

Supplementary Information

A Catalogue of Structural Variation across Ancestrally Diverse Asian Genomes

Joanna Hui Juan Tan^{1,*}, Zihui Li^{1,*}, Mar Gonzalez Porta^{1,2,*}, Ramesh Rajaby^{1,3,*}, Weng Khong Lim^{1,4,5,6}, Ye An Tan⁷, Rodrigo Toro Jimenez¹, Renyi Teo¹, Maxime Hebrard¹, Jack Ling Ow¹, Shimin Ang¹, Justin Jeyakani¹, Yap Seng Chong^{8,9}, Tock Han Lim¹⁰, Lih Ling Goh¹¹, Yih Chung Tham^{12,13}, Khai Pang Leong¹¹, Calvin Woon Loong Chin^{14,15}, SG10K_Health Consortium[^], Sonia Davila^{6,16,17,18}, Neerja Karnani^{19,20,21}, Ching-Yu Cheng^{12,22}, John Chambers^{23,24,25}, E Shyong Taj^{4,26,27,25}, Jianjun Liu^{28,29}, Xueling Sim⁷, Wing Kin Sung^{1,30,31}, Shyam Prabhakar^{1,32,33,#}, Patrick Tan^{1,4,5,25,#}, Nicolas Bertin^{1,#}

Table of Contents

| | |
|--|-----------|
| Supplementary Notes..... | 2 |
| Supplementary Note 1: Benchmarking of SV calling pipeline for deletions and insertions | 2 |
| Supplementary Note 2: Benchmarking of SV calling pipeline for duplications | 6 |
| Supplementary Note 3: Identifying novel variants with respect to gnomAD-SV catalogue | 7 |
| Supplementary Note 4: Identifying novel variants with respect to 1000G-SV catalogue | 8 |
| Supplementary Figures | 9 |
| Supplementary Fig. 1 Different types of structural variations detected in SG10K-SV. | 9 |
| Supplementary Fig. 2 Benchmarking of various SV callers for deletions and insertions. | 10 |
| Supplementary Fig. 3 True positive, false positive and false negative counts for Manta, Delly, Smoove and their combination for all classed of SVs..... | 11 |
| Supplementary Fig. 4 Comparison of SurVindel2 and Manta-SVimmer-Graphtyper2 pipeline for duplication identification. | 12 |
| Supplementary Fig. 5 Barplot showing the number of duplications detected by Manta-SVimmer-Graphtyper2 and SurVindel2 in different genomic regions. | 13 |
| Supplementary Fig. 6 Violin plot showing the number of events per genome for the Validation datasets..... | 14 |
| Supplementary Fig. 7 Allele distribution for the two validation datasets. | 15 |
| Supplementary Fig. 8 Samplot of a 9.16kb deletion event overlapping the PRKAG2 gene region..... | 16 |

| | | |
|----|--|-----------|
| 40 | Supplementary Fig. 9: Distribution of novel Asian-specific and known Asian- | |
| 41 | specific SVs across different allele frequency bins..... | 17 |
| 42 | Supplementary Fig. 10 Scatter plot of the top-2 principal components of a | |
| 43 | SG10K_Health dataset Single Nucleotide Variant based PCA analysis showing the | |
| 44 | population structure in the Singaporean population..... | 18 |
| 45 | Supplementary Fig. 11 PCA of variants in the discovery dataset showing the | |
| 46 | population structure in the SG10K-SV-r1.4. | 19 |
| 47 | Supplementary Fig. 12 Distribution of SVs shared among ethnic group across | |
| 48 | different allele frequency bins..... | 20 |
| 49 | <i>Supplementary references</i> | 21 |
| 50 | <i>SG10K_Health Consortium</i> | 22 |
| 51 | | |
| 52 | | |
| 53 | | |
| 54 | | |

55 **Supplementary Notes**

56 **Supplementary Note 1: Benchmarking of SV calling pipeline for** 57 **deletions and insertions**

58
59 We aimed to comprehensively assess the performance of our SV calling pipeline by
60 comparing the SG10K-SV (Manta¹) pipeline with three other popular SV detection
61 algorithms such as Delly² and Smoove³.

62
63 To accurately benchmark the performance of our SV detecting pipeline, we
64 downloaded a subset of 34 1000 Genome samples with both long and short read
65 whole genome sequencing (WGS) data. We retrieved the 30x short read WGS CRAM
66 files from <https://registry.opendata.aws/1000-genomes/>. Long-read sequencing data
67 have become the technique of choice for SV detections and hence it will serve as the
68 truth set for the comparison. We retrieved the comprehensive catalogue of SVs
69 detected using long-read sequencing from Ebert *et al.*⁴
70 ([https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGSVC2/release/v2.0/in-](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGSVC2/release/v2.0/integrated_callset/)
71 [tegrated_callset/](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGSVC2/release/v2.0/integrated_callset/)) to ascertain the sensitivity and precision of the short-read SVs
72 predicted using our SV detection pipeline and 2 other SV detection algorithms.

73

74 SV discovery using Manta

75 Manta¹ was executed in the single sample mode to identify deletions and insertions in
76 the 34 1000G samples using default parameters. We used SVimmer⁵ to cluster SVs
77 across samples using the default parameters and re-genotyped the SVs in each
78 sample using Graphtyper2⁶ with default parameters. We then merged the individual
79 re-genotyped VCF using Graphtyper2's *Vcfmerge* function. Lastly, we retained PASS
80 calls made under the aggregated genotyping model for downstream analysis. In
81 addition, we applied additional filters recommended by Graphtyper2.

82 For deletions, we filter the variants using bcftools with the following command:

```
83 bcftools 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="DEL" && QD > 9 &&  
84 (ABHet > 0.3 || ABHet < 0 ) && (AC/NUM_MERGED_SVS) < 25 && PASS_AC > 0 &&  
85 PASS_ratio > 0.1' ${vcf} | bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs . - |  
86 bcftools view -c1 -s ${meta} --output ${prefix}.DEL.vcf.gz -0z --threads $task.cpus  
87 -
```

88

89 For duplications, we retain variants which passed the following criteria:

```

90 bcftools view -i 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="DUP" &&
91 QD > 5 && (AC/NUM_MERGED_SVS) < 25 && PASS_AC >0 ' ${vcf} |
92 bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs . - | bcftools view -c1 -s ${meta} -
93 -output ${prefix}.DUPonly.vcf.gz -0z --threads $task.cpus -
94

```

95 Lastly, for insertions, we filtered the variants with bcftools using the following command:

```

96 bcftools view -i 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="INS" &&
97 PASS_AC >0 && (AC/NUM_MERGED_SVS) < 25 && PASS_ratio > 0.1 && (ABHet > 0.25 ||
98 ABHet < 0) && MaxAAS > 4' ${vcf} | bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs .
99 --threads $task.cpus - | bcftools view -c1 -s ${meta} --output
100 ${prefix}.INSONly.vcf.gz -0z --threads $task.cpus -
101

```

102 SV discovery using Delly

103 Delly² v1.2.6 was executed in the single sample mode to identify deletions and
104 insertions in the 34 1000G samples using default parameters. BCFtools⁷ was used to
105 convert the bcf output from Delly to VCF format before clustering SVs across samples
106 using SVimmer⁵. The SVs were re-genotyped in each sample using Graphtyper2⁶ with
107 default parameters. We merged the individual re-genotyped VCF using Graphtyper2's
108 *Vcfmerge* function. Lastly, we retained PASS calls made under the aggregated
109 genotyping model for downstream analysis. In addition, we applied additional filters
110 recommended by Graphtyper2.

111 For deletions, we filter the variants using bcftools with the following command:

```

112 bcftools 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="DEL" && QD > 9 &&
113 (ABHet > 0.3 || ABHet < 0) && (AC/NUM_MERGED_SVS) < 25 && PASS_AC > 0 &&
114 PASS_ratio > 0.1' ${vcf} | bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs . - |
115 bcftools view -c1 -s ${meta} --output ${prefix}.DEL.vcf.gz -0z --threads $task.cpus
116 -
117

```

118 For duplications, we retain variants which passed the following criteria:

```

119 bcftools view -i 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="DUP" &&
120 QD > 5 && (AC/NUM_MERGED_SVS) < 25 && PASS_AC >0 ' ${vcf} |
121 bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs . - | bcftools view -c1 -s ${meta} -
122 -output ${prefix}.DUPonly.vcf.gz -0z --threads $task.cpus -
123

```

124 Lastly, for insertions, we filtered the variants with bcftools using the following command:

```

125 bcftools view -i 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="INS" &&
126 PASS_AC >0 && (AC/NUM_MERGED_SVS) < 25 && PASS_ratio > 0.1 && (ABHet > 0.25 ||
127 ABHet < 0) && MaxAAS > 4' ${vcf} | bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs .
128 --threads $task.cpus - | bcftools view -c1 -s ${meta} --output
129 ${prefix}.INSONly.vcf.gz -0z --threads $task.cpus -
130

```

131

132 SV discovery using Smoove

133 Smoove was executed in the single sample mode to identify structural variations in the
134 34 1000G samples using Smoove Call function with default parameters. Variants were
135 merged across samples using the Smoove Merge function with default parameters.
136 Lastly, SVs were re-genotyped in each sample using the Smoove Genotype function
137 with default parameters.

