

1 Full characterization of unresolved structural variation through long- 2 read sequencing and optical genome mapping

3 Griet De Clercq^{1,2}, Lies Vantomme¹, Barbara Dewaele³, Bert Callewaert^{1,2}, Olivier Vanakker^{1,2}, Sandra
4 Janssens^{1,2}, Bart Loeys⁴, Mojca Strazisar^{5,6}, Wouter De Coster^{5,6}, Joris Robert Vermeesch^{3,7}, Annelies
5 Dheedene², Björn Menten^{*1,2}

6 1. Department of Biomolecular Medicine, Ghent University, Ghent, Belgium

7 2. Center for Medical Genetics Ghent, Ghent University Hospital, Ghent, Belgium

8 3. Center for Human Genetics Leuven, University Hospital Leuven, Leuven, Belgium.

9 4. Center for Medical Genetics Antwerp, University of Antwerp, Antwerp University Hospital,
10 Antwerp, Belgium.

11 5. VIB Center for Molecular Neurology, VIB, Antwerp, Belgium

12 6. Department of Biomedical Sciences, University of Antwerp, Antwerp, Belgium

13 7. Department of Human Genetics, KU Leuven, Leuven, Belgium.

14 ***Corresponding author:** Björn Menten (bjorn.menten@ugent.be)

15

16 Acknowledgements

17 We are grateful to all the individuals and their families who participated in this study. Furthermore,
18 we thank Professor S. Vergult and doctor H. Sryn for their fruitful and insightful discussions and
19 their support during manuscript preparation.

20 Abstract

21 Structural variants (SVs) are important contributors to human disease. Their characterization remains
22 however difficult due to their size and association with repetitive regions. Long-read sequencing (LRS)
23 and optical genome mapping (OGM) can aid as their molecules span multiple kilobases and capture
24 SVs in full. In this study, we selected six individuals who presented with unresolved SVs. We applied
25 LRS onto all individuals and OGM to a subset of three complex cases. LRS detected and fully resolved
26 the interrogated SV in all samples. This enabled a precise molecular diagnosis in two individuals.
27 Overall, LRS identified 100% of the junctions at single-basepair level, providing valuable insights into
28 their formation mechanisms without need for additional data sources. Application of OGM added
29 straightforward variant phasing, aiding in the unravelment of complex rearrangements. These results
30 highlight the potential of LRS and OGM as follow-up molecular tests for complete SV characterization.
31 We show that they can assess clinically relevant structural variation at unprecedented resolution.
32 Additionally, they detect (complex) cryptic rearrangements missed by conventional methods. This
33 ultimately leads to an increased diagnostic yield, emphasizing their added benefit in a diagnostic
34 setting. To aid their rapid adoption, we provide detailed laboratory and bioinformatics workflows in
35 this manuscript.

36 **Keywords:** long-read sequencing, optical genome mapping, structural variation, clinical genomics,
37 chromothripsis, complex genomic rearrangements

38 Introduction

39 Structural variants (SVs) are genomic rearrangements of 50 base pairs (bp) or larger that are
40 categorized into deletions, duplications, insertions, inversions, and translocations [1]. In addition, more
41 complex genomic rearrangements (CGRs) exist, where multiple SV types are combined in a single
42 event. Recent research indicates that one individual's genome harbors 26,000 SVs on average,
43 accounting for nearly 30 Mb of reorganized DNA [2]. Besides being part of normal variation between
44 individuals, they also contribute to disorders such as cancer [3], autism [4], and (syndromic) intellectual
45 disability [5].

46
47 Routine diagnostic techniques employed to detect structural variation consist of conventional
48 karyotyping, fluorescence in situ hybridization (FISH), microarrays, and more recently (shallow) whole
49 genome sequencing (WGS) [6]. These techniques are however unable to interrogate the full spectrum
50 of structural variation and often fail to pinpoint exact breakpoints. In a diagnostic setting, precise
51 identification and full characterization of SVs is however crucial as it allows for a conclusive clinical and
52 molecular diagnosis. This can provide insight into the prognosis of a disease and its possible
53 therapeutic management [7]. Short-read WGS has the potential to cater to these shortcomings, but
54 due to the small read size inherent to this technology, the detection of SVs that are only a few kbs in
55 size or reside within repeat regions is often challenging [8]. Recently, long-read sequencing (LRS) and
56 optical genome mapping (OGM) have been on the rise for accurate SV detection [9,10]. LRS is a
57 sequencing technology that reads DNA and RNA at the single-base level, while OGM is a non-
58 sequencing-based imaging technique that maps patterns of a fluorescent labelled DNA motif [11,12].
59 Reads generated by LRS typically measure between 10 to 100 kb, while OGM molecules start at 150
60 kb. Both techniques can however analyze fragments up to several Mbs in length [13]. Due to these
61 long fragment lengths, SVs can be fully captured in unprecedented detail and resolution.

62

63 In this study we selected six individuals with a previously identified SV that remained unresolved
64 through routine molecular techniques. In three cases no clear causal molecular diagnosis could be
65 made through previous tests. We sequenced all individuals with nanopore LRS (Oxford Nanopore
66 Technologies (ONT)), and additionally applied OGM (Bionano Genomics) to (highly) complex cases that
67 benefit from further haplotype phasing. The aim of this study was to investigate the ability of LRS and
68 OGM to fully characterize unresolved structural variation, their potential to identify missed cryptic
69 rearrangements, and to evaluate their added value in a clinical setting.

70 Results

71 Simple structural variants are identified at single base-pair level, enabling a new 72 molecular diagnosis

73 Individual S1 presented with severe ID and dysmorphic facial features. Through karyotyping, an
74 apparently balanced *de novo* reciprocal translocation between chromosome bands 9q21.2 and
75 10p15.2 was identified (table 1; sup. fig. S1a,b). As all other diagnostic tests came back normal, this
76 translocation was highly suggestive as causal for the observed phenotype. However, exact breakpoint
77 coordinates could not be further determined [14]. Application of LRS identified and characterized this
78 variant at single-basepair level as chr9:g.pter_cen_82209673::chr10:g.11203945_pter (der9) and
79 chr9:g.qter_82209674::chr10:g.11203946_cen_qter (der10) (fig. 1a,b; sup. table S3). Through
80 sequence analysis, a 1 bp guanine insertion at the breakpoint of derivative 10 was found (sup. text S1).
81 Additional confirmation by Sanger sequencing revealed the final variant to be a true event and 100%
82 concordant with the LRS consensus. Microhomology analysis near the breakpoints revealed no
83 microhomology stretches. The breakpoint on chromosome 9 falls within a LINE repeat, while on
84 chromosome 10 the *CELF2* gene is disrupted (fig. 1a, sup. fig. S8a). Itai et al. (2021) [15] recently
85 described deleterious variants in *CELF2* that lead to developmental and epileptic encephalopathy
86 (OMIM #619561), and ID and autistic features through a loss-of-function mechanism. This led to a
87 conclusive molecular diagnosis in this individual.

