# 1 STIX: Long-reads based Accurate Structural
# 2 Variation Annotation at Population Scale

3 Xinchang Zheng[1], Murad Chowdhury[2], Behzod Mirpochoev[3], Aaron Clauset[2,3], Ryan M
4 Layer[2,3], Fritz J Sedlazeck[1,4,5]

5

6 1: Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX, USA
7 2: Biofrontiers Institute, University of Colorado, Boulder, CO, USA
8 3: Computer Science Department, University of Colorado, Boulder, CO, USA
9 4: Department of Molecular and Human Genetics, Baylor College of Medicine, TX, USA
10 5: Department of Computer Science, Rice University, Houston, TX, USA

11

12 Corresponding: ryan.layer@colorado.edu, fritz.sedlazeck@bcm.edu

# 13 Abstract

14 The prioritization of Structural Variants (SV), which is needed to rank and identify
15 potential pathogenic alleles, is still in its infancy. This is exemplified over gnomAD only
16 being able to annotate 33.5% of GIAB SVs. To overcome this, we present the first
17 long-read based annotation resource for both GRCh38 and CHM13-T2T reference
18 using STIX. In contrast to previous methods, STIX indexes SV-informative long-reads
19 themselves, can thus be easily extended and accurately annotate all SV types
20 including insertions. STIX successfully annotated 95.9% of GIAB Tier1 SVs. STIX
21 further improved cancer based SV prioritization by highlighting 3,563 SV from
22 COSMIC being common in the population. We further showcase that mosaic SV can
23 be independently gained and may be widely spread throughout the population. This
24 highlights the need for accurate SV population frequency annotation to further facilitate
25 the adoption of SV via long-read sequencing in medical research and clinical
26 applications.

# 27 Introduction

28 Structural variants (SVs) are genomic alterations larger than 50 bp, typically including
29 deletions, duplications, insertions, inversions, translocations, and their
30 combinations[1,2]. SVs impact mendelian, neurological, and cardiovascular diseases, [3–
31 5], and are the driver mutations in many cancers including leukemia[6], lung cancer[7] and
32 triple-negative breast cancer[8]. SVs often occur in the cancer genomes at lower
33 frequencies but dynamically evolve under selection pressure, resulting in highly
34 heterogeneous tumor SV profiles[8–11]. SV detection is improving, with clear benefits
35 from using long-read sequencing[2,12–15]. However, previously published genomic
36 annotation tools such as VEP[16], ANNOVAR[17], or SnpEff[18] mainly focus on small

1

37  variants, while SV prioritization and interpretation are lacking behind, which limits our
38  ability to determine the role of SVs in disease studies.

39  Population frequency annotation is an effective approach to prioritize SVs since a
40  variant that is common in the population is less likely to be pathogenic[19]. Current
41  approaches such as dbVar[20] and gnomAD[21,22] report SV population frequencies, but
42  are limited for two main reasons. First, they are based on short-read sequencing,
43  which includes only a fraction of the SVs when compared to using long-reads[12].
44  Second, they are built upon merged and heavily filtered VCFs that do not depict a
45  harmonized representation[23]. Merging based on the rough overlap (e.g., 70%)[24] can
46  add significant bias when the caller and sequencing technologies (e.g insert size) vary.
47  Over-merging can occur due to imprecise thresholds, which leads to a wrong
48  interpretation of the allele itself[25]. While the identification of SVs is constantly
49  improving, the population frequency annotation remains challenging, resulting in
50  ineffective prioritization and impact prediction[26].

51  Long-read sequencing has significantly improved our understanding of the occurrence
52  and impact of SVs[1,27]. It has better alignments in repetitive regions such as tandem
53  repeats where 70% or more of SV are located[1,28]; is able to span the alleles; and has
54  generally less sequencing biases than short read sequencing[1]. Despite the increased
55  sequencing costs and sample requirements[1] , several current efforts in population
56  long-read sequencing have emerged to provide novel SV insights[2,29–31], and the time
57  is right to generate a harmonized SV annotation resource to improve SV population
58  frequency estimation. The challenge is the harmonization and the future readiness of
59  such a resource as current annotation methods still require version specific callers and
60  other filtering heuristics. To overcome some of these challenges, we recently
61  introduced STIX for short-read sequencing that overcomes limitations of merged VCF
62  files, by indexing and searching SV informative reads directly[26]. While the initial
63  version of STIX improved SV prioritization, it did not support insertions nor was it able
64  to utilize long-read sequencing data.

65  In this study, we expanded STIX to long-reads, allowing for the annotation of SV with
66  population-scale evidence across different long-read sequencing platforms. In this
67  update, we fully support all types of SVs, including insertions, and benchmarked the
68  performance of STIX not only for germline SV annotation but also demonstrated that
69  STIX can also improve mosaic SV annotation. To accomplish this, we generated the
70  first long-read annotation resource across both GRCh38 and CHM13-T2T[32] reference
71  utilizing 1,108 samples spanning all ethnicities. We demonstrate the utility of this
72  catalog together with STIX for variant prioritization across cancer and somatic variant
73  calling experiments. STIX has now solved the long-read annotation issues and
74  enables a future ready resource for the genomic community that can be easily
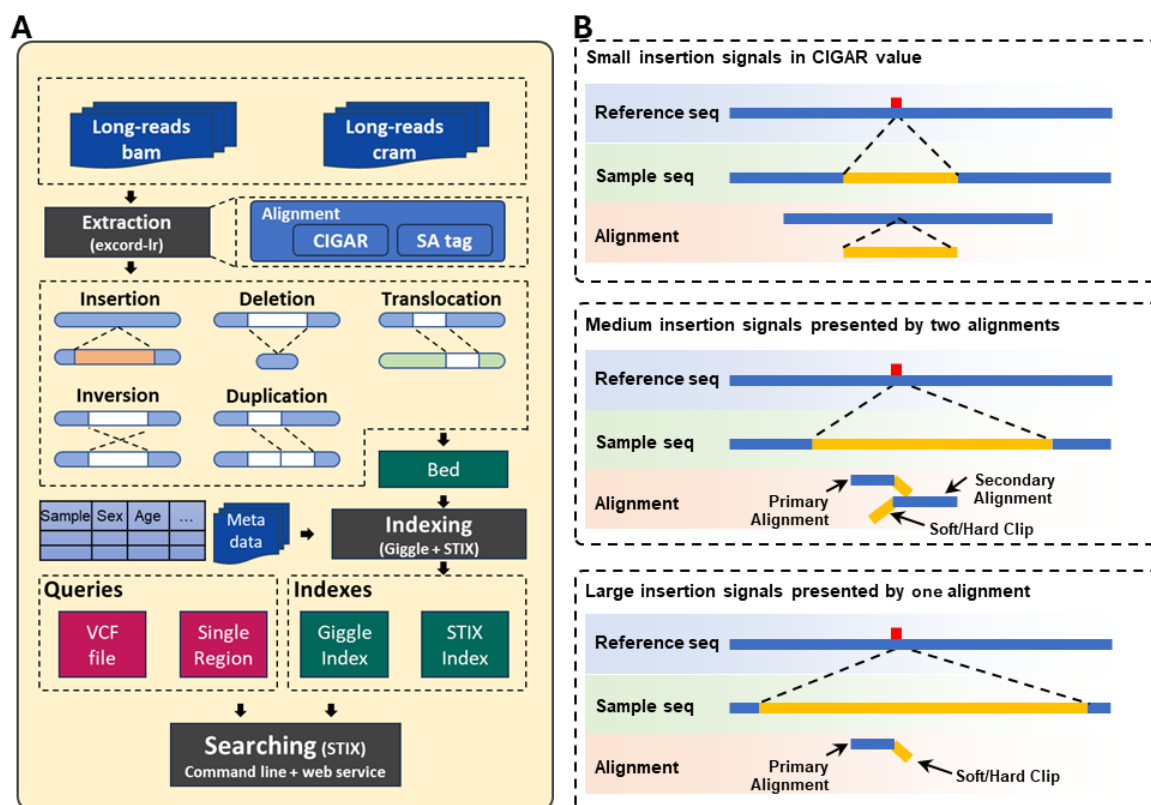75  accessed, extended and updated.

# Results

## Comprehensive and accurate SV annotation using STIX

SV annotation is important for rare disease research, medical applications, and genotype-phenotype related studies[33–35]. Our extensions of STIX to long-reads builds on our existing short-read platform[26] by supporting all SV types. **Fig. 1A** summarizes the three key steps of STIX. In the first step, STIX screens a bam file using a new implementation of excord called excord-lr (**see methods**). Excord-lr scans the bam file's CIGAR string and split-read signal to find potential SV information. The results are then stored in a bed file for each sample, reducing the storage burden by 99.85%. We have expanded excord-lr to include support for insertions, which are encoded in long-reads in three different ways based on the size of the event (**Fig. 1B**). Small insertions are included in the CIGAR string, split read events represent mid-sized insertions, and insertions larger than the read length generate a single primary alignment with large unaligned segments.

In the second step, STIX creates a compressed index[36] of the previously extracted SV data. Sample-level metadata, such as ethnicity, gender, and phenotype can be added to the index and later retrieved. To enable the resource's future-ready growth, we have implemented a sharding approach that allows the simultaneous handling of independent indexes and queries.

The third step annotates novel SVs by querying the previously generated indexes. Here STIX assesses if a SV is in the indexed population based on SV type, position, and length. In this step, we can integrate over short- and long-reads and annotate all SV types, including insertions. See methods for a detailed description about the individual steps. STIX is open source and available at https://github.com/ryanlayer/stix together with the annotation index at https://stix.colorado.edu.
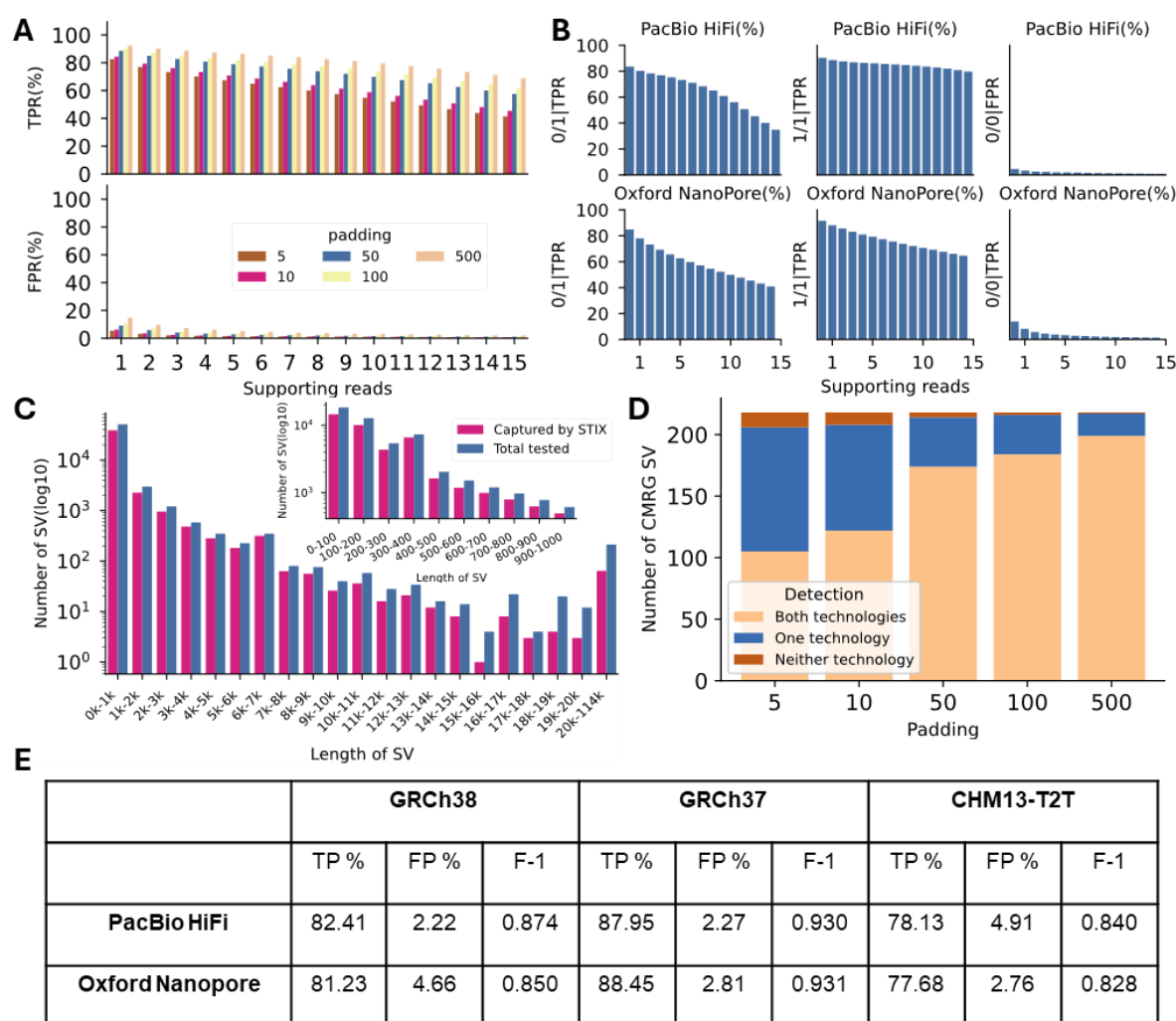
3