138

139 Combining SVs detected across the three algorithms

140 We obtained the single sample calls from each of the algorithms (Manta, Smoove,
141 Delly) and clustered across all samples and algorithm using SVimmer⁵ with the default
142 parameters. Lastly, we re-genotyped SVs in each sample using Graphtyper2⁶ with
143 default parameters and merged the individual re-genotyped VCF using Graphtyper2's
144 *Vcfmerge* function. Lastly, we retained PASS calls made under the aggregated
145 genotyping model for downstream analysis. In addition, we applied additional filters
146 recommended by Graphtyper2.

147 For deletions, we filter the variants using bcftools with the following command:

```
148 bcftools 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="DEL" && QD > 9 &&  
149 (ABHet > 0.3 || ABHet < 0 ) && (AC/NUM_MERGED_SVS) < 25 && PASS_AC > 0 &&  
150 PASS_ratio > 0.1' ${vcf} | bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs . - |  
151 bcftools view -c1 -s ${meta} --output ${prefix}.DEL.vcf.gz -0z --threads $task.cpus  
152 -
```

153

154 For duplications, we retain variants which passed the following criteria:

```
155 bcftools view -i 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="DUP" &&  
156 QD > 5 && (AC/NUM_MERGED_SVS) < 25 && PASS_AC > 0 ' ${vcf} |  
157 bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs . - | bcftools view -c1 -s ${meta} -  
158 -output ${prefix}.DUPonly.vcf.gz -0z --threads $task.cpus -
```

159

160 Lastly, for insertions, we filtered the variants with bcftools using the following command:

```
161 bcftools view -i 'INFO/SVMODEL="AGGREGATED" && FILTER="PASS" && SVTYPE="INS" &&  
162 PASS_AC > 0 && (AC/NUM_MERGED_SVS) < 25 && PASS_ratio > 0.1 && (ABHet > 0.25 ||  
163 ABHet < 0) && MaxAAS > 4' ${vcf} | bcftools filter -i 'FMT/FT ="PASS" ' --set-GTs .  
164 --threads $task.cpus - | bcftools view -c1 -s ${meta} --output  
165 ${prefix}.INSonly.vcf.gz -0z --threads $task.cpus -
```

166

167 Calculating precision, recall and F1-Score

168 To evaluate the performance of different SV algorithm, we focus the test on the
169 presence and absence of the variants in the long read dataset. We calculate the
170 precision, recall and F1-Score using Truvari⁸ with the SV calls from long read data

171 from Ebert *et al.*⁴ as the truth set. A variant is defined as a true positive (TP) if the
172 variant is found in both short-read and long-read dataset. A variant is defined as a
173 false positive (FP) if it is not found in the long read dataset.

174

175 Precision is defined as:

$$176 \quad \textit{Precision} = \frac{TP}{TP+FP} \quad (1)$$

177

178 Recall is defined as:

179

$$180 \quad \textit{Recall} = \frac{TP}{TP+FN} \quad (2)$$

181

182 F1-score is defined as:

$$183 \quad F1 = 2 * \frac{\textit{Precision} * \textit{Recall}}{\textit{Precision} + \textit{Recall}} \quad (3)$$

184

185

186 Evaluation of SV detection pipelines using 34 1000G project WGS data

187 The precision and recall for SV detection varied depending on the method. Fig 1c, d
188 and e show the precision, recall and F1 (combined statistics of precision and recall) of
189 the three different SV calling pipeline for calling structural variations.
190 Delly+Graphtyper2 has the highest precision in terms of SV detection for both 30x and
191 15x sequencing libraries. However, Manta+Graphtyper2 has a higher recall and F1
192 score compared to the rest of the pipelines. Combining SVs detected by all three
193 pipelines did not improve the precision, recall and F1-Score compared to running the
194 Manta+Graphtyper2 pipeline.

195

196 When analyzing the SVs separately based on the SV classes, Delly+Graphtyper2 has
197 the highest precision in terms of deletion and insertion detection for both 30x and 15x
198 sequencing libraries (Supplementary Fig. 2). In terms of recall and F1-score,
199 Manta+Graphtyper2 outperforms the other pipelines for both deletions and insertions.

200

201 Next, we evaluate how sequencing read depth affects the performance of the SV
202 pipelines. We down-sampled the 30x CRAM files to 15x using Sambamba⁹ and we
203 evaluate the performance of the four different approaches to detect SVs. Differences

204 in sequencing depth affect precision and recall (Fig. 1c, d, e, and Supplementary Fig.
205 2 and 3). Interestingly, the 30x dataset has a higher recall for all four approaches. In
206 addition, data with a higher sequencing depth (30x) have a lower precision compared
207 to the sequencing data with a lower coverage (15x). The lower precision and higher
208 recall of the 30x data could be attributed to higher number of misaligned reads leading
209 to spurious SV calling¹⁰.

210

211 To estimate the number of variants that are missed or incorrectly called using the 15x
212 samples as compared to 30x, we obtained the true positives (TP), false positives (FP),
213 and false negative (FN) counts for each pipeline across different sequencing depth.
214 Across all SV pipelines, 15x libraries have a lower FP count compared to the 30x
215 libraries (Supplementary Fig. 3). This could be attributed to the higher number of
216 misaligned reads in the 30x libraries which could lead to spurious SV calls.

217

218 **Supplementary Note 2: Benchmarking of SV calling pipeline for** 219 **duplications**

220

221 As the Manta-SVimmer-Grphtyper2 SV pipeline relies solely on discordant read pairs
222 and split-read alignments, it has inherent limitations to accurately detect duplication
223 events created by the presence of tandem repeat sequences (e.g., microsatellites and
224 minisatellites)^{11,12}. We thus complemented the above algorithms with SurVIndel2¹³, an
225 in-house developed algorithm that can detect duplication events at high sensitivity
226 (Supplementary Fig. 4).

227

228 To demonstrate the robustness of SurVindel2, we assessed false discovery rate (FDR)
229 and true positive (TP) statistics for duplications relative to Manta-SVimmer-
230 Grphtyper2, against a truth set of high quality SVs obtained by haplotype-resolved
231 long-read sequencing of a selected subset of 1000 Genomes Project analyzed
232 samples⁴.

233

234 We downloaded CRAM files at 30x coverage are available for all the samples¹⁴. We
235 randomly selected 10 samples for our benchmarking effort and down-sampled the
236 CRAM files to 15x using samtools⁷ we ran our pipeline on a dataset comprising 5,487

237 discovery samples plus the 10 benchmarking samples. we ran our pipeline on a
238 dataset comprising 5,487 discovery samples plus the 10 benchmarking samples. to
239 mimic our discovery dataset.

240

241 For this benchmarking, we ran our pipeline on a dataset comprising 5,487 discovery
242 samples plus the 10 benchmarking samples. we ran our pipeline on a dataset
243 comprising 5,487 discovery samples plus the 10 benchmarking samples. Finally, we
244 obtained a call set for each sample by retaining SVs with an allele count of at least 1
245 and an FS value of PASS.

246

247 We used an in-house tool (<https://github.com/Mesh89/SVComparator>) to compare, for
248 each sample, the predicted SVs with the set of SVs reported in HGVC2. Our pipeline
249 reports tandem duplications and insertions separately, while HGVC2 only reports
250 deletions and insertions; tandem duplications are considered insertions. For this
251 reason, we could not measure the sensitivity of our duplications and insertions
252 separately.

253

254 We measured an average per-sample duplication identification FDR of 12% and 36%
255 for SurVindel2 and Manta-SVimmer-Grphtyper2, respectively. SurVindel2 yielded a
256 better sensitivity than Manta-SVimmer-Grphtyper2 (Supplementary Fig. 3,
257 Supplementary Table 2). Furthermore, the gains in sensitivity were more pronounced
258 for tandem repeats (Supplementary Fig. 5).

259

260 One of the significant challenges when generating a dataset of SVs for a large
261 population is maintaining a low level of noise. Our benchmarking efforts show that our
262 call set is precise (average precision is 0.91 for deletions, 0.88 for duplications and
263 0.72 for insertions) (Supplementary Table 3). Unsurprisingly, long reads can discover
264 far more SVs compared to 15x Illumina paired-end reads. However, the number of
265 deletions, duplications and insertions we discovered is comparable to recent studies
266 such as gnomAD¹⁵ while using a lower sequencing depth. Coupled with the good
267 precision, we conclude that our pipeline is in line with the state of the art in the field.

