million DNA bases in each individual, a 23% increase over previous estimates. Notably, we find that regions of low-mappability are approximately 8 times more likely to harbor CNVs than the rest of the genome, which contrasts with somatic CNVs that are nearly uniformly distributed. In addition to the known enrichment in segmental duplication, we also observe that CNVs are enriched near centromeres and telomeres, in specific types of satellite and short tandem repeats, and in some of the most recent families of transposable elements. Although CNVs are found to be depleted in protein-coding genes, we identify 7206 genes with at least one exonic CNV, 682 of which harbored CNVs in low-mappability regions that would have been missed by other methods. Our results provide the most exhaustive map of CNVs across the human genome to date and demonstrate the broad functional impact of this type of genetic variation including in regions of low-mappability.

1 Introduction

Structural variants (SVs) are defined as genetic mutations affecting more than 100 base pairs and encompasses several types of rearrangements: deletion, duplication, novel insertion, inversion and translocation¹. Deletions and duplications, which affect DNA copy number, are also collectively known as copy number variants (CNVs). SVs arise from a broad range of mechanisms and show a heterogeneous distribution of location and size across the genome^{1,2,3}. In healthy individual, SVs are estimated to cumulatively affect a higher proportion of the genome as compared to single nucleotide polymorphisms (SNPs)⁴. Numerous diseases including Crohn's Disease⁵, schizophrenia⁶, obesity⁷, epilepsy⁸, autism⁹, cancer¹⁰ and other inherited diseases^{11,12}, harbor SVs with a demonstrated detrimental effect^{13,14,15}.

While cytogenetic approaches and array-based technologies have been used to identify large SVs, whole-genome sequencing (WGS) could in theory uncover the full range of SVs both in terms of size and type¹⁶. Numerous methods have been implemented to detect SVs from WGS data using either paired-end information^{17,18}, read-depth (RD) variation^{19,20,21}, breakpoints detection through split-read approach²² or de novo assembly²³. However, the presence of technical bias in WGS is an important challenge. Indeed, it has been shown that various features of sequencing experiments, such as mappability^{24,25}, GC content²⁶ or replication timing²⁷, have a negative impact on the uniformity of the coverage²⁸. Unfortunately, this variability is difficult to fully correct for as it involves different factors, some of which are unknown, that vary from one experiment to another. This issue particularly impairs the detection of SV with weaker signal, which is inevitable in regions of low-mappability, for smaller SVs or in cancer samples with stromal contamination or cell heterogeneity.

As a result, existing approaches suffer from limited sensitivity and specificity^{3,16}, especially in specific regions of the genome, including regions of low-complexity and low-mappability^{24,25}. Some methods^{29,30} try to model ambiguous mapping and repeat structure but address only specific situations. Another strategy to improve the accuracy of SV detection has been to use an ensemble approach that combines information from different methods relying on different types of reads. Large re-sequencing projects such as the 1000 Genome Project^{3,31} and the Genomes of Netherlands (GoNL) project^{32,33} have adopted this strategy and have successfully identified many SVs using an extensive panel of detection methods combined with low-throughput validation. Such a strategy increases the specificity of the calls but is less sensitive. In these studies, as in many others, repeat-rich regions and other problematic regions are frequently removed or smoothed at some step of the analysis, to improve the accuracy of the calls. Thus, low-mappability regions are just scarcely covered in some of the most recent CNV catalogs³¹. This is unfortunate given that CNVs in such regions have already been associated with various diseases^{34,35,36,12} and that these regions are also more likely variable. Indeed, CNVs are known to be enriched in segmental duplications², short and long tandem repeats can be highly polymorphic^{37,38} and SV formation can be facilitated by repeat templates³⁹.

In this work, we start by showing that technical variation challenges the uniformity of coverage assumption despite state-of-the-art intra-sample normalization. To correct for this, we propose a new method, PopSV, an approach that relies on RD but uses a set of reference samples to control for technical variation and detect abnormal read coverage. Our approach differs from previous RD methods, such as RDXplorer⁴⁰ or CNVnator²⁰, that scan the genome horizontally and look for regions that diverge from the expected global average. Even when approaches rely on a ratio between an aberrant sample and a control, such as FREEC¹⁹ or BIC-seq⁴¹, we show that they do

not sufficiently control for experiment-specific noise as compared to PopSV. Glusman et al.⁴² go further by normalizing the RD with pre-computed RD profiles that fit the GC-fingerprint of a sample but this approach excludes regions with extreme RD and doesn't integrate the variance in the pre-computed RD profiles which is essential to robustly deal with mappability bias. PopSV is also different from approaches such as cn.MOPS²¹ and Genome STRIP⁴³ that scan simultaneously the genome of several samples and fit a Bayesian or Gaussian mixture model in each region. Those methods have more power to detect SVs present in several samples but may miss sample-specific events. Moreover, their basic normalization of coverage and fully parametric models forces them to conceal a sizable portion of the genome and variants with weaker signal.

To demonstrate the utility of PopSV in characterizing CNVs across the genome, we apply the method to 640 WGS individuals from three human cohorts: a twin study with 45 individuals⁴⁴, a renal cell carcinoma datasets with 95 tumor and control pairs⁴⁵ and 500 unrelated individuals from the GoNL dataset³². Using this data we compare the performance of PopSV with existing CNV detection methods and validate the quality of the predictions. We also characterize the patterns of CNVs across the human genome and show that CNVs are enriched in regions of low-mappability and in different classes of repeats. Finally, we look at the functional significance of these structural variants and show that CNVs overlap thousands of genes, with hundreds of them hit by CNVs associated with regions of low-mappability.

2 Results

Intra-sample normalization does not remove coverage biases It is usually assumed that after correction for known biases such as GC content²⁶ and mappability^{24,25}, sequencing reads in a WGS experiment are uniformly distributed across the genome. To test this hypothesis, we first filtered and normalized the RD in the normal samples of the renal cancer dataset, following standard techniques. We computed the RD as the number of properly mapped read in non-overlapping genomic windows (bins) of size 5 kilo bases (Kb). Read counts in the bins were corrected for GC-bias and, to be conservative in this initial analysis, regions with extreme read coverage were removed (Methods). Bin counts were then quantile normalized to obtain the same distribution for all samples (Fig. S1). Unexpectedly, and in contrast to simulated datasets, the inter-sample mean coverage in each bin was observed to vary from one genomic region to the other, highlighting the presence of additional biases (Fig. 1a). Supporting this observation, the bin coverage variance across samples was lower than expected and also varied between genomic regions (Fig. S2). Such region-specific bias is overlooked when global estimates and genome-scanning methods are used to detect coverage differences. To further investigate this bias, we computed the proportion of the genome where a given sample had either the highest or the lowest coverage of all samples. Some samples looked more affected by this bias than others, as they consistently showed the highest, or the lowest, coverage across large portions of the genome (Fig. 1b). Similar patterns were also observed in the other two cohorts (Methods and Fig. S3 and S4). In short, we observed significant coverage biases even after intra-sample normalization and when focusing on the least problematic regions of the genome. Not surprisingly, this effect was even stronger when the whole genome was being evaluated (Fig. S5). This artificial variation has implications for CNV detection approaches that assume a uniform distribution of the RD technical variation across the genome or across samples^{19,20,41}. First, the rate of false positives will be higher as the coverage will artificially fluctuate. Moreover, this experimental noise will confuse the detection of weaker signal, e.g. in

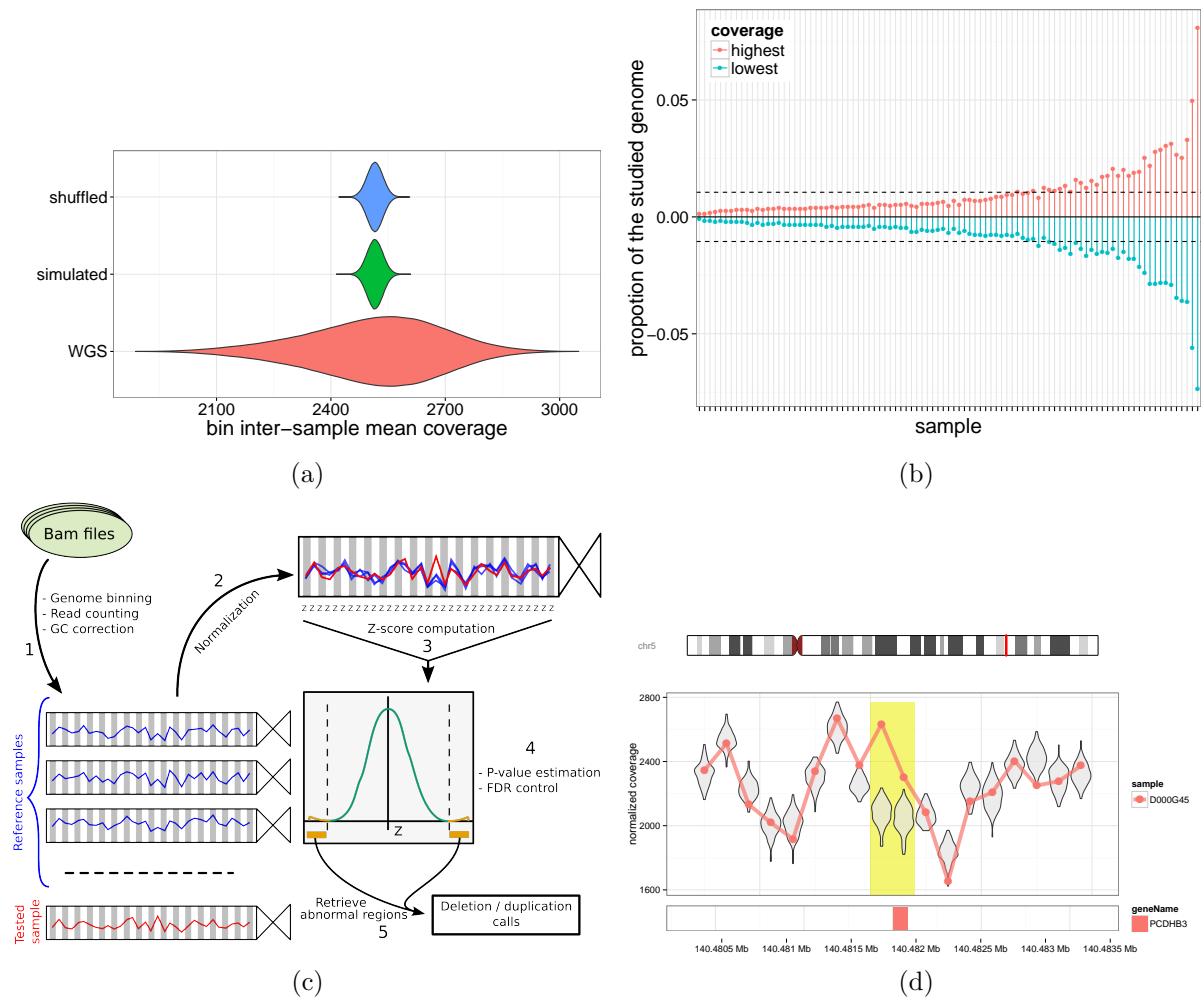


Figure 1: **Coverage bias in whole-genome sequencing and the PopSV approach.** a) Distribution of the bin inter-sample mean coverage (red) and null distribution (blue: bins shuffled, green: simulated normal distribution). b) Proportion of the genome in which a given sample (x-axis) has the highest (red) or lowest (blue) RD. In the absence of bias all samples should be the most extreme at the same frequency (dotted horizontal line). c) PopSV approach. First the genome is fragmented and reads mapping in each bin are counted for each sample and GC corrected (1). Next, coverage of the sample is normalized (2) and each bin is tested by computing a Z-score (3), estimating p-values (4) and identifying abnormal regions (5). d) The line and points represent the coverage of one sample with a duplication (highlighted in yellow); the violin plots represent the distribution of the coverage in the reference samples.

low-mappability regions, for smaller CNVs or in cancer samples with stromal contamination or cell heterogeneity.

A population-based normalization and CNV detection method The main idea behind PopSV is to assess whether the coverage observed in a given location of the genome diverges significantly from the coverage observed in a set of reference samples. In PopSV, the genome is first segmented into bins and RD is computed for each sample as the number of reads with proper mapping in each bin. In a typical design, the genome is segmented in non-overlapping consecutive windows of equal size, but custom designs could also be used. After normalization, the value observed in each bin is compared to the values observed in the reference samples and a Z-score is calculated (Fig. 1c, 1d and Methods). False Discovery Rate (FDR) is estimated based on these Z-score distributions and a bin is marked as abnormal based on a user-defined FDR threshold. Consecutive or nearby abnormal bins are merged and considered as one variant. Other segmentation approaches, such as the circular binary segmentation can also be used. The normalization step is critical here since we have shown that simple approaches will fail to give acceptable normalized RD scores (Fig. 1b). Moreover, with global median/variance adjustment or quantile normalization, the remaining subtle experimental variation impairs the abnormal RD test (Fig. S6a). With PopSV, we propose a new normalization procedure, which we call targeted normalization, that retrieves, for each bin, other genomic regions with similar profile across the reference samples and uses these regions to normalize RD (Methods). In contrast to other methods, targeted normalization shows better distribution features (Fig. S6b). It is important to note that it is critical for the success of this targeted normalization that the set of reference samples used is comparable to the tested samples. We have included in PopSV a set of exploratory tools to help assess this (Methods).

Sensitivity and specificity of PopSV To demonstrate the effectiveness of PopSV, we first applied it to the twin dataset (Methods). Using 5 Kb bins, we observe smooth normal-like Z-score distributions and overall consistency of the bin values in the twin pairs (Fig. 2a). Applying the same methodology to the normal/tumor cancer cohort lead to similar results and highlighted, as expected, a large number of duplications and deletions in the tumors (Fig. S7). Encouragingly, in regions of low-mappability, the Z-score distribution was found to be identical to the one in regions of normal mappability (Fig. S8). Next, we estimated the copy number of each bin by dividing the RD in a given sample by the average RD across the reference samples multiplied by two, to reflect the fact that reference set is assumed to be diploid in each bin. We anticipate the copy number estimate to be reliable if the detected event spans the entire bin but less accurate for smaller event or partial signal (e.g. contamination or cell heterogeneity in cancer). The distribution of these copy-number estimates further supported the quality of the PopSV calls, with clear peaks around integer values (Fig. 2b). It's important to note that, in contrast to some of the other methods^{21,43}, this aggregation around integer values is completely independent of the calling process which only marked bins with abnormal RD.

To evaluate the performance of PopSV, we compared it to FREEC¹⁹ and cn.MOPS²¹, two popular RD methods that can be applied to WGS datasets to identify CNVs. FREEC segments the RD values of a sample using a LASSO-based algorithm while cn.MOPS considers simultaneously several samples and detects copy number variation using a Poisson model and a Bayesian approach. First, in the twin study, we measured the number of CNVs identified in each twin that were also found in the matching twin (Methods). High frequency CNVs were removed to ensure that systematic

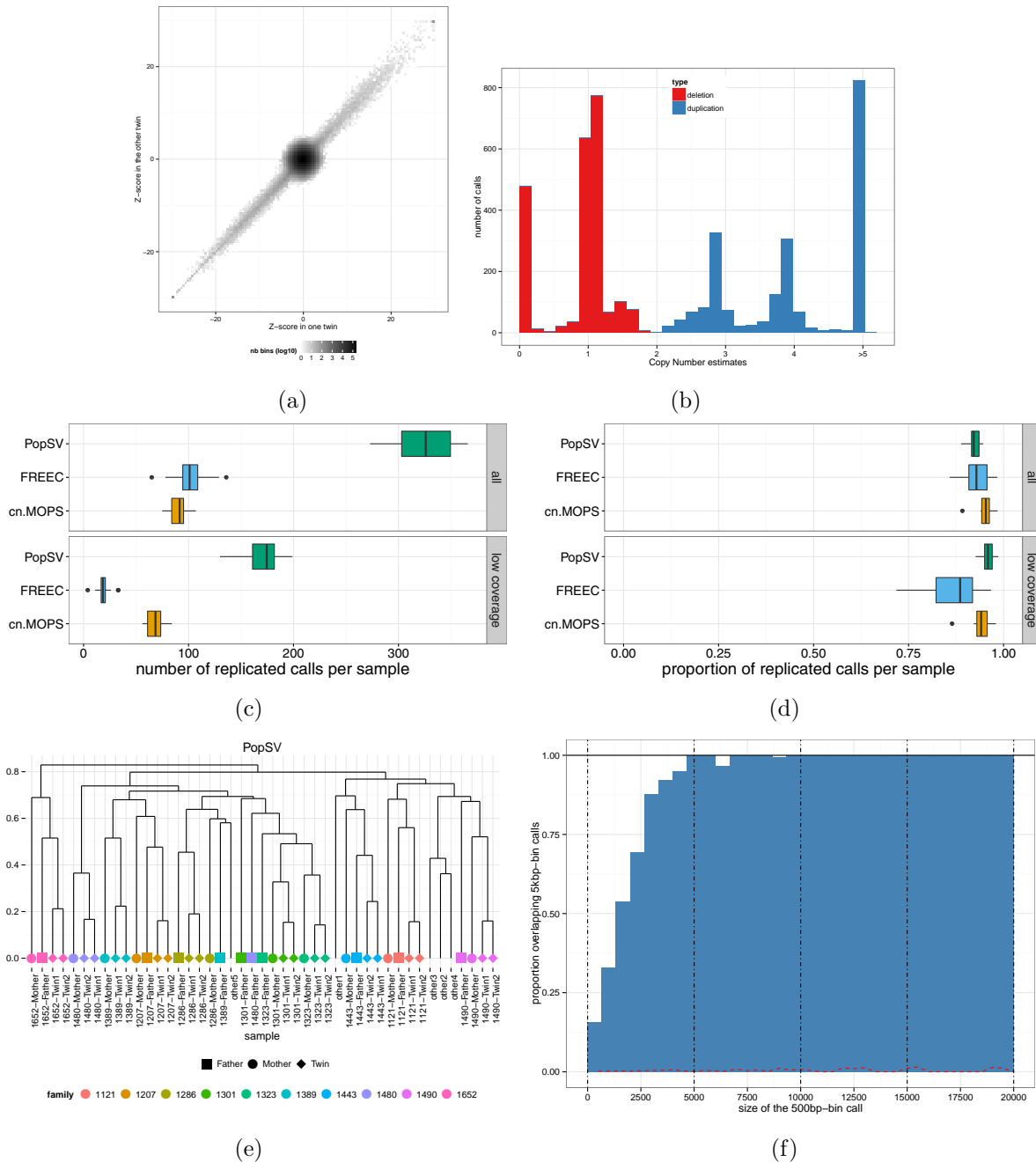


Figure 2: **Sensitivity, specificity and resolution of PopSV.** a) Z-scores for the genomic bins from two twins (x - y - axis). Non-zero positive or negative Z-score supports a duplication or a deletion, respectively. b) Focusing on large events, copy numbers can be estimated accurately and segregate close to integer values. c-d) Number (c) and proportion (d) of variants from a twin that was replicated in the other twin. e) Samples are clustered using PopSV calls in regions of extremely low coverage and recover almost perfectly the family structure. f) Proportion of 500 bp calls of different sizes (x -axis) overlapping a 5 Kbp call.

errors were not biasing our replication estimates. Using 5 Kb bins, PopSV recovered on average more replicated CNV events per sample, 324 versus 102 and 92 for FREEC and cn.MOPS respectively, while maintaining comparable specificity (Fig. 2c and 2d). Notably, focusing on the regions of low mappability, we found that PopSV also outperforms the other approaches with 174 replicated events per sample on average, while cn.MOPS and FREEC only found 70 and 19 respectively. In those regions, PopSV had a slightly higher specificity with 96% of the calls being replicated (Fig. 2d). To explore the quality of the CNV calls further, we clustered individuals according to the CNV calls and compared the result to the known pedigree for these samples (Methods). We found that PopSV shows better concordance as assessed by the Rand index (Fig. S9). Even using only the regions of extremely low coverage resulted in a clustering dendrogram mimicking almost perfectly the family relationships (Fig. 2e). Additionally, the distribution of CNV recurrence shows a clearer peak at 3-sample for PopSV (Fig. S10), which is expected due to the aggregation of CNVs present in both twins and one parent.

