

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, 35806,
16 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 considerably more
28 accurate 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 disease 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, eight of which
35 (8/96, 8.33%) harbored pathogenic or likely pathogenic variants. Newly identified variants were
36 visible in both srGS and lrGS in nine probands (~9.4%) and resulted from changes to
37 interpretation mostly from recent gene-disease association discoveries. Seven cases included
38 variants that were only interpretable in lrGS, including copy-number variants, an inversion, a
39 mobile element insertion, two low-complexity repeat expansions, and a 1 bp deletion. While
40 evidence for each of these variants is, in retrospect, visible in srGS, they were either: not called
41 within srGS data, were represented by calls with incorrect sizes or structures, or failed quality-
42 control and filtration. Thus, while reanalysis of older data clearly increases diagnostic yield, we
43 find that lrGS allows for substantial additional yield (7/96, 7.3%) beyond srGS. We anticipate that
44 as lrGS analysis improves, and as lrGS datasets grow allowing for better variant frequency
45 annotation, the additional lrGS-only rare disease yield will grow over time.

46

47

48

49 INTRODUCTION

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

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

73 facilitated by data sharing via the MatchMaker Exchange and GeneMatcher (Philippakis et al.
74 2015; Sobreira et al. 2017).

75 Here we discuss findings from lrGS on a cohort of 96 short-read-negative cases, drawn
76 from several studies focused on rare, suspected congenital diseases, especially early-onset
77 neurodevelopmental disorders. We describe 19 relevant or potentially clinically relevant variants
78 not previously evaluated or considered in 16 cases. We show that a combination of more
79 comprehensive variant detection from lrGS and updated reanalysis contribute to these
80 discoveries, supporting the value of a combination of lrGS and reanalysis as a strategy to
81 maximize rates of discovery of highly penetrant variation leading to rare disease.

82

83 **RESULTS**

84 We selected individuals with rare diseases who had undergone short-read exome sequencing
85 (srES; n=2) or srGS (n=94) in previous research studies yet had no pathogenic or likely
86 pathogenic variants (P/LP) nor variants of uncertain significance (VUS) identified (Bowling et al.
87 2017; East et al. 2021; Bowling et al. 2022). Most of our cohort consisted of children (89% were
88 <18 years of age at time of enrollment) with a neurodevelopmental disorder (NDD, 70%), multiple
89 congenital anomalies (MCA, 22%), or a suspected genetic myopathy (8%). Probands consisted
90 of 66% males (63/96); genetically inferred ancestries for probands revealed 72% European
91 (69/96), 21% African/African American (20/96), 3% Admixed American (3/96), 1% Southeast
92 Asian (1/96) and 3% unspecified admixture ancestries (Table 1). For these 96 cases, we
93 performed lrGS using Pacific Biosciences “HiFi” sequencing to a median depth of 27X (Table 1,
94 Supplemental Table S1). For a subset (10/96), we also performed lrGS on parents (median
95 parental HiFi depth of 22X, Supplemental Table S2). We also generated *de novo* assemblies for
96 each proband using hifiasm (Cheng et al. 2021), with parental srGS used for kmer-based binning
97 and phasing when available. The median N50 for all proband contigs was 29.05 Mb (Table 1).

98

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 |

99

100 HiFi reads were aligned to hg38, and variant calling was performed using DeepVariant
 101 (Poplin et al. 2018) and pbsv (<https://github.com/PacificBiosciences/pbsv>, see Methods). A
 102 median of 4.4 million SNVs and 970,031 indels were called in each proband (Table 1).

103

104 **Analysis of SVs**

105 We detected a median of 55,586 SVs of varying classes across the 96 probands using pbsv
 106 (Table 1, Supplemental Table S3). We sought to characterize how filtering SVs by frequency
 107 could reduce manual curation burden by considering five allele frequency resources. First, we
 108 created a set of “cohort” SVs by performing SV call merging across all 96 probands using Jasmine
 109 (Kirsche et al. 2023) and generating allele counts from the merged set (set 1). We then used a
 110 second Jasmine merge step to match cohort SVs with SV frequencies from: an in-house set of