268

269 **Supplementary Note 3: Identifying novel variants with respect to gnomAD-**
270 **SV catalogue**

271

272 To identify SVs that have a higher prevalence in Asian population within gnomAD-SV
273 catalogue, we first identify variants that overlap between SG10K-SV and gnomAD-SV
274 using SVimmer⁵. We identified 23,434 SVs in the SG10K-SV dataset which overlap
275 with gnomAD-SV. This includes 4,725 deletions, 7,458 duplications and 11,251
276 insertions.

277

278

279 **Supplementary Note 4: Identifying novel variants with respect to 1000G-**
280 **SV catalogue**

281

282 To identify SVs that have a higher prevalence in Asian population within 1000G-SV
283 catalogue, we first identify variants that overlap between SG10K-SV and 1000G-SV
284 using SVimmer⁵. We identified 9,668 SVs in the SG10K-SV dataset which overlap with
285 1000G-SV. This includes 3,105 deletions, 284 duplications and 6,279 insertions.

286

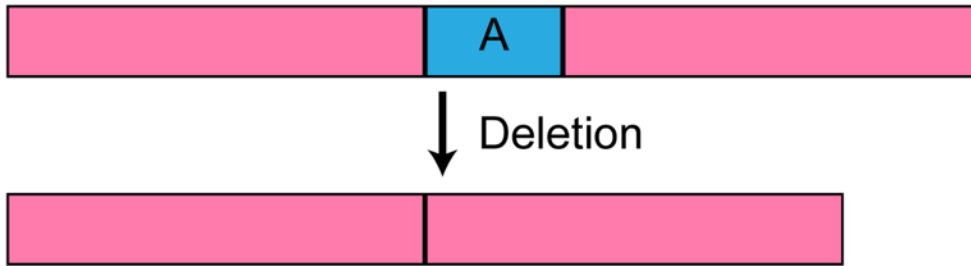
287

288

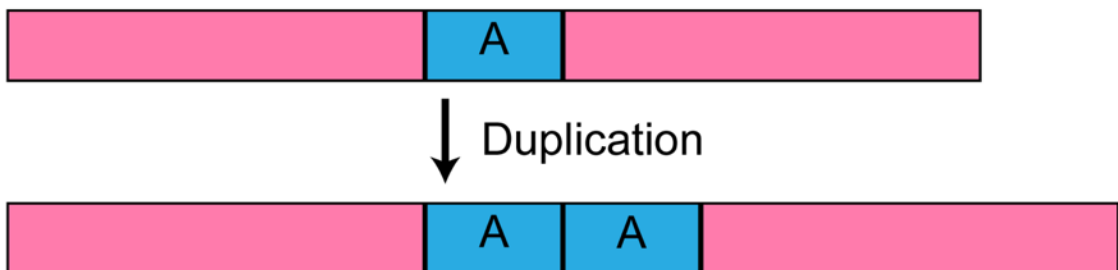
289 **Supplementary Figures**

290

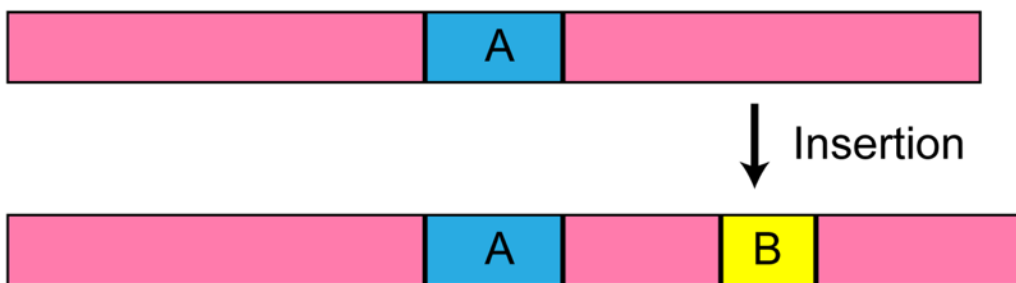
A



B



C



291

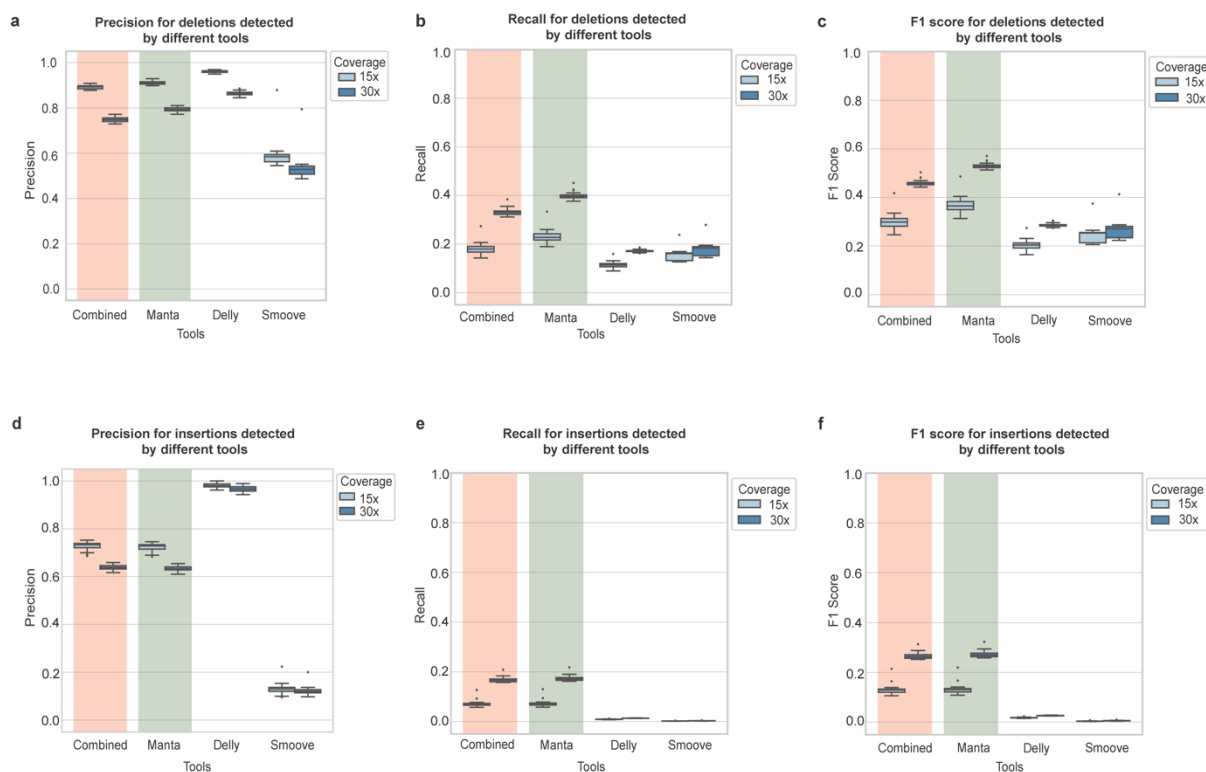
292

293

294

Supplementary Fig. 1 Different types of structural variations detected in SG10K-SV.