To further assess the performance of PopSV, we also tested the approach on the cancer dataset by comparing the agreement between germline events in tumor/normal pairs in a similar way as was done for the twin pairs. We observed comparable results with PopSV reporting on average 293 replicated CNV calls per sample while cn.MOPS and FREEC only detected 75 and 48 such events respectively (Table S1). Once again, the specificity of the different methods was comparable at around 88%. This was true overall as well as in low mappability regions where PopSV found twice as many replicated calls.

Resolution and validation of the PopSV calls To evaluate the performance of PopSV at different resolutions, we repeated the analysis of the twin dataset using 500 bp bins. With smaller bins there is more noise and long stretches of bins of low significance might be missed. For this reason, the 500 bp calls were combined with the 5 Kb calls (Methods). At this resolution, we observed that PopSV still found on average 1.7 and 6.3 times more replicated calls per sample compared to cn.MOPS and FREEC while maintaining similar specificity (Table S1). PopSV also detected on average 1.3 and 23.2 times more replicated variants in regions of low mappability compared to cn.MOPS and FREEC respectively, and had the highest specificity of all tested methods. Similar results were observed with 500 bp bins in the renal cancer data set (Table S1). We also compared the 5 Kb and the 500 bp individual calls and observed high consistency (see Methods). Remarkably, the results suggest that PopSV could detect 75% of the events as large as half the bin size (Fig. 2f). Additionally, we assessed the performance in each genomic bin individually (Methods). This analysis showed that PopSV is reliable across the range of mappability, GC and repeat content (Fig. S11), and, overall, across a larger fraction of the genome compared to FREEC and cn.MOPS (Table S2).

Finally, some variants were experimentally validated. First, we randomly selected 20 one-copy and two-copy deletions, among small (~ 700 bp) and large (~ 4 Kb) variants and visually inspected them to design PCR primers (Methods). In total, 18 out of 20 (90%) were successfully validated, close to our *in silico* estimates (Table S3). Next, we designed PCR primers for the validation of rare deletions in low-mappability regions. For these regions we had to perform local read re-assembly in order to better predict potential breakpoints (Methods). In this case, 11 out of the 18 (61%) deletions in low-coverage regions were successfully validated (Table S4). We note that designing primers in repeat-rich regions is challenging and that this might lead to an under-estimation of the true validation rates. Additionally, we observed that the majority of the non-validated CNVs

were predicted to be deletion smaller than 100 bp. If we focus on deletions larger than 100 bp, the validation rate in regions of low-mappability was increased to 77% (10/13) once again very close to our *in silico* estimates.

Global patterns of CNVs across the human genome Having demonstrated the sensitivity, specificity and resolution of PopSV, we wanted to characterize the global patterns of CNVs across the human genome. We started with an analysis of the twins and the normal samples in the renal cancer dataset, both of which have an average sequencing depth around 40X. We will be especially interested in looking at calls in regions of low-mappability which represents 12.6% of the human genome (Fig. S12 and Methods).

PopSV was used to make calls using 500 bp and 5 Kb bins, which were then merged to create a final set of variants as before. On average per genome, 7.4 Mb of the reference genome had abnormal read coverage, 4 Mb showing an excess of reads indicating duplications and 3.4 Mb showing a lack of reads indicating deletions (Table 1). In both datasets, the average variant size was around 4 Kb and 70% of the variants found were smaller than 3 Kb. We compared our numbers to equivalent CNVs detected in the recent human SV catalog from the 1000 Genomes Project³¹ (Methods). In that study, we calculated that 6.0 Mb of the reference genome was found to be variable on average in each genome (Table S5). We also notice that no variants except for a few deletions were identified in low mappability regions in this catalog. Similarly, small duplications (< 3 Kbp) are absent from this catalog. In contrast, the set of variants identified by PopSV included variants in low mappability regions as well as small deletions and duplications (Table 1), explaining in part the ~ 23% increase. While the study from the 1000 Genomes Project³¹ explored a wider range of SVs, our set of variant is likely more representative of the distribution of CNVs in a normal genome since a broader portion of the genome could be analyzed.

Next, we applied PopSV to the 500 unrelated samples from the GoNL cohort (Table 1). Due to a lower sequencing depth (~13X), we used bins of size 2 Kb and 5Kb that gave the best signal to noise ratio (Methods). Slightly fewer variants were found in these samples mainly because of the reduced sequencing depth, which limits the detection of smaller CNVs. Nevertheless, a large sample size helps better characterize the frequency patterns and provides a more comprehensive map of rare CNVs. In total, across these three cohorts, 326 Mb were found to be affected by a CNV with more duplications (325,602) detected than deletions (248,937). This contrasts with the CNVs reported by the 1000 Genomes Projects³¹ that were heavily skewed towards deletions (Table 1 and Table S5), likely due to the usage of different methods to detect various types of CNVs. The frequency distribution of deletions and duplications found using PopSV was also much more balanced compared with the ones from Sudmant et al.³¹ (Fig. S13). Of note, we observed the same when comparing PopSV with other methods: PopSV's frequencies are more similar between deletions and duplications compared to FREEC (Fig. S14). As expected, both deletions and duplications detected by cn.MOPS tend to be skewed towards more common events.

CNVs are enriched near centromeres and telomeres and in regions of low-mappability

Large CNVs have been shown to be enriched near centromeres, telomeres and assembly gaps (CTGs)⁴⁶. We were interested in exploring this observation further using the set of high resolution calls from PopSV. We compared the distribution of CNVs calls made across the 3 datasets to randomly distributed regions of similar sizes (Fig. S15 and Methods). In an average genome, we found that 33% of the CNVs calls were within 1 Mb of a CTG, while we would have expected 11%

by chance. To verify that these observations were not simply a consequence of the methodology used, we also looked at the somatic CNVs (sCNVs) that we could detect in the renal dataset. For this purpose, we extracted the variants found by PopSV in the tumor sample of an individual but missing from its paired normal sample (Methods). As expected, somatic CNVs were found to be significantly larger and to affect a much larger fraction of the genome (Table S6). Reassuringly, and in contrast to germline CNVs, sCNVs were not preferentially found near CTGs (Fig. S15), with only 14% of the sCNVs within 1 Mb of a CTG.

Notably, when looking at the genomic distribution of CNVs, we also observed a 8.2 fold-enrichment of variants in regions of low mappability (Fig. 3a). Segmental duplications (SD), DNA

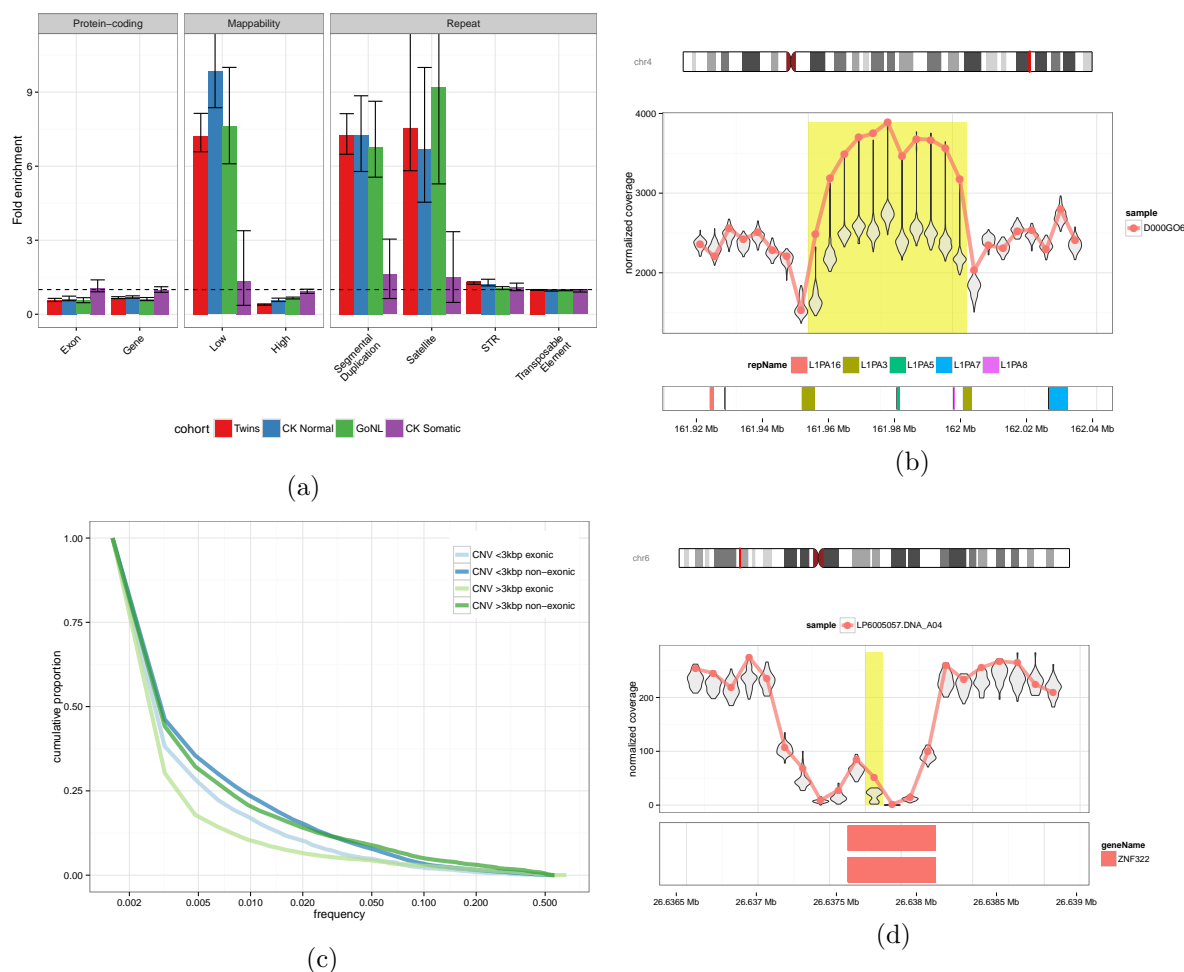


Figure 3: CNVs in normal genomes. a) Enrichment of CNVs in different genomic classes (x-axis) across different cohorts (colors). Bars show the fold enrichment compared to control regions. The error bars represent 80% of the samples. b) Example of CNV likely caused by non-allelic homologous recombination between two L1PA3 repeats. The line and points represent the coverage of one sample with a duplication (highlighted in yellow); the violin plots represent the distribution of the coverage in the reference samples. c) Frequency of CNVs of different sizes and overlap with coding exons. d) Example of a duplication in *ZNF322* exon, located in a challenging region (low coverage).

satellites and Short Tandem Repeats (STR) were also significantly enriched with fold-enrichment of 7.1, 7.8 and 1.2, respectively. The over-representation of CNVs in SDs has been described before² and in a recent study⁴⁷, half of the CNV base pairs were shown to overlap a SD. To investigate the contribution of low-mappability regions beyond SDs, random control regions were selected to have the same size distribution as the real CNVs, to have comparable overlap with SDs and to have similar distance to CTG (Methods). Even after controlling for these known enrichments, we found that CNVs overlapped low-coverage regions twice as much as expected (Fig. S16a). This two-fold enrichment is independent of the SD association and consistently observed in the 3 cohorts of normal genomes. In contrast to germline CNVs, sCNVs were once again found to be more uniformly distributed (Fig. 3a). These results suggest that the enrichments of germline CNVs near CTGs and in regions of low-mappability are trustworthy and are unlikely to be the result of a methodological artifact.

Various repeat families are more prone to harbor CNVs We wanted to characterize further the distribution of germline CNVs in relation to different genomic features, including looking at the contribution of different repeat classes. By comparing CNVs to the same control regions with matched overlap with SD and distance to CTGs we can look for patterns that are specific to repeat sub-families without the risk of being biased by the global enrichments. Using this approach, we found that CNVs were still significantly enriched in satellites repeats and in short tandem repeats (STRs) (P-value < 10⁻¹³, Fig. S16a), with fold-enrichments of 2.3 and 1.1 respectively. In contrast, and as expected, protein-coding genes and exons were found to be significantly under-represented.

Although it is known that satellites and simple repeats DNA are more unstable⁴⁸, the extent to which CNVs are found in these regions in humans had, to our knowledge, not been systematically explored. Satellite repeats are grouped into distinct families depending on their repeated unit and we found that not all satellite repeats were equally likely to overlap a CNV (Fig. S16b). In particular, Alpha satellites have the highest and most significant enrichment (P-value 5.8×10⁻¹⁴), with more than 3.5 times more CNVs than in the control regions. Satellite-like Repeat 1 (SATR1) was also enriched in CNVs. We noted that satellites tend to span completely CNVs (Fig. S17), suggesting that satellites are likely directly involved, in contrast to having a passive role, in the CNV formation. Short and long tandem repeats can be highly polymorphic^{37,38}. Constrained by read length, recent studies^{49,50} focused on variation of STRs smaller than 100 bp. In our analysis we found that CNVs were significantly enriched in the largest annotated STRs (>100 bp or >400 bp, Fig. S18). STR can be grouped by motif and we further tested the largest and most frequent families (Fig. S16c). Except for *TA(/AT)* tandem repeats, we conclude that the enrichment is driven by other STR families, most likely without specific involvement of sub-families. Here the repeats tend to overlap just a fraction of the variant, but a clear subset of the variants are fully covered by these tandem repeats (Fig. S17).

Finally, although transposable elements (TEs) as a whole did not show enrichment (Fig. 3a), the Other repeat class, which corresponds to SVA repeats, was found to be significantly enriched in the three datasets (Fig. S16d). Moreover, looking at TEs at the level of individual repeat families, we found a number of them to be enriched including SVA D-F or L1Hs. Surprisingly, a few older ERV families, including HERV-H that has been shown to be expressed and important in human embryonic stem cells^{51,52}, were also in the list of enriched TEs. Several families of older L1 repeats (e.g. L1PA2 to L1PA5) were also enriched and often implicated in what appears to be non-allelic homologous recombination (see examples in Fig. 3b and S19). Reassuringly, the somatic CNVs

once again did not show any of these enrichments (Fig. S16d).

Impact of CNVs in low-mappability regions on protein-coding genes Although both small and large CNVs were depleted in genes (Fig. 3c and Methods), 7206 protein-coding genes were found to have an exon overlapping an event in at least one of the 640 normal genomes studied (Table 2). Moreover, if we included the promoter regions (10 Kbp upstream of the transcription

| Set | CNVs | Genes with CNVs | | |
|------------------------|--------|-----------------|------------|----------|
| | | Exon | + Promoter | + Intron |
| Germline CNVs | | | | |
| All CNVs | 91733 | 7206 | 11341 | 13259 |
| Low coverage | 26888 | 682 | 1151 | 1977 |
| Extremely low coverage | 10010 | 347 | 465 | 521 |
| STR | 4286 | 45 | 286 | 748 |
| Satellite | 1822 | 2 | 21 | 33 |
| TE | 20491 | 164 | 1747 | 3998 |
| STR/Sat/TE | 22313 | 166 | 1760 | 4014 |
| Somatic CNVs | | | | |
| All CNVs | 331778 | 18121 | 18909 | 18969 |
| Low coverage | 9508 | 450 | 836 | 1384 |
| Extremely low coverage | 2476 | 173 | 295 | 328 |
| STR | 2829 | 39 | 267 | 662 |
| Satellite | 1423 | 2 | 24 | 35 |
| TE | 45205 | 400 | 3820 | 7137 |
| STR/Sat/TE | 46630 | 402 | 3839 | 7154 |

Table 2: Impact of CNVs on protein-coding genes. The *CNVs* number represents the number of different CNVs, after collapsing CNVs with more than 50% reciprocal overlap. Repeat CNV: more than 90% of the CNV is annotated as repeat. Genes are protein-coding genes and the promoter region is defined as the 10 Kbp region upstream of the transcription start site.

start site), at least 11341 protein-coding genes were potentially affected by at least one CNV (Methods). Focusing on regions of low-mappability, we found 4286 different CNVs that were completely included in regions annotated as STR. These STR-CNVs overlapped exonic regions of 45 protein-coding genes, and 286 genes when including the promoter region (Table 2). In contrast, for CNVs included in satellite regions, only 21 genes had an exon or the promoter region overlapping one of the 1822 Satellite-CNVs. Finally we focused on CNVs that completely spans regions of extremely low mappability (see Methods). Even there, 347 genes were found to have an exon overlapping such a CNVs and this number increased to 465 if we included the promoter regions. These CNVs are distinct from larger aberrations (see Fig. 3d for example) and could easily be missed by other approaches masking low mappability regions. Of note, we also found that 173 genes were affected by somatic exonic CNVs located within these extremely low coverage regions (Table 2).

3 Discussion

Why are SVs so difficult to detect in WGS data? We have answered this question by showing that the various experimental biases cannot be corrected for using basic intra-sample normalization and affect the uniformity of read coverage across the genome. It is important to note that the amplitude of these biases varied from one cohort to the next and did not appear to be strictly linked to the sequencing platform used but also to the way the samples were prepared (Fig. 1a, 1b, S3 and S4). With PopSV, samples that were sequenced with the same technology and protocols can now be analyzed jointly to control for these biases. When only a few samples are available this inter-sample normalization procedure might be less efficient but we estimate that with 20 reference samples or more PopSV will be preferable over methods working on single samples (or pairs of samples). We note that WGS is probably one of the most straightforward next-generation sequencing (NGS) protocol that only involves DNA extraction, shearing, sometimes amplification, and sequencing. It is likely that other NGS experiments, such as ChIP-Seq, are also similarly affected by sample preparation conditions and that these would also benefit from a similar inter-sample normalization procedure.

Comparing different calling methods is not straightforward, especially when different strategies are implemented. To begin, we compared PopSV, cn.MOPS and FREEC using the same large bin size (5 Kb) in order to assess their ability to detect different types of signal: full versus partial signal, single versus multiple bin support, normal versus low mappability. Next, we ran the methods with a smaller bin size (500 bp) to compare the methods in a situation with higher background noise. In each comparison we made sure that PopSV had similar specificity estimates compared to other methods, in order to reliably compare the sensitivity. We concluded that PopSV was more capable of detecting partial or single-bin signal (Fig. 2c and 2f), which is valuable to be able to observe smaller variants or variants in more challenging regions. Even when the background noise was significant, PopSV showed the best sensitivity and could reliably test a wider range of the genome (Table S1 and S2). In contrast to cn.MOPS, FREEC and ensemble methodologies^{32,31}, PopSV was also able to detect both deletions and duplications as efficiently (Fig. S13 and S14).

A notable strength of this new approach is that it enables the analysis of CNVs across the genome. Using PopSV on 140 normal genomes with high sequencing depth ($\sim 40X$) and 500 additional samples with medium coverage ($\sim 13X$), we found that regions of low mappability, which only represent 12.6% of the genome, overlap with 65% of the CNVs detected. The fact that this enrichment was observed for germline events and not somatic events was both reassuring and interesting because of the implications on the selection forces at play. Having a more complete CNV catalog also enabled an unbiased characterization of the CNV patterns across genome and potentially increases the power for trait-association studies. In particular, we were able for the first time to quantify the extent to which some regions in the genome are more prone to harbor such structural rearrangements. For example, we could describe genome-wide enrichment for different families of DNA satellites, simple repeats and several TE families, such as SVA, L1Hs and HERV-H. Although PopSV doesn't characterize fully STR variation it is able to detect CNVs in large STRs, something that cannot be done by STR detection methods using WGS. A recent analysis of an haploid human cell line⁵³ found a large number of novel SVs thanks to the use of long sequencing reads. Their study highlighted the variation involving complex repetitive DNA. Although the short reads in our study don't allow for a full characterization, we could detect the presence of such CNVs across a large population of normal genomes.

Because PopSV looks for abnormal read coverage in each bin independently, it does not require

the coverage to be uniform across the genome. For this reason, a natural extension of PopSV would be to apply it to targeted sequencing data, such as whole-exome sequencing data. In this context, the fragmented nature of the coverage and the differences in baseline from one region to another would seamlessly be integrated and corrected for by the set of samples used as a reference. Actually, several methods for CNV detection from whole-exome data that use information from other samples already exist^{54,55}, although they do not control for the biases described above the way PopSV does. Similarly, another logical extension of PopSV would be to apply it not only to correctly mapped reads but also to discordant reads to detect abnormal discordant coverage. Here, any type of discordant mapping, such as read pairs with incorrect insert size, orientation or with only one pair mapped could be counted together or separately. Discordant reads are intrinsically difficult to work with because they are usually ambiguous and found in regions of low-mappability. Issues of ambiguous mapping are context-specific and are exceedingly difficult to model directly. The advantage of working with a set of reference samples, as in the PopSV framework, is that we would have a way to control for this variability empirically. An additional advantage of incorporating the discordant reads in PopSV is that it would also allow for defining more precise breakpoints for the SVs detected, including in regions of low-mappability.