**Fig.1 Overview of STIX: A)** *An overview of the main steps in STIX SV annotation* **B)** *Insertion signals extracted by excord-lr in long-reads dataset are classified by size. Small insertions are in CIGAR string (top). Longer insertions are either one primary alignment and one supplementary alignment (middle), or a single primary alignment (bottom).*

# Performance assessment of germline structural variation annotation

STIX returns the population evidence depth for a query SV. To account for the imprecision in SV calling, STIX uses a *padding* parameter to determine how many bases up and downstream of a variant to consider. To derive this threshold and the minimum evidence depth for reporting an SV, we used the HG002 SV benchmark from the Genome in a Bottle (GIAB) project (Tier1GIAB SVs) and created an HG002 STIX index for long-read sequencing data from Oxford Nanopore (ONT) and PacBio HiFi (Pb-HiFi)[37]. SV calls in the benchmark and recovered by STIX were considered true positives (TP). Calls not recovered by STIX are false negatives (FN). Using the GIAB phased assembly SV pipeline, we generated calls for HG00733[28,25]. Calls that were specific to HG00733 and recovered by STIX in the HG002 indexes were considered false positives (FP) (**see methods**).

We tested the robustness of STIX using SV benchmarks from three reference genomes, including the well established HG002 (v0.6) on GRCh37 and two prototype

4

126    SV benchmarks from GIAB for GRCh38 and CHM13-T2T (V0.012-20231107). For FP,
127    we called SVs in HG00733 based on its assembly compared to GRCh37, GRCh38
128    and CHM13-T2T, excluded common SV with HG002 and queried these HG00733
129    private SV against the HG002 index.
130



131
132

***Fig.2 Performance assessment for Germline SV of STIX based on GIAB:***
*Performance metrics for correctly genotyping SVs in different platforms(Pb-HiFi and*
*ONT). All details can be found at **Supplementary Table S2 and S3**. **A)** comparison*
*of Tier1GIAB SV (TPR and FPR) with different support reads(x-axis), padding(color of*
*the bar). **B)** Impact of genotypes on STIX across different platforms. Each sub-figure*
*represents either TPR(first two columns) or FPR(third column) for Pb-HiFi(top row)*
*and ONT(bottom row). X-axis shows the supporting reads. Y-axis is TPR or FPR in*
*percent in each sub-figure. **C)** Assessment of capture rate (y-axis) of STIX for*
*Tier1GIAB SVs across different sizes(x-axis). SVs were grouped by their length. Sub-*
*figure in the right top corner shows the SVs no longer than 1kb with 100bp as bin width.*
*For each group, two bars represent tested(blue) and captured by STIX(red)*
*respectively. **D)** Performance evaluation of STIX for detecting CMRG SV(y-axis) with*
*the increase of padding(x-axis). Colors in each bar represent the detection*
*status(orange:detected in both Pb-HiFi and ONT; blue: detected in either Pb-HiFi or*
*ONT; yellow: not detected in all platforms). **E)** STIX performance with different*

5

148 *reference genome and sequencing platforms*

149

150 As expected, TP and FP rates decreased as the minimum number of supporting reads
151 and the padding increased(**Fig. 2A**). Interestingly, the padding parameter impacted
152 TP, but had minimal effect on FP. We observed slightly different impacts on parameter
153 choice across the two sequencing platforms(**Supplementary FIgure S5**). Pb-HiFi
154 exhibited consistent performance independent of the padding, while ONT benefited
155 from larger paddings. This difference may be due to higher error rates given that the
156 GIAB ONT SVs were produced with the older R9 basecaller with higher error rate[38,39].

157 Based on these results, we recommend a padding of 100 bp and a minimum read
158 support of 5 reads for germline variants. With these parameters, STIX achieved an
159 overall TP rate of 81.82% and a FP rate of 3.44%. These parameters further minimized
160 the differences between Pb-HiFi and ONT across all three reference genomes(**Fig.
161 2E**, additional comparisons in **Supplementary Section 1** and **Supplementary Table
162 S2**).

163

164 For insertions and deletions, the only two SV types reported in the GIAB HG002
165 benchmark, STIX had a higher TP and FP rate (TPR and FPR) as well as F-1 score
166 in deletions when compared with insertions (**Supplementary Table S1 and
167 Supplementary Fig. S6**). In addition to performing well in genome-wide benchmarks,
168 STIX recovered SVs in regions that are traditionally difficult to resolve. Among the 218
169 SVs in the Challenging Medical Relevant Genes(CMRG) GIAB benchmark[40], STIX
170 correctly annotated 89.45%(195) of the SVs based on ONT and 94.04%(205) of SVs
171 based on Pb-HiFi. When the two sequencing   technologies were combined, STIX
172 re-identified 216 out of the 218 SVs (99.08%)(**Fig. 2D**). Across the three references,
173 STIX (TPR: 79.7%, FPR:3.27%) also performed similarly to two state-of-art single
174 sample only SV genotypers, Sniffles2[41] (TPR: 87.2%, FPR: 4.65%) and cuteSV[42](TPR:
175 82.3%, FPR: 3.41%) (**Supplementary Table S1**). We further benchmarked STIX with
176 its short-read index that showed clear improvements across the long-read version
177 (**Supplementary Section 5 and Supplementary Figure S9-S15**).

178 When categorized by zygosity, Pb-HiFi and ONT showed similar TP rates, with
179 heterozygous SVs being lower than homozygous variants (**Fig. 2B**, left). We
180 anticipated this difference due to fewer reads supporting the variant allele in
181 heterozygous locations. The FP rate for Pb-HiFi was about half that of ONT (**Fig. 2B**,
182 right). Among 23,114 heterozygous variants, 5,473 homozygous variants from HG002,
183 and 6,413 wild-type variants from the HG00733 call set(see **Supplementary Section
184 1**), the TP and FP rates for Pb-HiFi were 87.32% and 2.18%, respectively, while for
185 ONT, they were 87.14% and 4.41%, respectively (**Supplementary Table S4**).

186 To assess the impact of SV length on the performance, we grouped the GIAB HG002
187 SVs in different size bins and calculated the recall per 1 kbp bin. As illustrated in **Fig.
188 2C**, STIX captures 78.33%(average of each length category) of SVs that are less than
189 10kb (73.09% of SVs less than 15kb). We further stratified those SVs between 50bp-
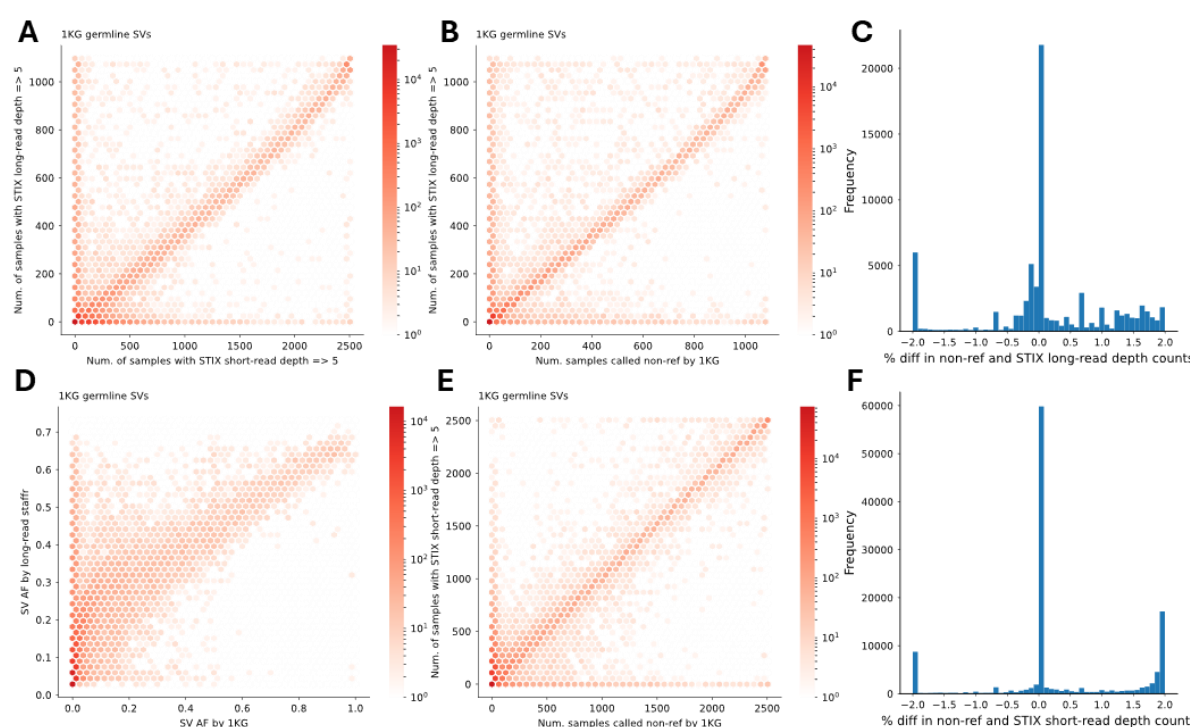
190    1kbp into 100bp bins as the SV in this size amount to the most number of SV. We
191    observed a good performance in SVs that overlap with Alu (300-400 bp) or LINE
192    elements (6,000-7,000 bp). STIX has limited performance for SVs larger than 15kb but
193    still capture 35.26% of them.

194    Overall, STIX performed well across both sequencing platforms, regardless of SV
195    zygosity or length. Additionally, STIX can now accurately annotate insertions, a feature
196    not available with the same level of precision as other population annotation methods.
197    We also conducted a benchmark assessment for the short-reads version of STIX to
198    demonstrate the consistency of the new version(**Supplementary Section 5**).

## STIX improves the SV annotation at population scale

200    While long-read sequencing has well-known advantages for detecting SVs compared
201    to short-read sequencing[1], as of now, there are no long-read based resources
202    available for annotating SVs. To address this gap, we integrated data from the The
203    1000 Genomes Project (1KG)[43] and The Human Pangenome Reference Consortium
204    (HPRC) project[44] (**Methods and Supplementary Section 3**) to produce a STIX index
205    with 1,108 samples(1,104 for the CHM13-T2T version) from 26 populations and 5
206    super-populations (**Supplementary Table S5 and Supplementary Fig. S16**)
207    sequenced by Pb-HiFi or ONT at 17x coverage or higher[45,46].