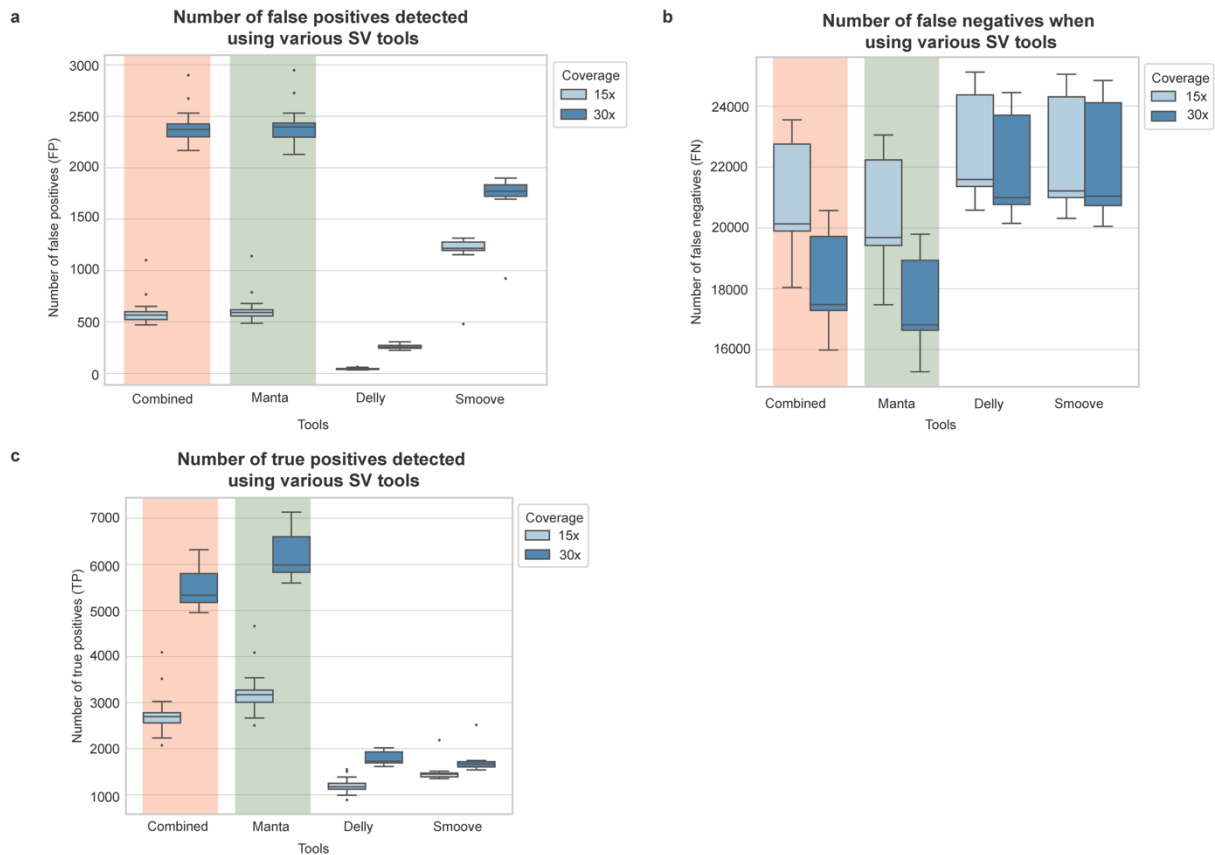
295
296



297
298
299
300
301
302
303
304
305
306
307
308
309

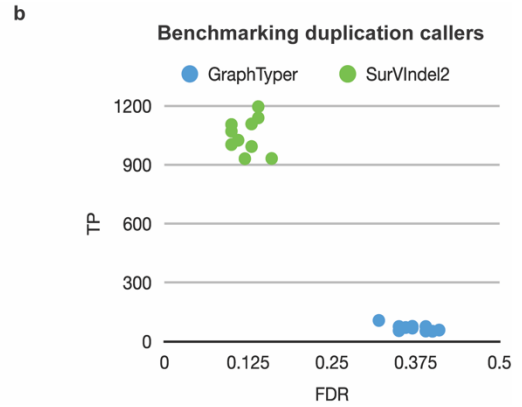
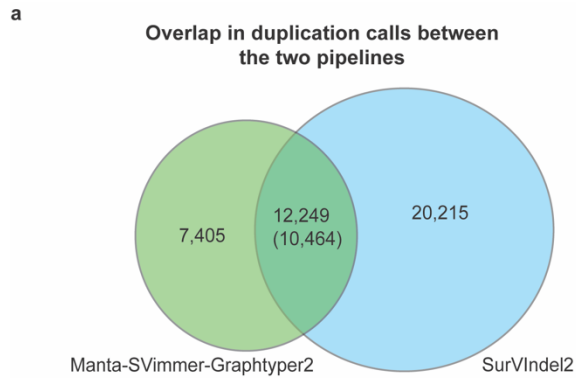
Supplementary Fig. 2 Benchmarking of various SV callers for deletions and insertions using 34 1000G samples with two different sequencing depths.

a Boxplot showing the precision for deletions between 15x and 30x coverage for each caller. Combined refers to variants that are detected in all three pipelines. **b** Boxplot showing the recall for deletions between 15x and 30x coverage for each caller. **c** Boxplot showing the F1-score for deletions between 15x and 30x coverage for each caller. **d** Boxplot showing the precision for insertions between 15x and 30x coverage for each caller. **e** Boxplot showing the recall for insertions between 15x and 30x coverage for each caller. **f** Boxplot showing the F1-score for insertions between 15x and 30x coverage for each caller. The boxplots in a-f display the median and first/third quartiles.



310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322

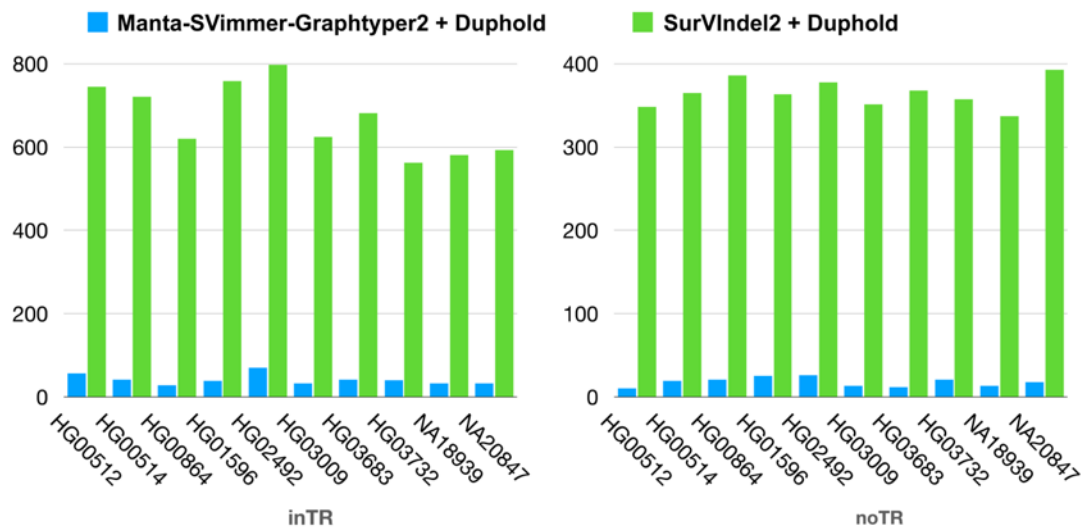
Supplementary Fig. 3 True positive, false positive and false negative counts for Manta, Delly, Smoove and their combination for all classed of SVs using 34 1000G samples with two different sequencing depth (15x and 30x coverage). **a** Boxplot showing the number of false positive counts between 15x and 30x coverage for each SV caller. Combined refers to variants that are detected in all three pipelines. **b** Boxplot showing the false negative counts between 15x and 30x coverage for each SV caller. **c** Boxplot showing the true positive counts between 15x and 30x coverage for each SV caller. The boxplots showed in a-c display the median and first/third quartiles.



323
324
325
326
327
328
329
330
331

Supplementary Fig. 4 Comparison of SurVindel2 and Manta-SVimmer-Graphtyper2 pipeline for duplication identification.

a Comparison of the number of duplications detected by Manta-Graphtyper2 and SurVindel2. **b** Scatterplot comparing the number of true positives detected duplication and FDR achieved with Manta-SVimmer-Graphtyper2 and SurVindel2 for a truth set of high quality SVs obtained by haplotype-resolved long-read sequencing of a selected subset of 1000 Genomes Project analyzed samples⁴.