111 266 HiFi genomes including all cohort probands and parents, samples from other internal projects,
112 and public HiFi data (set 2, see Methods); gnomAD v4 SV frequencies from 63,046 short read
113 samples (set 3, Collins et al. 2020) Human Genome Structural Variant Consortium phase 2
114 (HGSVC2) assembly-based calls from 18 HiFi samples (set 4, Ebert et al. 2021); and a PacBio-
115 provided set of pbsv calls from 103 HiFi samples from the Human Pangenome Reference
116 Consortium (HPRC) and Genome in a Bottle (GIAB) consortia (set 5, Ebert et al. 2021;
117 [https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGSVC2/release/v2.0/integrated_ca](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGSVC2/release/v2.0/integrated_callsset/variants_freeze4_sv_insdels.vcf.gz)
118 [llset/variants_freeze4_sv_insdels.vcf.gz](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGSVC2/release/v2.0/integrated_callsset/variants_freeze4_sv_insdels.vcf.gz);
119 [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)
120 [cf.gz](https://github.com/PacificBiosciences/svpack/raw/main/resources/HPRC_GIAB_GRCh38.pbsv.vcf.gz)). We note that these data are at least partially redundant with one another (e.g., set 1 is a
121 partial subset of set 2 and many public samples such as NA12878 are shared between sets 2-5).
122 However, our goal was to maximize filtering performance by including calls from as many datasets
123 and discovery methods as possible. When filtering SVs with this data, we found probands had a
124 median of 1,721 “rare” SVs, defined as having an allele frequency <1% in public SV databases
125 (sets 3-5) and an allele count <4 in the cohort or internal (sets 1 or 2). A subset of these (proband
126 median of 87 SVs) were within 50 bp of a RefSeq exon (Supplemental Table S3). We also filtered
127 to identify “private” SVs as those absent from sets 3-5, having a set 1 allele count of 1, and having
128 a set 2 allele count of ≤ 2 (allowing for parental inheritance). This filtering results in a median of
129 733 SVs per proband, only 40 of which are within 50 bp of a RefSeq exon (Supplemental Table
130 S3).

131

132 **Findings from IrGS**

133 IrGS SNVs/indels were annotated with features such as gene overlaps, coding consequences,
134 computational impact scores, and allele frequencies. They were then filtered and analyzed using
135 in-house software that is also used for srGS data (Hiatt et al. 2021). Rare SVs were assessed by
136 visualization of reads in IGV and prioritization and analysis using SvAnna (Danis et al. 2022). All

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 | IrGS required? | Orthogonal Validation |
|------------|-------------------------------|--|-------------------------|------------------------|---------------------------|--|----------------|----------------|--|
| 1 | <i>ZBTB20</i> | NC_000003.12:g.110477273_114639202_inv | de novo | P | Definitive Diagnostic | 2A(+1.00), 5A (+0.15) | SV | Yes | Yes, Research Sanger of breakpoints |
| 2 | <i>ALS2</i> | NC_000002.12:g.201720435-201725085_del; NC_000002.12:g.200115181-201739349_del | biparental | LP; P | Definitive Diagnostic | 2E (+0.90);2D-4 (+0.90), 3B (+0.45) | SVx2 | Yes; Yes | Pending, supported by srGS data |
| 3 | <i>HCFC1</i> | NC_000023.11:g.153948602_ins4902, NM_005334.3:c.*745_ins4902 | maternal (X-linked) | VUS | Uncertain | NA | SV | Yes | Yes, Research PCR amplification of breakpoints. |
| 4 | <i>ABAT, PMM2, USP7, etc.</i> | NC_000016.10:g.(8742452_9220783)ins[(8742452_8879961)_ (8879962_9000189)x2_(9000190_9220783)] | de novo | VUS | Uncertain | 2K (+0.30) | SV | Yes | 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 | Yes | Yes, Clinical Testing* |
| 6 | <i>AFF3</i> | NM_001386135.1:c.-64-281_-64-280insGGC[90] | unknown | VUS | Uncertain | NA | TRE | Yes | 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 | Yes | 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 | No | Yes, Clinical Sanger |
| 9 | <i>CSNK2B</i> | NM_001320.6:c.202C>T, p.(Gln68Ter) | paternal | LP | Likely Diagnostic | PVS1_VeryStrong, PM2_Moderate | SNV | No | 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 | No | Yes, Clinical Sanger |
| 11 | <i>MCF2</i> | NM_005369.5:c.2234G>T, p.(Gly745Val) | maternal (X-linked) | VUS | Uncertain | PM2 | SNV | No | Yes, Clinical Sanger |
| 12 | <i>NOTCH3</i> | NM_000435.3c.6409_6410delCT, p.(Leu2137GlyfsTer104) | paternal | LP | Likely Diagnostic | PVS1_Strong, PM2_Moderate | SNV | No | Yes, Clinical Sanger |
| 13 | <i>AFF4</i> | NM_014423.4, c.879delA, p.(His294IlefsTer5) | paternal | VUS | Uncertain | PM2 | SNV | No | Yes, Clinical Sanger |
| 14 | <i>KCNT2, KIF21A</i> | NC_000001.11:g.196329420-196344697_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) | No; No | 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 | No | Pending, supported by srGS data |
| 16 | <i>SCN1A</i> | NM_001165963.4:c.4003-603T>C | paternal | VUS | Uncertain | PM2_Moderate | SNV | No | 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. [#]Independent clinical testing indicated that the *PHOX2B* expansion was de novo; we only sequenced the proband in our research study.

