

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

***de novo* variant calling identifies cancer mutation profiles in the 1000 Genomes Project**

Jeffrey Ng<sup>1</sup>, Pankaj Vats<sup>2</sup>, Elyn Fritz-Waters<sup>3</sup>, Evin M. Padhi<sup>1</sup>, Zachary L. Payne<sup>1</sup>, Shawn Leonard<sup>3</sup>, Stephanie Sarkar<sup>1</sup>, Marc West<sup>2</sup>, Chandler Prince<sup>3</sup>, Lee Trani<sup>4</sup>, Marshall Jansen<sup>3</sup>, George Vacek<sup>2</sup>, Mehrzad Samadi<sup>2</sup>, Timothy T. Harkins<sup>2</sup>, Craig Pohl<sup>3</sup>, Tychele N. Turner<sup>1,\*</sup>

1. Department of Genetics, Washington University School of Medicine, St. Louis, MO 63110, USA
2. NVIDIA
3. Research Infrastructure Services, Washington University School of Medicine, St. Louis, MO 63110, USA
4. McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO 63110, USA

\*Correspondence to: [tychele@wustl.edu](mailto:tychele@wustl.edu)  
Tychele N. Turner, Ph.D.  
Washington University School of Medicine  
Department of Genetics  
4523 Clayton Avenue  
Campus Box 8232  
St. Louis, MO 63110

1 **ABSTRACT**

2 Detection of *de novo* variants (DNVs) is critical for studies of disease-related variation and  
3 mutation rates. We developed a GPU-based workflow to call DNVs, using 602 trios from the 1000  
4 Genomes Project as a control. We detected 445,711 DNVs, having a bimodal distribution, with  
5 peaks at 200 and 2000 DNVs. The excess DNVs are cell line artifacts that are increasing with cell  
6 passage. Reduction in DNVs at CpG sites and in percent of DNVs with a paternal parent-of-origin  
7 with increasing number of DNVs supports this finding. Detailed assessment of individual  
8 NA12878 across multiple genome datasets from 2012 to 2020 reveals increasing number of DNVs  
9 over time. Mutation signature analysis across the set revealed individuals had either 1) age-related,  
10 2) B-cell lymphoma, or 3) no prominent signatures. Our approach provides an important  
11 advancement for DNV detection and shows cell line artifacts present in lymphoblastoid cell lines  
12 are not always random.

13

14

## 1 INTRODUCTION

2 The 1000 Genomes Project is a data resource for the study of genetic variation that includes  
3 individuals from diverse genetic ancestries (1, 2). In this study, we present the first ever assessment  
4 of new mutations, termed *de novo* variants (DNVs) that are only found in children and not their  
5 parents, represented in this collection. DNVs are important for assessing mutation rates (3) and  
6 have been shown to contribute to human disease (autism (4-12), epilepsy (13, 14), intellectual  
7 disability (15-18), congenital heart disorders (19-21)). This is a critical reason to assess these  
8 variants in the 18 populations represented in this data (Figure S1). Moreover, the 1000 Genomes  
9 Project has been utilized in many applications as a control resource for filtering of genetic variation  
10 by allele frequency and/or variant presence-absence in the dataset (22).

11  
12 One complicating factor of DNV assessment in this resource is the fact that sequencing  
13 data is generated from DNA isolated from lymphoblastoid cell lines (LCLs) (2) as opposed to  
14 primary tissue. Epstein-Barr Virus is used to make these LCLs and passaging over time enables  
15 the accumulation of cell line artifacts. These artifacts can complicate variant filtration schemes and  
16 the utility of this data as a frequency control. As opposed to a random accumulation of mutations  
17 in each individual, we found that DNVs fit one of three possible mutation profiles: 1) age-related  
18 (similar to true DNVs); 2) B-cell lymphoma related; and 3) an absence of specific profile. The  
19 most problematic of these is the similarity to B-cell lymphomas and it would be imperative that  
20 this data not be used as a control in the context of the study of these and related cancers.

21  
22 A second aspect of this paper is the development of a rapid approach for identifying DNVs  
23 that will be applicable to any short-read whole-genome sequencing data. Typically calling of  
24 DNVs from raw sequence data to final calls can take days to weeks. Multiple *de novo* workflows  
25 exist that primarily rely on central processing unit (CPU)-based approaches (4-9, 11, 12, 14-17,  
26 19, 23-34). These DNV workflows employ different approaches including strict filtering, utilizing  
27 multiple variant callers as opposed to using only one, machine-learning, and incorporation of  
28 genotypic information at other sites around the genome. Overall, there is no community consensus  
29 on a standard method for detecting DNVs. It is imperative that this process be streamlined and  
30 flexible to enable broad adoption across the community. The approach we developed to accelerate  
31 this analysis was the utilization of graphics processing units (GPUs), enabling highly-parallelized  
32 analyses. Herein, we developed a DNV workflow and integrated it into NVIDIA Parabricks (35)  
33 software for a drastic acceleration and reduced run time from alignment files, for each trio, to final  
34 DNVs in less than one hour. This type of computational upgrade for specific workflows are  
35 essential for the future of genomics as projects become larger in scale.

36  
37

## 1 RESULTS

### 2 *Rapid DNV calling with GPUs*

3 Our DNV workflow utilized the existing features of the NVIDIA Parabricks software  
4 including fast GATK (36) HaplotypeCaller and Google's DeepVariant (37) gvcf generation. The  
5 run times are ~40 minutes per sample on a 4 GPU node and can be run in parallel on all three  
6 family members in the parent-child trio. Genotyping of the trio is quick through the use of GLnexus  
7 (38). Finally, our DNV workflow runs in ~1 hour with speedups at all steps with parallelization  
8 providing a clear advantage over CPU-based approaches (Figure 1A).

9  
10 To benchmark our DNV workflow, we tested it on a monozygotic twin pair with WGS data  
11 derived from blood DNA. These individuals should share the same DNVs from generation in the  
12 germline. However, they may differ at some sites if DNVs occur in a post-zygotic, somatic manner.  
13 The twins shared 75 DNVs and contained 83 and 81 DNVs, respectively (Figures 1B and 1C). The  
14 percent CpG was 19.3% and 17.2%, respectively and in line with estimates ~20% (Figures 1B and  
15 1C). As this monozygotic twin pair was discordant for the phenotype of autism, we also tested  
16 whether there were any protein-coding DNV differences between the two twins. These would  
17 potentially be relevant for autism, but there were no such differences.

18  
19 We next assessed DNVs with our workflow in four trios from the 1000 Genomes Project  
20 (Figure 1D). Two were chosen at random (i.e., HG00405, HG00408) and two were chosen because  
21 they were "famous" trios assessed in many other studies (i.e., NA12878 (26, 39), NA19240 (26)).  
22 One of these trios (HG00405) had 70 DNVs and a CpG percent of 21.4 as we would have expected  
23 from DNA derived from blood. To our surprise, the other trios had varying numbers of DNVs  
24 from 592 to 2,230 with NA12878 (arguably the most studied individual in the 1000 Genomes  
25 Project) having the most DNVs. With the increase in DNVs the CpG percent dropped considerably  
26 down to ~10%. We also assessed 3,598 of the DNVs from the four trios by visual inspection of  
27 the reads in each family member (Table S1) and found that 93.6% of the variants appeared to be  
28 true DNVs, 4.9% were inherited, and 1.5% were low confidence calls.

### 29 *Excess of DNVs observed in trio dataset*

30  
31 We applied our DNV workflow to all 602 trios (Table S2) in the 1000 Genomes Project  
32 and detected 445,711 total DNVs with  $740.4 \pm 968.0$  DNVs per individual (Tables S3 and S4).  
33 There was a clear bimodal distribution wherein some individuals contained an excess of DNVs  
34 (Figure 2A). We split the data into two groups: individuals having less than or equal to 100 DNVs  
35 ( $n = 123$ ) and individuals had greater than 100 DNVs ( $n = 479$ ). The individuals with less than or  
36 equal to 100 DNVs were in line with expectation for expected DNVs. While we expected there to  
37 be no difference by ancestry we sought to see if this could somehow be a factor (Figure 2B). The  
38 population with the most DNVs was the CEU having on average 1,688 DNVs. We hypothesized  
39 that this may be because the CEU is one of the oldest cohorts in the 1000 Genome Project dating  
40 back to the HapMap project (40) and these individuals may have cell lines that have been cultured  
41 more over time than other populations.

42  
43 Since we thought the individuals with excess of DNVs represented cell line artifacts, we  
44 assessed two main features of typical DNVs. These were the percent of DNVs at CpG locations  
45 and the percent of DNVs arising on the paternal chromosome. It has been well-established that the  
46 percent of DNVs at CpG should be ~20% (3, 8) and the percent of DNVs arising on the paternal

1 chromosome should be ~80% (41). We saw that overall,  $13.7 \pm 4.4\%$  of DNVs per individual  
2 occurred at CpG sites. In the individuals with less than or equal to 100 DNVs this rose to  $17.4 \pm$   
3  $5.2\%$  and in families with greater than 100 DNVs it fell to  $12.7 \pm 3.6\%$  (Figure 2C). The percent  
4 of DNVs that were phase-able for parent-of-origin was  $37.2 \pm 7.5\%$  (Figure S2). Of the phased  
5 variants,  $61.3 \pm 11.3\%$  were on the chromosome of paternal origin (Table S5). In the families with  
6 less than or equal to 100 DNVs this rose to  $72.0 \pm 8.5\%$  and in the families with greater than 100  
7 DNVs it fell to  $58.6 \pm 10.3\%$ . The drop leveled off to ~50% in the individuals with the most DNVs  
8 (Figure 2D). This showed that the individuals with DNVs at expectation for counts also behaved  
9 more like true DNVs in regard to CpG percentage and percent arising on the paternal chromosome.

10

### 11 *DNVs increase over time*

12 We utilized the fact that the 1000 Genomes Project individual NA12878 has been studied  
13 and sequenced multiple times over the past ten years by WGS (2) (SRA identifiers: SRR944138  
14 and SRR952827). Presumably, across time, the utilization of NA12878 has required additional  
15 culturing of this cell line, and potentially even by different laboratories. We aggregated five  
16 Illumina WGS datasets from this individual, downsampled them to ~30x coverage, and assessed  
17 them by our DNV workflow. The data for this individual ranged from the year 2012 to the year  
18 2020 and we found that the 2012 experiment had the least DNVs ( $n = 2,060$ ) and the 2020  
19 experiment had the most DNVs ( $n = 2,230$ ) (Figure 3A). Overall, the five replicates had a large  
20 overlap of DNVs ( $n = 1,820$ ) across all samples. These shared DNVs constitute what were present  
21 in the ancestor of all the cell line replicates. Mutations not shared by all five replicates are  
22 sometimes shared by a subset of the replicates and are sometimes unique to the replicate. To  
23 formally assess the ancestral state, we built a phylogenetic tree based only on the DNVs and saw  
24 that the farthest replicates from each other in the tree were the 2012 and 2020 replicates (Figure  
25 3B).

26

### 27 *Genomes with cancer mutation profiles*

28 We used mutation profile analysis (42) (Table S6) to determine whether the DNVs  
29 identified in individuals from the 1000 Genomes project had any certain characteristics. For this  
30 analysis, we utilized a method that would enable comparisons to known mutational profiles that  
31 are either age-related (reminiscent of true DNVs) or are seen in cancers (Figure 4A and Figure  
32 4B). There were 186 individuals that had a strong contribution of an age-related signature  
33 (Signature 1A, Signature 1B). To our surprise, the other contributing signatures in individuals were  
34 primarily those associated with B-cell lymphomas (Signature 5, Signature 9 and Signature 17) in  
35 241 individuals. This was intriguing because lymphoblastoid cell lines are generated from B-cells  
36 that are infected with Epstein Barr Virus and demonstrates that new mutations are not arising in a  
37 random manner. Rather they are being generated in a manner consistent with the development of  
38 cancer in the same cell type.

39

40 We further sought to determine what the mechanism was for the generation of a B-cell  
41 lymphoma-like state. First, we determined whether there was high rate of aneuploidies in the cell  
42 lines. By digital karyotyping (Table S7) we found that 595 individuals (98.8%) had a typical  
43 chromosome complement (46,XX or 46,XY), four were missing a sex chromosome (45,X0), one  
44 was 47,XXY, one had three chromosome 12 (47,XY), and one had three chromosome 9 (47,XY).  
45 This demonstrated that while these aneuploidies are occurring in some cell lines, they are probably  
46 not the main driving factor. Next, we looked at DNVs in genes involved in DNA repair and found

1 17 individuals contained a missense or loss-of-function in one of these genes (Table S8).  
2 Individuals with B-cell lymphoma profiles and disruptive mutations in DNA repair genes included  
3 mutations in the following genes *FANCF* (HG01126), *MUS81* (NA10838), *POLB* (NA10838),  
4 *POLD1* (NA19677), *POLE* (HG01096), *RAD18* (NA12864) (Figure 4C), *RAD51* (HG02683),  
5 *RPA4* (HG02630), and two individuals with mutations in *FANCA* (HG02841, HG03200) and *WRN*  
6 (HG04115, NA19161), respectively (Table 1). Third, we looked at Epstein Barr Virus load in each  
7 of the genomes (Table S9) and found that there was a weak, yet significant, correlation with the  
8 number of DNVs ( $p = 2.32 \times 10^{-5}$ ,  $r = 0.17$ ) (Figure S3). By visual inspection of phased variation  
9 in all individuals we also identified individuals with clusters of mutations (e.g., NA07048, Figure  
10 4D, Figure S4).

### 11 12 *Excess of DNVs in IGLL5*

13 We applied a multi-phase approach to determine if there were any genes with enrichment  
14 of protein-coding DNVs in individuals with greater than 100 DNVs. In the first phase, we tested  
15 whether there was genome-wide significance for enrichment of protein-coding DNVs (missense,  
16 loss-of-function) in any specific genes. By application of two methods (chimpanzee-human,  
17 denovolyzeR), we identified 29 significant genes (*ARMC3*, *BCL2*, *BCR*, *C6orf15*, *CCDC168*,  
18 *CSMD3*, *EGR3*, *EXO1*, *HLA-B*, *HLA-C*, *IGLL5*, *KMT2D*, *LINGO2*, *LTB*, *MEOX2*, *MUC16*,  
19 *MUC22*, *NPAP1*, *PCLO*, *PRPF40A*, *RUNX1T1*, *SGK1*, *STRAP*, *TMEM232*, *TNXB*, *TTN*, *WDFY4*,  
20 *XIRP2*, *ZNF488*) with excess of DNVs (Table S10). In the second phase, we tested these 29 genes  
21 to see whether there were significantly more protein-coding DNVs in individuals with greater than  
22 100 DNVs in comparison to individuals with less than or equal to 100 DNVs. Only *IGLL5* was  
23 significant in this comparison ( $1.79 \times 10^{-3}$ ) (Table S10, Table S11, Figure 4E).

### 24 25 *DNVs identified in clinically-relevant variants*

26 We tested whether any of the DNVs detected were already known to be pathogenic or  
27 likely-pathogenic in the Clinvar (43) database (Table 1). There were 15 mutations meeting these  
28 criteria (Table S12). We rescored these variants using Franklin software to assess their  
29 pathogenicity and found that 13 were also pathogenic or likely-pathogenic by this approach.  
30 Twelve of these variants were associated with described phenotypes in Clinvar. These included a  
31 missense variant in *SOS1* involved in Noonan syndrome, a missense variant in *SCN2A* involved in  
32 seizures, a stop gained variant in *UNC80* involved in a syndrome with hypotonia, intellectual  
33 disability, and characteristic facies, a missense variant in *THRB* involved in thyroid hormone  
34 resistance, a missense variant in *PKHD1* involved in polycystic kidney disease, a stop-gained in  
35 *ERCC6* involved in Cockayne syndrome, a stop-gained in *ANO5* involved in gnathodiaphyseal  
36 dysplasia, a stop-gained in *PHF21A* involved in inborn genetic disease, a missense in *MYO7A* in  
37 Usher syndrome type 1, a stop-gained in *ROBO3* in Gaze palsy with progressive scoliosis, a  
38 missense in *COL4A1* involved in inborn genetic disease, and a missense in *POLG* involved in  
39 POLG-related disorder.

