

1 **Mind the gap: the relevance of the genome reference to**
2 **resolve rare and pathogenic inversions**

3
4 Kristine Bilgrav Saether^{1,2}, Jesper Eisfeldt^{1, 2, 3*}, Jesse Bengtsson⁴, Ming Yin Lun⁴,
5 Christopher M. Grochowski⁵, Medhat Mahmoud^{5,6}, Hsiao-Tuan Chao^{5,7,8,9,10,11}, Jill A.
6 Rosenfeld⁵, Pengfei Liu^{5,12}, Jakob Schuy¹, Adam Ameer¹³, Undiagnosed Diseases
7 Network, James Paul Hwang⁶, Fritz J. Sedlazeck^{5,6,14}, Weimin Bi^{5,12}, Ronit Marom^{5,7},
8 Ann Nordgren^{1,3,15,16}, Claudia M.B. Carvalho^{4#}, Anna Lindstrand^{1,3#}

- 9 1. Department of Molecular Medicine and Surgery, Karolinska Institute, Stockholm,
10 Sweden.
- 11 2. Science for Life Laboratory, Stockholm, Sweden.
- 12 3. Department of Clinical Genetics and Genomics, Karolinska University Hospital,
13 Stockholm, Sweden.
- 14 4. Pacific Northwest Research Institute, Seattle, WA, US.
- 15 5. Department of Molecular and Human Genetics, Baylor College of Medicine,
16 Houston, TX 77030, USA.
- 17 6. Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX
18 77030, USA
- 19 7. Texas Children's Hospital, Houston, TX 77030, USA.
- 20 8. Cain Pediatric Neurology Research Laboratories, Jan and Dan Duncan
21 Neurological Research Institute, Houston, TX, USA.
- 22 9. Division of Neurology and Developmental Neuroscience, Department of
23 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)

39

40 **Abstract**

41 Both long-read genome sequencing (lrGS) 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
51 cytogenetically detected inversions. An INV6 and INV10 were resolved using GRCh38
52 and T2T-CHM13. Finally, an INV9 required optical genome mapping, *de novo* assembly
53 of lrGS 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 lrGS and optical genome mapping in solving rearrangements in challenging
61 regions of the genome. This is particularly important for inversions and may impact
62 clinical diagnostics.

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

68

69 Introduction

70 Inversions are defined as a copy-number neutral structural variants characterized by
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 false-
75 negative and false-positive rates in genome sequencing ¹. Moreover, recurrent
76 inversions formed by non-allelic homologous recombination (NAHR) use segmental
77 duplications (SDs) or other types of highly similar repeats as recombinant substrates ²⁻⁴
78 which adds to the challenge of detecting junctions mapping to poor quality regions of
79 the genome ^{1,5,6}.

80
81 We have previously shown that 28% of cytogenetically visible inversions are
82 undetected by short read genome sequencing (srGS) ⁷, suggesting that the breakpoint
83 junctions are located within large stretches of repetitive sequences. Long read genome
84 sequencing (lrGS) was shown to improve alignment and enable phasing and better
85 resolution across repetitive regions ⁸⁻¹⁰. Regardless, inversions with breakpoints
86 mapping to large repeats remain challenging to resolve even when applying lrGS ^{5,11}.

87
88 Previous studies show that the new T2T-CHM13 (T2T) reference provides increased
89 sensitivity in inversion detection due to increased resolution across repetitive sequences
90 ^{5,12}. This reference genome has filled gaps present in earlier reference genomes, adding
91 >200 Mb of sequence compared to GRCh38 ¹². In fact, both GRCh37 and GRCh38 lack

92 information across hundreds of mega-base pairs (Mb) of regions such as telomeres,
93 centromeres and other repetitive regions¹²⁻¹⁷. Often forgotten resources in human
94 genetic analysis are the closely related primate genomes chimpanzee¹⁸ and bonobo¹⁹
95 that have been fully sequenced, with up to 99% of gaps closed¹⁹. Many sequences
96 unmappable after srGS analysis may be present in primates^{17,20}.

97

98 We previously solved 72% (13/18) cytogenetically detected inversions using srGS and
99 GRCh37⁷. Here, we investigated six paracentric and pericentric large genomic
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, lrGS 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
108 the utility of applying long molecules to investigate the architecture of rare diseases.

109 **Results**

110 **Resolving inversions using lrGS and reference genomes**

111 The six rare inversions affect chromosomes 6, 9, 10, 11 and 18 in six unrelated
112 individuals (Table 1). They were aligned to GRCh37, GRCh38 and T2T.

113

114 **Table 1. Overview of included inversions.** Including the reference genome in which
 115 they could be detected and the type of genomic sequencing technology used (short-
 116 read (sr), linked-read (lir), long-read (lr) genome sequencing (GS) and optical genome
 117 mapping (OGM)). Results are indicated as detected (☑), unclear (-) or absent (X).