138 variants of interest were subject to curation using American College of Medical Genetics and
139 Genomics and Association for Molecular Pathology (ACMG/AMP) and ClinGen criteria to identify
140 potentially clinically relevant variation (Richards et al. 2015; Riggs et al. 2020). We identified 19
141 potentially “clinically relevant” variants, defined here as being pathogenic, likely pathogenic, or
142 variants of uncertain significance (P/LP/VUS), in 16 of the 96 cases (Table 2). Seven of these
143 have a case-level classification of Definitive Diagnostic or Likely Diagnostic, which we define as
144 P/LP variants that likely fully explain the reason for testing (Bowling et al. 2022). The remaining
145 nine cases have Uncertain case-level classifications, either due to the variants being VUSs or
146 being P/LP variants in genes whose associated phenotypes do not closely match the observed
147 phenotype. Findings in seven probands exemplify the unique benefits of IrGS and are highlighted
148 below.

149

150 IrGS-informed SVs

151 IrGS uncovered a *de novo*, 4 Mb, copy-neutral, paracentric inversion on chromosome 3
152 (NC_000003.12:g.110477273_114639202_inv) in Proband 1 (Figure 1). This inversion spans

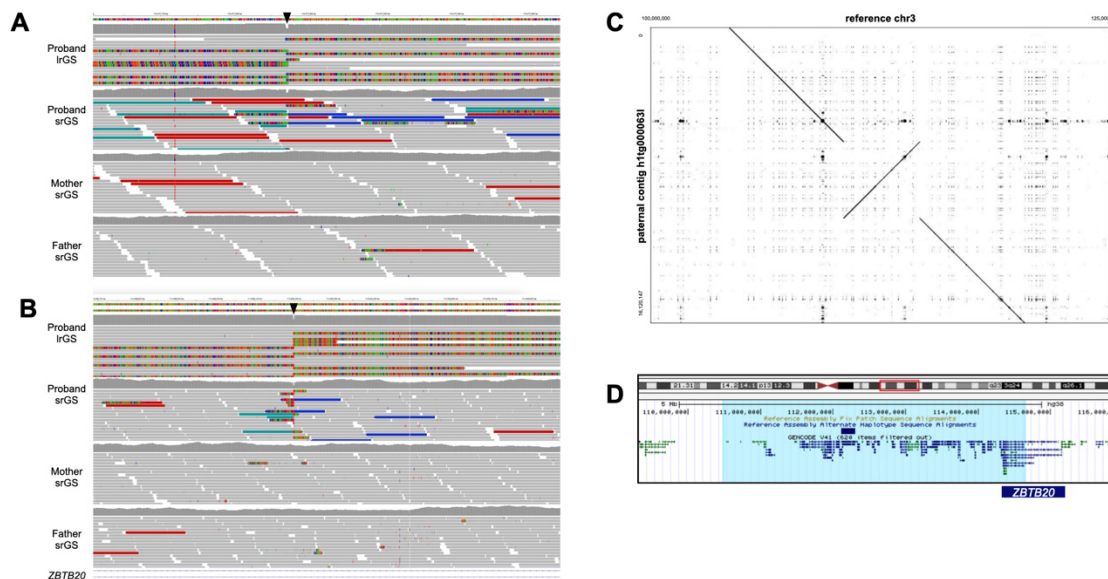


Figure 1. A *de novo*, 4 Mb paracentric inversion in proband 1, affecting *ZBTB20*. A, B. Visualization of a subset of proband and parent reads in IGV at the 5' (A) and 3' (B) breakpoints (black arrowheads) indicate a *de novo* event. C. Alignment of the proband's assembled paternal contig versus the reference genome supports the inversion. D. Visualization of the inverted region (highlighted in light blue) in the UCSC browser shows the inversion spans 35 protein-coding genes and likely disrupts the *ZBTB20* gene (dark blue bar).