88

89 Individual S2 manifested with mild ID, Down-like features, and autism spectrum disorder. Through
90 whole exome sequencing a heterozygous *de novo* deletion was detected of exon 9 of the *MYT1L* gene
91 (table 1). Variations in this gene are associated with autosomal dominant intellectual developmental
92 disorder 39 (OMIM #616521). As this deletion results in an out-of-frame transcript, it was considered
93 causal for the underlying phenotype. LRS was applied to further test its ability to characterize simple

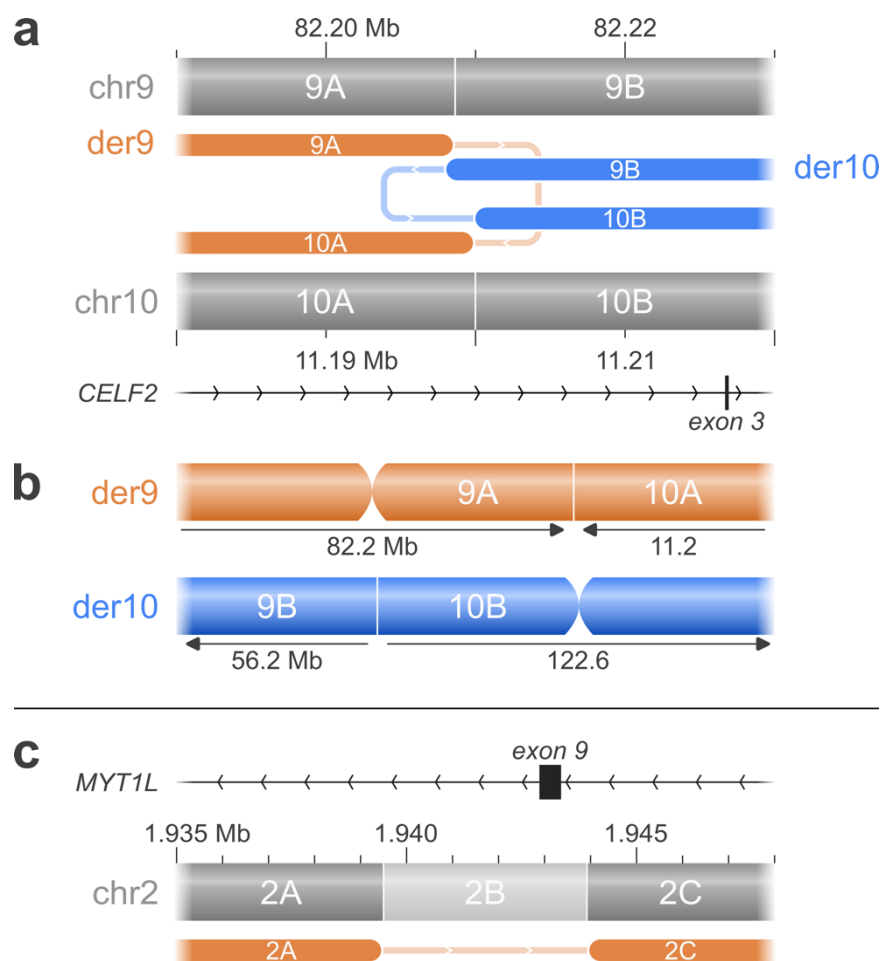


Fig. 1 Simple SVs in individuals S1 and S2. **(a)** Exact breakpoint determination of the t(9;10) variant in individual S1 through LRS allowed for a novel molecular diagnosis through the disruption of the *CEL F2* gene. **(b)** Local karyotype of the aberrant chromosomes of individual S1. **(c)** The exon 9 *MYT1L* deletion in individual S2 was delineated through LRS at single base pair resolution as a 4,574 bp deletion (fragment 2B) enclosing exon 9.

94 SVs. Variants were confined to the *MYT1L* region, which delineated the SV as a 4,574 bp deletion
 95 enclosing exon 9 (NC_000002.12:g.1939353-1943926del) (fig. 1c; sup. table S3; sup. fig. S4a,S8b).
 96 Microhomology analysis revealed a 1 bp microhomology at the junction breakpoint, and Sanger
 97 sequencing confirmed the junction sequence to be 100% concordant to the consensus sequence as
 98 determined by LRS (sup. text S2).

99 Conventional techniques underestimate the incidence and complexity of structural
 100 variation

101 Individual S3 was referred to the clinic for multiple phenotypic aberrations corresponding to a clinical
 102 diagnosis of Mowat-Wilson syndrome (OMIM #235730). Karyotyping revealed a balanced *de novo*

103 reciprocal translocation between chromosome bands 2q22 and 21q21 (table 1; sup. fig. S1c).
104 Haploinsufficiency of *ZEB2*, located on chromosome band 2q22, can lead to Mowat-Wilson syndrome
105 making this translocation likely causal for the observed phenotype. However, no exact breakpoints and
106 thus no precise molecular diagnosis could be established through other diagnostic methods. LRS was
107 applied and the variants filtered to retain translocations between the two long arms of chromosomes
108 2 and 21. This identified two distinct translocations in bands 2q22.3 and 21q21.2. These translocations
109 however do not disrupt *ZEB2* and their breakpoints on chromosome 2 hit respectively the protein
110 coding gene *GTDC1* and lncRNA gene *TEX41*. These genes are respectively lying up- and downstream
111 of *ZEB2*, suggesting a more complex disruption of the region. Further manual inspection revealed other
112 rearrangements which directly disrupt the *ZEB2* gene and involve chromosome 5 as well (sup. table
113 S3,S8). Reconstruction of the separate fragments led to the discovery of a CGR consisting of 23
114 breakpoints on chromosomes 2, 5, and 21, affecting 1.4 Mb in total (fig. 2a,b; sup. fig. S7a). Due to its
115 complexity, OGM was additionally applied to verify this reconstructed rearrangement. This confirmed
116 the CGR, although some smaller fragments could not be detected through OGM analysis. Eight
117 breakpoint junctions solely identified through LRS were therefore successfully confirmed through PCR
118 (sup. table S5,S8). The final rearrangement was found to be relatively balanced with no large CNVs,
119 except for a 25 kb deletion affecting exons 5 to 10 of the *ZEB2* gene (fig. 2a,b; sup. fig. S4g). This
120 deletion was picked up by both LRS and OGM SV analysis (sup. fig. S9). LRS CNV analysis however called
121 no copy number changing events in the affected regions (sup. table S2). Additionally, the vast majority
122 of breakpoints are flanked by minor losses of genetic material, identified as deleted fragments in the
123 CGR reconstruction (fig. 2a,b; sup. fig. S7a) and distinct drops in the LRS coverage data (sup. fig. S4b-
124 n). The discovery of this previously concealed complex rearrangement affecting *ZEB2* led to a
125 molecular diagnosis and the confirmation of the clinical diagnosis of Mowat-Wilson syndrome. Apart
126 from *GTDC1* and *ZEB2*, no other protein coding genes were directly disrupted. Investigation for
127 microhomology sequences around the LRS breakpoint junctions revealed 4 to 90 bp insertions in seven
128 variants (sup. text S3). In four out of five junctions with insertions over 20 bp, the insertion sequences

