



Long-read genome sequencing and variant reanalysis increase diagnostic yield in neurodevelopmental disorders

Susan M Hiatt, James MJ Lawlor, Lori H Handley, et al.

Genome Res. published online September 19, 2024
Access the most recent version at doi:[10.1101/gr.279227.124](https://doi.org/10.1101/gr.279227.124)

P<P	Published online September 19, 2024 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Long-read genome sequencing and variant reanalysis increase diagnostic yield in**
2 **neurodevelopmental disorders**

3
4 Susan M. Hiatt^{1*}, James M.J. Lawlor¹, Lori H. Handley¹, Donald R. Latner¹, Zachary T.
5 Bonnstetter¹, Candice R. Finnila¹, Michelle L. Thompson¹, Lori Beth Boston¹, Melissa Williams¹,
6 Ivan Rodriguez Nunez¹, Jerry Jenkins¹, Whitley V. Kelley¹, E. Martina Bebin², Michael A.
7 Lopez^{2,3,4}, Anna C. E. Hurst⁴, Bruce R. Korf⁴, Jeremy Schmutz¹, Jane Grimwood¹, Gregory M.
8 Cooper^{1*}

9
10 ¹HudsonAlpha Institute for Biotechnology, Huntsville, AL, 35806, USA

11 ²Department of Neurology, University of Alabama at Birmingham, Birmingham, AL, 35924, USA

12 ³Department of Pediatrics, University of Alabama at Birmingham, Birmingham, AL, 35924, USA

13 ⁴Department of Genetics, University of Alabama at Birmingham, Birmingham, AL, 35924, USA

14
15 [*shiatt@hudsonalpha.org](mailto:shiatt@hudsonalpha.org), gcooper@hudsonalpha.org, 601 Genome Way, Huntsville, AL,
16 35806, USA, 256-327-9490

17
18
19 **Running Title:**

20 Increased diagnostic yield with long-reads

21
22 **Key Words:**

23 LRS Special Issue, Long read sequencing, Clinical sequencing, neurodevelopmental disorder,
24 structural variation, repeat expansion

25

26 **ABSTRACT**

27 Variant detection from long-read genome sequencing (lrGS) has proven to be more accurate
28 and comprehensive than variant detection from short-read genome sequencing (srGS).
29 However, the rate at which lrGS can increase molecular diagnostic yield for rare disease is not
30 yet precisely characterized. We performed lrGS using Pacific Biosciences “HiFi” technology on
31 96 short-read-negative probands with rare diseases that were suspected to be genetic. We
32 generated hg38-aligned variants and *de novo* phased genome assemblies, and subsequently
33 annotated, filtered, and curated variants using clinical standards. New disease-relevant or
34 potentially relevant genetic findings were identified in 16/96 (16.7%) probands, nine of which
35 (8/96, ~9.4%) harbored pathogenic or likely pathogenic variants. Nine probands (~9.4%) had
36 variants that were accurately called in both srGS and lrGS and represent changes to clinical
37 interpretation, mostly from recently published gene-disease associations. Seven cases included
38 variants that were only correctly interpreted in lrGS, including copy-number variants, an
39 inversion, a mobile element insertion, two low-complexity repeat expansions, and a 1 bp
40 deletion. While evidence for each of these variants is, in retrospect, visible in srGS, they were
41 either not called within srGS data, were represented by calls with incorrect sizes or structures,
42 or failed quality-control and filtration. Thus, while reanalysis of older srGS data clearly increases
43 diagnostic yield, we find that lrGS allows for substantial additional yield (7/96, 7.3%) beyond
44 srGS. We anticipate that as lrGS analysis improves, and as lrGS datasets grow allowing for
45 better variant frequency annotation, the additional lrGS-only rare disease yield will grow over
46 time.

47

48

49

50 INTRODUCTION

51 Although genome and exome sequencing (GS/ES) are increasingly used to identify molecular
52 diagnoses for rare diseases, reported diagnostic rates range from 20-60% (Srivastava et al.
53 2019; Baxter et al. 2022), indicating that many conditions suspected to be genetic remain
54 refractory to genomic testing. While some tested individuals may have phenotypes resulting
55 from polygenic and/or environmental risk factors (e.g., Niemi et al. 2018), a subset of
56 undiagnosed cases likely result from genetic factors that we are as-yet unable to identify. It is
57 well-known that short-read genome sequencing (srGS) has poor sensitivity to many types of
58 variants, especially structural variants (SVs) and variants affecting repetitive sequences
59 (Wenger et al. 2019; Sanghvi et al. 2018; Mahmoud et al. 2024). Long-read genome
60 sequencing (lrGS), in contrast, has been shown to greatly improve sensitivity to many of the
61 variants missed by srGS (Logsdon et al. 2020), in addition to facilitating *de novo* assemblies to
62 allow for more effective evaluation of structural variation (Cheng et al. 2021). Accordingly, lrGS
63 has great potential to improve rare disease diagnostic testing and has been applied to several
64 rare disease cohorts (Cohen et al. 2022; Miller et al. 2021; Hiatt et al. 2021).

65 In addition to changes in sequencing technology, the scope of knowledge about genes
66 and our ability to annotate genetic variants has steadily increased. As such, systematic
67 reanalysis of GS/ES data can also lead to the discovery of previously overlooked clinically
68 relevant variants. Diagnostic yield increases from reanalysis have been reported to range from
69 4-31% depending on a variety of factors, most notably time since the previous analysis (Hiatt et
70 al. 2018; Liu et al. 2019; Schobers et al. 2022; Hartley et al. 2023). While a variety of factors
71 contribute to reanalysis discoveries, they often result from the discovery of new disease genes,
72 which contributes to 42-75% of reanalysis findings (Hiatt et al. 2018; Liu et al. 2019; Schobers et
73 al. 2022; Hartley et al. 2023). This reflects the rapid pace of discovery of new disease genes in

74 the rare disease research community, which has been facilitated by data sharing via the
75 MatchMaker Exchange and GeneMatcher (Philippakis et al. 2015; Sobreira et al. 2017).

76 Here we discuss findings from lrGS on a cohort of 96 short-read-negative cases, drawn
77 from several studies focused on rare, suspected genetic diseases, especially early-onset
78 neurodevelopmental disorders. We describe 19 relevant or potentially clinically relevant variants
79 not previously evaluated or considered in 16 cases. We show that both more comprehensive
80 variant detection from lrGS and variant reanalysis contribute to these discoveries. However,
81 lrGS clearly provides unique advantages, and these advantages are likely to grow in the future,
82 to maximize the rate of discovery of highly penetrant variation in any given individual with rare
83 disease.

84

85 **RESULTS**

86 We selected individuals with rare diseases who had undergone short-read exome sequencing
87 (srES; n=2) or srGS (n=94) in previous research studies yet had no pathogenic or likely
88 pathogenic variants (P/LP) nor variants of uncertain significance (VUS) identified (Bowling et al.
89 2017; East et al. 2021; Bowling et al. 2022). Most of our cohort consisted of children (89% were
90 <18 years of age at time of enrollment) with a neurodevelopmental disorder (NDD, 70%),
91 multiple congenital anomalies (MCA, 22%), or a suspected genetic myopathy (8%). Probands
92 consisted of 66% males (63/96); genetically inferred ancestries for probands revealed 72%
93 European (69/96), 21% African/African American (20/96), 3% Admixed American (3/96), 1%
94 Southeast Asian (1/96) and 3% unspecified ancestry admixture (Table 1). For these 96 cases,
95 we performed lrGS using Pacific Biosciences “HiFi” sequencing to a median depth of 27X
96 (Table 1, Supplemental Table S1). For a subset (10/96), we also performed lrGS on parents
97 (median parental HiFi depth of 22X, Supplemental Table S2). We also generated *de novo*
98 assemblies for each proband using hifiasm (Cheng et al. 2021), with parental srGS used for k-

99 mer-based binning and phasing when available. While the assemblies were not used for
100 structural variant calling in this study, they have proven useful in visualizing and evaluating
101 complex structural variation (Hiatt et al. 2021). The median N50 for assembled proband contigs
102 was 29.05 Mb (Table 1).

103 HiFi reads were aligned to hg38, and variant calling was performed using DeepVariant
104 (Poplin et al. 2018) and pbsv (<https://github.com/PacificBiosciences/pbsv>, see Methods). A
105 median of 4.4 million SNVs and 970,031 indels were called in each proband (Table 1). We
106 detected a median of 55,586 SVs of varying classes across the 96 probands using pbsv, of
107 which a median of 25,218 are greater than 50bp in length (Table 1, Supplemental Table S3).
108 These counts do not include any variant-quality filtration metrics and are expected to include an
109 undetermined fraction of false positive calls; only variant calls of potential disease interest were
110 evaluated for read support and quality (see Methods). Given that SVs are commonly arbitrarily
111 defined as >50bp in length, our counts generally are consistent with previous studies that have
112 detected a range of 22,000 to 33,000 SVs per genome (Pauper et al. 2021; Cohen et al. 2022;
113 Groza et al. 2024).

114

115 **Findings from IrGS**

116 IrGS SNVs/indels were annotated with features such as gene overlaps, coding consequences,
117 computational impact scores, and allele frequencies. They were then filtered and analyzed
118 using in-house software that is also used for srGS data (Hiatt et al. 2021). Rare SVs (see
119 Methods) were assessed by visualization of reads in IGV and prioritization and analysis using
120 SvAnna (Danis et al. 2022). All variants of interest were subject to curation using American
121 College of Medical Genetics and Genomics and Association for Molecular Pathology
122 (ACMG/AMP) and ClinGen criteria to identify potentially clinically relevant variation (Richards et
123 al. 2015; Riggs et al. 2020). We ultimately identified 19 potentially “clinically relevant” variants,
124 defined here as being pathogenic, likely pathogenic, or variants of uncertain significance

125 (P/LP/VUS), in 16 of the 96 cases (Table 2). Seven of these have a case-level classification of
126 Definitive Diagnostic or Likely Diagnostic, which we define as P/LP variants that likely fully
127 explain the reason for testing (Bowling et al. 2022). The remaining nine cases have Uncertain
128 case-level classifications, either due to the variants being VUSs or being P/LP variants in genes
129 whose associated phenotypes do not closely match the observed phenotype. Findings in seven
130 probands exemplify the unique benefits of IrGS and are highlighted below.

131

132 **IrGS-informed SVs**

133 IrGS uncovered a *de novo*, 4 Mb, copy-neutral, paracentric inversion on Chromosome 3
134 (NC_000003.12:g.110477273_114639202_inv) in Proband 1 (Figure 1). This inversion spans
135 about 35 protein-coding genes, and one breakpoint of this inversion lies within an intron of
136 *ZBTB20* (MIM: 606025), between exons 6 and 7 in the 5' UTR of NM_001348800.3. The event
137 breakpoints are visible in both long and short reads for this proband, but not in short reads for
138 either parent, suggesting it arose *de novo*. Phasing in long-read data indicates the inversion is
139 on the proband's paternal allele (Figure 1A). While evidence of the breakpoints are visible in the
140 proband's srGS data (Figure 1A, 1B) and the breakpoints are called by manta (v1.6.0, Chen et
141 al. 2016), the event only survived filtration and curation in IrGS (see discussion and
142 Supplemental Figure S1, which compares inversion and breakend calls between srGS and IrGS
143 for this proband). The variant is private to the proband and is predicted to disrupt *ZBTB20*.
144 While the breakpoint lies between two UTR exons, it moves and flips the transcript's promoter
145 and the first six exons away from the remaining exons (Figure 1C). This inversion is easily
146 visualized in alignment of the proband's paternal contig vs. the reference genome in this region
147 (Figure 1D). Loss-of-function (LOF) variation in *ZBTB20* is associated with Primrose Syndrome
148 (MIM: 259050) and 3q13.31 Microdeletion Syndrome (Juven et al. 2020). Proband 1's reported
149 features include moderate intellectual disability (ID), delayed speech and language
150 development, muscular hypotonia, strabismus, and hypoplastic corpus callosum. She is also

151 non-ambulatory. Several of these features overlap Primrose Syndrome. We have classified this
152 variant as likely pathogenic and the case-level designation is Likely Diagnostic (Table 2).

153 Proband 2 was originally enrolled for sequencing at the age of ~20 years, and presented
154 with spasticity, ataxia, and leukodystrophy. srGS was negative, but lrGS identified two structural
155 variants (SVs) identified in *trans* in *ALS2* (MIM: 606352). These include a maternally-inherited
156 4.65 kb deletion (Chr 2:201720435-201725085_del) that removes exons 21-23 of
157 NM_020919.4, and a paternally-inherited ~1.6 Mb deletion (Chr 2:200115181-201739349_del)
158 that spans several genes, including the 3' end of *ALS2* (deletion of exons 12-34 of
159 NM_020919.4, Figure 2A, Table 2, Supplemental Figures S2-S5). These deletions are easily
160 visualized in alignment of the proband's maternal (Figure 2B) and paternal (Figure 2C) contigs
161 vs. the reference genome in this region. While these variants were called in short-read data, the
162 smaller deletion was called as a heterozygous deletion in the mother and as a homozygous
163 deletion in the proband, obfuscating the nature of the variation and raising quality-control
164 concerns (e.g., nested and conflicting copy-number states and Mendelian inconsistency) that
165 prevented effective curation (Supplemental Figures S6, S7). Given the results of the lrGS, the
166 srGS variant calls logically resulted from the small maternal deletion intersecting with the larger,
167 overlapping paternal deletion. This case highlights the difficulties in identification and analysis of
168 overlapping SVs of unknown phase. Variation in *ALS2* is associated with several AR conditions
169 (Juvenile amyotrophic lateral sclerosis 2, MIM: 205100; Juvenile Primary lateral sclerosis, MIM:
170 606353; and infantile onset ascending Spastic paralysis, MIM: 607225), each of which have
171 features that overlap the proband's presentation. These variants are classified as P/LP and
172 given the degree of overlap with expected phenotypes the case-level designation is Definitive
173 Diagnostic.

174 Proband 3 is a male, enrolled in our research study at age ~50 years, and has a strong
175 X-linked family history of ID (Figure 3A). srES and srGS were both negative, although in srES,
176 two neighboring SNVs of uncertain significance were called, and manually curated, in the 3'

177 UTR of *HCFC1* (MIM: 300019). While srGS resulted in no calls in this region, visualization of
178 reads in IGV suggested an insertion of unknown length and consequence (Figure 3C). Variation
179 (mostly missense and proposed regulatory variation) in *HCFC1* has been associated with X-
180 linked recessive Methylmalonic aciduria and homocysteinemia, cblX type (MIM: 309541).
181 However, most affected individuals present with severely delayed psychomotor development,
182 seizures, and methylmalonic aciduria. Proband 3's family reported neither of the latter two
183 features. Given the uncertainties in the identity, structure, and consequence of the potential
184 variants seen in srES and srGS, and the lack of clear phenotypic relevance for the gene, this
185 region was curated but not originally considered to be a strong candidate for clinical relevance.
186 HiFi sequencing identified a 4,902 bp mobile element insertion (MEI) in the 3' UTR of *HCFC1*
187 (Chr X:153948602_ins4902), consisting of both SVA and L1 sequence (Figure 3B). This variant
188 was likely inherited from a heterozygous carrier mother, as indicated by srGS reads at the
189 breakpoints (Figure 3C). While this insertion does not affect protein-coding sequence, it is
190 predicted to increase the length of the 3' UTR from 1,791 nt to 6,693 nt. We subsequently
191 performed 3'-end RNA-seq on blood samples from both the proband and his father, generating
192 ~8 million reads for each sample (see Methods). *HCFC1* shows the greatest expression
193 decrease in the proband relative to his father, an ~8.8-fold reduction, among all genes with at
194 least 10 counts in each sample (Supplemental Figure S8). While these results are consistent
195 with the hypothesis that the insertion has a large effect on *HCFC1* expression and potential
196 activity, they are not definitive. Further, expression or segregation analyses in additional family
197 members could not be assessed. Given the uncertainty of the molecular consequence, the
198 differences between observed phenotypic features and those reported to associate with
199 *HCFC1*, and the lack of additional segregation data, we classified this as a VUS, with a case-
200 level designation of Uncertain.

