## 1 Mind the gap: the relevance of the genome reference to

# 2 resolve rare and pathogenic inversions

3

23

## Kristine Bilgrav Saether<sup>1,2</sup>, Jesper Eisfeldt<sup>1, 2, 3\*</sup>, Jesse Bengtsson<sup>4</sup>, Ming Yin Lun<sup>4</sup>, 4 Christopher M. Grochowski<sup>5</sup>, Medhat Mahmoud<sup>5,6</sup>, Hsiao-Tuan Chao<sup>5,7,8,9,10,11</sup>, Jill A. 5 Rosenfeld<sup>5</sup>, Pengfei Liu<sup>5,12</sup>, Jakob Schuy<sup>1</sup>, Adam Ameur<sup>13</sup>, Undiagnosed Diseases 6 Network, James Paul Hwang<sup>6</sup>, Fritz J. Sedlazeck<sup>5,6,14</sup>, Weimin Bi<sup>5,12</sup>, Ronit Marom<sup>5,7</sup>, 7 Ann Nordgren<sup>1,3,15,16</sup>, Claudia M.B. Carvalho<sup>4#</sup>, Anna Lindstrand<sup>1,3#</sup> 8 1. Department of Molecular Medicine and Surgery, Karolinska Institute, Stockholm, 9 10 Sweden. 2. Science for Life Laboratory, Stockholm, Sweden. 11 3. Department of Clinical Genetics and Genomics, Karolinska University Hospital, 12 Stockholm, Sweden, 13 4. Pacific Northwest Research Institute, Seattle, WA, US. 14 5. Department of Molecular and Human Genetics, Baylor College of Medicine, 15 16 Houston, TX 77030, USA. 6. Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 17 18 77030, USA 7. Texas Children's Hospital, Houston, TX 77030, USA. 19 8. Cain Pediatric Neurology Research Laboratories, Jan and Dan Duncan 20 21 Neurological Research Institute, Houston, TX, USA. 9. Division of Neurology and Developmental Neuroscience, Department of 22

Pediatrics, Baylor College of Medicine, Houston, TX, USA.

24	10. Department of Neuroscience, Baylor College of Medicine, Houston, TX, USA
25	11. McNair Medical Institute, The Robert and Janice McNair Foundation, Houston,
26	TX, USA.
27	12. Baylor Genetics Laboratory, Houston, TX, USA.
28	13. Science for Life Laboratory, Department of Immunology, Genetics and
29	Pathology, Uppsala, Sweden.
30	14. Department of Computer Science, Rice University, Houston, TX, USA.
31	15. Department of Laboratory Medicine, University of Gothenburg, Gothenburg,
32	Sweden
33	16. Department of Clinical Genetics and Genomics, Sahlgrenska University Hospital,
34	Gothenburg, Sweden
35	# Shared senior authors
36	* Corresponding author:
37	Jesper Eisfeldt (jesper.eisfeldt@scilifelab.se)and Anna Lindstrand,
38	(anna.lindstrand@ki.se)

## 40 **Abstract**

41 Both long-read genome sequencing (IrGS) and the recently published Telomere to 42 Telomere (T2T) reference genome provide increased coverage and resolution across 43 repetitive regions promising heightened structural variant detection and improved 44 mapping. Inversions (INV), intrachromosomal segments which are rotated 180° and 45 inserted back into the same chromosome, are a class of structural variants particularly 46 challenging to detect due to their copy-number neutral state and association with 47 repetitive regions. Inversions represent about 1/20 of all balanced structural 48 chromosome aberrations and can lead to disease by gene disruption or altering 49 regulatory regions of dosage sensitive genes in cis. 50 Here we remapped the genome data from six individuals carrying unsolved cytogenetically detected inversions. An INV6 and INV10 were resolved using GRCh38 51 52 and T2T-CHM13. Finally, an INV9 required optical genome mapping, de novo assembly 53 of IrGS data and T2T-CHM13. This inversion disrupted intron 25 of EHMT1, confirming 54 a diagnosis of Kleefstra syndrome 1 (MIM#610253). 55 These three inversions, only mappable in specific references, prompted us to 56 investigate the presence and population frequencies of differential reference regions 57 (DRRs) between T2T-CHM13, GRCh37, GRCh38, the chimpanzee and bonobo, and 58 hundreds of megabases of DRRs were identified. 59 Our results emphasize the significance of the chosen reference genome and the added 60 benefits of IrGS and optical genome mapping in solving rearrangements in challenging 61 regions of the genome. This is particularly important for inversions and may impact clinical diagnostics. 62

- 63
- 64 Keywords: Structural variant, Inversions, long-read genome sequencing, optical
- 65 genome mapping, genetic diagnostics, formation mechanism, rare genetic diseases,
- 66 non-homologous end joining, NHEJ, satellites

67

## 69 Introduction

Inversions are defined as a copy-number neutral structural variants characterized by 70 71 two breakpoint junctions in cis each mapping to the same (paracentric inversion) or 72 distinct chromosomal arms (pericentric inversion). Inversions larger than the resolution 73 limitation of the methodology used for screening will be challenging to detect due to the 74 need of phasing both junctions *in cis*: this feature make them prone to high falsenegative and false-positive rates in genome sequencing<sup>1</sup>. Moreover, recurrent 75 76 inversions formed by non-allelic homologous recombination (NAHR) use segmental duplications (SDs) or other types of highly similar repeats as recombinant substrates <sup>2-4</sup> 77 78 which adds to the challenge of detecting junctions mapping to poor guality regions of the aenome  $^{1,5,6}$ . 79

80

We have previously shown that 28% of cytogenetically visible inversions are undetected by short read genome sequencing (srGS) <sup>7</sup>, suggesting that the breakpoint junctions are located within large stretches of repetitive sequences. Long read genome sequencing (IrGS) was shown to improve alignment and enable phasing and better resolution across repetitive regions <sup>8-10</sup>. Regardless, inversions with breakpoints mapping to large repeats remain challenging to resolve even when applying IrGS <sup>5,11</sup>.

87

Previous studies show that the new T2T-CHM13 (T2T) reference provides increased
 sensitivity in inversion detection due to increased resolution across repetitive sequences
 <sup>5,12</sup>. This reference genome has filled gaps present in earlier reference genomes, adding
 >200 Mb of sequence compared to GRCh38 <sup>12</sup>. In fact, both GRCh37 and GRCh38 lack

information across hundreds of mega-base pairs (Mb) of regions such as telomeres,
centromeres and other repetitive regions <sup>12-17</sup>. Often forgotten resources in human
genetic analysis are the closely related primate genomes chimpanzee <sup>18</sup> and bonobo <sup>19</sup>
that have been fully sequenced, with up to 99% of gaps closed <sup>19</sup>. Many sequences
unmappable after srGS analysis may be present in primates <sup>17,20</sup>.