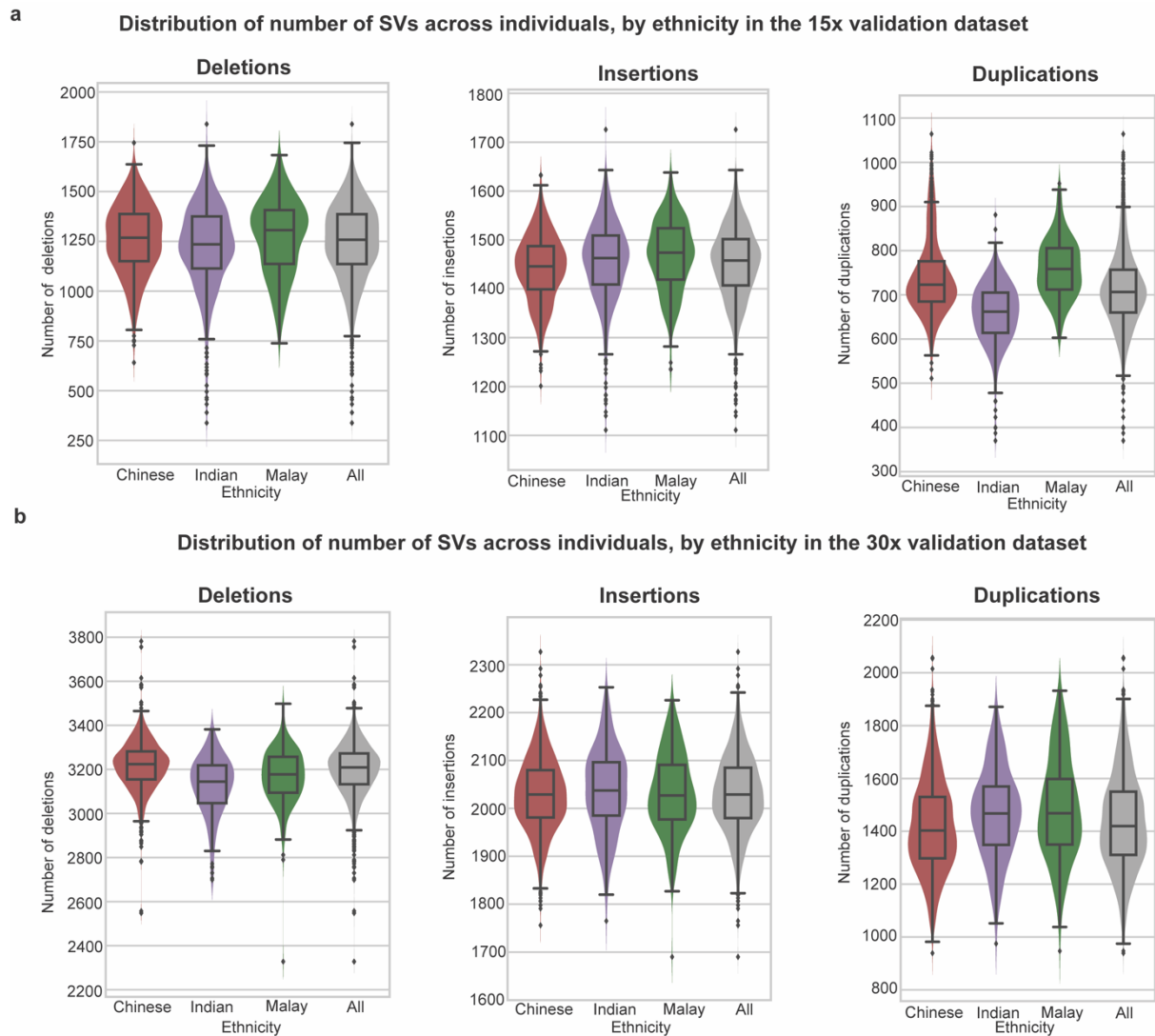
333

334 **Supplementary Fig. 5 Barplot showing the number of duplications detected by**
 335 **Manta-SVimmer-Graphtyper2 and SurVindel2 in different genomic regions.**

336 The Y-axis shows the number of SVs and the X-axis shows the sample name for each
 337 1KG sample. Blue bars indicate the number of duplications detected in each 1KG
 338 sample by the Manta-SVimmer-Graphtyper2 pipeline. Green bars indicate the number
 339 of duplications detected in each 1KG sample by SurVindel2.

340 The barplot on the left shows the number of SVs detected by Manta-SVimmer-
 341 Graphtyper2 and SurVindel2 in tandem repeat regions. The barplot on the right shows
 342 the number of SVs Manta-SVimmer-Graphtyper2 and SurVindel2 in non-tandem
 343 repeat regions. SurVindel2 detects more duplications in both tandem repeat and
 344 non-tandem repeat regions compared to Manta-SVimmer-Graphtyper2.

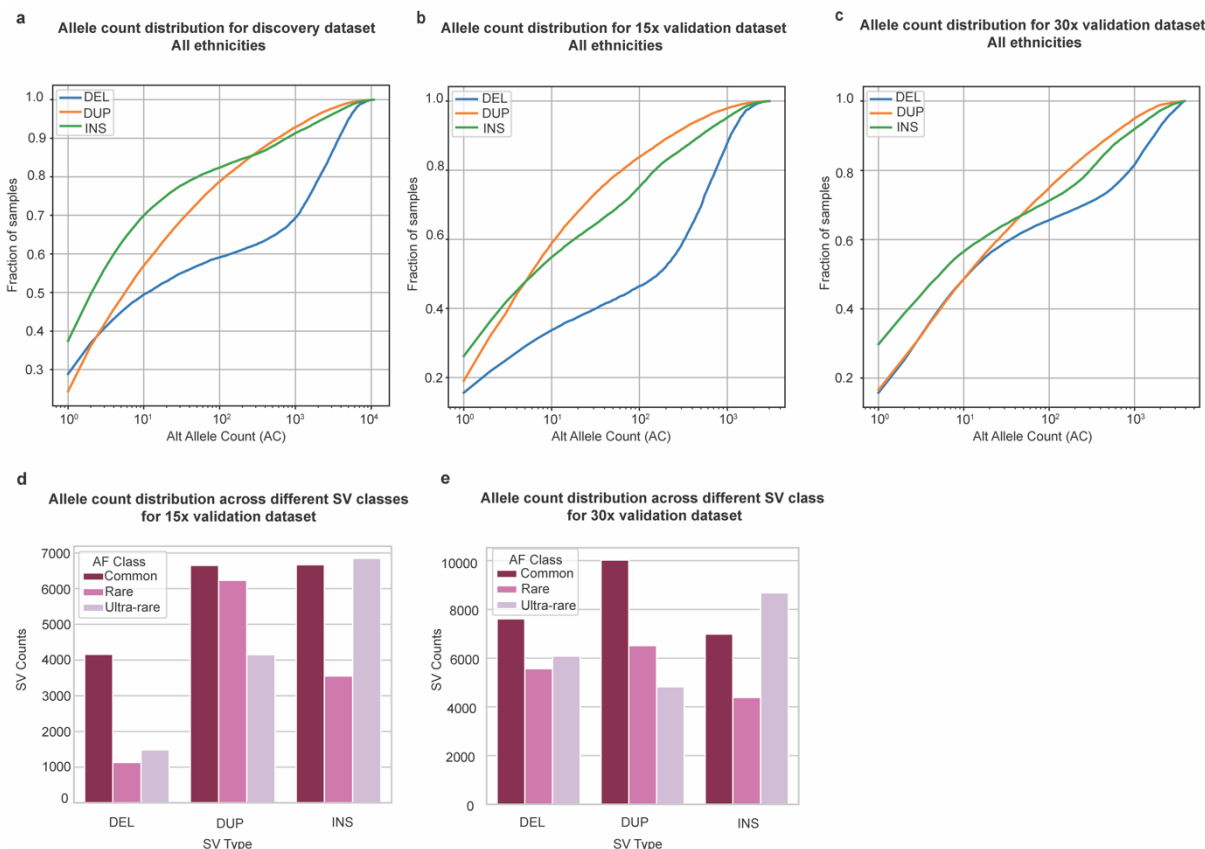
345



346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363

Supplementary Fig. 6 Violin plot showing the number of events per genome for the Validation datasets.

a Violin plots and boxplots showing the number of events per genome for each ethnic group (number of Chinese = 663, number of Malays = 278, number of Indians = 581). DEL, deletions; DUP, duplications; INS, insertions (including MEIs) in the 15x validation dataset. The boxplots display the minimum and maximum number of SVs as well as the median and the first/third quartile. **b** Violin plots and boxplots showing the number of events per genome for each ethnic group (number of Chinese = 1,433, number of Malays = 288, number of Indians = 198). DEL, deletions; DUP, duplications; INS, insertions (including MEIs) in the 30x validation dataset. The boxplots display the minimum and maximum number of SVs as well as the median and the first/third quartile.



365

366

Supplementary Fig. 7 Allele distribution for the two validation datasets.

367

a Distribution of alternate allele counts for different class of SVs in the discovery

368

dataset. **b** Distribution of alternate allele counts for different classes of SVs in the

369

SG10K 15x validation dataset. The majority of the SVs are rare variants ($AF < 1\%$).

370

c Distribution of alternate allele counts for different classes of SVs in the SG10K 30x

371

validation dataset. The majority of the SVs are rare variants ($AF < 1\%$).