201 Proband 4 has a complex *de novo* structural variant affecting 16p13.2
202 (NC_000016.9:g.(8742452_9220783)dup_ins[(8742452_8879961)_(9000190_9220783)], Table

203 2, Supplemental Figures S9-S18). Proband 4 first had trio srES and no variants were returned,
204 and this SV was not called. Trio lrGS identified the breakpoints, *de novo* status, phase, and
205 probable order, orientation, and copy number of segments in this SV. However, there remains
206 some uncertainty about the SV, as neither manual curation nor the proband's *de novo* assembly
207 could definitively resolve the full, exact structure of the locus. Two possible structures are shown
208 in Supplemental Figure S9. We believe that this is likely a limitation due to the length of the
209 reads, but it may also be influenced by low coverage in this proband (~17x), lack of parental
210 srGS data for hifiasm input, and a small number of informative variants near the breakpoints to
211 facilitate phasing.

212 Overlapping duplications in this region have been reported in gnomAD but are rare (Lek
213 et al. 2016). One individual in Decipher (Patient: 251349) has also been reported with a very
214 similar duplication of uncertain consequence (Deciphering Developmental Disorders 2015). This
215 region spans seven genes, three of which are associated with disease: *ABAT*, *PMM2*, and
216 *USP7*. The first (*ABAT*) is intersected by a duplication breakpoint in proband 4, but all other
217 breakpoints lie within intergenic regions (Supplemental Figure S9, S10). *USP7* is the only gene
218 associated with autosomal dominant disease (Hao-Fountain Syndrome, MIM:616863), but this
219 gene is expected to remain at a copy number of two in this proband whereas LOF is generally
220 the mechanism associated with disease (Hao et al. 2015). Some general features of Hao-
221 Fountain Syndrome overlap this proband, but it is not a strong phenotypic fit. The proband is
222 reported to have moderate ID, seizures, microcephaly, and facial dysmorphisms. Based on the
223 uncertain global structure and molecular consequence of the SV in this proband, in addition to
224 the clinical significance of variation in this region, we have classified this variant as a VUS, with
225 a case-level designation of Uncertain. Note that this variant is pending orthogonal validation.
226 Proband 4 was one of the two individuals who only previously had srES rather than srGS, and
227 this duplication is not easily visualized in srES data (Supplemental Figure S19).

228

229

230

231 IrGS-informed repeat expansions

232 Improved variant calling in repeat regions is also a benefit of IrGS (Nurk et al. 2022). In addition
233 to analysis of SNVs and SVs in our standard pipeline, we assessed variant calls in 66 tandem
234 repeat expansion (TRE) regions, including both known disease-associated and candidate
235 disease-associated loci (Supplemental Table S4). We intersected TRE regions of interest with
236 pbsv-called variants in each individual and compared these calls to known pathogenic
237 expansion sizes from the literature. We identified several large heterozygous insertions in
238 repeat regions that, while longer than the expected pathogenicity threshold, were predicted after
239 manual curation to be benign based on their sequence content (Nakamura et al. 2020,
240 Supplemental Figure S20). In two probands, we identified large heterozygous insertions in
241 *RFC1* (MIM: 102579), one of which is benign based on sequence content and one of which is
242 expected to be a pathogenic insertion. However, *RFC1*-associated disease (CANVAS, MIM:
243 614575) is caused by biallelic expansions, which we did not observe (Supplemental Figure
244 S20), suggesting the latter proband is merely a heterozygous carrier.

245 In Proband 5, we observed a *de novo* 18-bp alanine tract expansion in *PHOX2B* (MIM:
246 603851, NM_003924.4:c.741_758dup, p.(Ala255_Ala260dup), Table 2, Supplemental Figure
247 S21), associated with Central Hypoventilation Syndrome, with or without Hirschsprung Disease
248 (MIM: 209880). This disorder was clinically suspected, but variation in *PHOX2B* was missed by
249 initial clinical genetic testing and srGS, which was performed on a PCR+ srGS library. As
250 ExpansionHunter is intended to run on PCR-free srGS data, it was not run on this sample
251 (Dolzhenko et al. 2019). This variant is reported as pathogenic in the literature (Amiel et al.
252 2003), is a strong match to the proband's observed symptoms, and was confirmed
253 independently by additional clinical testing. We have classified this variant as pathogenic, with a
254 case-level classification of Definitive Diagnostic.

255 In proband 6, we identified a 270 bp insertion in *AFF3* (MIM:601464,
256 NM_001386135.1:c.-64-281_-64-280insGGC[90], Table 2, Supplemental Figure S22). While
257 missense variants in *AFF3* have been associated with KINSSHIP Syndrome (MIM:619297), an
258 expansion of a CGG-repeat in the promoter of this gene and subsequent hypermethylation of
259 the promoter has recently been reported to be associated with NDDs (Jadhav et al. 2023). Most
260 probands in our cohort (77/96) had both alleles matching hg38 in this region (Supplemental
261 Figure S23). Among the 19 probands harboring non-reference alleles, proband 6 had a 270 bp
262 insertion, and the remaining 18 probands had insertions ranging from 24-84 bp in length (8-28
263 triplet repeats). Similarly, when comparing to a larger database of in-house HiFi genomes
264 (n=266, “set 2”, see Methods), only 53/266 individuals harbor a non-reference allele; the longest
265 insertions outside of Proband 6 are 105 bp (35 repeats) and 93 bp (31 repeats), each in
266 different individuals, while the median non-reference insertion is 36 bp (12 repeats). Jadhav *et*
267 *al.* suggest that normal variation ranges to up to ~38 repeat units, with ≥ 61 repeats being a
268 likely pathogenic threshold. Thus, the 270 bp insertion (90 repeat units) in Proband 6 is well
269 above the normal range and pathogenicity threshold reported by Jadhav et al. and more than
270 twice as long as the second longest insertion in our sample of 266 individuals, which is in the
271 normal range reported by Jadhav et al. (Supplemental Table S4, Supplemental Figure S23).
272 Proband 6 was sequenced as a neonate and presented with intrauterine growth restriction
273 (IUGR) and hypoplastic left heart (HLH); Jadhav et al. reported a wide range of symptoms,
274 including intellectual disability/global developmental delays, seizures, behavioral disturbances,
275 and generalized hypotonia, but the specificity of phenotypic overlap with proband 6 is unclear
276 and neither supports nor refutes pathogenicity. While these results are consistent with
277 pathogenicity of the insertion in proband 6, there remains a need to further replicate and confirm
278 the spectrum of normal and pathogenic variation in *AFF3* repeat lengths. Given this uncertainty
279 and the uncertainty regarding the proband’s cognitive development, we have classified this TRE

280 as a VUS, and the case-level classification is Uncertain. Also note that this variant is pending
281 orthogonal validation.

282

283

284 **IrGS-informed SNV/indels**

285 A *de novo* *SHANK3* single-base deletion, predicted to lead to a frameshift
286 (NM_033517.1:c.3161delT, p.(Leu1054Argfs*10)), was identified by IrGS in proband 7. While
287 two reads in the srGS data support this deletion, the variant was not called by our srGS variant
288 calling pipeline (Supplemental Figure S24). LOF variation in *SHANK3* (MIM: 606230) is
289 associated with Phelan-McDermid syndrome (MIM: 606232). Features of this syndrome are
290 consistent with the proband's features, and we classified this variant as pathogenic (case-level
291 Definitive Diagnostic). Coverage of this region in short read data does not indicate a systematic
292 coverage deficiency, as mean coverage within 50 bp of this variant in the srGS data for the
293 cohort is 27.9x (n=94), while it is 15.1x for proband 7. Further, this gene is not present in the list
294 of medically relevant genes that tend to be poorly covered by srGS (Wagner et al. 2022). These
295 observations suggest the no-call may have resulted from a stochastic loss of alternative allele
296 reads in the srGS data in this proband. Nevertheless, IrGS has been shown to provide better
297 overall sensitivity and specificity to SNVs and indels in Genome-In-A-Bottle (GIAB) gold-
298 standard datasets compared to srGS (Logsdon et al. 2020; Hiatt et al. 2021) and thus detection
299 failures to variants such as this *SHANK3* event are more likely in srGS data in general.

300

301 **Reinterpretation of SNVs**

302 The remaining 9 cases had relevant variation identified following IrGS that were equally well
303 detected through reanalysis of existing srGS variant data (Table 2, Supplemental Case
304 Reports). In four cases (Probands 8-11, *HNRNPU*, *CSNK2B*, *GNB2*, *MCF2*) we identified
305 variation in genes that had additional published support for association of the gene or the

306 variant with disease since the time of the most recent analysis. In another four cases (Probands
307 12-15, *NOTCH3*, *AFF4*, *KCNT2*, *KIF21A*, *NRXN1*) we identified variation in established disease
308 genes that conflicted with the published data regarding molecular mechanisms or expected
309 mode of inheritance. Lastly, we identified a variant of interest in *SCN1A* that resulted from a
310 targeted analysis of candidate “poison exon” variants (Proband 16, Felker et al. 2023). We note
311 that in 7 of these 9 cases, variants were identified in an unaffected or mildly affected parent,
312 which was somewhat unexpected in these cases due to suspicion of high penetrance (also see
313 Supplemental Case Reports).

314

315 **Cohort Analysis of SVs**

316 Individual SV case analyses were performed on a rolling basis and variants were filtered and
317 prioritized using the HiFi variant data generated up to that point in time (see Methods,
318 Supplemental Figure S25). However, we also sought to characterize how filtering SVs by
319 frequency could reduce manual curation burden for future analyses by considering five allele
320 frequency resources (Supplemental Figure S26). First, we created a set of “cohort” SVs by
321 performing SV call merging across all 96 probands using Jasmine (Kirsche et al. 2023) and
322 generating allele counts from the merged set (set 1). We then used a second Jasmine merge
323 step to match cohort SVs with SV frequencies from: an in-house set of 266 HiFi genomes
324 including all cohort probands and parents, samples from other internal projects, and public HiFi
325 data (set 2, see Methods); gnomAD v4 SV frequencies from 63,046 short read samples (set 3,
326 Collins et al. 2020); Human Genome Structural Variant Consortium phase 2 (HGVC2)
327 assembly-based calls from 18 HiFi samples (set 4, Ebert et al. 2021,
328 [https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGVC2/release/v2.0/integrated_ca](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGVC2/release/v2.0/integrated_callset/variants_freeze4_sv_insdel_alt.vcf.gz)
329 [llset/variants_freeze4_sv_insdel_alt.vcf.gz](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGVC2/release/v2.0/integrated_callset/variants_freeze4_sv_insdel_alt.vcf.gz)); and a PacBio-provided set of pbsv calls from 103
330 HiFi samples from the Human Pangenome Reference Consortium (HPRC) and Genome in a
331 Bottle (GIAB) consortia (set 5, Liao et al. 2023; Zook et al. 2019;

332 https://github.com/PacificBiosciences/svpack/raw/main/resources/HPRC_GIAB.GRCh38.pbsv.v
333 [cf.gz](#)). We note that these data are partially redundant with one another (e.g., set 1 is a partial
334 subset of set 2 and some individual samples such as NA12878 are shared between sets 2-5).
335 We also note that gnomAD allele frequencies should be cautiously applied in the context of IrGS
336 SV annotation, as the srGS artifacts inherent in gnomAD variant calls may inappropriately mask
337 legitimate IrGS variants of interest, particularly given that gnomAD sample size is orders of
338 magnitude larger than the other datasets included in this analysis. However, our goal in this
339 analysis was to maximize filtering impact by including calls from as many datasets and
340 discovery methods as possible. Jasmine merging parameters were heuristically tuned until
341 further increased stringency demonstrated an inflection point in the unique variant count (see
342 Methods, Supplemental Figure S26).

343 When filtering SVs using these SV frequency datasets, we found that probands had a
344 median of 1,721 “rare” SVs, defined as having an allele frequency <1% in each of the public SV
345 databases (sets 3-5) and an allele count <4 in each internal cohort (sets 1 or 2). A small subset
346 of these, a median of 87 SVs per proband, overlapped a genomic position within 50 bp of a
347 RefSeq exon (Supplemental Table S3). We also filtered to identify “private” SVs as those absent
348 from sets 3-5, having a set 1 allele count of 1, and having a set 2 allele count of ≤ 2 (allowing
349 for uniparental inheritance). This filtering results in a median of 733 SVs per proband, only 40 of
350 which overlap a position 50 bp of a RefSeq exon (Supplemental Table S3). This analysis
351 includes all pbsv calls, and does not include quality, length, read-depth, or other filters that
352 might further reduce the numbers of variants needing curation.

353 DISCUSSION

354 There remains a substantial fraction of rare disease suspected to have genetic causes that is
355 refractory to genomic testing, a finding that has been repeatedly shown across many clinical
356 and research projects (e.g., Srivastava et al. 2019; Baxter et al. 2022). Several non-mutually
357 exclusive hypotheses exist to explain these observations. Environmental risk factors, such as
358 teratogenic exposure or infectious disease, may be relevant to some phenotypes. Multigenic
359 contributors are also likely to explain at least some cases. For example, a small but appreciable
360 fraction of “double diagnoses”, in which an affected individual is observed to harbor two distinct
361 diseases resulting from highly penetrant variation in two distinct genes, has been observed in
362 clinical genomic studies (e.g. Posey et al. 2017). Notably, such discoveries are typically only
363 made when the variation in both genes is independently amenable to a pathogenicity
364 determination (i.e., would be P/LP regardless of P/LP variation in the other gene), and it is
365 possible if not probable that at least some conditions result from combinations of variants in
366 different genes that are not pathogenic in isolation (Papadimitriou et al. 2019). At the further end
367 of this spectrum is a polygenic accumulation of many risk-factor alleles, which is known to be
368 relevant to many common, complex diseases and which may contribute to some rare
369 conditions, as has been suggested for at least a subset of NDDs (Niemi et al. 2018).

370 We find it likely that a substantial fraction of unexplained rare disease arises from highly
371 penetrant variation that we have not yet been able to precisely identify or confidently interpret.
372 The results from this study are consistent with that hypothesis, with over 15% of probands with
373 previously negative testing being now found to harbor a relevant or potentially relevant genetic
374 variant. Further, these observations are consistent with the general picture of rare disease
375 testing in recent years. With the advent of exome capture and sequencing ~15 years ago (Ng et
376 al. 2010, 2009) and subsequent improvements in cost and efficiency, genome-wide detection of
377 highly penetrant variation has greatly accelerated in recent years (Bamshad et al. 2019; Baxter
378 et al. 2022; Hamosh et al. 2022; Boycott et al. 2022). Long-read genome sequencing represents

379 the next phase of that acceleration by facilitating a substantial increase in variant
380 comprehensiveness and accuracy.

381 We note that the benefits of IrGS for detection of highly penetrant variation derive not
382 just from improved variant sensitivity, but specificity gains as well. Indeed, effective analysis
383 requires not just the ability to detect a variant but to do so accurately and with sufficient
384 specificity for routine and scalable use. For example, effectively all the variation we describe in
385 this study as being newly found within IrGS is, in fact, visible in srGS data, at least at the
386 breakpoint levels. However, being retrospectively visible, once the location and structure of a
387 variant is known to exist, is a much lower bar than the ability to prospectively detect, define,
388 filter, and curate such variation. For example, we describe here a 4.9 kb insertion of SVA and L1
389 sequence into the 3' UTR of *HCFC1*. In retrospective analysis, the Mobile-Element Locator Tool
390 (MELT, Gardner et al. 2017) detected a 1.2 kb SVA mobile element at this location in srGS.
391 However, this call is incorrect with respect to size and sequence composition and has only one
392 read flagged as supporting the right breakpoint. As another example, the breakpoints of the
393 inversion observed to overlap *ZBTB20* were correctly called as breakends by manta (Chen et al.
394 2016). However, each of these were retrospective analyses since MELT calls and breakend
395 calls produced by manta are too numerous for manual curation and thus not used as part of our
396 standard variant calling pipelines (Hiatt et al. 2021, Supplemental Figure S1). Related to this,
397 srGS CNV/SV callers are known to produce many false calls and require strict filtering to reduce
398 calls to a reasonable number for curation. For example, among the 94 samples that had srGS
399 in this study, the raw output from the four different CNV/SV callers that we currently run (see
400 Methods) produces a mean of 9,944 total calls per sample, of which only ~50 are subject to
401 manual curation.