97

We previously solved 72% (13/18) cytogenetically detected inversions using srGS and 98 GRCh37<sup>7</sup>. Here, we investigated six paracentric and pericentric large genomic 99 100 inversions (>10 Mb) detected by chromosomal karyotyping in individuals referred to 101 clinical studies, which remained unsolved after extensive genomic analysis. Using a 102 combination of srGS, IrGS and optical mapping together with remapping the GS data to 103 multiple reference genomes resolved a significant number of molecularly unsolved 104 inversions. Our results highlight a role for complex genomic regions in clinically relevant 105 structural variants with multiple breakpoint junctions in cis. Finally, we explore reference 106 genome differences using healthy Swedish individuals. Altogether, we demonstrate that 107 reference genomes have an impact on clinical structural variant calling and underscore the utility of applying long molecules to investigate the architecture of rare diseases. 108

## 109 **Results**

## 110 **Resolving inversions using IrGS and reference genomes**

111 The six rare inversions affect chromosomes 6, 9, 10, 11 and 18 in six unrelated

individuals (Table 1). They were aligned to GRCh37, GRCh38 and T2T.

114 **Table 1. Overview of included inversions.** Including the reference genome in which

they could be detected and the type of genomic sequencing technology used (short-

read (sr), linked-read (lir), long-read (lr) genome sequencing (GS) and optical genome

117 mapping (OGM)). Results are indicated as detected (2), unclear (-) or absent (X).

Case ID	Karyotype	GRCh37	GRCh38	T2T	Sequencing data
P4855_501	46, XY, inv(6)(p12q16.3)	×	?	?	srGS, lirGS, lrGS
P5371_208	46, XY, inv(9)(p13q22)	×	×	×	srGS, lirGS
P4855_208	46, XY, inv(11)(p11.1q12)mat	-	-	-	srGS, lirGS, lrGS
P4855_106	46, XY, inv(10)(q11q23)pat	×	×	?	srGS, lirGS
P5370_201	46, XY, inv(18)(p11.3q11.2)	×	×	×	srGS, lirGS
BH16643-1	46,XX,inv(9)(q12q34.3)dn	×	×	?	srGS, IrGS, OGM

118

Two (P4855\_501, P4855\_106) inversions were detected by standard variant callers
after realigning the srGS to a new reference genome. One additional case (BH16643-1)
was resolved using a combination of new reference genomes, de novo assembly and
optical genome mapping (OGM).

123

124 The first case (P4855\_501), with a pericentric inversion on chromosome 6 initially 125 detected by karyotyping, was undetected using both srGS, lirGS and IrGS SV analysis 126 as well as *de novo* assembly in GRCh37<sup>7</sup>. Following realignment of the GS data to 127 GRCh38 and T2T, the exact inversion breakpoints were present in both srGS, lirGS, 128 IrGS as well as in the *de novo* assembly (Fig. 1A & 3, Supplementary Fig. 1). The 129 GRCh38 analysis pinpointed the breakpoint junction at 6p12 to position chr6:51190256. 130 whilst the 6q16.3 breakpoint was specified to 6q16.1; chr6:93164914. The GRCh38 131 analysis pinpointed the breakpoint junction at 6p12 to position chr6:51190256, whilst the 132 other 6q16.3 breakpoint was specified to 6q16.1: chr6:93164914. Detailed breakpoint

sequences reveal presence of 3 bp microhomology (Supplementary Fig. 2). The
inverted segment covered 42 Mb (24% of chromosome 6). Discordant reads pairs were
present in GRCh37 at the 6q16.1 breakpoint, partnering with multiple genomic locations
(Supplementary Fig. 1). The affected patient suffered hearing impairment, intellectual
disability, autistic features, diplopia, anosmia as well as hypogonadism. No genes were
interrupted by the inversion breakpoints, while 324 genes were located within the
inverted segment.



142	Fig 1: Detection of unresolved inversions by short and long read genome
143	sequencing. A) An inversion 6 visible in short read genome sequencing (srGS), linked
144	read genome sequencing (lirGS) and long read genome sequencing (lrGS) using
145	GRCh38. <b>B)</b> An inversion 10 visible in srGS and lirGS data using T2T.
146	C) An inversion 9 only visibly by IrGS de novo assembly using T2T. D) An inversion 11
147	with centromeric breakpoints. The IrGS de novo assembly call fitting with the
148	cytogenetic analysis is shown.
149	
150	The second inversion (P4855_106), affecting chromosome 10 in a healthy individual,
151	could only be resolved using T2T (Fig. 1B, Supplementary Fig. 3), where it was visible
152	in both srGS, lirGS and <i>de novo</i> lirGS assembly. The 10q11.21 breakpoint was
153	pinpointed to chr10:42292350, whilst the 10q23.32 breakpoint was pinpointed to
154	chr10:93143588 (Supplementary Fig. 2). The inverted segment covered 50.9 Mb (40%
155	of chromosome 10). The inversion interrupts intron 1 in the gene CPEB3, however
156	disruption of this transcript is unlikely to be pathogenic. 2879 genes were located within
157	the inverted segment.
158	
159	The third inversion (BH16643-1), affecting chromosome 9, was first identified by
160	chromosomal karyotyping in an individual with global developmental delay, hypotonia,
161	feeding difficulties, congenital heart defect and dysmorphic craniofacial features (Fig. 1
162	& 2, Supplementary Fig. 5, Supplementary Information 2). The inversion was

163 undetected in srGS, IrGS using GRCh37. Manual inspection of the OGM data indicated

164 a structural variant breakpoint junction at 9q34.3 supported by raw molecules and was

165	narrowed down to 150.05 – 150.1 Mb using T2T OGM de novo assembly
166	(Supplementary Fig. 4). Lack of informative motifs in the raw molecules hampered our
167	ability to find the location of breakpoint at 9q12. Using T2T, OGM, IrGS and de novo
168	assembly, we were able to pinpoint the 9q34.3 breakpoint to chr9:150,079,673. The
169	9q12 breakpoint was located in a 28 Mb region (chr9:48424795-77056693) consisting of
170	satellite and simple repeats not represented in reference genomes GRCh37 and
171	GRCh38. Due to this, the 9q12 breakpoint is ambiguously aligned in both OGM, IrGS
172	and <i>de novo</i> assembly contigs (Fig 2D). The inverted segment covers ~95 Mb (63% of
173	chromosome 9). The 9q34.3 breakpoint interrupts intron 25 of the gene EHMT1,
174	haploinsufficiency of which causes Kleefstra syndrome 1 (MIM#610253), a diagnosis
175	fitting the clinical phenotype (Supplementary Information 2).