372

d Allele count distribution across different SV classes segregated by allele frequency classes for 15x

373

validation dataset. Allele frequency (AF) bins: Common ($AF \geq 0.01$), rare ($0.01 > AF \geq$

374

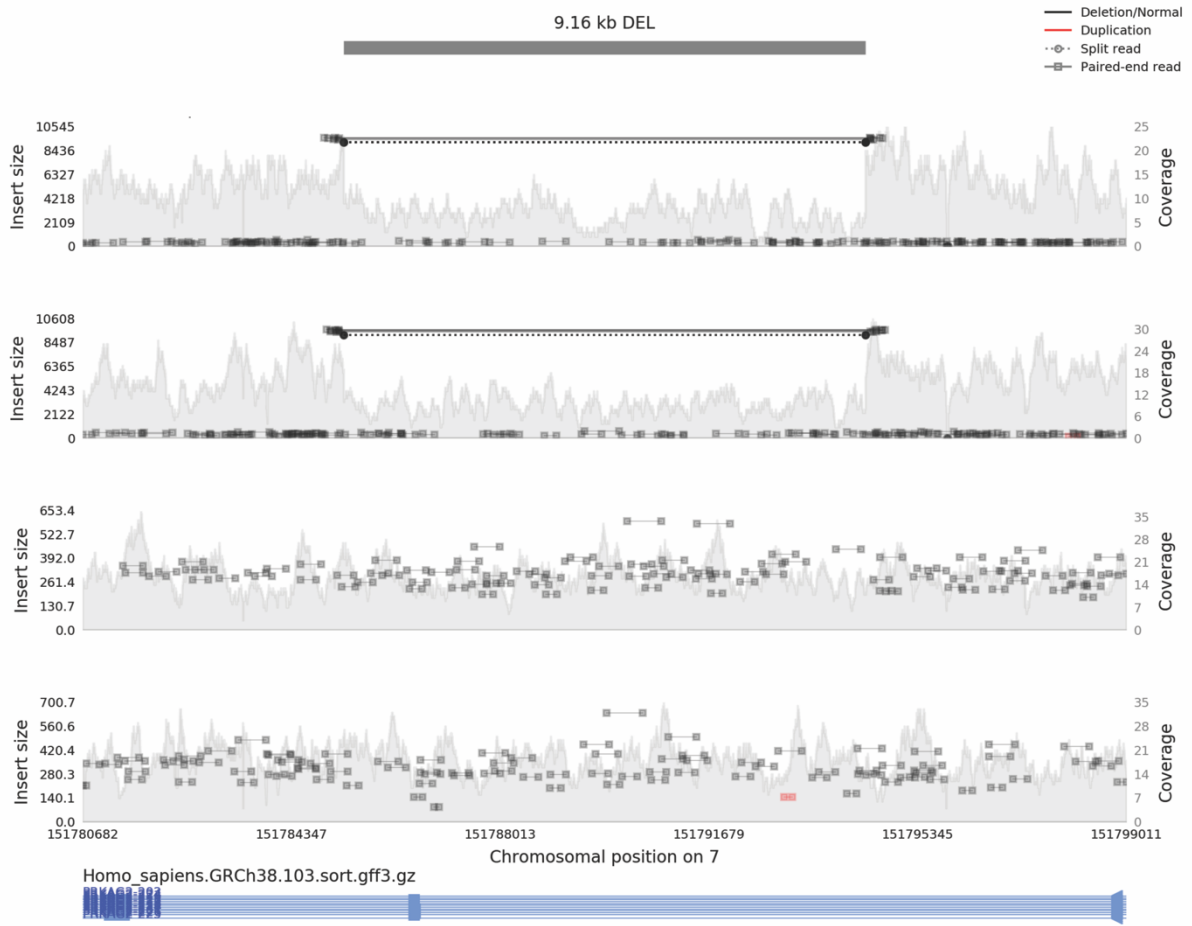
0.001) and ultra-rare ($AF < 0.001$).

375

e Allele count distribution across different SV classes segregated by allele frequency classes for 30x validation dataset.

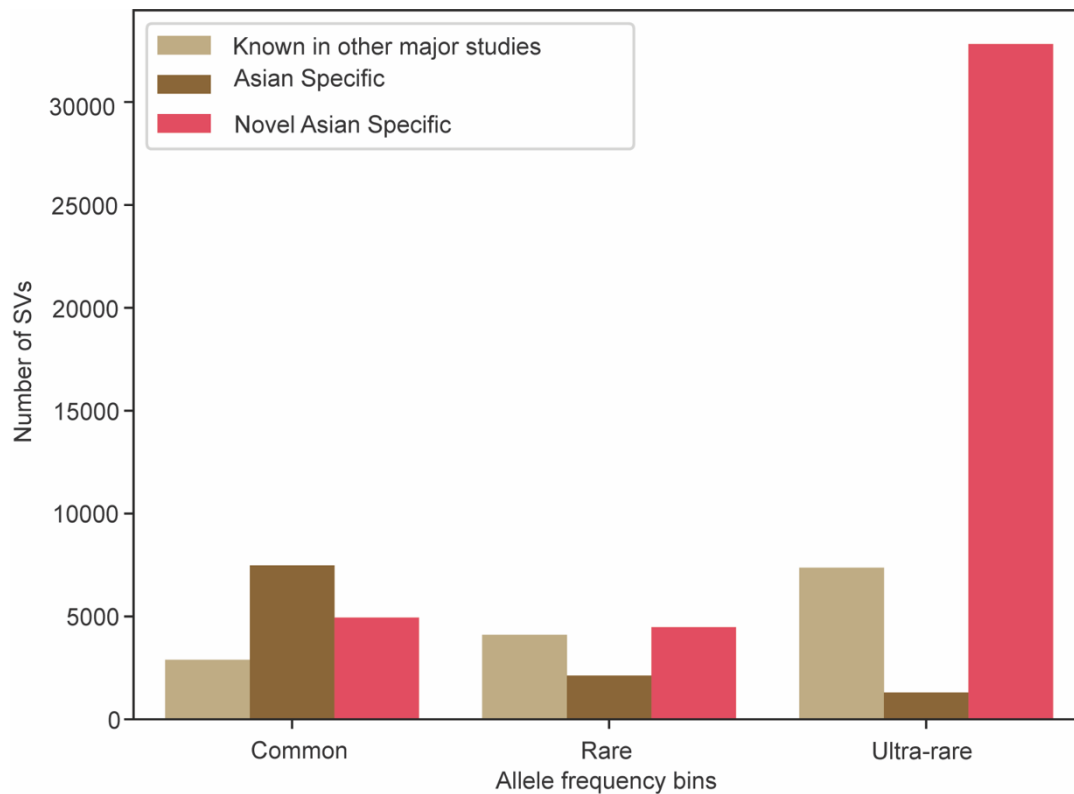
376

377



378
 379
 380
 381

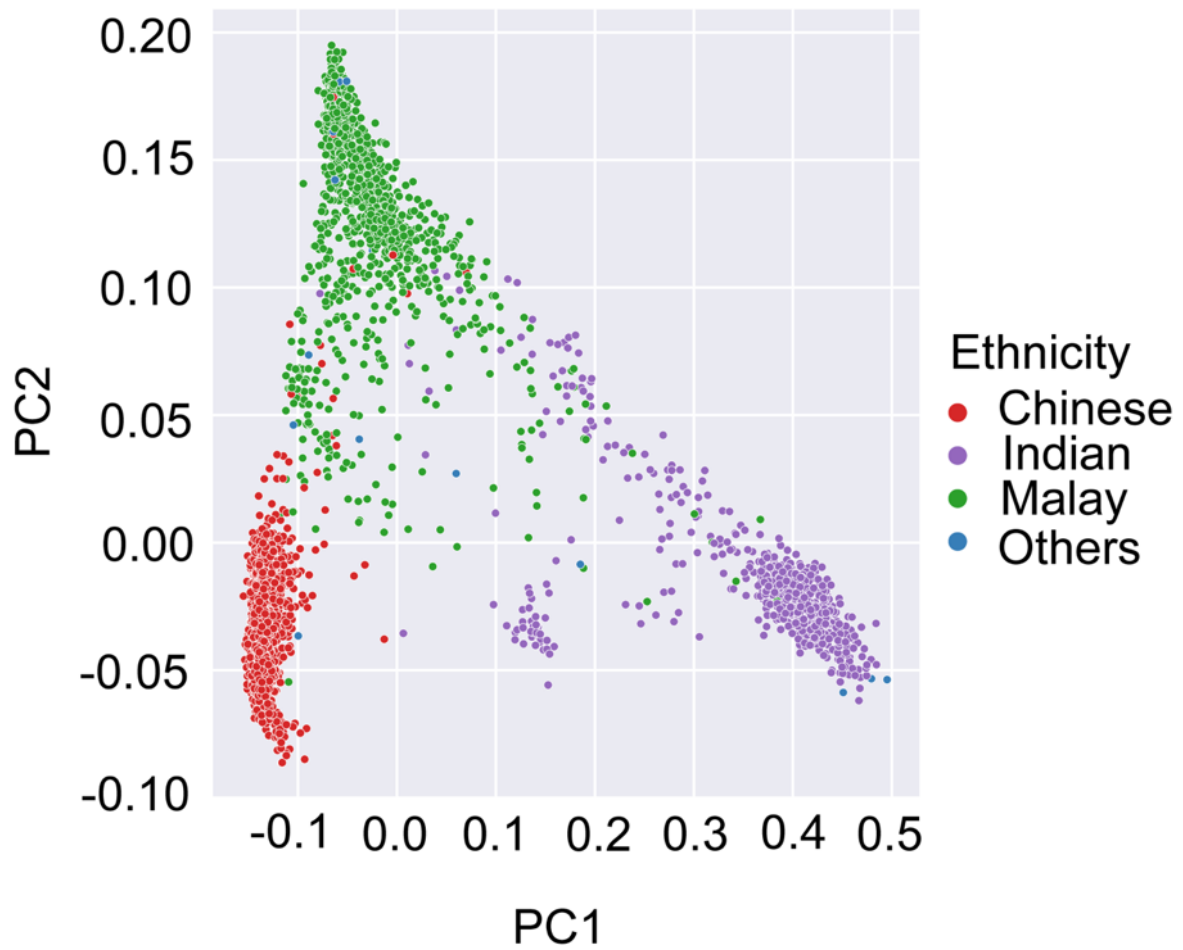
Supplementary Fig. 8 Samplot of a 9.16kb deletion event overlapping the *PRKAG2* gene region.