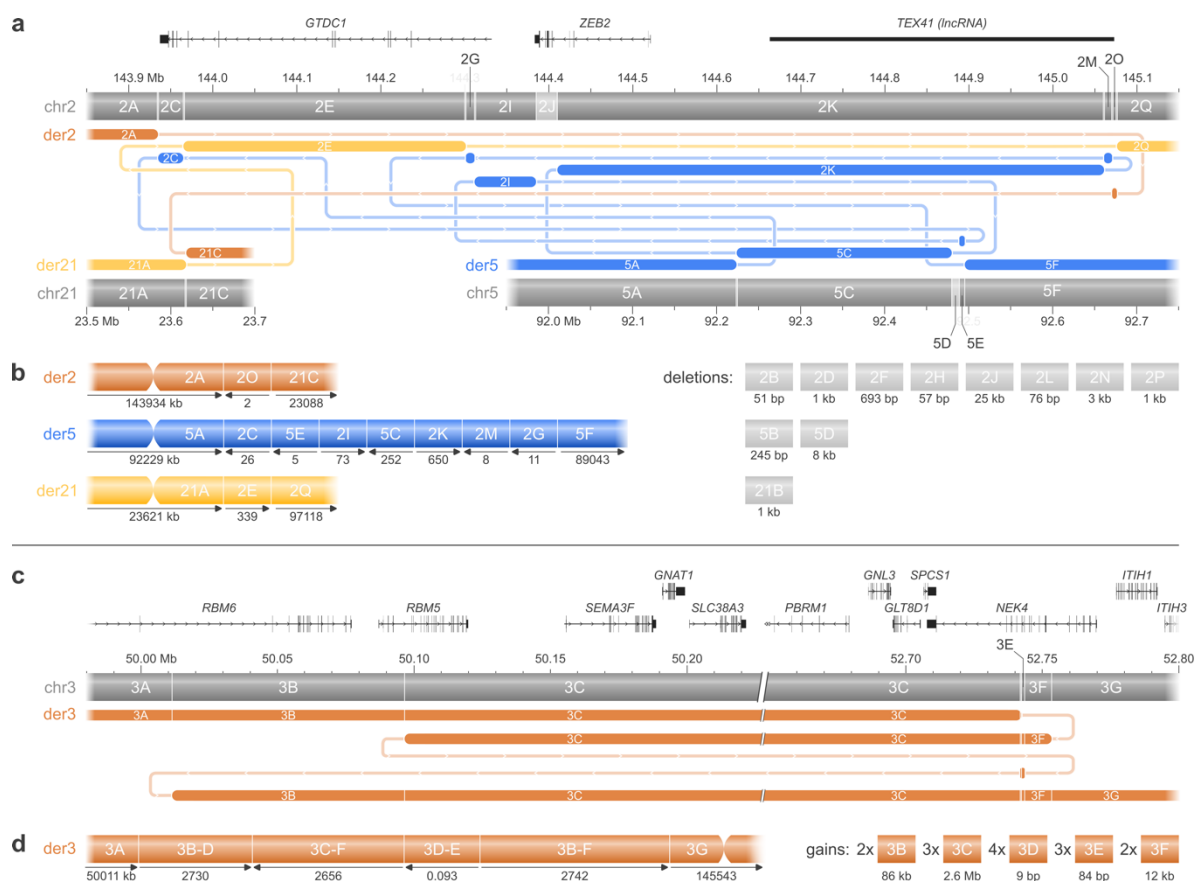


Fig. 2 Simple SVs in individuals S3 and S4 turned out to be more complex variants. Apart from the lncRNA gene *TEX41*, only protein coding genes are shown in the figures. **(a)** Through LRS and OGM a seemingly simple t(2;21) variant in individual S3 was delineated as a complex rearrangement additionally involving chr5 and consisting of 23 breakpoints. The rearrangement shows previously missed insertions of chr2 into chr5 (blue fragments), along with the translocation between chr2 and chr21 (yellow and orange fragments). This delineation allowed for a molecular diagnosis through the disruption of *ZEB2* and confirmed the clinical diagnosis of Mowat-Wilson syndrome. **(b)** Local karyotype of the aberrant chromosomes of individual S3. **(c)** LRS revealed that a previously identified triplication in individual S4 harbored a complex copy number changing rearrangement affecting 67 protein coding genes. **(d)** Local karyotype of the aberrant chromosome 3 of individual S4.

129 originated from regions within the CGR. Other variants presented with no to minimal microhomology
 130 (0 to 2 bp), except for one variant at the 2K-2M junction in derivative 5 where a microhomology stretch
 131 of 9 bp was detected.

132

133 Individual S4 was diagnosed with severe developmental delay, autism spectrum disorder,
 134 microcephaly, and facial dysmorphism. Shallow WGS revealed a *de novo* 2.6 Mb triplication at
 135 chromosome bands 3p21.31p21.1, flanked upstream by a *de novo* 90 kb duplication (table 1; sup. fig.
 136 S2a). This rearrangement was highly suggestive as causal, due to its size and large number of affected
 137 genes and additional diagnostic tests revealing no other molecular aberrations. However, no

138 immediate link could be established between this variant and the observed phenotype, and it was
139 therefore diagnosed as a variant of unknown significance. SV analysis through LRS identified a complex
140 copy-number changing rearrangement on chromosome 3 involving six breakpoints (sup. table S3,S9).
141 This CGR consists of an 86.1 kb duplication, followed by an inverted 2.6 Mb triplication, a 9 bp
142 quadruplication-inversion alongside an 84 bp triplication-inversion, and an inverted 11.6 kb duplication
143 (fig. 2c,d; sup. fig. S5a-e). This rearrangement was fully verified through PCR (sup. table S5). CNV
144 analysis through LRS was able to identify the 2.6 Mb triplication with identical genomic coordinates as
145 previously reported through shallow WGS (sup. table S2). The flanking 86.1 kb duplication was
146 however not called. Microhomology analysis detected no microhomology at the 3D-3F junction, a
147 templated 30 bp inserted sequence at the 3C-3E junction, and a 1 bp microhomology at the 3D-3B
148 junction (sup. text S4). The rearrangement covers 67 protein coding genes of which 17 have an
149 associated OMIM phenotype (sup. table S3,S4). Furthermore, it has breakpoints directly disrupting
150 non-coding regions of *RBM6*, *RBM5*, and *NEK4*, along with a breakpoint falling within exon 13 of *NEK4*.
151 None of these impacted genes can however be unambiguously linked to the observed phenotype.
152 OGM was not further applied as all variants could be easily phased and reconstructed through LRS
153 alone.

154 [Unresolved complex rearrangements are fully delineated, providing insight into their](#) 155 [formation mechanisms](#)

156 The phenotype for individual S5 consisted of ID and autism spectrum disorder. Chromosome analysis
157 revealed a balanced *de novo* reciprocal translocation between chromosome bands 2q36 and 7q32
158 (table 1; sup. fig. S1d). In addition, two heterozygous *de novo* CNVs were detected on chromosome 2
159 through microarray analysis, consisting of a 477 kb deletion in 2q24.2 and a 186 kb deletion in 2q36.3
160 (sup. fig. S2b,c). The 477 kb deletion in the 2q24.2 region encompasses the *TBR1* gene, described in
161 intellectual developmental disorder with autism and speech delay (OMIM #606053). This variant was
162 therefore considered causal for the associated phenotype. Furthermore, the translocation and 186 kb
163 deletion at chromosome band 2q36 were believed to be linked in a CGR. In regard to preconception

164 counseling of a family member, an attempt was made to delineate this region through further research
165 efforts, which however proved fruitless. LRS was therefore applied. Variants were filtered to retain
166 events lying within the long arms of chromosome 2 and 7, manually scanned at the previously defined
167 regions of interest, and further expanded upon finding cryptic variants. This approach generated
168 multiple variants of interest at chromosome band 2q24.2, consisting of one 190 kb inversion, one 839
169 kb deletion, and several translocations to the short arm of chromosome 7 (sup. table S10). The
170 translocations were identified as a novel 347 kb insertional inversion of chromosome 2 into
171 chromosome 7. One of its breakpoints on chromosome 2 coincides with the 839 kb deletion, and these
172 two variants were therefore considered to be part of the same event. Here, a 347 kb part of the 839
173 kb deletion was inserted back into chromosome 7, resulting in a 492 kb loss of material of chromosome
174 2. This is in line with the previous finding of a 477 kb deletion at location 2q24.2. The 190 kb inversion
175 lies 76 kb upstream of this complex variant. Phasing analysis through LRS could not be achieved due
176 to the absence of informative single nucleotide variants, but additional application of OGM confirmed
177 the CGR and identified all variants to be part of the same haplotype (event 1 in fig. 3a; sup. fig. S10a).
178 At location 2q36.3, a 222 kb deletion and several translocations to chromosome band 7q32.2 were
179 found through LRS. The translocations were deemed to be part of the same event that covers the
180 $t(2;7)(q36.3;q32)$ variant identified through karyotyping (event 2 in fig. 3a). The 222 kb deletion
181 present in the same chromosome band 2q36.3 coincides with the previously detected 186 kb deletion
182 at the same location. It was however considered separate from the translocation as it localizes nearly
183 2 Mb further downstream and phasing information could not be retrieved through LRS nor through

