1
2
3
4
5
6
7
8
9
10

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

Jeffrey K. Ng[1], Pankaj Vats[2], Elyn Fritz-Waters[3], Stephanie Sarkar[1], Eleanor I. Sams[1], Evin M. Padhi[1], Zachary L. Payne[1], Shawn Leonard[3], Marc A. West[2], Chandler Prince[3], Lee Trani[4], Marshall Jansen[3], George Vacek[2], Mehrzad Samadi[2], Timothy T. Harkins[2], Craig Pohl[3], Tychele N. Turner[1],*

1. Department of Genetics, Washington University School of Medicine, St. Louis, MO 63110, USA
2. NVIDIA Corporation, 2788 San Tomas Expressway Corporation, Santa Clara, CA 95051
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

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 rapidly call DNVs (HAT) and

4    demonstrated its effectiveness by applying it to 4,216 Simons Simplex Collection (SSC) whole-

5    genome sequenced parent-child trios from DNA derived from blood. In our SSC DNV data, we

6    identified $78 \pm 15$ DNVs per individual, $18\% \pm 5\%$ at CpG sites, $75\% \pm 9\%$ phased to the

7    paternal chromosome of origin, and an average allele balance of 0.49. These calculations are all

8    in line with DNV expectations. We sought to build a control DNV dataset by running HAT on

9    602 whole-genome sequenced parent-child trios from DNA derived from lymphoblastoid cell

10   lines (LCLs) from the publicly available 1000 Genomes Project (1000G). In our 1000G DNV

11   data, we identified $740 \pm 967$ DNVs per individual, $14\% \pm 4\%$ at CpG sites, $61\% \pm 11\%$ phased

12   to the paternal chromosome of origin, and an average allele balance of 0.41. Of the 602 trios,

13   $80\%$ had $> 100$ DNVs and we hypothesized the excess DNVs were cell line artifacts. Several

14   lines of evidence in our data suggest that this is true and that 1000G does not appear to be a static

15   reference. By mutation profile analysis, we tested whether these cell line artifacts were random

16   and found that 40% of individuals in 1000G did not have random DNV profiles; rather they had

17   DNV profiles matching B-cell lymphoma. Furthermore, we saw significant excess of protein-

18   coding DNVs in 1000G in the gene *IGLL5* that has already been implicated in this cancer. As a

19   result of cell line artifacts, 1000G has variants present in DNA repair genes and at Clinvar

20   pathogenic or likely-pathogenic sites. Our study elucidates important implications of the use of

21   sequencing data from LCLs for both reference building projects as well as disease-related

22   projects whereby these data are used in variant filtering steps.

23

1 **INTRODUCTION**

2       *de novo* variants (DNVs) are important for assessing mutation rates [1] and have been shown

3   to contribute to human disease (e.g., autism [2-10], epilepsy [11,12], intellectual disability [13-16],

4   congenital heart disorders [17-19]). Typically, the calling of DNVs from raw sequence data to final

5   calls can take days to weeks. Multiple DNV workflows exist that primarily rely on CPU-based

6   approaches [2-7,9,10,12-15,17,20-31]. These workflows employ different strategies including strict

7   filtering, utilizing multiple variant callers as opposed to using only one, machine-learning, and

8   incorporation of genotypic information at other sites around the genome. Overall, there is no

9   community consensus on a standard method for detecting DNVs. It is imperative that this process

10   be streamlined and flexible to enable broad adoption across the community. In this study, we

11   developed a rapid workflow to accelerate DNV calling using graphics processing units (GPUs)

12   that is integrated into NVIDIA Parabricks [32] software. We also developed an equivalent, freely

13   available open-source, CPU-based version of the workflow.  Together, the GPU-based workflow,

14   Hare, and the CPU-based workflow, Tortoise, make up HAT.

15

16       Our desire for a standardized, rapid DNV workflow stems from our interest in detecting

17   these DNVs in the large number of whole-genome sequencing (WGS) data in families with

18   neurodevelopmental disorders that has recently become available (https://anvilproject.org/data).

19   Studies assessing individuals with WGS data based on DNA derived from blood have provided

20   the field with our best estimates of DNV characteristics in humans [1]. One recent dataset, with DNA

21   derived from blood, consisting of 4,216 parent-child whole-genome sequenced trios from the

22   Simons Simplex Collection (SSC) has been extensively studied for DNVs [6,33-35]. We processed

23   this data with HAT and found that our method performed well.

1

2       This led us to assess the newly generated, publicly available, WGS dataset from a cohort

3       called the 1000 Genomes Project (1000G), where our initial goal was to build a control DNV

4       dataset.   Overall, 1000G is a data resource for the study of genetic variation that includes

5       individuals from diverse genetic ancestries [36,37]. Represented in the data are 602 trios from 18

6       worldwide populations (Figure S1). Moreover, as a field standard, 1000G has been utilized in

7       many applications as a control resource for filtering of genetic variation by allele frequency and/or

8       variant presence-absence in the dataset[38].

9

10       One complicating factor of DNV assessment in this resource is the fact that sequencing

11       data is generated from DNA isolated from lymphoblastoid cell lines (LCLs) [37] as opposed to

12       primary tissue. Epstein-Barr Virus is used to make these LCLs and passaging over time enables

13       the accumulation of cell line artifacts. These artifacts can complicate variant filtration schemes and

14       the utility of this data as a frequency control. As opposed to a random accumulation of mutations

15       in each individual, we found that 80% of 1000G individuals had an excess of DNVs and 40% of

16       all 1000G individuals had a profile matching a B-cell lymphoma. The similarity to this cancer is

17       problematic, and it would be imperative that this data not be used as a control in the context of the

18       study of these and related cancers. A secondary consequence of the excess DNVs is their presence

19       at disease-related sites whereby simple filtering schemes may accidentally remove sites of interest

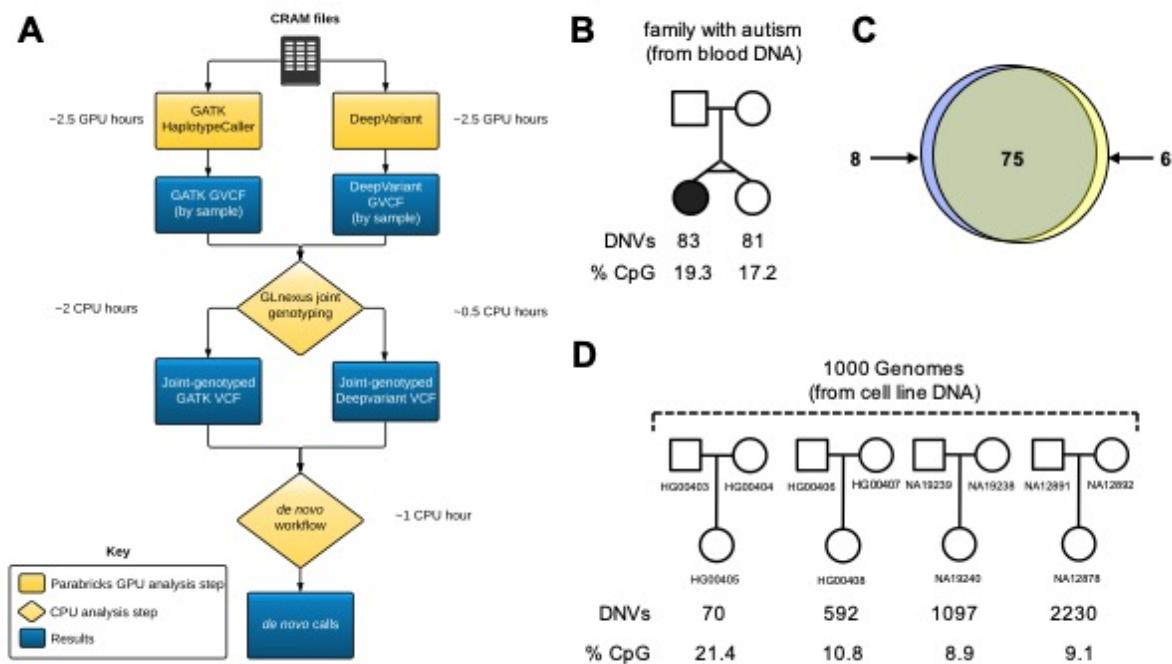20       in patients due to their presence in 1000G.