40  
41  
42

## 1 DISCUSSION

2 While the 1000 Genomes Project data has been extensively studied in the past, there has  
3 been no previous cross-cohort assessment of DNVs. This limitation is primarily because family-  
4 based sequencing was not available until 2020 when this cohort was sequenced by high-coverage  
5 short-read whole-genome sequencing ten years after the initial ground-breaking publication on the  
6 1000 Genomes Project (44). Determining DNV profiles across this dataset of diverse individuals  
7 is critical for assessment of mutation rates in the human population, while also providing a more  
8 complete catalog of all genetic variants within these individuals. The decision to sequence these  
9 individuals using DNA derived from lymphoblastoid cell lines was a practical one. However, it  
10 opened the door to the possibility of cell line artifacts, while simultaneously introducing a dynamic  
11 aspect to this extensive set of controls. As control samples, the cell lines that were used as the  
12 inputs for the 1000 Genome Project are still actively used across laboratories, acting as matched  
13 controls for workflows to known sets of variants. The large distribution of DNVs across the 1000  
14 Genomes Project suggest that a subset of the control source inputs are dynamic, and in some cases,  
15 harbor a spectrum of genetic variants associated with B-cell lymphomas or named clinical  
16 syndromes. Laboratories using control samples from the 1000 Genomes Project should account  
17 for both the presence and dynamic nature of the reported DNVs and in some cases may consider  
18 changing which control samples to use within the laboratory to avoid any of the associated issues  
19 with the presence of DNVs. Additionally, other public efforts to establish reference data sets using  
20 cell lines should consider the impacts of DNVs on their project design.

21  
22 We utilized a novel and accelerated analysis workflow to detect DNVs from short-read,  
23 whole-genome sequencing data. In total, we identified 445,711 DNVs in the 602 children assessed  
24 in this study. We provide family-level VCFs, DNV calls, and phased DNV results for the 602 trios  
25 in this study as a public community resource (Globus endpoint: “Turner Lab at WashU - DNV in  
26 1000 Genomes Paper”, direct link: [https://app.globus.org/file-manager?origin\\_id=3eff453a-88f4-11eb-954f-752ba7b88ebe&origin\\_path=%2F](https://app.globus.org/file-manager?origin_id=3eff453a-88f4-11eb-954f-752ba7b88ebe&origin_path=%2F)). Originally, it was assumed that the DNV’s across  
28 the 1000 Genomes Project would have been random and minimal, and yet only 20% of the  
29 offspring (123 children) have a number of DNVs around expectation (< 100) and the remainder  
30 have an excess of DNVs with the most extreme case being an individual (HG02683) having 11,219  
31 DNVs. We hypothesized that the excess DNVs were cell line artifacts and found multiple lines of  
32 evidence to support this hypothesis, including a reduction in the percent of DNVs at CpG as well  
33 as the reduction in percent phased to the paternal parent-of-origin chromosome with increasing  
34 DNVs, respectively. A detailed analysis of individual NA12878, who has been studied various  
35 times over the years, revealed increasing DNVs in the more recently sequenced samples also  
36 supporting this hypothesis. The changes in the DNVs for NA12878 suggest the dynamic nature  
37 of the DNVs, demonstrating that the number is increasing over time.

38  
39 When mutational signature analysis was performed on this new set of DNVs, the most  
40 common mutation signatures were those seen in B-cell lymphomas. This is important as the  
41 lymphoblastoid cell lines are generated from B-cells and points to a non-random accumulation of  
42 mutations that are in line with the development of cancer in this cell type. In particular, we  
43 identified mutations in key DNA repair genes as well as a statistically significant excess of DNVs  
44 in *IGLL5* (45, 46). This gene is found to be mutated in B-cell lymphomas and protein-coding  
45 DNVs are identified in 27 individuals in this cohort; all of which have >100 overall DNVs. From  
46 our work, we identify two contributing factors causing these higher levels of DNVs, one is the

1 mutation of DNA repair genes while the second is an excess of Epstein-Barr Viral load. Future  
2 work using long-read sequencing and *de novo* assemblies will be imperative to identify complete  
3 viral integration in these genomes as integration sites can have impacts on cell line stability. One  
4 unexpected consequence of B-cell lymphoma mutation signatures in some individuals from the  
5 1000 Genomes Project would be a new pathway to study the mechanisms and biology of the  
6 development of this cancer.

7  
8 In addition to the DNA repair gene DNVs, we identified fifteen pathogenic or likely-  
9 pathogenic DNVs that had already been implicated in a database of clinical variation (Clinvar).  
10 This calls into question the use of the 1000 Genomes Project data as a control for both B-cell  
11 lymphomas and more generally for DNVs identified in clinical patients. More importantly, the  
12 extensive spectrum of DNVs that can appear in a cell line call into question the use of control  
13 samples derived from lymphoblastoid cell lines. Currently, to our knowledge the Genome in a  
14 Bottle and Human Pangenome Reference Consortium (HPRC) are building reference databases  
15 and pangenomes using DNA from lymphoblastoid cell lines. Although it does seem that the use  
16 of blood for some samples was at least initially discussed for the HPRC  
17 ([https://www.genome.gov/Pages/Research/Sequencing/Meetings/HGR\\_Webinar\\_Summary\\_Mar](https://www.genome.gov/Pages/Research/Sequencing/Meetings/HGR_Webinar_Summary_March1_2018.pdf)  
18 [ch1\\_2018.pdf](https://www.genome.gov/Pages/Research/Sequencing/Meetings/HGR_Webinar_Summary_March1_2018.pdf)), it does appear the project has defaulted to using lymphoblastoid cell lines. We find  
19 it is imperative that these efforts consider utilizing native DNA isolated from blood as the source  
20 or utilize a family-based design to identify and remove DNVs. In this way, the highest quality  
21 references can be built that will stand the test of time. Finally, we recommend that much like the  
22 Simons Simplex Collection, that studies assessing DNVs in individuals with a particular phenotype  
23 of interest, also sequence DNA from blood cells and not DNA post-culturing of lymphoblastoid  
24 cell lines.



## 1 **Online Methods**

### 2 *Trio dataset*

3 A total of 602 trios from the 1000 Genomes Project were sequenced at the New York  
4 Genome Center as described previously (2). The aligned data files (crams) are located at  
5 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000G\\_2504\\_high\\_coverage/1000G\\_](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/1000G_2504_high_coverage.sequence.index)  
6 [2504\\_high\\_coverage.sequence.index](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/1000G_2504_high_coverage.sequence.index) and  
7 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000G\\_2504\\_high\\_coverage/1000G\\_](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/1000G_698_related_high_coverage.sequence.index)  
8 [698\\_related\\_high\\_coverage.sequence.index](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/1000G_698_related_high_coverage.sequence.index). Description of the 602 trios is found in Supplemental  
9 Table S2.

10

### 11 *Single-nucleotide variant and insertion/deletion calling*

12 The NVIDIA Parabricks program version 3.0.0 was utilized to call single-nucleotide  
13 variants (SNVs) and small insertions/deletions (indels) with GATK (36) version 4.1.0 and  
14 Google's DeepVariant (37) version 0.10. For each individual, a GVCF was generated for these  
15 two variant callers. The GVCFs were then genotyped, on a per trio basis, using the GLnexus (38)  
16 version 1.2.6 genotyper. Post-calling, we checked the counts of all variants and heterozygous  
17 variants per chromosome in each individual (Figure S5).

18

### 19 *de novo variant calling*

20 DNVs were called by identifying all DNVs in GATK and DeepVariant based on the parent  
21 and child genotypes, respectively. The intersection of these DNVs was then identified from the  
22 two callers and filtering was carried out as follows: genotype quality greater than 20, depth of  $\geq$   
23 10, an allele balance  $> 0.25$ , and no presence of the DNV in any reads in the parents. DNVs in low  
24 complexity regions, centromeres, and recent repeats were removed from further analysis. To assess  
25 the DNVs, we manually scored 3980 sites with SAMtools version 1.9 tview. To score these sites,  
26 we looked at the first column (variant location in tview images) of both parents and the proband  
27 sample to see what mutations were present. If there was any mutation in the first column of the  
28 mother or father, regardless of quality, that matched the main mutation in the proband's first  
29 column, then we denoted the mutation as maternal, paternal, or both depending on whether it was  
30 the mother's mutation that matched the proband or the fathers or both parents. If the main mutation  
31 in the first column of the parental samples did not match the proband's mutation, then we knew  
32 this sample would be de novo, thus verifying our results.

33

### 34 *Phasing of de novo variants*

35 We utilized Unfazed v1.0.2 (<https://github.com/jbelyeu/unfazed>) (47) to phase the *de novo* variants  
36 in our study with regard to the parent-of-origin chromosome. First, a bed file containing *de novo*  
37 variants was generated for each individual. Second, the *de novo* bed file, DeepVariant full genome  
38 trio VCF, and the alignment files for all trio members were run through Unfazed. Since Unfazed  
39 uses different approaches to phasing on the X chromosome in males and females, we only focus  
40 on phased variants on the autosomes in this study.

41

### 42 *NA12878 additional datasets*

43 We identified additional high-coverage whole-genome sequencing data from NA12878  
44 from the SRA (<https://www.ncbi.nlm.nih.gov/sra>) and other sources. These included SRA data  
45 SRR944138 and SRR952827 both from 2013, McDonnell Genome Institute data  
46 gerald\_HFKWMDSXX and H\_IJ-NA12878 both from 2018, and the high-coverage data from

1 2020. To avoid differences due to coverage, we downsampled all datasets to 30x using SAMtools.  
2 All data was re-mapped to build 38 using SpeedSeq (48) version 0.1.2 and run through the de novo  
3 workflow using the NA12891 and NA12892 parental WGS data from 2020. We again did a count  
4 check for total and heterozygous variants per chromosome (Figure S6).

#### 5 6 *Phylogenetic tree of de novo variants*

7 To assess the differences between different NA12878 replicates we built a multi-sequence  
8 fasta file where each fasta represents the aggregate of all possible DNVs identified in this  
9 individual. The specific steps to build the tree were as follows: 1) we first merged the samples  
10 together and converted the genotypes for each de novo mutation from 0/0 or 0/1 to the nucleotide  
11 counterparts (e.g., AA, CG, TC) for all of the NA12878 samples; 2) next we converted these  
12 genotype symbols to their IUPAC code; 3) we then collapsed the IUPAC symbols into a sequence  
13 per sample and placed them into a fasta file. We also included a reference "sample", which was  
14 just the reference allele at each de novo mutation and 4) we used MEGAX (49) version 10.2.4 to  
15 create a maximum likelihood phylogenetic tree.

#### 16 17 *Mutation profile assessment*

18 We utilized the deconstructSig (42) software version-1.9.0 inside of Parabricks to perform  
19 mutation signature analysis. The prominent signature was chosen for an individual and if there was  
20 not one prominent signature than the weights of two signature was equal to or greater than ( $\geq$   
21 0.31) both signatures were represented in the tables and figures.

#### 22 23 *Karyotype analysis*

24 Read-depth based karyotypes were generated by assessment of the aligned sequence data. First,  
25 the number of reads per chromosome was calculated using SAMtools (50) in each individual.  
26 Second, the size of each chromosome was generated using the reference genome data and by  
27 removing locations of gaps from the reference. Third, the copy number of each of the  
28 chromosomes was calculated as follows: ((fold coverage per chromosome) / (fold coverage of  
29 chromosome 1))\*2.

#### 30 31 *Viral analysis*

32 We ran SAMtools idxstats on all individuals to determine the number of mapped reads to  
33 each chromosome. We then calculated the copy number of EBV in each individual as follows:  
34  $EBV \text{ copy number} = ((\text{mapped reads to EBV} * 150 \text{ base pairs per read}) / \text{length of EBV}) / ((\text{mapped}$   
35  $\text{reads to chromosome 1} * 150 \text{ base pairs per read}) / \text{length of chromosome 1})$

#### 36 37 *DNV enrichment in genes*

38 To test for DNV enrichment in genes we utilized two methods: chimpanzee-human and  
39 denovolyzeR. These were run as previously described (10, 51).

#### 40 41 *Annotation of protein-coding DNVs*

42 We uploaded the DNV calls to the open-cravat program (<https://opencravat.org/>) and  
43 specifically identified Clinvar as one of the annotation categories. Rescoring of DNVs in Franklin  
44 was performed using Franklin (<https://franklin.genoox.com>).