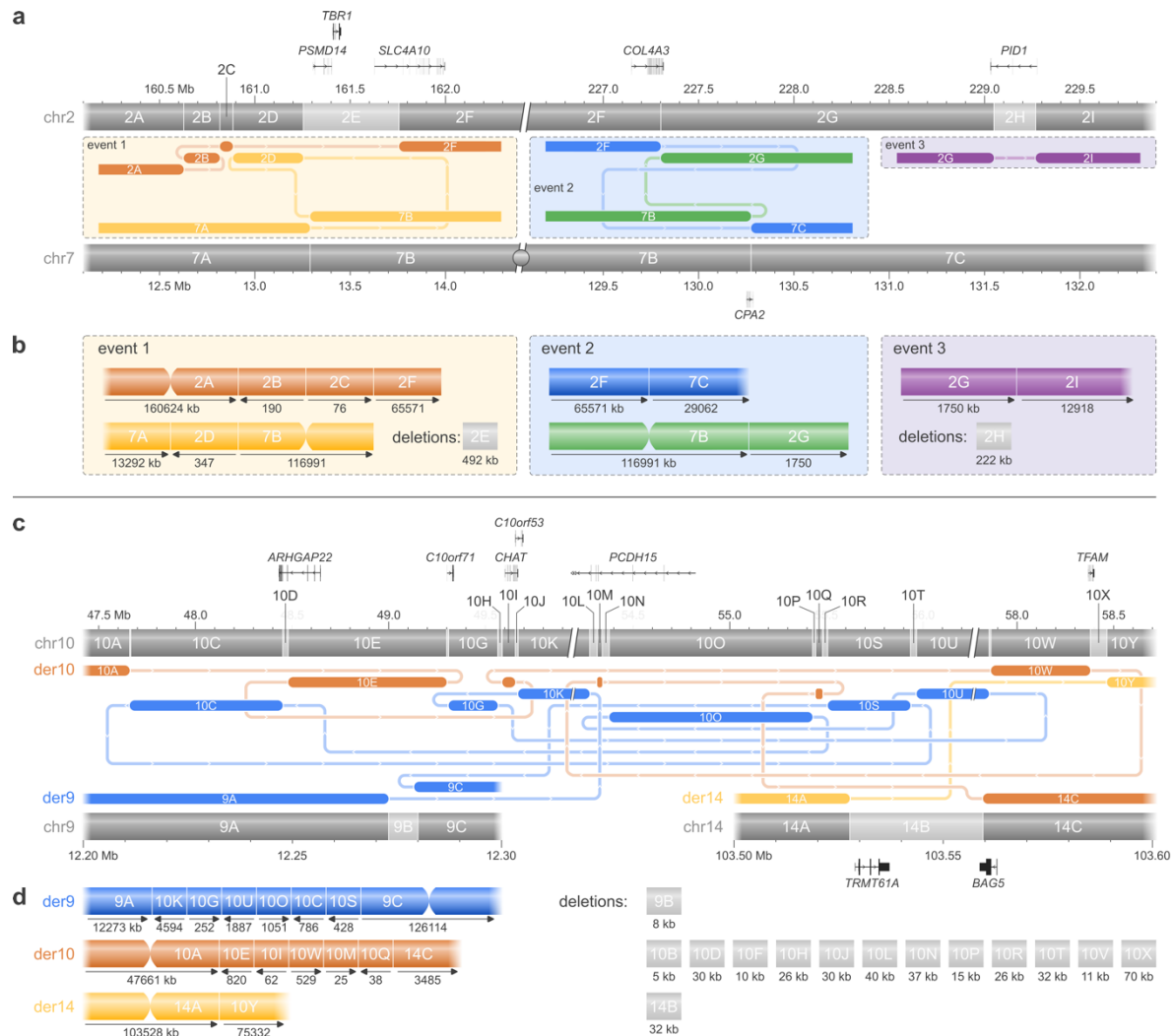


Fig. 3 Delineation of previously identified CGRs in individuals S5 and S6. Only genes directly affected by a breakpoint or deletion are shown in the figures. **(a)** In individual S5 conventional techniques detected two deletions at respectively chromosome bands 2q24.2 and 2q36.3, along with a t(2;7)(q36;q32) reciprocal translocation. The 2q36.3 deletion and the translocation were thought to be connected to each other and arisen through one single complex event. Application of LRS confirmed all previously detected variants, and identified an additional ins(2;7) variant connected to the 2q24.2 deletion and an upstream cryptic inversion (event 1). The t(2;7)(q36;q32) variant and the 2q36.3 deletion were found to be localizing nearly 2 Mb from each other and were thus considered as two distinct events (respectively events 2 and 3). **(b)** Local karyotypes of the three events identified in individual S5. **(c)** Individual S6 presented with a complex karyotype consisting of an ins(9;10) and a t(10;14) variant. LRS and OGM further delineated this CGR up until single base pair resolution and identified a much more intricate rearrangement involving 28 breakpoints on chr9, 10, and 14. It features 10 deletions larger than 15 kb and disrupts in total 10.8 Mb of genetic material, causal for the observed fertility problems. **(d)** Local karyotype of the aberrant chromosomes of individual S6.

184 OGM (event 3 in fig. 3a). OGM confirmed the variants in event 2 and 3 (fig 3a). Delineation of the SVs
 185 previously identified within this sample thus characterized them as three events, with the inversion
 186 and inverted insertion at region 2q24.2 as novel finds (fig. 3a,b; sup. fig. S7b). The two deletions at
 187 chromosome bands 2q24.2 and 2q36.3 were detected through copy number calling as well (sup. table
 188 S2). Microhomology analysis revealed only limited microhomology (0 to 3 bp) at the LRS breakpoint

189 junctions for most variants, except for the 2F-7C junction where a 7 bp microhomology stretch was
190 detected (sup. text S5). At the 2C-2F junction a thymine insertion was identified. The detected SVs
191 disturb four genes (*SLC4A10*, *COL4A3*, *PID1*, *CPA2*) directly through a breakpoint and two other protein
192 coding genes (*PSMD14*, *TRB1*) through a complete loss of genetic material (sup. table S3).

193

194 Individual S6 was referred to the clinic after repeated miscarriages and recurrent implantation failure.
195 Chromosome analysis showed an aberrant complex karyotype consisting of a translocation between
196 chromosome bands 10q21.2 and 14q32.3, and an insertion of chromosome region 10q11.2 to 10q21.2
197 into 9p22 (table 1; sup. fig. S1e-g). These variants were considered causal for the observed phenotype
198 as CGRs can lead to unbalanced gametes. Additional shallow WGS identified five deletions on
199 chromosome 10 ranging from 30 to 120 kb in length and one 30 kb deletion on chromosome 14 (sup.
200 fig. S3). LRS was applied to delineate the full situation. SV analysis exposed additional cryptic SVs (sup.
201 table S11), which were reconstructed back into a 10.8 Mb CGR consisting of 28 breakpoints involving
202 chromosomes 9,10, and 14 (fig. 3c,d; sup. table S3; sup. fig. S7c). Through this manual reconstruction,
203 14 deleted fragments were identified, of which 10 larger than 15 kb (sup. fig. S6). Six of these coincide
204 with the six previously detected deletions. CNV analysis through LRS was however not able to identify
205 any copy number changes in the regions of interest (sup. table S2). Breakpoint junction analysis
206 revealed no inserted sequences between fragments and limited microhomology (1 to 5 bp) in most
207 junctions (sup. text S6). Seven protein coding genes are directly disturbed by a breakpoint (*ARHGAP22*,
208 *C10orf71*, *CHAT*, *C10orf53*, *PCDH115*, *TFAM*, *BAG5*), while an eighth gene (*TRMT61A*) falls entirely
209 within a deleted fragment on chromosome 14 (sup. table S3). OGM was additionally applied to verify
210 the reconstruction of this highly complex rearrangement. This fully confirmed the LRS findings, yet
211 some smaller fragments could not be delineated. Four junctions seen solely through LRS analysis were
212 subsequently successfully confirmed by PCR (sup. table S5).

213

214 Long-read sequencing identifies 100% of breakpoint junctions

215 The majority of SV junctions in this study were detected through use of the Sniffles2 LRS SV caller (87%,
216 62/71, n=6), with a mean deviation to the actual breakpoint location of 3 ± 5 bp (n=76) (sup. table S6-
217 12). The junctions that were not detected through Sniffles2 could still be retrieved from the LRS data
218 through manual curation (sup. fig. S10b), resulting in a 100% pick-up rate through LRS (71/71, n=6).
219 OGM on the other hand, was able to retrieve 67% (41/61) of the junctions directly through automated
220 data analysis, and 74% (45/61) after manual curation (n=3). The average deviation from the actual
221 breakpoint coordinate through OGM was 5.0 ± 5.5 kb (n=66).