In summary, we have presented a novel method that enables the systematic detection of CNVs across the genome. Applying this method to a set of 640 WGS datasets, we were able to produce the most exhaustive map of CNVs across the human genome, including regions that were not well covered by the most recent CNV catalog³¹. We also highlighted the broad potential impact of this type of genetic variation including in regions of low mappability. In the future, we anticipate that population-based methods, such as PopSV, will facilitate the identification not only of CNVs but also of other types of SVs in both normal and cancer genomes.

4 Data and code availability

The PopSV R package and documentation are available at <http://jmonlong.github.io/PopSV/>. The scripts used to produce the graphs and numbers in this study have been deposited on <https://figshare.com/s/ba79730bb87a1322480d>. It also contains the necessary data to reproduce our results. The raw sequences of the different datasets have already been deposited by their respective consortium (Methods).

5 Acknowledgments

This work was supported by a grant from the National Sciences and Engineering Research Council (NSERC-448167-2013) and a grant from the Canadian Institute for Health Research (CIHR-MOP-115090). SLG and GB are supported by the Fonds de Recherche Santé Québec (FRSQ-29493 and FRSQ-25348). Data analyses were enabled by compute and storage resources provided by Compute Canada and Calcul Québec. We are grateful to the team of the Québec Study of Newborn twins who provided the twin dataset and the Cagekid consortium who provided the renal cancer dataset. This study also made use of data generated by the Genome of the Netherlands Project. A full list of the investigators is available from www.nlgenome.nl. Funding for the project was provided by the Netherlands Organization for Scientific Research under award number 184021007, dated July 9, 2009 and made available as a Rainbow Project of the Biobanking and Biomolecular Research Infrastructure Netherlands (BBMRI-NL). The sequencing was carried out in collaboration with

the Beijing Institute for Genomics (BGI). Finally, we would like to thank Simon Gravel, Mathieu Blanchette, Mathieu Bourgey and Toby Dylan Hocking for helpful discussions.

6 Author Contributions

JM and GB conceived and designed the study. JM implemented the method and performed the analyses. JM, CM and SLG designed and performed experimental validation. GR, PC and SLG contributed reagents/materials. Finally, JM and GB wrote the manuscript.

Abbreviation

CNV Copy-Number Variation or Copy Number Variant.

Kb Kilo base.

RD Read-Depth, also called read coverage or depth of coverage.

SV Structural Variation or Structural Variant.

WGS Whole-Genome Sequencing.

References

- [1] I. M. Hall and A. R. Quinlan. Detection and interpretation of genomic structural variation in mammals. *Methods in molecular biology (Clifton, N.J.)*, 838:225–48, 2012.
- [2] A. J. Sharp, Z. Cheng, and E. E. Eichler. Structural variation of the human genome. *Annual review of genomics and human genetics*, 7:407–442, 2006.
- [3] R. E. Mills *et al.* Mapping copy number variation by population-scale genome sequencing. *Nature*, 470(7332):59–65, 2011.
- [4] A. W. Pang *et al.* Towards a comprehensive structural variation map of an individual human genome. *Genome biology*, 11(5):R52, 2010.
- [5] S. a. McCarroll *et al.* Deletion polymorphism upstream of IRGM associated with altered IRGM expression and Crohn’s disease. *Nature genetics*, 40(9):1107–1112, 2008.
- [6] J. L. Stone *et al.* Rare chromosomal deletions and duplications increase risk of schizophrenia. *Nature*, 455(7210):237–241, 2008.
- [7] E. G. Bochukova *et al.* Large, rare chromosomal deletions associated with severe early-onset obesity. *Nature*, 463(7281):666–670, 2010.
- [8] H. C. Mefford *et al.* Rare copy number variants are an important cause of epileptic encephalopathies. *Annals of Neurology*, 70:974–985, 2011.
- [9] H. Stefansson *et al.* CNVs conferring risk of autism or schizophrenia affect cognition in controls. *Nature*, 505(7483):361–6, 2014.

- [10] R. Beroukhi *et al.* The landscape of somatic copy-number alteration across human cancers. *Nature*, 463(7283):899–905, 2010.
- [11] F. Balzola, C. Bernstein, G. T. Ho, and C. Lees. Genome-wide association study of CNVs in 16,000 cases of eight common diseases and 3,000 shared controls: Commentary. *Inflammatory Bowel Disease Monitor*, 11(1):26–27, 2010.
- [12] S. Ayarpadikannan and H.-S. Kim. The Impact of Transposable Elements in Genome Evolution and Genetic Instability and Their Implications in Various Diseases. *Genomics & Informatics*, 12(3):98, 2014.
- [13] H. V. Firth *et al.* DECIPHER: Database of Chromosomal Imbalance and Phenotype in Humans Using Ensembl Resources. *American Journal of Human Genetics*, 84(4):524–533, 2009.
- [14] D. F. Conrad *et al.* Origins and functional impact of copy number variation in the human genome. *Nature*, 464(7289):704–712, 2010.
- [15] M. Spielmann and E. Klopocki. CNVs of noncoding cis-regulatory elements in human disease. *Current opinion in genetics & development*, 23(3):1–8, 2013.
- [16] C. Alkan, B. P. Coe, and E. E. Eichler. Genome structural variation discovery and genotyping. *Nature reviews. Genetics*, 12(5):363–76, 2011.
- [17] K. Chen *et al.* BreakDancer: an algorithm for high-resolution mapping of genomic structural variation. *Nature methods*, 6(9):677–81, 2009.
- [18] M. R. Lindberg, I. M. Hall, and A. R. Quinlan. Population-based structural variation discovery with Hydra-Multi. *Bioinformatics (Oxford, England)*, pages 4–6, 2014.
- [19] V. Boeva *et al.* Control-free calling of copy number alterations in deep-sequencing data using GC-content normalization. *Bioinformatics (Oxford, England)*, 27(2):268–9, 2011.
- [20] A. Abyzov, A. E. Urban, M. Snyder, and M. Gerstein. CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs from family and population genome sequencing. *Genome research*, 21(6):974–84, 2011.
- [21] G. Klambauer *et al.* cn.MOPS: mixture of Poissons for discovering copy number variations in next-generation sequencing data with a low false discovery rate. *Nucleic acids research*, 40(9):e69, 2012.
- [22] K. Ye, M. H. Schulz, Q. Long, R. Apweiler, and Z. Ning. Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics (Oxford, England)*, 25(21):2865–71, 2009.
- [23] A. Rimmer *et al.* Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nature genetics*, 46(8):912–918, 2014.
- [24] T. J. Treangen and S. L. Salzberg. Repetitive DNA and next-generation sequencing: computational challenges and solutions. *Nature reviews. Genetics*, 13(1):36–46, 2012.

- [25] S. M. Teo, Y. Pawitan, C. S. Ku, K. S. Chia, and A. Salim. Statistical challenges associated with detecting copy number variations with next-generation sequencing. *Bioinformatics (Oxford, England)*, 28(21):2711–8, 2012.
- [26] Y. Benjamini and T. P. Speed. Summarizing and correcting the GC content bias in high-throughput sequencing. *Nucleic acids research*, 40(10):e72, 2012.
- [27] A. Koren *et al.* Genetic Variation in Human DNA Replication Timing. *Cell*, 159(5):1015–1026, 2014.
- [28] M.-S. Cheung, T. a. Down, I. Latorre, and J. Ahringer. Systematic bias in high-throughput sequencing data and its correction by BEADS. *Nucleic acids research*, 39(15):e103, 2011.
- [29] F. Hormozdiari *et al.* Next-generation VariationHunter: combinatorial algorithms for transposon insertion discovery. *Bioinformatics*, 26(12):i350–i357, 2010.
- [30] D. He, F. Hormozdiari, N. Furlotte, and E. Eskin. Efficient algorithms for tandem copy number variation reconstruction in repeat-rich regions. *Bioinformatics (Oxford, England)*, 27(11):1513–20, 2011.
- [31] P. H. Sudmant *et al.* An integrated map of structural variation in 2,504 human genomes. *Nature*, 526(7571):75–81, 2015.
- [32] L. C. Francioli *et al.* Whole-genome sequence variation, population structure and demographic history of the Dutch population. *Nature Genetics*, 46(8):818–825, 2014.
- [33] W. P. Kloosterman *et al.* Characteristics of de novo structural changes in the human genome. *Genome Research*, 25(6):792–801, 2015.
- [34] M. E. MacDonald *et al.* A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington’s disease chromosomes. *Cell*, 72(6):971–983, 1993.
- [35] S. M. Mirkin. Expandable DNA repeats and human disease. *Nature*, 447(7147):932–940, 2007.
- [36] J. Rich, V. V. Ogryzko, and I. V. Pirozhkova. Satellite DNA and related diseases, 2014.
- [37] M. Gymrek, D. Golan, S. Rosset, and Y. Erlich. lobSTR: A short tandem repeat profiler for personal genomes. *Genome Research*, 22(6):1154–1162, 2012.
- [38] P. E. Warburton *et al.* Analysis of the largest tandemly repeated DNA families in the human genome. *BMC genomics*, 9:533, 2008.
- [39] S. K. Sen *et al.* Human genomic deletions mediated by recombination between Alu elements. *American journal of human genetics*, 79(1):41–53, 2006.
- [40] S. Yoon, Z. Xuan, V. Makarov, K. Ye, and J. Sebat. Sensitive and accurate detection of copy number variants using read depth of coverage. *Genome research*, 19(9):1586–92, 2009.
- [41] R. Xi *et al.* Copy number variation detection in whole-genome sequencing data using the Bayesian information criterion. *Proceedings of the National Academy of Sciences of the United States of America*, 108(46):E1128–36, 2011.

- [42] G. Glusman *et al.* Identification of copy number variants in whole-genome data using Reference Coverage Profiles. *Frontiers in Genetics*, 6(February):1–13, 2015.
- [43] R. E. Handsaker *et al.* Large multiallelic copy number variations in humans. *Nature Genetics*, 47(3):296–303, 2015.
- [44] M. Boivin *et al.* The Quebec Newborn Twin Study into adolescence: 15 years later. *Twin research and human genetics : the official journal of the International Society for Twin Studies*, 16(1):64–9, 2013.
- [45] G. Scelo *et al.* Variation in genomic landscape of clear cell renal cell carcinoma across Europe. *Nature communications*, 5(May):5135, 2014.
- [46] D.-Q. Nguyen, C. Webber, and C. P. Ponting. Bias of selection on human copy-number variants. *PLoS genetics*, 2(2):e20, 2006.
- [47] P. H. Sudmant *et al.* Global diversity, population stratification, and selection of human copy number variation. *Science*, pages 1–16, 2015.
- [48] K. A. Eckert and S. E. Hile. Every microsatellite is different: Intrinsic DNA features dictate mutagenesis of common microsatellites present in the human genome. *Molecular Carcinogenesis*, 48(4):379–388, 2009.
- [49] T. F. Willems, M. Gymrek, G. Highnam, D. Mittelman, and Y. Erlich. The landscape of human STR variation. *Genome Research*, pages 1894–1904, 2014.
- [50] A. Functammasan *et al.* Accurate typing of short tandem repeats from genome-wide sequencing data and its applications. *Genome Research*, 25(5):736–749, 2015.
- [51] D. Kelley and J. Rinn. Transposable elements reveal a stem cell-specific class of long noncoding RNAs. *Genome biology*, 13(11):R107, 2012.
- [52] X. Lu *et al.* The retrovirus HERVH is a long noncoding RNA required for human embryonic stem cell identity. *Nature Structural & Molecular Biology*, 21(4):423–425, 2014.
- [53] M. J. P. Chaisson *et al.* Resolving the complexity of the human genome using single-molecule sequencing. *Nature*, 2014.
- [54] Y. Shi and J. Majewski. FishingCNV: a graphical software package for detecting rare copy number variations in exome-sequencing data. *Bioinformatics (Oxford, England)*, 29(11):1461–2, 2013.
- [55] C. Wang *et al.* PatternCNV: a versatile tool for detecting copy number changes from exome sequencing data. *Bioinformatics*, pages 1–3, 2014.
- [56] K. R. Rosenbloom *et al.* The UCSC Genome Browser database: 2015 update. *Nucleic Acids Research*, 43(D1):D670–D681, 2015.

7 Methods

7.1 Data

Twin study All patients gave informed consent in written form to participate in the Quebec Study of Newborn Twins⁴⁴. Ethic boards from the Centre de Recherche du CHUM, from the Université Laval and from the Montreal Neurological Institute approved this study. Sequencing was done on an Illumina HiSeq 2500 (paired-end mode, fragment length 300 bp). The reads were aligned using a modified version of the Burrows-Wheeler Aligner (bwa version 0.6.2-r126-tpx with threading enabled). The options were 'bwa aln -t 12 -q 5' and 'bwa sampe -t 12'. The aligned reads are available on the European Nucleotide Archive under [ENA PRJEB8308](#). The 45 samples had an average sequencing depth of 40x (minimum 34x / maximum 57x).

Renal cell carcinoma WGS data from renal cell carcinoma is presented in details in the CageKid paper⁴⁵. In short, 95 pairs of normal/tumor tissues were sequenced using GAIIX and HiSeq2000 instruments. Paired-end reads of size 100 bp totaled an average sequencing depth of 54x (minimum 26x / maximum 164x). Reads were trimmed with FASTX-Toolkit and mapped per lane with BWA backtrack to the GRCh37 reference genome. Picard was used to adjust pairs coordinates, flag duplicates and merged lane. Finally realignment was done with GATK. Raw sequence data have been deposited in the European Genome-phenome Archive, under the accession code [EGAS00001000083](#).

Genome of the Netherlands WGS data from the GoNL project is described in details in Francioli et al.³². This data have been derived from different sample collections:

- The [LifeLines Cohort Study](#), supported by the Netherlands Organization of Scientific Research (NWO, grant 175.010.2007.006), the Dutch government's Economic Structure Enhancing Fund (FES), the Ministry of Economic Affairs, the Ministry of Education, Culture and Science, the Ministry for Health, Welfare and Sports, the Northern Netherlands Collaboration of Provinces (SNN), the Province of Groningen, the University Medical Center Groningen, the University of Groningen, the Dutch Kidney Foundation and Dutch Diabetes Research Foundation.
- The [EMC Ergo Study](#).
- The LUMC Longevity Study, supported by the Innovation-Oriented Research Program on Genomics (SenterNovem IGE01014 and IGE05007), the Centre for Medical Systems Biology and the National Institute for Healthy Ageing (Grant 05040202 and 05060810).
- [VU Netherlands Twin Register](#).

In short, samples were sequenced on an Illumina HiSeq 2000 instrument (91-bp paired-end reads, 500-bp insert size). We downloaded the aligned read sequences (BAM) for the 500 parents in the data set. We further performed indel realignment using GATK 3.2.2, adjusted pairs coordinates with Samtools 0.1.19, marked duplicates with Picard 1.118, and performed base recalibration (GATK 3.2.2). The average sequencing depth was 14x (minimum 9x / maximum 59x).

Genomic annotations Gencode annotation (V19) was directly downloaded from the consortium FTP server at ftp://ftp.sanger.ac.uk/pub/gencode/Gencode_human/release_19/gencode.v19.annotation.gtf.gz. Other genomic annotations were downloaded from the UCSC database⁵⁶ server at <http://hgdownload.soe.ucsc.edu/goldenPath/hg19/database>. The file names of the corresponding annotations are

| | |
|--------------------------------------|-------------------------------------|
| Mappability | wgEncodeCrgMapabilityAlign100mer.bw |
| Cytogenetic bands | cytoBandIdeo.txt.gz |
| Centromere, telomere, assembly gap | gap.txt.gz |
| Segmental duplication | genomicSuperDups.txt.gz |
| Simple repeat / Short Tandem Repeats | simpleRepeat.txt.gz |
| RepeatMasker | rmsk.txt.gz |

7.2 Technical variation in Read-Depth from Whole-Genome Sequencing

To investigate the bias in RD we first fragmented the genome in non-overlapping bins of 5 Kbp. The number of properly mapped reads was used as RD measure, defined as read pairs with correct orientation and insert size, and a mapping quality of 30 (Phred score) or more. In each sample, GC bias was corrected by fitting a Loess model between the bin's RD and the bin's GC content. Using this model, the correction factor for each bin was estimated from its GC content. Bins with extreme coverage were identified when deviating from the median coverage by more than 3 standard deviation. After these conventional intra-sample corrections, RD across the different samples were combined and quantile normalized. At that point the different samples had the same global RD distribution and no bins with extreme coverage or GC bias.

Two control RD datasets were constructed to represent our expectation when no bias is present. One was derived from the original RD by shuffling the bins' RD in each sample. In the second, RD was simulated from a Normal distribution with mean and variance fitted to the real distribution. Simulation or shuffling ensures that no region-specific or sample-specific bias remains. To investigate region-specific bias, we computed the mean and standard deviation of the RD in each bin across the different samples. The same was performed in the control datasets. If there is no bias, the distribution of these estimators should be similar in the original, shuffled and simulated RD.

Next, to investigate experiment-specific bias, we retrieved which sample had the highest coverage in each bin. Then we computed, for each sample, the proportion of the genome where it had the highest coverage. If no bias was present, e.g. in the shuffled and simulated datasets, each sample should have the highest coverage in $\frac{100}{N}\%$ of the genome (with N the number of samples). If some experiment are more affected by technical bias it would be more often extreme. The same analysis was performed monitoring lowest coverage.

Finally, the same analyses were repeated with the challenging regions. Instead of excluding any bin with an extreme coverage in a sample, we kept any bin that was extreme in at least one sample. Hence it is the exact complement of the bins analyzed previously.

7.3 PopSV a population-based approach

Binning and coverage measure The genome is fragmented in non-overlapping consecutive bins of fixed size. We ran two separate analysis on the three datasets. Bin sizes of 5 Kbp and 500 bp were used on the Twin study and renal cell carcinoma. Because of its lower sequencing depth, the 500 bp run on GoNL gave only partial results. More precisely, we observed a truncated distribution

of the copy-number estimates, with most of the 1 and 3 copy number variants missing. It means that at this resolution many one-copy variation cannot be differentiated from background noise. For this reason we finally ran GoNL analysis using 2 Kbp and 5 Kbp bins.

In each bin and each sample the number of reads that overlap the bin and are properly mapped are counted to get a measure of coverage. Here proper mapping means read pairs with correct orientation and insert size, and a mapping quality of 30 (Phred score) or more. The bin counts were then corrected for GC bias. In each sample, a LOESS model was fitted between the bin's count and bin's GC content. A normalization factor was then defined for each bin from its GC content.

Constructing the set of reference samples In each dataset we choose the reference samples as follows: in the renal cancer dataset from the normal samples, in the Twins dataset from all the samples, in GoNL from a subset of 200 samples (see below). For each dataset, a Principal Component Analysis (PCA) was performed across samples on the counts normalized globally (median/variance adjusted). The resulting first two principal components are used to verify the homogeneity of the reference samples. In the presence of extreme outliers or clear sub-groups, a more cautious analysis is recommended. For example, outliers can stay included in the set of reference samples keeping in mind it might harbor more false calls later. Independent analysis in each of the identified sub-group is also a solution, especially when the same samples are to be used as reference. Although our three datasets showed different levels of homogeneity, we didn't need to exclude samples or split the analysis. The effect of weak outlier samples was either corrected by the normalization step or integrated in the population-view.

In GoNL, we decided to use only 200 of the 500 samples as reference. They were selected to span a maximum of the space defined by the principal components. In contrast to random selection, this ensures that weak outliers are included in the final set of reference samples, hence maximizing the technical variation integrated in the population-view.

Moreover, the principal components were used to select one control sample from the final set of reference samples. This sample is used in the normalization step as a baseline to normalize other samples against. We picked the sample closest to the centroid of the reference samples in the Principal Component space.

Normalization Although uniformity of the coverage across the genome is not required for our approach, RD values must be comparable across samples. When a particular region of the genome is tested, sample specific variation of technical origin must be minimized. This is done through a normalization step.

Naive global normalization approaches like the Trimmed-Mean M(TMM) or quantile normalization have been first implemented and tested. The TMM normalization robustly aligns the mean RD value in the samples. Quantile normalization forces the RD distribution to be exactly the same in each samples. After witnessing the presence of un-characterized sample-specific variation, we implemented a more suited normalization.

