







# PopDel identifies medium-size deletions simultaneously in tens of thousands of genomes

Sebastian Niehus <sup>1,2,3</sup>, Hákon Jónsson <sup>4</sup>, Janina Schönberger<sup>1,2</sup>, Eythór Björnsson<sup>4,5,6</sup>, Doruk Beyter<sup>4</sup>, Hannes P. Eggertsson <sup>4</sup>, Patrick Sulem <sup>3</sup>, Kári Stefánsson<sup>4,5</sup>, Bjarni V. Halldórsson <sup>4,7</sup> & Birte Kehr <sup>1,2,3,8</sup>✉

Thousands of genomic structural variants (SVs) segregate in the human population and can impact phenotypic traits and diseases. Their identification in whole-genome sequence data of large cohorts is a major computational challenge. Most current approaches identify SVs in single genomes and afterwards merge the identified variants into a joint call set across many genomes. We describe the approach PopDel, which directly identifies deletions of about 500 to at least 10,000 bp in length in data of many genomes jointly, eliminating the need for subsequent variant merging. PopDel scales to tens of thousands of genomes as we demonstrate in evaluations on up to 49,962 genomes. We show that PopDel reliably reports common, rare and de novo deletions. On genomes with available high-confidence reference call sets PopDel shows excellent recall and precision. Genotype inheritance patterns in up to 6794 trios indicate that genotypes predicted by PopDel are more reliable than those of previous SV callers. Furthermore, PopDel's running time is competitive with the fastest tested previous tools. The demonstrated scalability and accuracy of PopDel enables routine scans for deletions in large-scale sequencing studies.

<sup>1</sup>Regensburg Center for Interventional Immunology (RCI), Regensburg, Germany. <sup>2</sup>Berlin Institute of Health (BIH), Berlin, Germany. <sup>3</sup>Charité—Universitätsmedizin Berlin, Berlin, Germany. <sup>4</sup>deCODE genetics/Amgen Inc., Reykjavík, Iceland. <sup>5</sup>Faculty of Medicine, School of Health Sciences, University of Iceland, Reykjavík, Iceland. <sup>6</sup>Department of Internal Medicine, Landspítali—The National University Hospital of Iceland, Reykjavík, Iceland. <sup>7</sup>School of Science and Engineering, Reykjavik University, Reykjavík, Iceland. <sup>8</sup>Univeristät Regensburg, Regensburg, Germany. ✉email: [birte.kehr@ukr.de](mailto:birte.kehr@ukr.de)

Comprehensive and reliable collections of genetic variation are a foundation for research on human diversity and disease<sup>1</sup>. They facilitate a wide range of studies investigating mutation rates<sup>2–4</sup>, mutational mechanisms<sup>5–7</sup>, functional consequences of variants<sup>8–10</sup>, ancestry relationships<sup>11</sup>, disease risks<sup>12</sup>, or treatment options<sup>13</sup>. Due to increased throughput and decreased cost, whole-genome sequencing is now performed on cohorts of thousands of individuals, for example, sequencing at the population level in Iceland<sup>14</sup>, the United Kingdom<sup>15</sup>, or Crete<sup>16</sup>, as well as sequencing of large cohorts for specific diseases, such as autism<sup>17</sup> or asthma<sup>18</sup>, and in the general health research context in projects like GnomAD<sup>19</sup> or TopMed<sup>20</sup>. To create meaningful collections of genetic variation, the data from these large numbers of individuals need to be integrated. The most direct way of achieving this is done in joint variant calling approaches: detecting variant positions and inferring variant genotypes from data of many individuals together.

For single-nucleotide variants (SNVs) and small insertions/deletions (indels), joint calling has become the state of the art with tools that scale to tens of thousands of individuals<sup>21,22</sup>. For structural variants (SVs), the analysis of increasingly large numbers of individuals remains a major bioinformatic challenge<sup>23</sup>. Jointly detecting SVs in up to hundreds of individuals is a great achievement of previous projects and tools<sup>24,25</sup>. However, for larger cohorts, catalogs of SVs are generally created in a multistep approach by first analyzing the data of each individual separately or in small subsets of individuals, subsequently merging the resulting call sets and, finally, determining genotypes for all individuals on the merged call set<sup>26,27</sup>. Merging of SV call sets across individuals is often problematic and arbitrary when the same SV is detected with shifted positions in several individuals<sup>28,29</sup>. In addition, variants that are only weakly supported by the data may not be discovered using this multistep approach. Furthermore, the aligned read data is accessed at least twice in the process, for detecting and for genotyping SVs, requiring substantial computational resources. A joint SV detection approach simplifies the calling process, is computationally more efficient if accessing the large amounts of input data only once, eliminates the need for an error-prone variant merging step, and may reveal weakly supported variants if carried by several individuals as the support can be accumulated across individuals.

In this work, we overcome these limitations of current SV callers by introducing a joint calling approach, PopDel, for deletions of a few hundred up to tens of thousands of base pairs (bp) in length. We specifically designed the approach to scale to large cohorts and demonstrate that it jointly discovers SVs across tens of thousands of individuals, thereby directly creating joint call sets. Nevertheless, it can also be applied to a single genome or small numbers of genomes, where it achieves comparable or even better accuracy than widely used deletion callers. We demonstrate the high accuracy on simulated data, public benchmark data from the Genome-in-a-Bottle (GIAB) consortium, and on data of parent-child trios, which allows us to validate predicted genotypes using the genetic laws of inheritance. Our results demonstrate that the joint deletion detection approach can yield reliable call sets across very many genomes.

## Results

Our main result is a computational approach that can detect and genotype deletions in tens of thousands of genomes jointly. We have implemented the approach in the tool PopDel and evaluate PopDel's performance in comparison to previous popular tools for SV calling.

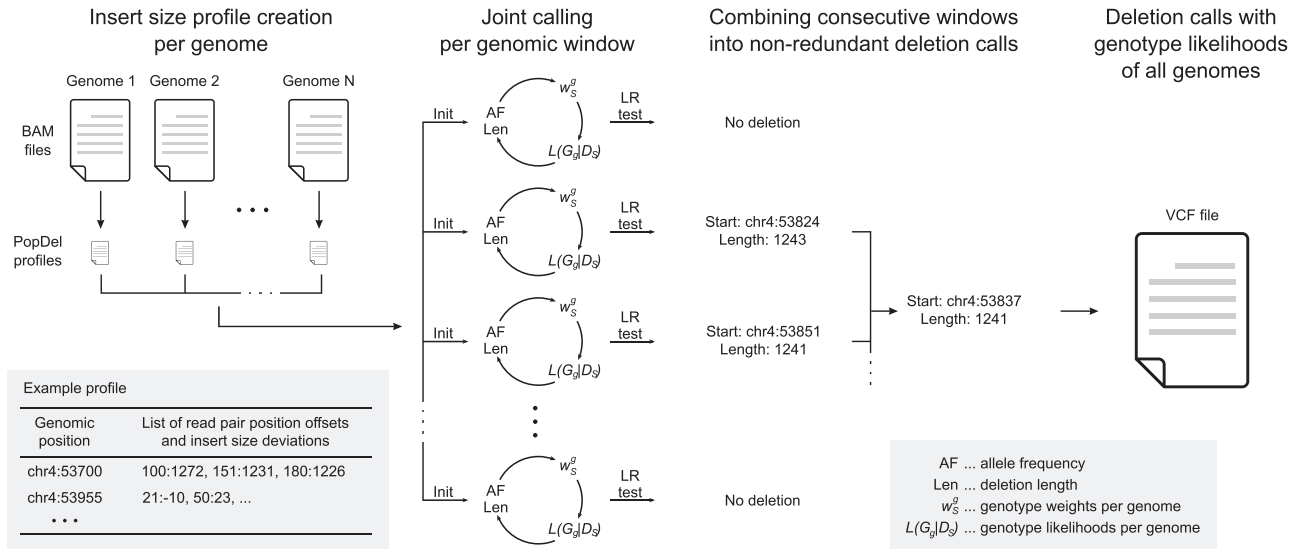
**Computational approach for joint deletion calling.** Deletions can manifest themselves in the reference alignment of short-read sequences as local drops in read depth, aberrations in the distance between the alignment of two reads in a pair, and split-aligned reads<sup>30</sup>. To detect deletions, PopDel focuses on local changes in the read pair alignment distance compared to the genome-wide distribution of read pair alignment distances. Split-aligned reads are used to infer a precise start position of detected deletions and read depth is used implicitly during genotyping.