176



178 Fig. 2: Inversion affecting chromosome 9. A) Pedigree displaying inheritance pattern

179 for inversion 9. B) G-banded chromosome analysis showed a paracentric inversion in

180 the long arm of one chromosome 9 between bands 9q12 and 9q34.3 in the proband. 181 The abnormal chromosome 9 is indicated by a blue arrow. Parental chromosome 182 analysis revealed no evidence of this inversion in either parent, suggesting that this is a 183 de novo event. C) Chromosome 9 inversion disrupted intron 25 out of 26 of 184 EHMT1/NM\_024757 at 9q34.3. D) Nucleotide sequence alignment of inversion 185 breakpoint junctions 1 (top) and 2 (bottom). 186 187 For one inversion, affecting chromosome 11 (P4855\_208) identified by karyotyping in a 188 patient suffering from neurodevelopmental delay (Table 1, Figure 1D) a potential 189 inversion call was suggested after IrGS de novo assembly. The suggested breakpoints 190 were present in all the assessed genomes (GRCh37, GRCh38 and T2T) and one of 191 them was verified using breakpoint PCR and Sanger sequencing (Supplementary Fig. 192 5). However, both breakpoints were in highly repetitive regions consisting of centromeric 193 satellite repeats and similar signals were present in unrelated controls rendering it 194 uncertain whether the true inversion breakpoints were detected or not (Supplementary 195 Fig. 6).

196

## 197 Reference genomes influence SV discovery

Reference genome analysis revealed that inversion breakpoint sequences were missing in reference genomes GRCh37 (inv6) and GRCh38 (inv9, inv10), making it impossible to solve them using these reference genomes (Fig. 3, Supplementary Fig. 1 and 3). In total 127 kb of sequence at 6p12.3 was present in GRCh38 but missing from GRCh37. The region, located at chr6:51102785-51230413 did not contain any known genes, and

203 consists of 51% repeat sequence, of which 49% interspersed repeats and 2% simple

repeats (Fig. 5C). The sequence aligned correctly in T2T, chimpanzee and bonobo,

205 concluding that the inversion was in fact visible in srGS except for when using GRCh37

- 206 (Supplementary Fig. 1 and 2).
- 207





in a 127kb region in GRCh38 missing from GRCh37. B) The chromosome 6p inversion
breakpoint in GRCh38 and T2T. C) The chromosome 10q breakpoint is located in a
69kb region missing in GRCh38, with a surrounding 4kb duplication which occurs only
once in T2T. D) The chromosome 9q12 breakpoint is located in a 28 Mb region missing
in GRCh38 shaded in blue.

218

For the inversion on chromosome 10, the 10:q11 breakpoint was located in a 69 kb

region of simple repeats only present in T2T (Fig. 3, 4, Supplementary Fig. 1 and 3).

The region, spanning from 10:42282056-42351085, consists of 99% simple repeats and

is surrounded by other regions of simple repeats. It does not contain any known genes.

223 The 9q12 breakpoint of inversion on chromosome 9 was located in a 28 Mb region of

224 79% simple and 19% satellite repeats which was not present in GRCh37, GRCh38,

bonobo or chimpanzee (Fig. 3, 4).

226

#### 227 Comparing variable sequences in human and primate reference genomes

228 Next, we evaluate the abundance of such Differential Reference Regions (DRRs), i.e. a 229 sequence larger than 10kb that is present in one reference and missing in another 230 during pairwise comparison. We compared three human (GRCh37, GRCh38, T2T) and 231 two primate (Chimpanzee and Bonobo) reference genomes pairwise. In comparing the 232 human references to each other, the longest uninterrupted DRR was detected in 233 GRCh38-GRCh37 (10kb-47Mb, median 50kb), whilst the most fragmented DRRs were detected in T2T-GRCh38 (10kb-34Mb, median 30kb). The chimpanzee-T2T (range 234 235 10kb-14Mb, median 40kb) and bonobo-T2T had similar ranges of DRRs (range 10kb-

- 19Mb and median 35kb) (Supplementary Table 2). In total, we uncovered 203 regions
- and 260.6 Mb present in T2T and missing from GRCh37 (T2T-GRCh37). Finally, T2T-
- 238 GRCh37 contains the highest total Mb of DRR (Table 1).
- 239
- 240 **Table 1: Differential reference regions between reference genomes.** For each
- template on the top row (grey) the total amount of sequence in Mb and on the second

row (white) the total number of DRRs is given in comparison with the query reference.

			Qı	uery		
	DRR (Mb)	GRCh37	GRCh38	T2T	Chimpanzee	Bonobo
	GRCh37	0	8.5	12.6	52.5	59.7
Template		0	84	130	686	717
	GRCh38	81.03	0	39.44	117.8	125.5
		340	0	814	870	885
	тот	260.6	216.9	0	289.36	295.3
	121	203	687	0	878	922
	Chimp	333.2	325.1	315.1	0	263.9
	Chinip	743	827	845	0	849
	Bonobo	408.9	400.7	392.48	336.1	0
	BOHODO	829	877	855	992	0

243

When comparing all DRRs where a sequence was present in T2T and missing from the query genome (T2T DRRs), we observe clustering of DRRs located in centromeric and telomeric regions as well as segmental duplications, the acrocentric p-arms and chr Y (Fig. 4). Of all T2T DRRs, 200 Mb of sequence was missing from all query reference

- genomes (Fig. 4C-D). For all GRCh38 DRRs, only 33 Mb of sequence was missing in
- all query reference genomes including T2T (Fig. 4B, Supplementary Fig. 7).



Fig 4: Shared DRR in T2T and GRCh38. A) Bar plot of all T2T DRRs B) Venn diagram
of Mb overlap between all GRCh38 DRRs, C) Venn diagram of Mb overlap between all
T2T DRRs.