382
383
384
385
386
387
388
389
390
391
392
393

Supplementary Fig. 9: Distribution of novel Asian-specific and known Asian-specific SVs across different allele frequency bins.

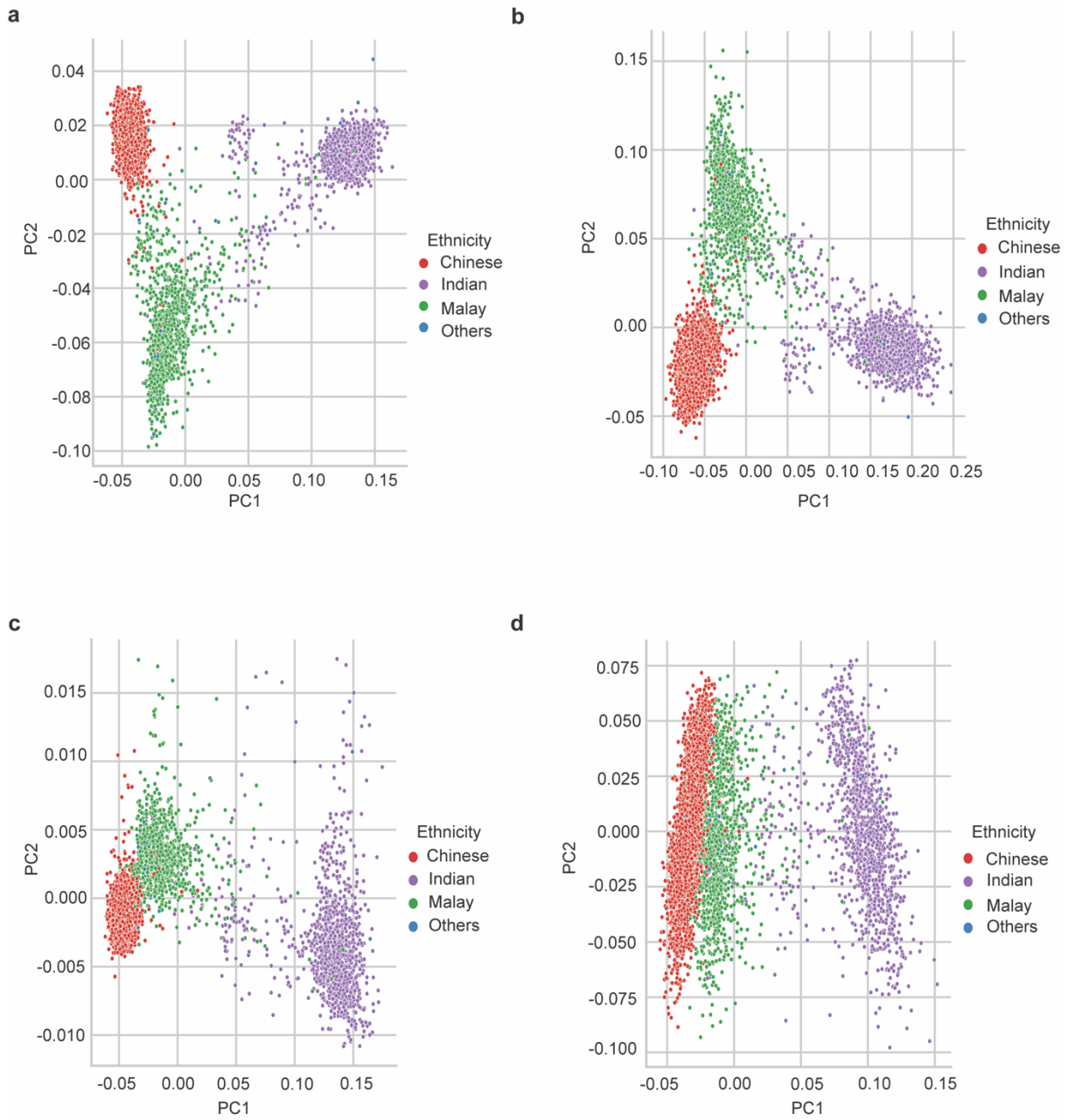
Light-brown bars indicate SVs in SG10K-SV that overlap either gnomAD-SV or 1000G-SV and do not have significant F_{st} . Brown bars indicate SVs in SG10K-SV which overlap variants in either gnomAD-SV and 1000G and have significant F_{st} , and therefore, they are termed as “Asian-specific”. Red bars indicate SVs that are only found in SG10K and have a call rate of ≥ 0.5 in Chinese, Malay, or Indians. These SVs are referred to as “Novel Asian Specific” SVs. The SVs are further partitioned into three different within SG10K-SV allele frequency (AF) bins: Common ($AF \geq 0.01$), rare ($0.01 > AF \geq 0.001$) and ultra-rare ($AF < 0.001$).



394
395
396
397
398

Supplementary Fig. 10 Scatter plot of the top-2 principal components of a SG10K_Health dataset Single Nucleotide Variant based PCA analysis showing the population structure in the Singaporean population.

399
400

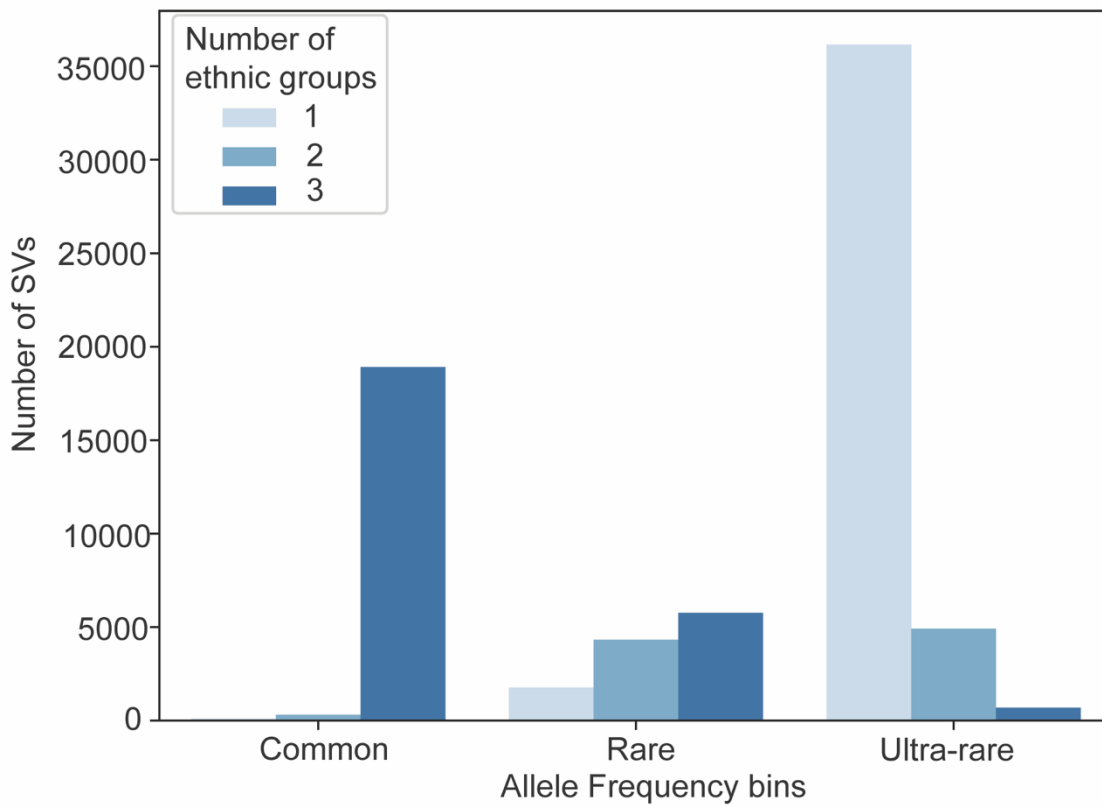


401
402
403
404
405
406
407

Supplementary Fig. 11 PCA of variants in the discovery dataset showing the population structure in the SG10K-SV-r1.4.

a PCA of all variants in the discovery dataset. **b** PCA using deletions only. **c** PCA using insertions only. **d** PCA using duplications only.

408
409
410



411
412
413
414
415
416
417
418
419

Supplementary Fig. 12 Distribution of SVs shared among ethnic group across different allele frequency bins.