To achieve scalability to very large numbers of individuals, PopDel has two steps: a profiling step, which reduces the aligned input sequencing read data per individual into a small read pair profile, followed by a joint calling step, which takes as input the read pair profiles and outputs deletion calls with genotypes across all individuals (Fig. 1; see “Methods”). While the two individual steps are novel, the enclosing two-step design is reminiscent of joint calling of small variants in the GATK HaplotypeCaller<sup>21</sup> and copy-number variant calling approaches that are based on read depth profiles<sup>31,32</sup>. Read pair profiles contain the sequencing experiment's distribution of read pair distances as well as alignment start positions and distances of all read pairs that match certain quality criteria (Supplementary Table 1). The joint calling step processes the profiles of all individuals together in small genomic windows (default 30 bp) to discover and genotype deletions. For all windows, likelihood ratio tests are performed to test if deletions overlap the window in any of the jointly analyzed individuals. In the likelihood computation, we use weighted genotype likelihoods to ensure that rare deletions can be found by boosting the signal in carriers and down-weighting the contribution of noncarriers dependent on the allele frequency. Finally, adjacent windows that support the same deletion are aggregated and output together with genotype likelihoods of all individuals.

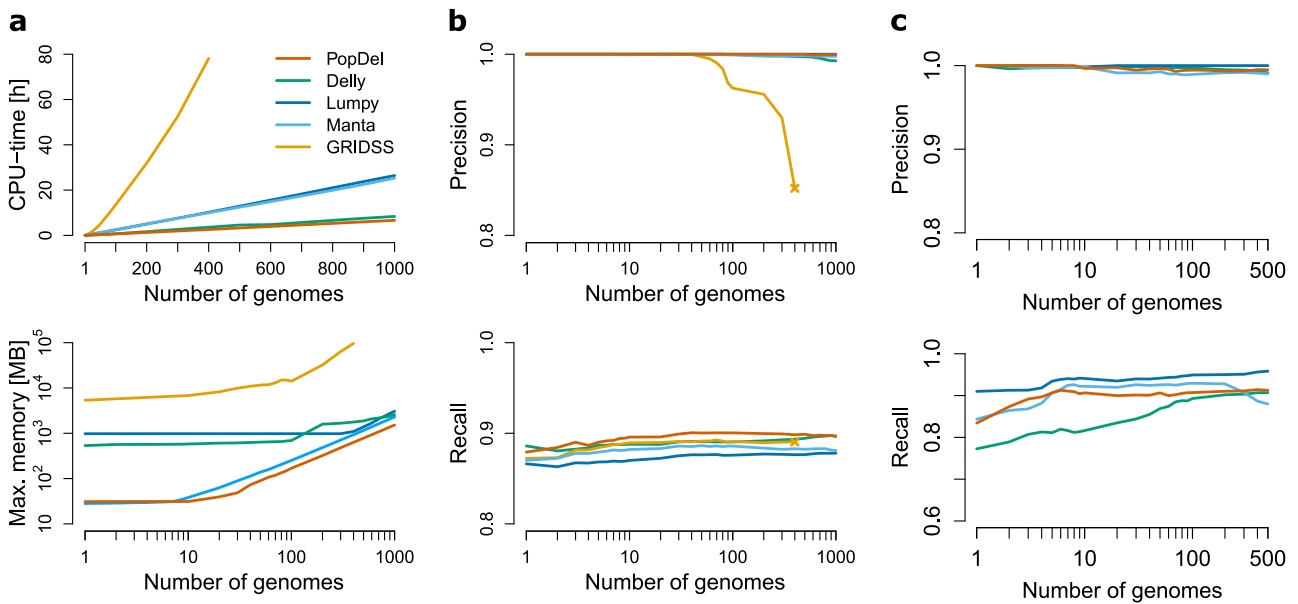
Most parameters of PopDel are calculated from the input data. The input parameters for each likelihood ratio test are iteratively estimated (Fig. 1; see “Methods”): the deletion length, allele frequency, genotype weights, and genotype likelihoods for all individuals for the three genotypes (noncarrier, heterozygous carrier, and homozygous carrier). The minimum length of deletions that can be identified with our likelihood ratio test derives from the standard deviation of the read pair distances.

**Assessment of scalability on simulated data.** For an initial assessment of PopDel's precision and recall, we simulated two cohorts of sequencing data: the first consisting of 1000 individuals carrying random deletions and the second consisting of 500 individuals with deletions reported in the 1000 Genomes Project (see “Methods”). Individuals in the first cohort carry on average 659 heterozygous and 673 homozygous deletions that are placed densely on chromosome 21, while individuals in the second cohort carry on average 167 heterozygous and 64 homozygous deletions that are more sparsely distributed on chromosomes 17–22. On these data, we compared the recall, precision, running time, and memory consumption of PopDel to that of four popular SV callers that can identify SVs jointly in a limited number of individuals or provide a pipeline for single-genome calling followed by merging and genotyping (Delly<sup>33</sup>, Lumpy<sup>34</sup> via the recommended Smoove pipeline (see URLs), Manta<sup>35</sup>, and GRIDSS<sup>36</sup>). We note that PopDel currently only reports deletions while other callers also report other types of SVs.

The precision and recall of PopDel and most other tools is high for both cohorts (Fig. 2) reflecting that the simulated data is easy to analyze. Only GRIDSS' performance in precision drops significantly with increasing numbers of individuals, which is why we excluded it from further joint calling comparisons. The recall and precision on the 1000 Genomes Project deletions



**Fig. 1 Overview of the approach implemented in PopDel.** First, the BAM file of one individual at a time is reduced into a small profile. The profiles of all individuals are processed together by sliding a window (of size 30 bp by default) over the genome and assessing the likelihood of each window to overlap a deletion in any individual. Sizes and allele frequencies of the deletions are estimated iteratively. Consecutive windows are combined into a single variant call and genotypes of all individuals are output to a VCF file. Init initialization, AF allele frequency, Len deletion length,  $w_s^g$  genotype weights per genome,  $L(G_j|D_j)$  genotype likelihoods per genome, LR likelihood ratio.



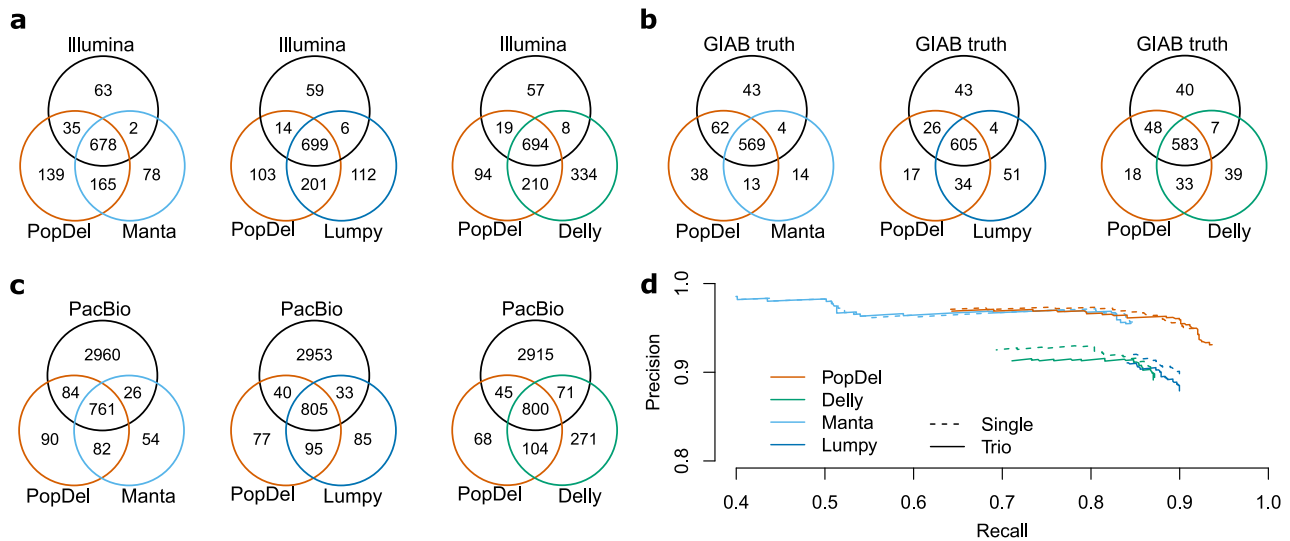
**Fig. 2 Performance of our approach on simulated data.** We compared PopDel to Delly, Lumpy (via Smoove), Manta, and GRIDSS for increasing numbers of genomes. **a** Running time and memory on data with deletions randomly placed on human chromosome 21. **b** Recall and precision on data with random deletions. **c** Recall and precision with deletions reported by the 1000 genomes project simulated on chromosomes 17–22.

fluctuate more due to the smaller number of simulated deletions per individual. Nevertheless, all tools show a clear trend of higher recall with increasing numbers of individuals on the 1000 Genomes Project deletions. This suggests that these deletions tend to be more difficult to identify than random deletions and, here, deletion callers benefit from integrating data of several individuals. Surprisingly, the recall of Lumpy is higher for the 1000 Genomes Project deletions than for the randomly generated deletions.