254

## 255 DRRs introduce repetitive sequences

- 256 Repeat analysis of all DRRs in GRCh38-GRCh37 and T2T-GRCh38 reveal most to be
- repeat regions, and ~10% to be unique sequence (Fig. 5A). As an example, the 127 kb
- 258 DRR affected by the inversion on chromosome 6 consisted of 49% unmasked









- 270 GRCh37 and T2T-GRCh38. C) Pie chart displaying repeat content in the GRCh38-
- 271 GRCh37 DRR sequences affected by the inversion 6 at the 6p12 junction
- 272 (chr6:51102785-51230413) in GRCh38.
- 273

## 274 DRR sequences in the general population

- Next, we aligned srGS data from 100 Swedish individuals <sup>21</sup> to the five references and
- assessed the presence of DRR across the population (Fig 6, Supplementary Fig. 9 and
- 277 10).
- 278







281 GRCh38 and their presence in Swedish individuals. Blue indicating absent (<8X) and

red present (>8X and <100X). B) Violin plot of coverage across the respective DRRs in</li>
100 Swedish individuals. C) The distribution of population frequencies of the detected
GRCh38-GRCh37 (green) and T2T-GRCh38 (orange) DRRs.

285

286 Of the DRRs in T2T-GRCh38, 68% were classified as not detected, meaning that the 287 average coverage per individual was blow the cutoff of 8X (Fig. 6A-C; Supplementary 288 Table 1). Of the 32% that were detected, 42% were observed in <5% (rare), 58% were 289 found in >5% (common) and 30% in over 90% (Fig. 6C). Across the 100 individuals, an 290 average of 1.3% of reads spanning DRRs were multimapping reads, meaning they map 291 to several locations in the genome. We also assessed the mapping quality of reads from 292 5 individuals across DRRs, (Supplementary Fig. 8) where 20% of reads had a mapping 293 quality above 20.

294

In comparison, for the GRCh38-GRCh37 DRRs, 60% were not detected

(Supplementary Fig. 10), whilst of the 40% detected, 13% were rare and 86% common
(Fig. 6C). Furthermore, 53% were found in over 90% of the queried individuals. The
mapping quality of reads from 5 individuals across DRRs was assessed, where 25% of
reads had a mapping quality above 20 (Supplementary Fig. 8). The violin plot confirms
that most DRRs lack aligned reads (Fig. 6B-C).

301

In Chimpanzee-T2T DRRs 80% were missing in the Swedish individuals, whilst of the
20% present, 6% were rare, 93% common and 75% in >90% of the individuals.

Of the Bonobo-T2T DRRs, 93% were not detected, whilst of the remaining 7%, 15%
were rare 85% common and 50% in >90% of the Swedish individuals (Supplementary
Fig. 9).

#### 307 Discussion

308 The availability of long read sequencing and the new reference T2T-CHM13

incentivized us to revisit six previously unsolved cytogenetically visible inversions. Three

inversions were solved by realigning the srGS data to GRCh38 and/or T2T. This

311 illustrates how reference genome variability may influence the accuracy of clinical

312 diagnostic SV calling and that IrGS in itself is not the sole answer.

313

For inversion 9, a *de novo* assembly proved necessary to pinpoint the breakpoints.

315 Highly repetitive DNA, LINE1 elements and centromeric sequences were involved in the

316 breakpoints. Resolving inversions with this level of repeat is challenging with srGS.

However, the inversion on chromosome 10 was resolved using srGS even though it was

318 located in a region consisting of 99% repeats, highlighting that detection of a true

319 positive SV call is highly dependent on completeness of the reference as well as the

320 representation of normal variation, even when using srGS. This is important from a

321 clinical perspective, where IrGS, which improves resolution of repeats, is not yet broadly

available clinically. Two inversions affecting chromosomes 9 and 18 remain unresolved,

323 both with at least one breakpoints positioned in or close to centromeric regions. For

inversion 11, IrGS nor *de* novo assembly was sufficient at pinpointing the breakpoints.

325 For the two remaining, IrGS was unfortunately not possible. Unfortunately, IrGS was not

326 possible for these cases.

327

328	For inversion 11, the IrGS de novo assembly generated a call matching the cytogenetic
329	findings (Fig. 1D). One breakpoint was verified by breakpoint PCR and Sanger
330	sequencing (Supplementary Fig. 5). However, further analysis revealed that similar
331	patters were identified in other individuals, this call may therefore represent normal
332	variation, or the true inversion call, formed through NAHR (Supplementary Fig. 6).
333	Overall, the case is an example of the challenge of pinpointing and resolving
334	breakpoints involving centromeric regions and indicates a need for new standards for
335	validating IrGS findings, as well as large scale population genomics databases for
336	filtering common variation. As supported by results from inversion 9, OGM can provide
337	further resolution in these examples.

338

We and others have previously suggested that a part of cytogenetically visible 339 inversions may have been formed through non-allelic homologous recombination 340 (NAHR) explaining why some remain undetected even after srGS analysis <sup>2,7 3</sup>. 341 342 Nonetheless, the breakpoint junction analysis of the three inversions resolved here 343 shows a distinct picture where none of the unsolved ones were mediated by ectopic 344 recombination between paralogous sequences. No matching repeats are detected, and 345 the junctions contain blunt ends or microhomology without additional copy-number variants or other concomitantly alterations suggesting canonical non-homologous end 346 joining (c-NHEJ) as the underlying mechanism of formation <sup>7,22</sup>. Even so, the inversion 347 348 breakpoint DRRs on chromosomes 6p12.3, 9q12 and 10q11 are highly repetitive. The 349 127 kb DRR on 6p contained 51% simple repeats, the 28 Mb DRR on 9g consisted of

99% satellite and simple repeats and the 69 kb DRR on 10q consisted of 99% simple
repeats. This result supports that copy-number neutral inversions, similarly to balanced
translocations, may result from an error prone repair of processed double-strand breaks
(DSBs) <sup>23</sup>.