402 Our results are also consistent with other studies of the benefits of IrGS for discovering
403 genetic contributors to disease. In 2021, for example, we showed in a small pilot project that two

404 of six previously srGS-negative probands harbored clinically relevant variation uniquely
405 interpretable by lrGS (Hiatt et al. 2021). Since that time, several other studies have also used
406 lrGS for molecular diagnosis of rare disease. For example, Cohen and colleagues showed
407 increased yield of GS (both lrGS and srGS) in exome-negative cases (Cohen et al. 2022).
408 Unique discoveries for lrGS included detection of novel repeat expansions of *STARD7* and a
409 compound heterozygous SNV/deletion that was easier to detect in lrGS. However, the increase
410 in diagnostic yield (~13%) was mainly from variants interpretable in either lrGS or srGS. Other
411 studies have also shown the diagnostic value of lrGS, especially in small, well-phenotyped
412 cohorts or families (Sakamoto et al. 2024; Kilich et al. 2024; Audet et al. 2023; Del Gobbo et al.
413 2023; Fukuda et al. 2023; Miller et al. 2021). Our results are similar to these studies at a high-
414 level. However, some important differences are worth noting. One difference is that we have
415 shown the value of lrGS in singletons. While srGS data was available for most of the probands'
416 parents (79/96, ~82%), only 10/96 probands had parent lrGS data. Our lrGS variant filtering
417 strategies were sufficient to allow curation of proband variants without parental lrGS data.
418 However, once flagged for interest, inheritance data could often be assessed by looking at
419 parental srGS reads. Based on this experience, prioritization of proband lrGS with targeted
420 validation in parents is an effective way to increase diagnostic yield with lrGS while reducing
421 lrGS sequencing needs.

422 We have also provided a direct, systematic comparison of lrGS to contemporaneously
423 analyzed srGS in previously negative cases. Our results thus allow for the separation and
424 description of clinically relevant variants that are “new” by virtue of benefitting from the unique
425 advantages of lrGS versus those that are “new” by virtue of reanalysis, and which could be
426 detected and accurately interpreted via lrGS or srGS-reanalysis. In that context, we note that
427 the benefits of reanalysis of older srGS data remain considerable. In the results described here,
428 the “new” discoveries in 9 of 16 cases in lrGS were called correctly and interpretable within
429 srGS data. These observations reflect the rapid pace of gene discovery in rare disease (Baxter

430 et al. 2022; Boycott et al. 2022; Hamosh et al. 2022), and are consistent with other studies. For
431 example, we previously showed that the probability of a negative srGS dataset harboring a
432 clinically relevant variant increased from 1% within one year of a previous analysis to ~22% if
433 more than three years have passed since a previous analysis (Hiatt et al. 2018). Several other
434 studies found similar results, with many reanalysis findings being due to recent publications of
435 new gene-disease associations (Liu et al. 2019; Schobers et al. 2022; Hartley et al. 2023).

436 One possible optimal path to maximizing overall yield in previously srGS-negative
437 individuals is to include srGS-reanalysis prior to lrGS. However, this reflects a cost/benefit ratio
438 that depends on the cost of the analysis step in relation to the costs of sequencing. While lrGS
439 costs currently remain higher than srGS, as lrGS costs decrease there may reach a point where
440 the cost of variant analysis alone (which is non-trivial and requires both software and compute
441 resources as well as manual curation) is substantial relative to sequencing costs per se and
442 which might favor a process of simply performing lrGS. Further, evaluation of a study design or
443 testing process needs to consider sequencing as only one of several parts of the overall
444 process (e.g., participant recruitment and phenotyping). Again, especially as lrGS costs decline,
445 these factors are likely to increasingly favor use of lrGS.

446 One caveat to this study is that the results were generated over a period of time with
447 considerable change in lrGS protocols. For example, most probands in this study were
448 sequenced on Sequel IIE machines (n=64) before the Revio (n=32) became available. Given
449 the costs of data production, the 64 Sequel IIE samples were covered at lower-depth (median
450 coverage 24.65X) than the 32 Revio-sequenced samples (median coverage 30.09X,
451 Supplemental Figure S27), which may have reduced our sensitivity to variants in the earlier
452 samples. Further, methylation data were not available in the early period of this study, although
453 methylation calls are now being generated and are available for the most recent 44 probands,
454 which may also impact diagnostic yield. For example, evaluation of the *AFF3* expansion in
455 Proband 6 (Supplemental Figure S22, S23) would benefit from an assessment of methylation

456 levels at this locus, as hypermethylation, in addition to insertion length, is likely associated with
457 disease risk (Jadhav et al. 2023). Additionally, some probands may harbor methylation
458 alterations that are clinically relevant even in the absence of a pathogenic genetic variant (Aref-
459 Eshghi et al. 2021).

460 In addition to changes in sequencing, informatic changes have also been considerable
461 over the course of the data generation for this study. While we present and describe a uniform
462 set of variant-calls, assemblies, and annotations (see Methods), analysis of individual samples
463 took place simultaneously with the optimization of variant-calling and annotation pipelines. One
464 particularly relevant change is variant-frequency annotations. While the key strength of IrGS is
465 its ability to see variants that are invisible or poorly detected in srGS, the ability to discriminate
466 genuine highly penetrant variation from the background of benign alleles depends on the ability
467 to annotate and remove alleles that are common in the general population. As the main allele
468 frequency resources are built from srGS data (e.g., Lek et al. 2016), we have limited ability to
469 filter away likely benign alleles among the variants uniquely detected by IrGS. This ability,
470 however, improved as more IrGS datasets were produced over the course of this study (Ebert et
471 al. 2021; Nurk et al. 2022). Projects like the COLORs consortium (<https://colorsdb.org/>, currently
472 with IrGS data from almost 1400 individuals) are likely to improve frequency annotation in the
473 future and continue to improve variant curation efficacy. Accumulating IrGS data from as many
474 samples and studies as possible is critical for the long-term maximization of IrGS benefits.
475 Further, assembly-based SV calling is an exciting opportunity that will likely improve SV
476 detection and prioritization (Groza et al. 2024).

477 In sum, rare disease genetics continues to be a rapidly advancing field. With data-
478 sharing (Sobreira et al. 2015; Philippakis et al. 2015; Muenzen et al. 2022) and technology
479 improvement (Wenger et al. 2019), the overall diagnostic yield for individuals with rare disease
480 is increasing at a considerable pace each year. In that context, IrGS has clear benefits over
481 srGS, providing substantial gains to variant specificity and sensitivity, especially for complex and

482 repeat-associated variants (Schuy et al. 2022). We anticipate that the degree of improvement
483 will widen over time, as sequencing and analysis pipelines mature and as IrGS datasets grow.

484 **METHODS**

485 *Short-read sequencing and variant calling*

486 Probands, their parents, and, when appropriate, affected siblings, were enrolled in one of four
487 research studies aimed at identifying genetic causes of rare disease (Bowling et al. 2017; East
488 et al. 2021; Bowling et al. 2022 and Pediatric Genomics (PGEN), unpublished). These studies
489 were monitored by Western IRB and UAB IRB (WIRB 0071, UAB IRB protocols 170303004,
490 300000328, and 130201001). Short read exome (srES) or short-read genome sequencing
491 (srGS) was performed as described (Bowling et al. 2022; East et al. 2021; Hiatt et al. 2021;
492 Bowling et al. 2017). Briefly, whole blood genomic DNA was isolated using the QIAasympyphony
493 (Qiagen), and sequencing libraries were constructed by the HudsonAlpha Genomic Services
494 Lab or the Clinical Services Laboratory, LLC, using a standard protocol that generally included
495 PCR amplification (86/96). Genomes were sequenced at an approximate mean depth of 30X,
496 with at least 80% of base positions reaching 20X coverage. Exomes were sequenced to a mean
497 depth of 71X. For short read reanalysis, srES and srGS reads were aligned to hg38 and small
498 variants and CNVs were called with a uniform pipeline. SNVs/indels and CNVs were curated
499 using an in-house software tool, as previously described (Hiatt et al. 2021). Expected sample
500 relatedness was confirmed with Somalier (v. 0.2.10, (Pedersen et al. 2020) and predicted major
501 genetic ancestries were calculated with peddy (v. 0.4.1, Pedersen and Quinlan 2017). srGS
502 data for participants who consented to controlled-access sharing in NIH-funded studies are
503 available via dbGaP/AnVIL (CSER1: [https://www.ncbi.nlm.nih.gov/projects/gap/cgi-](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001089.v4.p1)
504 [bin/study.cgi?study_id=phs001089.v4.p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001089.v4.p1), dbGaP accession phs001089;
505 SouthSeq: <https://anvilproject.org/data/studies/phs002307>, dbGaP accession phs002307).
506 Supplemental Table S6 provides the applicable srGS sample ID and dbGaP accession of all

507 individuals who consented to sharing, and the dbGaP/AnVIL metadata may be used to locate
508 relevant parent and familial srGS data.

509

510 *Long-Read sequencing, variant calling, analysis and de novo assemblies*

511 Long-read sequencing was performed using HiFi (CCS) mode on either a PacBio Sequel II or
512 Revio instrument (Pacific Biosciences of California, Inc.). Libraries were constructed using a
513 SMRTbell Template Prep Kit (V1.0, 2.0 or 3.0) and tightly sized on a SageELF or BluePippin
514 instrument (Sage Science, Beverly, MA, USA). Sequencing was performed using a 2 hour pre-
515 extension with either 24 or 30 hour movie times. The resulting raw data was processed using
516 either the CCS3.4 or CCS4 algorithm, as the latter was released during the course of the study.
517 Comparison of the number of high-quality indel events in a read versus the number of passes
518 confirmed that these algorithms produced comparable results. Probands were sequenced on 2-
519 3 Sequel II or one Revio SMRT cell. Top off sequencing was performed if the sequencing did
520 not meet the desired coverage (>20X). This resulted in an average estimated HiFi read depth of
521 26.1X (range 16.7-41) of raw, unaligned sequence data for probands. For 10 families, parents
522 were also sequenced on 2-3 Sequel II SMRT cells, with an average estimated depth of 21.6x
523 (range 14-28) of raw, unaligned sequence data for parents. Aligned sequencing metrics are
524 shown in Supplemental Tables S1 and S2. HiFi reads were aligned to the hg38 no-alt analysis
525 set

526 (https://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/001/405/GCA_000001405.15_GRCh38/seqs_for_alignment_pipelines.ucsc_ids/GCA_000001405.15_GRCh38_no_alt_analysis_set.fna.gz)

528 using pbmm2 (v. 1.10.0, <https://github.com/PacificBiosciences/pbmm2>). SNVs and small indels,
529 were called with DeepVariant (v. 1.5.0) and used to haplotag the aligned reads with whatshap
530 (v. 1.7)(Martin et al.). Structural variants were called using pbsv v2.9.0
531 (<https://github.com/PacificBiosciences/pbsv>), which discovers SVs based on read alignment and
532 split reads. Read-depth-based SV discovery was not performed. For the 10 cases with parent

533 lrGS, candidate de novo SVs required a proband genotype of 0/1 and parent genotypes of 0/0,
 534 with ≥ 6 alternate reads in the proband along with 0 alternate reads and ≥ 5 reference reads in
 535 the parents.

536

537 De novo assemblies were generated for all probands using hifiasm (v. 0.15.2 (r334 or v.0.16.1-
 538 r375, Cheng et al. 2021). Hifiasm was used to create two assemblies. First, the default
 539 parameters were used, followed by two rounds of Racon (v1.4.10) polishing of contigs (Vaser et
 540 al. 2017). For cases with parent GS data, trio-binned assemblies were built using k-mers
 541 (srGS). The k-mers were generated using yak (v0.1) using the suggested parameters for
 542 running a hifiasm trio assembly (k-mer size=31 and Bloom filter size of 2^{37}). Maternal and
 543 paternal contigs went through two rounds of Racon (v1.4.10) polishing. Individual parent
 544 assemblies were also built with hifiasm (v0.15.2) using default parameters. The resulting contigs
 545 went through two rounds of Racon (v1.4.10) polishing.

546

547 Coordinates of breakpoints were defined by a combination of assembly-assembly alignments
 548 using minimap2 (Li 2018) followed by use of `bedtools bamToBed` (Quinlan and Hall 2010),
 549 visual inspection of CCS read alignments, and BLAT. Dot plots illustrating sequence differences
 550 were created using Gepard (Krumtsiek et al. 2007).

551

552 *Structural Variant Merging and Filtering for Case Analysis*

553 Overall structural variant counts in Supplemental Table S1 were generated from the proband
 554 pbsv VCFs using BCFtools (v1.15.1)(Danecek et al. 2021) and awk (e.g., `bcftools filter`
 555 `-i 'SVTYPE=="DEL" <input.vcf> | {a[i++]=$1; sum+=$1} END{asort(a);`
 556 `min=a[1]; max=a[i]; if(i%2==1) median=a[int(i/2)+1]; else`
 557 `median=(a[i/2]+a[i/2+1])/2; mean=sum/i; print mean, median, min, max;}`).

558 Counts of SVs > 50 bp in length in Table 1 were created by addition of the BCFtools filters -e
559 'SVLEN>-50 && SVLEN<50'. No quality filtration or read support filtration was used.

560 An internal allele frequency catalog was constructed and periodically updated using all
561 available internally sequenced HiFi pbsv variant calls and pbsv calls generated from public HiFi
562 sequencing data (n=266 at the end of this analysis, “set 2”). Samples included: a majority of
563 participants from this cohort and our previous NDD pilot cohort (Hiatt et al. 2021, 120/266);
564 samples from HudsonAlpha non-NDD projects (35/266); HudsonAlpha-sequenced data from
565 HG001, HG003, HG004, HG006, and HG007 (5/266); HG00514, HG00731, HG00732,
566 NA19240 from HGSVC2 (Ebert et al. 2021) (4/266); CHM13 (1/266) (Nurk et al. 2022),
567 (https://www.ncbi.nlm.nih.gov/sra/?term=SRX789768*+CHM13); and the HPRC Year 1 and
568 HPRC_PLUS releases (101/266)(Liao et al. 2023). pbsv calls were merged naively using
569 `bcftools merge` (v. 1.15.1), i.e., merging only variants that were identical in terms of location,
570 reference sequence, and alternate sequence. The internal pbsv allele frequency catalog
571 included affected probands as well as parent/child trios. Supplemental Figures 25 and 26 detail
572 the pipeline for the creation and use of this catalog.

573

574 *Structural Variant Annotation and Curation*

575 For individual case structural variant analysis, a frequency-filtered subset of the proband’s pbsv
576 calls was generated using `bcftools annotate` and `bcftools filter`, requiring that calls
577 overlap a genomic position within 50bp of a RefSeq exon and have an allele count of < 4 in the
578 HudsonAlpha internal allele frequency catalog (set 2, described above). Exon regions were
579 defined as RefSeq exons +/- 50bp (calculated from `bedtools slop`) and were used to restrict
580 output from BCFtools to only calls overlapping those regions (the `-R` command line option). For
581 SVs that span multiple bases of the reference genome (duplications, deletions, and inversions),
582 this filter requires only at least one bp of overlap between the SV span and the target regions.

583 For SVs that exist at a single base of the reference genome (insertions and breakends), the
584 position must be within the target regions. Each call was visualized using a custom pipeline to
585 automatically generate IGV screenshots (see Supplemental Code for implementation,
586 <https://github.com/HudsonAlpha/igv-grapher>) (Robinson et al. 2011). Additionally, these filtered
587 pbsv variants were annotated, prioritized, and visualized with SvAnna (v1.0.4, annotations
588 v.2204 or v.2304, Danis et al. 2022) based on manually curated HPO terms for each case.
589 Supplemental Figure S25 details the SV case analysis pipeline from variant calling through
590 curation.

591
592 *Variant interpretation and orthogonal confirmation*
593 Variant interpretation was performed using ACMG and ClinGen (Richards et al. 2015; Riggs et
594 al. 2020). Variants of interest were either clinically confirmed by the HudsonAlpha Clinical
595 Services Lab, confirmed within a research lab, and/or were supported by short-read data,
596 except where noted in the text (Table 2).