21

22

1 **RESULTS**

2 *Rapid DNV calling with GPUs*

3        HAT consists of three main steps: GVCF generation, family-level genotyping, and filtering

4    of variants to get final DNVs. We utilized existing features of the NVIDIA Parabricks software

5    for rapid GVCF generation from GPU accelerated versions of  GATK [39] HaplotypeCaller and

6    Google DeepVariant [40] . The run times for GVCF generation are ~40 minutes per sample on a 4

7    GPU node and can be run in parallel on all three family members in the parent-child trio. Post-

8    GVCF generation, the trio is genotyped using the GLnexus joint genotyper [41]. Finally, our post-

9    genotyping custom DNV filtering workflow runs in ~1 hour with speedups at all steps with

10   parallelization providing a clear advantage over CPU-based approaches (Figure 1A).

11

## Figure 1



12

1    **Figure 1:** *de novo* variant calling in short-read whole-genome sequencing data. A) *de novo*

2    workflow for detection of *de novo* variants (DNVs) from aligned read files (crams); B and C)

3    Benchmarking DNV workflow in a monozygotic twin pair sequenced from DNA derived from

4    blood; D) DNV detecting in four trios in the 1000 Genomes Project.

5

6         To benchmark HAT, we tested it on a monozygotic twin pair with WGS data derived from

7    blood DNA. These individuals should share the same DNVs from generation in the germline.

8    However, they may differ at some sites if DNVs occur in a post-zygotic, somatic manner. The

9    twins shared 75 autosomal DNVs and contained 83 and 81 autosomal DNVs, respectively (Figures

10   1B and 1C). The percent CpG was 19.3% and 17.2%, respectively and in line with previous

11   published estimates of ~20% [1,6] (Figures 1B and 1C). As this monozygotic twin pair was discordant

12   for the phenotype of autism, we also tested whether there were any protein-coding DNV

13   differences between the two twins. These would potentially be relevant for autism, but there were

14   no such differences.

15

16        To establish a DNV callset from the 1000G data as a control, we started with the assessment

17   of DNVs with HAT in four trios from the 1000G (Figure 1D). Two were chosen at random (i.e.,

18   HG00405, HG00408) and two were chosen because they were "famous" trios assessed in many

19   other studies (i.e., NA12878 [23,42], NA19240 [23]). One of these trios (HG00405) had 70 DNVs and

20   a CpG percent of 21.4 as we would have expected from DNA derived from blood. To our surprise,

21   the other trios had varying numbers of DNVs from 592 to 2,230 with NA12878 (arguably the most

22   studied individual in 1000G) having the most DNVs. With the increase in DNVs the CpG percent

23   dropped considerably down to ~10%. We also assessed 3,598 of the DNVs from the four trios by
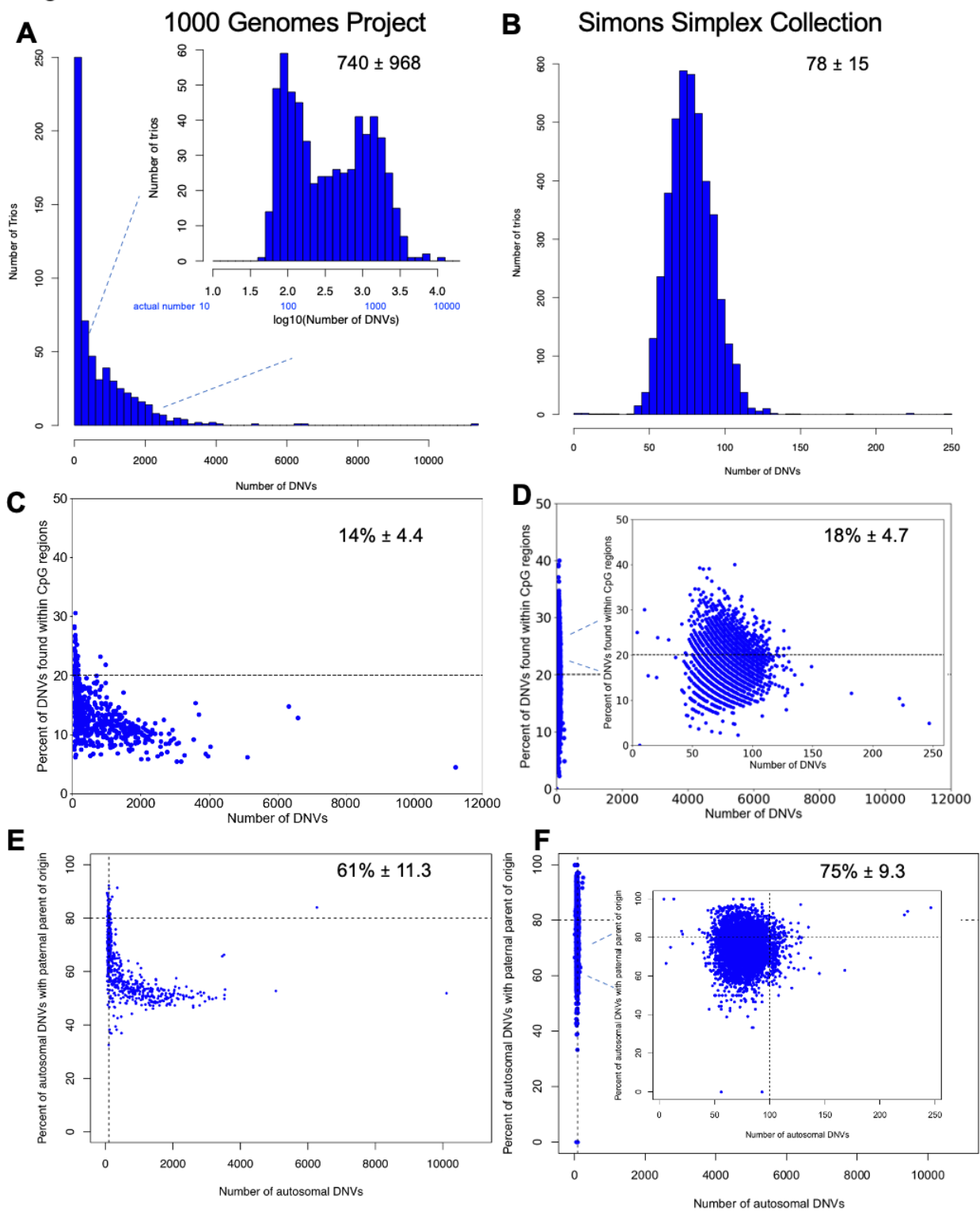
1    manual visual inspection of the underlying reads in each family member (Table S1) and found that

2    93.6% of the variants appeared to be true DNVs, 4.9% were inherited, and 1.5% were low

3    confidence calls.

4

5    *Differences in DNVs in blood and LCLs*

6         Our initial observations led us to focus on two main cohorts: the 602 trios from 1000G

7    (Table S2) with DNA derived from LCLs and 4,216 trios from the Simons Simplex Collection

8    (SSC) with DNA derived from blood. In the 1000G data we detected 445,711 total DNVs in the

9    cohort (Table S3). There were $740 \pm 968$ DNVs per individual (Table S4) with a clear bimodal

10   distribution (Hartigan's dip test: D = 0.033, p-value = $1.32 \times 10^{-4}$) wherein some individuals

11   contained an excess of DNVs (Figure 2A). In the SSC, we identified 329,589 total DNVs in the

12   cohort. There were $78 \pm 15$ DNVs per individual (Figure 2B, Table S5).  The values derived from

13   the SSC data are in line with expectation and highlight the effectiveness of our DNV workflow.

14   However, the values in the 1000G are higher than expected and we estimated the number of

15   individuals with appropriate numbers of DNVs by splitting the 1000G data into two groups:

16   individuals having less than or equal to 100 DNVs (n = 123) and individuals with greater than 100

17   DNVs (n = 479). This estimate suggests that only 20.4% of trios in the 1000G have the correct

18   number of DNVs and we thought those with excess DNVs may have cell line artifacts due to

19   culturing of LCLs.