222 Discussion

223 In this study we used LRS as the main method for SV characterization and utilized OGM as a subsequent
224 verification method for more complex cases. Through this, we show their potential as a follow-up
225 diagnostic test to delineate unresolved SVs. This is of importance in a diagnostic setting as it enables a
226 molecular diagnosis, influences future disease management, and can expose undiscovered disease-
227 associated genes. Through the application of LRS we successfully characterized multiple clinically
228 relevant rearrangements in six individuals up until single basepair level. This led to a conclusive
229 molecular diagnosis in two individuals, which was previously unattainable through standard diagnostic
230 techniques alone. A third individual remained without a precise diagnosis, although application of LRS
231 may give new insights into the underlying disease mechanism. The rearrangement disturbs exon 13 of
232 *NEK4* and several introns of *RBM6*, *RBM6*, and *NEK4*. None of these genes can currently be linked to
233 intellectual disability yet may constitute novel candidate genes. Alternatively, the newly discovered
234 inverted triplication could explain an alternative pathogenic mechanism through disturbed expression
235 patterns of its impacted genes or regulatory regions [16]. Future functional work to investigate this is
236 however outside the scope of this study.

237

238 The ability of LRS to characterize SVs at single base pair resolution and reveal their exact sequence
239 provides insight into their formation origins without the need for additional molecular methods.
240 Probands S1 and S2 both manifested with a simple SV. Their respective variants lack considerable
241 microhomology at the breakpoints and are lying outside repeat elements, indicating non-homologous
242 end joining (NHEJ) as the most likely formation mechanism. Slightly more complex cases were seen in
243 individuals S4 and S5. The CGR in participant S4 consists of multiple copy-number changes with several
244 junctions lying within repeat elements, hinting towards a replication defect based on the fork stalling
245 and template switching/microhomology-mediated break-induced replication pathway. Indeed,
246 extensive reports exist within literature describing similar DUP-TRIP/INV-DUP events driven by the
247 aforementioned mechanisms [17,18]. A less clear situation was seen in case S5 as not all events (fig.

248 3a,b) can be explained through the same formation mechanism. Non-allelic homologous
249 recombination is typically observed in recurrent SVs yet can also drive non-recurrent events between
250 smaller repetitive sequences [19]. Most fused junctions identified in this case indeed lie within repeat
251 elements of the same family. Event 3 however does not fit this mechanism as only one junction resides
252 within a repetitive region. Microhomology-mediated end joining (MMEJ) could then serve as an
253 alternative explanation, although here again, not all identified junctions adhere to its general
254 characteristics. This may point to two distinct processes driving the formation of the observed variants
255 in this individual instead of one.

256 Individuals S3 and S6 presented with highly intricate rearrangements which seem to originate from
257 separate mechanisms considering their differing characteristics. The CGR in proband S3 has templated
258 inserted sequences between breakpoints, microhomologies up to 7 bp, and no large CNVs apart for
259 one 25 kb deletion. We here propose chromothripsis through MMEJ as a formation mechanism. MMEJ
260 can act as a repair mechanism for double strand breaks such as those generated in mass during
261 chromothripsis yet is less preferred over the more utilized NHEJ pathway [20]. Most commonly it only
262 produces deletions and features loss of nucleotides at the breakpoints to reveal microhomology
263 necessary for mediating its mechanism. Indeed, flanking deletions of 1 kb or smaller can be observed
264 between junctions (fig. 2b). Case S6 has limited microhomologies at variant junctions, no insertions
265 between breakpoints, and several deletions larger than 15 kb. These observations are consistent with
266 chromothripsis through NHEJ [21].

267

268 The complex rearrangements detected in S3-S6 were followed up with OGM to confirm their
269 reconstruction achieved through LRS and to aid with haplotype phasing. Both techniques act
270 complementary to each other and enhance their respective qualities. LRS reads offer single base pair
271 resolution, yet only typically capture one SV or one SV junction per read. On the other hand, OGM
272 molecules often span multiple variants in one read and as such, can naturally phase variants together.
273 Due to this difference in technology, LRS and OGM detect rearrangements in a significantly different

274 manner. LRS unveils breakpoints and breaks the affected region up in fragments (sup. fig. S9a) that
275 have to be manually puzzled back together, making CGR reconstruction a feasible yet laborious task.
276 However, OGM detects variants in an aberrant haplotype approach (sup. fig. S9b,c). OGM thus allows
277 for easy phasing and reconstruction of variants, while LRS adds single base pair resolution. As
278 demonstrated in this study, their parallel application is therefore especially powerful for CGR
279 characterization and allows for the easy unraveling of highly complex rearrangements. Note however
280 that OGM glosses over intricate details within the more complex CGRs due to its resolution of 5 kb.
281 While OGM detects 74% of the junctions within its applied complex cases, this was 100% with LRS
282 alone. Application of each technology as a stand-alone method may thus manifest as a trade-off
283 between detailed resolution versus straightforward variant reconstruction in complex situations.
284 Historically, CGRs have been considered as rare constitutional events, yet recent studies indicate they
285 are more common than previously anticipated [18]. This is reflected in our study as we find highly
286 complex rearrangements in four out of six included individuals. Especially of interest is case S6. This
287 individual is phenotypically normal apart from reproductive issues, despite the presence of a
288 rearrangement consisting of 28 breakpoints affecting eight protein coding genes. CGRs are indeed
289 occasionally encountered in individuals with fertility problems who are otherwise healthy [22], and
290 might be more common than anticipated around reciprocal translocations [23]. A similar observation
291 was recently made by Eisfeldt *et al.* (2021) [24], who unraveled a CGR involving 137 breakpoints
292 through OGM and LRS, among other technologies. Their application may thus lead to more informed
293 decisions in settings where unidentified SVs can result in undesirable outcomes, such as in fertility
294 assistance. Preimplantation embryos have a tendency to form micronuclei, which in turn are
295 susceptible to the formation of complex rearrangements through chromoanagenesis [25]. Current
296 clinical technologies however underestimate both the incidence and complexity of CGRs and may thus
297 fail to identify these events [18], as demonstrated here as well. LRS and OGM can aid in more
298 accurately assessing the implantation success of these embryos and provide more substantiated
299 reproductive decisions. Especially LRS could prove beneficial in this setting, considering its ability to

300 assess the formation origins of several chromoanagenesis cases in this study without the need for
301 additional data sources.

302

303 In this manuscript, we furthermore include detailed laboratory and bioinformatics protocols that can
304 be utilized as an initial framework to implement LRS and OGM in a clinical setting. These can be applied
305 to identify genome-wide structural variation in LRS and OGM data. We show our protocols can be
306 employed in heterogeneous conditions, as they were successfully applied on divergent coverages,
307 instruments, and input materials. However, the differing experimental states mean we cannot
308 mutually compare the included samples to each other considering their overall quality statistics and
309 variant calls. In addition, separate CNV detection through LRS data failed to pick up a considerable
310 number of verified variants. This inability to detect CNVs through a separate analysis on the same LRS
311 dataset highlights an important hiatus in current SV detection tools. Due to the lack of a dedicated
312 LRS CNV calling algorithm, we used a CNV caller originally designed for short reads. This could fail to
313 detect signals inherent to LRS alone, resulting in missed variants. We therefore recommend cautious
314 result interpretation of these tools and propose a prior benchmarking on a verified set of CNVs to
315 assess their performance before employment in a diagnostic setting.

316

317 A limitation of this study is its small cohort size. While LRS and OGM were able to assess the clinically
318 relevant SV in all cases, the limitations of the isolated use of these techniques could still interfere with
319 successfully making a molecular diagnosis. Studies leveraging larger sample sizes however show similar
320 promising results for the application of both LRS and OGM in a clinical setting [9],[10]. Furthermore,
321 the current study heavily benefits from a priori knowledge obtained through previous tests. It has to
322 be taken into account that if LRS or OGM were to be applied blind, then stringent filtering needs to be
323 in place to identify the causal variant out of the myriad of SVs present within the human genome. A
324 crucial aspect of implementing a diagnostic test is namely its ability to distinguish between benign and
325 malignant genomic aberrations. Limited filtering based on an internal database is already in place

326 within the Bionano pipeline, yet its impact remains to be evaluated as studies applying this approach
327 in a blind diagnostic setting are lacking. Especially in the context of LRS, these studies would need to
328 rely on internal and external databases of common SVs present within the population to filter out
329 irrelevant variation. Few such databases exist [26] and ideally need to be extended with data from
330 emerging technologies that excel at SV characterization. We therefore envision that wide-scale
331 employment of these technologies in a diagnostic setting would further facilitate the sharing of these
332 events and thus their filtering and interpretation.