45  
46

## 1 ACKNOWLEDGMENTS

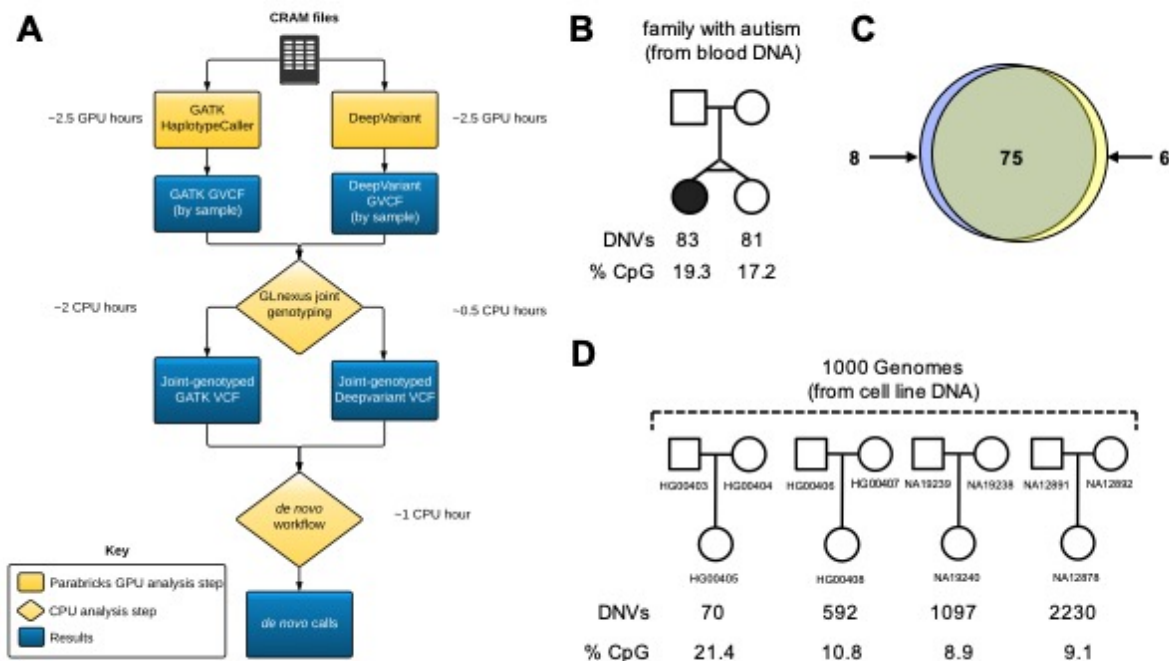
2 Funding: This work was supported by grants from the National Institutes of Health  
3 (R00MH117165 to T.N.T.). Author Contributions: J.N., T.H., C.P., and T.N.T. designed the study;  
4 J.N., P.V., E.F., E.M.P., Z.P., S.L., S.S., C.P., L.T., M.J., M.S., C.P. and T.N.T. performed  
5 analyses; J.N. and T.N.T. wrote the paper; and all authors reviewed and edited the paper.  
6 Competing interests: P.V., M.W., G.V. And T.T.H are full time employees of NVIDIA; and Data  
7 and materials availability: We provide family-level VCFs, DNV calls, and phased DNV results for  
8 the 602 trios in this study as a public community resource (Globus endpoint: “Turner Lab at  
9 WashU - DNV in 1000 Genomes Paper”, direct link: [https://app.globus.org/file-](https://app.globus.org/file-manager?origin_id=3eff453a-88f4-11eb-954f-752ba7b88e8e&origin_path=%2F)  
10 [manager?origin\\_id=3eff453a-88f4-11eb-954f-752ba7b88e8e&origin\\_path=%2F](https://app.globus.org/file-manager?origin_id=3eff453a-88f4-11eb-954f-752ba7b88e8e&origin_path=%2F)). The raw  
11 alignment files for the high-coverage 1000 Genomes Project data are available at the following  
12 website:  
13 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000G\\_2504\\_high\\_coverage/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/). The  
14 following cell lines/DNA samples were obtained from the NIGMS Human Genetic Cell  
15 Repository at the Coriell Institute for Medical Research: NA06984, NA06985, NA06986,  
16 NA06989, NA06994, NA07000, NA07037, NA07048, NA07051, NA07056, NA07347,  
17 NA07357, NA10847, NA10851, NA11829, NA11830, NA11831, NA11832, NA11840,  
18 NA11843, NA11881, NA11892, NA11893, NA11894, NA11918, NA11919, NA11920,  
19 NA11930, NA11931, NA11932, NA11933, NA11992, NA11994, NA11995, NA12003,  
20 NA12004, NA12005, NA12006, NA12043, NA12044, NA12045, NA12046, NA12058,  
21 NA12144, NA12154, NA12155, NA12156, NA12234, NA12249, NA12272, NA12273,  
22 NA12275, NA12282, NA12283, NA12286, NA12287, NA12340, NA12341, NA12342,  
23 NA12347, NA12348, NA12383, NA12399, NA12400, NA12413, NA12414, NA12489,  
24 NA12546, NA12716, NA12717, NA12718, NA12748, NA12749, NA12750, NA12751,  
25 NA12760, NA12761, NA12762, NA12763, NA12775, NA12776, NA12777, NA12778,  
26 NA12812, NA12813, NA12814, NA12815, NA12827, NA12828, NA12829, NA12830,  
27 NA12842, NA12843, NA12872, NA12873, NA12874, NA12878, NA12889, NA12890,  
28 NA12376, NA10838, NA12329, NA10852, NA10840, NA12386, NA12864, NA12801,  
29 NA12344, NA10861, NA07029, NA12753, NA12832, NA12485, NA12802, NA12739,  
30 NA10856, NA10845, NA12818, NA10831, NA12766, NA10864, NA10843, NA12877,  
31 NA12335, NA12817, NA12752, NA12767, NA10855, NA12707, NA10857, NA10839,  
32 NA12740, NA10837, NA10836, NA07348, NA11993, NA12057, NA11839, NA06993,  
33 NA07014, NA06995, NA12146, NA12865, NA10859, NA06991, NA12336, NA10860,  
34 NA12145, NA07045, NA07349, NA07031, NA07345, NA12891, NA07055, NA07435,  
35 NA10835, NA12274, NA12875, NA10842, NA12239, NA10830, NA12056, NA11917,  
36 NA12892, NA06997, NA07022, NA12264, NA11891, NA07034, NA12248, NA10865,  
37 NA10863, NA10854, NA11882, NA07346, NA07019, NA12343, NA10846. These data were  
38 generated at the New York Genome Center with funds provided by NHGRI Grant  
39 3UM1HG008901-03S1. We are grateful to all of the families at the participating SSC sites, as well  
40 as the principal investigators (A. Beaudet, R. Bernier, J. Constantino, E. Cook, E. Fombonne, D.  
41 Geschwind, R. Goin-Kochel, E. Hanson, D. Grice, A. Klin, 25D. Ledbetter, C. Lord, C. Martin,  
42 D. Martin, R. Maxim, J. Miles, O. Ousley, K. Pelphrey, B. Peterson, J. Piggot, C. Saulnier, M.  
43 State, W. Stone, J. Sutcliffe, C. Walsh, Z. Warren, and E. Wijsman). We appreciate obtaining  
44 access to phenotypic and genetic data for the monozygotic twin pair on SFARI Base. Approved  
45 researchers can obtain the SSC population dataset described in this study  
46 (<https://www.sfari.org/resource/simons-simplex-collection/>) by applying at <https://base.sfari.org>.



1 **FIGURES**

2

Figure 1



3

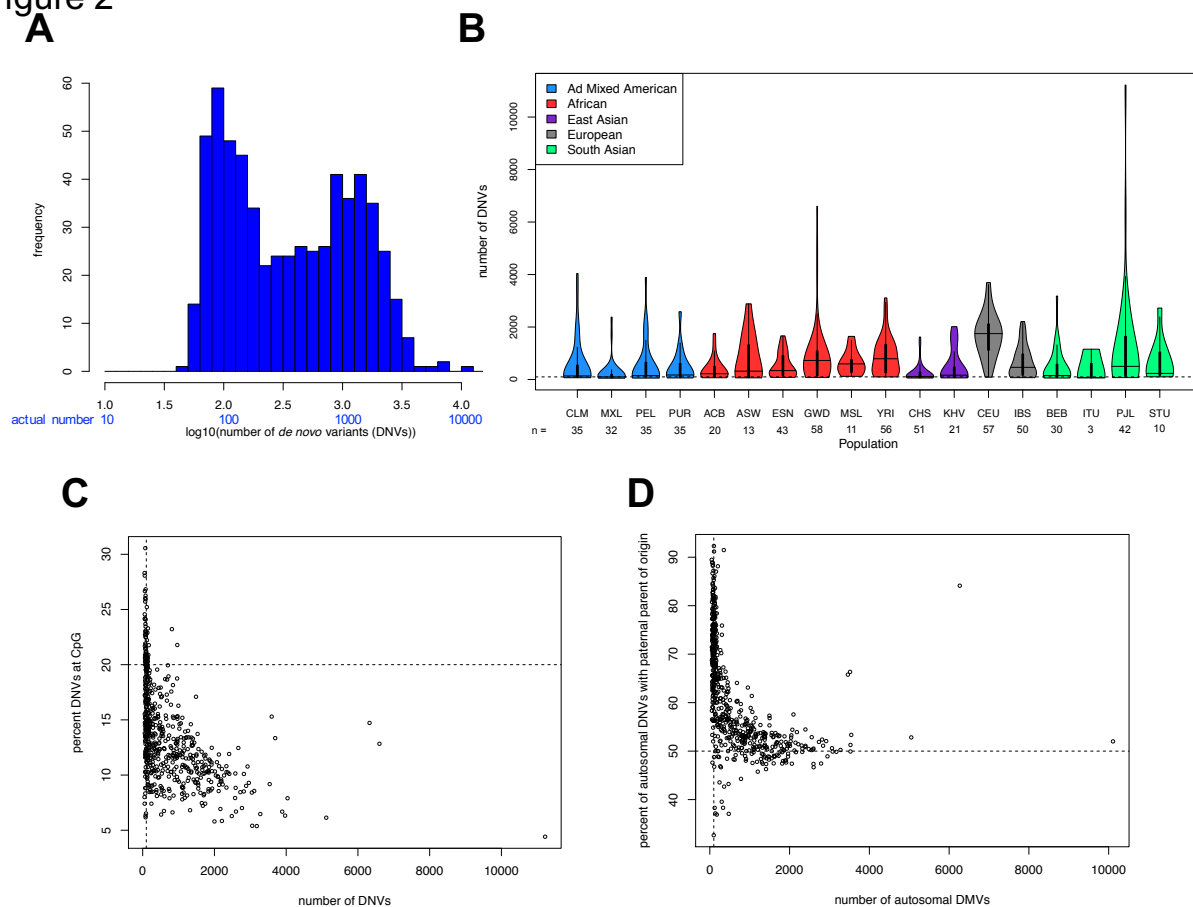
4 **Figure 1:** DNV calling in short-read whole-genome sequencing data. A) *de novo* workflow for  
 5 detection of DNVs from aligned read files (crams); B and C) Benchmarking DNV workflow in a  
 6 monozygotic twin pair sequenced from DNA derived from blood; D) DNV detecting in four trios  
 7 in the 1000 Genomes Project.

8

9

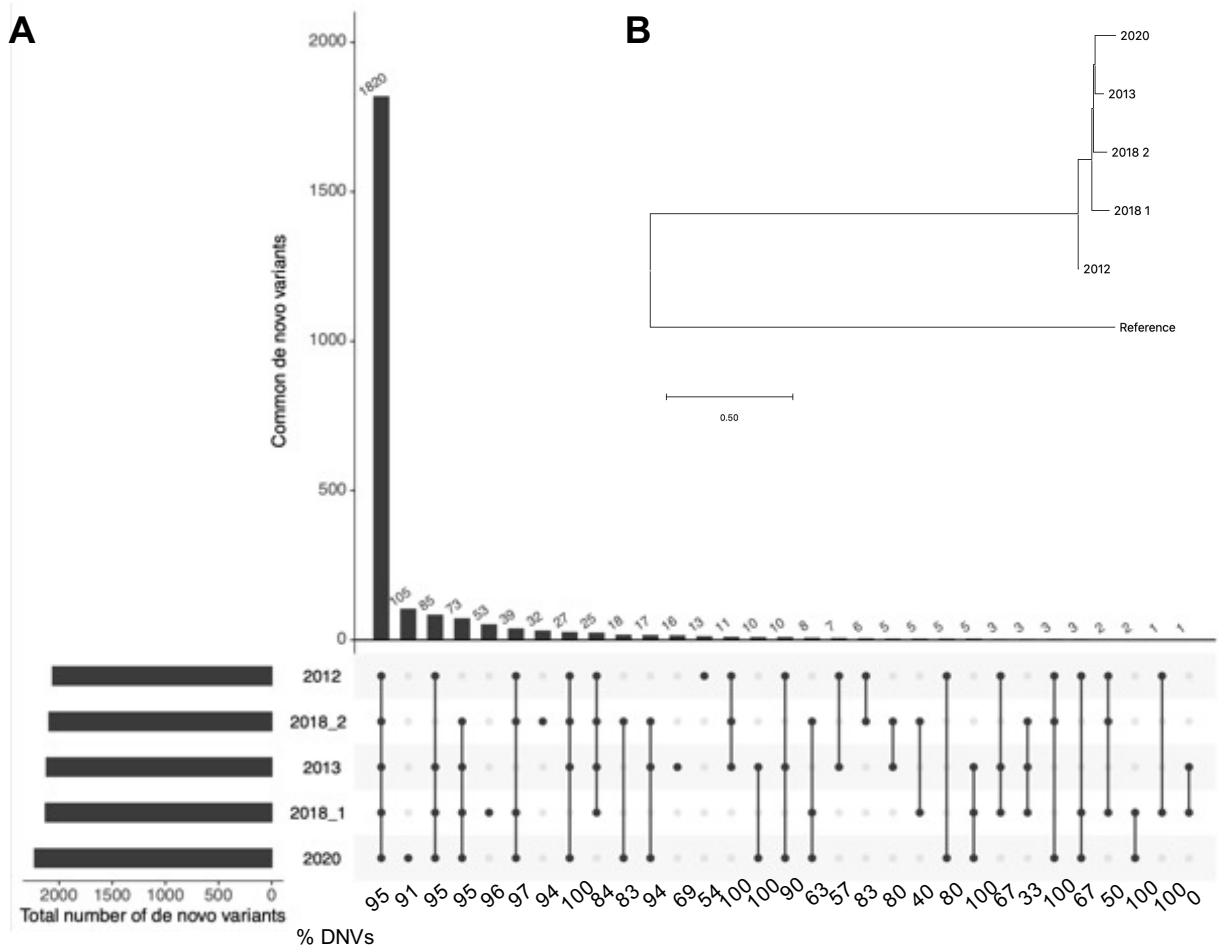
10

Figure 2



1  
 2 **Figure 2:** Characteristics of DNVs detected in 602 trios. A) Histogram of DNV counts in 602  
 3 trios; B) DNV counts by population; C) Percent of DNVs at CpG sites versus the total number of  
 4 DNVs; D) Percent of DNVs phased to have a paternal parent-of-origin versus the total number of  
 5 DNVs.

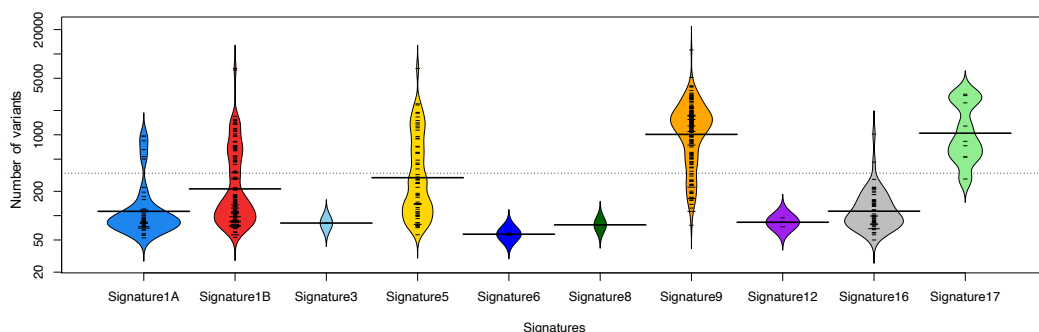
Figure 3



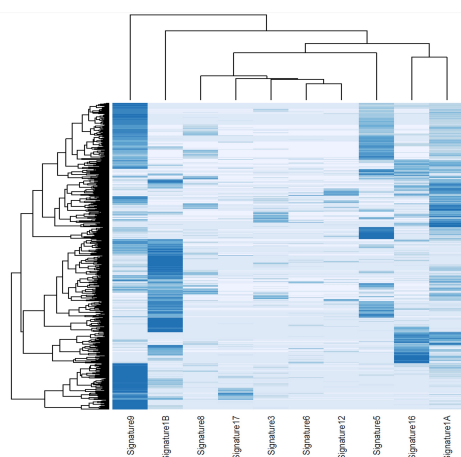
1  
2  
3  
4  
5  
6

**Figure 3:** Assessment of five replicates of NA12878. A) UpSet plot demonstrating the number of variants detected in the replicates (at the bottom of the plot the percent of true DNVs is listed for each category); B) Phylogenetic tree of the five replicates

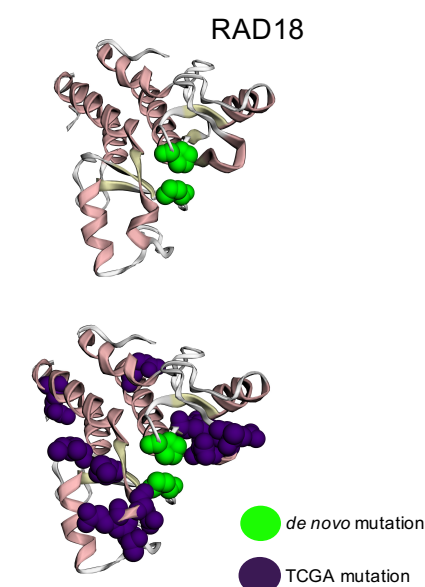
Figure 4  
A



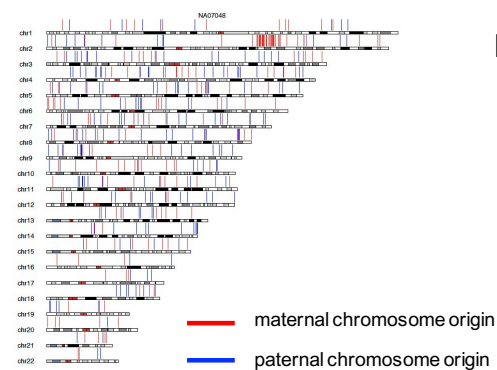
B



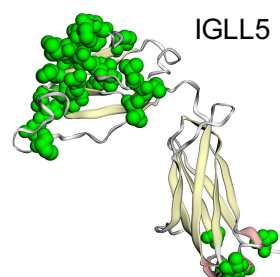
C



D



E



1  
2  
3

4 **Figure 4:** Mutational properties of DNVs. A) Mutation signature analysis showing the total  
5 number of DNVs and the individuals with each signature type; B) Heatmap of individuals based  
6 on their mutational signatures; C) Mutations in the DNA repair gene RAD18 shown on their 3D  
7 structure (and modeled using mupit). Also, shown are known cancer mutations from The Cancer  
8 Genome Atlas (TCGA); D) Location of DNVs based on their phased parent-of-origin in