208    To assess how accurately STIX can determine the frequency of SVs in a population,
209    we firstly evaluated the consistency of population evidence of 1KG published SV
210    between the long-read index and short-read index with sufficient evidence depth (at
211    least 5) and observed a significant positive correlation(R=0.8, p-value <0.01), as
212    shown in **Fig. 3A**. Moreover, we compared the number of non-reference samples for
213    the same 1KG SVs annotated by the STIX long-read index to those originally reported
214    by the 1KG. There was a significant correlation between the number of samples found
215    by 1KG and STIX that harbored an SV. (R=0.8, p-value<0.001, **Fig. 3B**). About 10%
216    of SVs were more frequent according to 1KG (less than -1 in **Fig. 3C**). These SVs
217    were evenly distributed among SV types (**Supplementary Table S10**) and tended to
218    be smaller (median length 271 bp vs 626 bp for all 1KG SVs). About 20% were more
219    frequent according to STIX (greater than 1 in **Fig. 3C**) and were enriched for inversions
220    (**Supplementary Table S10**), and tended to be longer (median length 6,113 bp). For
221    the remaining SVs, the difference between the STIX and 1KG frequencies was on
222    average 2.2%. Using a Hardy-Weinberg based Gaussian mixture model
223    (**Supplemental Section 6**), we estimated the allele frequency of common (> 10%
224    frequency) SVs that were highly correlated with the 1KG allele frequencies (R=0.73,
225    p-value<0.001, **Fig. 3D**), with a mean difference of 0.6% between our estimated allele
226    frequencies and those reported by 1KG. The STIX long-read results were consistent
227    with STIX short-read results (**Fig. 3E and 3F**), with the notable exceptions that the
228    short-read database had nearly 2X the number of samples (2,504), the long-read index
229    supports insertions, and the long-read mixture model allele frequency estimations
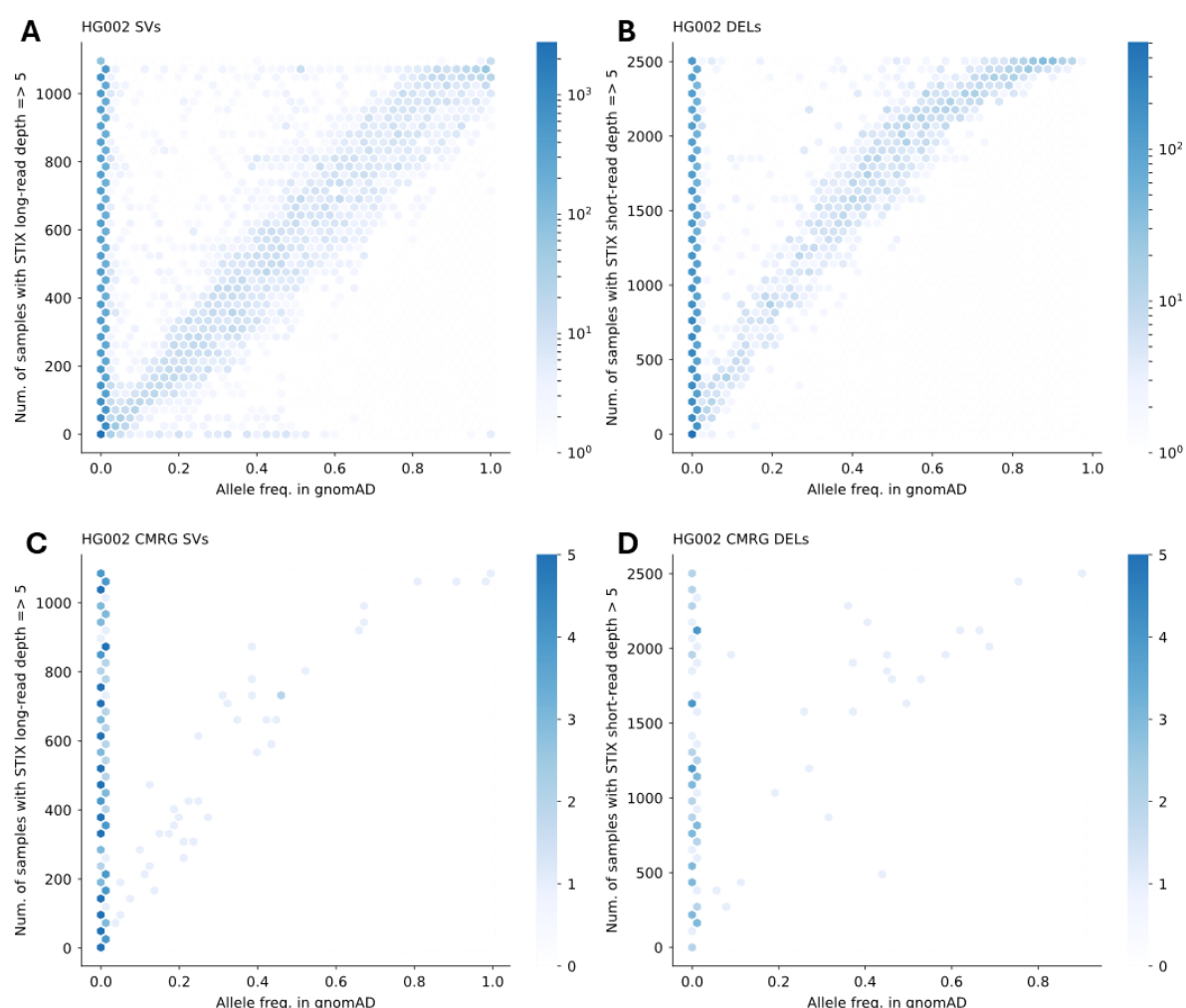230    were better calibrated (**Supplementary Fig. S17**).

7

231



232
233

234 **Fig. 3 Comprehensive resource of large cohorts: A)** *Comparison between the long-*
235 *read index and short-read index.* **B)** *The number of samples with at least 5 STIX long-*
236 *read hits versus the number of samples assigned a non-reference genotype by 1KG.*
237 **C)** *The distribution of average differences between the STIX long-read and 1KG*
238 *frequency estimates.* **D)** *The STIX long-read allele frequency estimated by mixture*
239 *model versus the allele frequencies published by 1KG. R=0.82.* **E)** *The number of*
240 *samples with at least 5 STIX short-read hits versus the number of samples assigned*
241 *a non-reference genotype by 1KG. R=0.64* **F)** *The distribution of average differences*
242 *between the STIX short-read and 1KG frequency estimates.*

243 STIX offers more comprehensive SV frequency annotation than standard catalogs. To
244 measure this improvement we compared the variant frequency annotations from the
245 10,847 samples in the gnomAD[24,47] catalog to the STIX long- and short-read 1KG
246 indices, which had 1,108 and 2,504 samples respectively (the long-read index
247 included 11 PacBio/ONT technical replicates), using the GIAB benchmark sample
248 HG002. **Supplementary Fig. S7** shows the frequency of GIAB SV across the STIX
249 index, which follows the expected distribution. While 99.8% of single nucleotide
250 variants (SNVs) in HG002 appeared in the gnomAD catalog[40], only 33.4% of structural
251 variants (SVs) were found. In contrast the STIX long-read index annotated 95.9% of
252 the HG002 SVs. Among the gnomAD annotations, there was a notable depletion in
253 the number of insertions. While there was evidence in the STIX long-read index for
254 94.6% of the HG002 insertions (the short-read index does not support insertions), only
255 17.3% were in the gnomAD catalog. For deletions, the STIX long- and short-read index
256 annotated 97.8% and 96.8%, respectively, while only 58.6% were in the gnomAD

257  catalog. When considering just the 218 HG002 SVs that overlapped the challaning
258  medically relevant genes (CMRG), the gnomAD catalog only found 69 (31.6%), while
259  the STIX long-read index found all 218, and the STIX-short read index found 97 (41%).
260  These differences are notable considering how many more samples the gnomAD
261  catalog included, and how many of the SVs missing from gnomAD are at high
262  frequency in the 1KG cohort (**Figure 4**).



263
264  ***Fig. 4: A comparison between frequency estimates annotation of the GIAB***
265  ***HG002 benchmark SVs using the gnomAD catalog and STIX long- and short-***
266  ***read indices. A)*** *The number of samples with at least 5 STIX long-read hits versus*
267  *the gnomAD allele frequency.* ***B)*** *The number of samples with at least 5 STIX short-*
268  *read hits versus the gnomAD allele frequency.* ***C and D)*** *The same comparison with*
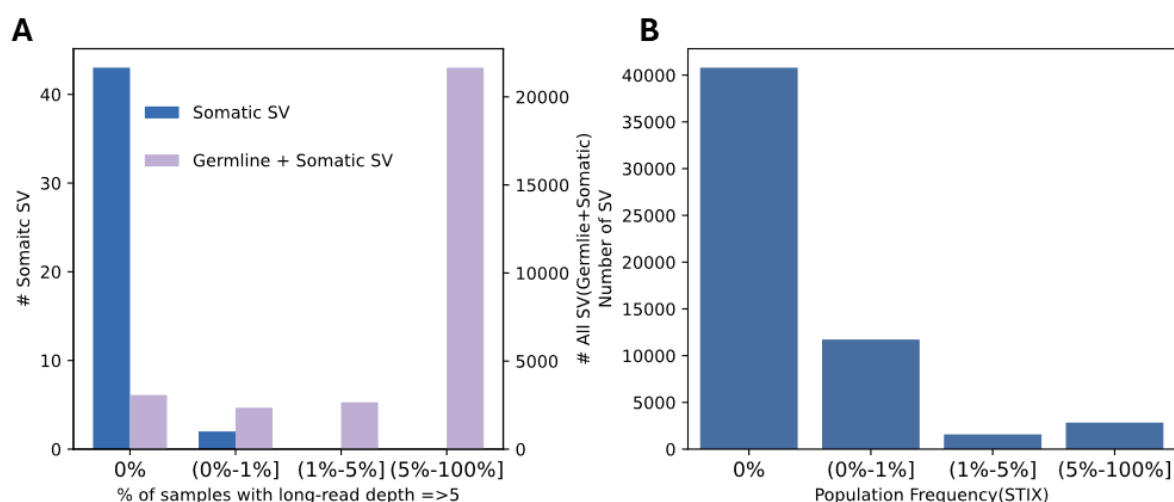269  *the subset of SVs that were in the challaning medically relevant genes (CMRG).*

## Improving somatic cancer SV prioritization

271  SV prioritization, especially in tumor samples, can be complicated by false positives
272  driving the wrong identification of somatic SV[48]. This can be mitigated by improved SV
273  comparison methods, but also by population frequency annotation. For the latter,

9

274 somatic SV that are not driving the cancer can occur naturally throughout the
275 population and thus might be commonly observed[8,49]. To assess the ability of STIX to
276 improve somatic cancer driver SV detection, we first tested this approach at
277 COLO829/COLO829BL, which are well characterized tumor-normal cell lines[50]. We
278 recently postulated somatic SV for COLO829/COLO829BL[50] and were interested in
279 how well STIX could identify potentially cancer only SV. We initially started with all
280 SVs (including 29,674 germline SV and 45 somatic SV) from the tumor
281 sample(COLO829). STIX was able to annotate 89.67%(26,649) of the germline SVs.
282 We use 1% AF as the threshold for common variants and found 81.77%(24,302) of
283 the germline COLO829 SVs can be indeed annotated as common in the population
284 and thus were likely non cancer drivers (**Fig. 5A**, purple bars). This is a great result
285 when for example one does not have a normal matched sample at hand at the right
286 quality or quantity to perform e.g. long-read sequencing. We next focused on the
287 previously postulated 45 SV that were identified as somatic for the COLO829 (cancer
288 sample) to further showcase the benefit of STIX[50]. **Fig. 5A**(blue bars) highlights these
289 45 somatic SV and their proportion of population scale evidence based on STIX.
290 Interestingly, STIX assigned evidence to two of the 45 somatic SV but none of the
291 somatic SV has a higher AF than 1% in our SV index. Thus, highlighting that despite
292 tumor-normal comparison a population annotation such as STIX might further narrow
293 down potential cancer driver mutations compared to likely benign SV.