Targeted normalization uses information across the set of reference samples to identify similar bins across the genome and normalize their counts separately (see Fig. S20). For each bin, the top 1000 bins with similar coverage patterns across the reference samples are used to normalize the coverage of the bin. TMM normalization is used on these top 1000 bins to derive the correct normalization factor for the bin to normalize. Similarity between two bins is measured using

Pearson correlation between the counts across the reference samples. Hence the top 1000 bins are most similar in term of relative coverage across the samples to the coverage in the bin to normalize. If some bias is present in some samples, the top 1000 bins should also harbor this bias. Hence other regions with similar bias patterns are used to correct for it. In this targeted approach, each genomic region is normalized independently. The 1000 supporting bins are saved and used to normalized new samples (e.g. case sample). Although computationally expensive it ensures that all bins are normalized with the same effort. In contrast global normalization or even PCA-based approaches corrects for the most common or spread bias, but a subset of regions with specific bias might not be corrected.

In order to compare the performance of the different normalization approaches we computed a set of quality metrics. The normalized RD will need to be suited for testing abnormal pattern across samples: under the null hypothesis, i.e. for normal bins, the RD should be relatively normally distributed and the samples rank should vary randomly from one bin to the other. The first metric is the *proportion of bins with non-normal RD* across the samples. Shapiro test was performed on each bin and a P-value lower than 0.01 defined non-normal RD. Then the randomness of the sample ranks was tested by comparing the RD of each sample a region with the median across all samples. In regions of 100 consecutive bins, we counted how many times the RD in a sample was higher than the median across sample. If the ranks are random this value should be around 0.5. The probability under the Binomial distribution is computed for each sample and corrected for multiple testing using Bonferroni correction. If any sample has an adjusted P-value lower than 0.05 we consider that the region has non-random ranks. The resulting QC metric is simply the *proportion of regions with non-random sample ranks*. This QC is specifically testing how much sample-specific bias remains. The remaining QC metrics look at the Z-score distribution in each sample. The proportion of non-normal Z-scores is computed by comparing the density curves of the Z-scores and simulated Normal Z-scores. We compute the proportion of the area under the density curve that doesn't overlap the Normal density curve. This estimate of the *proportion of non-normal Z-scores* is computed in each sample. The final metrics are the *average and maximum* across the samples.

Abnormal RD test and Z-score computation The test is based on Z-scores computed for each bin, corrected afterward for multiple testing. The Z-score represents how different the read count in the tested sample is from the reference samples. It is simply:

$$z = \frac{BC_t^b - \overline{BC_{ref}^b}}{sd(BC_{ref}^b)}$$

where BC_t^b is the bin count, i.e. the number of reads, in bin b and sample t .

Inevitably some samples are hosting common CNVs. We observed that just a couple of samples hosting a CNVs could be enough to bias the standard deviation used in the score computation and mask these CNVs in the coming tests. In many cases the RD signal was clearly showing several groups of samples with proportional read counts. To improve the Z-score computation in those regions, a simple approach was used: the samples were stringently clustered using their RD and the group with higher number of samples was chosen as reference and used to compute the mean and standard deviation for the Z-score computation. In practice, this clustering affects only bins with clear clusters but would remove just a few or no samples in most situations. Furthermore, a

median-based estimator was used for the standard deviation as it is less sensitive to outlier removal. A trimmed mean was also preferred over normal mean for its robustness to outliers.

Significance and multiple testing correction The Z-scores for all the bins of a sample are pooled and significance is estimated. Under the null hypothesis of normally distributed read counts, the Z-scores should also follow a normal distribution. For multiple testing correction, the Z-score empirical distribution is used to fit a normal and estimate the P-value and Q-value of each test. This step is performed using `fdrtool` R package.

By default the null distribution fitting for P-value computation assumes that only a low proportion of bins violates the null hypothesis. In aberrant genomes, e.g. in tumor samples, it is often an unrealistic assumption. We devised a new strategy to set the proportion of the empirical distribution, later used to estimate the null distribution variance. Here the null Z-score distribution is assumed to be centered on 0 and its variance is estimated by trimming the tails of the empirical distribution. To find a correct trimming factor, an iterative approach started from a low trimming factor and increased its value until reaching a plateau for the variance estimator. Indeed, once the plateau is reached, additional trimming doesn't change the estimated variance because there is no more abnormal Z-scores, only the central part of the null distribution. Samples with an important proportion of abnormal genome, e.g. tumor samples, showed more appropriate fit.

Of note, the P-values for positive Z-scores (duplication) and negative Z-scores (deletion) are estimated separately. Thus imbalance in the deletion to duplication ratio, or large aberration that lead to asymmetrical Z-score distribution doesn't affect the P-value estimation. Multiple testing correction is performed after pooling all the P-values.

Copy number estimation and other metrics Following the significance estimation, consecutive bins with abnormal coverage are merged into a call. In addition to the Z-score, P-value, Q-value and number of bins of each call, PopSV retrieves the average coverage in the reference samples and the fold change in the sample tested.

Copy number is also estimated by dividing the coverage in a region by the average coverage across the reference samples, multiplied by 2 (as diploidy is expected). In our bin setting, the estimation is correct if the bin spans completely the variant. For this reason we trust the copy number estimate for calls spanning 3 or more consecutive bins, as it is computed using the middle bin(s) which completely span the variant. In other cases we expect the copy number estimate to be under-estimated.

All this additional information can be used to order or retrieve high confidence calls. For examples, several consecutive bins or a copy number estimate around an integer value increases our confidence in a call. In our validation and analysis however, we used the entire set of calls.

The ZZ plots are computed directly from the Z-score of each bin in two different samples (e.g. paired normal/tumor samples, twins). The global distribution of the Z-score is also compared to the mappability estimate of the bins. At this point, we use the mappability track available from UCSC⁵⁶ (see Genomic annotations) and compute the mean level across the bin.

Coverage tracks For each run, we constructed coverage tracks based on the average coverage in the reference samples. Bins where the reference samples had, on average, the expected coverage were classified as *expected coverage*. Bins with a coverage higher or lower than 3 standard deviation from the median were classified as *high coverage* or *low coverage* respectively. To ensure robustness,

the standard deviation was derived from the Median Absolute Deviation. We use *low coverage* regions to define low-mappability regions, as the low coverage is a result of the lower mappability of a region.

Eventually, we also defined *extremely low coverage* region which have an average coverage close to 0. These region are defined by the peaks around 0 in the distribution of average coverage (see Figure S12). This sub-class of *low coverage* region is used in a few of the following analysis to highlight the most challenging regions.

7.4 Validation and benchmark

Running FREEC and cn.MOPS FREEC was run on each sample separately, starting from the BAM file. FREEC internally corrects the RD for GC and mappability bias. In order to compare its performance in low-mappability region, the minimum “*telocentromeric*” distance was set to 0. The remaining parameters were set to default. Of note an additional run with slightly looser parameter (`breakPointThreshold=0.6`) was performed to get a larger set of calls used in some parts of the *in silico* validation analysis to deal with borderline significant calls.

cn.MOPS was run on the same GC-corrected bin counts used for PopSV. All the samples are analyzed jointly. Of note an additional run with slightly looser parameter (`upperThreshold=0.32` and `lowerThreshold=-0.42`) was performed to get a larger set of calls used in some parts of the *in silico* validation analysis to deal with borderline significant calls.

Clustering the Twins samples A distance between two samples A and B is defined as : $1 - 2 \frac{|V_A \cap V_B|}{|V_A| + |V_B|}$ where V_A represents the variants found in sample A, $V_A \cap V_B$ the variants found in both A and B, and $|V|$ the cumulative size of the variants. Hence the similarity between two samples is represented by the amount of sequence found in both divided by the average amount of sequence called. This distance is used for hierarchical clustering of the samples. Different linkage criteria (*average*, *complete* and *Ward*) were used for the exploration. In our dendrograms we used the *average* linkage criterion. The same clustering was performed using only calls in regions with extremely low coverage (reference average ≤ 10 reads, see Coverage track section).

Frequency peak in Twins The frequency at which a region is affected by a CNV was compared between the different methods. In the Twins dataset, we expect a peak around frequency of 3 samples : the two twins and one parent. To compare the different methods the height of the peak, in the frequency distribution, represent the proportion of the affected genome called at each frequency.

Replication in twins For each twin, a CNV call was defined as *replicated* if also found in the other twin. In order to avoid missing calls with borderline significance, we used slightly less confident calls for the second twin. We removed calls present in more than 50% of the samples as they could be systematic errors that would look replicated. Hence a *replicated* call is most likely true as it is present in a minority of samples but consistently in the twin pair. The proportion of *replicated* calls per sample gives an estimate of specificity. The level of sensitivity is represented by the number of *replicated* calls. Even if we removed systematic calls, the most frequent calls in the cohort are more likely to look *replicated* by chance, compared to rare calls. To normalize for this effect, we use the frequency distribution to compute the number of *replicated* calls expected

by chance. In practice the null concordance for each call is simulated by a Bernoulli distribution of parameter the frequency of the call. This number of *replicated* calls by chance is subtracted to the original number of *replicated* calls to give a adjusted measure of sensitivity. Although we don't know the true number of variant, this number of *replicated* calls is used to compare the different methods. The same analysis is also performed using only calls located in low-mappability regions in order to get an estimate on challenging regions. A call was considered in a low mappability region if more than 90% of its sequence was annotated as so.

In addition to this *per-sample* concordance, we compute a *per-region* concordance estimate by pooling all the calls from all the samples. When more than 90% of the twins called in a specific bin have the other twin called too, the bin is classified as *reliable*. Then the bins can be grouped according to their GC content or repeat content to test that the quality of the calls is stable. This approach is particularly useful to verify that the proportion of *reliable* bins is similar in bin with extreme GC content or different repeat content. Finally we compute a null distribution with the same approach but using randomly selected samples instead of the sample called in each bin. Dividing the proportion of *reliable* bin by its null equivalent gives an idea of the significance of the observation. This fold-enrichment from the null is used to compare the different methods. Figure S21 shows PopSV's robustness and superior performance even in challenging regions. In addition, we use this *per-region* metric to estimate the amount of the genome that can be correctly tested. Here the genome is fragmented in 1 Mbp windows and we count how many show more *reliable* regions than *reliable* by chance. The 1 Mbp fragmentation is used in order to avoid biases from segmentation behavior. If the regions were used as-is, a segmentation that tend to locally call longer segments will look largely superior even-though it calls the same variants. The fragmentation of the genome in large windows limits this bias and allows for fair comparison between the different methods. By counting how many 1 Mbp windows can be called correctly, we estimate how much of the genome can be correctly tested by each method. We observed that a higher fraction of the genome is reliably called with PopSV compared to cn.MOPS and FREEC (1.5 and 2.7 more, respectively, Table S2 and Methods).

Replication between paired normal and tumor samples The same approach as described previously when comparing pairs of twins was applied in the renal cancer dataset, on pairs of normal/tumor samples. Here true germline calls should also be found in the paired tumor sample, and replication estimates is computed using normal samples only. Again both *per-sample* and *per-region* estimates are computed and compared between methods. Reliability is defined based on the proportion of the normal samples with consistent calls in their paired tumor (Table S1 and S2).

Concordance between different bin sizes We compared the calls using small bins (500 bp) and the calls using larger bins (5 Kbp). In theory, calls from the 5 Kb analysis should be supported by many 500 bp calls. We also expect large stretches of 500 bp calls to be detected in the 5 Kb analysis. This comparison is informative as it explores the quality of the calls, the size of detectable events and the resolution for different bin size. First we counted how many small bin calls supported any large bin call. These metrics were separated according to the size of the large bin call. Overall, we find that 5 Kb calls are well supported by 500 bp calls, with only 14% of the 5 Kb bins not supported by any 500 bp bin (Fig. S22a). To investigate large bin calls with no supporting small bin call, we display the average Z-scores in the small bins overlapping large bin calls. This is useful to test if the lack of support is due to lower confidence or real discordance between the different

runs. If the Z-scores in the small bins deviates from 0 in the correct direction we conclude that they support the large bin call. Even for these unsupported 5 Kb calls, we find that the 500 bp bins RD was consistently enriched (or depleted) although not enough to be called with confidence (Fig. S22b and S22c). This is expected given the higher background noise in the 500 bp analysis that will reduce the power to call these variants. Next, we looked at the proportion of 500 bp calls, grouped by size, that were found in the 5 Kb calls. More specifically, we grouped them by size to verify that large enough small bin calls are present in the large bin calls. This analysis is used to both test the sensitivity of PopSV with a particular bin size, and its resolution to variants smaller than the bin size. Indeed this framework allow us to ask questions such as: how much of the variants spanning only half a bin are detected? We find that the concordance gradually increases until the 500 bp calls reach 5 Kb in size where the concordance rises to nearly 100% (Fig. 2f). This suggests that PopSV is able to detect approximately 75% of the events as large as half its bin size, and almost all events larger than its bin size. As expected, only a small proportion of the small 500 bp calls overlap 5 Kb calls and they likely corresponds to fragmented larger calls. Considering the trade-off between bin size and noise, this suggests to run PopSV with a few bin sizes to better capture variants of different sizes.

Experimental validation Experimental validation was performed on samples from the Twin study. The 20 variants chosen for experimental validation were randomly selected among both one-copy and two-copy deletions. We selected both small (~ 700 bp) and large (~ 4 Kbp) variants in each class. The coverage at base pair resolution was visually inspected for each deletion in order to map the breakpoints. PCR primers were designed to target the whole deleted region. We performed long-range PCR followed by gel electrophoresis. We then compared the size of the amplified fragment in affected and control samples. If the affected sample showed a lower band than a control with a predicted 2 copies, the deletion was considered validated. On the other hand if affected sample and controls had one similar band, the deletion was considered non-validated. Of note, the validation rate might be under-estimated because visual prediction of the breakpoint is not always accurate and could lead to non-validation when the variant is actually present.

We then randomly selected deletions overlapping low mappability regions and detected in 6 samples or less. We chose to test rare variants because they are likely enriched in false-positives. Hence, this batch of validation represents the most challenging regions to call and validate, and enriched in false-positives. Here we couldn't use the base-pair coverage to fine-tune the breakpoints because the low-mappability blurs any clear signal. Instead we retrieved the reads (and their pairs) mapping to the region and assembled them. With this approach we could sometimes get a better breakpoint resolution and design PCR primers that would amplify the deleted region. In addition to gel electrophoresis, the amplified DNA of some regions was sequenced using Sanger sequencing.

7.5 Genomic patterns of CNVs

Merging results using two different bin sizes Small bins gives better resolution for smaller variant. Large bins gives better sensitivity. For this reason we merged the calls from the 500 bp bin and 5 Kbp bin runs. Variant supported by both sets of calls were merged into one. To decide which set to use to define breakpoints and other information (e.g. copy number estimate), the proportion of overlap was used. If call(s) using small bins overlapped more than a third of a call from the large bin run, it was considered fully recovered by the small bin call which was then used to define breakpoints and other information. If not, the large bin run was considered more appropriate to

define the final breakpoints and additional information. Calls unique to each run were simply added to the final set of calls.

Computing global estimates of copy number variation In Table 1, a call in low coverage region is completely located within *extremely low coverage* regions (as defined by our coverage tracks). The amount of sequence affected in a genome is computed by merging all the variants (e.g. if several samples are combined) and counting the number of bases in this merged set. After the merging step, each base of the genome either overlapped a merged variant or not. Hence each affected base is counted only once, even if it overlaps CNVs in several samples, or with large copy number differences.

Comparing with 1000 Genomes SV set The SV catalog from Sudmant et al.³¹ was downloaded and parsed into our preferred BED-like format. We first checked that we could reproduce the numbers in the main SV paper. Then we retrieved the set of autosomal deletion, duplication and CNVs. We removed deletions smaller than 300 bp as well as variants with high frequency (> 80%). This sub-set of SV represent CNVs that could in theory be detected by PopSV's approach. Using this sub-set, we derived the number of variants, number of variants smaller than 3 Kbp, number of variants in *extremely low coverage* regions, and amount of genome affected. These number are computed exactly as the one presented in Table 1 for PopSV's results, and hence can be compared.

Distance to centromere, telomere and assembly gaps The centromeres, telomeres and assembly gaps (CTGs) are annotated in the gap track from UCSC⁵⁶. However some chromosomes were missing telomere annotations. We defined them as the 10 Kbp region at the ends of chromosomes derived from the cytogenetic bands track.

The distance from each variant to the nearest CTG was computed and represented as a cumulative proportion, meaning the proportion of variants located at a distance d or closer to a CTG.

Because this distribution changes with the size of the variants, we sampled random regions in the genome with similar sizes and computed the same distance distribution. Thanks to this null distribution we are able to see if variants are closer/further to CTG than we would expect by chance.

Simulating control regions Control regions are simulated to have the same size distribution and same overlap with specified genomic features. In practice, this was used to control for the distance to centromere, telomeres and assembly gaps, as well as the overlap with segmental duplications. Hence the patterns observed afterwards are not caused by the over-representation of region in segmental duplication or their proximity to CTGs.

First, thousands of bases are randomly chosen in the genome. The distance between each base and the genomic features is then computed. At this point, simulating a region of a specific size and with specific overlap profile can be done by randomly choosing as center one of the bases that fit the profile :

$$\{b, \forall \text{ feature } f, O_f(d_f^b - \frac{S_r}{2}) < 0\}$$

with O_f equals 1 if the original region overlaps with feature f , -1 if not; d_f^b is the distance between base b and feature f ; and S_r is the size of the original region.

Hence for each input region, a control region is selected as described and will have by construction the exact same size and overlap profile (Fig. S23) as the input.

The number of random bases is important: a low number might result in duplicated regions in the output, but a high number is more computationally demanding. In practice, we perform the simulation twice with 10^5 random bases. The second run is used to simulate again all the duplicated regions from the first run.

Of note, to control for the distance to CTG, we select the base that fit the defined profile and with the closest distance to CTG as desired. Although it doesn't result in exactly the same distribution, it gave satisfactory results (Fig. S23).

Enrichment in genomic features CNVs and control regions are overlapped with genomic features. We then compute the proportion of regions overlapping each feature.

The proportion of overlap and the control regions are computed separately for each sample. Hence the control region fits perfectly the profile of the variants in each sample and is not simply a reflection of the majority of the samples. For each sample (and each feature), the enrichment measure is the difference between the proportion in the original and control regions. A Wilcoxon test on this measure assesses how significant is the potential deviation from 0. The fold-enrichment is the ratio between overlap proportion between original and control regions.

Eventually, we display how much of a variant overlap a feature of interest. This distribution is useful to get a sense if the genomic feature overlap completely the variants or just a small fraction of them.

Somatic variant definition Somatic variants were defined as variant in a tumor samples with no or low overlap with variant in the paired normal sample. In CageKid data, overlapping tumor variant with the ones from the paired normal showed almost only two peaks, at 0 and 100% overlap. A tumor variant was defined as somatic if it overlapped less than 10% of any variant in the paired normal.

Frequency distribution The frequency at which a region is affected by a CNV is computed using calls from the 640 samples. The copy-number change is not taken into account in the computation and the frequency is derived for all the nucleotide that overlaps at least one CNV. The cumulative proportion of affected genome is shown for each frequency in the frequency curve. In addition, frequency curves are computed using small or large variants, exonic or non-exonic variants, and deletions or duplications.

Eventually, we perform the same analysis with the set of comparable CNVs extracted from the 1000 Genomes catalog. Of note, the CNV set was down-sampled to 640 random samples in order to give comparable frequency curves.

CNV impact Exons of protein-coding genes and promoter regions (10 Kbp upstream of the transcription start site) were extracted from the Gencode annotation v19. We counted how many different genes had their exons, exons + promoter and exon + promoter + introns hit by a CNV, in a sample or in the entire dataset.

We did the same for CNVs that overlapped more than 90% of specific classes of repeats. These numbers are shown in table 2. For example, STR-CNVs are CNVs with more than 90% of the region annotated as STR.