7

Figure 2

1    **Figure 2:** Comparison of characteristics of DNVs detected in1000 Genomes Project (1000G)

2    and Simons Simplex Collection (SSC) callsets. A)    Histogram of DNV counts from 1000G in

3    602 trios; B) Histogram of DNV counts from SSC in 4216 trios; C) ; C) Percent of DNVs found

4    within CpG sites versus the total number of DNVs for 1000G; D) Percent of DNVs found within

5    CpG sites versus the total number of DNVs found for SSC; E) Percent of autosomal DNVs with

6    paternal parent of origin  versus the total number of DNVs for 1000G; F) Percent of autosomal

7    DNVs with paternal parent of origin  versus the total number of DNVs for SSC.

8

9         We assessed two main features of typical DNVs to investigate our hypothesis that the

10    excess DNVs found in individuals were cell line artifacts. These features are DNVs at CpG

11    locations and the percent of DNVs arising on the paternal chromosome. As mentioned previously,

12    the percent of DNVs at CpG should be ~20% and the percent of DNVs arising on the paternal

13    chromosome should be ~80% [43]. We saw that in the 1000G trios $14 \pm 4.4\%$ of DNVs per individual

14    occurred at CpG sites (Figure 2C) with individuals with less than or equal to 100 DNVs having17.4

15    $\pm 5.2\%$ DNV at CpG and in families with greater than 100 DNVs $12.7 \pm 3.6\%$ DNV at CpG. The

16    difference in DNVs at CpG sites between these two groups was significant (Wilcoxon rank sum:

17    p-value $< 2.2 \times 10^{-16}$). In the SSC, the percent of DNVs at CpG was $18 \pm 4.7\%$ and in line with

18    expectation (Figure 2D). In the 1000G, the percent of DNVs that were phase-able for parent-of-

19    origin was $37.2 \pm 7.5\%$ (Figure S2). Of the phased variants, $61 \pm 11.3\%$ were on the chromosome

20    of paternal origin (Figure 2E, Table S6). In the families with less than or equal to 100 DNVs this

21    rose to $72.0 \pm 8.5\%$ and in the families with greater than 100 DNVs it fell to $58.6 \pm 10.3\%$. This

22    difference in percent phased variants of paternal origin was found to be significantly different

23    (Wilcoxon rank sum: p-value $< 2.2 \times 10^{-16}$). The drop leveled off to ~50% in the individuals with

1    the most DNVs (Figure 2D). In the SSC, we were able to phase 37% of DNVs (Figure S3) with

2    the percent of DNVs phased to paternal chromosome of origin was 75% ± 9.24 and that was also

3    in line with expectation (Figure 2F).

4

5    We also tested whether there was a difference in the allele balance (AB) in the child at

6    DNV sites in the 1000G and the SSC (Figure S4). We found that the 1000G had a mean AB of

7    0.42 and in the SSC it was nearly perfect at an AB of 0.49 in line with expectation of 0.5. This

8    indicated a lower average AB level in 1000G from newly arising mutations from cell line artifacts.

9

10    Overall, these comparisons showed that the individuals in the 1000G with less than or equal

11    to 100 DNVs behaved more like true DNVs in regard to CpG percentage, percent arising on the

12    paternal chromosome, and allele balance. This also was true for the SSC trios where DNA was

13    derived from blood. However, individuals in the 1000G with > 100 DNVs did not have statistical

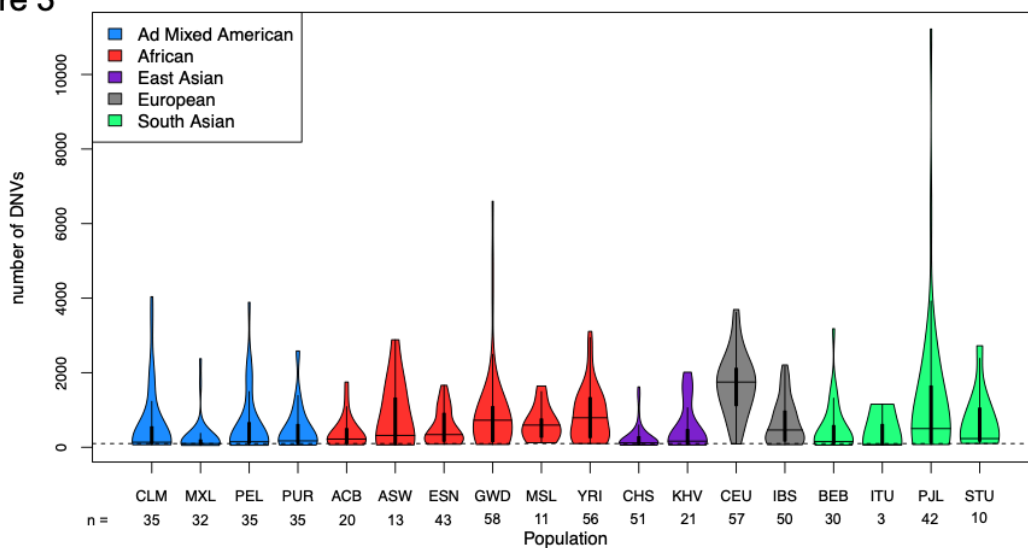14    properties of true DNVs providing evidence they may be cell line artifacts.

15

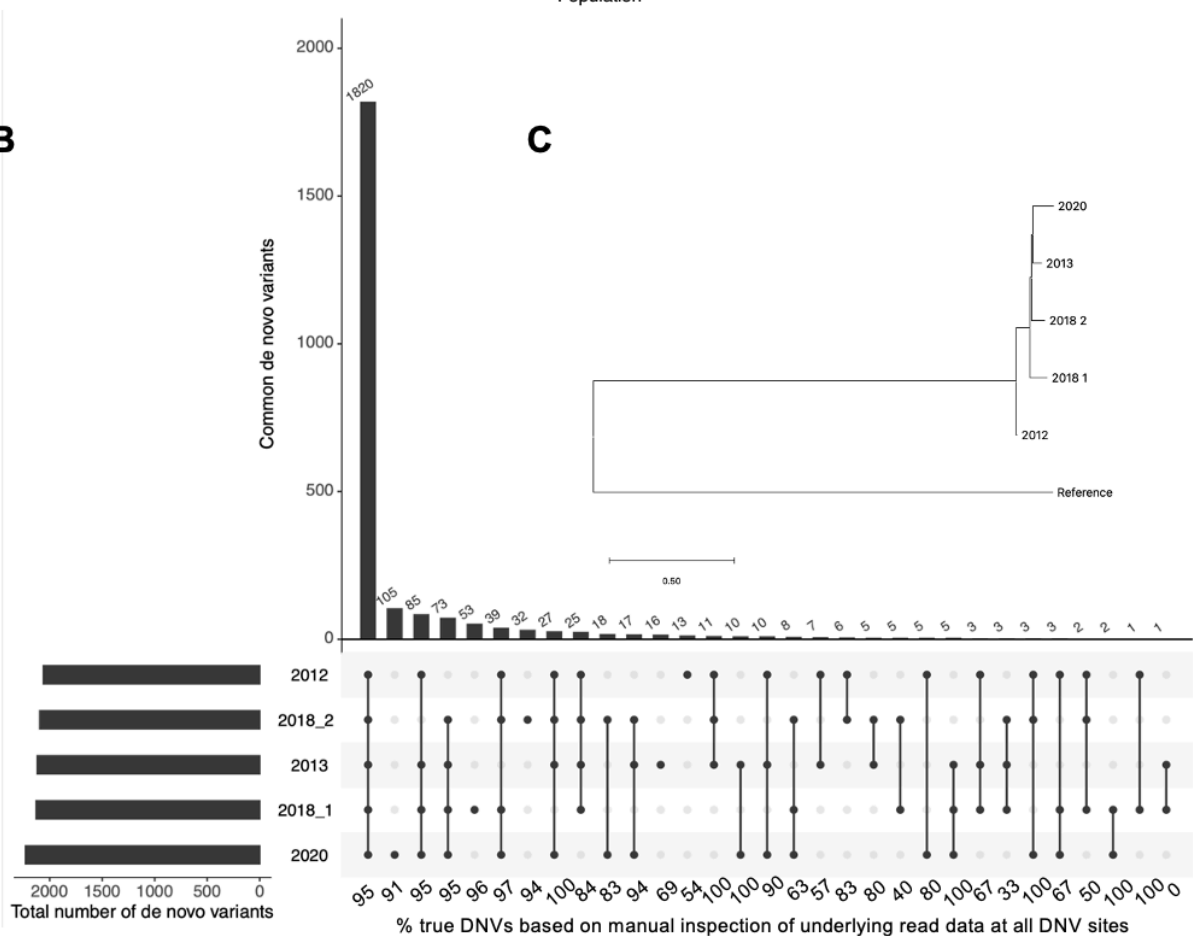16    *DNVs by 1000G population*

17    While we expected there to be no difference in DNV counts per individual by ancestry we

18    sought to see if there were any populations with excess DNVs (Figure 3A). The population with

19    the most DNVs was the CEU having on average 1,688 DNVs per individual. We hypothesized

20    that this may be because the CEU is one of the oldest cohorts in the 1000 Genome Project dating

21    back to the HapMap project [44] and these individuals may have cell lines that have been cultured

22    more over time than other populations.

**Figure 3:** Assessment of five replicates of NA12878. A) Population distribution of 1000G

dataset. B) UpSet plot demonstrating the number of variants detected in the replicates (at the

11

1    bottom of the plot the percent of true DNVs is listed for each category); C) Phylogenetic tree of

2    the five replicates.