153 about 35 protein-coding genes, and one breakpoint of this inversion lies within an intron of
154 *ZBTB20* (MIM: 606025). The event breakpoints are visible in both long and short reads for this
155 proband, but it was only called as an inversion in IrGS. The variant is private to the proband and
156 is predicted to disrupt *ZBTB20*. Loss-of-function (LOF) variation in *ZBTB20* is associated with
157 Primrose Syndrome (MIM: 259050) and 3q13.31 Microdeletion Syndrome (Juven et al. 2020).
158 Proband 1's reported features include moderate intellectual disability (ID), delayed speech and
159 language development, muscular hypotonia, strabismus, and hypoplastic corpus callosum. She
160 is also non-ambulatory. Several of these features overlap Primrose Syndrome. We have classified
161 this variant as pathogenic and the case-level designation is Definitive Diagnostic (Table 2).

162 Proband 2 presented with spasticity, ataxia, and leukodystrophy. srGS was negative, but
163 IrGS identified two structural variants (SVs) identified in *trans* in *ALS2* (MIM: 606352). These
164 include a maternally-inherited 4.65 kb deletion (chr2:201720435-201725085_del) that removes
165 exons 21-23 of NM_020919.4, and a paternally-inherited ~1.6 Mb deletion (chr2:200115181-
166 201739349_del) that spans several genes, including the 3' end of *ALS2* (deletion of exons 12-34
167 of NM_020919.4, Figure 2, Table 2). While these variants are visible in short-read data, the
168 smaller deletion was called as a heterozygous deletion in the mother and as a homozygous
169 deletion in the proband, obfuscating the nature of the variation and raising quality-control
170 concerns. Given the results of the IrGS, the srGS variant calls logically resulted from the small
171 maternal deletion intersecting with the larger, overlapping paternal deletion. This case highlights
172 the difficulties in identification and analysis of overlapping SVs of unknown phase. Variation in
173 *ALS2* is associated with several AR conditions (Juvenile amyotrophic lateral sclerosis 2, MIM:
174 205100; Juvenile Primary lateral sclerosis, MIM: 606353; and infantile onset ascending Spastic
175 paralysis, MIM: 607225), each of which have features that overlap the proband's presentation.

176 These variants are classified as P/LP and, given the degree of overlap with expected phenotypes,
177 the case-level designation is Definitive Diagnostic.

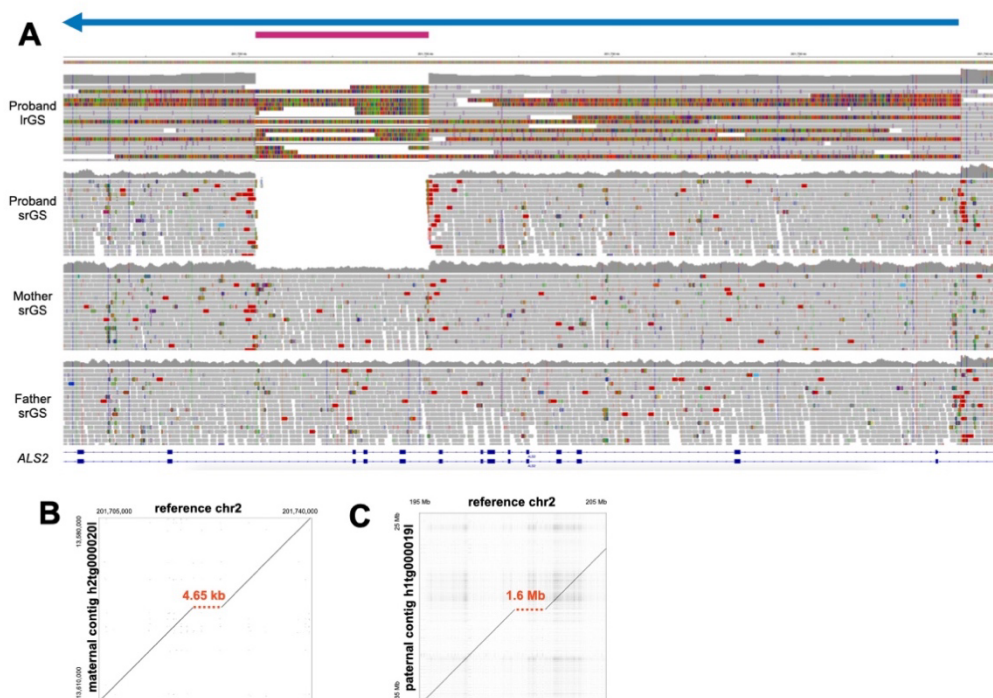
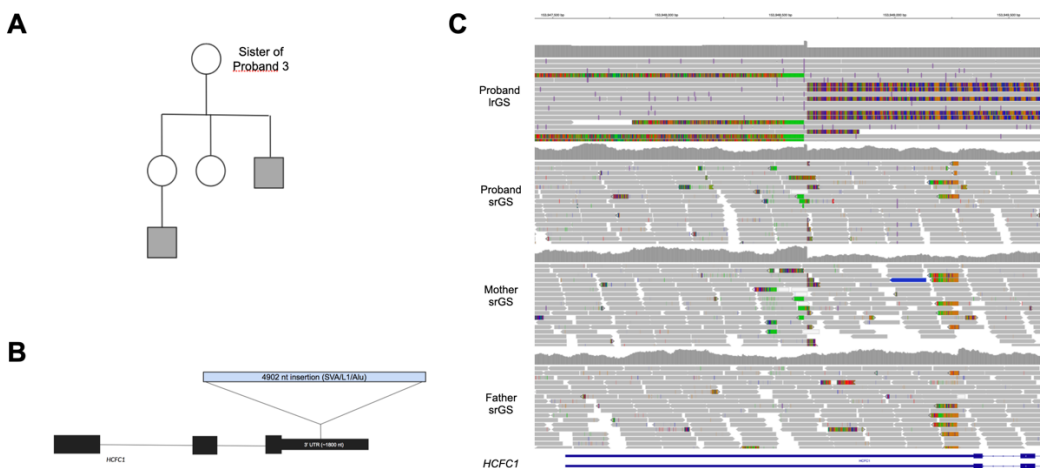