PopDel is the only tool that we could run on all 1000 individuals using our system’s default settings. For all other tools, we had to increase the ulimit (maximum number of open file handles) in

order to complete the calling on >500 individuals jointly, indicating that the other tools were primarily designed to jointly analyze smaller numbers of individuals. With a running time of 397 min and a peak memory of 1.5 GB for profiling and joint calling on the first cohort of all 1000 individuals, PopDel is the fastest tool and among the tools that require the least memory (Fig. 2).

**Comparison to reference deletion sets from the GIAB consortium.** Next, we assessed PopDel’s performance compared to Delly, Lumpy (via Smoove), and Manta on short-read whole-genome sequencing data of the well-studied HG001/NA12878 and HG002/NA24385 genomes and their parent’s genomes (see



**Fig. 3 Accuracy of our approach on Genome-in-a-Bottle (GIAB) benchmark genomes.** We compared PopDel on HG001 and HG002 to Manta, Lumpy (via Smoove), and Delly using a minimum reciprocal overlap criterium of 50%. Call set overlap for HG001 with **a** the Genome-in-a-Bottle short-read reference set of deletions and **c** a deletion set called from PacBio long read data. **b** Call set overlap for HG002 with deletions from the Genome-in-a-Bottle preliminary variant set. **d** Precision-recall curve when calling deletions from HG002 alone (Single) or jointly with the parental genomes (Trio) for different genotype quality thresholds.

“Methods”). For this assessment, we used reference sets of deletion calls prepared by the GIAB consortium<sup>37</sup>: the short-read-based reference set and a set of deletions called from PacBio long read data for HG001, and the preliminary variant set for HG002 (see “Data availability”).

PopDel is competitive with the three other tools on the data of HG001 and HG002 (Fig. 3 and Supplementary Figs. 1 and 2). All tools succeed to identify the majority of deletions reported in the short-read HG001 reference set (661/778, 85.0%) with PopDel identifying marginally more deletions (713, 91.6%) than Lumpy (705, 90.6%), Delly (702, 90.2%), and Manta (680, 87.4%). The fraction of PacBio deletions identified by all three tools is much lower (731/3,831, 19.1%). This is expected as the long PacBio reads reveal variants involving repeats that are invisible or hard to detect in short-read data. PopDel identifies a similar number of PacBio deletions (847, 22.1%) as Lumpy (838, 21.9%), Delly (871, 22.7%), and Manta (786, 20.5%). Including those deletions that are not part of the two HG001 reference sets, PopDel reports fewer deletions than Delly, a similar number of deletions as Lumpy, and more deletions than Manta. As the NA12878 reference sets do not claim to be complete, the additional deletions can either be true or false positives.

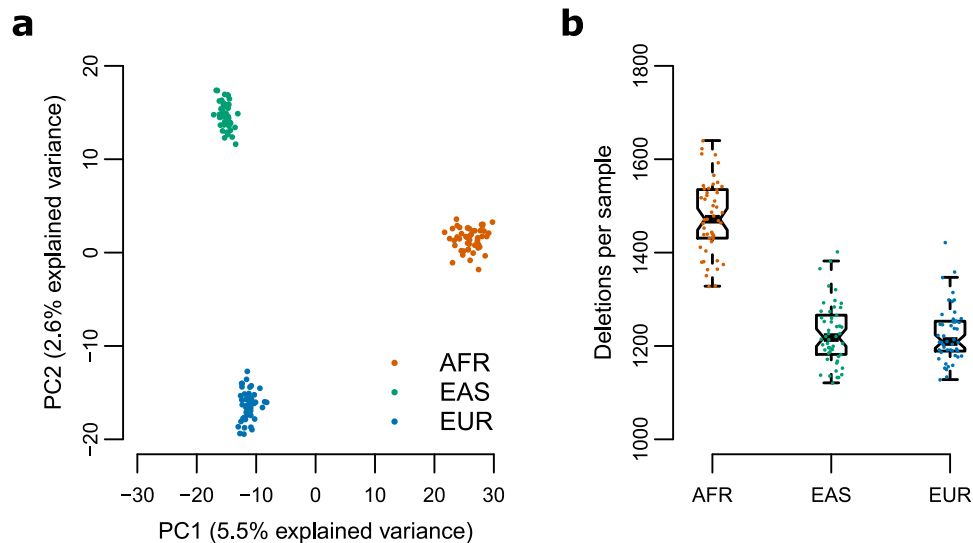
The preliminary HG002 deletion set has been released as the first reference set that is near complete within defined high-confidence regions of the genome and, hence, allows us to evaluate the precision and recall of PopDel compared to the other tools (Fig. 3). The call set of PopDel, when run on the data of HG002 and its parental genomes jointly, comprises 629 of the 678 reference deletions (93.7% recall) with a precision of 93.1% resulting in an  $F_1$  score of 93.4%. Only Manta’s precision is higher (95.5%) at the cost of a much lower recall (84.5%), resulting in an  $F_1$  score of only 89.7%, which is similar to that of Delly (88.1%) and Lumpy (88.9%). Thus, PopDel outperforms the other tools by 3.7 to 5.3 percentage points in terms of  $F_1$  score. PopDel’s recall is higher when adding the parental genomes compared to running it on the data of HG002 alone (91.6%), indicating a benefit of joint calling. While Manta’s performance is hardly affected by adding the parental genomes to the calling, Delly and Lumpy lose precision on our data without the recall being affected.

### Analysis of population structure based on deletions in the Polaris Diversity cohort.

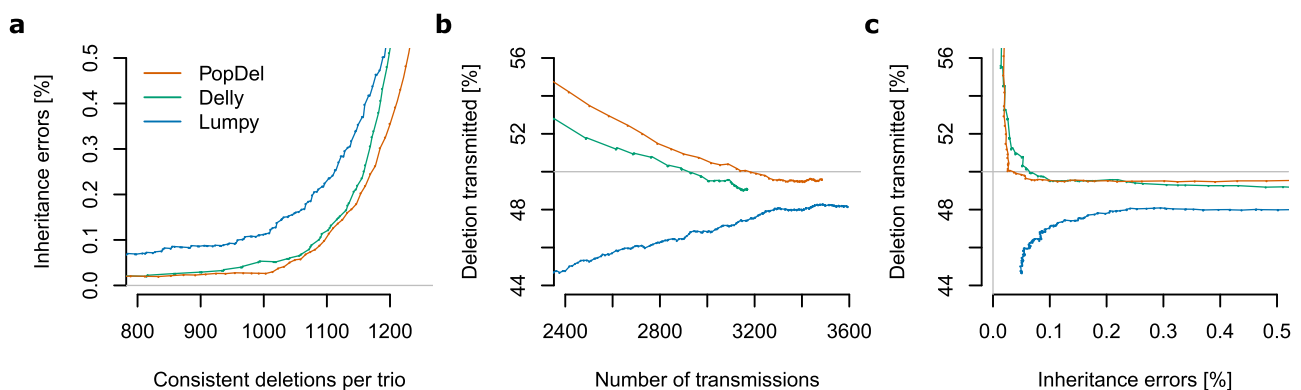
We applied PopDel to the 150 genomes in the Polaris HiSeq X Diversity cohort (BioProject accession PRJEB20654) to evaluate if the deletions identified by PopDel reflect population structure (see “Methods”). The cohort consists of three continental groups: Africans, East Asians, and Europeans. PopDel identifies an average of 969 heterozygous and 340 homozygous deletions per individual overall. Consistent with previous studies, Africans carry significantly ( $P$  value  $< 2.2 \times 10^{-16}$ , two-sided  $t$  test) more deletions than Europeans and East Asians (Fig. 4a). Principal component analysis of PopDel’s deletion calls shows a clear separation between the three continental groups (Fig. 4b) mirroring the well-known clustering resulting from small variants<sup>24,38</sup>. In particular, the first principal component separates the African genomes from the other continental groups, while the second principal component additionally pulls apart the European and East Asian genomes. These findings indicate that the deletions detected and genotyped by PopDel well reflect the biological differences between the continental groups. Similar results were obtained for Delly and Lumpy via Smoove (Supplementary Fig. 3).