354

355 One inversion revealed a breakpoint disrupting *EHMT1* likely leading to loss of function 356 of the gene, consistent with the expected underlying biological mechanism for Kleefstra 357 syndrome 1. The clinical phenotype of the individual that includes hypotonia and global 358 developmental delay, congenital heart defect, recurrent respiratory infections and visual 359 impairment is also consistent with the syndrome. The individual presented with 360 dysmorphic features, including midface retrusion, everted lower lip and prognathism, 361 that fit the Kleefstra syndrome's characteristic facial gestalt. Recently, we have reported 362 a patient with multiple paracentric and pericentric copy-neutral inversions affecting 363 chromosome 6 that disrupted ARID1B in an individual with neurodevelopmental phenotype<sup>22</sup>. All together, these results underscore the relevance of inversions to 364 365 unsolved rare disease, often undetected by current clinical genome sequencing. 366

We proceeded to evaluate DRRs differing between reference genomes across GRCh37, GRCh38, T2T, chimpanzee and bonobo (Fig. 4). Our results (216 Mb and 260 Mb DRR in T2T compared to GRCh38 and GRCh37 respectively) are comparable to previous work showing that T2T introduce >200 Mb compared to GRCh38<sup>12,14</sup>. T2T has the highest amount of DRR (200 Mb) sequence not present in any of the other human or primate references, indicating that the T2T reference is more complete. Assessing

the repeat content in all sequences, we find around 10% of DRR sequences in GRCh38
and T2T to be unique and the remaining to be repeats, where satellite repeats is the
major contributor. Interestingly, T2T add around 20% simple repeats (Fig. 5).

377 Next, we analyzed the variability of DRR sequences in 100 healthy Swedish individuals 378 (Fig. 6). Of note, most of the srGS sequences that align to T2T DRRs have a very low 379 mapping quality (20% with a mapping quality >20) indicating that short read technology 380 is not the best option for analyzing these regions. This is likely due to those DRRs 381 mainly consist of repeat and satellite sequences resulting in ambiguous alignment of short reads, in addition to these regions being highly variable between individuals<sup>24</sup> 382 383 (Fig. 6 and Supplementary Fig. 7). Regardless, some DRR sequences are present in 384 most individuals (32% of T2T-GRCh38 and 40% of GRCh38-GRCh37 DRRs are found 385 at >8X in >50% of the Swedish individuals).

386

387 Although we now have an almost 100% fully resolved human reference genome, no 388 single genome can represent the full genetic diversity in humans. To address these 389 shortcomings, the pangenome consortium made a reference genome representing 47 diploid assemblies represented as a graph <sup>25</sup>. This assembly is able to represent large 390 391 genomic variation, complex loci and increased number of SVs per haplotype. With 392 future refinement and aspects of including >700 haplotypes, providing a better 393 representation of the human genome, which provides better alignment and variant calling. 394

395

396 In conclusion, we show that for solving rearrangements in variable genomic regions, the 397 success rate highly depends on the reference genome and its completeness, and novel 398 IrGS databases and verification methods are needed. To fully understand the IrGS 399 findings and be able to offer digital karyotyping as a first line test we need to understand 400 the limits of the analysis. Furthermore, our results highlight that to improve clinical 401 genomic analysis genomic diversity needs to be considered. The available human and primate genomes are a useful resource to improve our understanding of repetitive and 402 403 complex regions which have previously been understudied.

404

## 405 **Methods**

#### 406 Study participants

IDs used in this article are not known to anyone outside the research group. Subjects
carrying Inversions 6, 10, 11 and 18 were enrolled at Karolinska University Hospital,
Stockholm, Sweden<sup>7</sup>. Patient BH16643-1 was enrolled using research protocol H47281/Pacific Northwest Research Institute WIRB #20202158 and 15HG0130 with the
NIH IRB as part of the Undiagnosed Diseases Network (UDN). Whole blood samples
(3-10mL) were collected from the patient and parents. DNA was isolated from blood
according to standard procedures.

414

The SweGen dataset (n=1000)<sup>21</sup>, consists of 1000 unrelated Swedish individuals
representing the genetic variation in the Swedish population. In brief, the individuals
were selected from the Swedish Twin Registry, a nationwide cohort of 10,000 Swedish-

born individuals. The samples were sequenced using Illumina short-read sequencing to
an average of 30X coverage. From these, we selected 100 random, unrelated samples
for further use in this study.

421

## 422 Genome sequencing

For samples (P4855\_501, P5371\_208, P4855\_208, P4855\_106, P5370\_106) srGS and
10X genomics linked read sequencing of the samples was performed at the national
genomics infrastructure (NGI) at Science for Life laboratory (SciLifeLab) Stockholm as
previously mentioned <sup>7</sup>. Analysis for structural variants was performed using FindSV as
described previously <sup>7</sup>.

428

IrGS was performed on P4855\_501 and P4855\_208 using Pacific Biosciences (PacBio)
Sequel II (NGI SciLifeLab Uppsala, Sweden).

431

432 For the BH16643 family, short-read trio genome sequencing was performed at the

433 Baylor College of Medicine Human Genome Sequencing Center (HGSC) with KAPA

434 Hyper PCR-free reagents on the NovaSeq 6000 to an average of 40X coverage. Post-

435 sequencing data analysis was performed using the HGSC HgV analysis pipeline  $^{26}$ .

436 IrGS of the proband (BH16643-1) was done on the PacBio Sequel II instrument using

437 two SMRTcells.

438

439 Genome analysis

440 The srGS data was aligned to reference genomes GRCh37, GRCh38, T2T,

- 441 Chimpanzee and Bonobo using BWA mem for the srGS.
- 442 bwa mem -p -t 16 <ref> <fastq>
- 443 The lirGS was aligned using:
- 444 longranger wgs --id <id> --reference <ref> --fastq <fastq> --
- 445 vcmode freebayes
- 446
- 447 The IrGS data was aligned to GRCH37, GRCh38 and T2T. Analysis of was done using
- 448 an in house developed pipeline LOMPE (<u>https://github.com/kristinebilgrav/LOMPE</u>).
- LOMPE uses minimap2 for alignment and combines Sniffles (v1) <sup>27</sup> and CNVpytor <sup>28</sup> for
- 450 structural variant calling, and produces a single VCF file which is annotated using
- 451 Variant Effect Predictor (VEP)<sup>29</sup>. The resulting IrGS data had a read depth of 10 (inv
- 452 11), 13 (inv 6) and 26X (inv9) and an average read length of 18kb.
- 453

#### 454 *De novo* assembly

- 455 De novo assembly on IrGS from samples P4855\_501, P4855\_208 and BH16643-1 was
- 456 performed using hifiasm <sup>30</sup>. Quality control was performed using quast <sup>31</sup>. Alignment to
- <sup>457</sup> reference genomes GRCh37, GRCh38 and T2T was performed using minimap2 <sup>32</sup>, and
- 458 variant calling was performed using sniffles (v1) <sup>27</sup> and htsbox
- 459 (https://github.com/lh3/htsbox). On lirGS from sample P4855\_106 a *de novo* assembly
- 460 was performed using 10X Genomics Supernova  $^{33}$ .
- 461
- 462 **Optical genome mapping**