8 Supplementary Figures

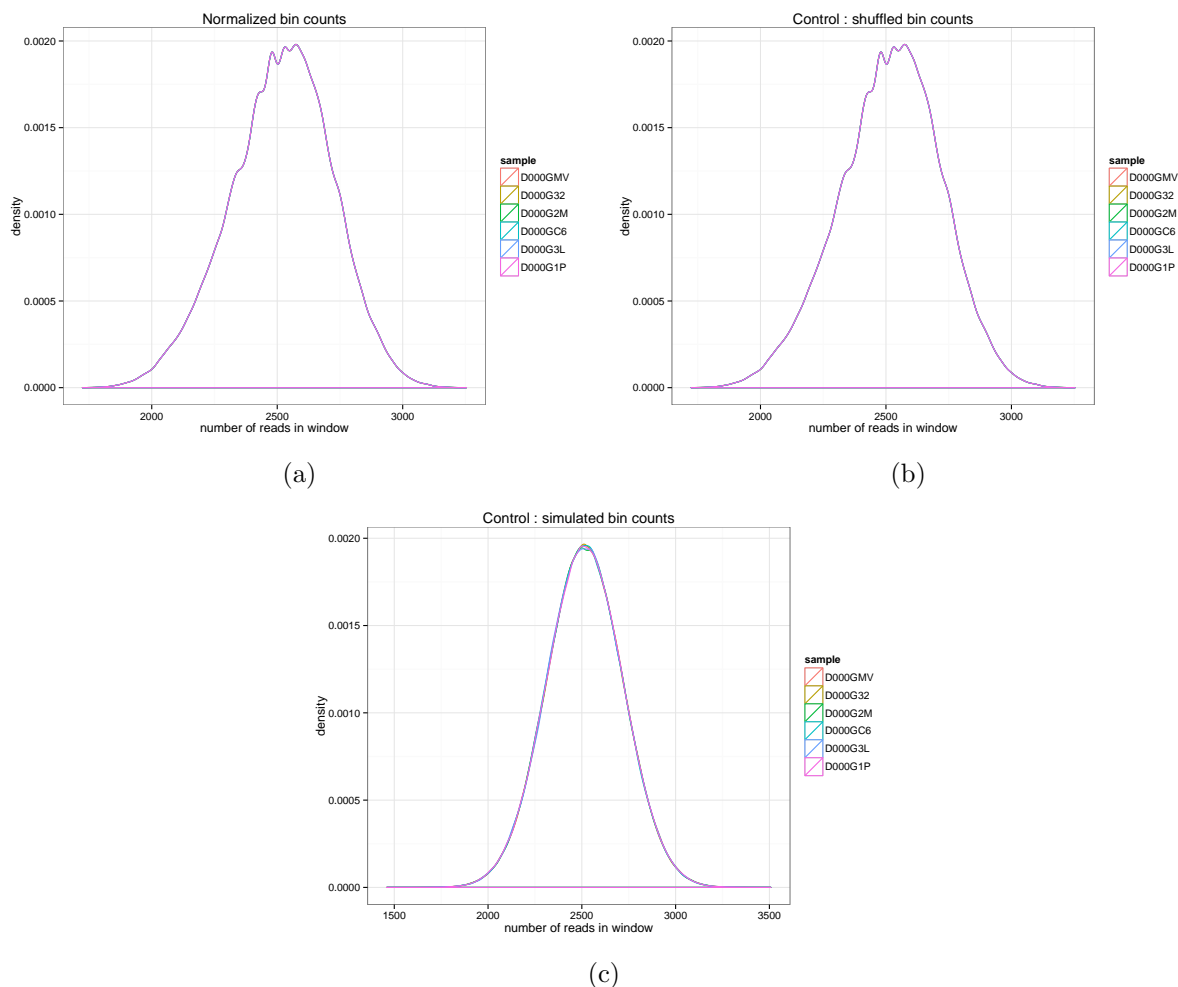


Figure S1: **Distribution of the bin counts after removal of regions of extreme coverage and normalization.** a) All samples have exactly the same RD distribution after quantile normalization. We build the distribution under the null hypothesis (i.e. uniform coverage) by shuffling the bins (b) or simulating RD from a Normal distribution (c).

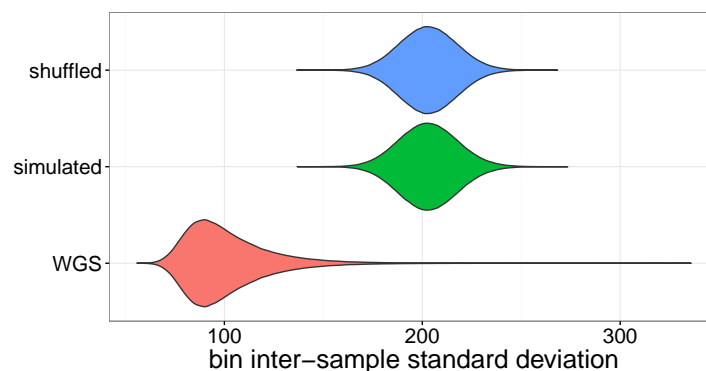


Figure S2: **Variation and bias in whole-genome sequencing experiments.** Distribution of the bin inter-sample standard deviation coverage (red) and null distribution (blue: bins shuffled, green: simulated normal distribution).

| Dataset | Region | Bin size | Number of concordant calls | | | Fold change PopSV vs | | Proportion of concordant calls | | |
|--------------|--------------|------------|----------------------------|-------|---------|----------------------|---------|--------------------------------|-------|---------|
| | | | PopSV | FREEC | cn.MOPS | FREEC | cn.MOPS | PopSV | FREEC | cn.MOPS |
| Twin study | whole genome | 5kbp | 324.5 | 101.5 | 91.5 | 3.20 | 3.55 | 0.92 | 0.93 | 0.95 |
| | | 500bp-5kbp | 883.0 | 140.0 | 506.5 | 6.31 | 1.74 | 0.89 | 0.92 | 0.88 |
| | low coverage | 5kbp | 173.5 | 19.0 | 69.5 | 9.13 | 2.50 | 0.96 | 0.89 | 0.94 |
| | | 500bp-5kbp | 546.0 | 23.5 | 407.5 | 23.23 | 1.34 | 0.94 | 0.90 | 0.87 |
| Renal cancer | whole genome | 5kbp | 293.0 | 48.0 | 75.0 | 6.10 | 3.91 | 0.88 | 0.80 | 0.88 |
| | | 500bp-5kbp | 949.0 | 80.0 | 564.0 | 11.86 | 1.68 | 0.79 | 0.72 | 0.75 |
| | low coverage | 5kbp | 107.0 | 6.0 | 52.0 | 17.83 | 2.06 | 0.91 | 0.78 | 0.87 |
| | | 500bp-5kbp | 445.0 | 2.0 | 267.0 | 222.50 | 1.67 | 0.83 | 0.62 | 0.72 |

Table S1: **Concordance in different datasets, methods and bin size.** The numbers are the average across samples.

| Dataset | Bin size | Number of reliable 1 Mb bins | | | Fold change PopSV vs | |
|--------------|------------|------------------------------|-------|---------|----------------------|---------|
| | | PopSV | FREEC | cn.MOPS | FREEC | cn.MOPS |
| Twin study | 5kbp | 1260 | 753 | 353 | 1.67 | 3.57 |
| | 500bp-5kbp | 2034 | 762 | 1360 | 2.67 | 1.50 |
| Renal cancer | 5kbp | 2107 | 808 | 484 | 2.61 | 4.35 |
| | 500bp-5kbp | 2699 | 1149 | 2106 | 2.35 | 1.28 |

Table S2: **Amount of genome reliably tested in different datasets, methods and bin size.**

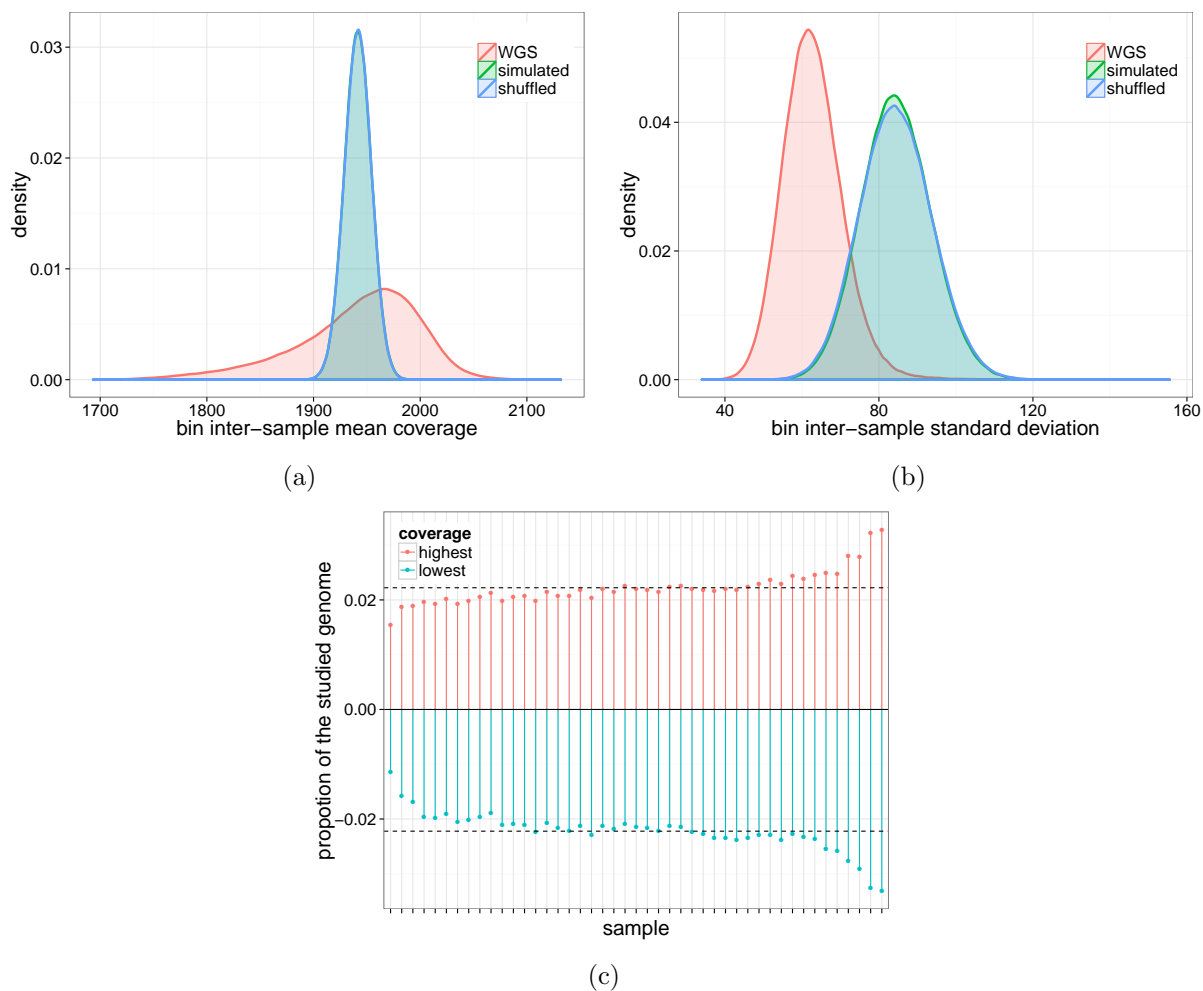


Figure S3: **Variation and bias in whole-genome sequencing in the *Twins* dataset.** a) Average bin RD across the samples (red) and null distribution (blue: bins shuffled, green: simulated normal distribution). b) Same with standard deviation. c) Proportion of the genome in which a sample (x-axis) has the highest(red) or lowest(blue) RD. In the absence of bias all samples should be the extreme one with the same frequency (dotted horizontal line).

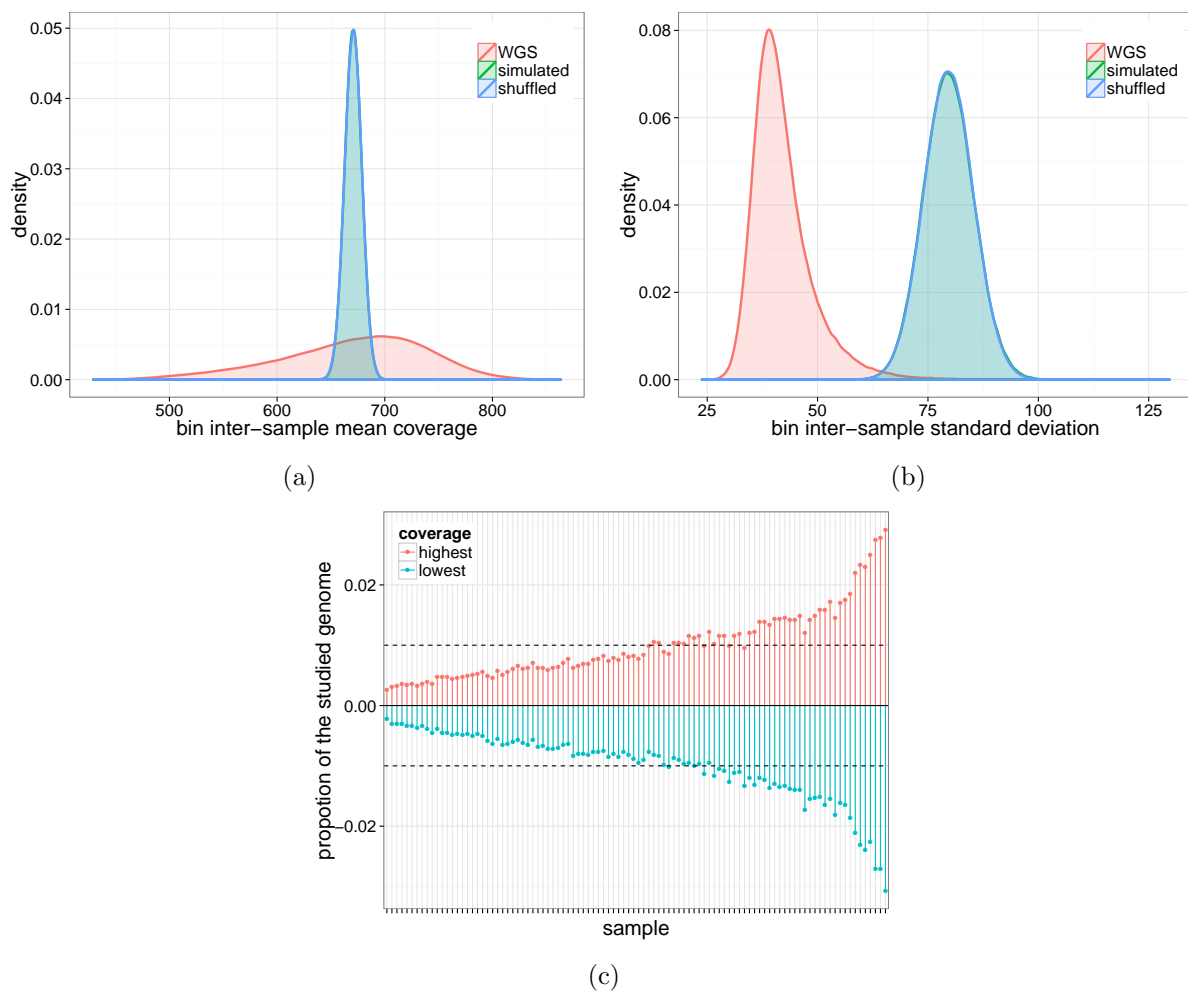


Figure S4: **Variation and bias in whole-genome sequencing in the GoNL dataset.** a) Average bin RD across the samples (red) and null distribution (blue: bins shuffled, green: simulated normal distribution). b) Same with standard deviation. c) Proportion of the genome in which a sample (x-axis) has the highest(red) or lowest(blue) RD. In the absence of bias all samples should be the extreme one with the same frequency (dotted horizontal line).

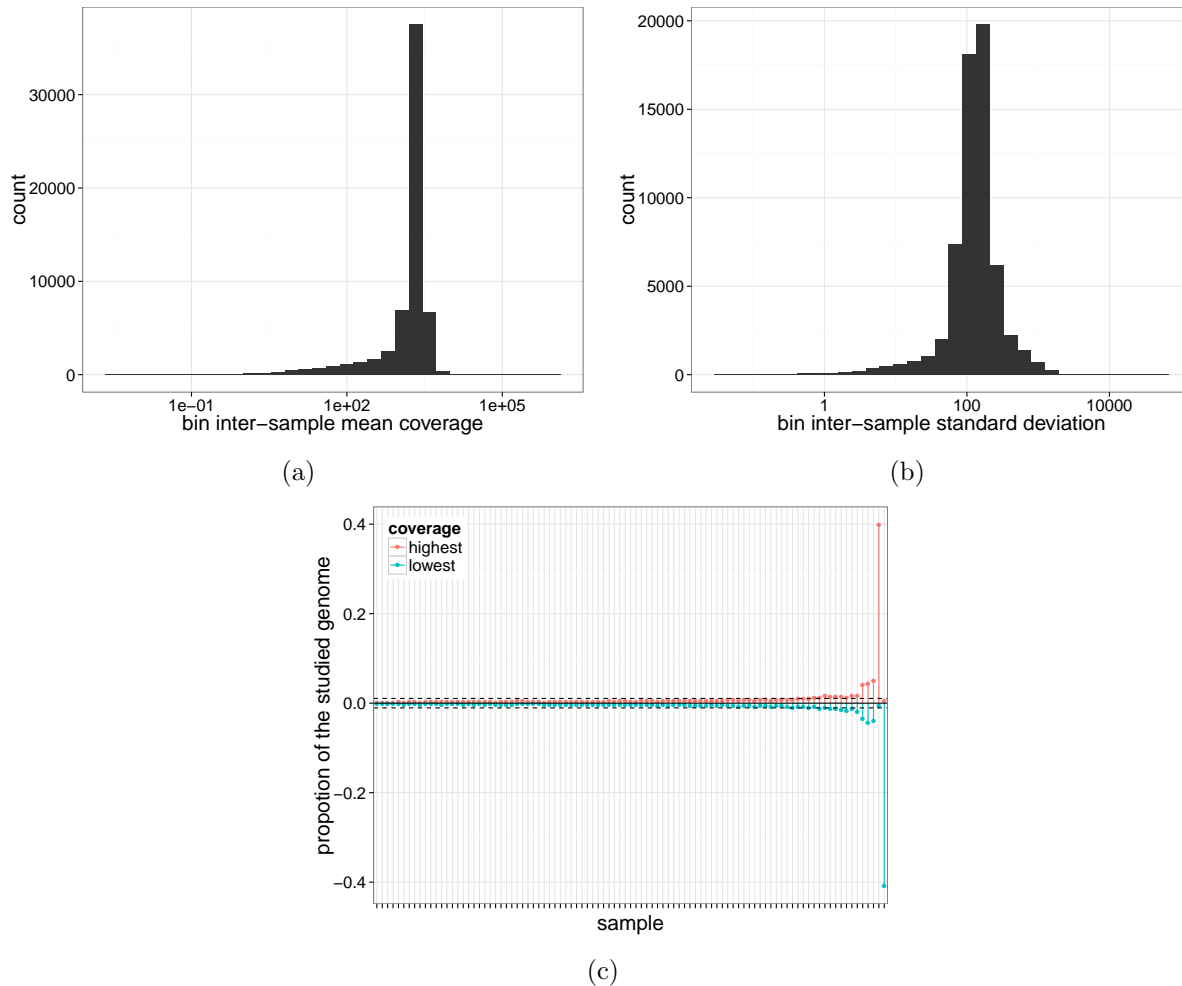


Figure S5: **RD bias is stronger when including all genomic regions.** In renal cancer normals, the same analysis as summarized in Fig. 1a and 1b is performed using all genomic regions, i.e. without filtering for extreme coverage. Quantile normalization is used again to force the same RD distribution in all samples. Of note, in a) and b) the distribution of the mean and variance across samples is shown on a log-scale as it spans several orders of magnitude.