- 1 NA07048. Most notable there are a cluster of mutations on the maternal chromosome on
- 2 chromosome 2; E) DNVs in IGLL5 shown on their 3D structure (and modeled using mupit).
- 3

1 **TABLES**

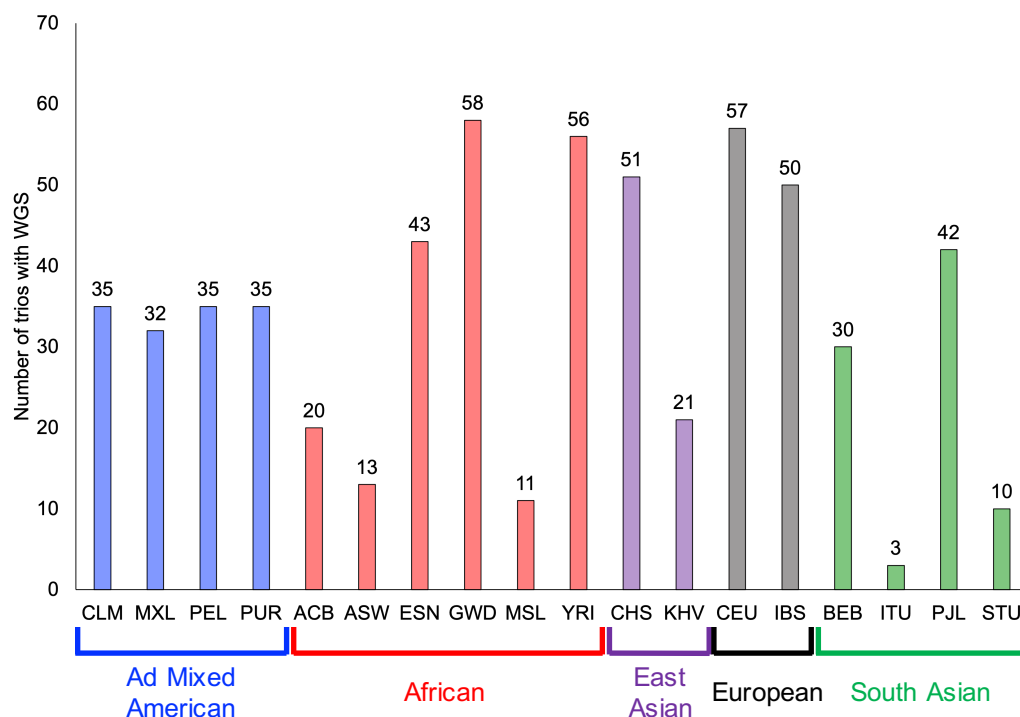
2 **Table 1.** DNVs in DNA damage repair genes and clinically relevant variants

Category	individual	<i>de novo</i> variant	variant type	gene
DNA damage repair gene	HG01074	chr3_48447050_C_G	missense	<i>ATRIP</i>
	HG02841	chr16_89799603_A_G	splice_donor	<i>FANCA</i>
	HG03200	chr16_89762010_C_T	missense	<i>FANCA</i>
	HG01126	chr11_22625482_T_G	missense	<i>FANCF</i>
	NA18875	chr5_80654794_G_A	missense	<i>MSH3</i>
	HG02650	chr6_31759121_C_T	missense	<i>MSH5</i>
	NA10838	chr11_65865247_C_T	missense	<i>MUS81</i>
	NA10838	chr8_42357362_AT_A	frameshift	<i>POLB</i>
	NA19677	chr19_50407375_G_A	missense	<i>POLD1</i>
	HG01096	chr12_132634327_C_T	missense	<i>POLE</i>
	HG01755	chr15_89321792_C_T	missense	<i>POLG</i>
	NA12864	chr3_8958938_G_T	missense	<i>RAD18</i>
	HG02683	chr15_40729853_T_C	missense	<i>RAD51</i>
	HG02630	chrX_96884884_G_A	missense	<i>RPA4</i>
	NA19919	chr3_133644039_A_G	missense	<i>TOPBP1</i>
	HG04115	chr8_31120294_C_T	missense	<i>WRN</i>
	NA19161	chr8_31124967_G_T	missense	<i>WRN</i>
Clinvar pathogenic / likely pathogenic	HG03795	chr11_22274728_C_T	stop_gained	<i>ANO5</i>
	NA10854	chr13_110179298_C_T	missense	<i>COL4A1</i>
	NA10842	chr10_49530737_G_A	stop_gained	<i>ERCC6</i>
	HG02668	chr1_111787063_C_T	missense	<i>KCND3</i>
	HG02466	chr1_39485559_G_A	missense	<i>MACF1</i>
	HG02129	chr11_77206108_G_A	missense	<i>MYO7A</i>
	HG03122	chr11_45949458_G_A	stop_gained	<i>PHF21A</i>
	NA12707	chr6_52058438_C_T	missense	<i>PKHD1</i>
	HG01755	chr15_89321792_C_T	missense	<i>POLG</i>
	HG02892	chr11_124875581_C_T	stop_gained	<i>ROBO3</i>
	HG03635	chr2_165310406_G_A	missense	<i>SCN2A</i>
	NA10830	chr2_39023106_C_T	missense	<i>SOS1</i>
	NA10831	chr3_24143512_G_A	missense	<i>THRB</i>
	HG01629	chr2_209775898_C_T	stop_gained	<i>UNC80</i>
	HG00558	chr16_88435401_G_A	missense	<i>ZNF469</i>

3  
4  
5  
6

## 1 SUPPLEMENTARY FIGURES

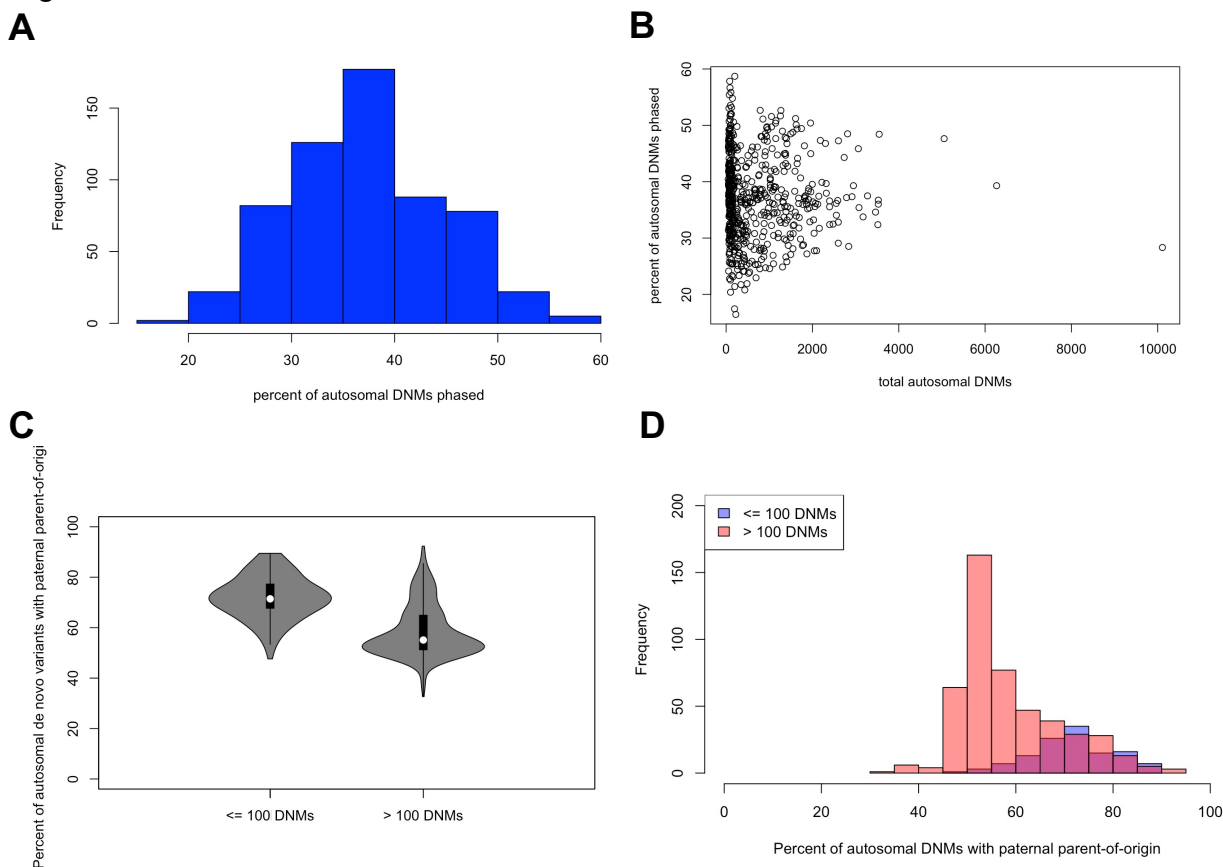
### Figure S1



2  
3 **S1 Sample population distribution.** A distribution of the populations by super and sub  
4 populations defined by the 1000 Genomes Project. Blue represents Ad Mixed American super  
5 population, red represents African super population, purple is East Asian super population, grey  
6 represents European super population, and green represents South Asian super population.

7  
8

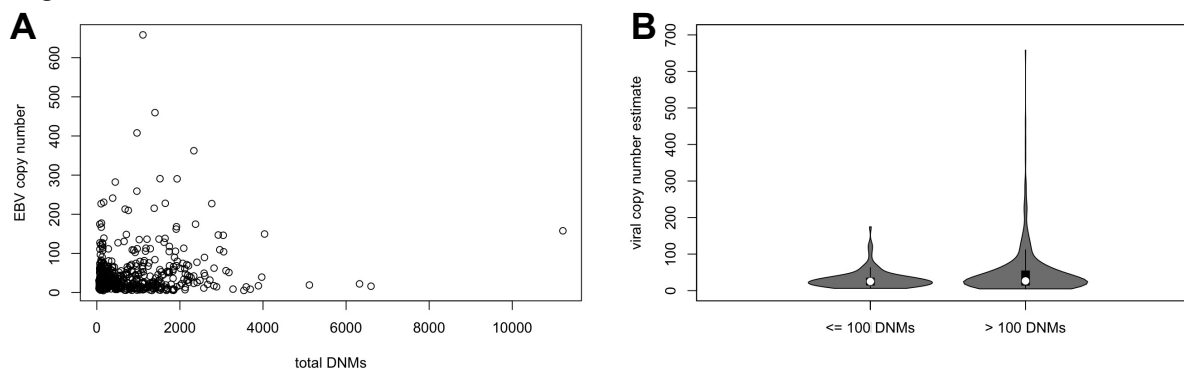
Figure S2



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11

**S2 Phase distribution of DNM.** Histogram of the percent autosomal DNMs that was fixed A) distribution of percent autosomal DNMs phased vs. total autosomal DNMs.; B) Violin plot of percent of autosomal DNMs that with paternal parent-of-origin.; C) Distribution of DNMs with paternal parent-of-origin. The pink graph represents samples that had greater than 100 DNMs. The blue graph represents samples that had less than or equal to 100 DNMs. There is a trend of higher percent of DNMs with paternal parent-of-origin compared in the group that had less than or equal 100 DNMs compared to those with greater than 100 DNMs, which would be expected if the DNMs are real.