Optical genome mapping was performed as described previously <sup>34</sup>. Briefly, ultra-high 463 464 molecular weight (UHMW) genomic DNA for use in genomic optical mapping was extracted from blood using Bionano Prep<sup>TM</sup> Blood and Cell Culture DNA Isolation Kit 465 466 (Bionano Genomics) with an input of 1.5 million cells. Subsequent DNA quantity and size was confirmed using Qubit<sup>™</sup> dsDNA BR Assay Kit. A total of 0.75 µg of HMW DNA 467 468 was then labeled by DLE-1 using the Bionano Prep direct label and stain (DLS) method 469 (Bionano Genomics) and loaded onto a flow cell to run on the Saphyr optical mapping 470 system (Bionano Genomics). Raw optical mapping molecules in the form of BNX files 471 were run through a preliminary bioinformatic pipeline that filtered out molecules less 472 than 150 kb in size and with less than 9 motifs per molecule to generate a *de novo* 473 assembly of the genome maps. The data collected provided 1637 Gb of data greater 474 than 150 kb, with at least 9 labels per molecule. Data was then aligned to an *in-silico* 475 reference genome (GRCh37, GRCh38, and T2T-CHM13) using the Bionano Solve v3.7 476 RefAligner module. Structural variant calls were generated through comparison of the 477 reference genome using a custom Bionano SV caller. Manual inspection of proposed 478 breakpoint junctions were then visualized in the Bionano Access software program 479 v1.7.2.

480

## 481 Breakpoint verification by Sanger sequencing

Breakpoint verification of breakpoints identified in P4855\_208 was performed as
previously described<sup>7</sup>.

484

#### 485 **Reference genome analysis**

486	Reference genomes GRCh37 (GCF_000001405.13), GRCh38 (GCF_000001405.26),
487	T2T-CHM13 (v2.0, GCF_009914755.1), bonobo (GCF_029289425.1) and chimpanzee
488	(GCF_028858775.1) were downloaded from National Center for Biotechnology
489	Information (NCBI) <sup>35</sup> . Alternative sequences were excluded for all reference genomes.
490	They were aligned to one another using minimap2 using the settings for cross-species
491	full genome alignment and overlap between long reads (2.24-r1122) <sup>32</sup> . This enables
492	sequence comparison between the two reference genomes.
493	minimap2 -cx asm5 template.fa query.fa > aln.paf
494	
495	minimap2 -ax asm5 template.fa query.fa   samtools view -Sbh -
496	samtools sort -m 4G -@1 - > aln.bam
497	samtools index aln.bam
498	
499	Coverage analysis of the resulting pairwise compared reference genomes was
500	performed using TIDDIT v.3.6.0 $^{36}$ , producing a bed file with gaps between the query
501	and template. Files with known gap regions were downloaded from UCSC
502	TableBrowser <sup>37</sup> and these regions were excluded from the coverage analysis. A
503	differential reference region (DRR) was identified as a region of template genome which
504	was not covered by the query genome.
505	
506	Differential reference regions in SweGen
507	100 SweGen samples were aligned to each of the 5 reference genomes and coverage
508	analysis across the genome was performed as described above. Coverage across
509	DRRs identified above was extracted. A DRR was considered present in SweGen if the

510	coverage across the DRR >8X and <100X, and absent if the coverage was <8X.
511	Regions with coverage >100X were not considered. The thresholds were set based on
512	coverage experienced to support the presence of one or multiple genomic copies
513	(Supplementary table 1). On a populational level, a DRR was considered common if it
514	was present in >5% of the population and absent if none had it.
515	
516	For the VENN diagrams, a DRR was considered overlapping if the region was missing
517	in all query genomes, but present in the template genome.
518	
519	Multimapping reads were identified by extracting the number of times a read was
520	aligned in the bam file. Mapping quality was assessed by extracting the mapping quality
521	of all reads in the bam file.
521 522	of all reads in the bam file.
521 522 523	of all reads in the bam file. Data access
521 522 523 524	of all reads in the bam file. Data access The reference genomes can be downloaded from <i>NCBI</i> <sup>35</sup> . The clinical samples are not
521 522 523 524 525	of all reads in the bam file. <b>Data access</b> The reference genomes can be downloaded from <i>NCBI</i> <sup>35</sup> . The clinical samples are not available due to ethical permissions. The SweGen dataset is available at
521 522 523 524 525 526	of all reads in the bam file. <b>Data access</b> The reference genomes can be downloaded from <i>NCBI</i> <sup>35</sup> . The clinical samples are not available due to ethical permissions. The SweGen dataset is available at <u>https://swefreq.nbis.se/</u> upon signing a data agreement. The srGS analysis pipeline
521 522 523 524 525 526 527	of all reads in the bam file. <b>Data access</b> The reference genomes can be downloaded from <i>NCBI</i> <sup>35</sup> . The clinical samples are not available due to ethical permissions. The SweGen dataset is available at <u>https://swefreq.nbis.se/</u> upon signing a data agreement. The srGS analysis pipeline FindSV is available on GitHub at <u>https://github.com/J35P312/FindSV</u> . The IrGS analysis
521 522 523 524 525 526 527 528	of all reads in the bam file. <b>Data access</b> The reference genomes can be downloaded from <i>NCBI</i> <sup>35</sup> . The clinical samples are not available due to ethical permissions. The SweGen dataset is available at <u>https://swefreq.nbis.se/</u> upon signing a data agreement. The srGS analysis pipeline FindSV is available on GitHub at <u>https://github.com/J35P312/FindSV</u> . The IrGS analysis pipeline LOMPE is available at <u>https://github.com/kristinebilgrav/LOMPE</u> . Tools TIDDIT
521 522 523 524 525 526 527 528 529	of all reads in the bam file. <b>Data access</b> The reference genomes can be downloaded from <i>NCBI</i> <sup>35</sup> . The clinical samples are not available due to ethical permissions. The SweGen dataset is available at <u>https://swefreq.nbis.se/</u> upon signing a data agreement. The srGS analysis pipeline FindSV is available on GitHub at <u>https://github.com/J35P312/FindSV</u> . The IrGS analysis pipeline LOMPE is available at <u>https://github.com/kristinebilgrav/LOMPE</u> . Tools TIDDIT and SVDB are available at <u>https://github.com/SciLifeLab/TIDDIT</u> and
521 522 523 524 525 526 527 528 529 530	of all reads in the bam file. <b>Data access</b> The reference genomes can be downloaded from <i>NCBI</i> <sup>35</sup> . The clinical samples are not available due to ethical permissions. The SweGen dataset is available at <u>https://swefreq.nbis.se/</u> upon signing a data agreement. The srGS analysis pipeline FindSV is available on GitHub at <u>https://github.com/J35P312/FindSV</u> . The IrGS analysis pipeline LOMPE is available at <u>https://github.com/kristinebilgrav/LOMPE</u> . Tools TIDDIT and SVDB are available at <u>https://github.com/SciLifeLab/TIDDIT</u> and <u>https://github.com/J35P312/SVDB</u> .