### Evaluation of genotyping using data of 49 Polaris trios.

By combining the Polaris Diversity cohort with the Polaris HiSeq X Kids cohort (BioProject accession PRJEB25009), we obtain a set of 49 trios that allow a thorough evaluation of the genotype predictions. After running PopDel, Delly, and Lumpy via Smoove (see “Methods”), we analyzed inheritance patterns of deletions and their genotypes in the 49 trios. In particular, we calculated the Mendelian inheritance error rate and transmission rates for each tool (see “Methods”). The Mendelian inheritance error rate effectively assesses the genotyping of common variants. The transmission rate is also meaningful for rare variants measuring how often a deletion allele is inherited from a heterozygous parent (Supplementary Fig. 4), in particular when it is calculated for deletions unique to one trio and the second parent being a noncarrier. As we noted an overabundance of heterozygous deletions in all call sets, we removed deletions that are not in Hardy–Weinberg equilibrium ( $P$  value 0.01) before all other calculations.



**Fig. 4 Analysis of PopDel's deletion calls on the Polaris Diversity cohort.** The Polaris Diversity cohort consists of 50 individuals from three continental groups each. **a** Principal component analysis. **b** Number of deletions per genome. Africans carry significantly ( $P$  value  $<2.2 \times 10^{-16}$ , two-sided  $t$  test) more deletions than Europeans and East Asians. Each point represents an individual. Center line denotes median. Boxes limit upper and lower quartiles. 1.5x interquartile ranges are given as whiskers. AFR, African; EAS, East Asian; EUR, European.



**Fig. 5 Assessment of genotype inheritance patterns in the Polaris Kids cohort.** We filtered call sets for Hardy-Weinberg equilibrium ( $P$  value 0.01) and varying genotype quality thresholds. **a** Mendelian inheritance error rate by the number of deletion sites per trio that are consistent with Mendelian inheritance. **b** Percentage of transmitted deletion alleles for all deletions unique to one trio and one parent by the number of possible transmissions, that is, the number of deletions unique to one trio and one parent. **c** Percentage of transmitted deletion alleles by Mendelian inheritance error rate. Gray lines indicate the ideal values of the Mendelian inheritance error rate and transmission rate.

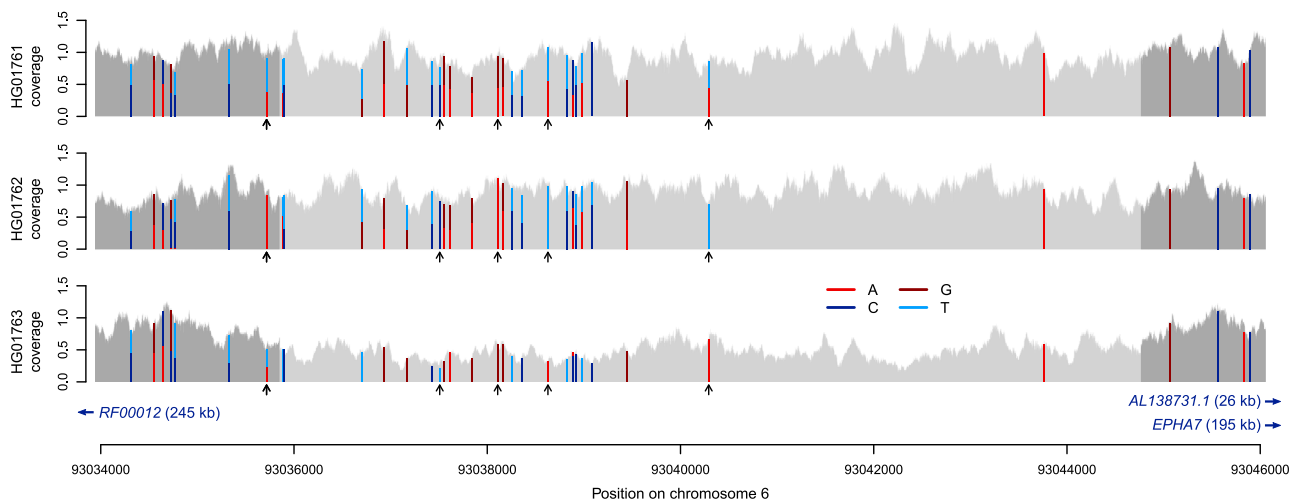
The genotyping of PopDel is very well calibrated in our comparison to Delly and Lumpy (Fig. 5). The deletions called by all three tools can be filtered to Mendelian inheritance error rates below 0.1% using reported genotype quality values. Notably, PopDel reports a larger number of deletions consistent with Mendelian inheritance than Delly and Lumpy when filtering to any Mendelian inheritance error rate. For example, PopDel reports an average of 1177 consistent deletions per trio compared to 1161 (Delly) and 1128 (Lumpy) when filtering to an error rate just below 0.3%. Similarly, PopDel reports more deletions unique to one trio than Delly when filtering by genotype quality to any given transmission rate, for example, 3167 PopDel deletions compared to 2935 Delly deletions when filtering to the best transmission rate closest to 50%. Furthermore, PopDel has a lower Mendelian inheritance error rate than Delly when filtering to the expected transmission rate of 50%. Lumpy's transmission rate shows a pattern that may be indicative of false positives (Supplementary Note 2) and never reaches 50%.

We further assessed the genotyping performance based on the rate of transmission from parents to children including scenarios where more than one parent is heterozygous and the deletion was found in several trios. For each tool, we chose genotype quality thresholds that filter the deletions to a Mendelian inheritance error rate just below 0.3% (Table 1). Using these filters, the rate of deletions transmitted from parent to child that are private to a single trio is not significantly different from 50% for PopDel and Delly (two-sided binomial test,  $P$  value threshold 0.05). When considering deletions found in several trios where only one parent is heterozygous, the transmission rate of PopDel is not significantly different from the expected 50% (two-sided binomial test,  $P$  value threshold 0.05), whereas it is different for Delly when one parent is a homozygous carrier. When both parents are heterozygous, all three tools show an overabundance of heterozygous calls in the child, indicating that this is the most challenging configuration for genotyping. The transmission rate of Lumpy is significantly different from the expected value for all considered genotype configurations. This analysis suggests that

**Table 1 Transmission of deletions in the 49 Polaris trios after filtering for Hardy-Weinberg equilibrium ( $P$  value 0.01) and genotype quality (GQ) to a Mendelian inheritance error rate just below 0.3%.**

Genotype of			PopDel (GQ = 26)	Delly (GQ = 28)	Lumpy (GQ = 78)			
Parent 1	Parent 2	Child						
Deletions present in any number of trios								
0/0	0/1	0/0	17,153	<u>50.02%</u>	16,584	<u>49.57%</u>	17,376	51.06%
		0/1	17,135	<u>49.97%</u>	16,863	<u>50.41%</u>	16,652	48.94%
		1/1	2	0.01%	8	0.02%	0	0.00%
0/1	0/1	0/0	2057	<u>24.47%</u>	2189	24.10%	2153	26.10%
		0/1	4474	53.22%	4961	54.62%	4479	54.30%
		1/1	1875	22.31%	1932	21.27%	1616	19.59%
0/1	1/1	0/0	1	0.02%	8	0.12%	2	0.04%
		0/1	3080	<u>50.80%</u>	3392	51.94%	2853	51.92%
		1/1	2982	<u>49.18%</u>	3130	47.93%	2640	48.04%
Deletions present in only a single trio								
0/0	0/1	0/0	1714	<u>50.50%</u>	1573	<u>50.63%</u>	1714	51.91%
		0/1	1680	<u>49.50%</u>	1534	<u>49.37%</u>	1588	48.09%
		1/1	0	0.00%	0	0.00%	0	0.00%

Underlined percentages are not significantly less or greater than the expected value 50% or 25%, respectively, according to a two-sided binomial test with a  $P$  value threshold of 0.05.



**Fig. 6 De novo deletion identified by PopDel in one individual of the Polaris Kids cohort.** The child (HG01763) is a carrier of a heterozygous deletion (indicated by the light gray interval) that is not present in either parent. Black arrows mark SNPs that allow phasing of the haplotypes in the child.

PopDel's genotyping accuracy is superior to that of Delly and Lumpy.

**Application to population-scale data from Iceland.** We applied PopDel to whole-genome data of 49,962 Icelanders, including 6794 parent-offspring trios (see "Methods" and Supplementary Note 2). The average number of deletions PopDel reports per Icelandic autosomes is 1504 heterozygous and 209 homozygous deletions (genotype quality threshold 25). The Mendelian inheritance error rate in the 6794 trios is 1.4% (1963 consistent deletions on average per trio). The transmission rate for 4256 deletions unique to a single trio is 49.2%, which is again not significantly different from the expected 50% (two-sided binomial test,  $P$  value 0.32). This implies that the majority of errors appear as common deletions shared by several individuals.