Different shades of blue indicate the number of ethnic groups in which the SV is detected in. The SVs are furthered partition into three different allele frequency bins. Common indicates variants with allele frequency ≥ 0.01 ; rare indicates variants with allele frequency ≥ 0.001 and allele frequency < 0.01 ; ultra-rare variants refers to variants with allele frequency < 0.001 .

420 Supplementary references

- 421
- 422 1. Chen, X. *et al.* Manta: rapid detection of structural variants and indels for
423 germline and cancer sequencing applications. *Bioinformatics* **32**, 1220-2
424 (2016).
- 425 2. Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end
426 and split-read analysis. *Bioinformatics* **28**, i333-i339 (2012).
- 427 3. Pedersen, B.S., Layer, R., Quinlan, A. R. smooove: structural-variant calling
428 and genotyping with existing tools. 0.2.8 edn
429 (<https://github.com/brentp/smooove>, 2020).
- 430 4. Ebert, P. *et al.* Haplotype-resolved diverse human genomes and integrated
431 analysis of structural variation. *Science* **372**(2021).
- 432 5. Eggertsson, H.P. Structural Variant Merging Tool.
433 (<https://github.com/DecodeGenetics/svimmer>, 2021).
- 434 6. Eggertsson, H.P. *et al.* GraphTyper2 enables population-scale genotyping of
435 structural variation using pangenome graphs. *Nature Communications* **10**,
436 5402 (2019).
- 437 7. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience*
438 **10**(2021).
- 439 8. English, A.C., Menon, V.K., Gibbs, R.A., Metcalf, G.A. & Sedlazeck, F.J.
440 Truvari: refined structural variant comparison preserves allelic diversity.
441 *Genome Biology* **23**, 271 (2022).
- 442 9. Tarasov, A., Vilella, A.J., Cuppen, E., Nijman, I.J. & Prins, P. Sambamba: fast
443 processing of NGS alignment formats. *Bioinformatics* **31**, 2032-2034 (2015).
- 444 10. Kosugi, S. *et al.* Comprehensive evaluation of structural variation detection
445 algorithms for whole genome sequencing. *Genome Biology* **20**, 117 (2019).
- 446 11. Rajaby, R. & Sung, W.K. SurVIndel: improving CNV calling from high-
447 throughput sequencing data through statistical testing. *Bioinformatics* **37**,
448 1497-1505 (2021).
- 449 12. Dashnow, H. *et al.* STRetch: detecting and discovering pathogenic short
450 tandem repeat expansions. *Genome Biology* **19**, 121 (2018).
- 451 13. Rajaby, R. & Sung, W.-K. SurVIndel2: improving local CNVs calling from next-
452 generation sequencing using novel hidden information. *bioRxiv*,
453 2023.04.23.538018 (2023).
- 454 14. Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the
455 expanded 1000 Genomes Project cohort including 602 trios. *Cell* **185**, 3426-
456 3440.e19 (2022).
- 457 15. Abel, H.J. *et al.* Mapping and characterization of structural variation in 17,795
458 human genomes. *Nature* **583**, 83-89 (2020).
- 459
- 460

461 **SG10K_Health Consortium**

462 Khung Keong Yeo³⁴, Stuart Alexander Cook³⁴, Chee Jian Pua³⁴, Chengxi Yang³⁴,
463 Tien Yin Wong¹², Charumathi Sabanayagam^{12,35}, Lavanya Raghavan¹², Tin
464 Aung^{12,35}, Miao Ling Chee¹², Miao Li Chee¹², Hengtong Li^{12,13}, Jimmy Lee^{36,37}, Eng
465 Sing Lee^{38,39}, Joanne Ngeow^{23,40}, Paul Eillot⁴¹, Elio Riboli⁴¹, Hong Kiat Ng²³,
466 Theresia Mina²³, Darwin Tay²³, Nilanjana Sadhu²³, Pritesh Rajesh Jain²³, Dorrain
467 Low²³, Xiaoyan Wang²³, Jin Fang Chai⁷, Rob M Van Dam^{7,42,43}, Yik Ying Teo⁷, Chia
468 Wei Lim¹¹, Pi Kuang Tsai¹¹, Wen Jie Chew⁴⁴, Wey Ching Sim¹¹, Li-xian Grace Toh¹¹,
469 Johan Gunnar Eriksson^{19,45}, Peter D Gluckman^{46,47}, Yung Seng Lee^{19,48}, Fabian
470 Yap⁴⁹, Kok Hian Tan⁵⁰

471

472 ¹Genome Institute of Singapore, Agency for Science, Technology and Research,
473 Singapore

474 ²Duke-NUS Medical School, Singapore

475 ³SingHealth Duke-NUS Institute of Precision Medicine, Singapore Health Services,
476 Duke-NUS Medical School, Singapore

477 ⁴SingHealth Duke-NUS Genomic Medicine Centre, Duke-NUS Medical School,
478 Singapore

479 ⁵Saw Swee Hock School of Public Health, National University of Singapore and
480 National University Health System, Singapore

481 ⁶Department of Obstetrics & Gynaecology, Yong Loo Lin School of Medicine,
482 National University of Singapore, Singapore

483 ⁷Institute for Human Development and Potential (IHDP), Agency for Science,
484 Technology and Research (A*STAR), Singapore

485 ⁸NHG Eye Institute, Tan Tock Seng Hospital, National Healthcare Group, Singapore

486 ⁹Personalised Medicine Service, Tan Tock Seng Hospital, Singapore

487 ¹⁰Singapore Eye Research Institute, Singapore National Eye Centre, Singapore

488 ¹¹Centre for Innovation and Precision Eye Health, Yong Loo Lin School of Medicine,
489 National University of Singapore, Singapore

490 ¹²Department of Cardiology, National Heart Centre Singapore, Singapore

491 ¹³Cardiovascular ACP, Duke-NUS Medical School, Singapore

492 ¹⁴SingHealth Duke-NUS Institute of Precision medicine, Singapore Health Services,
493 Singapore

494 ¹⁵Cardiovascular and Metabolic Disorders Program, Duke-NUS Medical School,
495 Singapore

496 ¹⁶Current affiliation: Translational Medicine, Sidra Medicine, Qatar

497 ¹⁷Human Development, Singapore Institute for Clinical Sciences, Singapore

498 ¹⁸Clinical Data Engagement, Bioinformatics Institute, Agency for Science,
499 Technology and Research, Singapore

500 ¹⁹Department of Biochemistry, Yong Loo Lin School of Medicine, National University
501 of Singapore, Singapore

502 ²⁰Centre for Innovation and Precision Eye Health, Yong Loo Lin School of Medicine,
503 National University of Singapore, Singapore

504 ²¹Population and Global Health, Nanyang Technological University, Lee Kong Chian
505 School of Medicine, Singapore

506 ²²Department of Epidemiology and Biostatistics, Imperial College London, UK

507 ²³Precision Health Research, Singapore

508 ²⁴Department of Medicine, Yong Loo Lin School of Medicine, National University of
509 Singapore, Singapore

510 ²⁵Saw Swee Hock School of Public Health, National University of Singapore and
511 National University Health System, Singapore
512 ²⁶Laboratory of Human Genomics, Genome Institute of Singapore (GIS), Agency for
513 Science, Technology and Research (A*STAR), Singapore
514 ²⁷Yong Loo Lin School of Medicine, National University of Singapore, Singapore
515 ²⁸Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore
516 ²⁹Cancer Science Institute of Singapore, National University of Singapore, Singapore
517 ³⁰National Heart Research Institute Singapore, National Heart Centre Singapore,
518 Singapore
519 ³¹Ophthalmology & Visual Sciences Academic Clinical Program, Duke-NUS Medical
520 School, Singapore
521 ³²Department of Psychosis, Institute of Mental Health, Singapore
522 ³³Nanyang Technological University, Singapore
523 ³⁴National Healthcare Group, Singapore
524 ³⁵Nanyang Technological University, Lee Kong Chian School of Medicine, Singapore
525 ³⁶National Cancer Centre, Singapore
526 ³⁷School of Public Health, Imperial College London, UK
527 ³⁸Department of Nutrition, Harvard T.H. Chan School of Public Health, US
528 ³⁹Exercise and Nutrition Sciences, Milken Institute School of Public Health, The
529 George Washington University, US
530 ⁴⁰Clinical Research & Innovation Office, Tan Tock Seng Hospital, Singapore,
531 ⁴¹Department of Obstetrics & Gynaecology, Yong Loo Lin School of Medicine, NUS,
532 Singapore
533 ⁴²Singapore Institute for Clinical Sciences, Singapore
534 ⁴³Liggins Institute, University of Auckland, New Zealand
535 ⁴⁴Department of Paediatrics, National University Hospital, Singapore
536 ⁴⁵Pediatrics, KK Women's and Children's Hospital, Singapore
537 ⁴⁶Department of Obstetrics & Gynaecology, KK Women's and Children's Hospital,
538 Singapore
539