532 Ethics statement

Ethics approval for analysis of participant samples was given by the Regional Ethical 533 534 Review Board in Stockholm, Sweden (ethics permit numbers 2012/222-31/3). This 535 ethics permit allows for use of clinical samples for analysis of scientific importance as 536 part of clinical development. The IRB approval does not require us to get written 537 consent for clinical testing. The research conformed to the principles of the Helsinki 538 Declaration. Patient BH16643-1 was enrolled using research protocol H-47281/Pacific 539 Northwest Research Institute WIRB #20202158 and 15HG0130 with the NIH IRB as 540 part of the Undiagnosed Diseases Network (UDN). Informed consent was obtained from 541 the legal guardians. 542 543 **Competing interest statement** 544 AL has received honoraria from Illumina and PacBio. The Department of Molecular and 545 Human Genetics at Baylor College of Medicine receives revenue from clinical genetic 546 testing conducted at Baylor Genetics Laboratories. The remaining authors have nothing 547 to declare. 548 549 Funding 550 Research reported in this publication was supported by the Swedish Research Council 551 2019-02078, the Swedish Brain Fund FO2022-0256, the Stockholm Regional Council 552 ALF funding, the Swedish Rare Diseases Research foundation (AL) and the National

- 553 Institute of General Medical Sciences NIGMS R01 GM132589 (CMBC). Additional
- support was provided through the National Institute of Neurological Disorders and
- 555 Stroke of the National Institutes of Health (U01HG007709 and U01HG007942) and the

National Institute of Health (NIH S10 1S10OD028587). The content is solely the
responsibility of the authors and does not necessarily represent the official views of the
National Institutes of Health. The funders had no role in study design, data collection
and analysis, decision to publish, or preparation of the manuscript.

560

## 561 Acknowledgements

562 The authors thank the families and individuals enrolled in this research. Gratefully

563 acknowledge the support from the National Genomics Infrastructure (NGI) Stockholm at

564 Science for Life Laboratory in providing assistance in massive parallel sequencing.

565 Thank you to Davut Pehlivan, from Baylor College of Medicine and Texas Children's

566 Hospital for helping with patient enrollment. The computations were performed on

resources provided by SNIC through Uppsala Multidisciplinary Center for Advanced

568 Computational Science (UPPMAX) under Project SNIC sens2017106 and

sens2020021. Two of the authors of this publication are members of the European

570 Reference Network on Rare Congenital Malformations and Rare Intellectual Disability

571 ERN-ITHACA [EU Framework Partnership Agreement ID: 3HP-HP-FPA ERN-01-

572 2016/739516].

573

## 574 **Author contributions:**

- 575 Study design: KBS, JE, AL, CMBC
- 576 Clinical information: HTC, LB, JAR, RM
- 577 Bench work experiments: JB, CMG, JPH, WB, JS
- 578 Bioinformatic analysis: KBS, JE, AA, MYL, FS,

579 Manuscript: KBS, JS, JE, AN, AL, CMBC

580 Figures, tables, visualizations: KBS

581 Supervision of the manuscript process: JE, AL, CMBC, AA

582

583

## 584 **References**

585 1. Chaisson MJP, Sanders AD, Zhao X, et al. Multi-platform discovery of haplotype-586 resolved structural variation in human genomes. Nat Commun. Apr 16 2019;10(1):1784. 587 doi:10.1038/s41467-018-08148-z 588 Stankiewicz P, Lupski JR. Genome architecture, rearrangements and genomic 2. 589 disorders. Trends Genet. Feb 2002;18(2):74-82. doi:10.1016/s0168-9525(02)02592-1 590 Carvalho CM, Lupski JR. Mechanisms underlying structural variant formation in genomic 3. 591 disorders. Nat Rev Genet. Apr 2016;17(4):224-38. doi:10.1038/nrg.2015.25 592 Dittwald P, Gambin T, Gonzaga-Jauregui C, et al. Inverted low-copy repeats and 4. 593 genome instability--a genome-wide analysis. Hum Mutat. Jan 2013;34(1):210-20. 594 doi:10.1002/humu.22217 595 Porubsky D, Harvey WT, Rozanski AN, et al. Inversion polymorphism in a complete 5. 596 human genome assembly. Genome Biol. Apr 30 2023;24(1):100. doi:10.1186/s13059-023-597 02919-8 598 Kidd JM, Graves T, Newman TL, et al. A human genome structural variation sequencing 6. 599 resource reveals insights into mutational mechanisms. Cell. Nov 24 2010;143(5):837-47. 600 doi:10.1016/j.cell.2010.10.027 601 Pettersson M, Grochowski CM, Wincent J, et al. Cytogenetically visible inversions are 7. 602 formed by multiple molecular mechanisms. Hum Mutat. Nov 2020;41(11):1979-1998. 603 doi:10.1002/humu.24106 604 8. Kronenberg ZN, Fiddes IT, Gordon D, et al. High-resolution comparative analysis of 605 great ape genomes. Science. Jun 8 2018;360(6393)doi:10.1126/science.aar6343 606 Logsdon GA, Vollger MR, Eichler EE. Long-read human genome sequencing and its 9. 607 applications. Nat Rev Genet. Oct 2020;21(10):597-614. doi:10.1038/s41576-020-0236-x 608 10. Sudmant PH, Rausch T, Gardner EJ, et al. An integrated map of structural variation in 609 2,504 human genomes. Nature. 2015/10/01/ 2015;526(7571):75-81. doi:10.1038/nature15394 610 11. Porubsky D, Vollger MR, Harvey WT, et al. Gaps and complex structurally variant loci in phased genome assemblies. Genome Res. Apr 2023;33(4):496-510. 611 doi:10.1101/gr.277334.122 612 613 Nurk S, Koren S, Rhie A, et al. The complete sequence of a human genome. Science. 12. 614 Apr 2022;376(6588):44-53. doi:10.1126/science.abj6987 615 Pan B, Kusko R, Xiao W, et al. Similarities and differences between variants called with 13. 616 human reference genome HG19 or HG38. BMC Bioinformatics. Mar 14 2019;20(Suppl 2):101. 617 doi:10.1186/s12859-019-2620-0 618 Schneider VA, Graves-Lindsay T, Howe K, et al. Evaluation of GRCh38 and de novo 14. 619 haploid genome assemblies demonstrates the enduring quality of the reference assembly. 620 Genome Res. May 2017;27(5):849-864. doi:10.1101/gr.213611.116