**Identification of a de novo deletion in the Polaris data.** We searched for de novo deletions (see "Methods") in the Polaris

Kids cohort and identified 12 candidate de novo variants. Manual inspection suggests that three of them are true de novo events: a 8901 bp deletion at chr6:93,035,858–93,044,759 in the Spanish female HG01763 (Fig. 6), an exonic 984 bp deletion at chr6:27,132,732–27,133,716 in the Spanish male HG01683 in the *H2BC11* gene, and an exonic 769 bp deletion at chr7:105,505,500–105,506,269 in the Chinese female HG00615 in the *PUS7* gene. The 8901 bp deletion in HG01763 is flanked and overlapped by SNVs that allow us to phase and confirm the de novo event. An SNV that overlaps with read pairs supporting the deletion indicates that the deletion haplotype was inherited from the mother (HG01762). Further evidence for this to be a true de novo event is given by 25 SNVs within the deletion that confirm the child to carry a single haplotype where both parents are heterozygous. All three individuals are heterozygous for numerous SNVs upstream of the deletion. Four of the SNVs within the deletion confirm that the event happened on a maternal haplotype. Given that this deletion is intergenic and HG01763 is part of a cohort of healthy individuals<sup>1</sup>, we expect the de novo deletion not to be of medical relevance. The closest transcript annotations

**Table 2 Running times (wall clock time, CPU hours in parentheses) on NA12878 and the Polaris Diversity cohort.**

	NA12878 (single individual)	Polaris Diversity cohort (150 individuals)
PopDel	0:25 (0:58)	56:17 (111:12)
Delly	1:42 (1:40)	389:43 (371:27) <sup>a</sup>
Lumpy (via Smoove)	0:18 (0:30)	87:58 (103:21) <sup>a</sup>
Manta	7:09 (7:09)	—

Note that Delly, Lumpy, and Manta report other types of SVs apart from deletions. These SVs were excluded after single-sample calling before merging to reduce running times.  
<sup>a</sup>Single-sample calling with subsequent variant merging and sample-wise genotyping.

in Gencode v29<sup>39</sup> are the lncRNA *AL138731.1* at a distance of 25.6 kb and the *EPHA7* gene at a distance of 195.3 kb.

**Running times on public benchmarking data.** We assessed the running time of PopDel compared to Delly and Lumpy via Smoove on the data of the NA12878 genome and the 150 genomes in the Polaris Diversity cohort (see “Methods” and Table 2). With a total wall clock running time of ~25 min for profile creation and deletion calling of the NA12878 genome, PopDel is almost as fast as Lumpy via the Smoove pipeline and four times faster than Delly. A similar behavior can be observed on the 150 genomes of the Polaris Diversity cohort confirming scalability: PopDel completes deletion calling within <2 days and 9 h, which is similar to the running time of Lumpy via Smoove (3 days and 16 h) and several times faster than Delly (16 days, 6 h). PopDel can be trivially parallelized by creating profiles of different individuals in parallel and splitting the joint calling by genomic region.

## Discussion

Identification and genotyping of structural variation in large sequencing cohorts is a major computational challenge. To enable the joint analysis of the increasingly large cohorts that are being sequenced, we developed a deletion calling approach implemented in the tool PopDel. Compared to existing approaches, the joint calling approach in PopDel greatly simplifies the analysis workflow, shows comparable if not better accuracy, and has a competitive running time. PopDel scales to very large cohorts as our tests on population-scale data from Iceland substantiate.

PopDel has high accuracy independent of the number of jointly analyzed individuals and across the deletion allele frequency spectrum. On data of a single genome, PopDel shows higher recall than other tools at high precision. On the Polaris Diversity cohort, the deletions called by PopDel recapitulate previous population genetic results showing that Africans carry on average more deletions than other continental groups and confirming that joint calling can be used to identify population structure. On the Polaris Kids cohort, PopDel identifies more deletions at a better transmission rate compared to other tools and reports a de novo deletion of about 9 kb. On Icelandic data, PopDel identifies deletions jointly in almost 50,000 genomes maintaining an excellent transmission rate for rare variants. All results confirm that the joint calling approach in PopDel is accurate across the allele frequency spectrum and the number of individuals analyzed.

The de novo deletion in the Polaris Kids cohort together with the good transmission rate of rare variants in a large number of Icelandic genomes demonstrates that PopDel’s joint calling approach provides a basis for studying rarely observed de novo deletion events. A previous study<sup>40</sup> verified in 258 healthy trios seven de novo deletions that fall into the size range addressed by

us. Given their rate of de novo deletions, we expect to observe 1.33 de novo deletions of medium size in the 49 Polaris trios. This is well in line with our finding of three candidate de novo events including one that we could confirm based on nearby SNVs.

When we tested an early version of PopDel on a selected 54 kb region covering the *LDLR* gene in 43,202 Icelanders, we identified a previously unknown 2.5 kb deletion in three closely related Icelanders shown to affect LDL levels<sup>41</sup>. This finding shows that PopDel is able to identify variants of biomedical interest even if they are present at a very low allele frequency in a population-scale cohort, and showcases the importance of SVs in human health.

PopDel consists of two steps: creation of read pair profiles per individual and joint deletion calling. The computational advantage of this two-step design is that the large input BAM files containing aligned read data need to be processed only once. The joint calling step takes the information needed for deletion detection and genotyping from the small read pair profiles. This implies that additional genomes, the  $N + 1$ st genome, can be added to the analysis without the need to access all input BAM files reducing the computational burden considerably. PopDel is currently limited to deletions. However, the two-step design and the likelihood ratio test can be generalized to junctions of other types of SVs. We are aiming for extending our approach accordingly in the near future.

## Methods

**Read pair profile creation.** PopDel reduces a coordinate-sorted BAM file of each sample into a read pair profile in a custom binary format (Supplementary Note 1). This profile stores positions and insert sizes of read pairs that align confidently (Supplementary Table 1) to the reference genome. In addition, the profile file contains meta information, including distribution of insert sizes across the sample and an index, which allows for jumping to genomic positions in the profile. We define the insert size as the distance between the leftmost alignment position of the forward read to the rightmost alignment position of the reverse read in the pair extended by any clipped bases (Supplementary Fig. 5). The null distribution of insert sizes is estimated by sampling the BAM file using pre-defined but user-configurable genomic regions with good mappability (Supplementary Note 1). If more than one library has been sequenced for a sample, PopDel writes separate profile data per read group to the profile file. An excerpt of an example profile is shown in Supplementary Table 2. The profiling vastly reduces I/O during joint calling as the size of the profiles is on average only 1.76% of the original BAM file size (Supplementary Figure 6).

**Likelihood ratio test for joint deletion calling.** For a given genomic window, our likelihood ratio test compares the relative likelihood that a deletion of a particular length  $l$  overlaps the window, against the relative likelihood of observing the reference haplotype:

$$\Lambda = \frac{\mathcal{L}(\text{no del})}{\mathcal{L}(\text{del of length } l)} \quad (1)$$

Our null hypothesis is that the data observed in a window is drawn from the reference model (numerator) rather than the deletion model (denominator). We reject the null hypothesis in PopDel using a cutoff for  $\Lambda$  calculated from  $-2 \log \Lambda \sim \chi^2$  with a  $P$  value threshold of 0.01 (one-tailed) and 1 degree of freedom in order to decide if the window overlaps a deletion of length  $l$ .

Let  $S \in \mathcal{S}$  be a single sample from the set of all samples  $\mathcal{S}$  and let  $I^S$  be the list of insert sizes for all the read pairs of  $S$  overlapping the given window (Supplementary Fig. 5). Furthermore, let  $\Delta^S = (i - \mu_S | i \in I^S)$  be the deviations of the insert sizes from the mean insert size of the sample  $\mu_S$ . We assume independence of samples and calculate the relative likelihood of the reference model as the product of the samples’ likelihoods  $\mathcal{L}(G_0 | \Delta^S)$  for the reference genotype  $G_0$

$$\mathcal{L}(\text{no del}) = (1 - \pi) \prod_{S \in \mathcal{S}} \mathcal{L}(G_0 | \Delta^S) \quad (2)$$

where  $\pi$  is the prior probability of observing a deletion (default  $10^{-4}$ ). For the likelihood of the deletion model, we use the weighted sums of all three genotype likelihoods in a similar product

$$\mathcal{L}(\text{del of length } l) = \pi \prod_{S \in \mathcal{S}} a_0^S \mathcal{L}(G_0 | \Delta^S) + a_1^S \mathcal{L}(G_1 | \Delta^S) + a_2^S \mathcal{L}(G_2 | \Delta^S) \quad (3)$$

where the  $a_g^S$  are sample- and genotype-specific weights (see below) with genotypes  $g \in \{0, 1, 2\}$  corresponding to 0, 1, or 2 variant alleles and  $a_0^S + a_1^S + a_2^S = 1$  for any  $S \in \mathcal{S}$ .