333

334 In conclusion, we show that nanopore LRS coupled with Bionano OGM offers several benefits
335 compared to the currently employed diagnostic techniques used for SV identification. Not only are
336 they able to detect previously identified SVs, they also allow for a complete characterization of
337 currently unresolved SVs up to base-level resolution and the detection of cryptic rearrangements
338 missed by conventional genetic tests. Furthermore, LRS provides valuable insights in the formation
339 mechanisms behind an SV and works in synergy with OGM to easily delineate CGRs, a currently
340 underestimated variant class. Most importantly, the ability of LRS to pinpoint exact breakpoint
341 locations enables molecular diagnoses, resulting in a modest higher diagnostic yield.

342 Materials and Methods

343 Study design

344 Six individuals (S1-6) were selected for whom an SV was identified through standard clinical diagnostic
345 methods, yet exact breakpoint coordinates could not be established. The structural variation events
346 were determined to be causal for the underlying phenotype, except for individuals S1, S3, and S4
347 where no conclusive molecular diagnosis could be made. Standard testing included a selection of G-
348 banded chromosome analysis, FISH, whole exome sequencing, shallow WGS, qPCR, and microarray
349 analysis. For four individuals a simple SV event consisting of maximum two breakends was identified,
350 including two translocations (S1,3), a deletion (S2), and a triplication (S4). For two individuals (S5,6) a
351 rearrangement consisting of three or more breakends, referred to as a CGR, was detected. All
352 individuals presented with ID or DD, except for S6 who came to the clinic due to reduced fertility. All
353 SVs occurred *de novo*, except for S6 for which the *de novo* status could not be further evaluated
354 through segregation analysis. All clinical diagnoses and SVs identified through standard clinical testing
355 can be consulted in table 1. More detailed information on the observed phenotypes and specific
356 conducted genetic tests per individual can be found in supplementary methods.

357 Long-read whole genome sequencing

358 Sample preparation and sequencing

359 High molecular weight (HMW) DNA was extracted from blood or lymphoblastoid cell lines in an
360 automated manner using respectively the Genomic DNA Large Volume Whole Blood Kit (1.200µl; RBC
361 Bioscience) and the Genomic DNA Cultured Cells Kit (RBC Bioscience). Quantification was done with a
362 Qubit fluorometer (Thermo Fisher) and the dsDNA HS Assay kit (Thermo Fisher). DNA was
363 subsequently sheared to 20-30 kb lengths by placing g-TUBES (Covaris) in a centrifuge at 3200 rpm for
364 four minutes, or size selected to enrich for sizes of 10 kb or greater using the Short Read Eliminator XS
365 kit (PacBio) (sup. table S1). Approximately 7 µg of processed DNA was used as input for the EXP-
366 BND104, SQK-LSK109, or the SQK-LSK114 library preparation kits (ONT) following the manufacturer's

367 **Table 1** Phenotypical features of the six individuals (S1-6) included in this study, along with their
 368 unresolved SVs identified through standard clinical diagnostic testing. Individuals S1-4 presented with
 369 a simple structural variation event, while in individuals S5 and S6 a CGR was identified. Molecular
 370 diagnoses were made for all participants, except for individuals S1, S3, and S4.

individual	phenotypical features	identified SVs
S1	severe ID ¹ facial dysmorphism	46,XY,t(9;10)(q21.2;p15.2)
S2	mild ID Down-like features ASD ²	NC_000002.12(NM_001303052.2):c.(152+1_153-1)_(504+1_505-1)del
S3	Mowat-Wilson syndrome severe DD ³	46,XY,t(2;21)(q22;q21)
S4	ASD mild microcephaly	sseq[GRCh38]3p21.31p21.1(50100001_52755000)x3 dn
S5	ID ASD	46,XY,t(2;7)(q36;q32) arr[GRCh37]2q24.2(162105681_162582744)x1 dn arr[GRCh37]2q36.3(229936820_230122602)x1 dn
S6	low egg count repeated miscarriages implantation failures	46,XX,ins(9;10)(p22;q11.2q21.2),t(10;14)(q?;q32.3)

371 ¹ ID: intellectual disability; ² ASD: autism spectrum disorder; ³ DD: developmental delay

372
 373 protocol with minor adjustments. Briefly, modifications consisted of longer incubation times and
 374 elution in higher volumes to maximize the retention of HMW DNA fragments. Libraries were then
 375 sequenced on a MinION or PromethION 24 device (ONT) for 80 hours, with flushing and reloading after
 376 24 and 48 hours to boost the final yield. Flow cells were of type R9.4.1 or R10.4.1 depending on the
 377 used library preparation kit. Loading quantities in fmol were calculated based on size profiles obtained
 378 from a TapeStation 4150 device (Agilent) using the Genomic DNA ScreenTape kit (Agilent). Additional
 379 library preparation and sequencing was done for CGR samples where the full rearrangement could not
 380 be reconstructed and the coverage was below 15x (coverage cut-off determined through internal data,
 381 paper in preparation). More detailed information on the processing and quality statistics of each
 382 sample can be found in supplementary table S1.

383 Sequencing analysis and structural variant detection

384 Raw sequencing data were basecalled with Guppy v6.3.7 (ONT) using the super accuracy model
385 without q-score filtering. Demultiplexing was done using the Guppy v6.3.7 barcoder (ONT). Reads were
386 aligned to the hg38 reference genome with minimap2 v2.24 [27], and split libraries were merged into
387 one alignment file using samtools v1.15 [28]. Coverage depth was calculated using Mosdepth v0.3.3
388 [29]. Other quality control metrics were checked with NanoPlot v1.40.0 [30] and PycoQC v2.5.2 [31].
389 SVs were called using Sniffles2 v2.0.7 [32] with a read support parameter of one, and further processed
390 to only retain variants with a read support of three or higher and a length of at least 50 bp. Resulting
391 alignment and variant files were filtered to relevant genomic regions based on previous diagnostic
392 tests using respectively samtools v1.15 [28] and vcftools v0.1.16 [33]. Both were then scanned
393 manually using IGV v2.13.2 [34] for SVs of interest, with special attention paid near identified
394 breakpoints to allow for potentially undetected CGRs. If no variants of interest were found or a CGR
395 was suspected, the filter criteria were loosened to allow for SVs at lower read supports and analysis
396 was repeated. The total number of SVs in a sample was calculated as all variants outputted by Sniffles2
397 with a length of at least 50 bp and with an aberrant read support of three or higher. Detailed
398 commands regarding this workflow can be found in supplementary methods.

399 Copy number variant detection

400 Large CNVs (≥ 15 kb) were called from the LRS data using WisecondorX v1.2.5 [35]. A custom reference
401 set was built from 36 LRS samples which were sequenced in-house on a PromethION 24 instrument
402 (ONT) on R9.4.1 flow cells. Reference samples were size selected using the Short Read Eliminator XS
403 kit (PacBio) and prepared for sequencing with the SQK-LSK109 kit (ONT). Reference samples were
404 further processed according to the workflow described in this study and had an average coverage of
405 $25.8x \pm 6.0x$ and an average N50 of $20.8 \text{ kb} \pm 5.1 \text{ kb}$. The bin size of the reference was set to 15 kb.
406 WisecondorX was then executed according to the manual instructions to call CNVs on the LRS data of
407 individuals S1-6. Further filtering was applied to only retain CNV events with a log₂ ratio lower than -
408 0.50 and higher than 0.35. Subsequently, events lying within centromeric regions were removed. CNVs

409 were further evaluated and solely reported if present within the region of interest as defined by
410 previous diagnostic tests. Exact commands can be found in supplementary methods. Intermediate
411 filtering results are listed in supplementary table S2.