294 Motivated by this observation, we next investigated if we could further annotate and
295 identify SV that are postulated to be somatic cancer mutations, but might actually be
296 common in the population. We use somatic SVs from the Catalogue of Somatic
297 Mutations in Cancer project(COSMIC)[51] as it summarizes one of the largest cancer
298 datasets. **Fig. 5B** shows the result with respect to the annotatable SV. We annotated
299 46,755 somatic cancer only SVs using STIX and were able to retrieve population
300 frequencies for 13,564 (29.01%). Among them, we identified 3,563 SV (26.27%) as
301 common in the population (>1% AF) indicating their potential non-pathogenic role.
302 These SVs could also represent non-cancer driving somatic SVs that were
303 accumulated during cancer initiation and progression and thus be potentially
304 harmless[52,53]. However, some SV could potentially also represent several falsely
305 identified somatic SV in COSMIC, which is hard to determine. The remaining 10,001
306 (73.73%) SV showed a low proportion of evidence(<1%) as may be expected from
307 somatic benign or cancer SV postulated by COSMIC.

308 Across both examples, STIX showed to be important in further prioritization of cancer
309 vs. benign somatic SV. This is obviously just one of the potential important use cases
310 of STIX for the prioritization of SV across many human genetic diseases.
311

***Fig. 5: STIX application to cancer structural variation prioritization. A)** Proportion of evidence distribution of all SV(purple bar,right y-axis) and somatic only SV(blue bar, left y-axis) from COLO829. STIX effectively narrows down the scope of potential pathogenic structural variants (SVs) to approximately 10%, and includes 96% of true somatic SVs, as identified in the matched normal sample COLO829BL. **B)** Proportion of evidence distribution of 46,755 COSMIC SV annotated by STIX. STIX found evidence for 29.01% of these SVs.*

## Exploring mosaic SVs across population level

Following the identification of potentially common somatic SVs in the population, explored STIX's ability to identify mosaic SVs (i.e., those with a low variant allele fraction (VAF)). To benchmark STIX's ability to assess and annotate mosaic SVs, we created STIX indexes from read sets sampled at different rates ranging from 0.05 to 0.45 from HG002 and HG00733. For each index, we searched for evidence of the HG002 benchmark SVs (**see methods, Supplementary Section 2**). Overall, STIX had an average of 90.61% TP for deletions and 64.59% TP for insertions(**Fig. 6A**). We speculate that the random choice of reads favor the more frequent shorter read length in the distribution of HG002 and thus reduce our performance on insertions. Similar to the germline ONT benchmark, we observed an impact of the padding parameter on the identification ability of STIX and did not observe a bias towards SV length with 95.27% of all SVs captured across the different lengths (**Fig. 6B**). Among the CMRG regions (**Supplementary Fig. S18**), STIX captured 182 SVs (83.49%) despite the SV being present at lower VAF. Overall, STIX showed great performance in annotating mosaic SV despite a greater challenge of lower number of reads present per SV.

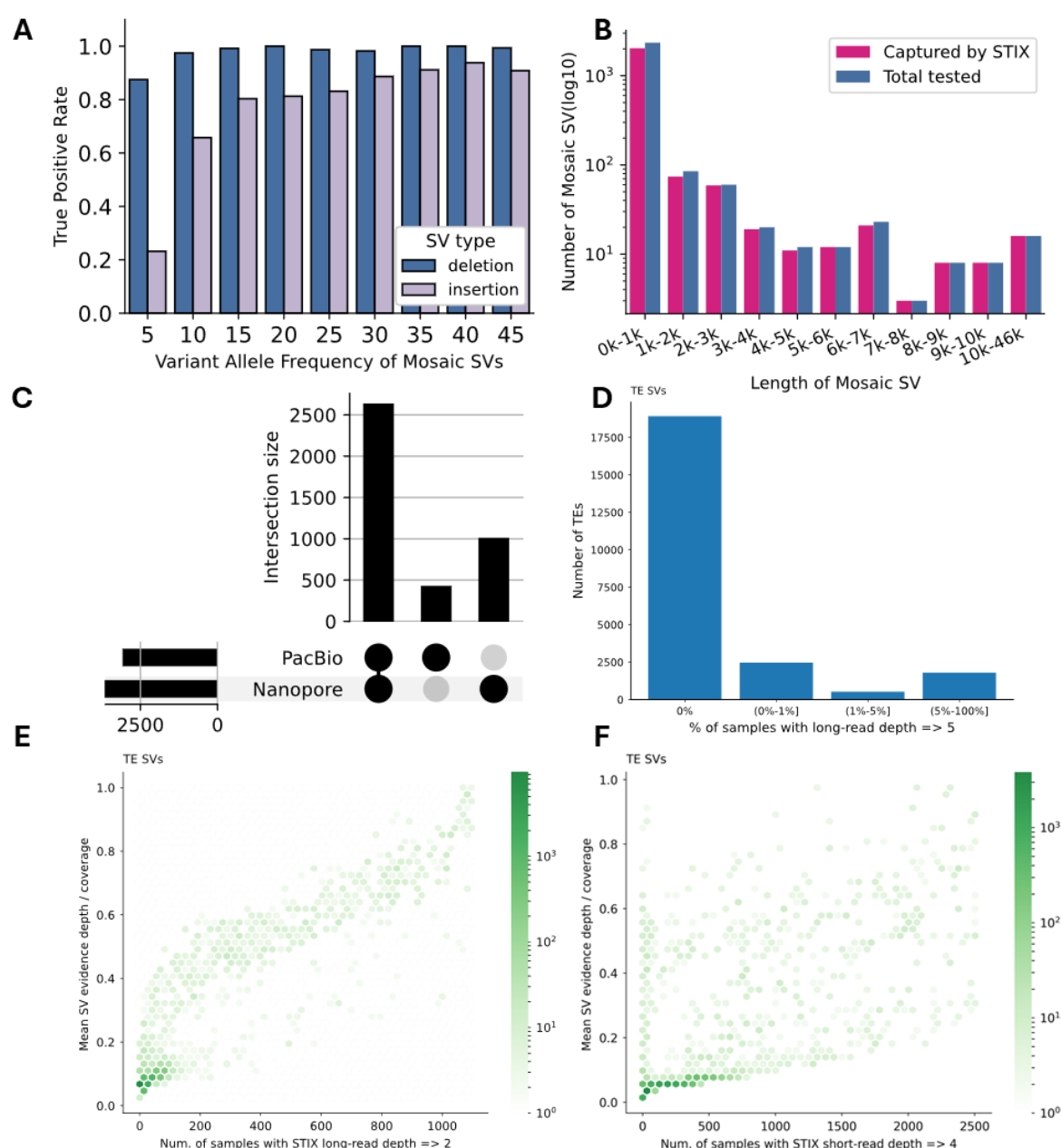Next investigated if a mosaic SV that was induced by a germline insertion identified in our recent work[41], was more common in the population than expected. In our previous work, we identified and validated a 6% VAF deletion that was due to an apparently common germline insertion of an ALU-Y element near an existing ALU-Y[41]. These two ALU-Y seemingly recombined in a certain number of cells and thus lead to a deletion

11

342  signal of ~6% VAF. Over PCR and sanger validation we could confirm these variants,
343  but also saw emerging bands in other brain samples that we tested. Thus, we
344  speculate that since the germline insertion is common in the population that the
345  recombinant mosaic allele might also be common in the population. We used STIX to
346  annotate both events reported and validated previously. STIX reported the ALU-Y
347  insertion in around 48.10% (533 samples) of the samples clearly confirming our
348  suspicion that this is a common germline insertion in the population. The resulting
349  recombinant between the two ALU elements was found at 0.2% proportion of evidence
350  in the population(2 samples). It's noteworthy that if we reduce the read threshold to 3,
351  given the overall coverage of 15-20x across the samples, we identified 23
352  samples(2.08% proportion of evidence). This indicates that while this is a
353  somatic/mosaic variant it is also commonly independently gained through the
354  population.

355  This observation spiked our interest in what other repeat recombinants might be
356  commonly shared in the population despite being independently gained throughout
357  the population. To further extend the search for common, but mosaic alleles in
358  individuals we next studied a recently reported set of mosaic transposable elements[54].
359  First, we focused on six samples where we had matched Pb-HiFi and ONT data
360  available in the STIX index. Thus, if we would observe the consistency of these alleles
361  across both technologies, it would strongly indicate that some of these postulated
362  mosaic SVs are shared in the population. We collected 18,331 mosaic SVs (3,551
363  deletion, 1,224 duplication and 13,348 inversion) from the CELL paper and removed
364  the duplicates as they were originally discovered at per-read level. We identified 9.28%
365  SVs being annotatable in at least one sample. Among them, 64.9% SV were identified
366  both in Pb-HiFi and ONT indicating the high-consistency across sequencing
367  technologies(**Fig. 6C**). Interestingly, there are 84.7% deletions and 69.0% duplications
368  having support by both Pb-HiFi and ONT, while only 6.32% inversions were observed
369  in both technologies. While there are many more inversions reported by the previous
370  study, we were curious on why there is such a discrepancy. Given this observation,
371  we investigated ONT and Pb-HiFi based indexes separately. Over manual inspection
372  (**Supplementary Fig. S8**) we speculate that these inversions are false positives in the
373  previous study caused by chimeras found in ONT data. We had reported something
374  similar for inversions before[41] but didn't expect to find this artifact in a published SV
375  call set. Interestingly, these inversions are typically smaller with a median size of 573
376  bp, which further matches our expectation of likely artificial chimeras at low frequency.
377  Thus, we concluded that many of these inversions are false positive SV calls in the
378  previous CELL publication[54], which further explains their high number of events (2.7
379  fold more inversions than deletions and insertions)

380  Next, we focus on mosaic TEs that are presented as deletions as they show high
381  consistency. Overall, 23,703 SVs(deletions) were included. Among them we can
382  identify 6,100 (25.7%) SVs in the healthy population and 2,324 (10.1%) are even
383  common(>1% population evidence). This suggests that mosaic SVs are more often

384    independently gained than we might expect. For those SVs that are shown positive
385    signals at population level, the majority of them exist at low VAF status while the others
386    tend to present in the population status (**Fig. 6D**). Despite those SVs with low VAF
387    and low proportion of evidence in the population(**left bottom part in the Fig. 6E**), we
388    observed a linear pattern showing a clear correlation between VAF and the proportion
389    of evidence in the population (**Fig. 6E**). This indicates that some of the postulated
390    mosaic repeat recombinants rose to germline SV. This observation is also supported
391    by the short-reads version of STIX (**Fig. 6F**).

392



393
394

395    ***Fig. 6 Investigation of mosaic SV at population level: A)*** *True positive rate of*
396    *mosaic SVs. SVs were spiked in with different fractions(x-axis) to simulate various*
397    *VAF. Each VAF has two bars representing deletions(blue) and insertions(yellow)*
398    *respectively.* ***B)*** *The impact of SV length. Mosaic SVs were classified using 1kb bin*

13

399 *size(x-axis). y axis represents the number of SVs with log 10 scale. The blue bars*
400 *represent SVs total assessed, red bars represent SVs captured by STIX. **C)** the*
401 *comparison of 6 samples sequenced by both Pb-HiFi and ONT. SVs were grouped by*
402 *their type(DEL,DUP, and INV) and the existence of two platforms. **D)** Proportion of*
403 *evidence distribution of mosaic TEs that was annotated by STIX. **E-F)** Distribution of*
404 *average variant allele frequency of a SV in samples(y-axis) versus number of samples*
405 *that harbor the SV(x-axis). Variant allele frequency of a SV in a sample was calculated*
406 *by the supporting reads / sample mean coverage. **E:** long-read index, **F:** short-read*
407 *index*

# Discussion

409 In this work, we present a long-reads based annotation approach by extending STIX
410 to long-reads datasets. This new version now also supports insertions and can utilize
411 either Pb-HiFi or ONT data. We validated the performance of STIX for long-reads with
412 the GIAB benchmark(HG002 sample) and assess its specificity using a negative
413 control sample(HG00733). We found that the overall true positive rate (TPR) is 81.82%
414 and false positive rate (FPR) is 3.44%. We observed higher performance of STIX and
415 other two state of the art SV genotypers on the GRCh37 version benchmark as it is
416 better established than the GRCh38 and CHM13-T2T benchmark. According to our
417 experiment, STIX shows robust performance with SVs with different length and
418 genotype. We further tested SITX in SVs that are located in challenging medical
419 relevant genes(CMRG SVs). We achieved 89.45% TPR in both sequencing platforms
420 (Pb-HiFi and ONT) and 99.08% TPR when considering at least one sequencing
421 platform. Furthermore, it is also important to highlight its future readiness as the index
422 is ready to be extended, can include meta information (e.g. ethnicities, diseases
423 background etc) and does not rely on a SV calling pipeline that ever needs to be
424 updated. Thus, STIX demonstrates a clear advantage when annotating SV with
425 population frequencies.