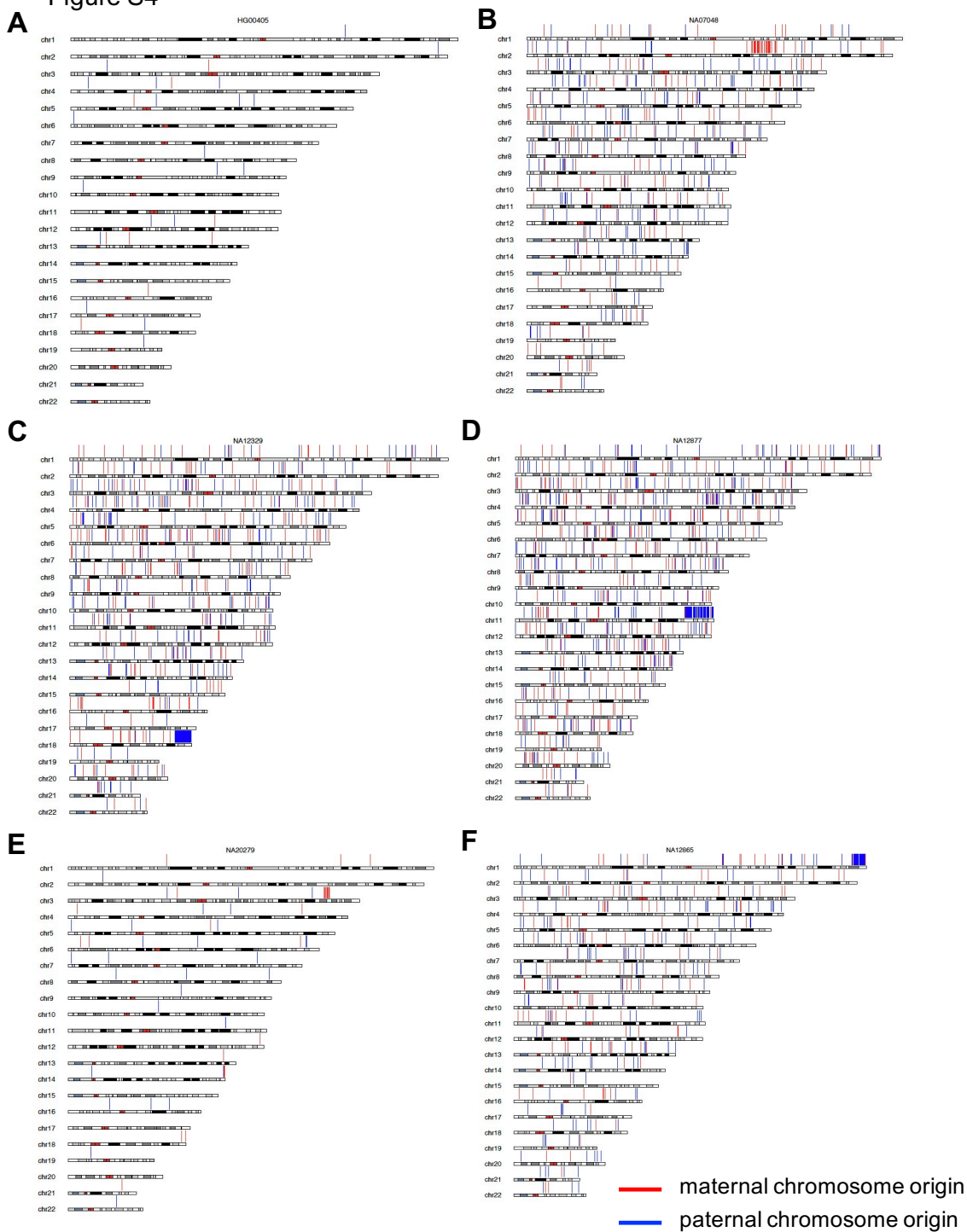
Figure S3



1  
2  
3  
4  
5  
6  
7

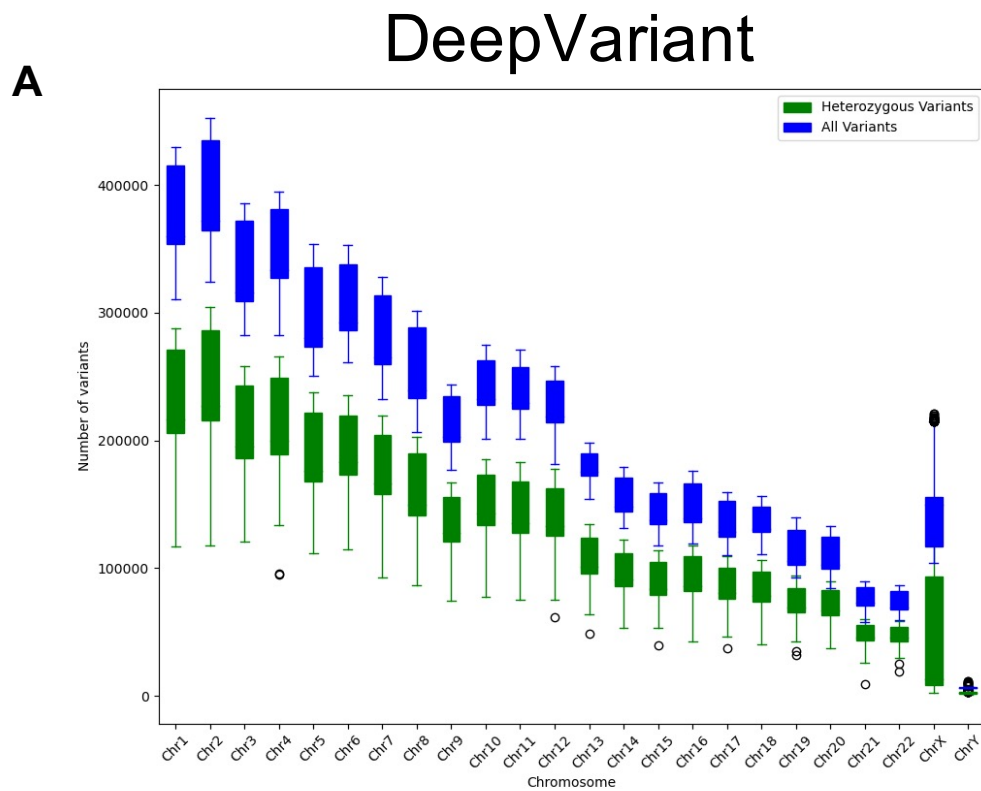
**S3 EBV Copy Number Estimation.** A) distribution of estimated EBV copy number vs total number of DNMs. There was a minor correlation between EBV copy number and total DNMs ( $p = 2.32e-05$ ,  $r = 0.17$ ).; B) violin plot comparing estimated EBV copy number and groups of less than or equal to 100 DNMs and greater than 100 DNMs.

Figure S4

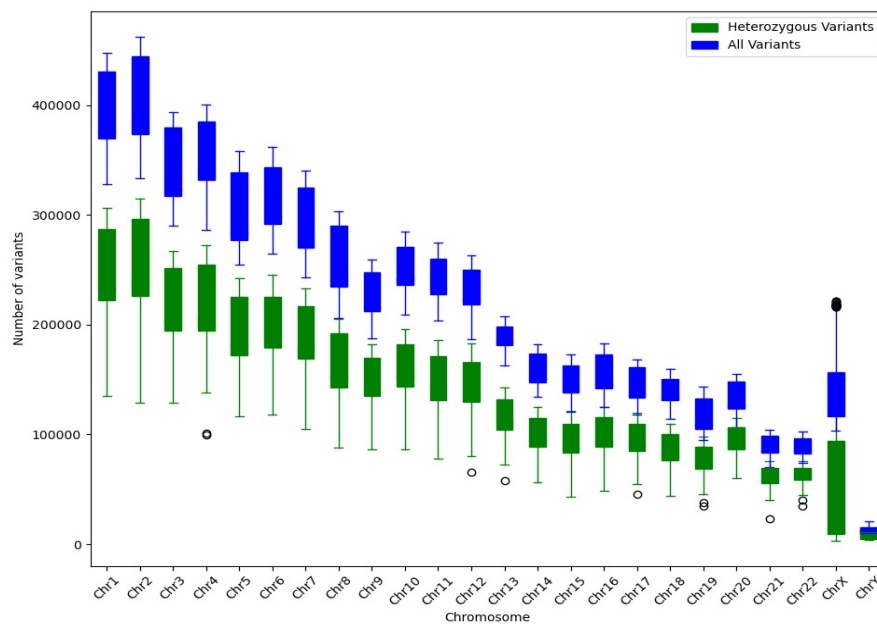


1  
2 **S4 Phase distributions across autosomal chromosomes.** Distribution of phased DNMs across  
3 6 different samples. The blue marks represent paternal parent-of-origin and the red mark  
4 represents the maternal parent-of-origin. A) represents a normal distribution, B) represents a  
5 clustering of maternal mutations, C), D) and F) represent paternal clustering of mutations, and E)  
6 represents a clustering of maternal mutations possible event that could lead to further DNVs.  
7

Figure S5



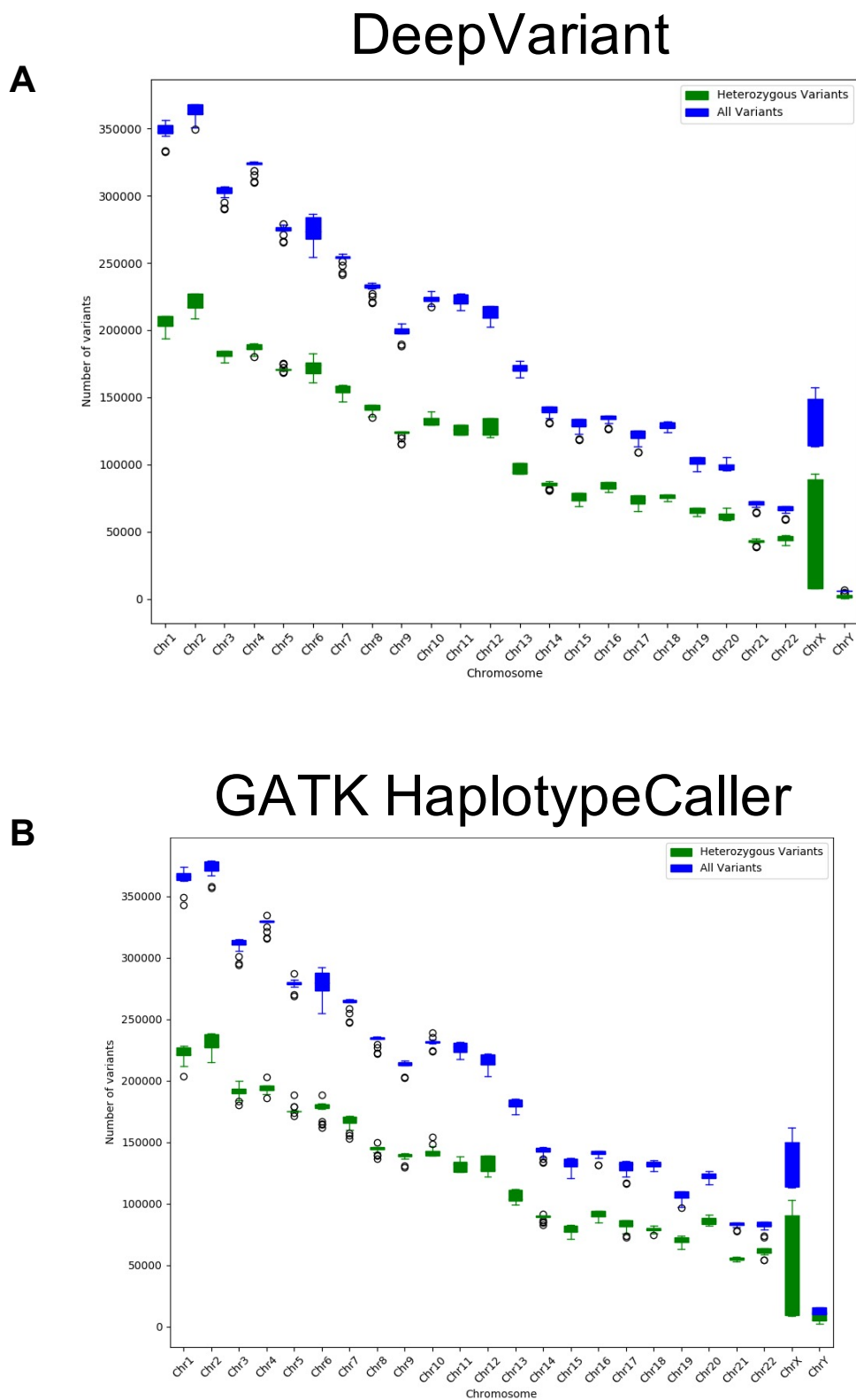
### B GATK HaplotypeCaller



1 **S5 Total counts of variants for DeepVariant and Haplotypcaller trio calls.** Distribution of  
2 variants by chromosome found in the trios from the 602 1000 Genomes Project families. The  
3 total number of variants in blue, and heterozygous variants, in green.  
4



Figure S6



1 **S6 Total counts of variants for DeepVariant and Haplotypcaller downsampled 30x**  
2 **coverage NA12878 trio calls.** Distribution of variants by chromosome found in the trios from  
3 the various NA12878 samples, downsampled to 30x coverage, as well as samples NA12891 and  
4 NA12892. The total number of variants in blue, and heterozygous variants, in green.  
5  
6

## 1 REFERENCES

- 2
- 3 1. A. Auton *et al.*, A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
- 4 2. M. Byrska-Bishop *et al.*, High coverage whole genome sequencing of the expanded 1000  
5 Genomes Project cohort including 602 trios. *bioRxiv*, 2021.2002.2006.430068 (2021).
- 6 3. L. Séguirel, M. J. Wyman, M. Przeworski, Determinants of mutation rate variation in the human  
7 germline. *Annual review of genomics and human genetics* **15**, 47-70 (2014).
- 8 4. P. Feliciano *et al.*, Exome sequencing of 457 autism families recruited online provides evidence  
9 for autism risk genes. *NPJ genomic medicine* **4**, 19 (2019).
- 10 5. I. Iossifov *et al.*, The contribution of de novo coding mutations to autism spectrum disorder.  
11 *Nature*, (2014).
- 12 6. B. J. O'Roak *et al.*, Exome sequencing in sporadic autism spectrum disorders identifies severe de  
13 novo mutations. *Nature genetics* **43**, 585-589 (2011).
- 14 7. B. J. O'Roak *et al.*, Sporadic autism exomes reveal a highly interconnected protein network of de  
15 novo mutations. *Nature* **485**, 246-250 (2012).
- 16 8. T. N. Turner *et al.*, Genomic Patterns of De Novo Mutation in Simplex Autism. *Cell* **171**, 710-  
17 722.e712 (2017).
- 18 9. T. N. Turner *et al.*, Genome Sequencing of Autism-Affected Families Reveals Disruption of  
19 Putative Noncoding Regulatory DNA. *American journal of human genetics* **98**, 58-74 (2016).
- 20 10. T. N. Turner *et al.*, Sex-Based Analysis of De Novo Variants in Neurodevelopmental Disorders.  
21 *American journal of human genetics* **105**, 1274-1285 (2019).
- 22 11. S. J. Sanders *et al.*, De novo mutations revealed by whole-exome sequencing are strongly  
23 associated with autism. *Nature* **485**, 237-241 (2012).
- 24 12. S. De Rubeis *et al.*, Synaptic, transcriptional and chromatin genes disrupted in autism. *Nature*  
25 **515**, 209-215 (2014).
- 26 13. K. L. Helbig *et al.*, Diagnostic exome sequencing provides a molecular diagnosis for a significant  
27 proportion of patients with epilepsy. *Genetics in medicine : official journal of the American*  
28 *College of Medical Genetics*, (2016).
- 29 14. A. S. Allen *et al.*, De novo mutations in epileptic encephalopathies. *Nature* **501**, 217-221 (2013).
- 30 15. J. Kaplanis *et al.*, Integrating healthcare and research genetic data empowers the discovery of 28  
31 novel developmental disorders. *bioRxiv*, 797787 (2020).
- 32 16. DDD, Prevalence and architecture of de novo mutations in developmental disorders. *Nature* **542**,  
33 433-438 (2017).
- 34 17. J. de Ligt *et al.*, Diagnostic exome sequencing in persons with severe intellectual disability. *The*  
35 *New England journal of medicine* **367**, 1921-1929 (2012).
- 36 18. C. Gilissen *et al.*, Genome sequencing identifies major causes of severe intellectual disability.  
37 *Nature* **511**, 344-347 (2014).
- 38 19. J. Homsy *et al.*, De novo mutations in congenital heart disease with neurodevelopmental and  
39 other congenital anomalies. *Science (New York, N.Y.)* **350**, 1262-1266 (2015).
- 40 20. S. Zaidi *et al.*, De novo mutations in histone-modifying genes in congenital heart disease. *Nature*  
41 **498**, 220-223 (2013).
- 42 21. A. Sifrim *et al.*, Distinct genetic architectures for syndromic and nonsyndromic congenital heart  
43 defects identified by exome sequencing. *Nature genetics* **48**, 1060-1065 (2016).
- 44 22. H. L. Rehm *et al.*, ACMG clinical laboratory standards for next-generation sequencing. *Genetics*  
45 *in medicine : official journal of the American College of Medical Genetics* **15**, 733-747 (2013).
- 46 23. S. Besenbacher *et al.*, Novel variation and de novo mutation rates in population-wide de novo  
47 assembled Danish trios. *Nature communications* **6**, 5969 (2015).
- 48 24. D. I. Boomsma *et al.*, The Genome of the Netherlands: design, and project goals. *European*  
49 *journal of human genetics : EJHG* **22**, 221-227 (2014).
- 50 25. A. Chesi *et al.*, Exome sequencing to identify de novo mutations in sporadic ALS trios. *Nature*  
51 *neuroscience* **16**, 851-855 (2013).