3

4    *DNVs increase over time*

5        We utilized the fact that the 1000G individual NA12878 has been studied and sequenced

6    multiple times over the past ten years by WGS [37] (SRA identifiers: SRR944138 and SRR952827).

7    Presumably, across time, the utilization of NA12878 has required additional culturing of this cell

8    line, and potentially even by different laboratories. We aggregated five Illumina WGS datasets

9    from this individual, downsampled them to ~30x coverage, and assessed them with HAT. The data

10   for this individual ranged from the year 2012 to the year 2020 and we found that the 2012

11   experiment had the least DNVs (n = 2,060) and the 2020 experiment had the most DNVs (n =

12   2,230) (Figure 3B). Overall, the five replicates had a large overlap of DNVs (n = 1,820) across all

13   samples. These shared DNVs constitute what were present in the ancestor of all the cell line

14   replicates. DNVs not shared by all five replicates are sometimes shared by a subset of the replicates

15   and are sometimes unique to the replicate. To formally assess the ancestral state, we built a

16   phylogenetic tree based only on the DNVs and saw that the farthest replicates from each other in

17   the tree were the 2012 and 2020 replicates (Figure 3C). To further assess the DNVs in NA12878,

18   we randomly sampled 25 DNVs from the union dataset from the five replicates. We performed

19   Sanger sequencing on DNA from NA12878 and her parents (NA12891, NA12892) (Figure S5-

20   S29, Table S7). We found that 24 of the 25 DNV sites gave clear results in the Sanger sequencing

21   with 23 confirming as real DNVs. Surprisingly, we found two sites, chr12_91353615_T_C and

22   chr13_81142986_T_A, that were determined to be true *de novo* variants but were previously

23   shown to be a false positive reading using Sanger sequencing [23]. In the Sanger experiments, one

1    site chr11_134531608_C_G showed subtle evidence for the variant allele in NA12878, so we

2    pursued deep amplicon sequencing of this region in the trio using Oxford Nanopore Technologies

3    (ONT) sequencing (Figure S30).  This resulted in a variant allele frequency of 11% in NA12878

4    suggesting this is a cell line artifact. This was elevated in comparison to the background rate of

5    1% in NA12891 and 0% in NA12892, that is in line with expectation for error rates from ONT.

6    Intriguingly, this DNV was only found in one of the NA12878 replicates (2018_1). Overall, this

7    indicates that 96% of DNVs called with HAT are real (24/25) and this estimate is close to the

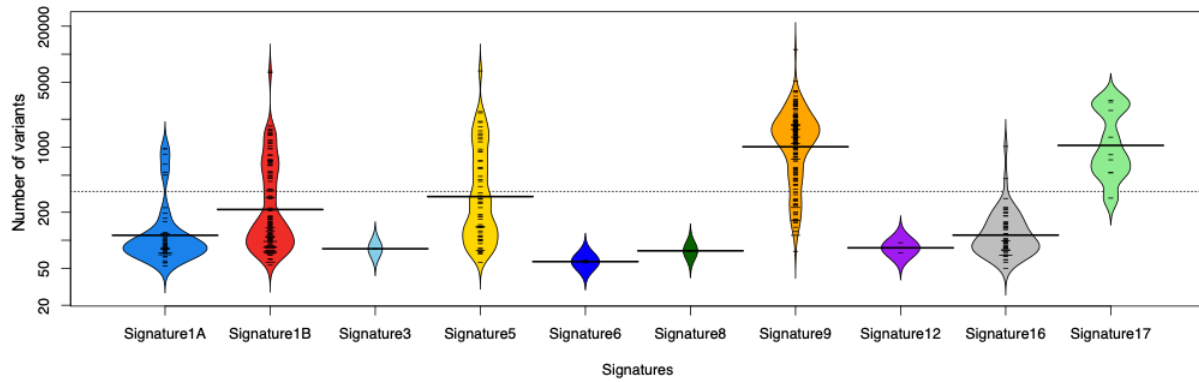8    93.6% we saw by manual inspection of underlying read data at 3,598 DNV sites (see above).

9

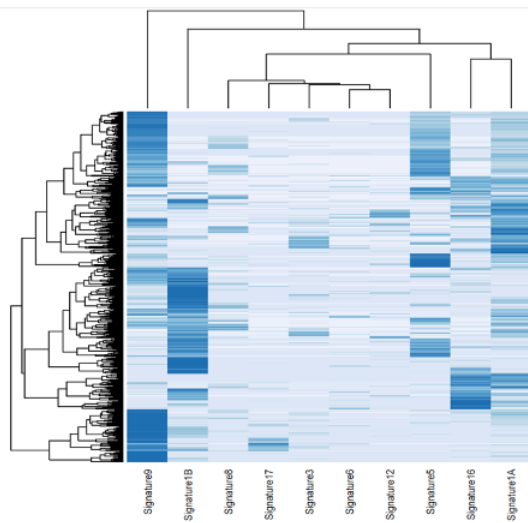10    *Genomes with cancer mutation profiles*

11          We used mutation profile analysis [45] (Table S8) to determine whether the DNVs identified

12    in individuals from the 1000G had any certain characteristics. For this analysis, we utilized a

13    method that would enable comparisons to known mutational profiles that are either age-related

14    (reminiscent of true DNVs) or are seen in cancers (Figure 4A and Figure 4B). There were 186

15    individuals (30%) that had a strong contribution of an age-related signature (Signature 1A,

16    Signature 1B). To our surprise, the other contributing signatures in individuals were primarily

17    those associated with B-cell lymphomas (Signature 5, Signature 9 and Signature 17) in 241

18    individuals (40%). This was intriguing because lymphoblastoid cell lines are generated from B-

19    cells that are infected with Epstein Barr Virus and demonstrates that new mutations are not arising

20    in a random manner. Rather they are being generated in a manner consistent with the development

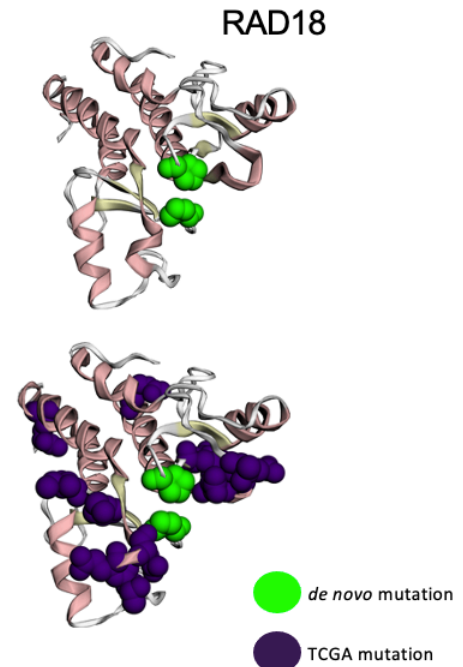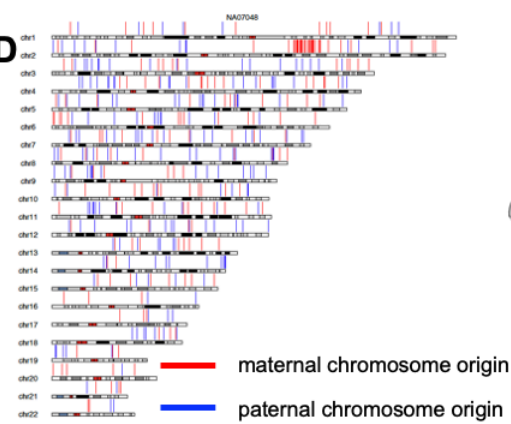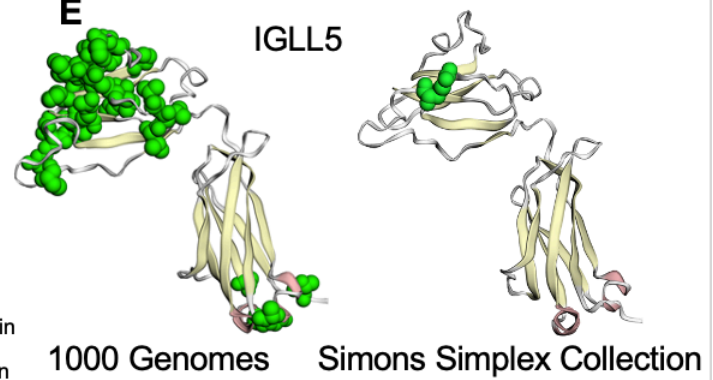21    of cancer in the same cell type.