412 [Phasing through single nucleotide detection](#)

413 Phasing was applied to samples where the SVs identified through LRS could not be accurately
414 reconstructed into one haplotype. Basecalled fastq files were filtered on a q-score of 10 or higher with
415 NanoFilt v2.8.0 [30] and subsequently mapped with minimap2 v2.24 [27] to the hg38 reference
416 genome. SNVs were then called through Clair3 v1.0.2 [36]. The resulting variant file was used as input
417 for WhatsHap v1.7 [37] to phase the long reads into haplotype groups. Aligned reads were haplotagged
418 with the same algorithm to allow visualization in IGV v2.13.2 [34]. Detailed information on the used
419 commands can be found in supplementary methods.

420 [Optical genome mapping](#)

421 [Sample preparation and molecule imaging](#)

422 OGM was applied to a subset of individuals (S3,5,6) to verify the arrangement detected through LRS
423 and to provide further phasing information. Ultra HMW DNA was extracted from frozen cell pellets (-
424 80°C) gathered from lymphoblastoid cell lines according to manufacturers' instructions in the SP
425 Frozen Cell Pellet Isolation Protocol and the SP Blood and Cell Culture DNA isolation kit (Bionano
426 Genomics). In short, frozen samples were thawed in a 37°C warm water bath and approximately 1.5
427 million cells were collected. Cells were centrifuged at 2200 g for 2 minutes at 4°C and subsequently
428 washed with cold DNA stabilizing buffer (Bionano Genomics). Cells were then digested in the presence
429 of proteinase K, and lysate was transferred to a nanobind disk. After washing, HMW DNA was eluted
430 and incubated overnight at 25°C to homogenize. For each sample 750 ng of extracted DNA was labelled
431 using the Direct Labeling Enzyme (DLE-1) following the Direct Label and Stain kit protocol (Bionano
432 Genomics). This DLE-1 enzyme tags the 6 bp CTTAAG sequence that occurs approximately every 5 kb
433 in the human genome with a green fluorophore. The sample was further stained for backbone
434 visualization and subsequently loaded onto a Saphyr instrument with a G2.3 chip (Bionano Genomics)

435 for molecule imaging. The instrument had a run time between 8 to 16 hours, depending on if the
436 manufacturer's minimum aim of 320 Gb throughput was obtained, except for sample S5 which ran for
437 72 hours. DNA was quantified after extraction and labelling using respectively the dsDNA BR and HS
438 Assay kit (Thermo Fisher) on a Qubit fluorometer (Thermo Fisher). Mixing of samples was done using
439 a wide bore pipette tip to maximize the preservation of long DNA fragments. Quality metrics were
440 assessed according to the manufacturer's guidelines and can be found in supplementary table S1.

441 [De novo assembly and structural variant detection](#)

442 Molecule data from the Saphyr instrument were analyzed using the Bionano Solve v3.5.1 tool (Bionano
443 Genomics). The De Novo Assembly pipeline v10322 was run on all three OGM samples, where a diploid
444 assembly is constructed from the analyzed molecules and subsequently aligned to the hg38 reference
445 optical map. Structural aberrations are then detected by identifying differences between the tagged
446 sequence motifs of the two aligned genomes. Subsequent visualization, filtering, and interpretation of
447 the identified SVs was done with Bionano Access v1.7.2 (Bionano Genomics). Standard filter settings
448 were adjusted to select variants only present in the regions of interest as identified through previous
449 standard diagnostic tests and in less than 1% of control samples. The total number of SVs identified in
450 a sample was determined from the variant vcf-files as outputted by the De Novo Assembly Pipeline
451 without any additional data filtering.

452 [Variant confirmation](#)

453 The nanopore LRS data of variants of interest were manually inspected in IGV v2.13.2 [34] to
454 determine the exact breakpoint locations of each variant. Genomic coordinates were set to the first
455 nucleotide immediately upstream of where the majority of the LRS reads switched reference genome
456 mapping locations, i.e. from local mapping to the reference genome to mapping to a distal location.
457 Simple structural variation events (S1,2) were then confirmed with Sanger sequencing. For SVs within
458 a CGR and identified through LRS alone or without parallel strong OGM data support, an additional
459 confirmation through PCR and resulting fragment length was executed. SVs within a CGR and identified
460 through both LRS and OGM were considered concordant if their breakpoints were within 25 kb of each

461 other and had the same strand orientation. Concordant events were not further confirmed through an
462 additional technique. SVs identified through OGM without parallel LRS data support were not
463 encountered in this study. The aberrant haplotypes in samples with a CGR were reconstructed
464 manually by finding the path that leads to one aberrant chromosome and one structurally normal
465 chromosome. Resulting subway plots were drawn by hand using Inkscape v1.2 [38].

466 [Assessment of structural variant origins](#)

467 For each variant a 50 bp consensus sequence surrounding the SV breakpoint was extracted from the
468 nanopore LRS data. Similarly, a 50 bp sequence surrounding the proximal and distal SV breakpoint was
469 selected from the reference genome. These three sequences were aligned to each other using Clustal
470 Omega v1.2.4 [39] to determine microhomology patches. The regions were extended with another
471 100 to 200 bp if no microhomology between the three sequences could be detected. If insertions
472 larger than 20 bp were found between breakpoint junctions, this sequence was compared against
473 human reference genome hg38 through blastn [40] to find its initial origin. Furthermore, the genomic
474 region around the breakpoints was investigated for potential repeat elements using the UCSC Genome
475 Browser [41].

476 [Technology evaluation](#)

477 The accuracy on breakpoint locations of variants of interest was evaluated for both LRS and OGM. For
478 LRS the breakpoint coordinates as identified through Sniffles2 SV calling were compared to the
479 manually verified variant breakpoint coordinates. For OGM data the breakpoint locations as seen in
480 the Bionano Access software were compared to the nearest manually verified variant breakpoint
481 coordinates with equal strand orientation. Breakpoint locations were manually extracted as an OGM
482 read can span multiple variants, yet only reports two breakpoints per molecule. The manually
483 extracted OGM breakpoints were then set to the genomic position of the last fluorophore tag before
484 the successive reference mapping was disrupted. Comparison between nanopore LRS and Bionano