597
598 *Structural Variant Merging and Filtering for Cohort-level Analysis*
599 In order to provide a cohort-level descriptive analysis of SVs and assess the filtering efficacy of
600 combining multiple allele frequency resources, a more robustly merged cohort catalog (n=96)
601 was constructed using Jasmine. Jasmine allows for merging of similar structural variants that
602 may have non-identical representation in terms of genomic position or variant sequence via
603 spanning a structural variant proximity graph. First, heuristic trials were conducted to determine
604 a set of stringent Jasmine merge options that would still effectively reduce the resulting count of
605 distinct SVs. The centroid merging strategy was chosen instead of the Jasmine default to lessen
606 the risk of over-merging two dissimilar SVs connected through a third SV of intermediate
607 position and length. In a single trial of Jasmine merging strategies, we observed that using the
608 centroid merging strategy increased the unique SV count by ~8,000 while using the most-

609 stringent clique merging increased the unique SV count by ~99,000. Supplemental Figure S28
610 shows the effects of the minimum overlap percentage threshold on the unique deletion count
611 and the minimum sequence identity threshold on unique insertion count. Jasmine was ultimately
612 run with options `--centroid_merging --min_overlap=0.65 --`
613 `min_sequence_id=0.75 --output_genotypes` to create the merged catalog. This cohort
614 catalog was then annotated with five sets of structural variant frequency annotation. Cohort
615 allele frequencies were generated from the merged set with `bcftools +fill-tags` (set 1,
616 `n=96`). A second Jasmine merge was used to combine the cohort catalog allele count with
617 additional allele frequency resources: the HudsonAlpha internal pbsv callset (described above)
618 (set 2, `n=266`); gnomAD structural variants (Collins et al. 2020, v4.0,
619 <https://gnomad.broadinstitute.org/news/2023-11-v4-structural-variants>) (set 3, `n=63,046`);
620 HGVC2 structural variants (Ebert et al. 2021)
621 [https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGVC2/release/v2.0/integrated_ca](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGVC2/release/v2.0/integrated_callset/variants_freeze4_sv_insel_alt.vcf.gz)
622 [llset/variants_freeze4_sv_insel_alt.vcf.gz](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGVC2/release/v2.0/integrated_callset/variants_freeze4_sv_insel_alt.vcf.gz)) (set 4, `n=18`); pbsv calls for individuals from the
623 Human Pangenome Reference Consortium and Genome in a Bottle (from PacBio, Liao et al.
624 2023; Zook et al. 2019,
625 [https://github.com/PacificBiosciences/svpack/raw/main/resources/HPRC_GIAB.GRCh38.pbsv.v](https://github.com/PacificBiosciences/svpack/raw/main/resources/HPRC_GIAB.GRCh38.pbsv.vcf.gz)
626 [cf.gz](https://github.com/PacificBiosciences/svpack/raw/main/resources/HPRC_GIAB.GRCh38.pbsv.vcf.gz)) (set 5, `n=103`). The second Jasmine merge was run with the prior settings but without the
627 `--output_genotypes` option and records without genotypes were discarded. A custom
628 Python script was used to transfer allele frequency annotations from the annotation source
629 VCFs to the merged VCF based on unique variant identifiers (which were created for each allele
630 frequency resource as needed) and the IDLIST record from Jasmine (see Supplemental Code,
631 https://github.com/HudsonAlpha/jasmine_sv_annotations). The annotated catalog was split into
632 individual VCFs with `BCFtools` (`bcftools view -I -s <id>`) which were filtered with
633 `BCFtools` to generate the counts in Supplemental Table S2. The “rare” filter described in

634 Supplemental Table S3 was defined as excluding variants with any of the following: within-
635 cohort allele count > 3, in-house allele count > 3, gnomAD maximum population allele frequency
636 (POPMAX_AF) >= 1%, HPRC_GIAB allele frequency >=1%, or HGVC2 allele frequency >=
637 1%. The “proband-exclusive” filter described in Supplemental Table S3 was defined as
638 excluding variants with any of the following: within-cohort allele count > 1, in-house allele count
639 > 2, gnomAD allele count > 0, HPRC_GIAB allele count > 0, or HGVC2 allele count > 0. The
640 in-house allele count cutoff was set at 2 in order to account for the fact that some parental
641 samples were included in the in-house frequency database. Exon regions were defined as
642 RefSeq exons +/- 50bp (described above). Variant filtering and counting was performed with
643 BCFtools and command-line tools, e.g., `bcftools view -R`
644 `refseq_exons_plus50bp.bed.gz -e 'AC>3 | inhouse_pbsv_AC >3 |`
645 `gnomad4_POPMAX_AF >= 0.01 | hgsvc2_AF >= 0.01 | hprc_giab_pbsv_AF >=`
646 `0.01' -H <input_file> | wc -l` for the rare filter and `bcftools view -R`
647 `refseq_exons_plus50bp.bed.gz -e 'AC>1 | inhouse_pbsv_AC>2 |`
648 `gnomad4_AC>0 | hgsvc2_AC>0 | hprc_giab_pbsv_AC>0' -H <input_file> | wc`
649 `-l` for the exclusive filter. No filtering based on SV length was used; counts include insertions
650 and deletions less than 50bp in length that were called by pbsv. Supplemental Figure 26 details
651 the cohort SV merging and analysis pipeline.

652

653 *Sequencing Metrics*

654 Sequencing metrics were generated from the aligned BAMs using the Sentieon implementations
655 of Picard sequencing metrics (Kendig et al. 2019). Sentieon algorithms QualityYield,
656 AlignmentStat were run with default settings. Sentieon `WgsMetricsAlgo` was run with settings
657 `--include_unpaired true --min_map_qual 0 --min_base_qual 0 --`
658 `coverage_cap 5000`. Further coverage metrics were generated with `cramino` using the `--`

659 `phased option` (De Coster and Rademakers 2023). Supplemental Table S5 maps each tool
660 or command and output field name to the corresponding entry in Supplemental Tables S1 and
661 S2. SNV and indel counts and ratios were calculated with `rtg vcfstats`
662 (<https://github.com/RealTimeGenomics/rtg-tools>). Summary statistics and graphs were
663 calculated with R (v4.3.1), RStudio (v2023.9.1 build 494), and `ggplot2` (v3.4.3)(Wickham 2011;
664 Team 2020; R Core Team 2023). Coverage across regions of interest, such as *SHANK3*, was
665 calculated using `samtools bedcov` with default parameters.

666

667 *Repeat region detection and analysis*

668 We curated a BED file of disease-associated low-complexity repeat regions in 66 genes from
669 previous studies (Hiatt et al. 2021; Cohen et al. 2022 and references therein). Variant calls from
670 `pbsv` were extracted from these regions +/- 30bp (Supplemental Table S4). Reads were also
671 visualized using the Integrated Genomics Viewer (IGV). Coverage across low complexity repeat
672 regions was calculated using `samtools bedcov` with default parameters. A coverage of at
673 least 8x across the low-complexity region was required for inclusion in Supplemental Table S4
674 and Supplemental Figure S23. We also used TRGT (Dolzhenko et al. 2024) and companion tool
675 TRVZ for visualization of a subset of calls including those for display in Figure S3. TRGT was
676 fed an hg38 reference genome FASTA, BED catalog of tandem repeats, and a sample's BAM to
677 generate a VCF containing genotypes for each tandem repeat from the provided catalog in the
678 given sample and a BAM containing only reads that span the repeat sequences. Output VCFs
679 and BAMs were sorted, indexed, and fed into TRVZ with the same hg38 reference FASTA and
680 BED catalog of tandem repeats to generate pileup plots for any desired variants.

681

682 *3' mRNA-seq*

683 Total RNA was isolated from blood samples in PAXgene RNA tubes (PreAnalytiX #762165)
684 according to the manufacturer's instructions and stored short-term at -20°C. RNA was isolated
685 using the PAX gene Blood RNA Kit (Qiagen #762164) according to the manufacturer's
686 instructions. Isolated RNA was quantified by the Qubit RNA HS Assay Kit (Thermo Q32855).
687 425 ng of RNA was used as input for the QuantSeq 3' mRNA-Seq Library Prep Kit FWD for
688 Illumina and UMI Second Strand Synthesis Module for QuantSeq FWD (Illumina, Read 1) from
689 Lexogen (015.96 and 081.96, respectively). Libraries were quantified using the Qubit DNA HS
690 Assay Kit (Thermo Q32854) and visualized with the Bioanalyzer High Sensitivity DNA Analysis
691 kit (Agilent 5067–4626) and 2100 Bioanalyzer Instrument (Agilent). Sequencing was carried out
692 using Illumina NextSeq 75 bp single-end. UMIs were first extracted from the reads with UMI-
693 tools extract with regex. Reads were then trimmed with bbduk
694 (<https://sourceforge.net/projects/bbmap/>) and aligned to hg38-GENCODEv42 using STAR
695 (Dobin et al. 2013) with the Lexogen recommended parameters for QuantSeq. Bams were
696 deduplicated by UMI and mapping coordinates using UMI-tools dedup (Smith et al. 2017). Count
697 tables were generated using htseq-count with the intersection-nonempty method (Anders et al.
698 2015).

699

700 DATA ACCESS

701 82 out of 96 families (85%) consented to controlled-access genomic data sharing, including all
702 probands described in detail in the manuscript and listed in Table 2. Consented samples include
703 100 lrGS genomes (82 probands and 18 parents, comprising 9 full trios and 73 singletons). For
704 participants who consented to controlled-access sharing, the lrGS data generated in this study
705 will be available through AnVIL (<https://anvilproject.org/data/studies>) under accession number
706 phs003537. Researchers may apply for access via dbGaP (<https://www.ncbi.nlm.nih.gov/gap/>).
707 Supplemental Table S6 describes the samples and de-identified file identifiers and, when

708 applicable, their corresponding short-read dbGaP identifiers and other metadata. Custom scripts
709 for automated IGV plot generation and Jasmine structural variant merging and annotation are
710 available in Supplemental Code and on GitHub at <https://github.com/HudsonAlpha/igv-grapher>
711 <https://github.com/HudsonAlpha/igv-grapher>.

712

713 **COMPETING INTEREST STATEMENT**

714 The authors declare no competing interests.

715

716

717 **ACKNOWLEDGMENTS**

718 We thank the families who have contributed to our studies and our collaborating physicians and
719 clinical staff for recruitment and enrollment for these research studies. Some reagents were
720 provided by PacBio as part of an early-access testing program. The CSER1 project was
721 supported by a grant from the US National Human Genome Research Institute (NHGRI;
722 UM1HG007301). The SouthSeq project (U01HG007301) was supported by the Clinical
723 Sequencing Evidence-Generating Research (CSER2) consortium, which was funded by the
724 National Human Genome Research Institute with co-funding from the National Institute on
725 Minority Health and Health Disparities and the National Cancer Institute. The Alabama Genomic
726 Health Initiative is an Alabama-State earmarked project (F170303004) through the University of
727 Alabama in Birmingham. The PGEN cohort was funded by the Alabama Pediatric Genomics
728 Initiative. IrGS of some samples was supported by a Research Grant from the Muscular
729 Dystrophy Association (MDA 963255). SMH, JMJL, LHH, DRL, ZTB, CRF, MLT, LBB, MW, IRN,
730 and JJ generated, analyzed, and/or interpreted data. JMJL and CRF coordinated data sharing
731 to AnVIL. WVK, EMB, MAL, ACEH, and BRK recruited and enrolled patients and provided
732 clinical assessment. SMH, JMJL, JS, JG, and GMC designed the study and oversaw

733 interpretation of data. SMH, JMJL and GMC drafted and revised the manuscript. All authors
734 read and approved the final version of the manuscript.

735

736

737 **REFERENCES**

- 738 Amiel J, Laudier B, Attié-Bitach T, Trang H, De Pontual L, Gener B, Trochet D, Etchevers H,
739 Ray P, Simonneau M, et al. 2003. Polyalanine expansion and frameshift mutations of the
740 paired-like homeobox gene PHOX2B in congenital central hypoventilation syndrome. *Nat*
741 *Genet* **33**: 459–461. <https://pubmed.ncbi.nlm.nih.gov/12640453/> (Accessed February 5,
742 2024).
743
- 744 Anders S, Pyl PT, Huber W. 2015. HTSeq—a Python framework to work with high-throughput
745 sequencing data. *Bioinformatics* **31**: 166–169.
746 <https://dx.doi.org/10.1093/bioinformatics/btu638> (Accessed August 1, 2024).
747
- 748 Aref-Eshghi E, Kerkhof J, Pedro VP, France G DI, Barat-Houari M, Ruiz-Pallares N, Andrau JC,
749 Lacombe D, Van-Gils J, Fergelot P, et al. 2021. Evaluation of DNA Methylation
750 Episignatures for Diagnosis and Phenotype Correlations in 42 Mendelian
751 Neurodevelopmental Disorders. *Am J Hum Genet* **108**: 1161–1163.
752 <https://pubmed.ncbi.nlm.nih.gov/34087165/> (Accessed February 27, 2024).
753
- 754 Audet S, Triassi V, Gelinat M, Legault-Cadieux N, Ferraro V, Duquette A, Tetreault M. 2023.
755 Integration of multi-omics technologies for molecular diagnosis in ataxia patients. *Front*
756 *Genet* **14**: 1304711. <http://www.ncbi.nlm.nih.gov/pubmed/38239855> (Accessed February
757 13, 2024).
758
- 759 Bamshad MJ, Nickerson DA, Chong JX. 2019. Mendelian Gene Discovery: Fast and Furious
760 with No End in Sight. *Am J Hum Genet* **105**: 448–455.
761 <https://pubmed.ncbi.nlm.nih.gov/31491408/> (Accessed July 19, 2022).
762
- 763 Baxter SM, Posey JE, Lake NJ, Sobreira N, Chong JX, Buyske S, Blue EE, Chadwick LH,
764 Coban-Akdemir ZH, Doheny KF, et al. 2022. Centers for Mendelian Genomics: A decade
765 of facilitating gene discovery. *Genet Med* **24**: 784–797.
766 <https://pubmed.ncbi.nlm.nih.gov/35148959/> (Accessed April 16, 2023).
767
- 768 Bowling KM, Thompson ML, Amaral MD, Finnila CR, Hiatt SM, Engel KL, Cochran JN, Brothers
769 KB, East KM, Gray DE, et al. 2017. Genomic diagnosis for children with intellectual
770 disability and/or developmental delay. *Genome Med* **9**: 43.
771 <http://www.ncbi.nlm.nih.gov/pubmed/28554332>.
772
- 773 Bowling KM, Thompson ML, Finnila CR, Hiatt SM, Latner DR, Amaral MD, Lawlor JM, East
774 KM, Cochran ME, Greve V, et al. 2022. Genome sequencing as a first-line diagnostic test
775 for hospitalized infants. *Genet Med* **24**: 851–861.
776 <https://pubmed.ncbi.nlm.nih.gov/34930662/> (Accessed May 29, 2023).
777

- 778 Boycott KM, Azzariti DR, Hamosh A, Rehm HL. 2022. Seven years since the launch of the
779 Matchmaker Exchange: The evolution of genomic matchmaking. *Hum Mutat* **43**: 659–667.
780 <https://pubmed.ncbi.nlm.nih.gov/35537081/> (Accessed July 19, 2022).
781
- 782 Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Källberg M, Cox AJ, Kruglyak S,
783 Saunders CT. 2016. Manta: rapid detection of structural variants and indels for germline
784 and cancer sequencing applications. *Bioinformatics* **32**: 1220–2.
785 <https://github.com/Illumina/manta>. (Accessed June 22, 2020).
786
- 787 Cheng H, Concepcion GT, Feng X, Zhang H, Li H. 2021. Haplotype-resolved de novo assembly
788 using phased assembly graphs with hifiasm. *Nat Methods* **18**: 170–175.
789 <https://pubmed.ncbi.nlm.nih.gov/33526886/> (Accessed February 6, 2024).
790
- 791 Cohen ASA, Farrow EG, Abdelmoity AT, Alaimo JT, Amudhavalli SM, Anderson JT, Bansal L,
792 Bartik L, Baybayan P, Belden B, et al. 2022. Genomic answers for children: Dynamic
793 analyses of >1000 pediatric rare disease genomes. *Genetics in Medicine* **24**: 1336–1348.
794 <https://pubmed.ncbi.nlm.nih.gov/35305867/> (Accessed August 24, 2022).
795
- 796 Collins RL, Brand H, Karczewski KJ, Zhao X, Alföldi J, Francioli LC, Khera A V., Lowther C,
797 Gauthier LD, Wang H, et al. 2020. A structural variation reference for medical and
798 population genetics. *Nature* **581**: 444–451. <https://pubmed.ncbi.nlm.nih.gov/32461652/>
799 (Accessed February 15, 2024).
800
- 801 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
802 McCarthy SA, Davies RM. 2021. Twelve years of SAMtools and BCFtools. *Gigascience* **10**:
803 1–4. <https://dx.doi.org/10.1093/gigascience/giab008> (Accessed August 1, 2024).
804
- 805 Danis D, Jacobsen JOB, Balachandran P, Zhu Q, Yilmaz F, Reese J, Haimel M, Lyon GJ,
806 Helbig I, Mungall CJ, et al. 2022. SvAnna: efficient and accurate pathogenicity prediction of
807 coding and regulatory structural variants in long-read genome sequencing. *Genome Med*
808 **14**. <https://pubmed.ncbi.nlm.nih.gov/35484572/> (Accessed February 12, 2024).
809
- 810 De Coster W, Rademakers R. 2023. NanoPack2: population-scale evaluation of long-read
811 sequencing data. *Bioinformatics* **39**. <https://dx.doi.org/10.1093/bioinformatics/btad311>
812 (Accessed August 1, 2024).
813
- 814 Deciphering Developmental Disorders S. 2015. Large-scale discovery of novel genetic causes
815 of developmental disorders. *Nature* **519**: 223–228.
816 <https://www.ncbi.nlm.nih.gov/pubmed/25533962>.
817
- 818 Del Gobbo GF, Wang X, Couse M, Mackay L, Goldsmith C, Marshall AE, Liang Y, Lambert C,
819 Zhang S, Dhillon H, et al. 2023. Long-read genome sequencing reveals a novel intronic
820 retroelement insertion in NR5A1 associated with 46,XY differences of sexual development.