## Figure 4

1    **Figure 4:** Mutational properties of DNVs. A) Mutation signature analysis showing the total

2    number of DNVs and the individuals with each signature type; B) Heatmap of individuals based

3    on their mutational signatures; C) Mutations in the DNA repair gene RAD18 shown on their 3D

4    structure (and modeled using mupit). Also, shown are known cancer mutations from The Cancer

5    Genome Atlas (TCGA); D) Location of DNVs based on their phased parent-of-origin in

6    NA07048. Most notable there are a cluster of mutations on the maternal chromosome on

7    chromosome 2; E) DNVs in IGLL5 shown on their 3D structure (and modeled using mupit).

8    The image on the left is modelling variants discovered in 1000G, the image on the right is

9    modelling variants discovered in SSC.

10

11    We further sought to determine what the mechanism was for the generation of a B-cell

12    lymphoma-like state. First, we determined whether there was high rate of aneuploidies in the cell

13    lines. By digital karyotyping (Table S9) we found that 595 individuals (98.8%) had a typical

14    chromosome complement (46,XX or 46,XY), four were missing a sex chromosome (45,X0), one

15    was 47,XXY, one had three chromosome 12 (47,XY), and one had three chromosome 9 (47,XY).

16    This demonstrated that while these aneuploidies are occurring in some cell lines, they are probably

17    not the main driving factor. Next, we looked at DNVs in genes involved in DNA repair and found

18    17 individuals contained a missense or loss-of-function in one of these genes (Table S10).

19    Individuals with B-cell lymphoma profiles and disruptive mutations in DNA repair genes included

20    mutations in the following genes *FANCF* (HG01126), *MUS81* (NA10838), *POLB* (NA10838),

21    *POLD1* (NA19677), *POLE* (HG01096), *RAD18* (NA12864) (Figure 4C), *RAD51* (HG02683),

22    *RPA4* (HG02630), and two individuals with mutations in *FANCA* (HG02841, HG03200) and *WRN*

23    (HG04115, NA19161), respectively (Table 1). Third, we looked at Epstein Barr Virus load in each

15

1    of the genomes (Table S11) and found that there was a weak, yet significant, correlation with the

2    number of DNVs (p = 2.32 × 10$^{-5}$, r = 0.17) (Figure S31). By visual inspection of phased variation

3    in all individuals we also identified individuals with clusters of mutations (e.g., NA07048, Figure

4    4D, Figure S32).

5

6    *Excess of DNVs in IGLL5*

7         We applied a multi-phase approach to determine if there were any genes with enrichment

8    of protein-coding DNVs in individuals with greater than 100 DNVs. In the first phase, we tested

9    whether there was genome-wide significance for enrichment of protein-coding DNVs (missense,

10   loss-of-function) in any specific genes. By application of two methods (chimpanzee-human,

11   denovolyzeR), we identified 29 significant genes (*ARMC3*, *BCL2*, *BCR*, *C6orf15*, *CCDC168*,

12   *CSMD3*, *EGR3*, *EXO1*, *HLA-B*, *HLA-C*, *IGLL5*, *KMT2D*, *LINGO2*, *LTB*, *MEOX2*, *MUC16*,

13   *MUC22*, *NPAP1*, *PCLO*, *PRPF40A*, *RUNX1T1*, *SGK1*, *STRAP*, *TMEM232*, *TNXB*, *TTN*, *WDFY4*,

14   *XIRP2*, *ZNF488*) with excess of DNVs (Table S11). In the second phase, we tested these 29 genes

15   to see whether there were significantly more protein-coding DNVs in individuals with greater than

16   100 DNVs in comparison to individuals with less than or equal to 100 DNVs. Only *IGLL5* was

17   significant in this comparison (1.79 × 10$^{-3}$) (Table S12, Table S13, Figure 4E). To test whether

18   this finding was relevant only to LCLs, we looked for protein-coding DNVs in SSC and only found

19   one missense variant (Figure 4E). This gene did not have significant excess of DNVs in SSC.

20

21   *DNVs identified in clinically relevant variants*

22         We tested whether any of the DNVs detected were already known to be pathogenic or

23   likely-pathogenic in the Clinvar [46] database (Table 1). There were 15 mutations meeting these

1 criteria (Table S14). We rescored these variants using Franklin software to assess their

2 pathogenicity and found that 13 were also pathogenic or likely-pathogenic by this approach.

3 Twelve of these variants were associated with described phenotypes in Clinvar. These included a

4 missense variant in *SOS1* involved in Noonan syndrome, a missense variant in *SCN2A* involved in

5 seizures, a stop gained variant in *UNC80* involved in a syndrome with hypotonia, intellectual

6 disability, and characteristic facies, a missense variant in *THRB* involved in thyroid hormone

7 resistance, a missense variant in *PKHD1* involved in polycystic kidney disease, a stop-gained in

8 *ERCC6* involved in Cockayne syndrome, a stop-gained in *ANO5* involved in gnathodiaphyseal

9 dysplasia, a stop-gained in *PHF21A* involved in inborn genetic disease, a missense in MYO7A in

10 Usher syndrome type 1, a stop-gained in *ROBO3* in Gaze palsy with progressive scoliosis, a

11 missense in *COL4A1* involved in inborn genetic disease, and a missense in *POLG* involved in

12 POLG-related disorder.

13

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

| Category | individual | *de novo* variant | variant type | gene |
|---|---|---|---|---|
| | HG01074 | chr3_48447050_C_G | missense | *ATRIP* |
| | HG02841 | chr16_89799603_A_G | splice_donor | *FANCA* |
| | HG03200 | chr16_89762010_C_T | missense | *FANCA* |
| | HG01126 | chr11_22625482_T_G | missense | *FANCF* |
| | NA18875 | chr5_80654794_G_A | missense | *MSH3* |
| | HG02650 | chr6_31759121_C_T | missense | *MSH5* |
| | NA10838 | chr11_65865247_C_T | missense | *MUS81* |
| DNA damage repair gene | NA10838 | chr8_42357362_AT_A | frameshift | *POLB* |
| | NA19677 | chr19_50407375_G_A | missense | *POLD1* |
| | HG01096 | chr12_132634327_C_T | missense | *POLE* |
| | HG01755 | chr15_89321792_C_T | missense | *POLG* |
| | NA12864 | chr3_8958938_G_T | missense | *RAD18* |
| | HG02683 | chr15_40729853_T_C | missense | *RAD51* |
| | HG02630 | chrX_96884884_G_A | missense | *RPA4* |
| | NA19919 | chr3_133644039_A_G | missense | *TOPBP1* |
| | HG04115 | chr8_31120294_C_T | missense | *WRN* |

| | NA19161 | chr8_31124967_G_T | missense | *WRN* |
|---|---|---|---|---|
| | HG03795 | chr11_22274728_C_T | stop_gained | *ANO5* |
| | NA10854 | chr13_110179298_C_T | missense | *COL4A1* |
| | NA10842 | chr10_49530737_G_A | stop_gained | *ERCC6* |
| | HG02668 | chr1_111787063_C_T | missense | *KCND3* |
| | HG02466 | chr1_39485559_G_A | missense | *MACF1* |
| | HG02129 | chr11_77206108_G_A | missense | *MYO7A* |
| Clinvar pathogenic / likely pathogenic | HG03122 | chr11_45949458_G_A | stop_gained | *PHF21A* |
| | NA12707 | chr6_52058438_C_T | missense | *PKHD1* |
| | HG01755 | chr15_89321792_C_T | missense | *POLG* |
| | HG02892 | chr11_124875581_C_T | stop_gained | *ROBO3* |
| | HG03635 | chr2_165310406_G_A | missense | *SCN2A* |
| | NA10830 | chr2_39023106_C_T | missense | *SOS1* |
| | NA10831 | chr3_24143512_G_A | missense | *THRB* |
| | HG01629 | chr2_209775898_C_T | stop_gained | *UNC80* |
| | HG00558 | chr16_88435401_G_A | missense | *ZNF469* |

**DISCUSSION**

While the 1000G data has been extensively studied in the past, there has been no previous cross-cohort assessment of DNVs. This limitation is primarily because family-based sequencing was not available until 2020 when this cohort was sequenced by high-coverage short-read WGS ten years after the initial ground-breaking publication on the 1000G [47]. Determining DNV profiles across this dataset of diverse individuals is critical for assessment of mutation rates in the human population, while also providing a more complete catalog of all genetic variants within these individuals. The decision to sequence these individuals using DNA derived from lymphoblastoid cell lines was a practical one. However, it opened the door to the possibility of cell line artifacts, while simultaneously introducing a dynamic aspect to this extensive set of controls. As control samples, the cell lines that were used as the inputs for the 1000G are still actively used across laboratories, acting as matched controls for workflows to known sets of variants. The large distribution of DNVs across the 1000G suggest that a subset of the control source inputs are dynamic, and in some cases, harbor a spectrum of genetic variants associated with B-cell

18

1     lymphomas or named clinical syndromes. Laboratories using control samples from the 1000G

2     should account for both the presence and dynamic nature of the reported DNVs and in some cases

3     may consider changing which control samples to use within the laboratory to avoid any of the

4     associated issues with the presence of DNVs. Additionally, other public efforts to establish

5     reference data sets using cell lines should consider the impacts of DNVs on their project design.