**Iterative estimation of parameters for the likelihood ratio test.** The likelihood ratio test requires as input a deletion length  $l$ , genotype likelihoods  $\mathcal{L}(G_g|\Delta^S)$  for all samples  $S$  and sample- and genotype-specific weights  $a_g^S$ . PopDel estimates these values for each window iteratively from the profiles together with an allele frequency  $f$  that is needed for updating the weights (Fig. 1). For simplicity, the following assumes one read group per sample. Our implementation in PopDel also handles multiple read groups (Supplementary Note 1). To be able to detect deletions of different lengths from different haplotypes overlapping the same window, the iteration and likelihood ratio test are performed for several initializations of the deletion length. Initial lengths are estimated by identifying samples with similar third quartiles of  $\Delta^S$  via greedy clustering (Supplementary Note 1). The initial allele frequencies  $f$  are set to the fraction of deletion-supporting read pairs of all samples in the window (Supplementary Note 1). To calculate the genotype likelihoods of the three genotypes  $G_0$ ,  $G_1$ , and  $G_2$  of a single sample  $S$ , PopDel transforms the insert size histogram of  $S$  to reflect how many read pairs with a given insert size deviation  $\delta \in \Delta^S$  are expected to overlap a window of size  $w$  (Supplementary Note 1). We denote the resulting relative likelihood of observing a read pair with insert size deviation  $\delta$  as  $H^S(\delta)$ .

PopDel calculates the likelihoods  $\mathcal{L}(G_g|\Delta^S)$  as

$$L(G_0|\Delta^S) = \prod_{\delta \in \Delta^S} H^S(\delta - \epsilon_S) \tag{4}$$

$$L(G_1|\Delta^S) = \prod_{\delta \in \Delta^S} \frac{H^S(\delta - \epsilon_S) + H^S(\delta - l)}{2} \tag{5}$$

$$L(G_2|\Delta^S) = \prod_{\delta \in \Delta^S} H^S(\delta - l) \tag{6}$$

where  $\epsilon_S$  is a sample-specific reference shift (Supplementary Note 1) that accounts for local biases of the data such as GC (guanine–cytosine) content<sup>42,43</sup>.

Our sample- and genotype-specific weights  $a_g^S$  are designed to give low weight to samples with a small likelihood for the genotype and a high weight to those with a good likelihood for the genotype. Furthermore, the weights make it more likely to observe a carrier genotype when the allele frequency is high:

$$a_g^S = \frac{\mathcal{L}(G_g|\Delta^S)\mathcal{L}(f, G_g)}{\sum_{j=0}^2 (\mathcal{L}(G_j|\Delta^S)\mathcal{L}(f, G_j))} \tag{7}$$

with  $\mathcal{L}(f, G_g)$  as the expected genotype frequencies given the population allele frequency  $f$ :

$$L(f, G_0) = (1 - f)^2 \tag{8}$$

$$L(f, G_1) = 2f(1 - f) \tag{9}$$

$$L(f, G_2) = f^2 \tag{10}$$

Given the weights, we update the allele frequency  $f$  using:

$$f^{\text{new}} = \frac{1}{2|S|} \sum_{S \in S} (a_1^S + 2a_2^S) \tag{11}$$

To update the deletion length  $l$ , probabilities  $P_{l, \epsilon_S}^S(\delta)$  reflecting that a given insert size deviation  $\delta$  resulted from a distribution shifted by  $l$  rather than by  $\epsilon_S$  are calculated as

$$P_{l, \epsilon_S}^S(\delta) = a_1^S \frac{H^S(\delta - l)}{H^S(\delta - \epsilon_S) + H^S(\delta - l)} + a_2^S \tag{12}$$

and used to update  $l$  jointly across all samples as the weighted sum over all insert size deviations:

$$l^{\text{new}} = \frac{\sum_{S \in S} \sum_{\delta \in \Delta^S} \delta P_{l, \epsilon_S}^S(\delta)}{\sum_{S \in S} \sum_{\delta \in \Delta^S} P_{l, \epsilon_S}^S(\delta)} \tag{13}$$

The iteration for parameter estimation terminates when both the allele frequency and deletion length converge or additional termination conditions are met (Supplementary Note 1), for example, reaching the maximum number of iterations (default 15). A start position of the potential deletion is estimated during above calculations by keeping track of the rightmost aligned positions of the forward reads of read pairs whose  $\delta$  support the deletion estimate (Supplementary Note 1).

**Combining consecutive deletion windows.** The likelihood ratio tests are performed per initialized deletion length for each genomic window (of size 30 bp by default). The deletions identified by PopDel typically overlap several consecutive windows and in each window the null hypothesis of the likelihood ratio test may be rejected. To provide the user with a nonredundant list of deletion variants, adjacent windows that support the same deletion are combined. PopDel sorts all pairs of windows and deletion lengths, for which the null hypothesis of the likelihood ratio test can be rejected, in ascending order of the predicted deletion start position,

deletion length, and deletion likelihood ratio. Traversing this sorted list of windows  $w_0, w_1, \dots$ , a window  $w_i, i \geq 0$  is combined with another window  $w_{i+k}, k > 0$  if their start positions and deletion sizes are similar enough (Supplementary Note 1). When no more windows can be combined with  $w_i$ , a deletion is output with a start position and length calculated as the median over all combined windows. The algorithm continues with the next window  $w_{i+k+1}$  that has not been combined with any other window so far.

**Deletion output.** We report the mean PHRED-scaled genotype likelihoods across the combined windows of one sample in the output. Samples without sufficient data or much higher than average coverage at the locus are not genotyped (Supplementary Note 1). The allele frequency is estimated by counting the number of alleles predicted to carry the variant, divided by the total number of genotyped alleles. We report a genotype quality as the difference of the best and second best PHRED-scaled genotype likelihoods.

**Simulation of sequencing data with deletions.** We simulated two cohorts of sequencing data, the first consisting of 1000 diploid individuals carrying random deletions and the second consisting of 500 individuals carrying deletions reported in the 1000 Genomes Project (G1k). For the first cohort, we simulated a set of 2000 deletion variants with uniformly distributed lengths between 100 and 10,000 bp, uniformly distributed allele frequencies between 0 and 1, and uniformly distributed positions on chromosome 21 of GRCh38. Regions containing N's were excluded and deletion were required to be at least 1000 bp apart. For the second cohort, we downloaded deletions identified in the 1000 Genomes Project (see ‘‘Data availability’’) on chromosomes 17–22 of GRCh37 and their reported allele frequencies.

Using the random deletion set, we created 2000 haplotypes by sampling deletions according to their allele frequency and inserting them into GRCh38. Using the G1k deletions, we created 1000 haplotypes by sampling deletions according to their allele frequencies and inserting them into GRCh37. The haplotypes were combined into 1000 and 500 diploid samples, respectively. These samples were used to simulate NGS reads with art\_illumina<sup>44</sup> and the reads aligned to GRCh38 or GRCh37 using BWA-mem<sup>45</sup>.

**Setup of SV callers on simulated data.** PopDel (1.2.2) and Manta (1.6.0) were run with an option to limit the calling to chromosome 21. Smooove (0.2.4 with Lumpy 0.2.13 and SVtyper 0.7.0, obtained from <https://github.com/brentp/smooove>) was applied as recommended by its authors using the provided exclude regions for GRCh38 and GRCh37, and excluding mappings not on the simulated chromosomes. Delly (0.7.8) was applied without small indel realignment (option -n). GRIDSS (1.8.1) was provided a maximum heap size of 8 GB. All tools were applied on increasing numbers of BAM files (up to 1000 or until failure) in steps of 1 up to 10, steps of 10 up to 100, and steps of 100 up to 1000.

**Evaluation on simulated data.** Running time and memory consumption were measured on a dedicated workstation (Intel Xeon E5-1630v3 8 × 3.5 GHz, 64 GB RAM) using the Unix time command. For tools consisting of multiple steps, the running times are the sum of the time taken by all steps from input BAM files to output variant call format (VCF) file. The memory consumption is stated as the maximum memory consumption of all steps. As GRIDSS (Genomic Rearrangement Identification Software Suite) produces two break-ends per deletion, corresponding pairs of break-ends were collapsed into a single call and LOW\_QUAL variants were removed. The calls of Delly were not filtered for variants that have the filter field set to PASS as this has a negative impact on its performance. A call is considered to match a simulated variant in case of a reciprocal overlap of at least 50%. Each simulated variant is allowed to be matched with only one predicted variant. See Supplementary Note 2 for results using alternative match criteria.