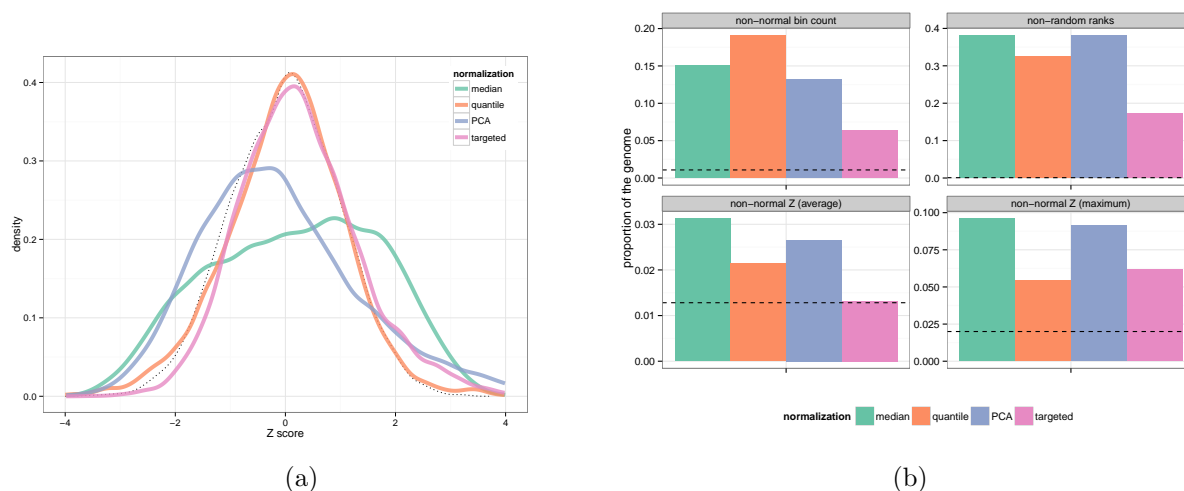


Figure S6: Comparison of different normalization approaches. a) For each normalization approach, the sample with the least normal Z-score distribution is shown. b) After targeted normalization, a lower proportion of the genome looks problematic for the analysis. Less bins have non-normal bin counts (top-left), the sample ranks are more random pointing at less sample-specific bias (top-right), and Z-scores fit better a Normal distribution on average (bottom-left) and in the worst sample (bottom-right). The dotted line is computed from simulated bin counts.

| Validated | Chr. | Start | End | Class | Left PCR primer | Right PCR primer |
|-----------|------|-----------|-----------|------------|---------------------------------|---------------------------------|
| V | 3 | 6649794 | 6654897 | large CN 0 | CCTTAGTATTTCAGTGGTTTCTGTGAGGTAT | ATAAATATCAGTGCCTCAACTTGGACTT |
| V | 5 | 127407030 | 127411341 | large CN 0 | TATTCATATTAACCTATCCTCACAGAAAGA | TTTTTAAGAGATTTGAACTAAAATTCCAC |
| V | 3 | 5535139 | 5539535 | large CN 0 | TACTTTTTGAATTTGTAAATTTCCTTTGTA | GAAATCAGAAAATCAAGATCATACTGAAG |
| V | 1 | 116229111 | 116233162 | large CN 0 | GTGTTACAGAATTAGTTTTACTGAGTGGTC | ATCTATAAAGAAGCTTTTCCAAAATAAACCA |
| V | 1 | 158961082 | 158966958 | large CN 1 | GTAGAATGAGCTGTGTATGAGATGGT | ATGACTTTCATTGTTTGAATGTAGTGAC |
| V | 15 | 26748887 | 26752614 | large CN 1 | CAATTTATCTATCAAGTTATTTACGGTAG | AGTGAGATTTCAATTTAAGCTTGTCTTC |
| V | 6 | 33937344 | 33942846 | large CN 1 | ACATTGTAGCCTGATGACCTTGTTTC | TGTGTTCTGAGGTTTACTTTATAATCTAGG |
| V | 12 | 82095501 | 82099389 | large CN 1 | ACCTATAACTAAGTGTAGCTGCTGTAACCTG | TCAGTAAAAATGATTACTACAGTGGAAAAAT |
| V | 5 | 8255604 | 8260914 | large CN 1 | TGAACATACATTCATACACATAATACAA | TACATCACTGAACAAACCTCTATAGTCATA |
| V | 20 | 7398397 | 7403743 | large CN 1 | AATAAACATTCTCTATAAACCTAAAATGG | CTTTGTACCATATTTTATAAACGTTAGATC |
| V | 18 | 40053822 | 40057873 | large CN 1 | TAACCTTCTTTTCTAAAGCTTTTGAGATAT | GTGAATTAAGATTCATGCTCTGCTAATA |
| V | 16 | 48904951 | 48906510 | small CN 0 | TCTTATTTATTTTGACAGCTCTTTACTCTG | AGATAATCAACTCTTTGTTTATCTTTTCAG |
| V | 2 | 241086647 | 241087801 | small CN 0 | ATCAACATTTAGCCAGTGTGTCTTAG | GTCTCTTGTGCTCTACTTTGGCTT |
| V | 13 | 110221621 | 110222631 | small CN 0 | ACCTCAGGAGAACTACTTCATACATTTCTA | GTATGAAAAACACTCATGGATATCATTCT |
| V | 11 | 60571017 | 60572170 | small CN 0 | AATGTTGAAGTGTGCTTTCTGTAATATCT | GTGTTTTGTGTGCTATTTGTTTAGTA |
| V | 5 | 166402295 | 166404219 | small CN 0 | TCACCTTATTCATAACATTTTCAGTGTAGAG | GATCATATGCTTAAAATGCTAATGAGG |
| N | 3 | 160126422 | 160127288 | small CN 1 | TAAGATACAAGAAATAGAGATAACACTGGG | TCTGAACACTTATTTAAGAAAATGAAAAA |
| N | 17 | 10612674 | 10613775 | small CN 1 | AATTTAGCAGTCTCTTACATTTCTTCTACC | TCTCTTCTATAAAAATAAATGGCTAAAAGC |
| V | 10 | 70253713 | 70255155 | small CN 1 | AATAAAATCAAAGGTATATTAAGTACAGAG | ATATACTCTTTAACTTTTGACCATTTTGG |
| V | 8 | 53700635 | 53702050 | small CN 1 | TAAGGAAAATTTAGTATAGTCTGGACCTGT | ATGAAAATATATCTGTATGGGTGAC |

Table S3: Experimental validation results. Location of the validated (V) and non-validated (N) CNVs for different classes. The last two columns show the primer sequences used for PCR amplification.

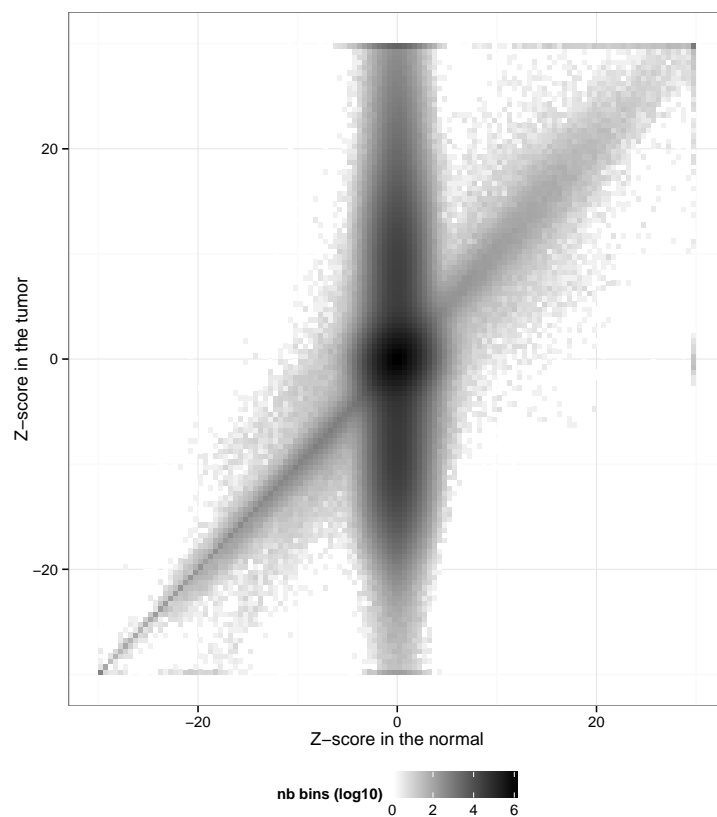


Figure S7: **ZZ plots between normal and tumor pairs.** In renal cancer, Z-scores from each normal samples (x-axis) is plotted against Z-scores from its tumor samples (y-axis). This graph is an aggregation of all normal/tumor pairs. Z-scores are winsorized at -30 and 30 for visibility purpose.

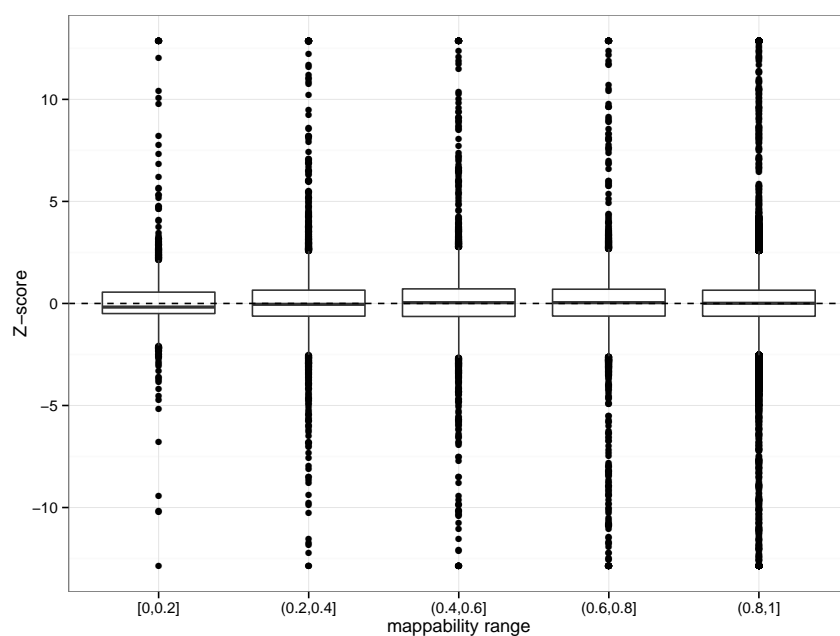


Figure S8: **Z-score distribution versus the mappability of the bin.** One randomly selected sample from the Twin dataset. At this stage, mappability was extracted from the UCSC track (Methods).

| Chr. | Start | End | CN | PCR product size | PCR product size when deletion | Validated | Gel | Sanger Sequencing |
|------|-----------|-----------|------|------------------|--------------------------------|-----------|------------------|-------------------|
| 14 | 40098378 | 40100213 | 0 | 2586 | 751 | Yes | Different bands | Yes: confirmed |
| 5 | 85559864 | 85564846 | 1.05 | 5690 | 708 | Yes | Different bands | Yes: confirmed |
| 6 | 14299746 | 14299801 | 0.79 | 755 | 700 | Yes | Double bands | No |
| 7 | 153000055 | 153000246 | 1.76 | 1137 | 946 | Yes | Double bands | Yes: confirmed |
| 4 | 96401034 | 96401460 | 1.13 | 745 | 319 | Yes | Double bands | No |
| 16 | 34230052 | 34230512 | 1 | 1139 | 679 | Yes | Double bands | No |
| 16 | 8688137 | 8689592 | 1.02 | 2121 | 666 | Yes | Double bands | Yes: confirmed |
| 2 | 12018994 | 12022932 | 1.02 | 4291 | 353 | Yes | Double bands | Yes: confirmed |
| 3 | 121051576 | 121060845 | 1.14 | 9485 | 216 | Yes | Double bands | No |
| 3 | 54433855 | 54433912 | 0 | 952 | 895 | Yes | One band | Yes: insertion |
| 2 | 151031059 | 151038246 | 1.11 | 7485 | 298 | Yes | Small band only | No |
| 9 | 45462450 | 45462522 | 1.1 | 530 | 458 | No | One band | No |
| 7 | 63233184 | 63233261 | 1.33 | 390 | 313 | No | One band | Yes: nothing |
| 9 | 106371251 | 106371330 | 1.28 | 484 | 405 | No | One band | No |
| 16 | 20466400 | 20466487 | 1.27 | 393 | 306 | No | One band | No |
| 5 | 85559864 | 85564842 | 0.78 | 5690 | 712 | No | One band | No |
| 10 | 65703860 | 65708900 | 1.64 | 5430 | 390 | No | One band | No |
| 7 | 159117395 | 159122761 | 1.09 | 5909 | 543 | No | One band | No |
| 2 | 83066824 | 83068234 | 0.57 | 2097 | 687 | NA | No amplification | No |
| 13 | 35996202 | 35996254 | 1.13 | 546 | 494 | NA | Non-specific | No |
| 4 | 159799983 | 159801372 | 1.03 | 2313 | 924 | NA | Non-specific | Yes: not clear |
| 7 | 52963172 | 52964911 | 1.48 | 2316 | 577 | NA | Non-specific | No |
| 10 | 69323932 | 69326507 | 1.62 | 2795 | 220 | NA | Non-specific | Yes: not clear |
| 6 | 58618198 | 58624080 | 1.04 | 6518 | 636 | NA | Non-specific | No |

Table S4: **Experimental validation in low-coverage regions.** The result of the PCR validation was either concordant with PopSV call (Yes), discordant (No) or inconclusive (NA). In some cases, Sanger sequencing was performed. The *CN* column is the estimated copy-number of the deleted allele.

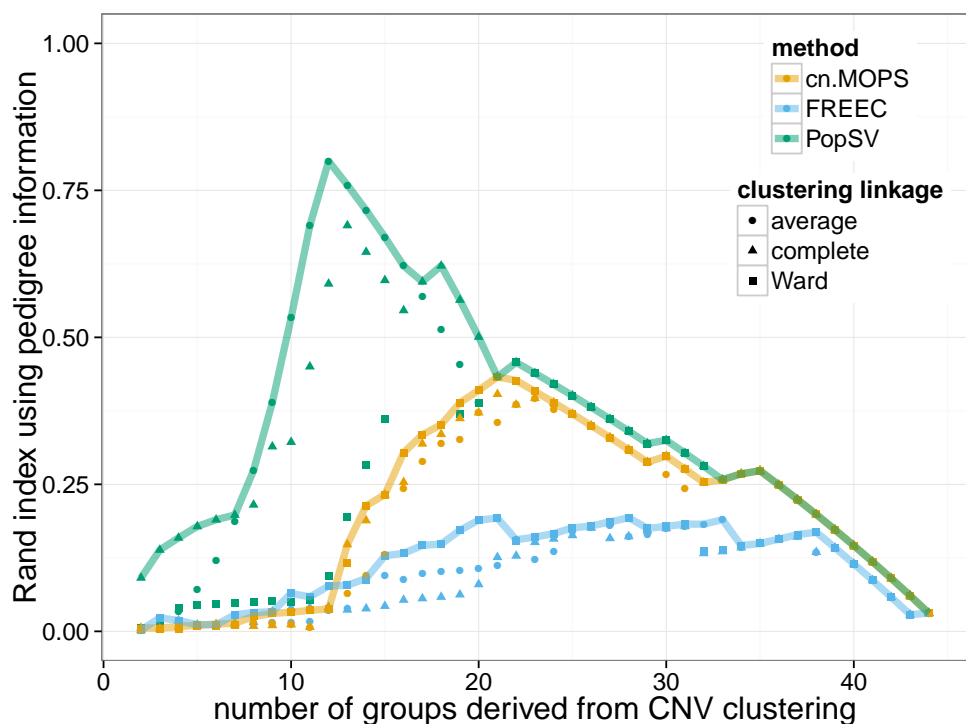


Figure S9: **Twin dataset: sample clustering and pedigree.** Samples are clustered using the CNV calls from the different methods (*colors*). The amount of genomic sequence called in only one of two samples defines the distance used for clustering. After cutting the hierarchical cluster tree (e.g. Fig. 2e) at different levels (*x-axis*), cluster groups are compared to the known pedigree using the Rand index (*y-axis*). Different clustering linkage criterion (*point style*) are used and the one showing the best Rand index is highlighted by the line.

| Sample | Type | Total variants | | Variants per sample | | Variants < 3 Kbp | | Affected genome (Mb) | |
|--------|------|----------------|-----------|---------------------|-----------|------------------|------------|----------------------|------------|
| | | <i>WG</i> | <i>LC</i> | <i>WG</i> | <i>LC</i> | Total | Per sample | Total | Per sample |
| 2504 | All | 2382489 | 3628 | 924 | 2 | 1420566 | 551 | 581.08 | 6.04 |
| | CNV | 312401 | 0 | 124 | 0 | 0 | 0 | 85.05 | 2.74 |
| | DEL | 2041543 | 3628 | 787 | 2 | 1420566 | 551 | 298.70 | 3.13 |
| | DUP | 28545 | 0 | 11 | 0 | 0 | 0 | 264.09 | 0.32 |

Table S5: **1000 Genomes deletions, duplications and CNVs.** We removed variants with high frequency (> 80%), variants in the chromosome X, and variants smaller than 300 bp in order to compare with PopSV's numbers (Table 1). *WG*: whole genome; *LC*: low-coverage regions. *Affected genome* represents the amount of the reference genome that overlaps at least one CNV.

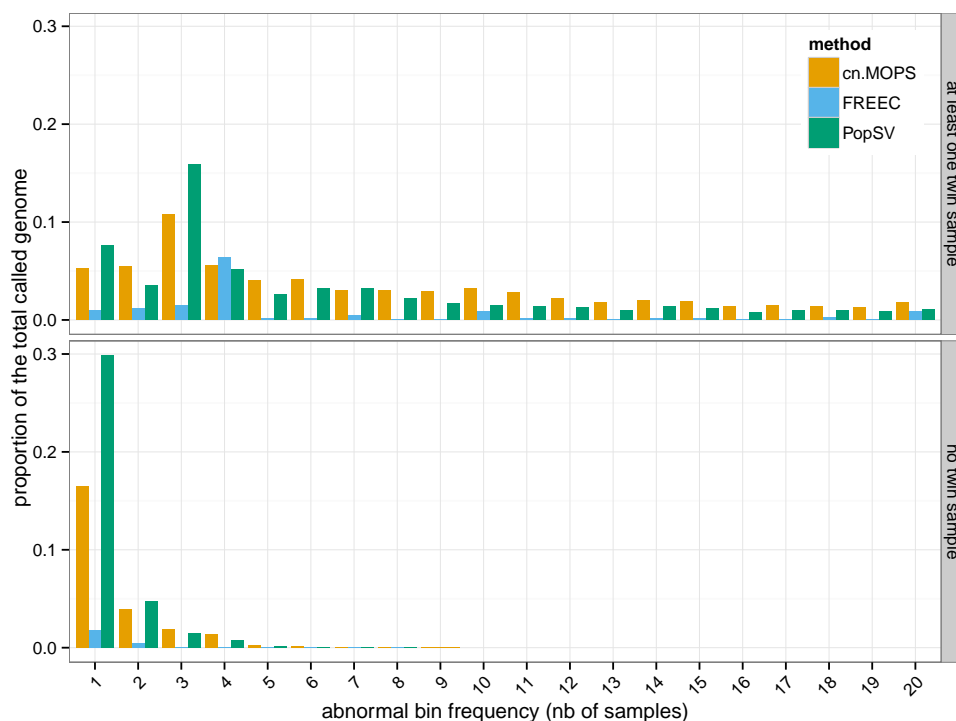


Figure S10: **Twin dataset: recurrence distribution.** The distribution of the event frequency shows a nice peak at 3-samples frequency when focusing on regions involving at least one twin (top). Using regions with no twin involved (bottom), the 3-samples peak should disappear.

| Set | Sample | Variants | | | Avg Size (Kbp) | Variants <3 Kbp | | Affected genome (Mb) | | | | |
|----------------------|--------|--------------------|------------|-----------|----------------|-----------------|------------|----------------------|-------------|------------|--------|--------|
| | | Total | Per sample | | | Proportion | Per sample | Total | Per sample | | | |
| | | | <i>WC</i> | <i>LC</i> | | | | <i>min</i> | <i>mean</i> | <i>max</i> | | |
| Renal cancer somatic | 95 | 391860 | 4124.84 | 44.40 | 58.54 | 0.48 | 1966.36 | 2455.18 | 4.16 | 232.83 | 664.86 | |
| | | <i>deletion</i> | 194181 | 2044.01 | 2.72 | 70.81 | 0.42 | 865.64 | 1695.56 | 0.01 | 136.35 | 413.66 |
| | | <i>duplication</i> | 197679 | 2080.83 | 43.68 | 46.50 | 0.53 | 1100.72 | 1464.00 | 0.12 | 96.48 | 367.53 |

Table S6: **Somatic CNVs in renal cancer dataset.** Same as Table 1 and S5. WG: whole genome; LC: low-coverage regions. *Affected genome* represents the amount of the reference genome that overlaps at least one CNV.