426 Variant annotation and thus prioritization are of utmost importance to identify potential
427 pathogenic variants for the medical and biological sciences[55,56]. The principle still
428 holds that common variants in the population for the most part are not pathogenic[19].
429 Therefore the identification of rare variants based on the population represents still
430 one of the best practices to rank or prioritize variants. Over the past decade multiple
431 advancements in SNV annotation have been made, which greatly improved medical
432 and clinical studies[57–60]. These advances have been mainly achieved by short-reads
433 as long-reads so far have been cost prohibitive to build up significant annotation
434 resources. Thus, despite the fact that long-reads improve SV detection, short reads
435 have been utilized to annotate SV[24]. This has led to the issue that many SVs are not
436 annotatable and thus the prioritization was not possible. Therefore the identification of
437 pathogenic SV has been hindered. This is demonstrable by gnomAD only being able
438 to annotate 33.5% of the HG002 GIAB SV catalog despite the sample likely being
439 present in gnomAD as the SNV are annotatable by more than 99.7%. To overcome
440 this, we have extended STIX to incorporate long-read information for SV annotation.

14

441   With this we are now able to annotate and thus prioritize SV. Despite the significantly
442   smaller indexed population size of 1,108 genomes, STIX was able to demonstrate
443   significant concordance with large short-read based catalogs for SV that were
444   annotatable in both data sets. This together with low false positives rate really
445   improves the SV annotation. Another important aspect about STIX is also the accuracy
446   it annotates variants avoiding the commonly used reciprocal overlap of 50% or 70%[24].
447   It is important to also discuss the behavior of STIX when annotating false positive SV
448   themselves. As highlighted with inversions, this typically results in high population
449   evidence annotations from STIX. This typically would discard these false positive SV
450   or artifacts from subsequent studies. Of note this is in contrast to the FPR of STIX,
451   which represents the annotation of a similar but different, close by SV. In this work we
452   have carefully benchmarked and concluded that actually a 100bp padding and min 5
453   reads are sufficient to accurately annotate SV and avoid FPR with close by SV.

454   We were able to demonstrate STIX ability and utility in cancer data sets and mosaic
455   SV call sets. For cancer we identified 3,563 SV that are reported in COSMIC[51], but are
456   actually common (>1% proportion of population evidence) in the population. While it
457   remains unclear if these are errors in COSMIC or maybe common genome instabilities,
458   it shows that STIX is an important method to prioritize cancer mutations or avoid
459   potential common SV. Thus, also speeding up the prioritization of SV for cancer
460   studies. Based on our experiment, STIX can narrow down the scope of potential
461   pathogenic SV to approximately 10% and retain 94% of true somatic SVs. This
462   demonstrates clearly that STIX is useful in the case of tumor only or even paired
463   sample analysis. To investigate the potential for common mosaic SV, we have
464   annotated a previous validated mosaic deletion, which was caused by an Alu-Y repeat
465   recombinant [41]. STIX was able to annotate both events and indeed showed that the
466   germline Alu-Y insertion at this location was highly common in the population (48.10%
467   of the individuals). In addition, the resulting mosaic deletion was still detected in two
468   samples (increases to 23 when reducing the thresholds) despite the general lower
469   coverage of the indexed data sets (15-25x). We expanded this study by investigating
470   somatic repeat recombinations previously published. Here STIX was able to find 10.1%
471   of the SVs are commonly shared(>1%) within different individuals of our index. Thus
472   highlighting that mosaic SV despite being rare in an individual can be independently
473   gained thought the population by likely common genome instabilities (ie. repeat
474   recombination sites). This result is significant as it demonstrates the plausibility of the
475   3,563 COSMIC somatic SV that are commonly shared not being false positive, but
476   rather benign variants. This is consistent with previous studies[52,53,61,62]. Thus, clearly
477   highlighting the importance of STIX to be used to annotate SV to streamline
478   pathogenic SV detection itself. This finding further highlights the need to study and
479   understand mosaic SV further as some fraction seems to be commonly shared in the
480   population and even arise to higher zygosity.

481   Lastly, we compared the new long-read version of STIX presented in this paper with
482   the short-read version, finding a high level of consistency between them. This

15

483 demonstrates STIX's robustness across both long-read and short-read data, paving
484 the way for joint SV annotation using long-read and short-read datasets
485 simultaneously. A gaussian mixture model was developed to provide the estimation of
486 the population allele frequency of query SVs based on the amount of evidence (the
487 number of alignments) that support the presence of a query SV in a sample(and
488 population of samples) that was returned by STIX. We tested this model with both
489 long-read and short-read indices and observed a high correlation with the original
490 allele frequency from 1KG. As expected, the correlation was higher with long-reads
491 (**Supplementary Fig. S17**) given its lower noise profile. Since the model was based
492 on the Hardy-Weinberg equilibrium, it struggles to classify SVs with allele frequencies
493 greater than 0.7, and most samples have two copies of the variant.

494 Overall, STIX has shown to be a versatile and accurate SV annotation methodology
495 that represents a significant improvement over other resources. Its indexing of raw
496 reads makes it updatable and future ready for improvements in SV identification over
497 the next few years.

# Methods

## Extending STIX to long-reads datasets

500 STIX extends the ability to extract, index, and search all types of SV signals from datasets
501 sequenced by ONT and Pb-HiFi platforms**(Figure 1A)**. To accomplish this, we developed a
502 new program named excord-lr, which is the long-read version of excord[26] to extract signals
503 from long-reads datasets. This program was written in the Rust programming language.
504 Precompiled binaries are available under releases in its GitHub repository. Briefly, the Excord-
505 lr program takes the BAM file as input and extracts SV signals from both the CIGAR value and
506 SA tags. It incorporates a variety of empirically-based filters to minimize false positive signals.
507 This includes 1) controlling the maximum number of supplementary alignments for a single
508 read, and 2) limiting the proportion of overlap between two supplementary alignments. A
509 function was also embedded in Excord-lr in order to extract insertions with different situations.
510 The output is compatible with short-version, which generates a file with BEDPE format.
511 Excord-lr also supports running with multiple threads to accelerate the speed. Debug
512 information and reads id of a SV signal are also available under the debug mode. Excord-lr is
513 specifically designed to manage the distinct features of long-read sequencing datasets and to
514 ensure they are compatible with further searching and indexing processes.
515
516 One of the major updates in STIX is supporting the extracting, indexing and searching of the
517 insertions with different lengths. To accomplish this, STIX uses excord to extract three types
518 of insertions encompassed in long-read sequencing datasets according to its length,mapping
519 status, and the fields present in the BAM file.
520     1. Shorter insertions will be encoded as "I" tags in CIGAR values.
521     2. Large insertion will be encoded as supplementary alignment in the SA tag. But the
522        insertion is shorter than the read length, but can not be presented in CIGAR value. It

523    may generate a primary alignment and a supplementary alignment. Each alignment
524    will have a soft-clip or hard-clip on the left and right respectively.

525  3.  When the large insertion is longer than read length, it will only result in a primary
526    alignment with one large soft-clip or hard-clip only.

527

528  Taking into account those conditions, excord-lr implies three different approaches to extract
529  insertion signals. For the insertions in CIGAR value, it will extract the position of insertion point
530  as well as the length of insertion. For insertions present in SA tag that generate two alignment,
531  excord-lr will extract them by using the criteria as follows:

532  1.  must have one primary alignment and one supplementary alignment.
533  2.  Each alignment must have a soft-clip or hard-clip(no less than 1kb by default).
534  3.  two alignments should overlap to each other.
535  4.  both alignment should have the same chromosome name.
536  5.  For insertions that only generate one primary alignment, excord-lr extract them by
537    including the following filters:
538    a.  only one primary alignment, no supplementary alignment.
539    b.  primary alignment must have soft-clip or hard clip more than 1kb.

540

541  Excord-lr encodes insertions as 0 base pair intervals in BEDPE format. The insert point is
542  presented at the end of the left region and the start of the right region. In other words, the end
543  position of the left region will equal the start position of the right region if the record encodes
544  an insertion. This 0 base pair interval will be used to STIX searching step to distinguish
545  insertions from other types of SVs. if the insertions are accompanied with length
546  information(insertions that are extracted from CIGAR value), the length will be encoded as the
547  length of the right region, otherwise, the start position is equal to the end position in the right
548  region.

549

550  STIX kept the same steps to index SV for signals from long-reads datasets with the short-
551  reads datasets, as described[26]. Briefly, Giggle, a fast genomics search engine , is used as a
552  dependency to index SV signals[36]. The output bed file from excord-lr will be sorted according
553  to the position of the left and right region and be compressed with gzip. For example:

```
cat HG002.bed \
 | LC_ALL=C sort --buffer-size 2G -k1,1 -k2,2n -k3,3n \
      | bgzip -c > ./data/HG002.sorted.bed.gz
```

554  Giggle index then can be built using command as following:

```
giggle index \
 -i "./data/*.bed.gz" \
 -o ./giggle_idx \
 -s -f
```

555  STIX index will be created based on top of sorted bed file, giggle index, and metadata.

```
stix \
 -i ./giggle_idx/ \
 -d stix_idx.db \
 -p HG002.meta.ped \
 -c 5
```

17

556 STIX requires a file that describes sample-related information, which can be specified with the
557 -p option. The -c option is used to indicate the column of the file names.
558

559 The search step in the current STIX version is mainly inherited from the old version but with a
560 few new options to enhance its ability. Concisely, STIX can take a single region pair or a VCF
561 file as input. It can also accept a tabular format as input in the latest version in case users find
562 it difficult to generate a VCF when they want to query more than one region pair.
563

564 When searching for an insertion. Users need to provide two base pair intervals for left region
565 and right region respectively as well as an additional parameter (-L) to specify the length of
566 the insertion. The steps for how STIX searches for an insertion are briefly outlined as follows:
567     1. searching the potential records that overlap with the query regions.
568     2. For each record, mark them as an insertion record if the end position of the left region
569        equals the start position at the right region and the chromosomes of both regions are
570        identical.
571     3. Compare the insertion point between the query and each result, report a hit if they are
572        identical, otherwise report a non-hit.
573     4. If the length of query insertion is provided by the user, STIX will try to further compare
574        the length between the query insertion and the potential insertions from the index. A
575        relative error will be calculated to describe the similarity of the length between the
576        query and we set 0.2 as default cutoff.
577     5. If the user does not provide the length information for query insertion. STIX will search
578        the potential insertions located in the padding region.
579

580 There are examples for running STIX in different search mode:
581     1. Single query:

```
stix \
        -i ./giggle_idx/ \
        -d ./stix_idx.db \
        -s 100 \
        -t DEL \
        -l 1:934054-934074 \
        -r 1:934895-934915
```

582     2. Batch query using VCF as input:

```
stix \
        -i ./giggle_idx/ \
        -d ./stix_idx.db \
        -s 100 \
        -f HG002.high_confidencce.vcf.gz
```

583     3. Batch query using table as input:

```
stix \
        -i ./giggle_idx/ \
        -d ./stix_idx.db \
        -s 100 \
        -Q HG002.high_confidencce.SV.tsv
```