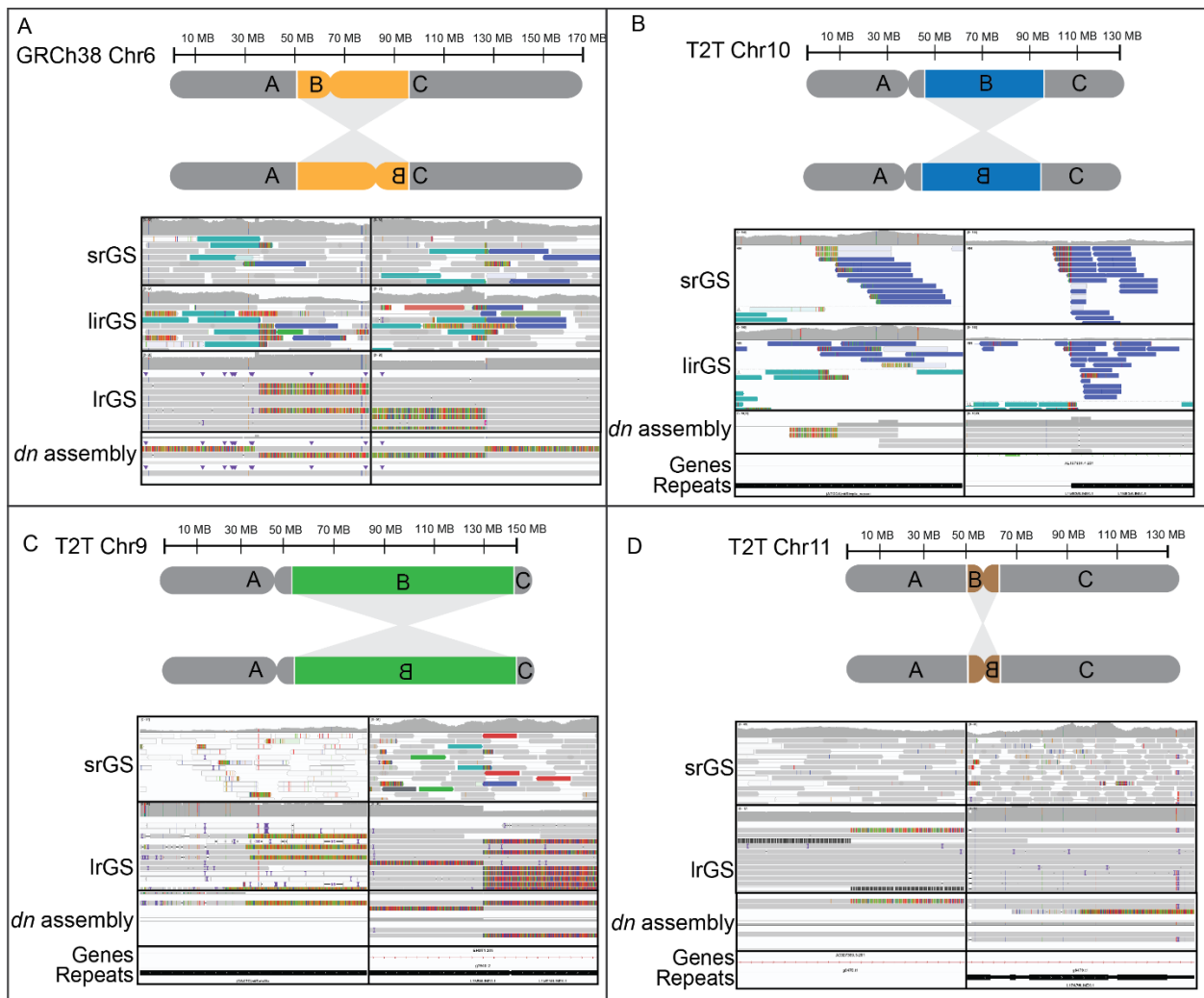
Case ID	Karyotype	GRCh37	GRCh38	T2T	Sequencing data
P4855_501	46, XY, inv(6)(p12q16.3)	×	☑	☑	srGS, lirGS, lrGS
P5371_208	46, XY, inv(9)(p13q22)	×	×	×	srGS, lirGS
P4855_208	46, XY, inv(11)(p11.1q12)mat	-	-	-	srGS, lirGS, lrGS
P4855_106	46, XY, inv(10)(q11q23)pat	×	×	☑	srGS, lirGS
P5370_201	46, XY, inv(18)(p11.3q11.2)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

118
 119 Two (P4855_501, P4855_106) inversions were detected by standard variant callers
 120 after realigning the srGS to a new reference genome. One additional case (BH16643-1)
 121 was resolved using a combination of new reference genomes, de novo assembly and
 122 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 lrGS SV analysis
 126 as well as *de novo* assembly in GRCh37⁷. Following realignment of the GS data to
 127 GRCh38 and T2T, the exact inversion breakpoints were present in both srGS, lirGS,
 128 lrGS 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

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

140



141

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)** An inversion 10 visible in srGS and lirGS data using T2T.
146 **C)** An inversion 9 only visibly by lrGS de novo assembly using T2T. **D)** An inversion 11
147 with centromeric breakpoints. The lrGS *de novo* assembly call fitting with the
148 cytogenetic analysis is shown.