- 1 26. D. F. Conrad *et al.*, Variation in genome-wide mutation rates within and between human families. *Nature genetics* **43**, 712-714 (2011).
- 2
- 3 27. S. Dimassi *et al.*, Whole-exome sequencing improves the diagnosis yield in sporadic infantile
- 4 spasm syndrome. *Clinical genetics* **89**, 198-204 (2016).
- 5 28. S. Dong *et al.*, De novo insertions and deletions of predominantly paternal origin are associated
- 6 with autism spectrum disorder. *Cell reports* **9**, 16-23 (2014).
- 7 29. M. Fromer *et al.*, De novo mutations in schizophrenia implicate synaptic networks. *Nature* **506**,
- 8 179-184 (2014).
- 9 30. I. Iossifov *et al.*, De novo gene disruptions in children on the autistic spectrum. *Neuron* **74**, 285-
- 10 299 (2012).
- 11 31. H. Jonsson *et al.*, Parental influence on human germline de novo mutations in 1,548 trios from
- 12 Iceland. *Nature* **549**, 519-522 (2017).
- 13 32. S. E. McCarthy *et al.*, De novo mutations in schizophrenia implicate chromatin remodeling and
- 14 support a genetic overlap with autism and intellectual disability. *Molecular psychiatry* **19**, 652-
- 15 658 (2014).
- 16 33. G. Narzisi *et al.*, Accurate de novo and transmitted indel detection in exome-capture data using
- 17 microassembly. *Nature methods* **11**, 1033-1036 (2014).
- 18 34. A. Ramu *et al.*, DeNovoGear: de novo indel and point mutation discovery and phasing. *Nature*
- 19 *methods* **10**, 985-987 (2013).
- 20 35. K. R. Franke, E. L. Crowgey, Accelerating next generation sequencing data analysis: an
- 21 evaluation of optimized best practices for Genome Analysis Toolkit algorithms. *Genomics Inform*
- 22 **18**, e10 (2020).
- 23 36. A. McKenna *et al.*, The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
- 24 generation DNA sequencing data. *Genome research* **20**, 1297-1303 (2010).
- 25 37. R. Poplin *et al.*, A universal SNP and small-indel variant caller using deep neural networks.
- 26 *Nature biotechnology* **36**, 983-987 (2018).
- 27 38. T. Yun *et al.*, Accurate, scalable cohort variant calls using DeepVariant and GLnexus. *bioRxiv*,
- 28 2020.2002.2010.942086 (2020).
- 29 39. J. M. Zook *et al.*, Integrating human sequence data sets provides a resource of benchmark SNP
- 30 and indel genotype calls. (2014).
- 31 40. R. A. Gibbs *et al.*, The International HapMap Project. *Nature* **426**, 789-796 (2003).
- 32 41. T. A. Sasani *et al.*, Large, three-generation human families reveal post-zygotic mosaicism and
- 33 variability in germline mutation accumulation. *Elife* **8**, (2019).
- 34 42. R. Rosenthal, N. McGranahan, J. Herrero, B. S. Taylor, C. Swanton, DeconstructSigs: delineating
- 35 mutational processes in single tumors distinguishes DNA repair deficiencies and patterns of
- 36 carcinoma evolution. *Genome biology* **17**, 31 (2016).
- 37 43. M. J. Landrum *et al.*, ClinVar: improving access to variant interpretations and supporting
- 38 evidence. *Nucleic acids research* **46**, D1062-d1067 (2018).
- 39 44. G. R. Abecasis *et al.*, A map of human genome variation from population-scale sequencing.
- 40 *Nature* **467**, 1061-1073 (2010).
- 41 45. F. Chen, Y. Zhang, C. J. Creighton, Systematic identification of non-coding somatic single
- 42 nucleotide variants associated with altered transcription and DNA methylation in adult and
- 43 pediatric cancers. *NAR Cancer* **3**, zcab001 (2021).
- 44 46. L. Pedrosa *et al.*, Proposal and validation of a method to classify genetic subtypes of diffuse large
- 45 B cell lymphoma. *Scientific reports* **11**, 1886 (2021).
- 46 47. J. R. Belyeu, T. A. Sasani, B. S. Pedersen, A. R. Quinlan, Unfazed: parent-of-origin detection for
- 47 large and small *de novo* variants. *bioRxiv*, 2021.2002.2003.429658 (2021).
- 48 48. C. Chiang *et al.*, SpeedSeq: ultra-fast personal genome analysis and interpretation. *Nature*
- 49 *methods* **12**, 966-968 (2015).
- 50 49. S. Kumar, G. Stecher, M. Li, C. Knyaz, K. Tamura, MEGA X: Molecular Evolutionary Genetics
- 51 Analysis across Computing Platforms. *Molecular biology and evolution* **35**, 1547-1549 (2018).

- 1 50. H. Li *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford,*  
2 *England)* **25**, 2078-2079 (2009).
- 3 51. B. P. Coe *et al.*, Neurodevelopmental disease genes implicated by de novo mutation and copy  
4 number variation morbidity. *Nature genetics* **51**, 106-116 (2019).
- 5

# Figure 1

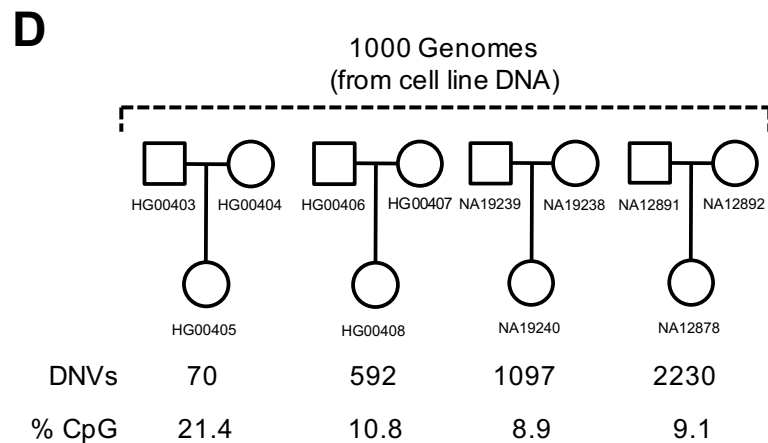
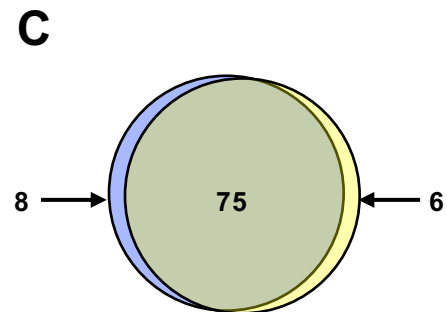
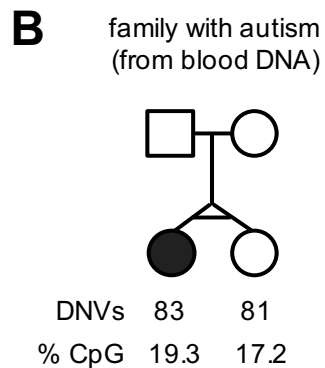
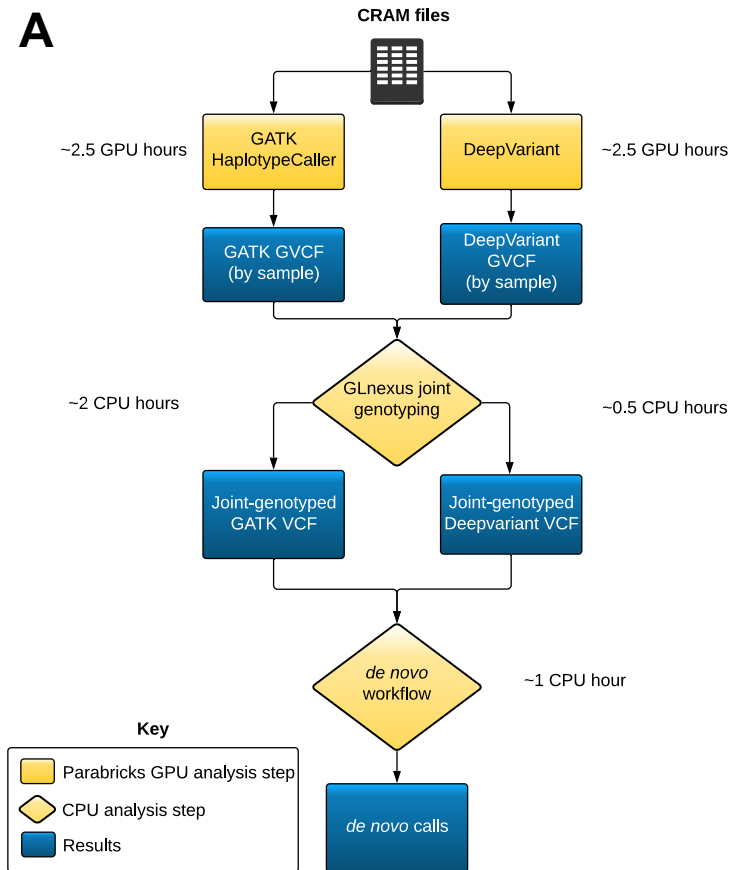
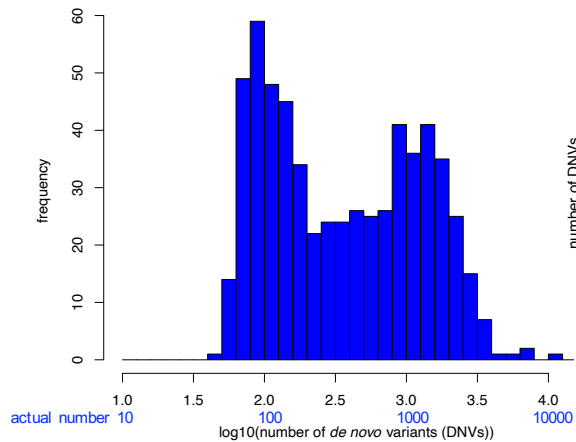
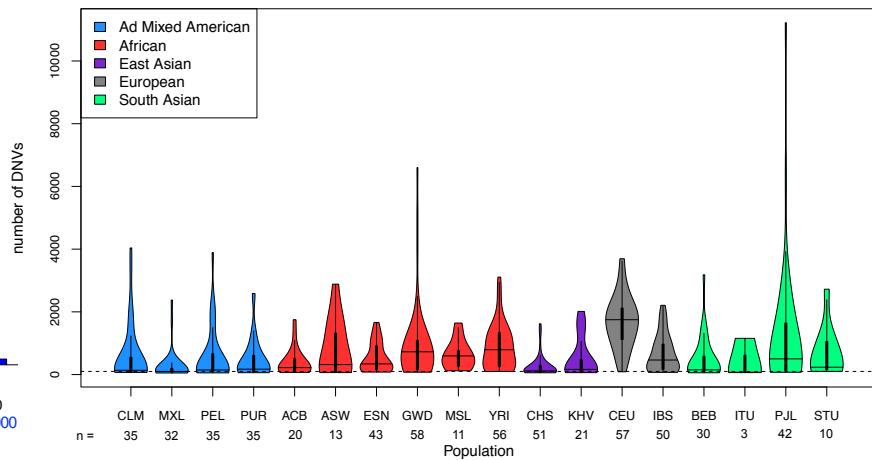


Figure 2

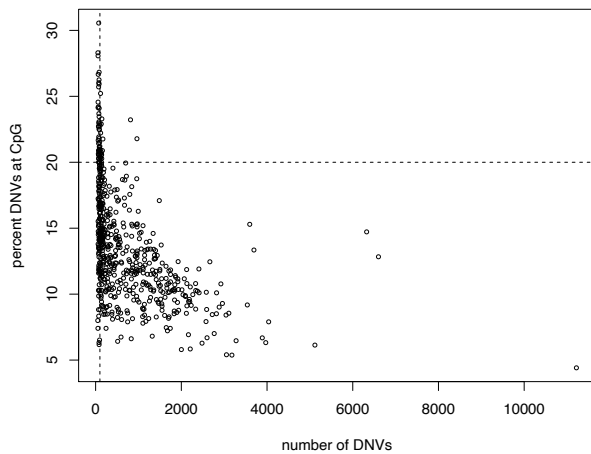
**A**



**B**



**C**



**D**

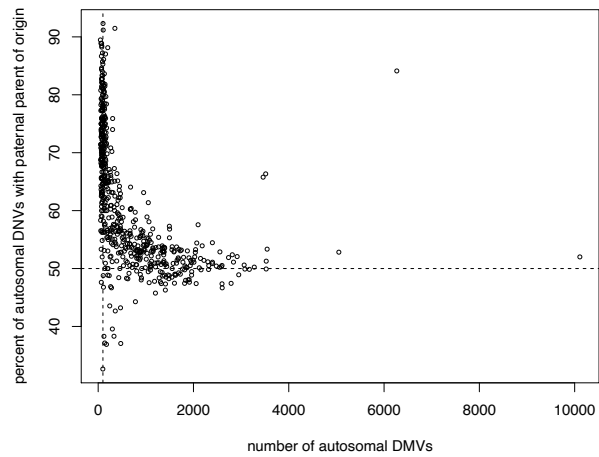
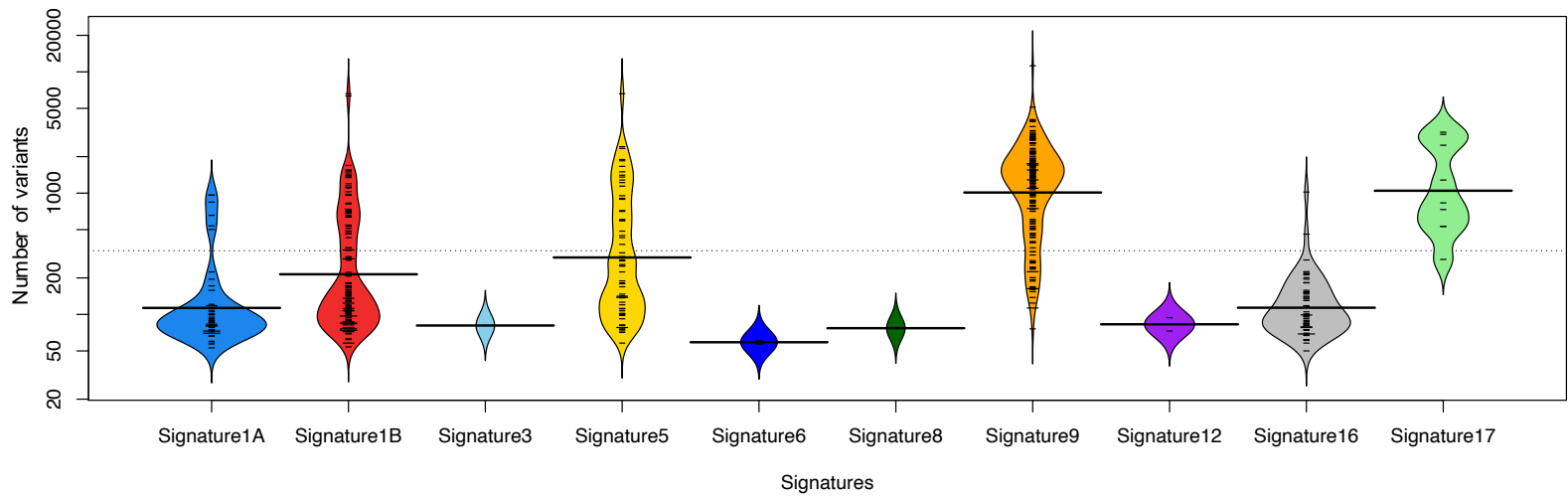