6

7          We utilized a novel and accelerated analysis workflow to detect DNVs from short-read,

8     whole-genome sequencing data. We showed this new workflow is of high-quality by running it on

9     4,216 trios with WGS, from the SSC, on DNA derived from blood. This analysis revealed expected

10    number of DNVs, percent of DNVs at CpG sites, percent of DNVs phased to the paternal

11    chromosome of origin, and average allele balance of the DNV. This was an important analysis and

12    was in contrast to our DNV analysis of the 1000G. In total, we identified 445,711 DNVs in the

13    602 children from 1000G assessed in this study. We provide a cross-cohort joint-genotyped VCF,

14    family-level VCFs, DNV calls, and phased DNV results for the 602 trios in this study as a public

15    community resource (Globus endpoint: "Turner Lab at WashU - DNV in 1000 Genomes Paper",

16    direct      link:      https://app.globus.org/file-manager?origin_id=3eff453a-88f4-11eb-954f-

17    752ba7b88ebe&origin_path=%2F). Originally, it was assumed that the DNVs across the 1000G

18    would have been random and minimal, and yet only 20% of the offspring (123 children) have a

19    number of DNVs around expectation (< 100) and the remainder have an excess of DNVs with the

20    most extreme case being an individual (HG02683) having 11,219 DNVs. We hypothesized that

21    the excess DNVs were cell line artifacts and found multiple lines of evidence to support this

22    hypothesis, including a reduction in the percent of DNVs at CpG as well as the reduction in percent

23    phased to the paternal parent-of-origin chromosome with increasing DNVs, respectively. A

19

1  detailed analysis of individual NA12878, who has been studied various times over the years,

2  revealed increasing DNVs in the more recently sequenced samples also supporting this hypothesis.

3  The changes in the DNVs for NA12878 suggest the dynamic nature of the DNVs, demonstrating

4  that the number is increasing over time.

5

6        When mutational signature analysis was performed on this new set of DNVs, the most

7  common mutation signatures were those seen in B-cell lymphomas. This signature was found in

8  40% of individuals in the 1000G. This is important as the lymphoblastoid cell lines are generated

9  from B-cells and points to a non-random accumulation of mutations that are in line with the

10  development of cancer in this cell type. In particular, we identified mutations in key DNA repair

11  genes as well as a statistically significant excess of DNVs in *IGLL5* [48,49]. This gene is found to be

12  mutated in B-cell lymphomas and protein-coding DNVs are identified in 27 individuals in this

13  cohort; all of which have >100 overall DNVs. From our work, we identify two contributing factors

14  causing these higher levels of DNVs, one is the mutation of DNA repair genes while the second is

15  an excess of Epstein-Barr Viral load. Future work using long-read sequencing and *de novo*

16  assemblies will be imperative to identify complete viral integration in these genomes as integration

17  sites can have impacts on cell line stability.  One unexpected consequence of B-cell lymphoma

18  mutation signatures in some individuals from the 1000G would be a new pathway to study the

19  mechanisms and biology of the development of this cancer.

20

21        In addition to the DNA repair gene DNVs, we identified fifteen pathogenic or likely-

22  pathogenic DNVs that had already been implicated in a database of clinical variation (Clinvar).

23  This calls into question the use of the 1000G data as a control for both B-cell lymphomas and more

1   generally for DNVs identified in clinical patients. More importantly, the extensive spectrum of

2   DNVs that can appear in a cell line call into question the use of control samples derived from

3   lymphoblastoid cell lines. Currently, to our knowledge the Genome in a Bottle and Human

4   Pangenome Reference Consortium (HPRC) are building reference databases and pangenomes

5   using DNA from lymphoblastoid cell lines. Although it does seem that the use of blood for some

6   samples      was      at      least      initially      discussed      for      the      HPRC

7   (https://www.genome.gov/Pages/Research/Sequencing/Meetings/HGR_Webinar_Summary_Mar

8   ch1_2018.pdf), it does appear the project has defaulted to using lymphoblastoid cell lines. We find

9   it is imperative that these efforts consider utilizing native DNA isolated from blood as the source

10  or utilize a family-based design to identify and remove DNVs. In this way, the highest quality

11  references can be built that will stand the test of time. Finally, we recommend that much like the

12  Simons Simplex Collection, that studies assessing DNVs in individuals with a particular phenotype

13  of interest, also sequence DNA from blood cells and not DNA post-culturing of lymphoblastoid

14  cell lines.

15

1    **METHODS**

2    *Software Code availability*

3    The description for HAT (**H**are **A**nd **T**ortoise) can be found at

4    https://github.com/TNTurnerLab/HAT.  Hare, which was used for analyses in this paper are

5    present at https://github.com/TNTurnerLab/GPU_accelerated_de_novo_workflow, v1.0. We also

6    developed a fully open-source CPU-based version of the code that does not require the NVIDIA

7    Parabricks license, Tortoise,  and it is available at  https://github.com/TNTurnerLab/Tortoise.  We

8    found that Tortoise is just as accurate as Hare, with high level of overlap between the two versions

9    when tested on NA12878 and the monozygotic twin pair (Figure S33).

10

11    *1000 Genomes trio whole-genome sequencing dataset*

12    As described previously [37], a total of 602 trios from the 1000 Genomes Project (1000G)

13    were whole-genome sequenced, from lymphoblastoid cell line DNA, at the New York Genome

14    Center. We downloaded the publicly available aligned data files (crams), totaling around 27TB,

15    onto the Washington University Information Technology's Research Infrastructure Services (RIS),

16    a LSF based, high compute server for further analysis described below.  The download locations

17    are                                          described                                          here

18    http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/1000G_

19    2504_high_coverage.sequence.index                           and                           here

20    http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/1000G_

21    698_related_high_coverage.sequence.index. Details of the 602 trios are found in Supplemental

22    Table S2.

23

1    *Simons Simplex Collection whole-genome sequencing dataset*

2    We downloaded Simons Simplex Collection whole-genome sequencing alignment files

3    (crams) from SFARI Base using Globus, totaling around 239TB, onto the RIS. Importantly, these

4    genomes were sequenced, from DNA derived from blood, at the New York Genome Center [35]. We

5    utilized the crams as the starting point for running in HAT. In total, we assessed 8,922 individuals

6    from both quad (unaffected father, unaffected mother, one child with autism, one child without

7    autism) and trio (unaffected father, unaffected mother, one child with autism) families resulting in

8    a total of 4,216 parent-child sequenced trios. The following individuals were not present in the

9    Globus link and were excluded from the study: SSC03147, SSC03138, SSC03133, SSC03146,

10    SSC06708, SSC06703, SSC06699.

11

12    *Single-nucleotide variant and insertion/deletion calling*

13    The NVIDIA Parabricks program version 3.0.0 was utilized to call single-nucleotide

14    variants (SNVs) and small insertions/deletions (indels) with GATK [39] version 4.1.0 and Google's

15    DeepVariant [40] version 0.10 using default parameters (note for DeepVariant the model_type

16    utilized is WGS). The reference genome utilized for these analyses was

17    GRCh38_full_analysis_set_plus_decoy_hla.fa as the data was originally mapped to this reference

18    genome [37]. For each individual, a GVCF was generated for these two variant callers. The GVCFs

19    were then genotyped, on a per trio basis, using the GLnexus [41] version 1.2.6 joint genotyper using

20    prebuilt configs for each respective caller. Post-calling, we checked the counts of all variants and

21    heterozygous variants per chromosome in each individual as a quality check (Figure S34).

22

23

1

*de novo variant calling*

DNVs were called by identifying all putative DNVs in GATK and DeepVariant based on the parent and child genotypes, respectively. Specifically, the parent genotypes had to be homozygous for the reference allele (i.e., 0/0) and the child had to be, at a minimum, heterozygous for the alternate allele (e.g., 0/1, 1/1). DNVs identified in both GATK and DeepVariant (intersection of the two callers) were then identified and further filtering was carried out as follows: depth, in each trio member, at the DNV position had to be $>= 10$, the genotype quality of the DNV had to be $> 20$, the DNV had to have an allele balance $> 0.25$, and there could be no presence of the DNV allele present in any reads in the parents. Finally, we removed DNVs in low complexity regions, centromeres, and recent repeats from further analysis.