- 821 *Am J Med Genet A*. <http://www.ncbi.nlm.nih.gov/pubmed/38131126> (Accessed February
822 13, 2024).
- 823
- 824 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras
825 TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15.
826 [/pmc/articles/PMC3530905/](https://pubmed.ncbi.nlm.nih.gov/24752237/) (Accessed August 1, 2024).
- 827
- 828 Dolzhenko E, Deshpande V, Schlesinger F, Krusche P, Petrovski R, Chen S, Emig-Agius D,
829 Gross A, Narzisi G, Bowman B, et al. 2019. ExpansionHunter: a sequence-graph-based
830 tool to analyze variation in short tandem repeat regions. *Bioinformatics* **35**: 4754–4756.
831 <https://pubmed.ncbi.nlm.nih.gov/31134279/> (Accessed June 25, 2024).
- 832
- 833 Dolzhenko E, English A, Dashnow H, De Sena Brandine G, Mokveld T, Rowell WJ, Karniski C,
834 Kronenberg Z, Danzi MC, Cheung WA, et al. 2024. Characterization and visualization of
835 tandem repeats at genome scale. *Nature Biotechnology* **2024** 1–9.
836 <https://www.nature.com/articles/s41587-023-02057-3> (Accessed February 15, 2024).
- 837
- 838 East KM, Kelley W V., Cannon A, Cochran ME, Moss IP, May T, Nakano-Okuno M, Sodeke SO,
839 Edberg JC, Cimino JJ, et al. 2021. A state-based approach to genomics for rare disease
840 and population screening. *Genet Med* **23**: 777–781.
841 <https://pubmed.ncbi.nlm.nih.gov/33244164/> (Accessed July 26, 2022).
- 842
- 843 Ebert P, Audano PA, Zhu Q, Rodriguez-Martin B, Porubsky D, Bonder MJ, Sulovari A, Ebler J,
844 Zhou W, Mari RS, et al. 2021. Haplotype-resolved diverse human genomes and integrated
845 analysis of structural variation. *Science (1979)* **372**.
- 846
- 847 Felker SA, Lawlor JMJ, Hiatt SM, Thompson ML, Latner DR, Finnila CR, Bowling KM,
848 Bonnstetter ZT, Bonini KE, Kelly NR, et al. 2023. Poison exon annotations improve the
849 yield of clinically relevant variants in genomic diagnostic testing. *Genet Med* **2023**: 100884.
850 <https://pubmed.ncbi.nlm.nih.gov/37161864/> (Accessed May 29, 2023).
- 851
- 852 Fukuda H, Mizuguchi T, Doi H, Kameyama S, Kunii M, Joki H, Takahashi T, Komiya H, Sasaki
853 M, Miyaji Y, et al. 2023. Long-read sequencing revealing intragenic deletions in exome-
854 negative spastic paraplegias. *J Hum Genet* **68**: 689–697.
- 855
- 856 Gardner EJ, Lam VK, Harris DN, Chuang NT, Scott EC, Stephen Pittard W, Mills RE, Devine
857 SE. 2017. The mobile element locator tool (MELT): Population-scale mobile element
858 discovery and biology. *Genome Res* **27**: 1916–1929.
- 859
- 860 Groza C, Schwendinger-Schreck C, Cheung WA, Farrow EG, Thiffault I, Lake J, Rizzo WB,
861 Evrony G, Curran T, Bourque G, et al. 2024. Pangenome graphs improve the analysis of
862 structural variants in rare genetic diseases. *Nat Commun* **15**: 657.
863 <https://pubmed.ncbi.nlm.nih.gov/38253606/> (Accessed February 13, 2024).
- 864

- 865 Hamosh A, Wohler E, Martin R, Griffith S, Rodrigues E da S, Antonescu C, Doheny KF, Valle D,
866 Sobreira N. 2022. The impact of GeneMatcher on international data sharing and
867 collaboration. *Hum Mutat* **43**: 668–673. <https://pubmed.ncbi.nlm.nih.gov/35170833/>
868 (Accessed July 19, 2022).
869
- 870 Hao YH, Fountain MD, Fon Tacer K, Xia F, Bi W, Kang SHL, Patel A, Rosenfeld JA, Le Caignec
871 C, Isidor B, et al. 2015. USP7 Acts as a Molecular Rheostat to Promote WASH-Dependent
872 Endosomal Protein Recycling and Is Mutated in a Human Neurodevelopmental Disorder.
873 *Mol Cell* **59**: 956–969. <https://pubmed.ncbi.nlm.nih.gov/26365382/> (Accessed February 11,
874 2024).
875
- 876 Hartley T, Soubry É, Acker M, Osmond M, Couse M, Gillespie MK, Ito Y, Marshall AE, Lemire
877 G, Huang L, et al. 2023. Bridging clinical care and research in Ontario, Canada:
878 Maximizing diagnoses from reanalysis of clinical exome sequencing data. *Clin Genet* **103**:
879 288–300.
880
- 881 Hiatt SM, Amaral MD, Bowling KM, Finnila CR, Thompson ML, Gray DE, Lawlor JMJ, Cochran
882 JN, Bebin EM, Brothers KB, et al. 2018. Systematic reanalysis of genomic data improves
883 quality of variant interpretation. *Clin Genet* **94**.
884
- 885 Hiatt SM, Lawlor JMJ, Handley LH, Ramaker RC, Rogers BB, Partridge EC, Boston LB,
886 Williams M, Plott CB, Jenkins J, et al. 2021. Long-read genome sequencing for the
887 molecular diagnosis of neurodevelopmental disorders. *HGG Adv* **2**: 100023.
888 <http://www.ncbi.nlm.nih.gov/pubmed/33937879> (Accessed June 9, 2021).
889
- 890 Jadhav B, Garg P, van Vugt JJ, Ibanez K, Gagliardi D, Lee W, Shadrina M, Mokveld T,
891 Dolzhenko E, Martin-Trujillo A, et al. 2023. A phenome-wide association study of
892 methylated GC-rich repeats identifies a GCC repeat expansion in *AFF3* as a significant
893 cause of intellectual disability. *medRxiv*. <https://pubmed.ncbi.nlm.nih.gov/37205357/>
894 (Accessed February 6, 2024).
895
- 896 Juven A, Nambot S, Piton A, Jean-Marçais N, Masurel A, Callier P, Marle N, Mosca-Boidron AL,
897 Kuentz P, Philippe C, et al. 2020. Primrose syndrome: a phenotypic comparison of patients
898 with a *ZBTB20* missense variant versus a 3q13.31 microdeletion including *ZBTB20*. *Eur J*
899 *Hum Genet* **28**: 1044–1055. <https://pubmed.ncbi.nlm.nih.gov/32071410/> (Accessed
900 February 15, 2024).
901
- 902 Kendig KI, Baheti S, Bockol MA, Drucker TM, Hart SN, Heldenbrand JR, Hernaez M, Hudson
903 ME, Kalmbach MT, Klee EW, et al. 2019. Sentieon DNaseq Variant Calling Workflow
904 Demonstrates Strong Computational Performance and Accuracy. *Front Genet* **10**.
905 <https://pubmed.ncbi.nlm.nih.gov/31481971/> (Accessed February 15, 2024).
906
- 907 Kilich G, Hassey K, Behrens EM, Falk M, Vanderver A, Rader DJ, Cahill PJ, Raper A, Zhang Z,
908 Westerfer D, et al. 2024. Kagami Ogata syndrome: a small deletion refines critical region

- 909 for imprinting. *NPJ Genom Med* **9**: 5. <http://www.ncbi.nlm.nih.gov/pubmed/38212313>
910 (Accessed February 13, 2024).
911
- 912 Kirsche M, Prabhu G, Sherman R, Ni B, Battle A, Aganezov S, Schatz MC. 2023. Jasmine and
913 Iris: population-scale structural variant comparison and analysis. *Nat Methods* **20**: 408–
914 417. <https://pubmed.ncbi.nlm.nih.gov/36658279/> (Accessed February 15, 2024).
915
- 916 Krumsiek J, Arnold R, Rattei T. 2007. Gepard: a rapid and sensitive tool for creating dotplots on
917 genome scale. **23**: 1026–1028. <http://mips.gsf.de/services/analysis/gepard> (Accessed May
918 18, 2020).
919
- 920 Lek M, Karczewski KJ, Minikel E V, Samocha KE, Banks E, Fennell T, O'Donnell-Luria AH,
921 Ware JS, Hill AJ, Cummings BB, et al. 2016. Analysis of protein-coding genetic variation in
922 60,706 humans. *Nature* **536**: 285–291. <https://www.ncbi.nlm.nih.gov/pubmed/27535533>.
923
- 924 Li H. 2018. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094–
925 3100. <https://github.com/ruanjue/smartdenovo>; (Accessed September 7, 2020).
926
- 927 Liao WW, Asri M, Ebler J, Doerr D, Haukness M, Hickey G, Lu S, Lucas JK, Monlong J, Abel
928 HJ, et al. 2023. A draft human pangenome reference. *Nature* **617**.
929
- 930 Liu P, Meng L, Normand EA, Xia F, Song X, Ghazi A, Rosenfeld J, Magoulas PL, Braxton A,
931 Ward P, et al. 2019. Reanalysis of Clinical Exome Sequencing Data. *New England Journal*
932 *of Medicine* **380**: 2478–2480.
933
- 934 Logsdon GA, Vollger MR, Eichler EE. 2020. Long-read human genome sequencing and its
935 applications. *Nat Rev Genet* **21**: 597–614. [https://www.nature.com/articles/s41576-020-](https://www.nature.com/articles/s41576-020-0236-x)
936 [0236-x](https://www.nature.com/articles/s41576-020-0236-x) (Accessed November 4, 2020).
937
- 938 Mahmoud M, Huang Y, Garimella K, Audano PA, Wan W, Prasad N, Handsaker RE, Hall S,
939 Pionzio A, Schatz MC, et al. 2024. Utility of long-read sequencing for All of Us. *Nat*
940 *Commun* **15**: 837. <https://pubmed.ncbi.nlm.nih.gov/38281971/> (Accessed February 5,
941 2024).
942
- 943 Martin M, Patterson M, Garg S, Fischer SO, Pisanti N, Klau GW, Schönhuth A, Marschall T.
944 WhatsHap: fast and accurate read-based phasing. <https://doi.org/10.1101/085050>
945 (Accessed August 1, 2024).
946
- 947 Miller DE, Sulovari A, Wang T, Loucks H, Hoekzema K, Munson KM, Lewis AP, Fuerte EPA,
948 Paschal CR, Walsh T, et al. 2021. Targeted long-read sequencing identifies missing
949 disease-causing variation. *Am J Hum Genet* **108**: 1436–1449.
950 <https://pubmed.ncbi.nlm.nih.gov/34216551/> (Accessed February 13, 2024).
951

- 952 Muenzen KD, Amendola LM, Kauffman TL, Mittendorf KF, Bensen JT, Chen F, Green R, Powell
953 BC, Kvale M, Angelo F, et al. 2022. Lessons learned and recommendations for data
954 coordination in collaborative research: The CSER consortium experience. *Human Genetics
955 and Genomics Advances* **3**: 100120.
956 <http://www.cell.com/article/S2666247722000367/fulltext> (Accessed August 1, 2024).
957
- 958 Nakamura H, Doi H, Mitsuhashi S, Miyatake S, Katoh K, Frith MC, Asano T, Kudo Y, Ikeda T,
959 Kubota S, et al. 2020. Long-read sequencing identifies the pathogenic nucleotide repeat
960 expansion in RFC1 in a Japanese case of CANVAS. *J Hum Genet* **65**: 475–480.
961 <https://www.nature.com/articles/s10038-020-0733-y> (Accessed June 27, 2021).
962
- 963 Ng SB, Buckingham KJ, Lee C, Bigham AW, Tabor HK, Dent KM, Huff CD, Shannon PT, Jabs
964 EW, Nickerson DA, et al. 2010. Exome sequencing identifies the cause of a mendelian
965 disorder. *Nat Genet* **42**: 30–35. <https://pubmed.ncbi.nlm.nih.gov/19915526/> (Accessed
966 February 27, 2024).
967
- 968 Ng SB, Turner EH, Robertson PD, Flygare SD, Bigham AW, Lee C, Shaffer T, Wong M,
969 Bhattacharjee A, Eichler EE, et al. 2009. Targeted capture and massively parallel
970 sequencing of 12 human exomes. *Nature* **461**: 272–276.
971 <https://pubmed.ncbi.nlm.nih.gov/19684571/> (Accessed February 13, 2024).
972
- 973 Niemi MEK, Martin HC, Rice DL, Gallone G, Gordon S, Kelemen M, McAloney K, McRae J,
974 Radford EJ, Yu S, et al. 2018. Common genetic variants contribute to risk of rare severe
975 neurodevelopmental disorders. *Nature* **562**: 268–271.
976 <https://pubmed.ncbi.nlm.nih.gov/30258228/> (Accessed February 11, 2024).
977
- 978 Nurk S, Koren S, Rhie A, Rautiainen M, Bizikadze A V., Mikheenko A, Vollger MR, Altemose N,
979 Uralsky L, Gershman A, et al. 2022. The complete sequence of a human genome. *Science*
980 **376**: 44–53. <https://pubmed.ncbi.nlm.nih.gov/35357919/> (Accessed January 28, 2024).
981
- 982 Papadimitriou S, Gazzo A, Versbraegen N, Nachtegaele C, Aerts J, Moreau Y, Van Dooren S,
983 Nowé A, Smits G, Lenaerts T. 2019. Predicting disease-causing variant combinations. *Proc
984 Natl Acad Sci U S A* **116**: 11878–11887. <https://pubmed.ncbi.nlm.nih.gov/31127050/>
985 (Accessed February 14, 2024).
986
- 987 Pauper M, Kucuk E, Wenger AM, Chakraborty S, Baybayan P, Kwint M, van der Sanden B,
988 Nelen MR, Derks R, Brunner HG, et al. 2021. Long-read trio sequencing of individuals with
989 unsolved intellectual disability. *European Journal of Human Genetics* **29**: 637–648.
990
- 991 Pedersen BS, Bhetariya PJ, Brown J, Kravitz SN, Marth G, Jensen RL, Bronner MP, Underhill
992 HR, Quinlan AR. 2020. Somalier: rapid relatedness estimation for cancer and germline
993 studies using efficient genome sketches. *Genome Med* **12**.
994 <https://pubmed.ncbi.nlm.nih.gov/32664994/> (Accessed February 15, 2024).
995