584

18

585  Additionally, we introduced a sharding mode to facilitate indexing of large datasets. This has
586  a wrapper script for generating separate Giggle indexes and an -B option in STIX, which
587  accepts a two-column TSV file. Each row in the file specifies a Giggle index and the
588  corresponding STIX index.

```
stix \
      -B shards.tsv
      -s 100 \
      -Q HG002.high_confidencce.SV.tsv
```

589

# Data collection and preprocessing

591  The original links of all public datasets can be found at **Supplementary Table S1.** Briefly,the
592  SVs of HG002 in the Variant Call Format were downloaded from the Genome In A Bottle
593  project(GIAB). The accompanying high-confidence regions were also downloaded from the
594  same path. The GRCh37 version of the SV callset of HG00733 was collected from the Truvari
595  project[25]. We generated the GRCh38 version and CHM13-T2T version callset for HG00733
596  using an assembly-based pipeline which replicated the GIAB TR process for variant mapping[63].
597  We initially collected two bam files of HG002 which were sequenced by Pb-HiFi and ONT
598  respectively, and a bam file of HG00733 that sequenced by ONT. All bam files were aligned
599  to GRCh37 reference and were remapped to GRCh38 reference and CHM13-T2T reference
600  respectively using samtools and minimap2. The bam files for long-reads index were
601  downloaded from The 1000 Genomes Project and The Human Pangenome Reference
602  Consortium project. Raw reads were extracted from HRPC bam files and realigned to the
603  GRCh38 and the CHM13-T2T reference using minimap2. Callsets from COLO829 and
604  COLO829BL were collected from https://zenodo.org/records/10819636. CMRG SVs callset
605  were downloaded from the GIAB project and intersected with the high-confidence regions. IDs
606  for each variant were added using perl -pe 'if(!/^#/){$c++; $_=~s/\t\./\tGRCh38_CMRG_$c/}' to
607  make it meet the input requirement of SVAFotate(0.0.1). High-quality SVs from high-coverage
608  short-reads dataset in the 1KG project were collected from this paper[64]. The raw vcf was
609  downloaded from 1KG FTP site and retained samples that only overlap with STIX LR index
610  by using bcftools view -S sample_id.unique.txt --force-samples -o 1KGP.subset.vcf
611  1KGP_3202.gatksv_svtools_novelins.freeze_V3.wAF.vcf.gz . Mosaic transposable elements
612  (TEs) were collected as requested from the authors[54]. The raw format is tabular format and
613  was transformed into a 5-columns tabular format (left-region \t right-region \t length \t svtype
614  \t ID) using python. The projection of svtype is listed below: "del" --> "DEL", "dupl" --> "DUP",
615  "Inter" --> "BND", "inv" --> "INV". Raw fastq file of TP10Plus was downloaded from under
616  SRA(PRJNA636606) and remapped to GRCh38 using minimap2(2.26-r1175).

617

# Benchmarking for germline SV

619  Raw SV signals were extracted using excord-lr(version 0.1.17) and build index by using STIX
620  for HG002 BAM files for both Pb-HiFi and ONT platform respectively as described above and
621  the previous paper[26] and https://github.com/ryanlayer/stix. To collect high-quality HG002
622  SVs and ensure the consistent existence of those SVs in HG002 bam files, we only consider
623  the SVs located in the high-confidence regions provided by GIAB. SVs that were larger than

19

624    50bp were used for downstream analysis. HG00733 specific SVs were generated by excluding
625    the common SVs from HG002 callset by using SURVIVOR(1.0.7)[65] with parameters '1000 1
626    1 1 0 50'. This merges the VCF files with 1000bp padding enforcing the same type and strand
627    for variants larger than 50bp. Based on final VCFs for HG002 and HG00733, 5-columns tables
628    were generated for STIX one by one query mode. STIX took SVs from HG002 and HG00733
629    as input and query each SV on the HG002 index using command:

```
/workspace/stix/bin/stix \
        -i 03.1.giggle_idx/ \
        -d 03.1.stix_idx.db \
        -s ${padding} \
        -t ${svtype} \
        -l ${left_region} \
        -r ${right_region}
```

630
631    We use different padding: 5, 10, 50, 100, and 500 to test the performance and consider the
632    minimal supporting reads from 1 to 15. We consider the number of hits in HG002 SVs as true
633    positives and non-hits as false negatives, the number of hits in HG00733 as false positives
634    and non-hits as true negatives. Other statistics can be calculated using the following equations:
635

$$TPR = \frac{TP}{TP + FN}$$
636

$$FPR = \frac{FP}{FP + TN}$$
637

$$Recall = \frac{TP}{TP + FN}$$
638

$$Precision = \frac{TP}{TP + FP}$$
639

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$
640

641
642
643    Run Sniffles2[41] (2.2) genotyping mode:
644

```
sniffles --threads 8 --allow-overwrite --input <bam> --genotype-vcf
<input.vcf> --vcf <output.vcf>
```

645
646
647    Run cutesv[42](2.2.1) genotyping mode:
648

```
cuteSV <bam> <ref.fa> <output.vcf> ./ --max_cluster_bias_INS 1000 --
diff_ratio_merging_INS 0.9 --max_cluster_bias_DEL 1000 --
diff_ratio_merging_DEL 0.5 -Ivcf <input.vcf> -q 10 -L -1
```

20

## Benchmarking for mosaic SV

We have developed a framework for generating in-silico mosaic structural variations (SVs) with varying variant allele frequencies, in addition to creating corresponding callsets. Briefly, mixed BAM files that are used in the mosaic benchmark were generated by injecting raw reads from HG002 into HG00733 bam. We used samtools to subset the HG002 bam with the -f parameter to specify the percentage of downsampling, and merged the subset BAM with the HG00733 bam. Finally, the merged bam files were sorted and indexed. Sequencing depths of the BAM files are estimated using mosdepth[66]. SV callsets from HG002, which are used in the germline benchmark, were re-genotyped in the mixed BAM files using Sniffles (version 2.2) in force-call mode with the parameters --input mixture.resort.bam --genotype-vcf HG002.vcf --vcf mixture_regenotype.vcf --allow-overwrite. Only the SVs that were reported with supporting reads were retained for downstream analysis. Variant allele fraction(VAF) were calculated based on the AD and DP in the output VCF file. Raw SV signals were extracted using excord-lr and the STIX index was built based on top of the sorted bed files with command listed in the Methods .

## Build large index on long-reads datasets

We gathered publicly available long-read sequencing datasets from healthy individuals. In summary, we obtained 100 datasets generated using the ONT platform from the 1000 Genome Project. These datasets were already aligned to the GRCh38/CHM13-T2T reference genome using minimap2 [67]. Additionally, we acquired 100 datasets (in unaligned BAM format) generated using the Pb-HiFi platform from The Human Pangenome Reference Consortium project. We aligned those dataset to GRCh38 and CHM13-T2T reference with the following command: 'minimap2 -ax map-hifi <GRCh38.fa> <sample.fq.gz> | samtools view -b > aln.bam'. We downloaded 908 datasets from the 1000 Genome Project (Vienna project) with GRCh38 and CHM13-T2T reference individually. The depth of coverage for each dataset was calculated using mosdepth[66] with the following parameters: '-b 1000 -x -t 8 -Q 20 --no-per-base. The ethnicities of those individuals were annotated by manual curation. The raw SV signals were extracted by excord-lr and the sharded giggle indices and STIX indices were created as described above.

## Annotate SVs based on short-reads datasets

We use SVAFotate[24] to annotate SVs that are used in this study. SVAfotate software (version 0.0.1) was installed according to their instruction[24]. The corresponding libraries were downloaded from the same repository. We downloaded the SVs from Challenge Medical Relevant Genes from GIAB with the filename HG002_GRCh38_difficult_medical_gene_SV_benchmark_v0.01.vcf.gz. We selected the SVs located in high-quality regions by bedtools(v2.30.0) using command:

```
bedtools intersect -u -header -a
HG002_GRCh38_difficult_medical_gene_SV_benchmark_v0.01.vcf.gz -b
HG002_GRCh38_difficult_medical_gene_SV_benchmark_v0.01.bed | bgzip >
HG002_GRCh38_difficult_medical_gene_SV_benchmark_v0.01_trusted.vcf.gz
```

687  Unique IDs were manually added for each SVs since the SVAFotate needed them to make
688  the right output. The annotate command used in this study is listed below:

689

```
svafotate annotate -s gnomAD -f 0.8 -c 0 -a best -v input.vcf -o output.vcf -b
SVAFotate_core_SV_popAFs.GRCh38.bed.gz
```

690

# Analysis of Mosaic Repeat Elements

692

693  Raw data of TP10plus were downloaded from the SRA database and converted to fastq format
694  using SRAtools. We next align the raw data against the GRCh38 reference using minimap2
695  with parameter " -t 8 -ax map-ont -Y". Then use samtools to convert the SAM file to BAM
696  format. Excord-lr was used to extract SV signals from aligned bam and STIX index was
697  created by using the aligned BAM file. Mosaic repeat elements locus were downloaded from
698  the supplementary materials of previous published work[54]. We further convert the format for
699  STIX searching. In summary, for each original file, we extract the left and right breakpoints.
700  The duplicated records after rearrangement of left and right breakpoint were removed. We
701  projected the SV types from .out file to STIX accepted format by the following rules: del →
702  DEL; dupl → DUP; Inter → BND; inv → INV. a 5-column table was generated for each .out file
703  and further used for STIX annotation. SVs that were not located in the GIAB high-confidence
704  regions were excluded to avoid the noise.

705

# Hardy-Weinberg mixture model for estimated an SVs alternate allele frequency

708

709  STIX returns the amount of evidence (the number of alignments) that support the presence of
710  a query SV in a sample. While *evidence depth* is useful for reasoning about an SV's role in a
711  trait, it is not a standard population genetic metric and its unfamiliarity can limit its usefulness.
712  The ideal population genetic metric for estimating a variant's impact is population allele
713  frequency, or the proportion of haplotypes in the population that harbor the variant. The
714  primary difference between *evidence depth* and *allele frequency* is that the former is a binary
715  reporter–does a sample have evidence or not–while the latter sorts samples with evidence
716  into heterozygous or homozygous states. To quantify these additional states we leverage two
717  observations. First, is that when STIX is performed across a population, the distribution of
718  evidence depths for many SVs has 3 rough modes that correspond to the samples that have
719  zero, one, and two copies of the variant. Second is that past analyses indicate that 86% of
720  SVs are in Hardy-Weinberg equilibrium[22].

721

722  From these observations we developed STAFFR (structural variant allele frequency finder,
723  https://github.com/behzodcu/staffr), a probabilistic generative mixture model that embodies
724  the Hardy-Weinberg equilibrium ($p^2 + 2pq + q^2 = 1$), which can be used to estimate the
725  alternate allele frequency $q$ from STIX evidence. Our model interprets variabilities in
726  sequencing data, attributing different levels of evidence depths to specific genotypes: low or

22

727 very low depths correspond to a homozygous reference genotype where the structural variant
728 is absent, medium depths correspond to heterozygosity with one copy of the variant present,
729 and high depths correspond to a homozygous recessive genotype with two copies of the
730 variant present. This probabilistic generative mixture model is based on modeling variation in
731 evidence using Normal or Gaussian distributions, which have a probability density function of:

732
$$N(x|\mu,\sigma^2) \;=\; \frac{1}{\sigma\sqrt{*2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}} \, .$$

733 Here $\mu$ and $\sigma^2$ denote the sample mean and variance of the variation in evidence across
734 individuals. In general, such mixture models are defined by a collection of several Gaussian
735 distributions, each with a distinct mean and variance, which jointly represent the likelihood of
736 the observed data. The likelihood of each particular observation depends on its total probability
737 across each distribution or mode in the model, which has the form:

738
$$Pr(x_j|\Theta) = \sum_{i=1}^{k} \pi_i N(x_j \mid \mu_i, \sigma_i^2).$$

739 Here $\Theta = \{\pi, \mu, \sigma\}$ where $\pi_i$, $\mu_i$, and $\sigma_i^2$ are the probability of generating a value from the $i$th
740 mode, that mode's mean and its variance, respectively.