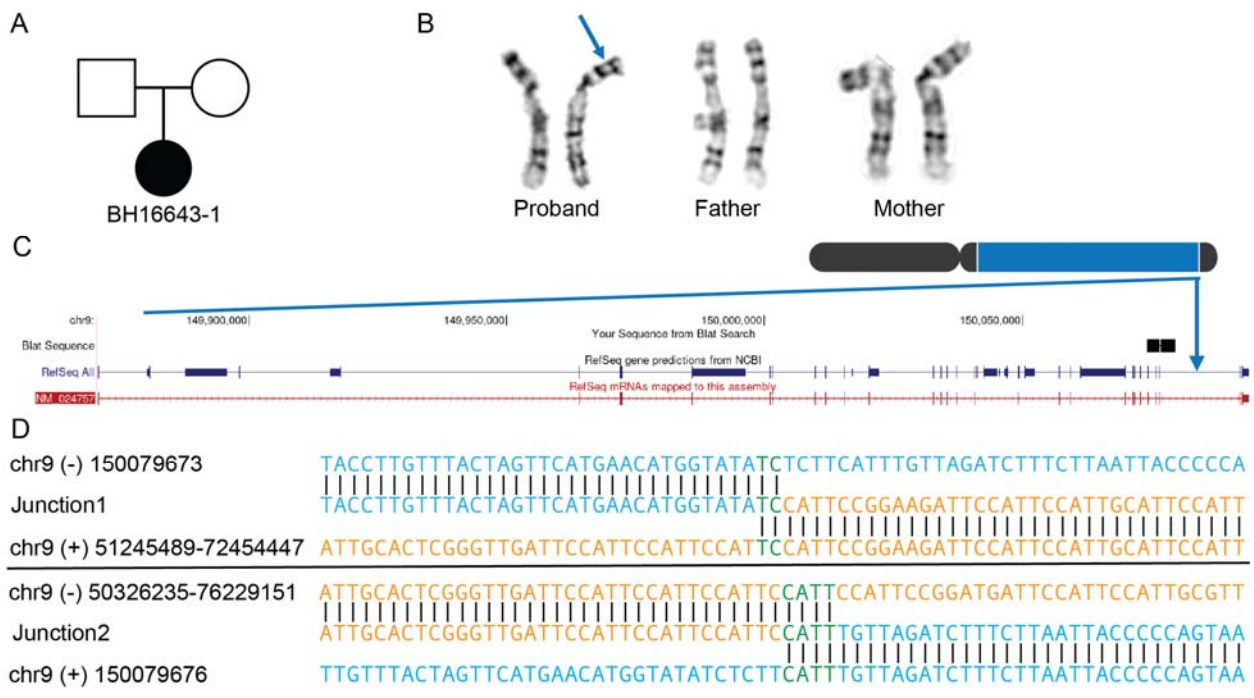
149

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 *de novo* 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, lrGS 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 *de novo* 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

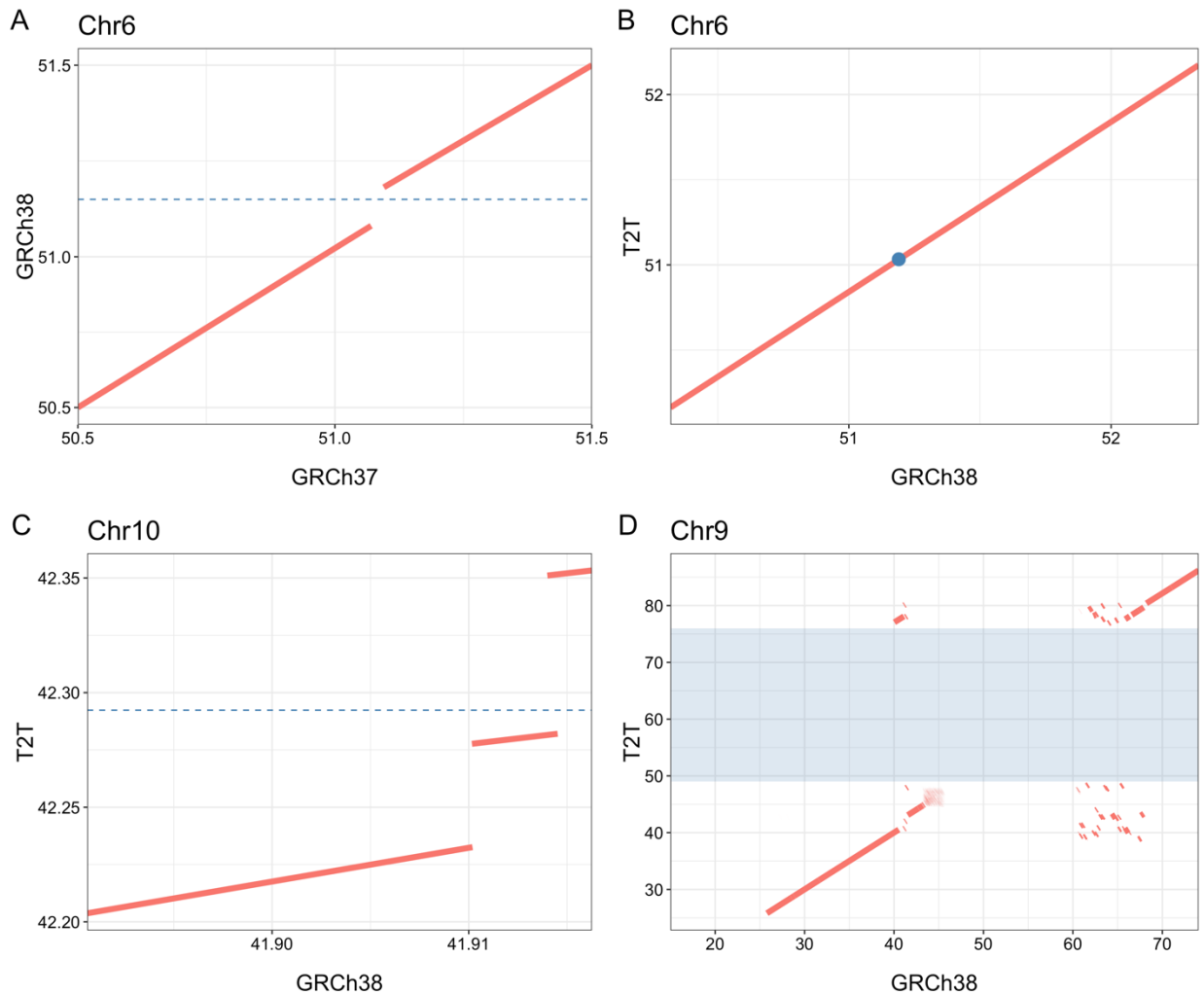
187 For one inversion, affecting chromosome 11 (P4855_208) identified by karyotyping in a
188 patient suffering from neurodevelopmental delay (Table 1, Figure 1D) a potential
189 inversion call was suggested after IrGS *de novo* assembly. The suggested breakpoints
190 were present in all the assessed 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**