- 996 Pedersen BS, Quinlan AR. 2017. Who's Who? Detecting and Resolving Sample Anomalies in
997 Human DNA Sequencing Studies with Peddy. *Am J Hum Genet* **100**: 406–413.
998 <https://pubmed.ncbi.nlm.nih.gov/28190455/> (Accessed February 15, 2024).
999
- 1000 Philippakis AA, Azzariti DR, Beltran S, Brookes AJ, Brownstein CA, Brudno M, Brunner HG,
1001 Buske OJ, Carey K, Doll C, et al. 2015. The Matchmaker Exchange: A Platform for Rare
1002 Disease Gene Discovery. *Hum Mutat* **36**: 915–921.
1003 <http://doi.wiley.com/10.1002/humu.22858> (Accessed November 25, 2018).
1004
- 1005 Poplin R, Chang PC, Alexander D, Schwartz S, Colthurst T, Ku A, Newburger D, Dijamco J,
1006 Nguyen N, Afshar PT, et al. 2018. A universal snp and small-indel variant caller using deep
1007 neural networks. *Nat Biotechnol* **36**: 983.
1008
- 1009 Posey JE, Harel T, Liu P, Rosenfeld JA, James RA, Coban Akdemir ZH, Walkiewicz M, Bi W,
1010 Xiao R, Ding Y, et al. 2017. Resolution of Disease Phenotypes Resulting from Multilocus
1011 Genomic Variation. *N Engl J Med* **376**: 21–31. <https://pubmed.ncbi.nlm.nih.gov/27959697/>
1012 (Accessed February 15, 2024).
1013
- 1014 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
1015 features. *Bioinformatics* **26**: 841–842. <https://pubmed.ncbi.nlm.nih.gov/20110278/>
1016 (Accessed August 1, 2024).
1017
- 1018 R Core Team. 2023. R Core Team 2023 R: A language and environment for statistical
1019 computing. R foundation for statistical computing. <https://www.R-project.org/>. *R Foundation*
1020 *for Statistical Computing*.
1021
- 1022 Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, Grody WW, Hegde M, Lyon E,
1023 Spector E, et al. 2015. Standards and guidelines for the interpretation of sequence
1024 variants: a joint consensus recommendation of the American College of Medical Genetics
1025 and Genomics and the Association for Molecular Pathology. *Genet Med* **17**: 405–424.
1026 <https://www.ncbi.nlm.nih.gov/pubmed/25741868>.
1027
- 1028 Riggs ER, Andersen EF, Cherry AM, Kantarci S, Kearney H, Patel A, Raca G, Ritter DI, South
1029 ST, Thorland EC, et al. 2020. Technical standards for the interpretation and reporting of
1030 constitutional copy-number variants: a joint consensus recommendation of the American
1031 College of Medical Genetics and Genomics (ACMG) and the Clinical Genome Resource
1032 (ClinGen). *Genetics in Medicine* **22**: 245–257. <https://pubmed.ncbi.nlm.nih.gov/31690835/>
1033 (Accessed September 8, 2020).
1034
- 1035 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP. 2011.
1036 Integrative genomics viewer. *Nat Biotechnol* **29**: 24–26.
1037 <https://pubmed.ncbi.nlm.nih.gov/21221095/> (Accessed February 15, 2024).
1038

- 1039 Sakamoto M, Kurosawa K, Tanoue K, Iwama K, Ishida F, Watanabe Y, Okamoto N, Tsuchida
1040 N, Uchiyama Y, Koshimizu E, et al. 2024. A heterozygous germline deletion within USP8
1041 causes severe neurodevelopmental delay with multiorgan abnormalities. *J Hum Genet* **69**:
1042 85–90. <http://www.ncbi.nlm.nih.gov/pubmed/38030753> (Accessed February 13, 2024).
1043
- 1044 Sanghvi R V., Buhay CJ, Powell BC, Tsai EA, Dorschner MO, Hong CS, Lebo MS, Sasson A,
1045 Hanna DS, McGee S, et al. 2018. Characterizing reduced coverage regions through
1046 comparison of exome and genome sequencing data across 10 centers. *Genet Med* **20**:
1047 855–866. <https://pubmed.ncbi.nlm.nih.gov/29144510/> (Accessed February 5, 2024).
1048
- 1049 Schobers G, Schieving JH, Yntema HG, Pennings M, Pfundt R, Derks R, Hofste T, de Wijs I,
1050 Wieskamp N, van den Heuvel S, et al. 2022. Reanalysis of exome negative patients with
1051 rare disease: a pragmatic workflow for diagnostic applications. *Genome Med* **14**.
1052
- 1053 Schuy J, Grochowski CM, Carvalho CMB, Lindstrand A. 2022. Complex genomic
1054 rearrangements: an underestimated cause of rare diseases. *Trends in Genetics* **0**.
1055 <http://www.cell.com/article/S0168952522001457/fulltext> (Accessed July 31, 2022).
1056
- 1057 Smith T, Heger A, Sudbery I. 2017. UMI-tools: Modelling sequencing errors in Unique Molecular
1058 Identifiers to improve quantification accuracy. *Genome Res* **27**: gr.209601.116.
1059 <https://genome.cshlp.org/content/early/2017/01/18/gr.209601.116> (Accessed August 1,
1060 2024).
1061
- 1062 Sobreira N, Schiettecatte F, Valle D, Hamosh A. 2015. GeneMatcher: a matching tool for
1063 connecting investigators with an interest in the same gene. *Hum Mutat* **36**: 928–30.
1064 <http://doi.wiley.com/10.1002/humu.22844> (Accessed November 25, 2018).
1065
- 1066 Sobreira NLM, Arachchi H, Buske OJ, Chong JX, Hutton B, Foreman J, Schiettecatte F, Groza
1067 T, Jacobsen JOB, Haendel MA, et al. 2017. Matchmaker Exchange. *Curr Protoc Hum*
1068 *Genet* **95**: 9.31.1-9.31.15. <http://doi.wiley.com/10.1002/cphg.50> (Accessed November 25,
1069 2018).
1070
- 1071 Srivastava S, Love-Nichols JA, Dies KA, Ledbetter DH, Martin CL, Chung WK, Firth H V.,
1072 Frazier T, Hansen RL, Prock L, et al. 2019. Meta-analysis and multidisciplinary consensus
1073 statement: exome sequencing is a first-tier clinical diagnostic test for individuals with
1074 neurodevelopmental disorders. *Genet Med* **21**: 2413–2421.
1075 <https://pubmed.ncbi.nlm.nih.gov/31182824/> (Accessed December 6, 2022).
1076
- 1077 Team Rs. 2020. RStudio: Integrated Development for R. *RStudio, Inc, Boston, MA*.
1078
- 1079 Vaser R, Sović I, Nagarajan N, Šikić M. 2017. Fast and accurate de novo genome assembly
1080 from long uncorrected reads. *Genome Res* **27**: 737–746.
1081 <https://pubmed.ncbi.nlm.nih.gov/28100585/> (Accessed August 1, 2024).
1082

- 1083 Wagner J, Olson ND, Harris L, McDaniel J, Cheng H, Fungtammasan A, Hwang YC, Gupta R,
 1084 Wenger AM, Rowell WJ, et al. 2022. Curated variation benchmarks for challenging
 1085 medically relevant autosomal genes. *Nat Biotechnol* **40**: 672–680.
 1086 <https://pubmed.ncbi.nlm.nih.gov/35132260/> (Accessed February 15, 2024).
 1087
- 1088 Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, Ebler J,
 1089 Fungtammasan A, Kolesnikov A, Olson ND, et al. 2019. Accurate circular consensus long-
 1090 read sequencing improves variant detection and assembly of a human genome. *Nat*
 1091 *Biotechnol* **37**: 1155–1162.
 1092
- 1093 Wickham H. 2011. ggplot2. *Wiley Interdiscip Rev Comput Stat* **3**.
 1094
- 1095 Zook JM, McDaniel J, Olson ND, Wagner J, Parikh H, Heaton H, Irvine SA, Trigg L, Truty R,
 1096 McLean CY, et al. 2019. An open resource for accurately benchmarking small variant and
 1097 reference calls. *Nat Biotechnol* **37**: 561–566.
 1098

1099 **FIGURE LEGENDS**

1100 **Figure 1. A *de novo*, 4 Mb paracentric inversion in proband 1, affecting *ZBTB20*.** A, B.

1101 Visualization of a subset of proband and parent reads in IGV at the 5' (A, Chr 3:110,477,108-
 1102 110,477,357) and 3' (B, Chr 3:114,639,120-114,639,284) breakpoints (black arrowheads)
 1103 indicate a *de novo* event. Reads in A show the SV is present on the proband's paternal allele.
 1104 C. Schematic of the proband's inversion. Section B-C is inverted, and junctions 1 and 2 (jct1,
 1105 jct2) are displayed. The inversion disrupts *ZBTB20*, with a promoter (P) and exons 1-6 of
 1106 NM_001348800.3 (yellow boxes) remaining at the distal end of the chromosome near jct2/D and
 1107 exons 7-12 of NM_001348800.3 (green boxes) moving upstream to the breakpoint at the more
 1108 proximal end of the chromosome (near jct1/A). Sequence level resolution of the inversion is also
 1109 shown. D. Alignment of the proband's assembled paternal contig versus the reference genome
 1110 supports the inversion.

1111 **Figure 2. Two *ALS2* deletions in trans in Proband 2.** A. Visualization of proband and parent
 1112 reads in IGV (Chr 2:201,719,400-201,741,000) indicate two overlapping deletions in *ALS2*; a
 1113 smaller maternal deletion (pink bar) and a larger paternal deletion (blue bar/arrow). Note that
 1114 the three breakpoints shown here overlap *Alu* sequence in the RepeatMasker track. Alignment

1115 of the proband's assembled maternal (B) or paternal contig (C) versus the reference genome
1116 support the two deletions (red dashed lines).

1117 **Figure 3. Proband 3 has a 4 kb insertion in the 3' UTR of HCFC1.** A. The proband's family
1118 has a history of X-linked intellectual disability, as the proband (arrow) and two other male
1119 relatives (gray squares) are affected. B. Model of the relative length of the insertion in the 3'
1120 UTR NM_005334.3. Only exons 24-26 are shown. C. The hemizygous insertion is likely
1121 inherited from a heterozygous carrier mother, as indicated by srGS reads. Note that some reads
1122 contain unaligned ends (multicolored bases) while some span the entire insertion (purple line at
1123 the 5' end of the TSD); all the proband's reads support the insertion. A target site duplication
1124 (TSD) is a hallmark of an L1-mediated insertion.

Table 1. Demographics, and sequencing and assembly metrics for the cohort

Sex	Male	63 (66%)			
	Female	33 (34%)			
Predicted Major Genetic Ancestry	European (EUR)	71.88%			
	African (AFR)	20.83%			
	Admixed American (AMR)	3.13%			
	Southeast Asian (SAS)	1.04%			
	Unspecified Admixed (UNKNOWN)	3.13%			
Sequencing Metrics		Median	Min	Max	
	Sequenced Bases (Gb)	80.9	53.5	131.5	
	Mean Read Length (Sequenced)	16,771	11,295	21,412	
	Median Coverage (X)	27	18	44	
	Percent Covered at 10x	97.80%	95.10%	99.00%	
	Percent Covered at 20x	84.90%	36.10%	97.50%	
Proband Assembly N50 (Mb)	All contigs	29.05	2.10	71.30	
	Maternal contigs	29.30	2.10	71.30	
	Paternal contigs	28.70	2.20	67.90	
	hap1 contigs	31.70	18.20	42.90	
	hap2 contigs	27.85	14.40	40.80	
Small Variant Metrics (DeepVariant)	SNVs	4,404,564	3,997,350	5,299,670	
	Total indels	970,031	926,286	1,153,616	
Structural Variant Metrics (pbsv)	Total Structural Variants	55,586	53,834	65,120	
	Deletion	24,757	23,799	29,531	
	Duplication	3,017	2,820	3,649	
	Insertion	27,594	26,441	32,437	
	Inversion	121	99	155	
	Breakend	193	124	312	
	Structural Variants >50bp	25,218	23,994	29,159	
	Deletion	10,303	9,804	12,306	
	Duplication	641	536	742	
	Insertion	14,012	13,344	16,159	
Inversion	121	99	155		
Breakend	193	124	312		

Table 2. Variants identified by long-read sequencing.

Proband ID	Gene(s) Affected	HGVS Nomenclature	Inheritance	Variant Classification	Case-Level Classification	ACMG/ClinGen Evidence Codes	SV, SNV or TRE	Step of srGS loss*	Orthogonal Validation
1	<i>ZBTB20</i>	NC_000003.12:g.1104772_73_114639202_inv	de novo	LP	Likely Diagnostic	2C(+1.00), 5A (+0.15)	SV	Filtering/QC	Yes, Research Sanger of breakpoints
2	<i>ALS2</i>	NC_000002.12:g.2017204_35-201725085_del; NC_000002.12:g.2001151_81-201739349_del	biparental	LP; P	Definitive Diagnostic	2E (+0.90);2D-4 (+0.90), 3B (+0.45)	SVx2	Filtering/QC	Pending, supported by srGS data
3	<i>HCFC1</i>	NC_000023.11:g.1539486_02_ins4902, NM_005334.3:c.*745_ins4_902	maternal (X-linked)	VUS	Uncertain	NA	SV	Variant Calling	Yes, Research PCR amplification of breakpoints. Pending
4	<i>ABAT, PMM2, USP7, etc.</i>	NC_000016.9:g.(8742452_9220783)dup_ins[(87424_52_8879961)_(9000190_9_220783)]	de novo	VUS	Uncertain	2K (+0.30)	SV	Variant Calling	Pending
5	<i>PHOX2B</i>	NM_003924.4:c.741_758dup, p.(Ala255_Ala260dup)	unknown [#]	P	Definitive Diagnostic	PS4_M, PM1_Strong, PM2_Moderate, PM6_Moderate [#]	TRE	Variant Calling	Yes, Clinical Testing*
6	<i>AFF3</i>	NM_001386135.1:c.-64-281_-64-280insGGC[90]	unknown	VUS	Uncertain	NA	TRE	Variant Calling	Pending
7	<i>SHANK3</i>	NM_033517.1:c.3161delT, p.(Lys1054Argfs*10)	de novo	P	Definitive Diagnostic	PVS1_VeryStrong, PS2_Strong, PM2_Moderate	SNV	Variant Calling	Yes, Clinical Sanger
8	<i>HNRNPU</i>	NM_031844.3:c.660_661dupAGGCGGCGGA, p.(Gly221ArgfsTer25)	de novo	P	Definitive Diagnostic	PVS1_VeryStrong, PS2_Strong, PM2_Moderate	SNV	Curation (Gene-Disease Association)	Yes, Clinical Sanger
9	<i>CSNK2B</i>	NM_001320.6:c.202C>T, p.(Gln68Ter)	paternal	LP	Likely Diagnostic	PVS1_VeryStrong, PM2_Moderate	SNV	Curation (Gene-Disease Association)	Yes, Clinical Sanger
10	<i>GNB2</i>	NM_005273.4:c.217G>A, p.(Ala73Thr)	maternal	LP	Uncertain	PS4_Moderate, PP2_Supporting, PP3_Supporting	SNV	Curation (Gene-Disease Association)	Yes, Clinical Sanger
11	<i>MCF2</i>	NM_005369.5:c.2234G>T, p.(Gly745Val)	maternal (X-linked)	VUS	Uncertain	PM2	SNV	Curation (Gene-Disease Association)	Yes, Clinical Sanger
12	<i>NOTCH3</i>	NM_000435.3:c.6409_641_0delCT, p.(Leu2137GlyfsTer104)	paternal	LP	Likely Diagnostic	PVS1_Strong, PM2_Moderate	SNV	Curation (unexpected mechanism)	Yes, Clinical Sanger
13	<i>AFF4</i>	NM_014423.4, c.879delA, p.(His294IlefsTer5)	paternal	VUS	Uncertain	PM2	SNV	Curation (unexpected mechanism)	Yes, Clinical Sanger
14	<i>KCNT2, KIF21A</i>	NC_000001.11:g.1963294_20-196344697_DUP, NM_017641.3:c.847C>T, p.(Arg283Cys); NM_017641.3:c.706C>T, p.(Gln236Ter)	paternal/biparental	VUS, VUS;LP	Uncertain	2I (+0.45); PM2, PP3; PVS1, PM2	SV; SNV(x2)	Curation (unexpected mechanism)	SV Pending, supported by srGS data; SNVs-Clinical Sanger
15	<i>NRXN1</i>	NC_000002.12:49922063_49928691del	de novo	VUS	Uncertain	2E (+0.30), 4C (+0.15), 4M (+0.30)	SV	Curation (unexpected mechanism)	Pending, supported by srGS data
16	<i>SCN1A</i>	NM_001165963.4:c.4003-603T>C	paternal	VUS	Uncertain	PM2_Moderate	SNV	Curation (noncoding variation)	Yes, Clinical Sanger

SV, Structural variant; SNV, single nucleotide variant; TRE, tandem repeat expansion. P, Pathogenic; LP, Likely Pathogenic; VUS, Variant of Uncertain Significance; NA, not applicable. *Filtering/QC, our filtering or prioritization strategy did not accurately present this for curation (or did not at all); Variant Calling, the variant was not called or was called incorrectly/inaccurately; Curation, manual curation did not result in flagging of the variant(s). [#]Independent clinical testing indicated that the PHOX2B expansion was de novo; we only sequenced the proband in our research study.