**Sample preparation and setup of SV callers on GIAB and Polaris data.** All samples were aligned to the human reference genome GRCh38<sup>46</sup> using BWA-mem<sup>47,48</sup>, except for the data of HG002/NA24385 and his parents, which was aligned to GRCh37. All SV callers were applied as recommended by the authors and, if possible, limited to the reference sequence of the 22 autosomes. Delly was run with the option -n. All tools could be run jointly on the GIAB trio data. Smooove, the recommended pipeline for running Lumpy, recommends a workflow consisting of single-sample calling, merging, and sample-wise re-genotyping for cohorts of 40 or more individuals and was applied accordingly. Joint calling with Delly and Manta did not finish within 4 weeks on the Polaris data. Therefore, Delly was run using single-sample calling with merging and sample-wise re-genotyping following the germline SV calling workflow described by the authors (Supplementary Note 2). For Manta, we could not find a description of a similar workflow, so we excluded it from the analysis of the Polaris data. Calls other than deletions were removed as early as possible in Delly's and Smooove's workflows to reduce the running time (Supplementary Note 2). The sample order of the Polaris Diversity and Kids cohorts was shuffled, but the same for all tools and no tool was provided pedigree information.



**Variant filtering on real data.** The reference call sets were filtered for deletion variants. All variants on contigs apart from chromosomes 1–22 and chromosome X were removed and VCF-to-BED conversion was performed for the PacBio reference set. Liftover from GRCh37 to GRCh38 was performed for HG001/NA12878 using the NCBI Genome Remapping Service (<https://www.ncbi.nlm.nih.gov/genome/tools/remap>).

Deletions identified by the tested tools in HG002/NA24385 and his parents were filtered using the high-confidence regions prepared by the GIAB consortium (see “Data availability”). All other deletion sets were filtered for centromeric regions. Centromeric regions were obtained through the UCSC table browser (group: Mapping and Sequencing, track: Centromeres) for GRCh38. Any deletions in a reference set or call set having any overlap with a centromeric region or a region outside the high-confidence regions were removed. Overlap was determined using bedtools intersect<sup>49</sup> (Supplementary Note 2).

Only deletions on the 22 autosomes were considered for analysis. Deletions were filtered to the size range from 500 to 10,000 bp. Two deletions were considered the same if they had a reciprocal overlap of 50% or more. The Polaris data was filtered to high-confidence deletions with genotype quality scores above a fixed threshold. This threshold was chosen once per tool on the Polaris Kids cohort such that the Mendelian inheritance error rate dropped below 0.3%<sup>37</sup>; 26 for PopDel, 28 for Delly, and 78 for Lumpy. To search for de novo deletions, the genotype threshold for PopDel was set to 50 and all 12 resulting candidate deletions were inspected manually in Integrative Genomics Viewer<sup>50</sup>. In all real data analyses, Delly variants were only considered if they had the FILTER field set to PASS.

**Principal component analysis.** Predicted genotypes of the Polaris Diversity cohort were converted into a variant/sample matrix containing deletion allele counts. Uninformative deletions and those in linkage disequilibrium were removed (Supplementary Note 2). Principal component analysis was computed using the R function `prcomp`.

**Mendelian inheritance error rate and transmission rate.** All deletions that are not in Hardy–Weinberg equilibrium ( $P$  value threshold 0.01) were removed. For all reported deletions, the three genotypes in each trio were inspected for Mendelian consistency (Supplementary Fig. 4). Trios with one or more missing genotypes and trios with all three samples genotyped as 0/0 were ignored. The transmission rate was calculated as the number of deletion alleles transmitted from the heterozygous parents to the children divided by the number of considered deletions. If indicated, we calculated the transmission rate considering only deletions that were called in a single trio, where one parent is a heterozygous carrier and the other parent carries the reference allele on both haplotypes.

**Sequence data from 49,962 Icelanders.** DNA was isolated from both blood and buccal samples. All participating subjects signed informed consent. The personal identities of the participants and biological samples were encrypted by a third-party system approved and monitored by the Data Protection Authority. The National Bioethics Committee and the Data Protection Authority in Iceland approved these studies. The Icelandic samples were whole-genome sequenced at deCODE genetics using Illumina GAIIX, HiSeq, HiSeq X, and NovaSeq sequencing machines, and sequences were aligned to the human reference genome (GRCh38)<sup>46</sup> using BWA-mem<sup>47,48</sup>. Details of the sample preparation, paired-end sequencing, read processing and alignment, and selection of the final set of BAM files have been previously described<sup>51</sup>.

**Running time measurements on real data.** The running time on NA12878 was measured using the Unix time command on a dedicated workstation (Intel Xeon E5-1630v3  $8 \times 3.5$  GHz, 64 GB RAM). Due to limited storage capacities on our workstations, running times for the Polaris Diversity cohort were measured on the BIH high-performance compute cluster as reported in SGE log files.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

The generated VCF files used for evaluation of all simulated and real data sets have been deposited in Zenodo (<https://doi.org/10.5281/zenodo.3992607>)<sup>52</sup>. The deletion set of the 1000 Genomes Project we used for simulating data was obtained from [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated\\_sv\\_map/ALL.wgs.mergedSV.v8.20130502.svs.genotypes.vcf.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated_sv_map/ALL.wgs.mergedSV.v8.20130502.svs.genotypes.vcf.gz). The short-read data from the HG001/NA12878 genome is publicly available under ERA run accession ERR194147. The short-read data from the HG002/NA24385 genome and parental genomes HG003/NA24149 and HG004/NA24143 are available under BioProject accession PRJNA200694. The precise list of run accession numbers we used for HG002 and his parents is given in the Supplementary Note 2. The GIAB short-read reference set for HG001/NA12878 was obtained from [ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/technical/svclassify\\_Manuscript/Supplementary\\_Information/Personalis\\_1000\\_Genomes\\_deduplicated\\_deletions.bed](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/technical/svclassify_Manuscript/Supplementary_Information/Personalis_1000_Genomes_deduplicated_deletions.bed) and the long read reference call sets from <ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/>

NA12878/NA12878\_PacBio\_MtSinai/NA12878.sorted.vcf.gz. The preliminary variant set for HG002 is available at [ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/NIST\\_SVs\\_Integration\\_v0.6/HG002\\_SVs\\_Tier1\\_v0.6.vcf.gz](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/NIST_SVs_Integration_v0.6/HG002_SVs_Tier1_v0.6.vcf.gz) and the high-confidence regions for HG002 at [ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/NIST\\_SVs\\_Integration\\_v0.6/HG002\\_SVs\\_Tier1\\_v0.6.bed](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/NIST_SVs_Integration_v0.6/HG002_SVs_Tier1_v0.6.bed). The short-read data of the Polaris HiSeq X Diversity cohort is publicly available under BioProject accession <https://www.ebi.ac.uk/ena/browser/view/PRJEB20654> and the Polaris HiSeq X Kids cohort under BioProject accession PRJEB25009. We obtained the human reference genome GRCh37 from [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2\\_reference\\_assembly\\_sequence/hs37d5.fa.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assembly_sequence/hs37d5.fa.gz) and downloaded GRCh38 on June 14th, 2017 from [ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/001/405/GCA\\_000001405.15\\_GRCh38/seqs\\_for\\_alignment\\_pipelines.ucsc\\_ids/GCA\\_000001405.15\\_GRCh38\\_full\\_analysis\\_set.fna.gz](ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/001/405/GCA_000001405.15_GRCh38/seqs_for_alignment_pipelines.ucsc_ids/GCA_000001405.15_GRCh38_full_analysis_set.fna.gz) using bwakit (0.7.15). Access to raw Icelandic sequence data is restricted by Icelandic state law and can be given through collaboration with KS.

## Code availability

PopDel is available for installation through Bioconda<sup>53</sup>. The source code is available at <https://github.com/kehrlab/PopDel> (v1.2.2, GNU GPLv3 license, <https://doi.org/10.5281/zenodo.4282041>)<sup>54</sup>. PopDel was implemented using the SeqAn C++ library<sup>55</sup>. Scripts used for evaluating PopDel are available at <https://github.com/kehrlab/PopDel-scripts>.