198 Reference genome analysis revealed that inversion breakpoint sequences were missing
199 in reference genomes GRCh37 (inv6) and GRCh38 (inv9, inv10), making it impossible
200 to solve them using these reference genomes (Fig. 3, Supplementary Fig. 1 and 3). In
201 total 127 kb of sequence at 6p12.3 was present in GRCh38 but missing from GRCh37.
202 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
 204 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



208 **Fig 3: Comparison of the inversion breakpoint region on chromosome 6p12.3,**
 209 **chromosome 10q11 and chromosome 9q12.** Reference sequences were aligned to
 210 each other and shown as dot-plots. The dashed line and dot represents the position of
 211 the breakpoint of the inversions. **A)** The chromosome 6p inversion breakpoint is located
 212

213 in a 127kb region in GRCh38 missing from GRCh37. **B)** The chromosome 6p inversion
214 breakpoint in GRCh38 and T2T. **C)** The chromosome 10q breakpoint is located in a
215 69kb region missing in GRCh38, with a surrounding 4kb duplication which occurs only
216 once in T2T. **D)** The chromosome 9q12 breakpoint is located in a 28 Mb region missing
217 in GRCh38 shaded in blue.

218

219 For the inversion on chromosome 10, the 10:q11 breakpoint was located in a 69 kb
220 region of simple repeats only present in T2T (Fig. 3, 4, Supplementary Fig. 1 and 3).

221 The region, spanning from 10:42282056-42351085, consists of 99% simple repeats and
222 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,
225 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

234 detected in T2T-GRCh38 (10kb-34Mb, median 30kb). The chimpanzee-T2T (range

235 10kb-14Mb, median 40kb) and bonobo-T2T had similar ranges of DRRs (range 10kb-

236 19Mb and median 35kb) (Supplementary Table 2). In total, we uncovered 203 regions
 237 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

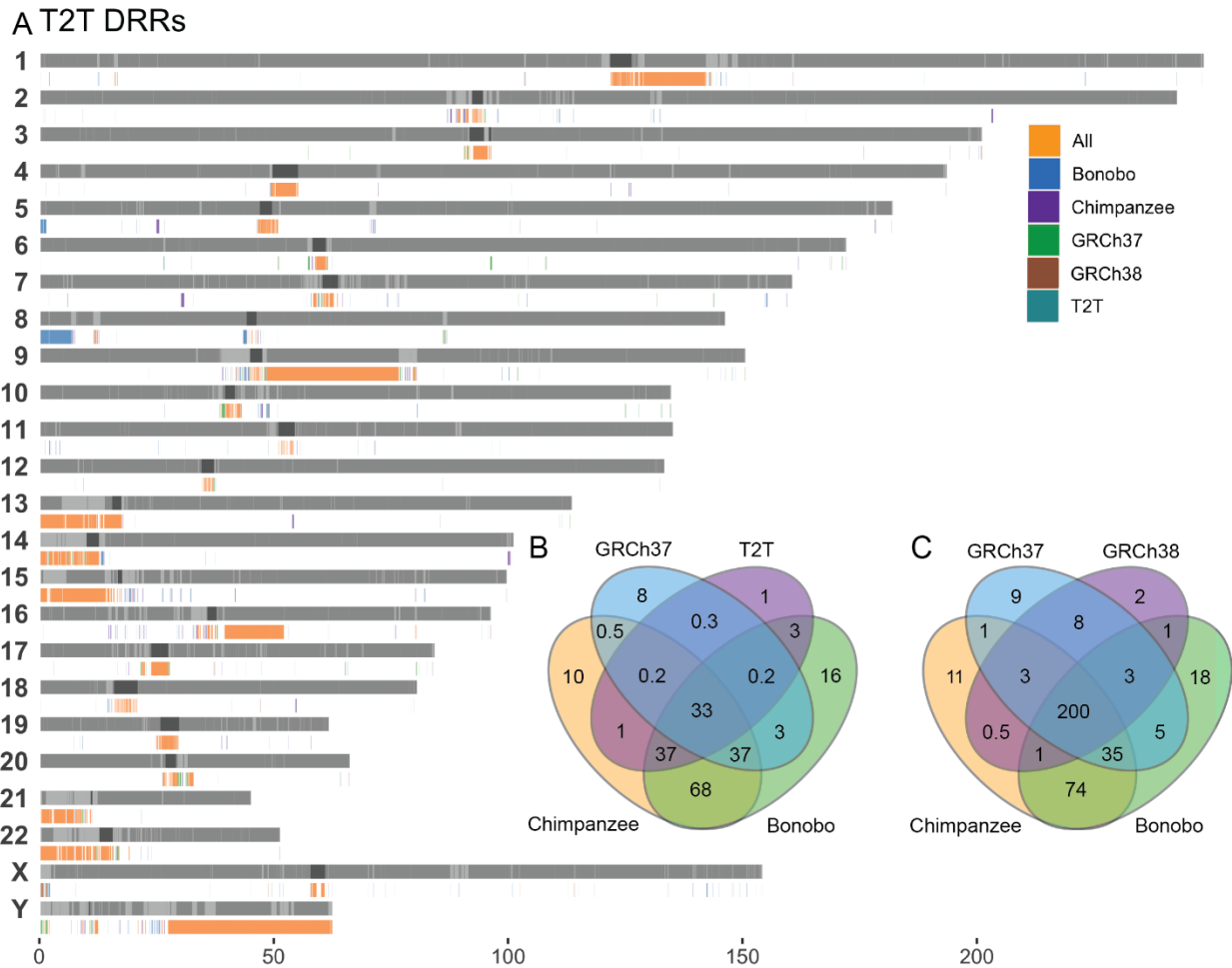
240 **Table 1: Differential reference regions between reference genomes.** For each
 241 template on the top row (grey) the total amount of sequence in Mb and on the second
 242 row (white) the total number of DRRs is given in comparison with the query reference.