To assess the quality of our DNVs, we manually scored 3,980 sites, by visualizing the underlying read data in each trio member, with SAMtools version 1.9 tview. To score these sites, we looked at the first column (variant location in the read data as seen in tview images) of both parents and the proband sample to see what variants were present (example shown in Figure S35). If there was any variant in the first column of the mother or father, regardless of quality, that matched the main variant in the proband's first column, then we denoted the variant as maternal, paternal, or both depending on whether it was the mother's variant that matched the proband or the fathers or both parents. If the main variant in the first column of the parental samples did not match the proband's variant, then we knew this sample would be a DNV, thus verifying our results.

24

1  As a second check of our DNVs, we randomly sampled 25 DNV sites in NA12878 and

2  performed Sanger sequencing in NA12878 and parents (NA12891, NA12892). Primers were

3  designed using Primer3Plus (https://primer3plus.com) to target each of the 25 variants. PCR

4  reactions were run using the primers, genomic DNA for individuals NA12878 (Coriell tube label

5  NA12878 * N44 12/02/2019), NA12891 (Coriell tube label NA12891 * H3 7/25/2019), and

6  NA12892 (Coriell tube label NA12892 * F3 8/6/2019), and Thermo Scientific Phusion High-

7  Fidelity PCR Master Mix with HF Buffer. All PCR products underwent PCR clean-up and Sanger

8  sequencing through Genewiz (https://www.genewiz.com). Trace files with the Sanger sequencing

9  data were assembled and visualized as chromatograms using Sequencher 5.4.6

10  (http://www.genecodes.com). For 24 of the variants, the result from Sanger sequencing was clear.

11  However, for site chr11_134531608_C_G we saw evidence of the alternate allele at a low

12  frequency. To test whether this signal was real, we pursued deep sequencing of the amplicon on

13  an Oxford Nanopore Technologies (ONT) MinION sequencer as follows. PCR products for

14  amplicon chr11_134531608_C_G, in each of the three individuals, underwent purification using

15  the QIAquick PCR Purification Kit. A library of the purified products was prepared using the

16  Oxford Nanopore Technologies (ONT) Rapid Barcoding Kit (SQK-RBK004). Sequencing of the

17  library was performed using the ONT MinION sequencer and the MinKNOW software. The fastq

18  output files containing the sequencing data for all three samples were mapped to the amplicon

19  reference sequence using minimap2 [50] (version 2.21) and all had coverage depth > 100x. A bam

20  file and indexed bam file were generated for each sample using SAMtools [51] (version 1.9). The

21  bam files were then visualized using the Integrated Genomics Viewer [52] to determine the count of

22  each nucleotide base at the variant position.

23

1

*Phasing of de novo variants*

We utilized Unfazed version 1.0.2 (https://github.com/jbelyeu/unfazed) [53] to phase the *de novo* variants in our study with regard to the parent-of-origin chromosome. First, a bed file containing *de novo* variants was generated for each individual. Second, the *de novo* bed file, DeepVariant full genome trio VCF, and the alignment files for all trio members were run through Unfazed. Since Unfazed uses different approaches to phasing on the X chromosome in males and females, we only focus on phased variants on the autosomes in this study.

*NA12878 additional datasets*

We identified additional high-coverage whole-genome sequencing data from NA12878 from the SRA (https://www.ncbi.nlm.nih.gov/sra) and other sources. These included SRA data SRR944138 from 2012 and SRR952827 from 2013, McDonnell Genome Institute data gerald_HFKWMDSXX and H_IJ-NA12878 both from 2018, and the high-coverage data from 2020. To avoid differences due to coverage, we downsampled all datasets to 30x using SAMtools. All data was re-mapped to build 38 using SpeedSeq [54] version 0.1.2 and run through the DNV workflow using the NA12891 and NA12892 parental WGS data from 2020 1000G. We again did a count check for total and heterozygous variants per chromosome (Figure S36).

*Phylogenetic tree of de novo variants*

To assess the differences between different NA12878 replicates we built a multi-sequence FASTA file where each FASTA represents the aggregate of all possible DNVs identified in this individual. The specific steps to build the tree were as follows: 1) we first merged the samples

26

1  together and converted the genotypes for each DNV site from 0/0 or 0/1 to the nucleotide

2  counterparts (e.g., AA, CG, TC) for all of the NA12878 samples; 2) next we converted these

3  genotype symbols to their IUPAC code; 3) we then collapsed the IUPAC symbols into a sequence

4  per sample and placed them into a FASTA file. We also included a reference "sample", which was

5  just the reference allele at each DNV site and 4) we used MEGAX [55] version 10.2.4 to create a

6  maximum likelihood phylogenetic tree.

7

8  *Mutation profile assessment*

9  We utilized the deconstructSig [45] software version-1.9.0 inside of Parabricks to perform

10  mutation signature analysis. The prominent signature was chosen for an individual and if there was

11  not one prominent signature than the weights of two signature was equal to or greater than (>=

12  0.31) both signatures were represented in the tables and figures.

13

14  *Karyotype analysis*

15  Read-depth based karyotypes were generated by assessment of the aligned sequence data. First,

16  the number of reads per chromosome was calculated using SAMtools [51] in each individual. Second,

17  the size of each chromosome was generated using the reference genome data and by removing

18  locations of gaps from the reference. Third, the copy number of each of the chromosomes was

19  calculated as follows: ((fold coverage per chromosome) / (fold coverage of chromosome 1))*2.

20

21  *Viral analysis*

22  We ran SAMtools idxstats on all individuals to determine the number of mapped reads to

23  each chromosome. We then calculated the copy number of EBV in each individual as follows:

27

1    EBV copy number = ((mapped reads to EBV * 150 base pairs per read) / length of EBV) / ((mapped

2    reads to chromosome 1 * 150 base pairs per read) / length of chromosome 1)

3

4    *DNV enrichment in genes*

5        To test for DNV enrichment in genes we utilized two methods: chimpanzee-human and

6    denovolyzeR. These were run as previously described [8,56].

7

8    *Annotation of protein-coding DNVs*

9        We uploaded the DNV calls to the open-cravat program (https://opencravat.org/) and

10    specifically identified Clinvar as one of the annotation categories. Rescoring of DNVs in Franklin

11    was performed using Franklin (https://franklin.genoox.com).

12

13

14    **ACKNOWLEDGMENTS**

manager?origin_id=3eff453a-88f4-11eb-954f-752ba7b88ebe&origin_path=%2F). 1000

Genomes Acknowledgement for deep coverage of the extended 3202 genomes (or subset thereof):

The following cell lines/DNA samples were obtained from the NIGMS Human Genetic Cell

Repository at the Coriell Institute for Medical Research: [NA06984, NA06985, NA06986,