Received: 22 August 2019; Accepted: 14 December 2020;

Published online: 01 February 2021

## References

- Auton, A. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- Francioli, L. C. et al. Genome-wide patterns and properties of de novo mutations in humans. *Nat. Genet.* **47**, 822–826 (2015).
- Jónsson, H. et al. Parental influence on human germline de novo mutations in 1,548 trios from Iceland. *Nature* **549**, 519–522 (2017).
- Halldórsson, B. V. et al. Characterizing mutagenic effects of recombination through a sequence-level genetic map. *Science* **363**, eaau1043 (2019).
- Carvalho, C. M. & Lupski, J. R. Mechanisms underlying structural variant formation in genomic disorders. *Nat. Rev. Genet.* **17**, 224–288 (2016).
- Abyzov, A. et al. Analysis of deletion breakpoints from 1,092 humans reveals details of mutation mechanisms. *Nat. Commun.* **6**, 7256 (2015).
- Goldmann, J. M. et al. Germline de novo mutation clusters arise during oocyte aging in genomic regions with high double-strand-break incidence. *Nat. Genet.* **50**, 487–492 (2018).
- GTEX Consortium. Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).
- Conrad, D. F. et al. Origins and functional impact of copy number variation in the human genome. *Nature* **464**, 704–712 (2010).
- Kircher, M. et al. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
- Novembre, J. et al. Genes mirror geography within Europe. *Nature* **456**, 98–101 (2008).
- Buniello, A. et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
- Ingelman-Sundberg, M., Mkrтчian, S., Zhou, Y. & Lauschke, V. M. Integrating rare genetic variants into pharmacogenetic drug response predictions. *Hum. Genomics* **12**, 26 (2018).
- Gudbjartsson, D. F. et al. Large-scale whole-genome sequencing of the Icelandic population. *Nat. Genet.* **47**, 435–444 (2015).
- Turro, E. et al. Whole-genome sequencing of patients with rare diseases in a national health system. *Nature* **583**, 96–102 (2020).
- Gilly, A. et al. Cohort-wide deep whole genome sequencing and the allelic architecture of complex traits. *Nat. Commun.* **9**, 4674 (2018).
- Yuen, R. K. C. et al. Whole genome sequencing resource identifies 18 new candidate genes for autism spectrum disorder. *Nat. Neurosci.* **20**, 602–611 (2017).
- Mak, A. C. Y. et al. Whole-genome sequencing of pharmacogenetic drug response in racially diverse children with asthma. *Am. J. Respir. Crit. Care Med.* **197**, 1552–1564 (2018).
- Karczewski, K. J. et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
- Taliun, D. et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. Preprint at *bioRxiv* <https://doi.org/10.1101/563866> (2019).
- Poplin, R. et al. Scaling accurate genetic variant discovery to tens of thousands of samples. Preprint at *bioRxiv* <https://doi.org/10.1101/201178> (2017).
- Eggertsson, H. P. et al. GraphTyper enables population-scale genotyping using pangenome graphs. *Nat. Genet.* **49**, 1654–1660 (2017).

23. Guan, P. & Sung, W.-K. Structural variation detection using next-generation sequencing data. *Methods* **102**, 36–49 (2016).
24. Sudmant, P. H. et al. An integrated map of structural variation in 2,504 human genomes. *Nature* **526**, 75–81 (2015).
25. Handsaker, R. E., Korn, J. M., Nemes, J. & McCarroll, S. A. Discovery and genotyping of genome structural polymorphism by sequencing on a population scale. *Nat. Genet.* **43**, 269–276 (2011).
26. Abel, H. J. et al. Mapping and characterization of structural variation in 17,795 human genomes. *Nature* **583**, 83–89 (2020).
27. Larson, D. E. et al. Svtools: population-scale analysis of structural variation. *Bioinformatics* **35**, 4782–4787 (2019).
28. Zarate, S. et al. Parliament2: accurate structural variant calling at scale. *GigaScience* **9**, gaaa145 (2020).
29. Mohiyuddin, M. et al. MetaSV: an accurate and integrative structural-variant caller for next generation sequencing. *Bioinformatics* **31**, 2741–2744 (2015).
30. Ho, S. S., Urban, A. E. & Mills, R. E. Structural variation in the sequencing era. *Nat. Rev. Genet.* **21**, 171–189 (2020).
31. Abyzov, A., Urban, A. E., Snyder, M. & Gerstein, M. CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs from family and population genome sequencing. *Genome Res.* **21**, 974–984 (2011).
32. Handsaker, R. E. et al. Large multiallelic copy number variations in humans. *Nat. Genet.* **47**, 296–303 (2015).
33. Rausch, T. et al. DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics* **28**, i333–i339 (2012).
34. Layer, R. M., Chiang, C., Quinlan, A. R. & Hall, I. M. LUMPY: a probabilistic framework for structural variant discovery. *Genome Biol.* **15**, R84 (2014).
35. Chen, X. et al. Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinformatics* **32**, 1220–1222 (2016).
36. Cameron, D. L. et al. GRIDSS: sensitive and specific genomic rearrangement detection using positional de Bruijn graph assembly. *Genome Res.* **27**, 2050–2060 (2017).
37. Zook, J. M. et al. A robust benchmark for detection of germline large deletions and insertions. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-020-0538-8>, 1–9 (2020).
38. Collins, R. L. et al. A structural variation reference for medical and population genetics. *Nature* **581**, 444–451 (2020).
39. Frankish, A. et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.* **47**, D766–D773 (2019).
40. Kloosterman, W. P. et al. Characteristics of de novo structural changes in the human genome. *Genome Res.* **25**, 792–801 (2015).
41. Björnsson, E. et al. Lifelong Reduction in LDL Cholesterol Due to a Gain-of-Function Mutation in *LDLR*. *Circulation: Genomic and Precision Medicine* (2020).
42. Benjamini, Y. & Speed, T. P. Summarizing and correcting the GC content bias in high-throughput sequencing. *Nucleic Acids Res.* **40**, e72–e72 (2012).
43. Iakovishina, D., Janoueix-Lerosey, I., Barillot, E., Regnier, M. & Boeva, V. SV-Bay: structural variant detection in cancer genomes using a Bayesian approach with correction for GC-content and read mappability. *Bioinformatics* **32**, 984–992 (2016).
44. Huang, W., Li, L., Myers, J. R. & Marth, G. T. ART: a next-generation sequencing read simulator. *Bioinformatics* **28**, 593–594 (2012).
45. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinformatics* **26**, 589–595 (2010).
46. Schneider, V. A. et al. Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. *Genome Res.* <https://doi.org/10.1101/gr.213611.116> (2017).
47. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
48. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint at <https://arxiv.org/abs/1303.3997> (2013).
49. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
50. Robinson, J. T. et al. Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
51. Jónsson, H. et al. Data descriptor: whole genome characterization of sequence diversity of 15,220 Icelanders. *Sci. Data* **4**, <https://doi.org/10.1038/sdata.2017.115> (2017).
52. Niehus, S. PopDel identifies medium-size deletions jointly in tens of thousands of genomes—variant call sets. *Zenodo* <https://doi.org/10.5281/zenodo.3992607> (2020).
53. Grünig, B. et al. Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nat. Methods* **15**, 475–476 (2018).
54. Niehus, S., Haldórsson, B. V. & Kehr, B. PopDel identifies medium-size deletions jointly in tens of thousands of genomes. *GitHub* <https://doi.org/10.5281/zenodo.4282041> (2020).
55. Reinert, K. et al. The SeqAn C++ template library for efficient sequence analysis: a resource for programmers. *J. Biotechnol.* **261**, 157–168 (2017).

## Acknowledgements

Computation has been performed on the HPC for Research cluster of the Berlin Institute of Health. The project was funded by the Federal Ministry of Education and Research (BMBF) through grant number FKZ 031L0180 to B.K.

## Author contributions

B.K. and B.V.H. conceived the approach. S.N. and B.K. developed PopDel, designed the experiments, and evaluated the results. J.S. assisted in the development of the approach. S.N. simulated data, performed analyses on simulated and public data, and tested on Icelandic data. B.V.H. applied PopDel on Icelandic data and assisted in testing. E.B., P.S., and K.S. studied the *LDLR* deletion. H.J. assisted in the computation of transmission rates. H.P.E. and D.B. assisted in the analyses of Icelandic data. B.K. and S.N. drafted the manuscript with feedback from B.V.H. All authors revised the draft and approved the final manuscript.

## Funding

Open Access funding enabled and organized by Projekt DEAL.

## Competing interests

H.J., E.B., H.P.E., D.B., P.S., K.S., and B.V.H. are employees of deCODE genetics/Amgen, Inc. The other authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-020-20850-5>.

**Correspondence** and requests for materials should be addressed to B.K.

**Peer review information** *Nature Communications* thanks Can Alkan and Kai Wang for their contribution to the peer review of this work.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021