Figure 1

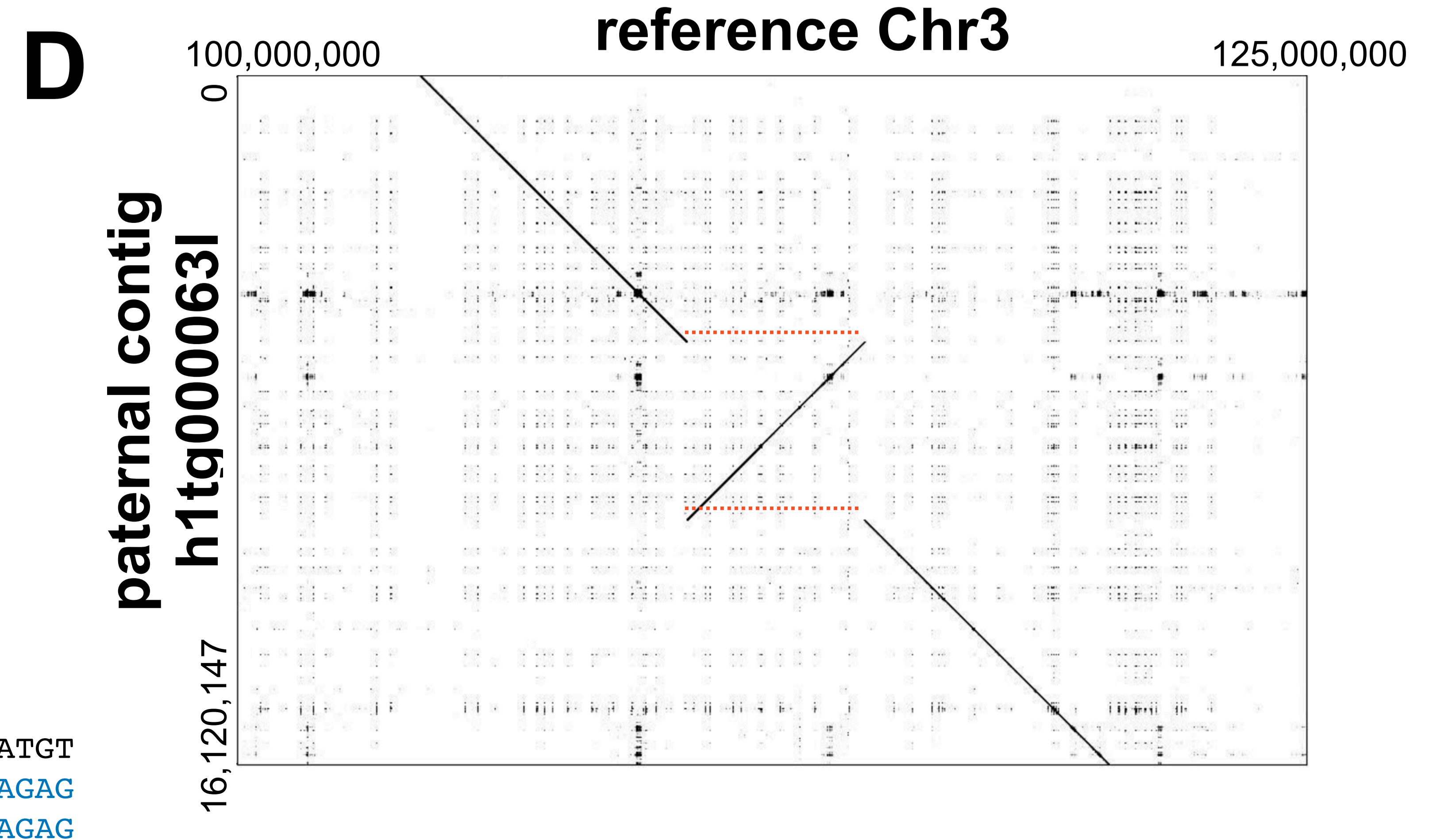
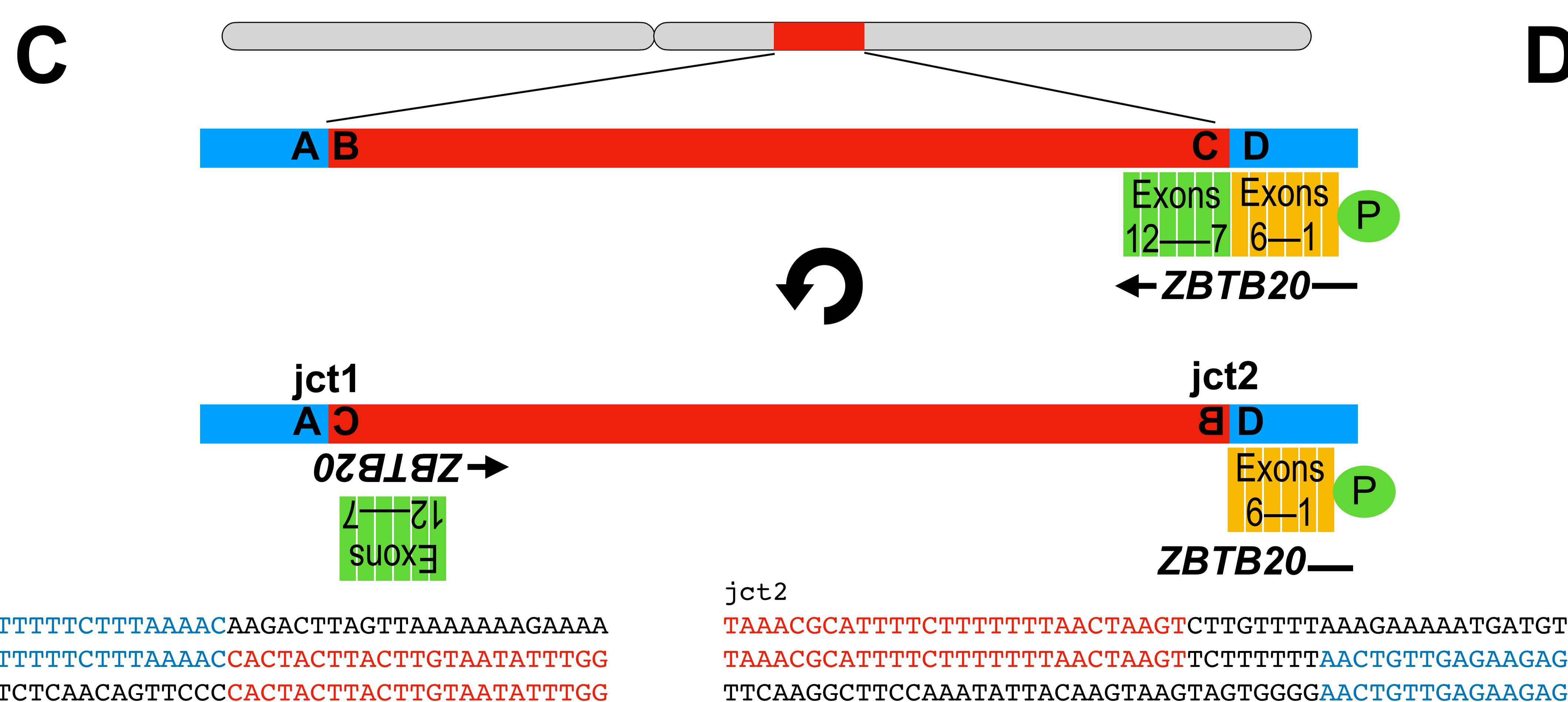
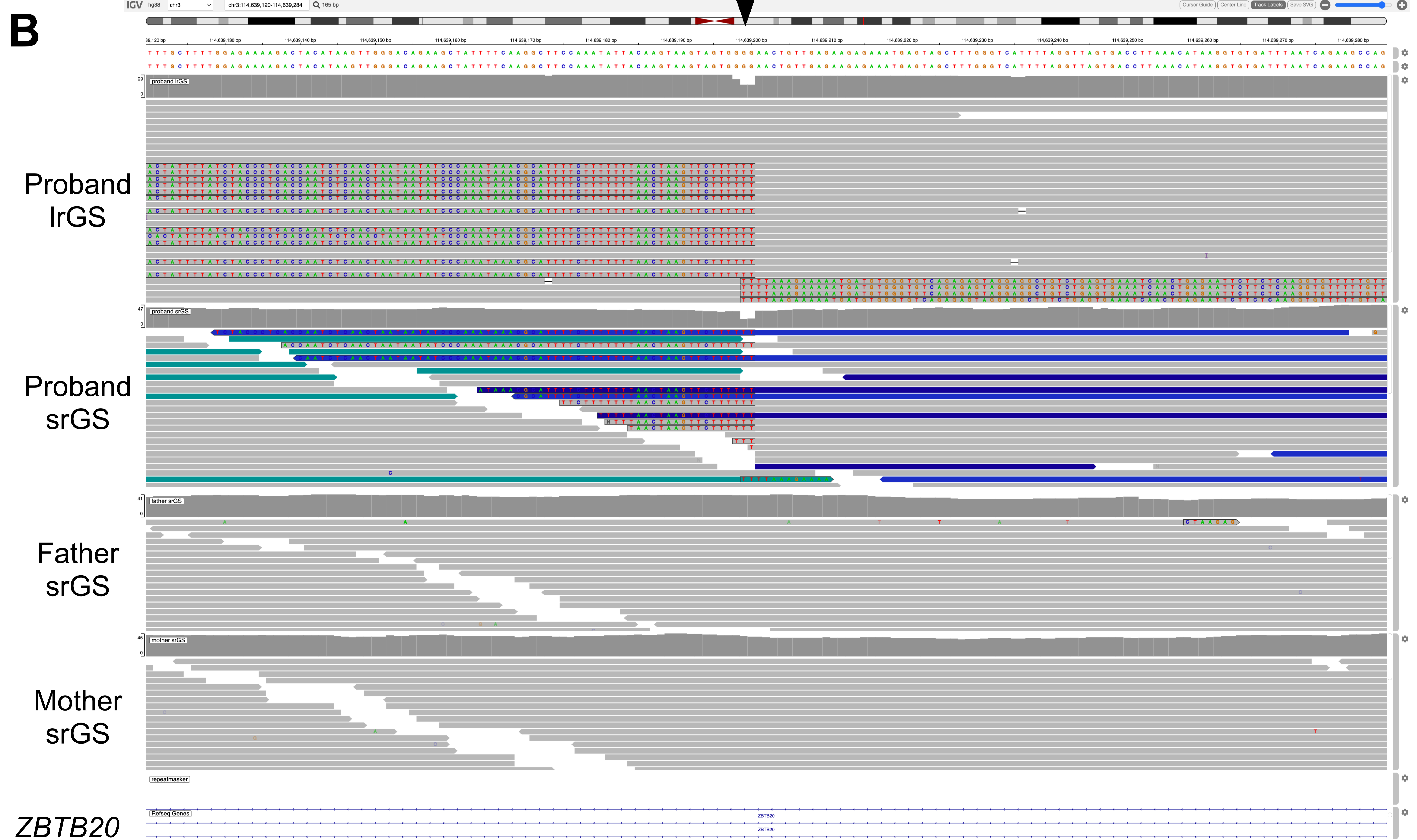
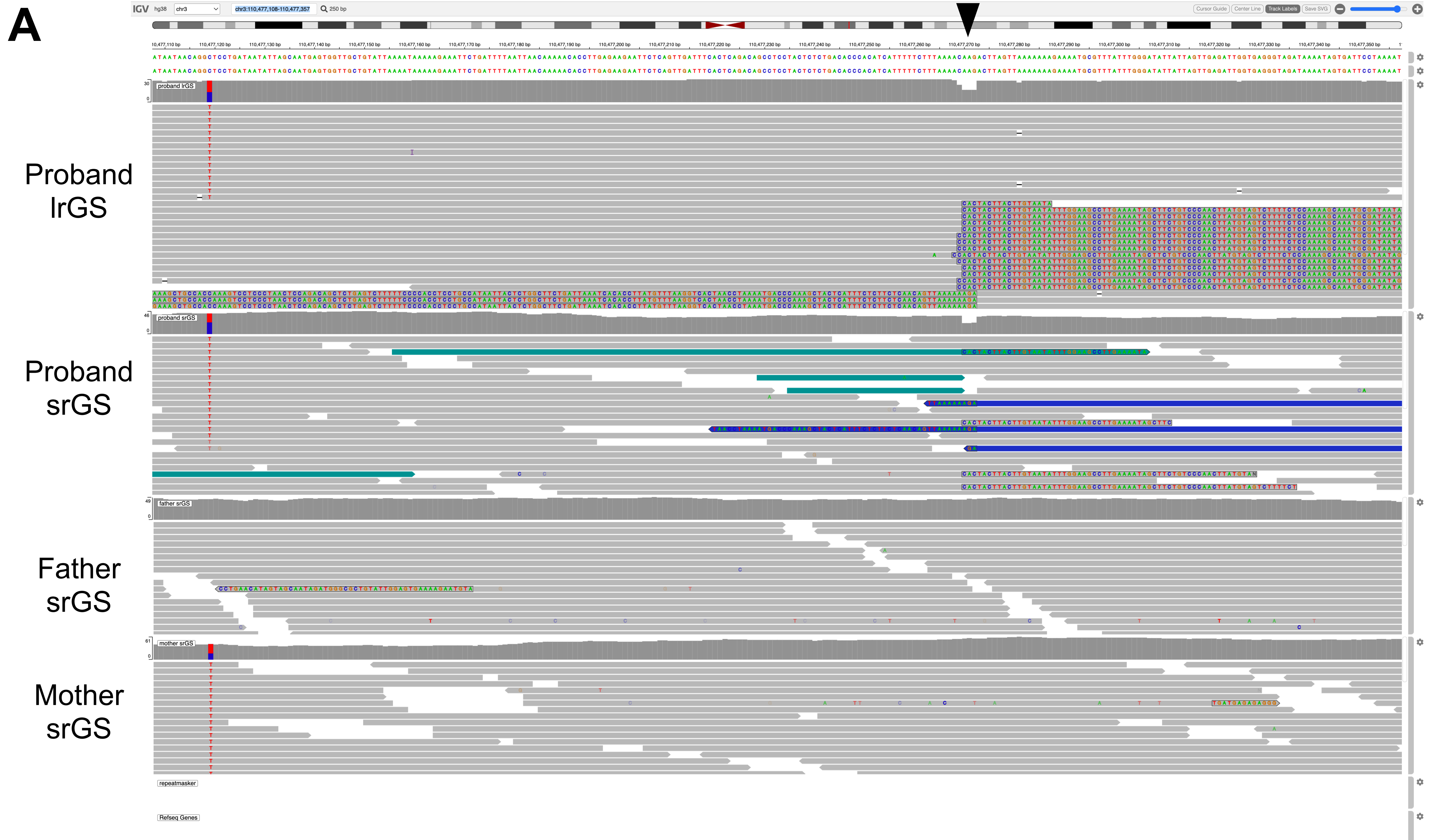