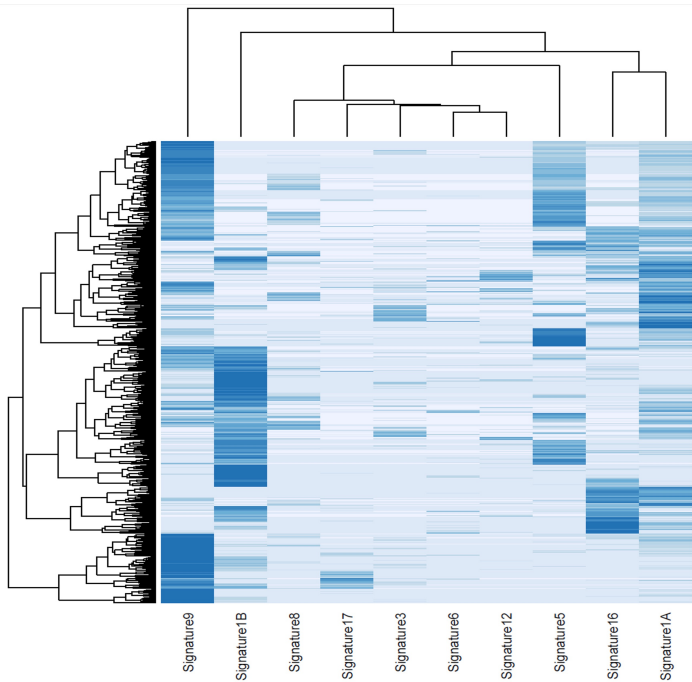


Figure 4

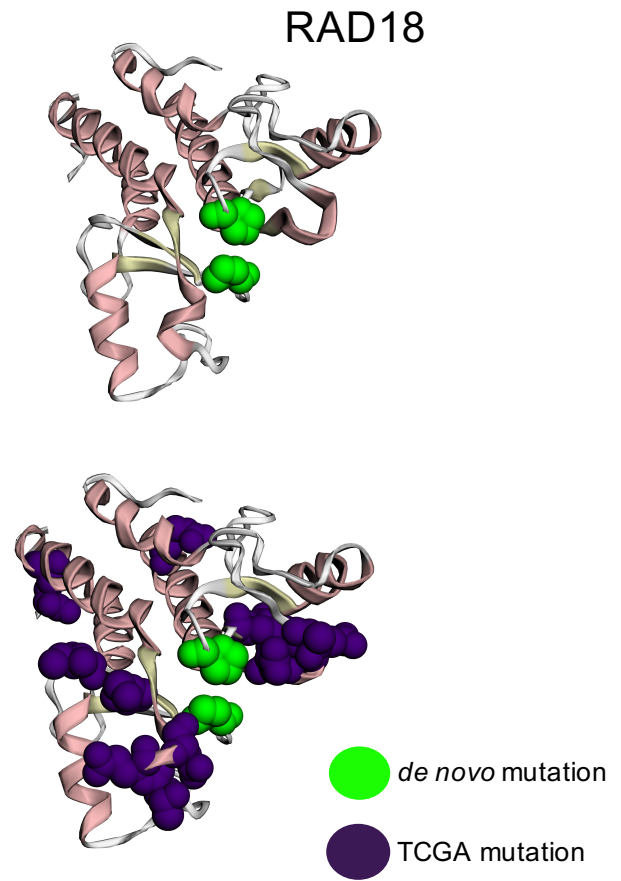
**A**



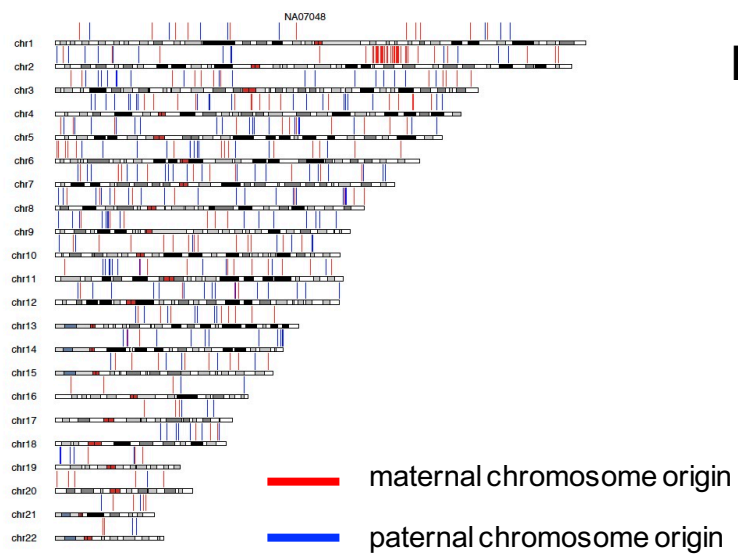
**B**



**C**



**D**



**E**

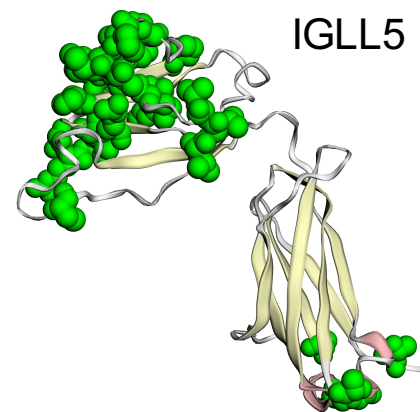


Figure S1

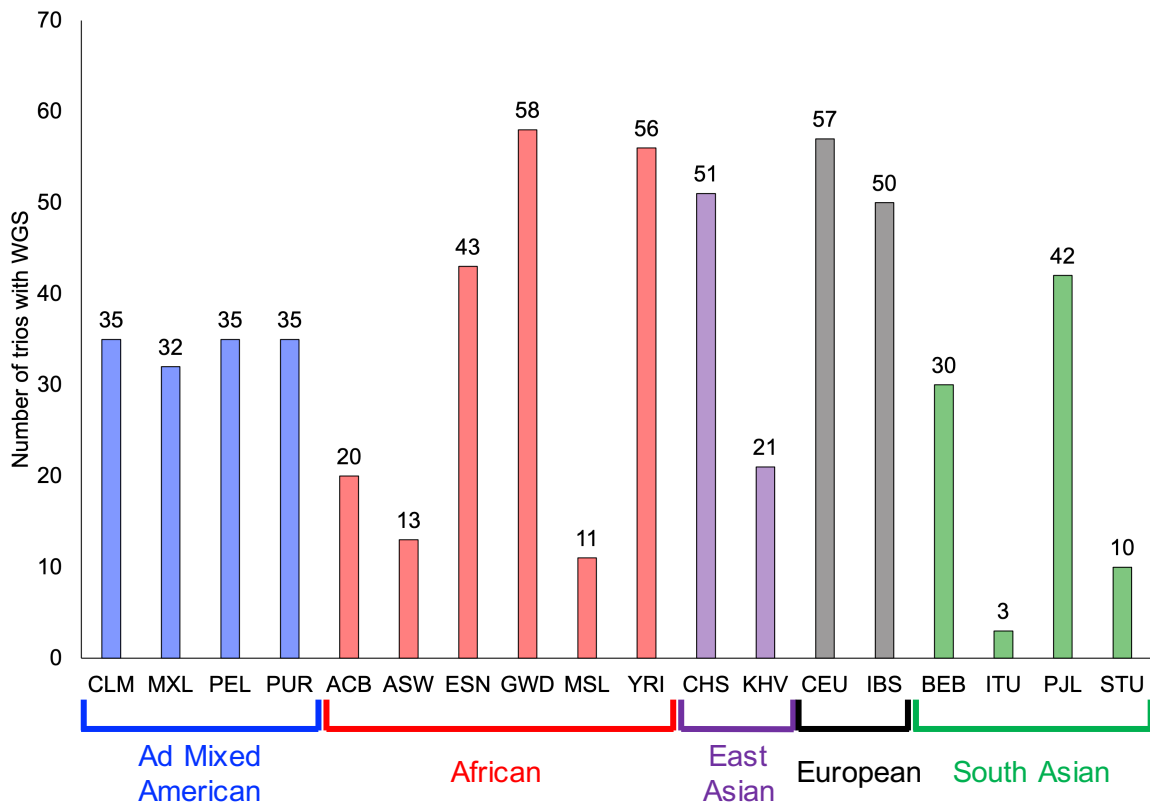
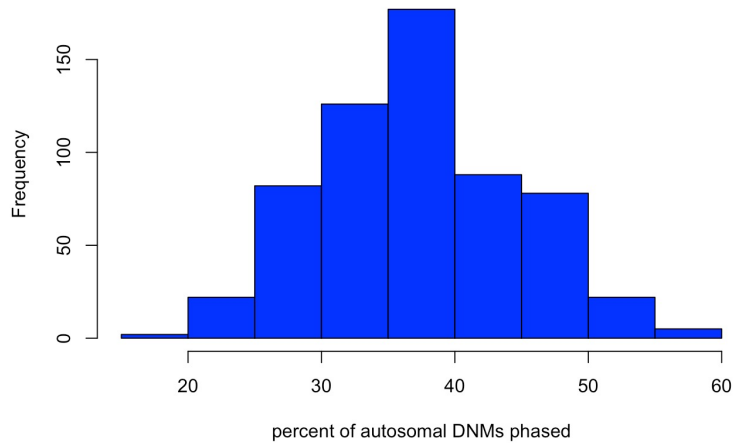
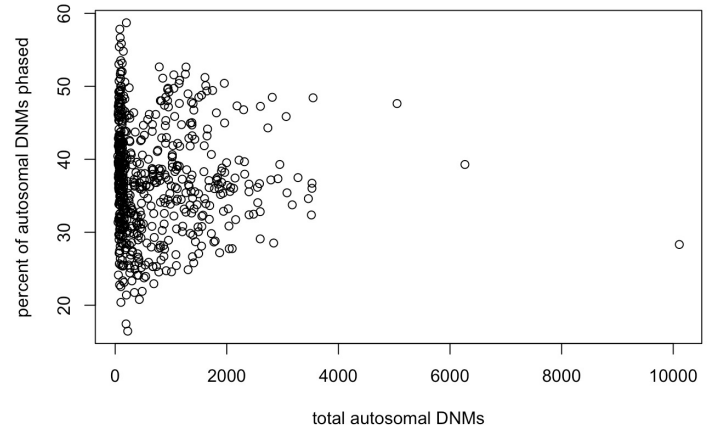