741

742 Integrating the Hardy-Weinberg equilibrium into the framework of a mixture model allows for
743 the incorporation of genetic principles directly into the model's estimation. Specifically, Hardy-
744 Weinberg equilibrium defines a mathematical relationship that governs the proportions of
745 genotypes in a population based on allele frequencies and their relative levels of evidence.
746 According to the equilibrium, the frequencies of homozygous reference, heterozygous, and
747 homozygous alternate genotypes occur as $p^2$, $2pq$, and $q^2$ respectively, with $p$ and $q$
748 representing the frequencies of the two alleles. By constraining the model's estimation to
749 conform to these expected genotype frequencies, a modified mixture model embodies the
750 foundational concepts of population genetics, ensuring that the estimated genotype
751 frequencies are consistent. In the above mixture model, we formalize this notion by associating
752 each genotype with the proportion of evidence depths there are, i.e., with their weight $\pi =$
753 $\{p^2, 2pq, q^2\}$.

754

755 Furthermore, a structural variant in Hardy-Weinberg equilibrium implies a constraint on the
756 model's estimated evidence depths based on the presence of one or two copies of the
757 structural variant. As $\mu$ represents the mean evidence depth associated with a single copy of
758 the variant under the model, the implied mean for the mode representing two copies of the
759 variant will be twice the mean of one copy, $\mu_3 = 2\mu_2$. And, considering that the reference
760 population exhibits no evidence of structural variant presence, we replace the first mode with
761 a simple proportion $\pi_1 = \alpha$ of data showing no variant, i.e. the proportion of zeros.

762

763 Putting these together, the total the likelihood of the observed data being generated from any
764 of the modes with these constraints, is given as:

765
$$L(X|\Theta) \;=\; \prod_{j=1}^{n} \left[ \alpha + \pi_2 N(x_J|\mu_2, \sigma_2^2) + \pi_3 N(x_J|\mu_3, \sigma_3^2) \right]$$

766 with log-likelihood

767
$$l(X|\Psi) \;=\; \sum_{j=1}^{n} ln \left[ \alpha + \pi_2 N(x_J|\mu_2, \sigma_2^2) + \pi_3 N(x_J|\mu_3, \sigma_3^2) \right] \, .$$

23

768  We identify divergences by the observed data from the Hardy-Weinberg equilibrium using a
769  null model approach[68]. In our null (non-HW) model, $q = 0$ and the noise from sampling the
770  reference genome is modeled as $\Psi = \{\alpha, \lambda\}$, where $\alpha$ is the proportion of zero evidence data
771  points (as in the HW model) and non-zero evidence values are modeled by a geometric tail
772  with an expected value $\lambda$:

$$Pr(x_j|\Psi) = \begin{cases} \alpha & \text{if } x_j = 0 \\ \frac{(1-\alpha)}{\lambda}\left(1 - \frac{1}{\lambda}\right)^{x_j-1} & \text{if } x_j > 0 \end{cases}$$

773  In the null (non-HW) model, the likelihood of the observed data is given as:

774
$$L(X|\Psi) = \prod_{j-1}^{n} \left[I(x_i = 0)(\alpha) + I(X_j > 0)\left(\frac{1-\alpha}{\lambda}\left(1 - \frac{1}{\lambda}\right)^{x_j-1}\right)\right],$$

775  where we use indicator variables $I$ to split the values of $x_j$ between zero and non-zero. The
776  log-likelihood is then

777
$$l(X|\Psi) = \sum_{j=1}^{n} ln\left[I(x_i = 0)(\alpha) + I(X_j > 0)\left(\frac{1-\alpha}{\lambda}\left(1 - \frac{1}{\lambda}\right)^{x_j-1}\right)\right].$$

778
779  The free parameters $\theta$ and $\Psi$ in the Hardy-Weinberg mixture model and null (non-HW) model
780  provide quantitative descriptions of the distribution of genotypes and evidence depths in a
781  given population from the perspective of a Hardy-Weinberg equilibrium vs. the alternative non-
782  HW scenario. We can estimate these parameters directly from a set of observed evidence
783  depths using a standard expectation-maximization (EM) algorithm[69]. The resulting estimates
784  of model parameters can be interpreted in ways consistent with concepts from population
785  genetics: allele frequencies ($\pi$ and $\alpha$), the positioning of these genotypes in our evidence
786  depth data ($\mu_2$), the width of these genotypes in our evidence depth data ($\sigma_2$ and $\sigma_3$), and our
787  noise measurement in our null model ($\lambda$).

788
789  To run the EM algorithm, we first derive maximum likelihood estimators (MLEs) for each model
790  parameter. In the HW model, we use standard MLEs for the Gaussian parameters, noting that
791  the location parameters of the 2nd and 3rd modes are coupled. In the non-HW (null) model,
792  the MLEs are

793
$$\hat{a} = \frac{\sum_{i=0}^{n} I(x_i=0)r_i}{\sum_{i=0}^{n} r_i},$$

794
$$\hat{\lambda} = \frac{\sum_{i=0}^{n} I(x_i \neq 0)x_i r_i}{(1-\hat{a})\sum_{i=0}^{n} r_i}.$$

795
796  In our analysis, we first determine whether the data fits more closely with the proposed Hardy-
797  Weinberg model or the non-Hardy-Weinberg (null) model using a numerical Kolmogorov-
798  Smirnov (KS) test to evaluate the null hypothesis that the observed data are plausibly
799  generated from a population with $q = 0$. We construct the null distribution of the test statistic
800  numerically in the following manner:
801     1.  Fit the null model to the dataset, $X$, from which we estimate the parameters $\hat{\alpha}$ and $\hat{\lambda}$.

24

2.  Given these estimated parameters, we compute the empirical distribution function (EDF) for $X$ and the cumulative distribution function (CDF), evaluated at each unique value in $X$, using the estimated model parameters $\hat{\alpha}$ and $\hat{\lambda}$.

3.  Calculate the empirical test statistic $D_*$ (KS distance) between the EDF and CDF.

4.  Then construct the null distribution of $D$. Generate >1000 datasets $\Psi_i$ from the null model with estimated parameters $\hat{\alpha}$ and $\hat{\lambda}$; fit the null model to each generated dataset; for each generated dataset, we calculate $D$ relative to its own fitted model.

5.  Calculate the p-value as the fraction of generated datasets with a test statistic $D$ that are at least as large as $D_*$.

Staffr can also characterize the degree of overlap of the evidence depths for the relevant genotypes under the HW equilibrium model. This test uses the *common language effect size* or *Mann-Whitney U* statistic (which ranges from 0.5 to 1.0) to quantify the separability of modes in the evidence depth data. The steps of the process are as follows:

1.  For each point, assign it to the mode (genotype) with the greatest responsibility $r_{ij}$ .

2.  For a large number of rounds, choose two points uniformly at random; let $y$ denote the value from the heterozygous mode and let $z$ denote the value from the homozygous alternate mode.

3.  Let $U$ denote the fraction of times that $y > z$.

If $U = 0.5$, then modes are statistically indistinguishable, and if $U = 1.0$ then the modes are completely statistically distinguishable, with intermediate values indicating some degree of overlap.

Staffr can also identify outliers relative to a non-HW distribution of evidence, which can be interpreted as evidence for non-reference genotypes. Its method for detecting such outliers is as follows:

1.  Over the non-zero evidence values in the data, the MLE of the non-HW model's geometric distribution parameter $\lambda$ is given by these data's mean value.

2.  Parameterize the geometric probability mass function

$$f(x|\lambda) = \frac{(1-\alpha)}{\lambda}\left(1 - \frac{1}{\lambda}\right)^{x-1},$$

and evaluate $f(x|\lambda)$ for each unique evidence value $x$.

3.  Let $b$ denote a user-selected minimum number of observations allowed for non-outliers. In practice, $b = 1$ is a relatively inclusive threshold, while $b = 0.1$ is more conservative.

4.  Declare all evidence values for which $f(x_i|\lambda) - b/n < 0$ to be outliers, i.e., values that are statistically unlikely to be observed under the null model.

# Statistics and Visualizations

We conducted all analysis and visualization using python with other packages: Scipy, pandas, numpy, matplotlib, and seaborn. We visualized the SV at read level using the IGV and samplot. Upset package(https://github.com/jnothman/UpSetPlot)was used to visualize the intersection of the mosaic TEs.

# Code Availability

846 exocrd-lr source code can be found at https://github.com/zhengxinchang/excord-lr . STIX
847 with a long-reads data compatible version can be found at https://github.com/ryanlayer/stix.

# Data Availability

849 GIAB HG002 PacBio HiFi data(27x) and Oxford Nanopore data(43x) were collected
850 at:https://ftp-
851 trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/HG002_NA24385_son/
852 PacBio_CCS_15kb/ and https://ftp-
853 trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/HG002_NA24385_son/
854 UCSC_Ultralong_OxfordNanopore_Promethion/. GIAB HG002 callset were collected from
855 https://ftp.ncbi.nlm.nih.gov/ReferenceSamples/giab/release/AshkenazimTrio/HG002_NA24385_son/N
856 IST_SV_v0.6/ and https://ftp-
857 trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/analysis/NIST_HG002_DraftBen
858 chmark_defrabbV0.012-20231107/
859 HG00733 callsets with three reference versions are hosted at
860 https://doi.org/10.5281/zenodo.13702318. CMRG bed was downloaded at https://ftp-
861 trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/analysis/NIST_HG002_medical_
862 genes_SV_benchmark_v0.01/. BAM and CRAM files from The 1000 Genomes Project were
863 downloaded from https://s3.amazonaws.com/1000g-
864 ont/index.html?prefix=FIRST_100_FREEZE/minimap2_2.24_alignment_data/,
865 https://s3.amazonaws.com/1000g-
866 ont/index.html?prefix=FIRST_100_FREEZE/minimap2_2.24_chm13_alignment_data/,
867 https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1KG_ONT_VIENNA/hg38/, and
868 https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1KG_ONT_VIENNA/t2t/. BAM files from
869 HPRC project were downloaded from https://human-
870 pangenomics.s3.amazonaws.com/index.html?prefix=working/HPRC/. The COLO829 callset was
871 collected from(https://zenodo.org/records/10819636.). The COSMIC SV were downloaded
872 from https://cancer.sanger.ac.uk/cosmic. The Mosaic TE callset were obtained with request
873 form authors. The raw reads of TP10plus were downloaded from the SRA
874 database(PRJNA636606). Details of the data link can also be found at **Supplementary**
875 **Table S7**.
876
877 All software used (with versions) is listed in **Supplementary Table S6.**

# Acknowledgments

26

887   https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1KG_ONT_VIENNA/README_1
888   KG_ONT_VIENNA_datareuse_statement_20240227.md. To comply with these restrictions,
889   all genome-wide population data and related datasets from this study have been deposited in
890   the original 1000 Genomes ONT archive at IGSR.

# Fundings

# Competing interests

F.J.S. receives research support from Illumina, Pacbio and Oxford Nanopore. All other authors declare no competing interests.

# Contributions

X.Z. developed the STIX LR and performed the analysis. M.C. performed the STIX SR analysis. B.M. & A.C. implemented the mixture model. R.M.L & F.J.S. supervised. All authors contributed to writing and reviewing the manuscript.

# Reference

1. Mahmoud, M. *et al.* Structural variant calling: the long and the short of it. *Genome Biol.* **20**, 246 (2019).

2. Mahmoud, M. *et al.* Utility of long-read sequencing for All of Us. *Nat. Commun.* **15**, 1–13 (2024).

3. Sanchis-Juan, A. *et al.* Complex structural variants in Mendelian disorders: identification and breakpoint resolution using short- and long-read genome sequencing. *Genome Med.* **10**, 95 (2018).