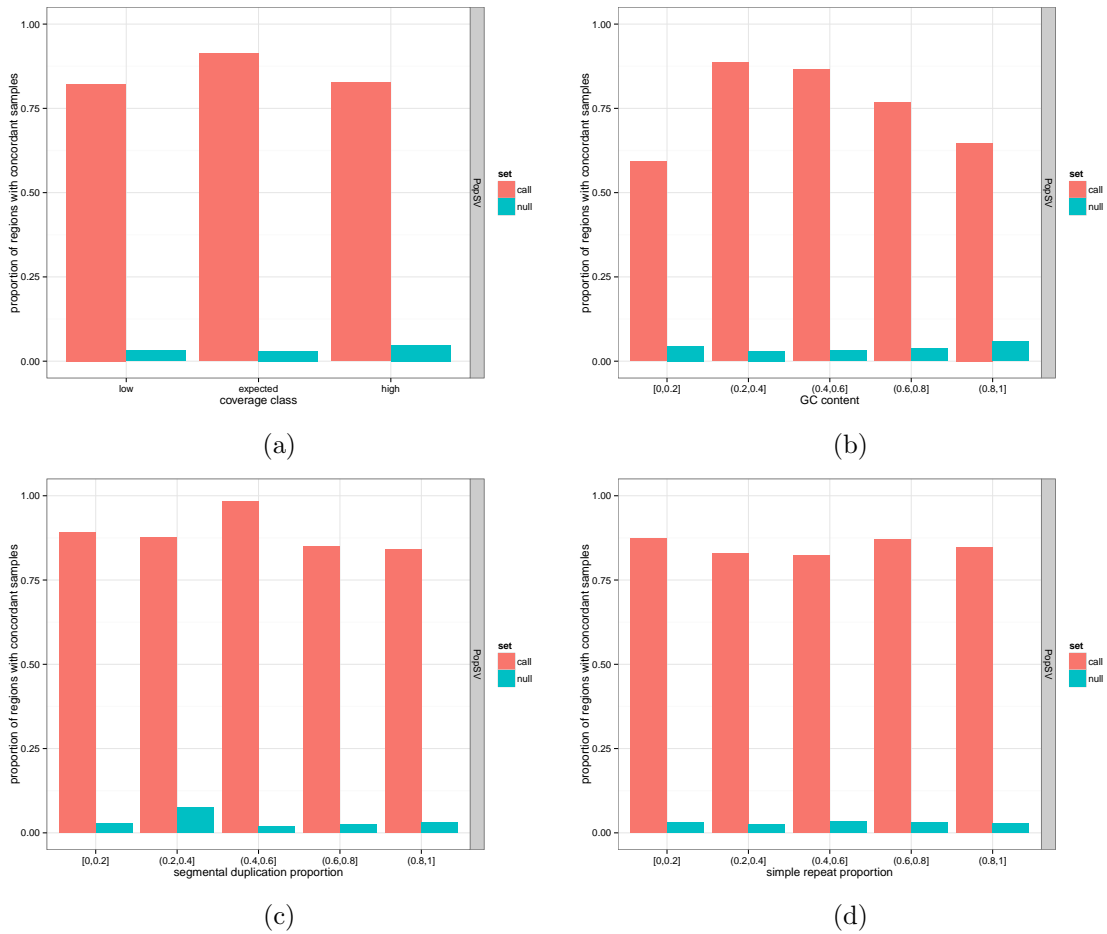


Figure S11: **Proportion of reliable bins.** Bins are grouped by coverage class (top-left), GC content (top-right), segmental duplication content (bottom-left) and simple repeat content (bottom-right). Reliable bins are defined by the concordance between twin pairs. The bars show the proportion of reliable bins. The null proportion (blue) represents the proportion expected by chance, computed by randomly selecting samples.

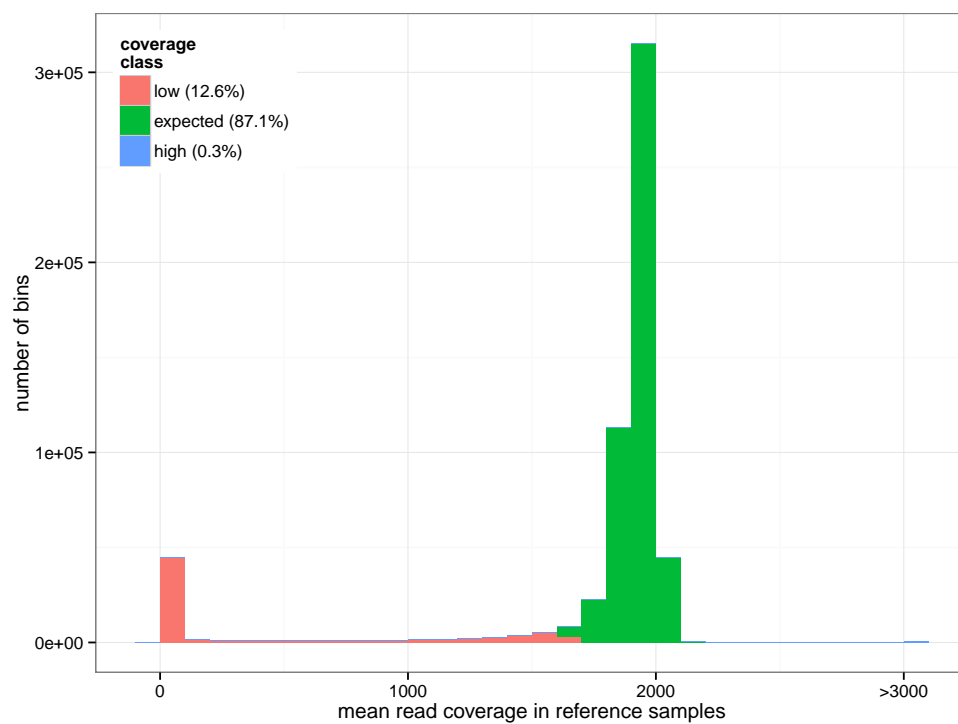


Figure S12: **Average coverage in 5 Kbp bins across reference samples in the *Twins* dataset.** We define low-mappability regions as the regions with consistently low RD across the reference samples (red).

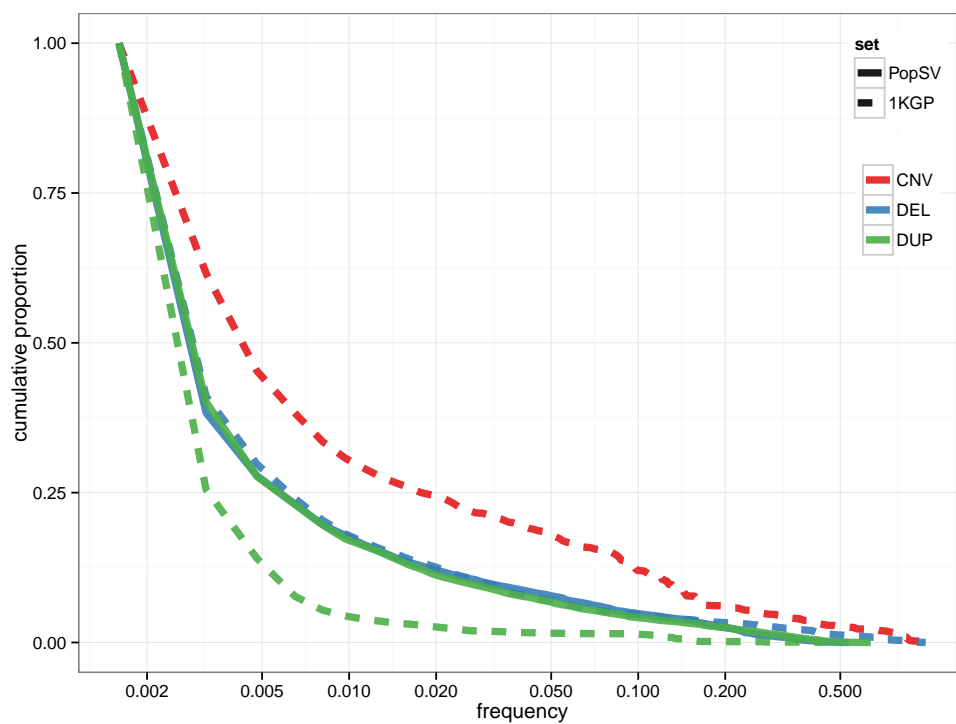


Figure S13: **Variant frequency compared to the 1000 Genomes.** The x-axis represents the frequency at which a genomic region is affected by a CNV. The y-axis represent the cumulative proportion of the affected genome.

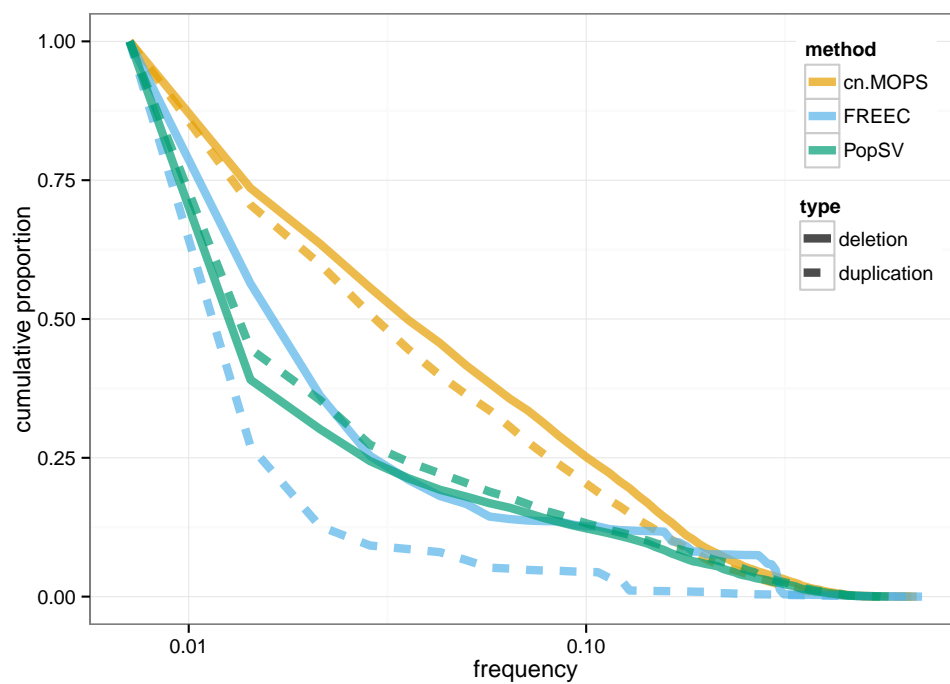


Figure S14: **Variant frequency with different methods.** The x-axis is log-scaled and represents the frequency at which a genomic region is affected by a CNV. The y-axis represent the cumulative proportion of the affected genome.

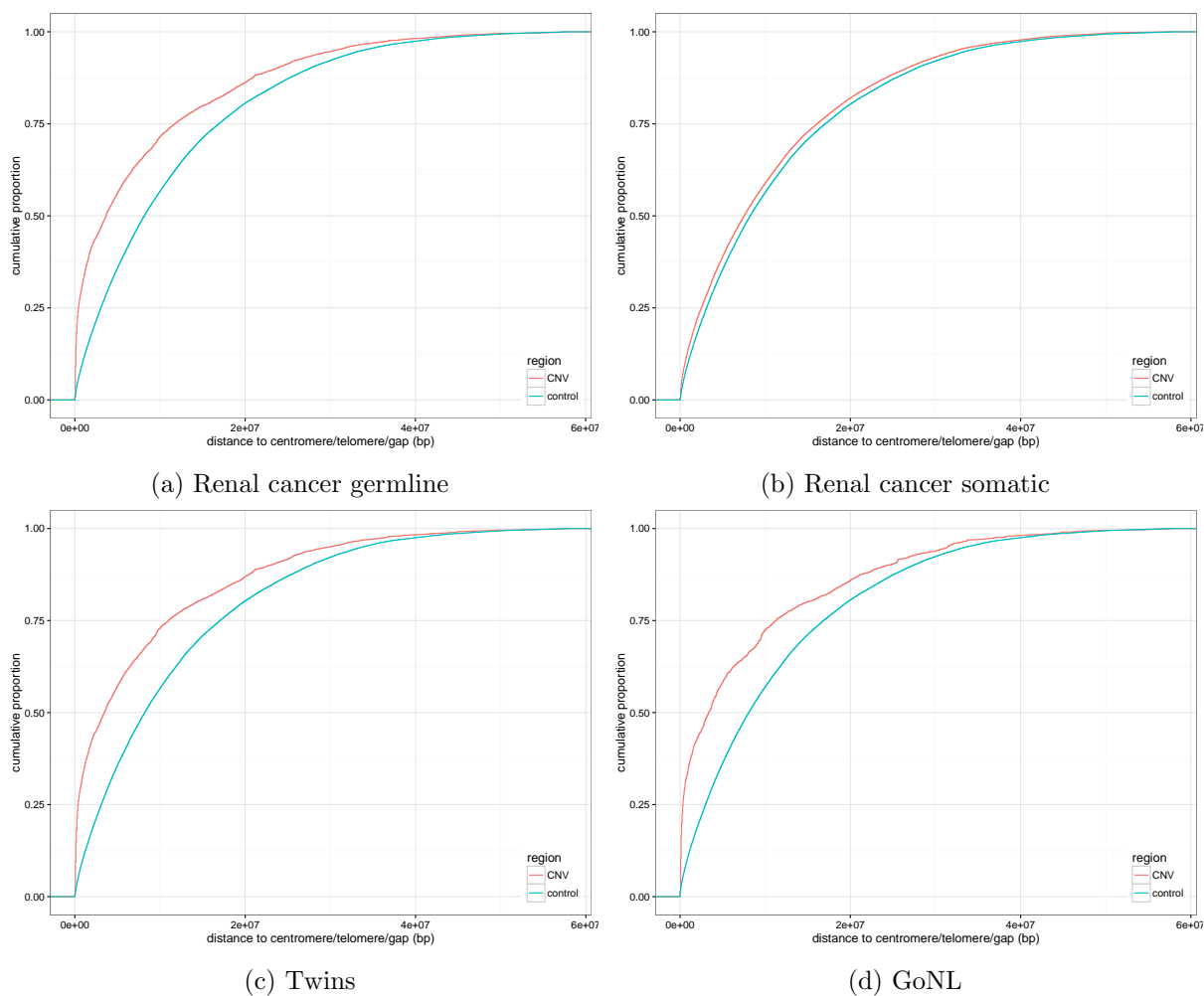


Figure S15: **Distance to a centromere, telomere or assembly gap (CTG)**. The y-axis represents the cumulative proportion of the affected genome. The *control* curve is computed from uniformly distributed genomic regions with matched size.

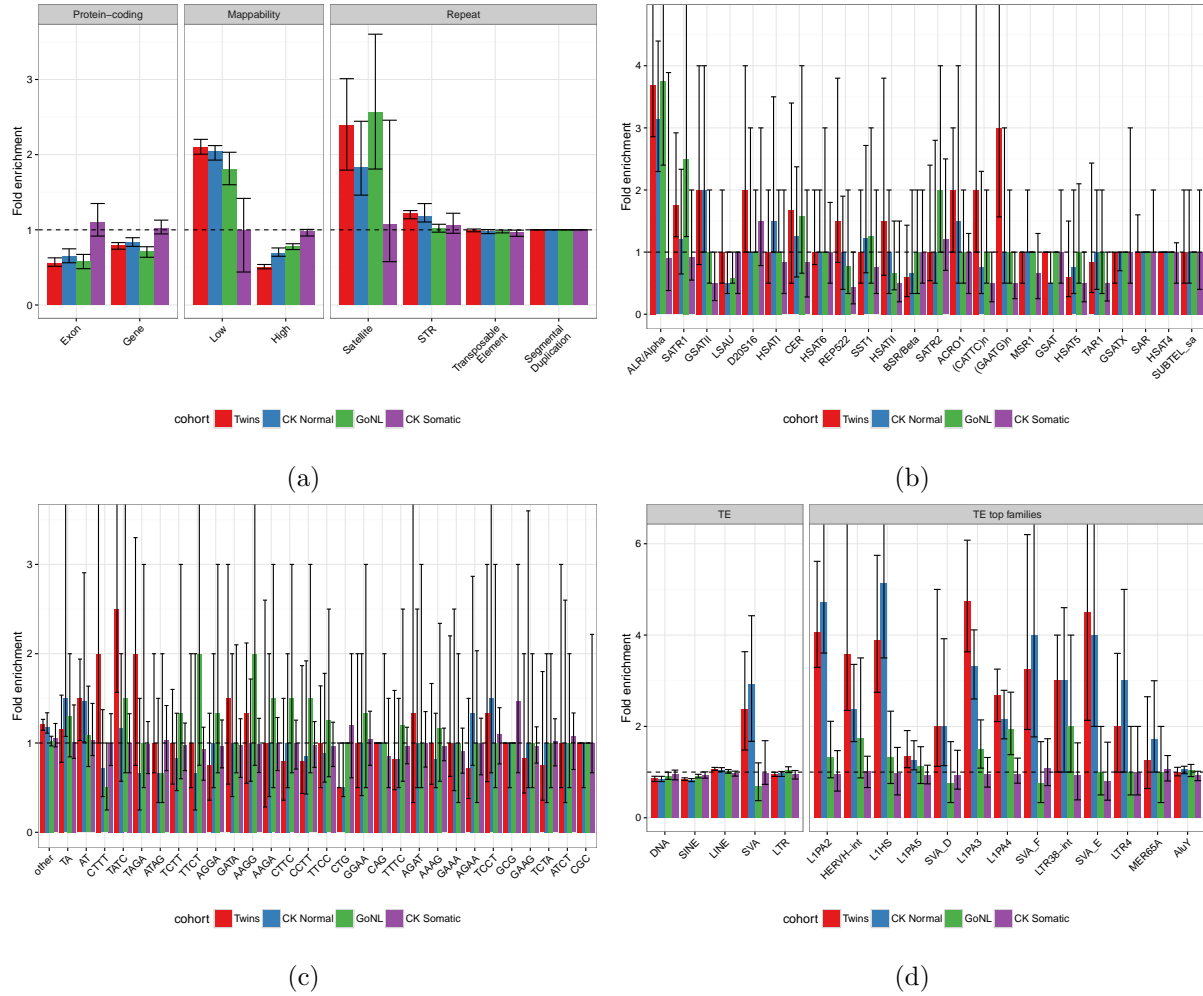


Figure S16: **CNVs enrichment after controlling for segmental duplication overlap and distance to CTG.** Enrichment of CNVs in a) different genomic features, b) satellite families, c) simple repeats, d) TE classes and top sub-families in the different cohorts (colors). Bars show the fold enrichment compared to control regions. The error bar represent 80% of the samples.

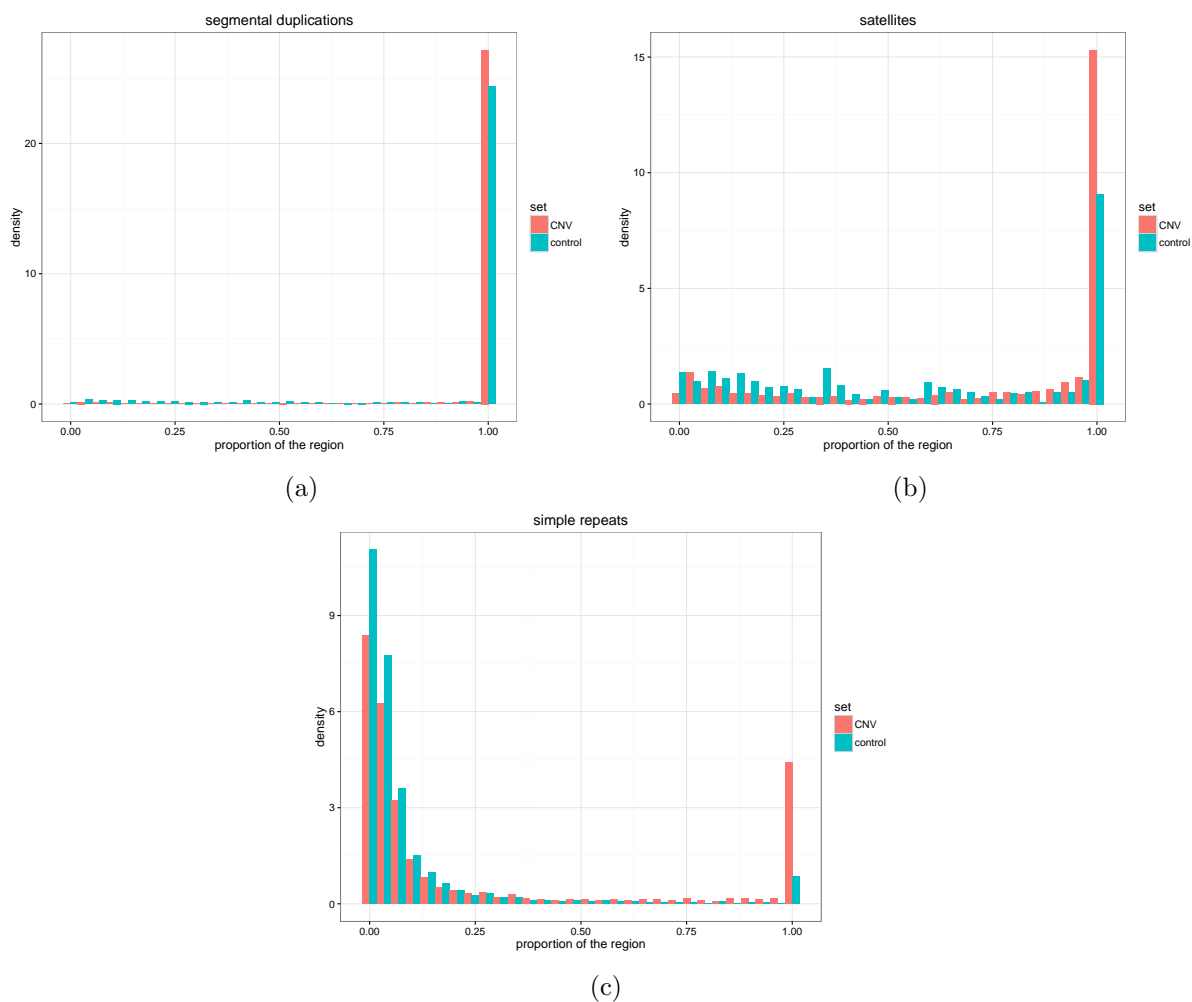


Figure S17: **Overlap between CNVs and repeats.** The histograms represent the proportion of the CNV region that overlaps a) a segmental duplication, b) satellite or c) a simple repeat, when they do overlap. The null distribution is computed from the same control regions used for the enrichment analysis (Methods).