Figure S2

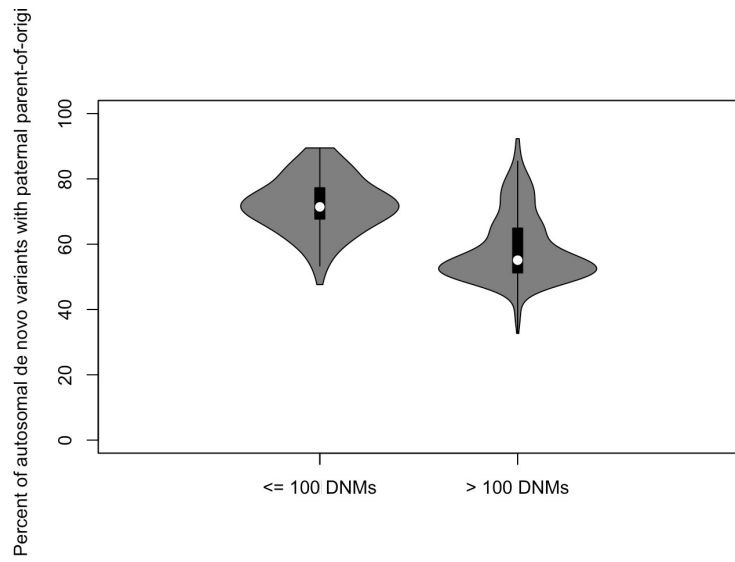
**A**



**B**



**C**



**D**

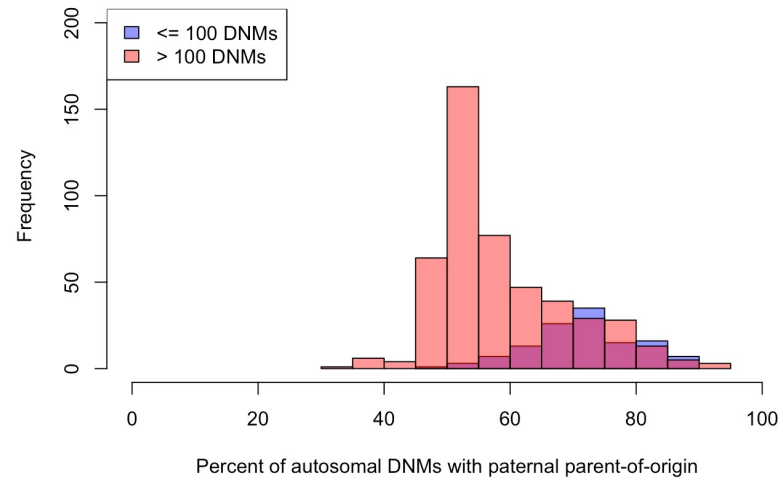


Figure S3

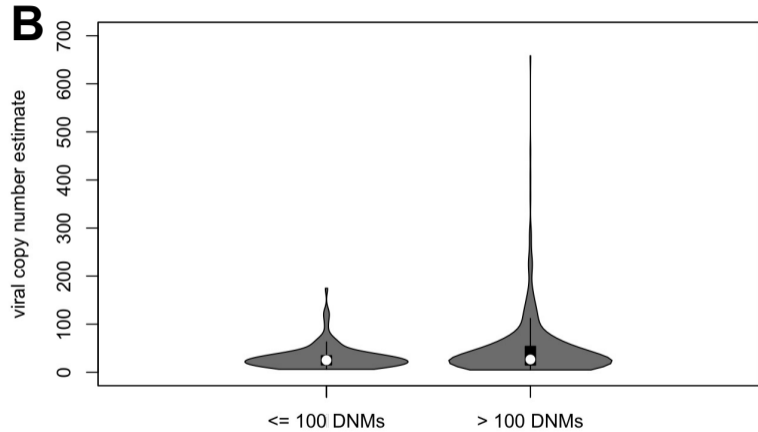
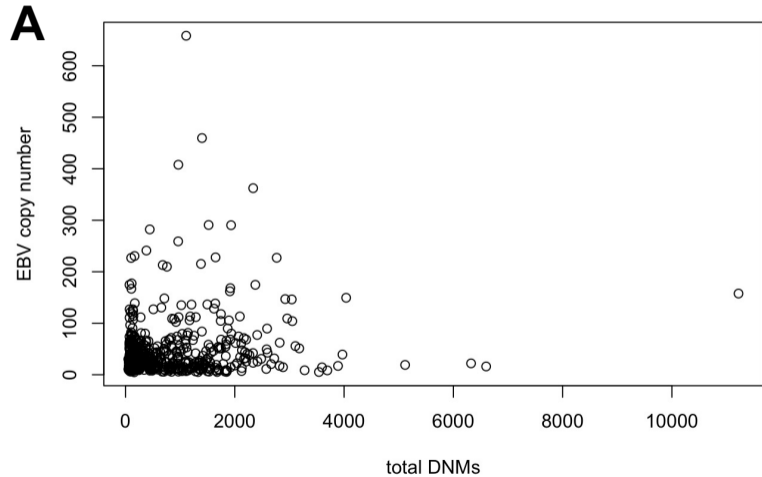
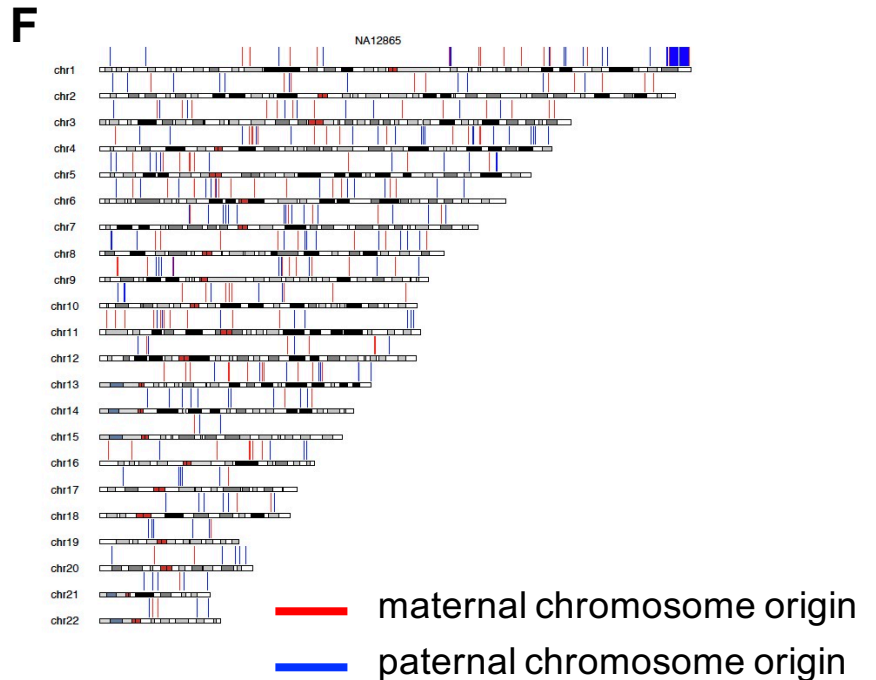
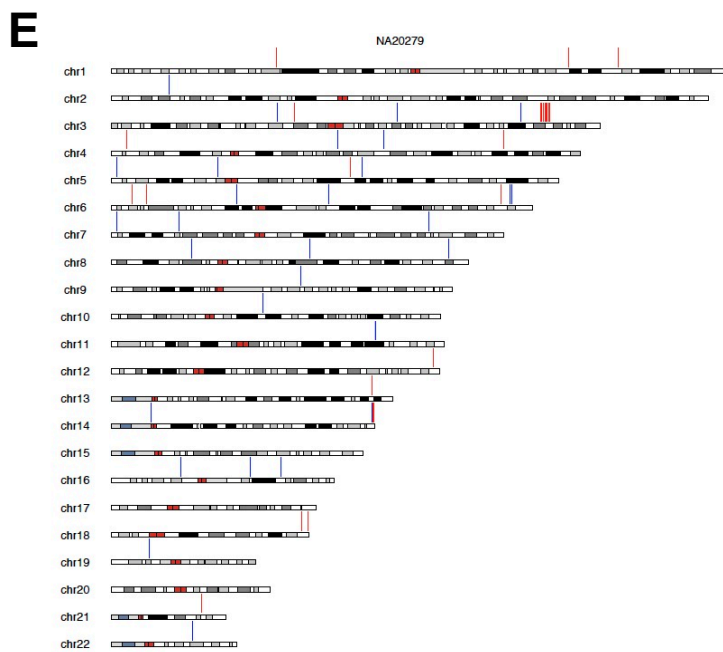
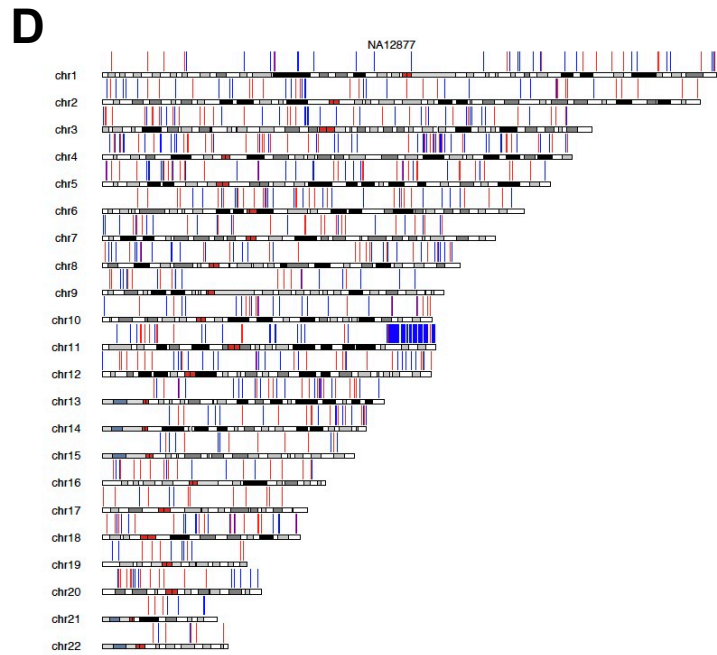
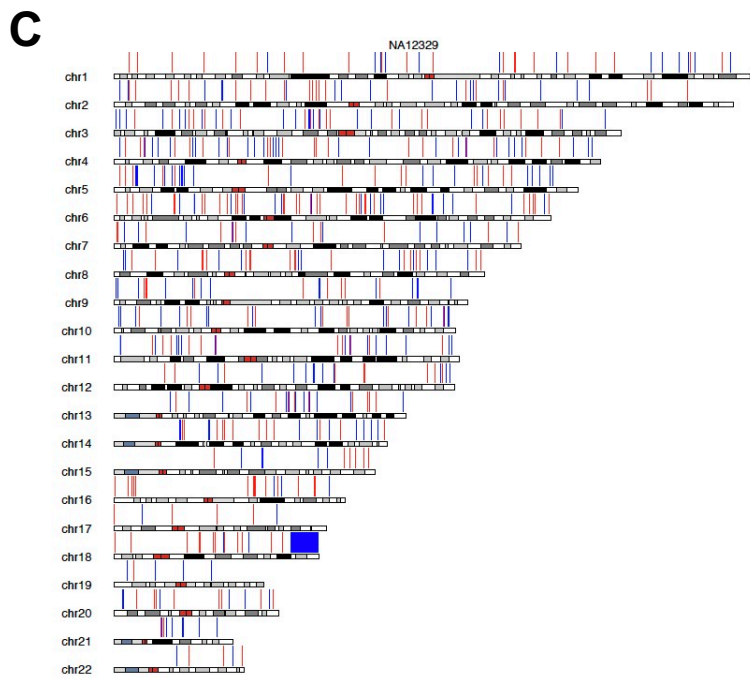
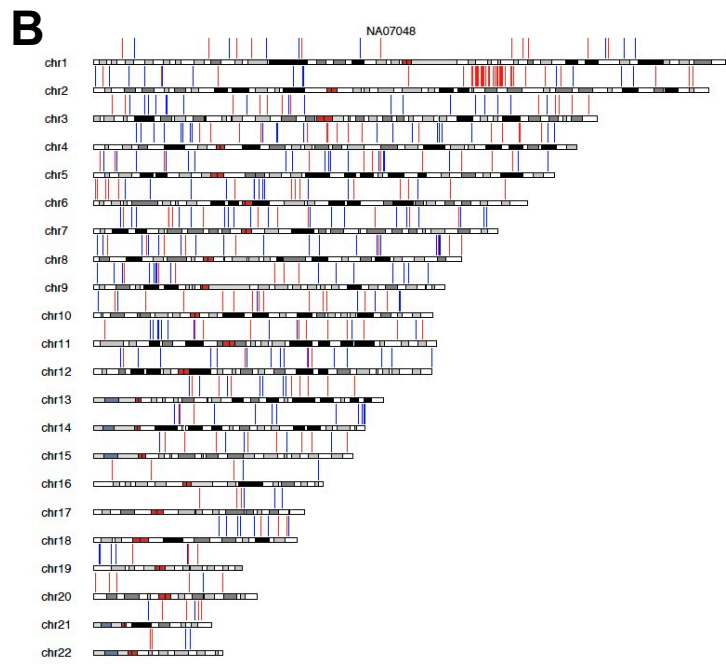
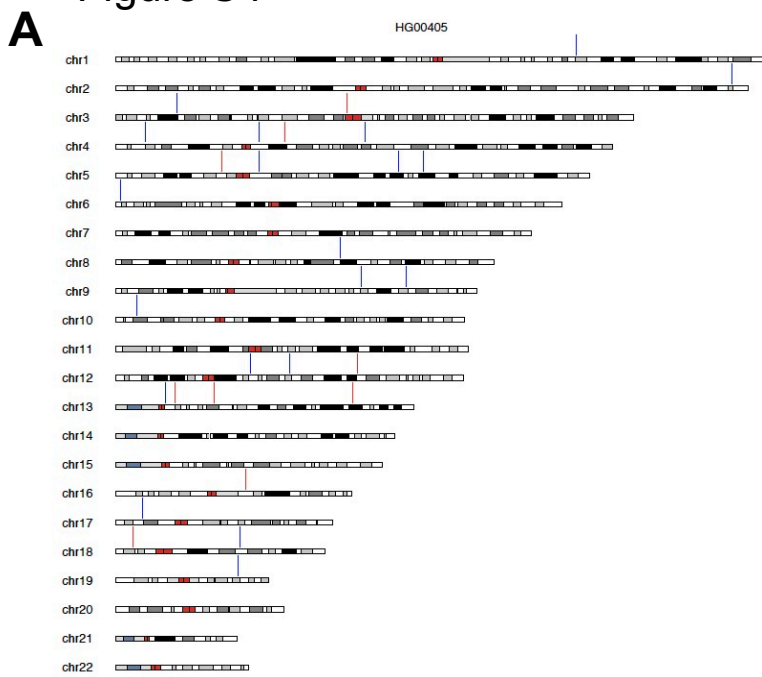
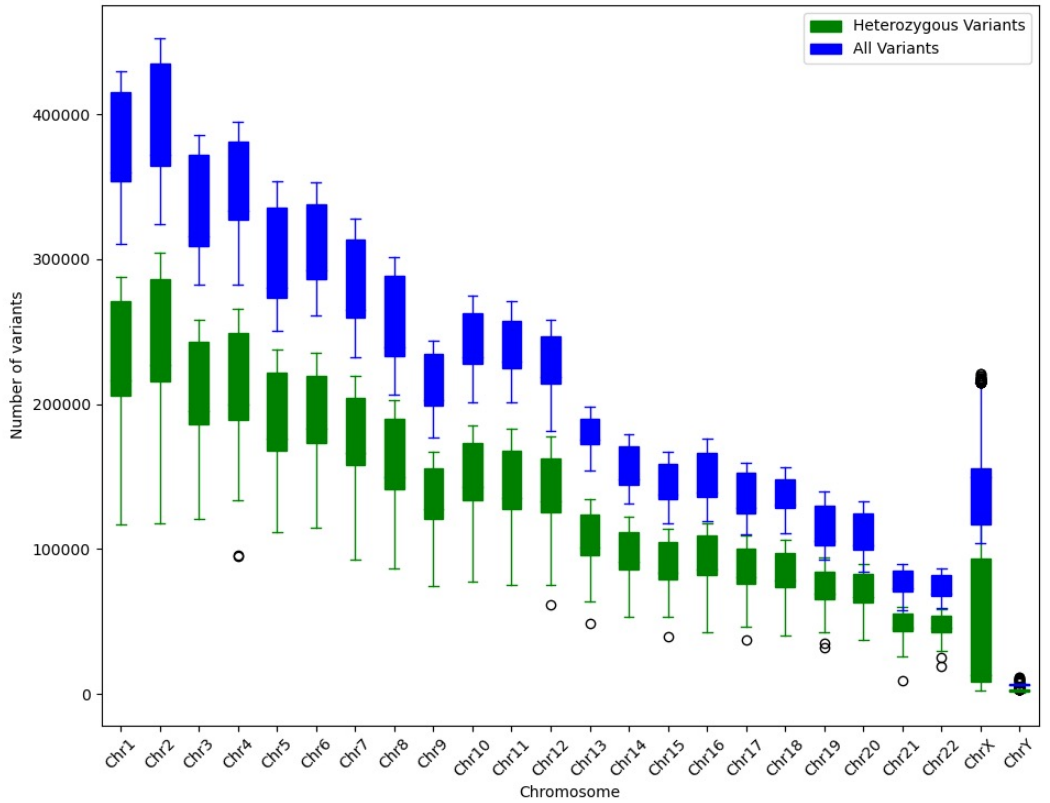


Figure S4



## DeepVariant

A



B

## GATK HaplotypeCaller

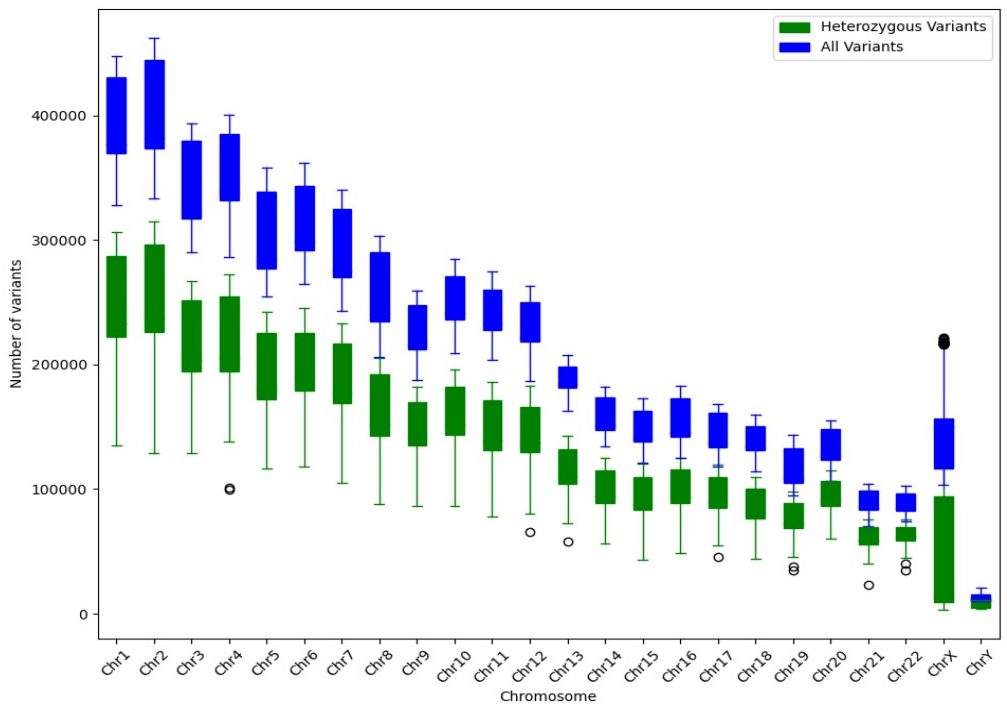
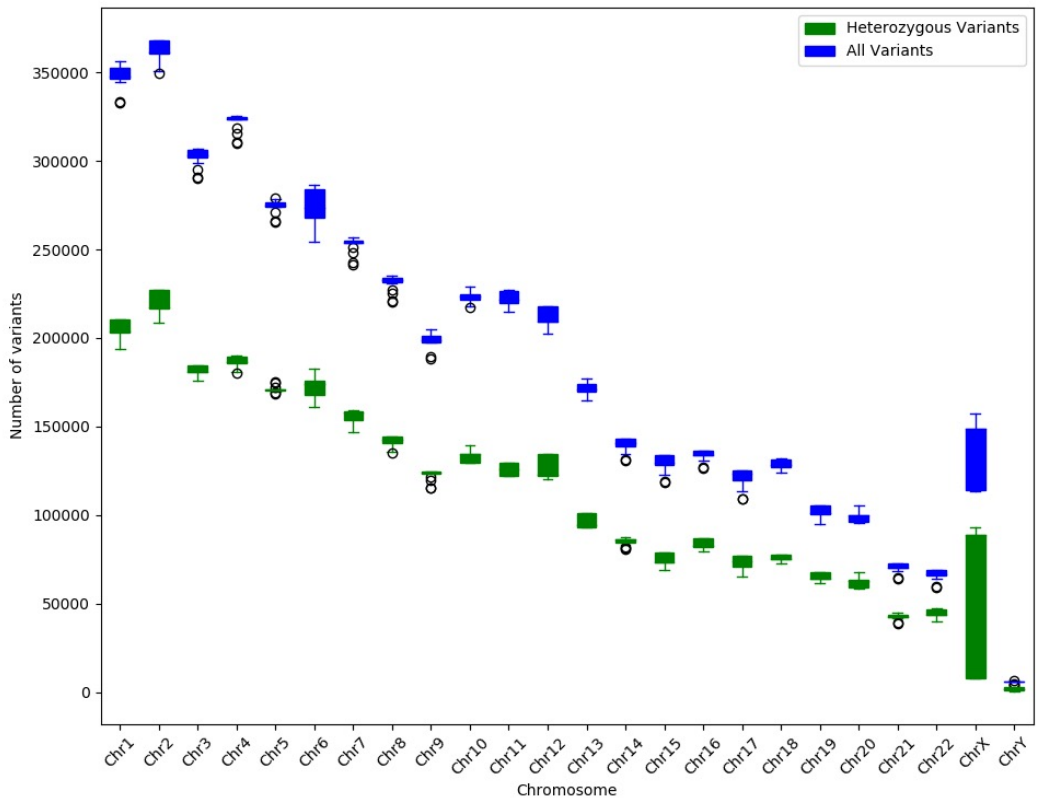


Figure S6

# DeepVariant

**A**



# GATK HaplotypeCaller

**B**