621 15. Church DM, Schneider VA, Graves T, et al. Modernizing reference genome assemblies. 622 PLoS Biol. Jul 2011;9(7):e1001091. doi:10.1371/journal.pbio.1001091 623 16. Ameur A, Che H, Martin M, et al. De Novo Assembly of Two Swedish Genomes Reveals 624 Missing Segments from the Human GRCh38 Reference and Improves Variant Calling of 625 Population-Scale Sequencing Data. Genes (Basel). Oct 9 626 2018;9(10)doi:10.3390/genes9100486 627 Eisfeldt J, Martensson G, Ameur A, Nilsson D, Lindstrand A. Discovery of Novel 17. 628 Sequences in 1,000 Swedish Genomes. Mol Biol Evol. Jan 1 2020;37(1):18-30. 629 doi:10.1093/molbev/msz176 630 Chimpanzee S, Analysis C. Initial sequence of the chimpanzee genome and comparison 18. 631 with the human genome. Nature. Sep 1 2005;437(7055):69-87. doi:10.1038/nature04072 632 Mao Y, Catacchio CR, Hillier LW, et al. A high-quality bonobo genome refines the 19. 633 analysis of hominid evolution. Nature. Jun 2021;594(7861):77-81. doi:10.1038/s41586-021-634 03519-x 635 20. Sherman RM, Forman J, Antonescu V, et al. Assembly of a pan-genome from deep 636 sequencing of 910 humans of African descent. Nat Genet. Jan 2019;51(1):30-35. 637 doi:10.1038/s41588-018-0273-y 638 Ameur A, Dahlberg J, Olason P, et al. SweGen: a whole-genome data resource of 21. 639 genetic variability in a cross-section of the Swedish population. Eur J Hum Genet. 2017/11// 640 2017;25(11):1253-1260. doi:10.1038/ejhg.2017.130 641 Grochowski CM, Krepischi ACV, Eisfeldt J, et al. Chromoanagenesis Event Underlies a 22. 642 de novo Pericentric and Multiple Paracentric Inversions in a Single Chromosome Causing 643 Coffin-Siris Syndrome. Front Genet. 2021;12:708348. doi:10.3389/fgene.2021.708348 644 Nilsson D, Pettersson M, Gustavsson P, et al. Whole-Genome Sequencing of 23. 645 Cytogenetically Balanced Chromosome Translocations Identifies Potentially Pathological Gene 646 Disruptions and Highlights the Importance of Microhomology in the Mechanism of Formation. 647 Hum Mutat. Feb 2017;38(2):180-192. doi:10.1002/humu.23146 648 Thakur J, Packiaraj J, Henikoff S. Sequence, Chromatin and Evolution of Satellite DNA. 24. 649 Int J Mol Sci. Apr 21 2021;22(9)doi:10.3390/ijms22094309 650 25. Liao WW, Asri M, Ebler J, et al. A draft human pangenome reference. Nature. May 651 2023;617(7960):312-324. doi:10.1038/s41586-023-05896-x 652 Regier AA, Farjoun Y, Larson DE, et al. Functional equivalence of genome sequencing 26. 653 analysis pipelines enables harmonized variant calling across human genetics projects. Nat 654 Commun. Oct 2 2018;9(1):4038. doi:10.1038/s41467-018-06159-4 655 Sedlazeck FJ, Rescheneder P, Smolka M, et al. Accurate detection of complex structural 27. 656 variations using single-molecule sequencing. Nat Methods. Jun 2018;15(6):461-468. 657 doi:10.1038/s41592-018-0001-7 658 Suvakov M, Panda A, Diesh C, Holmes I, Abyzov A. CNVpytor: a tool for copy number 28. 659 variation detection and analysis from read depth and allele imbalance in whole-genome 660 sequencing. GigaSci. Nov 18 2021;10(11)doi:10.1093/gigascience/giab074 661 McLaren W, Gil L, Hunt SE, et al. The Ensembl Variant Effect Predictor. Genome Biol. 29. 662 2016/12// 2016;17(1):122. doi:10.1186/s13059-016-0974-4 663 Cheng H, Jarvis ED, Fedrigo O, et al. Haplotype-resolved assembly of diploid genomes 30. 664 without parental data. Nat Biotechnol. Sep 2022;40(9):1332-1335. doi:10.1038/s41587-022-665 01261-x 666 Mikheenko A, Prijbelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome 31. 667 assembly evaluation with QUAST-LG. Bioinformatics. Jul 1 2018;34(13):i142-i150. 668 doi:10.1093/bioinformatics/bty266 Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. Sep 15 669 32. 670 2018;34(18):3094-3100. doi:10.1093/bioinformatics/bty191

- 671 33. Weisenfeld NI, Kumar V, Shah P, Church DM, Jaffe DB. Direct determination of diploid 672 genome sequences. *Genome Res.* May 2017;27(5):757-767. doi:10.1101/gr.214874.116
- 673 34. Grochowski CM, Bengtsson JD, Du H, et al. Break-induced replication underlies
- formation of inverted triplications and generates unexpected diversity in haplotype structures.
   *bioRxiv*. Oct 3 2023;doi:10.1101/2023.10.02.560172
- 676 35. Sayers EW, Bolton EE, Brister JR, et al. Database resources of the national center for 677 biotechnology information. *Nucleic Acids Res.* Jan 7 2022;50(D1):D20-D26.
- 678 doi:10.1093/nar/gkab1112
- 679 36. Eisfeldt J, Vezzi F, Olason P, Nilsson D, Lindstrand A. TIDDIT, an efficient and
- 680 comprehensive structural variant caller for massive parallel sequencing data. *F1000Res*.
- 681 2017;6:664. doi:10.12688/f1000research.11168.2
- 682 37. Karolchik D. The UCSC Table Browser data retrieval tool. *Nucleic Acids Research*.
- 683 2004/01/01/ 2004;32(90001):493D-496. doi:10.1093/nar/gkh103

684