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

243

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

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

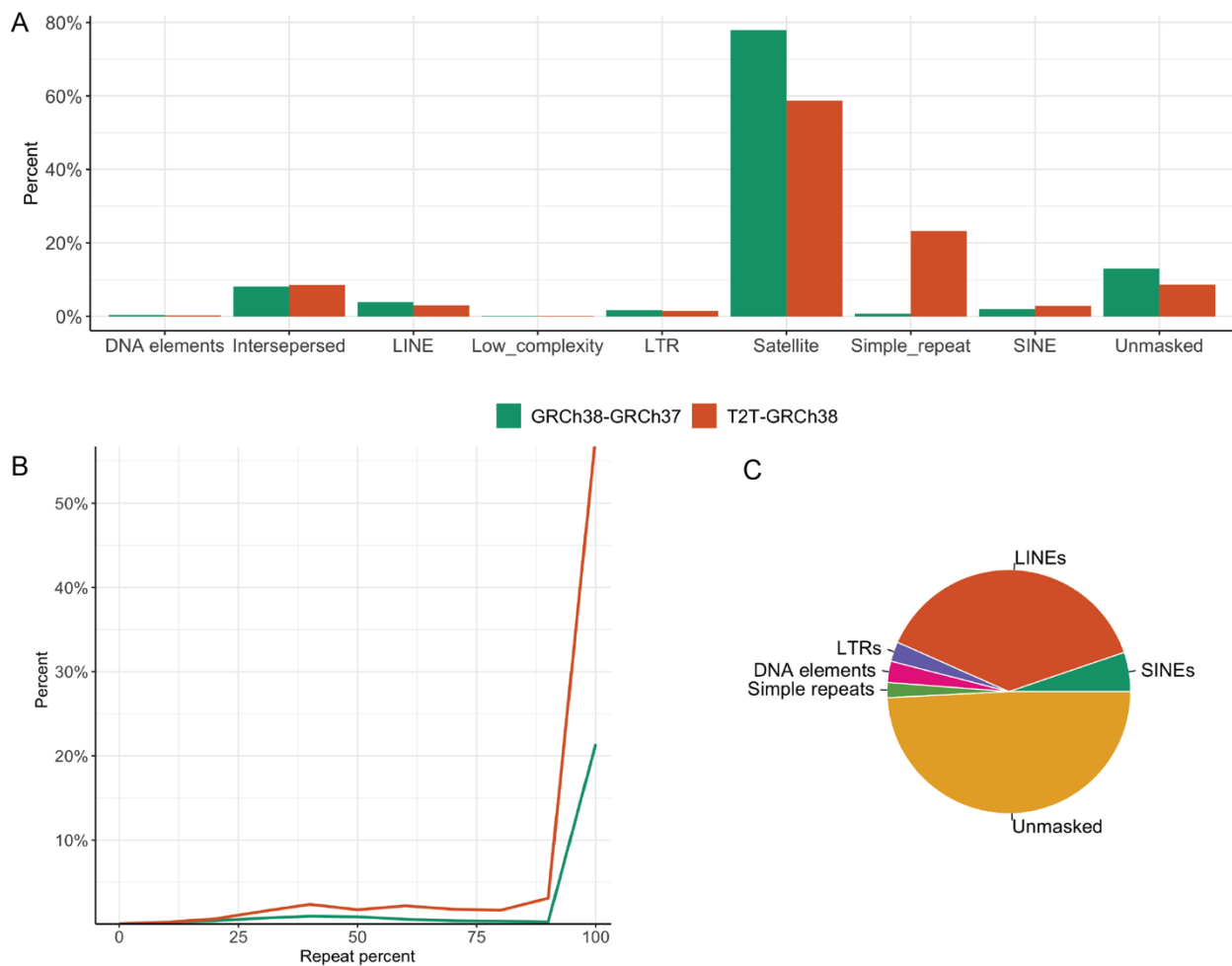


250 **Fig 4: Shared DRR in T2T and GRCh38. A)** Bar plot of all T2T DRRs **B)** Venn diagram
 251 of Mb overlap between all GRCh38 DRRs, **C)** Venn diagram of Mb overlap between all
 252 T2T DRRs.
 253

254
 255 **DRRs introduce repetitive sequences**

256 Repeat analysis of all DRRs in GRCh38-GRCh37 and T2T-GRCh38 reveal most to be
 257 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

259 sequence, 38% LINEs, 5.2% SINEs, 2% simple repeats, 2.6% LTRs and 2.8% DNA
 260 elements (Fig. 5C). Furthermore, of all T2T-GRCh38 DRR sequences, 55% consisted of
 261 100% repetitive DNA, 92% consisted of >50% repetitive DNA, 20% were located inside
 262 or within 10kb of centromeric regions and 30% within segmental duplications (Fig. 5B).
 263 Of GRCh38-GRCh37 DRR sequences, 20% consisted of 100% repetitive DNA, 89%
 264 consisted of >50% of repetitive DNA, 76% were located within 10 kb of centromeric
 265 regions and 20% within segmental duplications (Fig. 5B).