4. Carvalho, C. M. B. & Lupski, J. R. Mechanisms underlying structural variant formation in genomic disorders. *Nat. Rev. Genet.* **17**, 224–238 (2016).

5. Sekar, S. *et al.* Complex mosaic structural variations in human fetal brains. *Genome Res.* **30**, 1695–1704 (2020).

6. Papaemmanuil, E. *et al.* Genomic Classification and Prognosis in Acute Myeloid

27

915      Leukemia. *N. Engl. J. Med.* **374**, 2209–2221 (2016).

916   7.  Cancer Genome Atlas Research Network. Comprehensive genomic characterization of

917      squamous cell lung cancers. *Nature* **489**, 519–525 (2012).

918   8.  Gao, R. *et al.* Punctuated copy number evolution and clonal stasis in triple-negative

919      breast cancer. *Nat. Genet.* **48**, 1119–1130 (2016).

920   9.  Cosenza, M. R., Rodriguez-Martin, B. & Korbel, J. O. Structural Variation in Cancer:

921      Role, Prevalence, and Mechanisms. *Annu. Rev. Genomics Hum. Genet.* **23**, 123–152

922      (2022).

923  10. Dubois, F., Sidiropoulos, N., Weischenfeldt, J. & Beroukhim, R. Publisher Correction:

924      Structural variations in cancer and the 3D genome. *Nat. Rev. Cancer* (2024)

925      doi:10.1038/s41568-024-00738-y.

926  11. Jamal-Hanjani, M. *et al.* Tracking the Evolution of Non-Small-Cell Lung Cancer. *N. Engl.*

927      *J. Med.* **376**, 2109–2121 (2017).

928  12. Logsdon, G. A., Vollger, M. R. & Eichler, E. E. Long-read human genome sequencing

929      and its applications. *Nat. Rev. Genet.* **21**, 597–614 (2020).

930  13. Miller, D. E. *et al.* Targeted long-read sequencing identifies missing disease-causing

931      variation. *Am. J. Hum. Genet.* **108**, 1436–1449 (2021).

932  14. Wenger, A. M. *et al.* Accurate circular consensus long-read sequencing improves

933      variant detection and assembly of a human genome. *Nat. Biotechnol.* **37**, 1155–1162

934      (2019).

935  15. Ebert, P. *et al.* Haplotype-resolved diverse human genomes and integrated analysis of

936      structural variation. *Science* **372**, (2021).

937  16. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).

938  17. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants

939      from high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164 (2010).

940  18. Cingolani, P. *et al.* A program for annotating and predicting the effects of single

941      nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster

942      strain w1118; iso-2; iso-3. *Fly* **6**, 80–92 (2012).

943    19. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**,

944        285–291 (2016).

945    20. Lappalainen, I. *et al.* DbVar and DGVa: public archives for genomic structural variation.

946        *Nucleic Acids Res.* **41**, D936–41 (2013).

947    21. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in

948        141,456 humans. *Nature* **581**, 434–443 (2020).

949    22. Collins, R. L. *et al.* A structural variation reference for medical and population genetics.

950        *Nature* **581**, 444–451 (2020).

951    23. Danis, D. *et al.* SvAnna: efficient and accurate pathogenicity prediction of coding and

952        regulatory structural variants in long-read genome sequencing. *Genome Med.* **14**, 1–13

953        (2022).

954    24. Nicholas, T. J., Cormier, M. J. & Quinlan, A. R. Annotation of structural variants with

955        reported allele frequencies and related metrics from multiple datasets using SVAFotate.

956        *BMC Bioinformatics* **23**, 490 (2022).

957    25. English, A. C., Menon, V. K., Gibbs, R. A., Metcalf, G. A. & Sedlazeck, F. J. Truvari:

958        refined structural variant comparison preserves allelic diversity. *Genome Biol.* **23**, 271

959        (2022).

960    26. Chowdhury, M., Pedersen, B. S., Sedlazeck, F. J., Quinlan, A. R. & Layer, R. M.

961        Searching thousands of genomes to classify somatic and novel structural variants using

962        STIX. *Nat. Methods* **19**, 445–448 (2022).

963    27. Foord, C. *et al.* The variables on RNA molecules: concert or cacophony? Answers in

964        long-read sequencing. *Nat. Methods* **20**, 20–24 (2023).

965    28. English, A. C. *et al.* Analysis and benchmarking of small and large genomic variants

966        across tandem repeats. *Nat. Biotechnol.* (2024) doi:10.1038/s41587-024-02225-z.

967    29. Beyter, D. *et al.* Long-read sequencing of 3,622 Icelanders provides insight into the role

968        of structural variants in human diseases and other traits. *Nat. Genet.* **53**, 779–786

969        (2021).

970    30. Gong, J. *et al.* Long-read sequencing of 945 Han individuals identifies novel structural

29

971  variants associated with phenotypic diversity and disease susceptibility. *medRxiv*

972  2024.03.21.24304654 (2024) doi:10.1101/2024.03.21.24304654.

973 31. De Coster, W., Weissensteiner, M. H. & Sedlazeck, F. J. Towards population-scale

974  long-read sequencing. *Nat. Rev. Genet.* **22**, 572–587 (2021).

975 32. Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44–53 (2022).

976 33. Eilbeck, K., Quinlan, A. & Yandell, M. Settling the score: variant prioritization and

977  Mendelian disease. *Nat. Rev. Genet.* **18**, 599–612 (2017).

978 34. Han, L. *et al.* Functional annotation of rare structural variation in the human brain. *Nat.*

979  *Commun.* **11**, 1–13 (2020).

980 35. Demidov, G. *et al.* Structural variant calling and clinical interpretation in 6224 unsolved

981  rare disease exomes. *Eur. J. Hum. Genet.* 1–7 (2024).

982 36. Layer, R. M. *et al.* GIGGLE: a search engine for large-scale integrated genome

983  analysis. *Nat. Methods* **15**, 123–126 (2018).

984 37. Zook, J. M. *et al.* A robust benchmark for detection of germline large deletions and

985  insertions. *Nat. Biotechnol.* **38**, 1347–1355 (2020).

986 38. Rang, F. J., Kloosterman, W. P. & de Ridder, J. From squiggle to basepair:

987  computational approaches for improving nanopore sequencing read accuracy. *Genome*

988  *Biol.* **19**, 90 (2018).

989 39. Wang, Y., Zhao, Y., Bollas, A., Wang, Y. & Au, K. F. Nanopore sequencing technology,

990  bioinformatics and applications. *Nat. Biotechnol.* **39**, 1348–1365 (2021).

991 40. Wagner, J. *et al.* Curated variation benchmarks for challenging medically relevant

992  autosomal genes. *Nat. Biotechnol.* **40**, 672–680 (2022).

993 41. Smolka, M. *et al.* Detection of mosaic and population-level structural variants with

994  Sniffles2. *Nat. Biotechnol.* (2024) doi:10.1038/s41587-023-02024-y.

995 42. Jiang, T. *et al.* Regenotyping structural variants through an accurate force-calling

996  method. *bioRxiv* (2022) doi:10.1101/2022.08.29.505534.

997 43. Fairley, S., Lowy-Gallego, E., Perry, E. & Flicek, P. The International Genome Sample

998  Resource (IGSR) collection of open human genomic variation resources. *Nucleic Acids*

999     *Res.* **48**, D941–D947 (2020).

1000    44. Liao, W.-W. *et al.* A draft human pangenome reference. *Nature* **617**, 312–324 (2023).

1001    45. Gustafson, J. A. *et al.* Nanopore sequencing of 1000 Genomes Project samples to build

1002        a comprehensive catalog of human genetic variation. *medRxiv* (2024)

1003        doi:10.1101/2024.03.05.24303792.

1004    46. Wang, T. *et al.* The Human Pangenome Project: a global resource to map genomic

1005        diversity. *Nature* **604**, 437–446 (2022).

1006    47. gnomAD. https://gnomad.broadinstitute.org/help/sv-overview.

1007    48. van Belzen, I. A. E. M., Schönhuth, A., Kemmeren, P. & Hehir-Kwa, J. Y. Structural

1008        variant detection in cancer genomes: computational challenges and perspectives for

1009        precision oncology. *NPJ Precis Oncol* **5**, 15 (2021).

1010    49. Sun, R., Hu, Z. & Curtis, C. Big Bang Tumor Growth and Clonal Evolution. *Cold Spring*

1011        *Harb. Perspect. Med.* **8**, (2018).

1012    50. Paulin, L. F. *et al.* The benefit of a complete reference genome for cancer structural

1013        variant analysis. *medRxiv* (2024) doi:10.1101/2024.03.15.24304369.

1014    51. Sondka, Z. *et al.* COSMIC: a curated database of somatic variants and clinical data for

1015        cancer. *Nucleic Acids Res.* **52**, D1210–D1217 (2023).

1016    52. Jaiswal, S. & Ebert, B. L. Clonal hematopoiesis in human aging and disease. *Science*

1017        **366**, (2019).

1018    53. Lee-Six, H. *et al.* The landscape of somatic mutation in normal colorectal epithelial cells.

1019        *Nature* **574**, (2019).

1020    54. Pascarella, G. *et al.* Recombination of repeat elements generates somatic complexity in

1021        human genomes. *Cell* **185**, 3025–3040.e6 (2022).

1022    55. Weischenfeldt, J., Symmons, O., Spitz, F. & Korbel, J. O. Phenotypic impact of genomic

1023        structural variation: insights from and for human disease. *Nat. Rev. Genet.* **14**, 125–138

1024        (2013).

1025    56. Thibodeau, M. L. *et al.* Improved structural variant interpretation for hereditary cancer

1026        susceptibility using long-read sequencing. *Genet. Med.* **22**, 1892–1897 (2020).

1027   57.  Sherry, S. T., Ward, M. & Sirotkin, K. dbSNP-database for single nucleotide

1028        polymorphisms and other classes of minor genetic variation. *Genome Res.* **9**, 677–679

1029        (1999).

1030   58.  Chen, S. *et al.* A genomic mutational constraint map using variation in 76,156 human

1031        genomes. *Nature* **625**, 92–100 (2024).

1032   59.  UK10K Consortium *et al.* The UK10K project identifies rare variants in health and

1033        disease. *Nature* **526**, 82–90 (2015).

1034   60.  Cao, Y. *et al.* The ChinaMAP analytics of deep whole genome sequences in 10,588

1035        individuals. *Cell Res.* **30**, 717–731 (2020).

1036   61.  Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive selection of

1037        somatic mutations in normal human skin. *Science* **348**, 880–886 (2015).

1038   62.  Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age.

1039        *Science* **362**, 911–917 (2018).

1040   63.  English, A. *et al.* Benchmarking of small and large variants across tandem repeats.

1041        *bioRxiv* (2023) doi:10.1101/2023.10.29.564632.

1042   64.  Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the expanded

1043        1000 Genomes Project cohort including 602 trios. *Cell* **185**, 3426–3440.e19 (2022).

1044   65.  Jeffares, D. C. *et al.* Transient structural variations have strong effects on quantitative

1045        traits and reproductive isolation in fission yeast. *Nat. Commun.* **8**, 14061 (2017).

1046   66.  Pedersen, B. S. & Quinlan, A. R. Mosdepth: quick coverage calculation for genomes

1047        and exomes. *Bioinformatics* **34**, 867–868 (2018).

1048   67.  Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–

1049        3100 (2018).

1050   68.  Feng, C. X. A comparison of zero-inflated and hurdle models for modeling zero-inflated

1051        count data. *J Stat Distrib Appl* **8**, 8 (2021).

1052   69.  Moon, T. K. The expectation-maximization algorithm. *IEEE Signal Process. Mag.* **13**,

1053        47–60 (1996).

1054   70.  Olson, N. D. precisionFDA Truth Challenge V2: Calling variants from short- and long-

1055    reads in difficult-to-map regions. National Institute of Standards and Technology

1056    https://doi.org/10.18434/MDS2-2336 (2020)

1057