NA06989, NA06991, NA06993, NA06994, NA06995, NA06997, NA07000, NA07014,

NA07019, NA07022, NA07029, NA07031, NA07034, NA07037, NA07045, NA07048,

NA07051, NA07055, NA07056, NA07340, NA07345, NA07346, NA07347, NA07348,

NA07349, NA07357, NA07435, NA10830, NA10831, NA10835, NA10836, NA10837,

NA10838, NA10839, NA10840, NA10842, NA10843, NA10845, NA10846, NA10847,

NA10850, NA10851, NA10852, NA10853, NA10854, NA10855, NA10856, NA10857,

NA10859, NA10860, NA10861, NA10863, NA10864, NA10865, NA11829, NA11830,

NA11831, NA11832, NA11839, NA11840, NA11843, NA11881, NA11882, NA11891,

NA11892, NA11893, NA11894, NA11917, NA11918, NA11919, NA11920, NA11930,

NA11931, NA11932, NA11933, NA11992, NA11993, NA11994, NA11995, NA12003,

NA12004, NA12005, NA12006, NA12043, NA12044, NA12045, NA12046, NA12056,

NA12057, NA12058, NA12144, NA12145, NA12146, NA12154, NA12155, NA12156,

NA12234, NA12239, NA12248, NA12249, NA12264, NA12272, NA12273, NA12274,

NA12275, NA12282, NA12283, NA12286, NA12287, NA12329, NA12335, NA12336,

NA12340, NA12341, NA12342, NA12343, NA12344, NA12347, NA12348, NA12375,

NA12376, NA12383, NA12386, NA12399, NA12400, NA12413, NA12414, NA12485,

NA12489, NA12546, NA12707, NA12708, NA12716, NA12717, NA12718, NA12739,

NA12740, NA12748, NA12749, NA12750, NA12751, NA12752, NA12753, NA12760,

NA12761, NA12762, NA12763, NA12766, NA12767, NA12775, NA12776, NA12777,

14
15

## REFERENCES

1. Ségurel, L., Wyman, M.J. & Przeworski, M. Determinants of mutation rate variation in the human germline. *Annu Rev Genomics Hum Genet* **15**, 47-70 (2014).
2. Feliciano, P. *et al.* Exome sequencing of 457 autism families recruited online provides evidence for autism risk genes. *NPJ Genom Med* **4**, 19 (2019).
3. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum disorder. *Nature* (2014).
4. O'Roak, B.J. *et al.* Exome sequencing in sporadic autism spectrum disorders identifies severe de novo mutations. *Nat Genet* **43**, 585-9 (2011).
5. O'Roak, B.J. *et al.* Sporadic autism exomes reveal a highly interconnected protein network of de novo mutations. *Nature* **485**, 246-50 (2012).
6. Turner, T.N. *et al.* Genomic Patterns of De Novo Mutation in Simplex Autism. *Cell* **171**, 710-722.e12 (2017).
7. Turner, T.N. *et al.* Genome Sequencing of Autism-Affected Families Reveals Disruption of Putative Noncoding Regulatory DNA. *Am J Hum Genet* **98**, 58-74 (2016).
8. Turner, T.N. *et al.* Sex-Based Analysis of De Novo Variants in Neurodevelopmental Disorders. *Am J Hum Genet* **105**, 1274-1285 (2019).
9. Sanders, S.J. *et al.* De novo mutations revealed by whole-exome sequencing are strongly associated with autism. *Nature* **485**, 237-41 (2012).
10. De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in autism. *Nature* **515**, 209-15 (2014).
11. Helbig, K.L. *et al.* Diagnostic exome sequencing provides a molecular diagnosis for a significant proportion of patients with epilepsy. *Genet Med* (2016).
12. Allen, A.S. *et al.* De novo mutations in epileptic encephalopathies. *Nature* **501**, 217-21 (2013).
13. Kaplanis, J. *et al.* Integrating healthcare and research genetic data empowers the discovery of 28 novel developmental disorders. *bioRxiv*, 797787 (2020).
14. DDD. Prevalence and architecture of de novo mutations in developmental disorders. *Nature* **542**, 433-438 (2017).
15. de Ligt, J. *et al.* Diagnostic exome sequencing in persons with severe intellectual disability. *N Engl J Med* **367**, 1921-9 (2012).
16. Gilissen, C. *et al.* Genome sequencing identifies major causes of severe intellectual disability. *Nature* **511**, 344-7 (2014).
17. Homsy, J. *et al.* De novo mutations in congenital heart disease with neurodevelopmental and other congenital anomalies. *Science* **350**, 1262-6 (2015).
18. Zaidi, S. *et al.* De novo mutations in histone-modifying genes in congenital heart disease. *Nature* **498**, 220-3 (2013).
19. Sifrim, A. *et al.* Distinct genetic architectures for syndromic and nonsyndromic congenital heart defects identified by exome sequencing. *Nat Genet* **48**, 1060-5 (2016).
20. Besenbacher, S. *et al.* Novel variation and de novo mutation rates in population-wide de novo assembled Danish trios. *Nat Commun* **6**, 5969 (2015).
21. Boomsma, D.I. *et al.* The Genome of the Netherlands: design, and project goals. *Eur J Hum Genet* **22**, 221-7 (2014).
22. Chesi, A. *et al.* Exome sequencing to identify de novo mutations in sporadic ALS trios. *Nat Neurosci* **16**, 851-5 (2013).
23. Conrad, D.F. *et al.* Variation in genome-wide mutation rates within and between human families. *Nat Genet* **43**, 712-4 (2011).
24. Dimassi, S. *et al.* Whole-exome sequencing improves the diagnosis yield in sporadic infantile spasm syndrome. *Clin Genet* **89**, 198-204 (2016).
25. Dong, S. *et al.* De novo insertions and deletions of predominantly paternal origin are associated with autism spectrum disorder. *Cell Rep* **9**, 16-23 (2014).

26. Fromer, M. *et al.* De novo mutations in schizophrenia implicate synaptic networks. *Nature* **506**, 179-84 (2014).

27. Iossifov, I. *et al.* De novo gene disruptions in children on the autistic spectrum. *Neuron* **74**, 285-99 (2012).

28. Jonsson, H. *et al.* Parental influence on human germline de novo mutations in 1,548 trios from Iceland. *Nature* **549**, 519-522 (2017).

29. McCarthy, S.E. *et al.* De novo mutations in schizophrenia implicate chromatin remodeling and support a genetic overlap with autism and intellectual disability. *Mol Psychiatry* **19**, 652-8 (2014).

30. Narzisi, G. *et al.* Accurate de novo and transmitted indel detection in exome-capture data using microassembly. *Nat Methods* **11**, 1033-6 (2014).

31. Ramu, A. *et al.* DeNovoGear: de novo indel and point mutation discovery and phasing. *Nat Methods* **10**, 985-7 (2013).

32. Franke, K.R. & Crowgey, E.L. Accelerating next generation sequencing data analysis: an evaluation of optimized best practices for Genome Analysis Toolkit algorithms. *Genomics Inform* **18**, e10 (2020).

33. An, J.Y. *et al.* Genome-wide de novo risk score implicates promoter variation in autism spectrum disorder. *Science* **362**(2018).

34. Werling, D.M. *et al.* An analytical framework for whole-genome sequence association studies and its implications for autism spectrum disorder. *Nat Genet* **50**, 727-736 (2018).

35. Wilfert, A.B. *et al.* Recent ultra-rare inherited variants implicate new autism candidate risk genes. *Nat Genet* (2021).

36. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).

37. Byrska-Bishop, M. *et al.* High coverage whole genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. *bioRxiv*, 2021.02.06.430068 (2021).

38. Rehm, H.L. *et al.* ACMG clinical laboratory standards for next-generation sequencing. *Genet Med* **15**, 733-47 (2013).

39. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-303 (2010).

40. Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural networks. *Nat Biotechnol* **36**, 983-987 (2018).

41. Yun, T. *et al.* Accurate, scalable cohort variant calls using DeepVariant and GLnexus. *bioRxiv*, 2020.02.10.942086 (2020).

42. Zook, J.M. *et al.* Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls. (2014).

43. Sasani, T.A. *et al.* Large, three-generation human families reveal post-zygotic mosaicism and variability in germline mutation accumulation. *Elife* **8**(2019).

44. Gibbs, R.A. *et al.* The International HapMap Project. *Nature* **426**, 789-796 (2003).

45. Rosenthal, R., McGranahan, N., Herrero, J., Taylor, B.S. & Swanton, C. DeconstructSigs: delineating mutational processes in single tumors distinguishes DNA repair deficiencies and patterns of carcinoma evolution. *Genome Biol* **17**, 31 (2016).

46. Landrum, M.J. *et al.* ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res* **46**, D1062-d1067 (2018).

47. Abecasis, G.R. *et al.* A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061-73 (2010).

48. Chen, F., Zhang, Y. & Creighton, C.J. Systematic identification of non-coding somatic single nucleotide variants associated with altered transcription and DNA methylation in adult and pediatric cancers. *NAR Cancer* **3**, zcab001 (2021).

49. Pedrosa, L. *et al.* Proposal and validation of a method to classify genetic subtypes of diffuse large B cell lymphoma. *Sci Rep* **11**, 1886 (2021).

50. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-3100 (2018).

51. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-9 (2009).
52. Robinson, J.T. *et al.* Integrative genomics viewer. *Nat Biotechnol* **29**, 24-6 (2011).
53. Belyeu, J.R., Sasani, T.A., Pedersen, B.S. & Quinlan, A.R. Unfazed: parent-of-origin detection for large and small <em>de novo</em> variants. *bioRxiv*, 2021.02.03.429658 (2021).
54. Chiang, C. *et al.* SpeedSeq: ultra-fast personal genome analysis and interpretation. *Nat Methods* **12**, 966-8 (2015).
55. Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular Evolutionary Genetics Analysis across Computing Platforms. *Mol Biol Evol* **35**, 1547-1549 (2018).
56. Coe, B.P. *et al.* Neurodevelopmental disease genes implicated by de novo mutation and copy number variation morbidity. *Nat Genet* **51**, 106-116 (2019).