485 OGM variant data was achieved by comparing the breakpoint coordinates for each aberrant fragment
486 recombination.

487 References

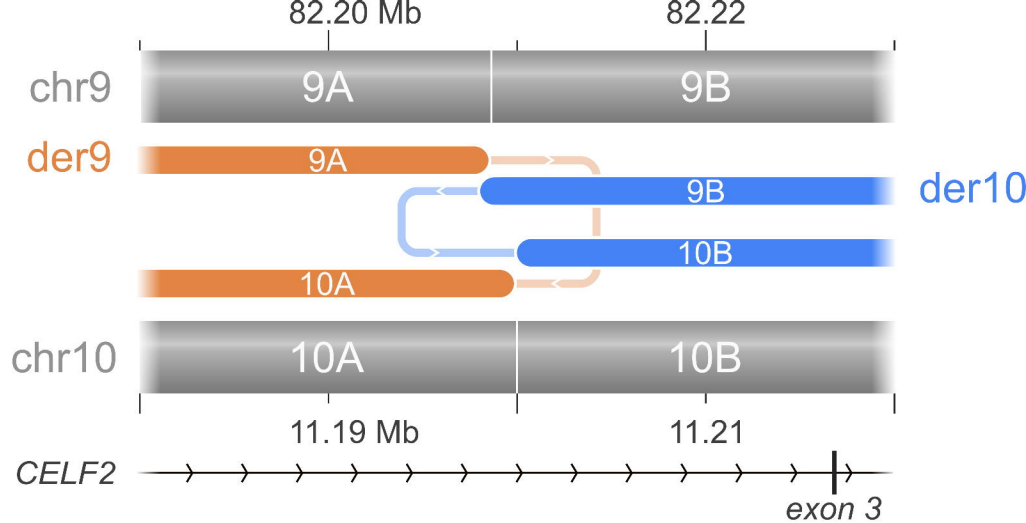
- 488 1. Weischenfeldt, J., Symmons, O., Spitz, F. & Korbel, J. O. Phenotypic impact of genomic
489 structural variation: Insights from and for human disease. *Nat. Rev. Genet.* **14**, 125–
490 138 (2013).
- 491 2. Porubsky, D. & Eichler, E. E. A 25-year odyssey of genomic technology advances and
492 structural variant discovery. *Cell* **187**, 1024–1037 (2024).
- 493 3. Gröbner, S. N. *et al.* The landscape of genomic alterations across childhood cancers.
494 *Nature* **555**, 321–327 (2018).
- 495 4. Sebat, J. *et al.* Strong association of de novo copy number mutations with autism.
496 *Science* **316**, 445–449 (2007).
- 497 5. Cooper, G. M. *et al.* A copy number variation morbidity map of developmental delay.
498 *Nat. Genet.* **43**, 838–846 (2011).
- 499 6. Balachandran, P. & Beck, C. R. Structural variant identification and characterization.
500 *Chromosome Res.* **28**, 31–47 (2020).
- 501 7. Wright, C. F., FitzPatrick, D. R. & Firth, H. V. Paediatric genomics: Diagnosing rare
502 disease in children. *Nat. Rev. Genet.* **19**, 253–268 (2018).
- 503 8. Ho, S. S., Urban, A. E. & Mills, R. E. Structural variation in the sequencing era. *Nat. Rev.*
504 *Genet.* **21**, 171–189 (2019).
- 505 9. Mantere, T. *et al.* Optical genome mapping enables constitutional chromosomal
506 aberration detection. *Am. J. Hum. Genet.* **108**, 1409–1422 (2021).
- 507 10. Miller, D. E. *et al.* Targeted long-read sequencing identifies missing disease-causing
508 variation. *Am. J. Hum. Genet.* **108**, 1436–1449 (2021).

- 509 11. Bocklandt, S., Hastie, A. & Cao, H. Bionano Genome Mapping: High-Throughput, Ultra-
510 Long Molecule Genome Analysis System for Precision Genome Assembly and Haploid-
511 Resolved Structural Variation Discovery. *Adv. Exp. Med. Biol.* **1129**, 97–118 (2019).
- 512 12. Logsdon, G. A., Vollger, M. R. & Eichler, E. E. Long-read human genome sequencing and
513 its applications. *Nat. Rev. Genet.* **21**, 597–614 (2020).
- 514 13. Wang, Y., Zhao, Y., Bollas, A., Wang, Y. & Au, K. F. Nanopore sequencing technology,
515 bioinformatics and applications. *Nat. Biotechnol.* **39**, 1348–1365 (2021).
- 516 14. Vergult, S. *et al.* Mate pair sequencing for the detection of chromosomal aberrations in
517 patients with intellectual disability and congenital malformations. *Eur. J. Hum. Genet.*
518 **22**, 652–659 (2014).
- 519 15. Itai, T. *et al.* De novo variants in CELF2 that disrupt the nuclear localization signal cause
520 developmental and epileptic encephalopathy. *Hum. Mutat.* **42**, 66–76 (2021).
- 521 16. Kosuthova, K. & Solc, R. Inversions on human chromosomes. *Am. J. Med. Genet.* **191**,
522 672–683 (2023).
- 523 17. Carvalho, C. M. B. & Lupski, J. R. Mechanisms underlying structural variant formation
524 in genomic disorders. *Nat. Rev. Genet.* **17**, 224–238 (2016).
- 525 18. Schuy, J., Grochowski, C. M., Carvalho, C. M. B. & Lindstrand, A. Complex genomic
526 rearrangements: an underestimated cause of rare diseases. *Trends Genet.* **38**, 1134–
527 1146 (2022).
- 528 19. Verdin, H. *et al.* Microhomology-Mediated Mechanisms Underlie Non-Recurrent
529 Disease-Causing Microdeletions of the FOXL2 Gene or Its Regulatory Domain. *PLoS*
530 *Genet.* **9**, 1003358; [10.1371/journal.pgen.1003358](https://doi.org/10.1371/journal.pgen.1003358) (2013).

- 531 20. Ottaviani, D., LeCain, M. & Sheer, D. The role of microhomology in genomic structural
532 variation. *Trends Genet.* **30**, 85–94 (2014).
- 533 21. Korbelt, J. O. & Campbell, P. J. Criteria for inference of chromothripsis in cancer
534 genomes. *Cell* **152**, 1226–1236 (2013).
- 535 22. Madan, K., Nieuwint, A. W. M. & Van Bever, Y. Recombination in a balanced complex
536 translocation of a mother leading to a balanced reciprocal translocation in the child.
537 Review of 60 cases of balanced complex translocations. *Hum. Genet.* **99**, 806–815
538 (1997).
- 539 23. Gajeccka, M. *et al.* Unexpected complexity at breakpoint junctions in phenotypically
540 normal individuals and mechanisms involved in generating balanced translocations
541 t(1;22)(p36;q13). *Genome Res.* **18**, 1733–1742 (2008).
- 542 24. Eisfeldt, J. *et al.* Hybrid sequencing resolves two germline ultra-complex chromosomal
543 rearrangements consisting of 137 breakpoint junctions in a single carrier. *Hum. Genet.*
544 **140**, 775–790 (2021).
- 545 25. Koltsova, A. S. *et al.* On the complexity of mechanisms and consequences of
546 chromothripsis: An update. *Front. Genet.* **10**, 446661; 10.3389/fgene.2019.00393
547 (2019).
- 548 26. Collins, R. L. *et al.* A structural variation reference for medical and population genetics.
549 *Nature* **581**, 444–451 (2020).
- 550 27. Li, H. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,
551 3094–3100 (2018).
- 552 28. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience* **10**, 33590861;
553 10.1093/gigascience/giab008 (2021).

- 554 29. Pedersen, B. S. & Quinlan, A. R. Mosdepth: quick coverage calculation for genomes
555 and exomes. *Bioinformatics* **34**, 867–868 (2018).
- 556 30. De Coster, W. & Rademakers, R. NanoPack2: population-scale evaluation of long-read
557 sequencing data. *Bioinformatics* **39**, (2023).
- 558 31. Leger, A. & Leonardi, T. pycoQC, interactive quality control for Oxford Nanopore
559 Sequencing. *J. Open Source Softw.* **4**, 1236 (2019).
- 560 32. Smolka, M. *et al.* Detection of mosaic and population-level structural variants with
561 Sniffles2. *Nat. Biotechnol.* **2**, 38168980; 10.1038/s41587-023-02024-y (2024).
- 562 33. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158
563 (2011).
- 564 34. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
- 565 35. Raman, L., Dheedene, A., De Smet, M., Van Dorpe, J. & Menten, B. WisecondorX:
566 improved copy number detection for routine shallow whole-genome sequencing.
567 *Nucleic Acids Res.* **47**, 1605–1614 (2019).
- 568 36. Zheng, Z. *et al.* Symphonizing pileup and full-alignment for deep learning-based long-
569 read variant calling. *Nat. Comput. Sci.* **2**, 797–803 (2022).
- 570 37. Martin, M. *et al.* WhatsHap: fast and accurate read-based phasing. Preprint at:
571 <https://www.biorxiv.org/content/10.1101/085050v2> (2016).
- 572 38. Inkscape Project. Inkscape. <https://inkscape.org>.
- 573 39. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence
574 alignments using Clustal Omega. *Mol. Syst. Biol.* **7**, 539 (2011).
- 575 40. Zhang, Z., Schwartz, S., Wagner, L. & Miller, W. A greedy algorithm for aligning DNA
576 sequences. *J. Comput. Biol.* **7**, 203–214 (2000).

- 577 41. Kent, W. J. *et al.* The Human Genome Browser at UCSC. *Genome Res.* **12**, 996–1006
578 (2002).
579

a**b****c**