266
 267 **Fig 5: Repeat characterization across DRRs. A)** Percentage of repeat elements
 268 (masked by repeat masker) in the DRR sequences from GRCh38-GRCh37 and T2T-
 269 GRCh38. **B)** Distribution of DRR sequences and their repeat percentage in GRCh38-

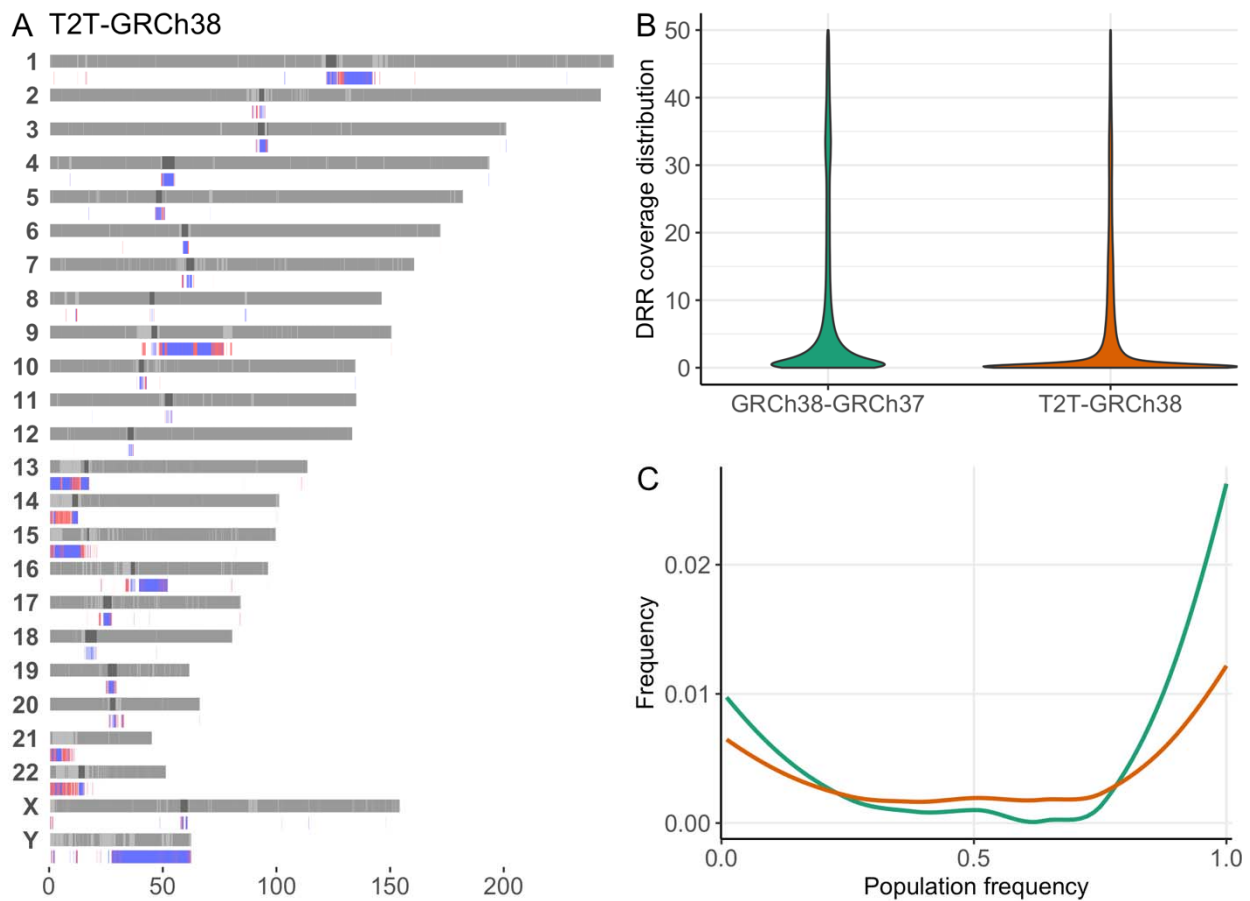
270 GRCh37 and T2T-GRCh38. **C)** Pie chart displaying repeat content in the GRCh38-
 271 GRCh37 DRR sequences affected by the inversion 6 at the 6p12 junction
 272 (chr6:51102785-51230413) in GRCh38.

273

274 **DRR sequences in the general population**

275 Next, we aligned srGS data from 100 Swedish individuals ²¹ to the five references and
 276 assessed the presence of DRR across the population (Fig 6, Supplementary Fig. 9 and
 277 10).