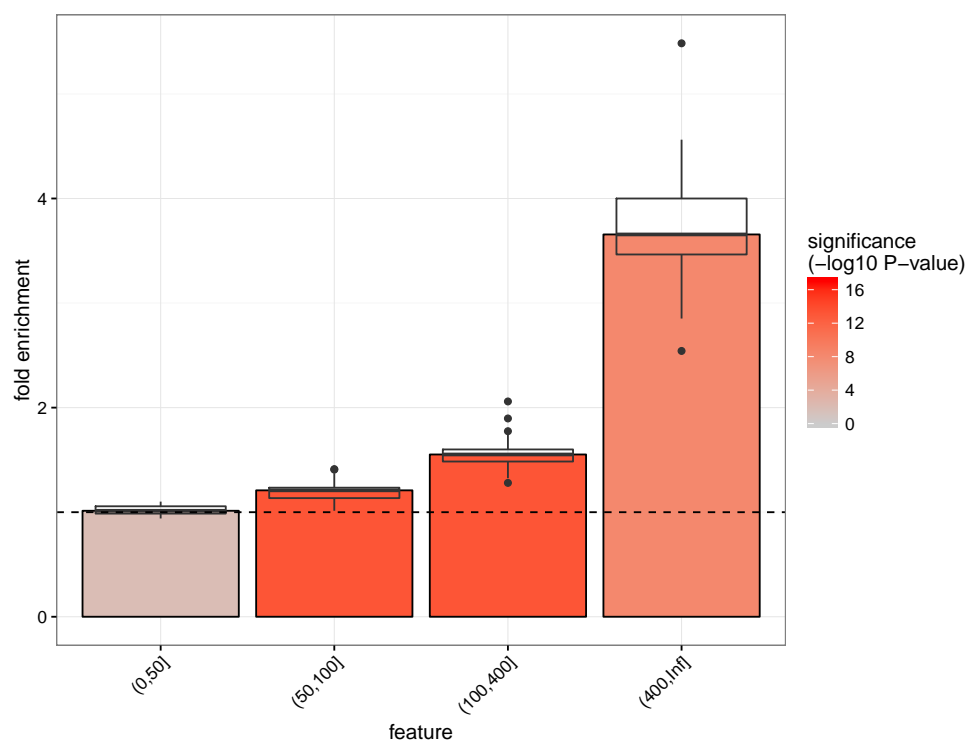


Figure S18: **Enrichment in simple repeats grouped by size.** Simple repeats are grouped by size of the annotated instance (x-axis). The fold enrichment between variant and control regions is shown in the y-axis. The boxplot show the distribution across all the samples, here from the *Twins* dataset.

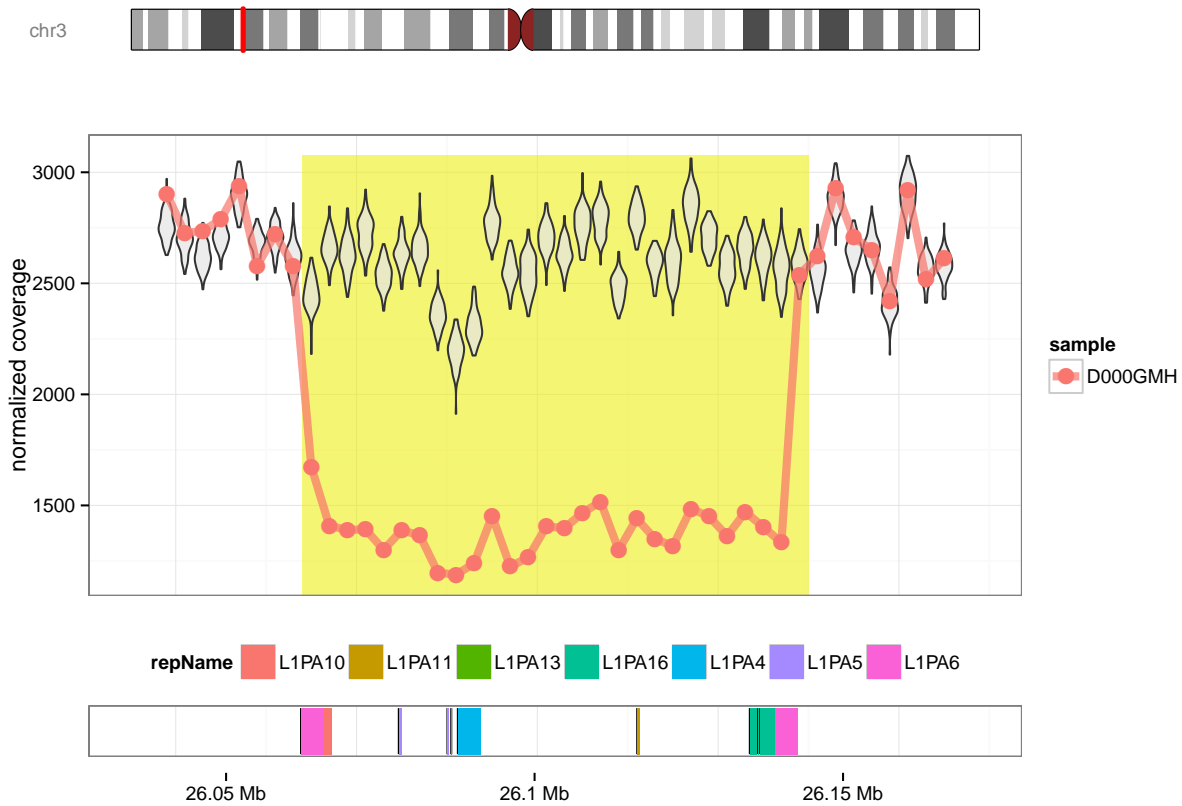


Figure S19: **Polymorphism likely caused by non-homologous allelic recombination between L1PA repeats.** Similar to Fig. 1d, violin plots represent the coverage in the reference samples, and the line and point the coverage in one sample. Here L1PA6 seems to serve as a template for a non-homologous allelic recombination resulting in a deletion.

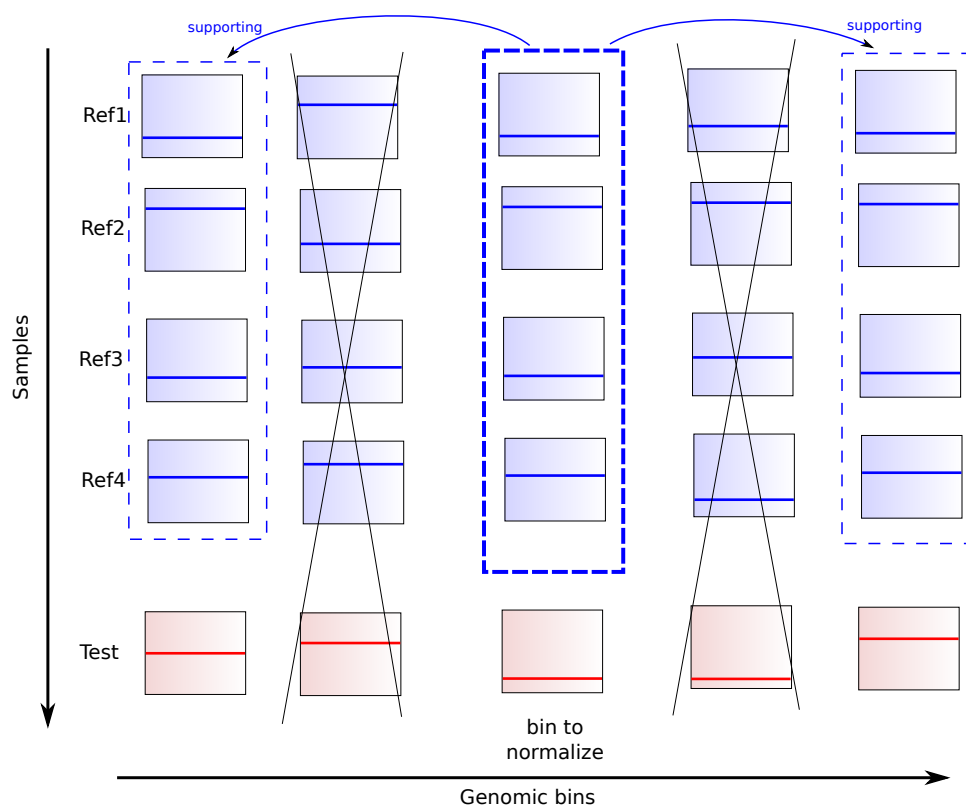


Figure S20: **Targeted normalization.** The coverage across the reference samples (blue) in the bin to normalize is used to find supporting bins across the genome. These supporting bins only are used to compute the normalization factor. The same supporting bins will be used to normalize the bin count in a test sample (red).

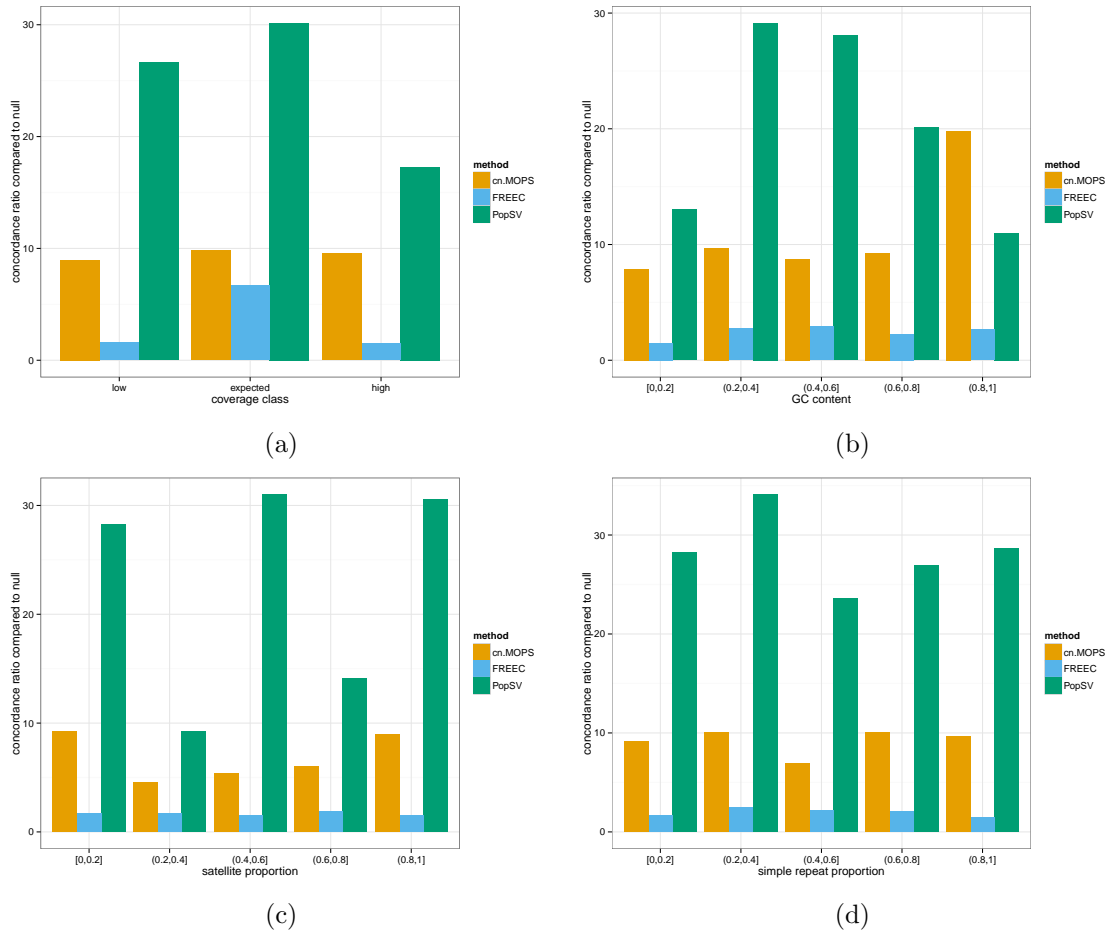


Figure S21: **Reliable bin enrichment in different methods.** Bins are grouped by coverage class (top-left), GC content (top-right), segmental duplication content (bottom-left) and simple repeat content (bottom-right). The y-axis represents the fold enrichment between the proportion of reliable bins and its expected value by chance (red versus blue in Figure S11).

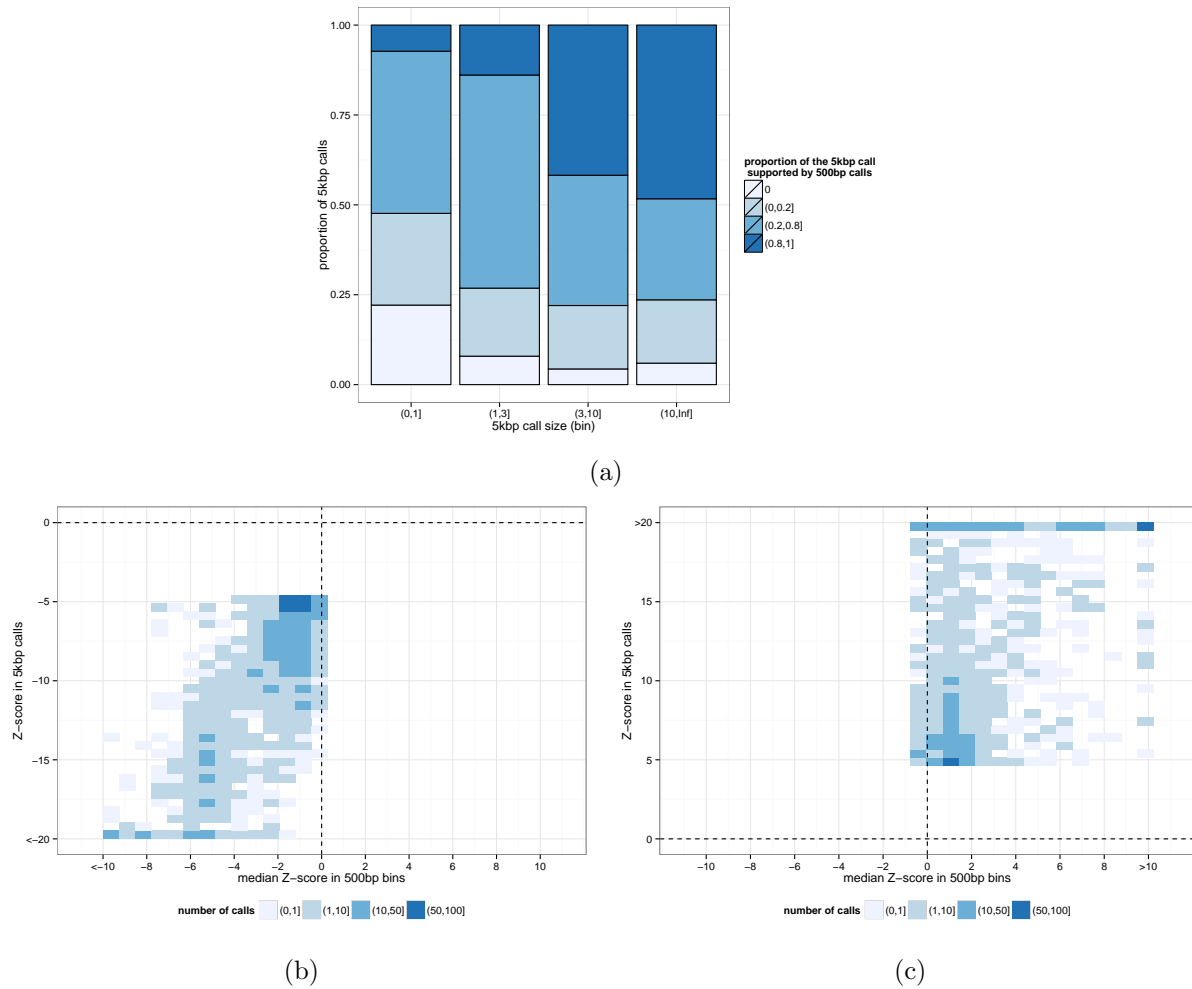


Figure S22: **5Kbp calls supported by 500bp calls.** a) 5 Kbp calls of different sizes (x-axis) are split according to the proportion of the call supported by 500 bp calls. The Z-score of 500 bp bins in 5 Kbp calls is consistent with the call for deletion b) and duplication c) signal. 5 Kbp calls with lower significance (e.g. single-bin calls) are less supported by 500 bp calls (a) but their Z-scores are in the consistent direction (b,c) although not always significant enough to be called.

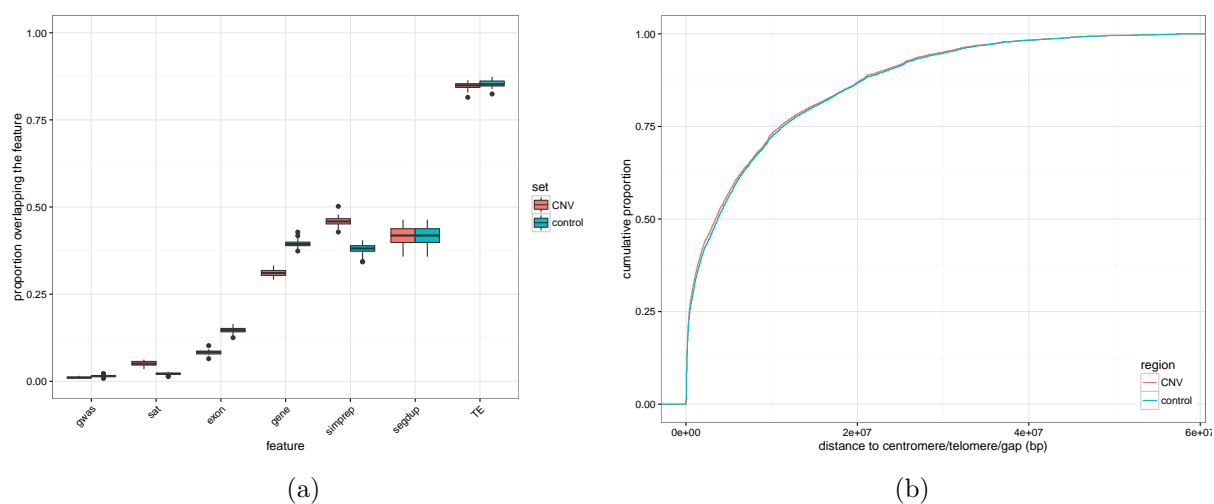


Figure S23: **Control regions quality control.** a) Control regions and CNVs have exactly the same proportion of overlap with the segmental duplications. b) When controlling for it, the distribution of the distance to a centromere, telomere or gap is very similar between CNVs and control.