Figure 2

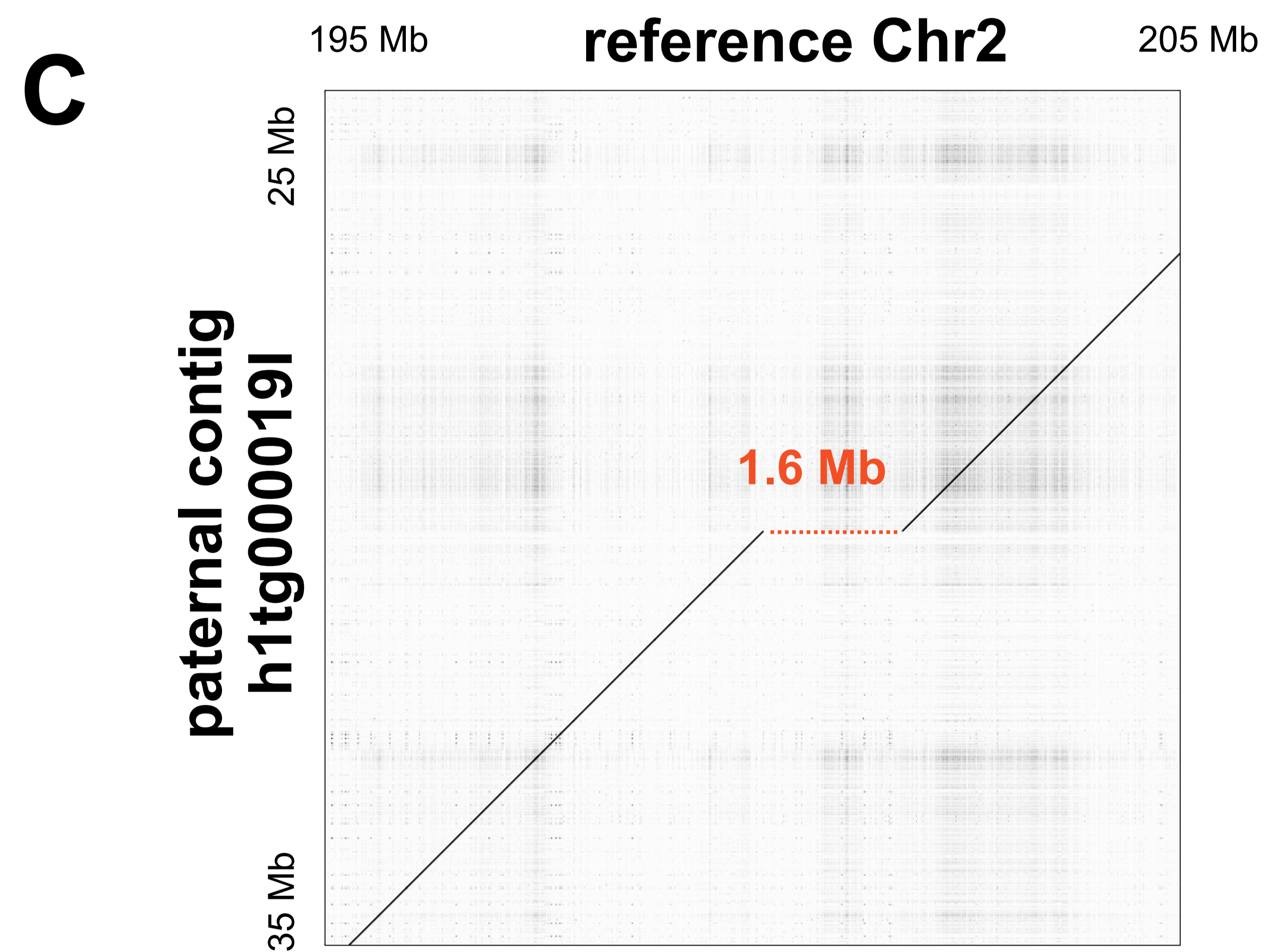
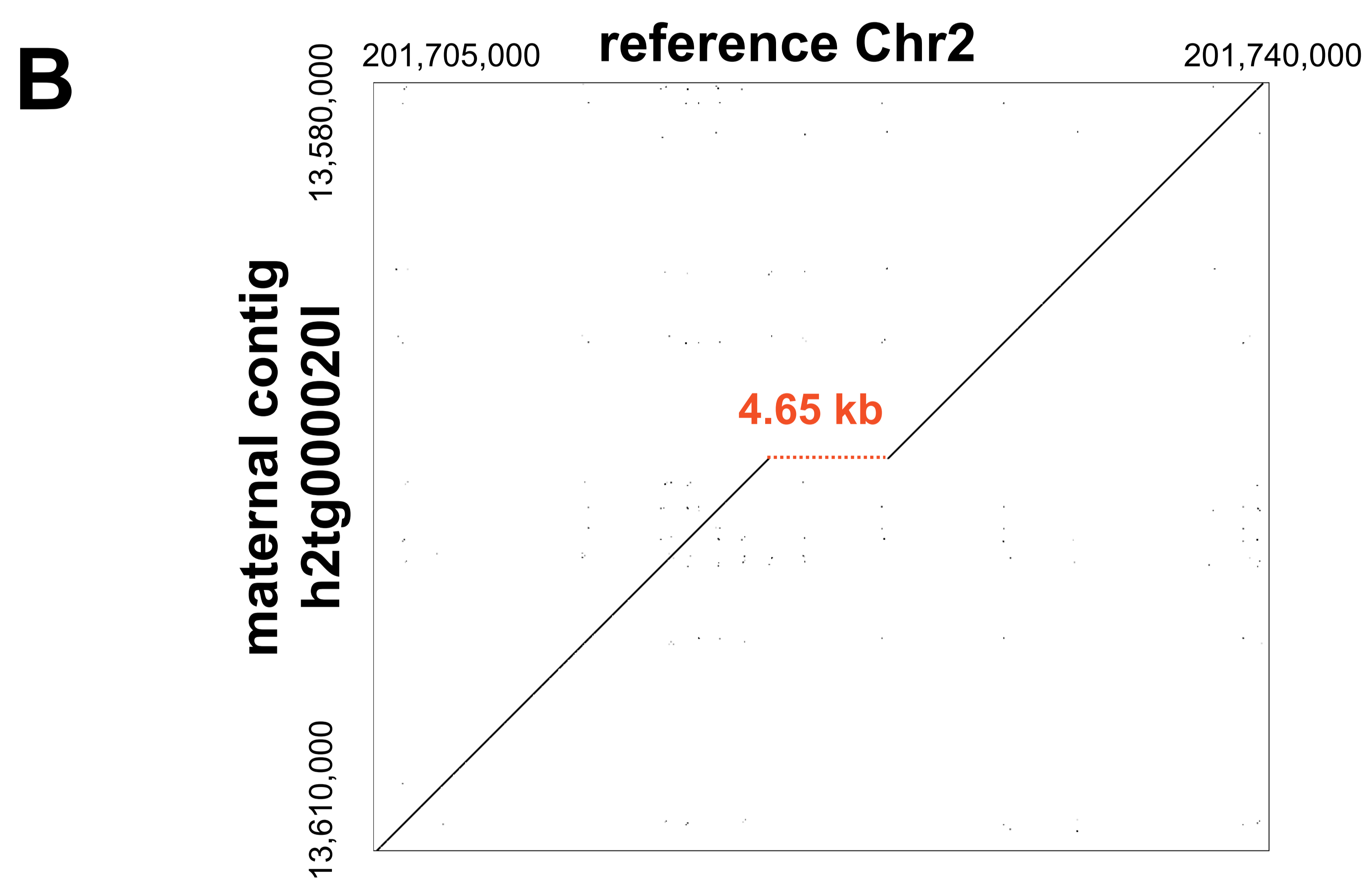
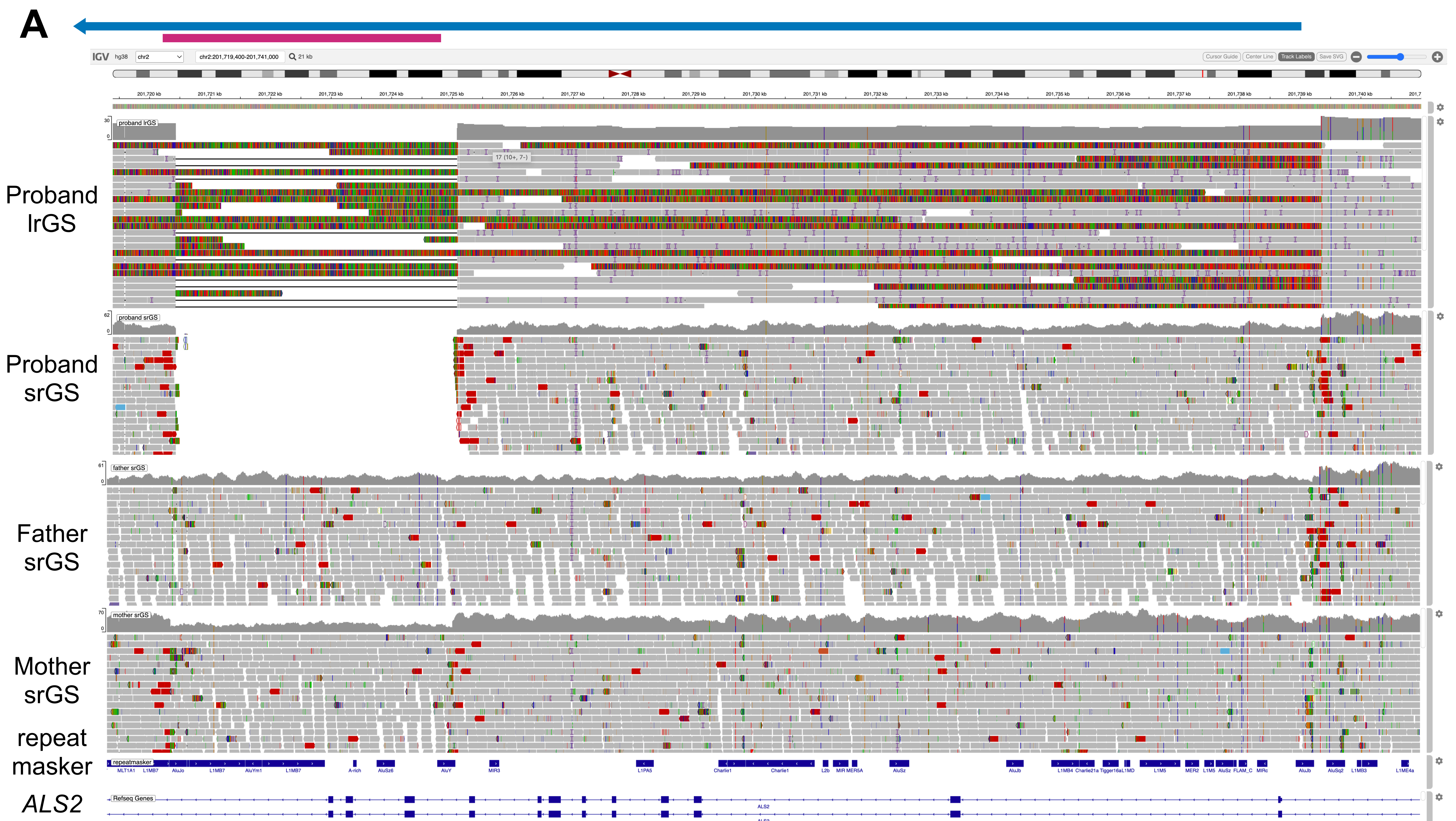
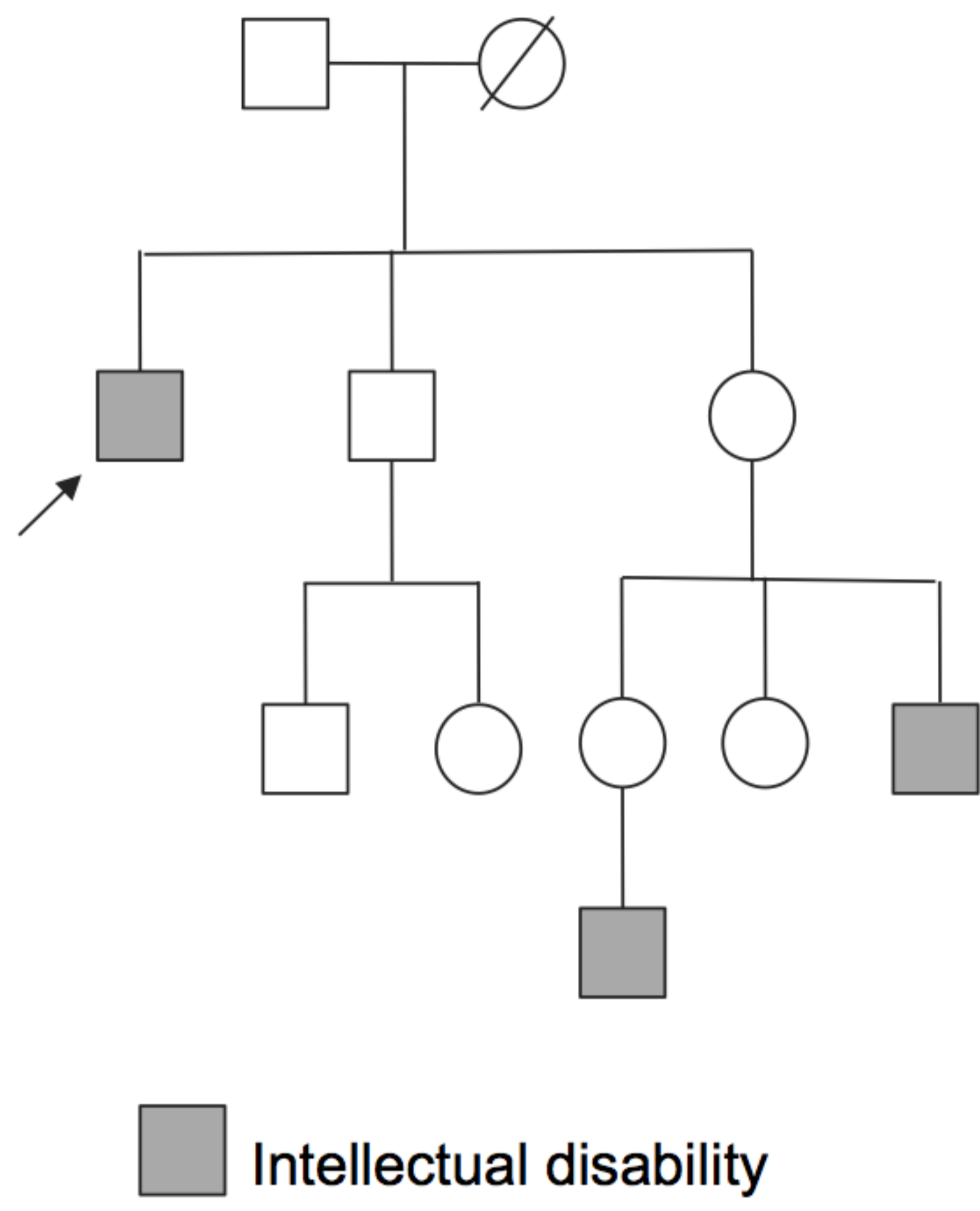
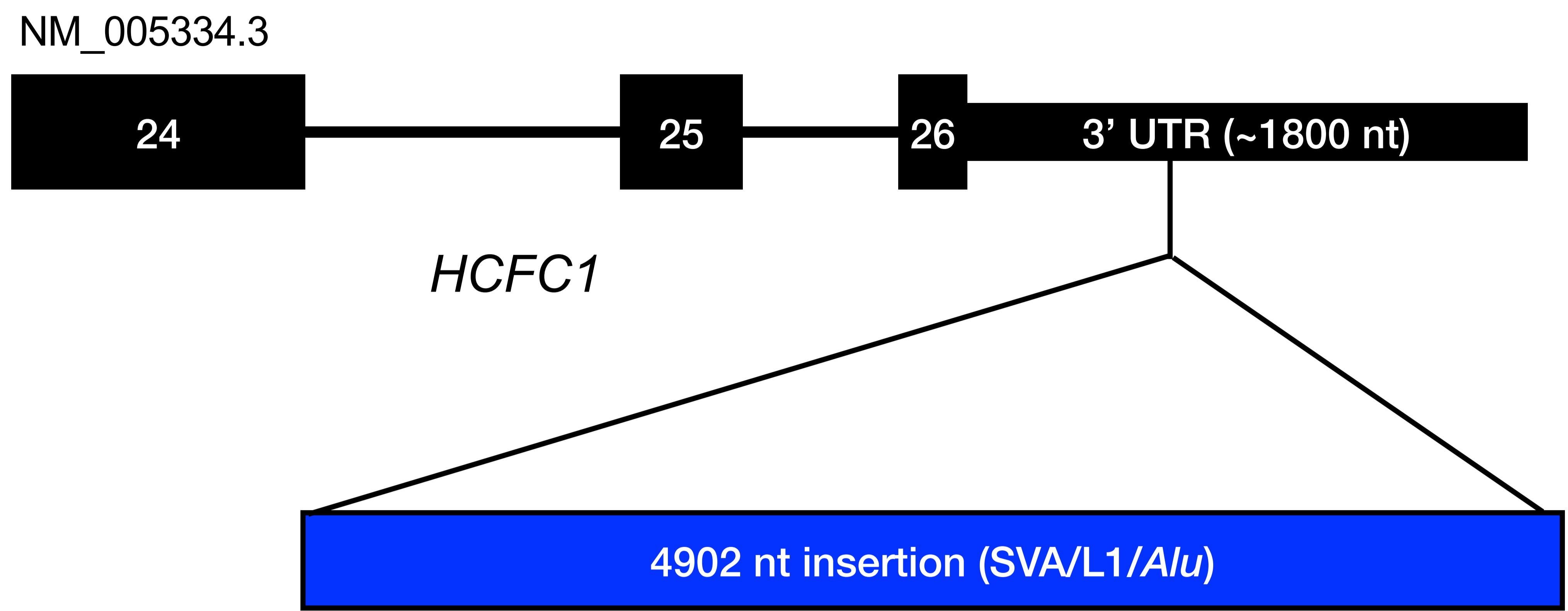


Figure 3

A



B



C