278



279

280 **Fig 6: Distribution of DRRs in 100 Swedish individuals. A)** DRRs between T2T-
 281 GRCh38 and their presence in Swedish individuals. Blue indicating absent (<8X) and

282 red present ($>8X$ and $<100X$). **B)** Violin plot of coverage across the respective DRRs in
283 100 Swedish individuals. **C)** The distribution of population frequencies of the detected
284 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 below 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

295 In comparison, for the GRCh38-GRCh37 DRRs, 60% were not detected
296 (Supplementary Fig. 10), whilst of the 40% detected, 13% were rare and 86% common
297 (Fig. 6C). Furthermore, 53% were found in over 90% of the queried individuals. The
298 mapping quality of reads from 5 individuals across DRRs was assessed, where 25% of
299 reads had a mapping quality above 20 (Supplementary Fig. 8). The violin plot confirms
300 that most DRRs lack aligned reads (Fig. 6B-C).

301

302 In Chimpanzee-T2T DRRs 80% were missing in the Swedish individuals, whilst of the
303 20% present, 6% were rare, 93% common and 75% in $>90\%$ of the individuals.

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

307 **Discussion**

308 The availability of long read sequencing and the new reference T2T-CHM13
309 incentivized us to revisit six previously unsolved cytogenetically visible inversions. Three
310 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 lrGS in itself is not the sole answer.

313

314 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.
317 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 lrGS, which improves resolution of repeats, is not yet broadly
322 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
324 inversion 11, lrGS nor *de novo* assembly was sufficient at pinpointing the breakpoints.
325 For the two remaining, lrGS was unfortunately not possible. Unfortunately, lrGS 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 patterns 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

339 We and others have previously suggested that a part of cytogenetically visible
340 inversions may have been formed through non-allelic homologous recombination
341 (NAHR) explaining why some remain undetected even after srGS analysis^{2,7 3}.
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
346 variants or other concomitantly alterations suggesting canonical non-homologous end
347 joining (c-NHEJ) as the underlying mechanism of formation^{7,22}. Even so, the inversion
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 9q consisted of

350 99% satellite and simple repeats and the 69 kb DRR on 10q consisted of 99% simple
351 repeats. This result supports that copy-number neutral inversions, similarly to balanced
352 translocations, may result from an error prone repair of processed double-strand breaks
353 (DSBs)²³.

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
364 phenotype²². All together, these results underscore the relevance of inversions to
365 unsolved rare disease, often undetected by current clinical genome sequencing.

366

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

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

376

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
382 short reads, in addition to these regions being highly variable between individuals²⁴
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
390 diploid assemblies represented as a graph²⁵. This assembly is able to represent large
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
394 calling.

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
402 primate genomes are a useful resource to improve our understanding of repetitive and
403 complex regions which have previously been understudied.

404

405 **Methods**

406 **Study participants**

407 IDs used in this article are not known to anyone outside the research group. Subjects
408 carrying Inversions 6, 10, 11 and 18 were enrolled at Karolinska University Hospital,
409 Stockholm, Sweden⁷. Patient BH16643-1 was enrolled using research protocol H-
410 47281/Pacific Northwest Research Institute WIRB #20202158 and 15HG0130 with the
411 NIH IRB as part of the Undiagnosed Diseases Network (UDN). Whole blood samples
412 (3-10mL) were collected from the patient and parents. DNA was isolated from blood
413 according to standard procedures.

414

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

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

421

422 **Genome sequencing**

423 For samples (P4855_501, P5371_208, P4855_208, P4855_106, P5370_106) srGS and
424 10X genomics linked read sequencing of the samples was performed at the national
425 genomics infrastructure (NGI) at Science for Life laboratory (SciLifeLab) Stockholm as
426 previously mentioned⁷. Analysis for structural variants was performed using FindSV as
427 described previously⁷.

428

429 lrGS was performed on P4855_501 and P4855_208 using Pacific Biosciences (PacBio)
430 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 pipeline²⁶.

436 lrGS 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 lrGS was aligned using:

```
444 longranger wgs --id <id> --reference <ref> --fastq <fastq> --  
445 vcmode freebayes
```

446

447 The lrGS data was aligned to GRCh37, GRCh38 and T2T. Analysis of was done using
448 an in house developed pipeline LOMPE (<https://github.com/kristinebilgrav/LOMPE>).

449 LOMPE uses minimap2 for alignment and combines Sniffles (v1)²⁷ and CNVpytor²⁸ for
450 structural variant calling, and produces a single VCF file which is annotated using
451 Variant Effect Predictor (VEP)²⁹. The resulting lrGS 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 lrGS from samples P4855_501, P4855_208 and BH16643-1 was
456 performed using hifiasm³⁰. Quality control was performed using quast³¹. Alignment to
457 reference genomes GRCh37, GRCh38 and T2T was performed using minimap2³², and
458 variant calling was performed using sniffles (v1)²⁷ and htsbox
459 (<https://github.com/lh3/htsbox>). On lrGS from sample P4855_106 a *de novo* assembly
460 was performed using 10X Genomics Supernova³³.

461

462 **Optical genome mapping**

463 Optical genome mapping was performed as described previously³⁴. Briefly, ultra-high
464 molecular weight (UHMW) genomic DNA for use in genomic optical mapping was
465 extracted from blood using Bionano Prep™ Blood and Cell Culture DNA Isolation Kit
466 (Bionano Genomics) with an input of 1.5 million cells. Subsequent DNA quantity and
467 size was confirmed using Qubit™ dsDNA BR Assay Kit. A total of 0.75 µg of HMW DNA
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**

482 Breakpoint verification of breakpoints identified in P4855_208 was performed as
483 previously described⁷.

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)*³⁵. 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)³². 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.3.6.0³⁶, producing a bed file with gaps between the query
501 and template. Files with known gap regions were downloaded from UCSC
502 TableBrowser³⁷ 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.

522

523 **Data access**

524 The reference genomes can be downloaded from *NCBI*³⁵. The clinical samples are not
525 available due to ethical permissions. The SweGen dataset is available at

526 <https://swefreq.nbis.se/> upon signing a data agreement. The srGS analysis pipeline

527 FindSV is available on GitHub at <https://github.com/J35P312/FindSV>. The lrGS analysis

528 pipeline LOMPE is available at <https://github.com/kristinebilgrav/LOMPE>. Tools TIDDIT

529 and SVDB are available at <https://github.com/SciLifeLab/TIDDIT> and

530 <https://github.com/J35P312/SVDB>.

531

532 **Ethics statement**

533 Ethics approval for analysis of participant samples was given by the Regional Ethical
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
554 support was provided through the National Institute of Neurological Disorders and
555 Stroke of the National Institutes of Health (U01HG007709 and U01HG007942) and the

556 National Institute of Health (NIH S10 1S10OD028587). The content is solely the
557 responsibility of the authors and does not necessarily represent the official views of the
558 National Institutes of Health. The funders had no role in study design, data collection
559 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
567 resources provided by SNIC through Uppsala Multidisciplinary Center for Advanced
568 Computational Science (UPPMAX) under Project SNIC sens2017106 and
569 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 2. Stankiewicz P, Lupski JR. Genome architecture, rearrangements and genomic
589 disorders. *Trends Genet.* Feb 2002;18(2):74-82. doi:10.1016/s0168-9525(02)02592-1
- 590 3. Carvalho CM, Lupski JR. Mechanisms underlying structural variant formation in genomic
591 disorders. *Nat Rev Genet.* Apr 2016;17(4):224-38. doi:10.1038/nrg.2015.25
- 592 4. Dittwald P, Gambin T, Gonzaga-Jauregui C, et al. Inverted low-copy repeats and
593 genome instability--a genome-wide analysis. *Hum Mutat.* Jan 2013;34(1):210-20.
594 doi:10.1002/humu.22217
- 595 5. Porubsky D, Harvey WT, Rozanski AN, et al. Inversion polymorphism in a complete
596 human genome assembly. *Genome Biol.* Apr 30 2023;24(1):100. doi:10.1186/s13059-023-
597 02919-8
- 598 6. Kidd JM, Graves T, Newman TL, et al. A human genome structural variation sequencing
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 7. Pettersson M, Grochowski CM, Wincent J, et al. Cytogenetically visible inversions are
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 9. Logsdon GA, Vollger MR, Eichler EE. Long-read human genome sequencing and its
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
611 phased genome assemblies. *Genome Res.* Apr 2023;33(4):496-510.
612 doi:10.1101/gr.277334.122
- 613 12. Nurk S, Koren S, Rhie A, et al. The complete sequence of a human genome. *Science.*
614 Apr 2022;376(6588):44-53. doi:10.1126/science.abj6987
- 615 13. Pan B, Kusko R, Xiao W, et al. Similarities and differences between variants called with
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 14. Schneider VA, Graves-Lindsay T, Howe K, et al. Evaluation of GRCh38 and de novo
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. Ameer 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 17. Einfeldt J, Martensson G, Ameer A, Nilsson D, Lindstrand A. Discovery of Novel
628 Sequences in 1,000 Swedish Genomes. *Mol Biol Evol.* Jan 1 2020;37(1):18-30.
629 doi:10.1093/molbev/msz176
- 630 18. Chimpanzee S, Analysis C. Initial sequence of the chimpanzee genome and comparison
631 with the human genome. *Nature.* Sep 1 2005;437(7055):69-87. doi:10.1038/nature04072
- 632 19. Mao Y, Catacchio CR, Hillier LW, et al. A high-quality bonobo genome refines the
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 21. Ameer A, Dahlberg J, Olason P, et al. SweGen: a whole-genome data resource of
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 22. Grochowski CM, Krepischi ACV, Einfeldt J, et al. Chromoanagenesis Event Underlies a
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 23. Nilsson D, Pettersson M, Gustavsson P, et al. Whole-Genome Sequencing of
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 24. Thakur J, Packiaraj J, Henikoff S. Sequence, Chromatin and Evolution of Satellite DNA.
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 26. Regier AA, Farjoun Y, Larson DE, et al. Functional equivalence of genome sequencing
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 27. Sedlazeck FJ, Rescheneder P, Smolka M, et al. Accurate detection of complex structural
656 variations using single-molecule sequencing. *Nat Methods.* Jun 2018;15(6):461-468.
657 doi:10.1038/s41592-018-0001-7
- 658 28. Suvakov M, Panda A, Diesch C, Holmes I, Abyzov A. CNVpytor: a tool for copy number
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 29. McLaren W, Gil L, Hunt SE, et al. The Ensembl Variant Effect Predictor. *Genome Biol.*
662 2016/12// 2016;17(1):122. doi:10.1186/s13059-016-0974-4
- 663 30. Cheng H, Jarvis ED, Fedrigo O, et al. Haplotype-resolved assembly of diploid genomes
664 without parental data. *Nat Biotechnol.* Sep 2022;40(9):1332-1335. doi:10.1038/s41587-022-
665 01261-x
- 666 31. Mikheenko A, Prijbelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome
667 assembly evaluation with QUAST-LG. *Bioinformatics.* Jul 1 2018;34(13):i142-i150.
668 doi:10.1093/bioinformatics/bty266
- 669 32. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* Sep 15
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
674 formation of inverted triplications and generates unexpected diversity in haplotype structures.
675 *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

685