Figure 2. Two ALS2 deletions in trans in Proband 2. A. Visualization of proband and parent reads in IGV indicate two overlapping deletions in ALS2; a smaller maternal deletion (pink bar) and a larger paternal deletion (blue bar/arrow). Alignment of the proband's assembled maternal (B) or paternal contig (C) versus the reference genome support the two deletions (red dashed lines).

178
179 Proband 3 is a male with a strong X-linked family history of ID (Figure 3A). srES and srGS
180 were both negative, although in srES, two neighboring SNVs were called, and manually curated,
181 in the 3' UTR of *HCFC1* (MIM: 300019). While srGS resulted in no calls in this region, visualization
182 of reads in IGV suggested an insertion of unknown length and consequence (Figure 3C). Variation
183 (mostly missense and proposed regulatory variation) in *HCFC1* has been associated with X-linked
184 recessive Methylmalonic aciduria and homocysteinemia, cblX type (MIM: 309541). However,
185 most affected individuals present with severely delayed psychomotor development, seizures, and
186 methylmalonic aciduria. Proband 3's family reported neither of the latter two features. Given the
187 uncertainties in the identity, structure and consequence of the variants in both srES and srGS and
188 the lack of clear phenotypic relevance for the gene, these variant calls were curated but not

189 originally considered to be good candidates for clinical relevance. HiFi sequencing identified a
190 4,902 bp mobile element insertion (MEI) in the 3' UTR of *HCFC1* (chrX:153948602_ins4902),
191 consisting of both SVA and L1 sequence; this variant was likely inherited from a heterozygous
192 carrier mother, as indicated by srGS reads at the breakpoints



193 **Figure 3. Proband 3 has a 4 kb insertion in the 3' UTR of *HCFC1*.** A. The proband's family has a history of X-linked intellectual disability, as the proband (not shown) and two other male relatives (gray squares) are affected. B. Model of the relative length of the insertion in the 3' UTR. C. The insertion is likely inherited from a heterozygous carrier mother, as indicated by srGS reads.

194 (Figure 3C). While this insertion does not affect protein-coding sequence, it is predicted to
195 increase the length of the 3' UTR by 3.75x (from ~1800 nt to about ~6800 nt). We subsequently
196 performed 3'-end RNA-seq on blood from both the proband and his father, generating $\sim 8 \times 10^6$
197 reads for each sample (see Methods). *HCFC1* shows the greatest expression decrease in the
198 proband relative to his father, an ~ 8.8 -fold reduction, among all genes with at least 10 counts in
199 each sample (Supplemental Figure S1). While these results are consistent with the hypothesis
200 that the insertion has a large effect on *HCFC1* expression and potential activity, they are not
201 definitive. Further, expression or segregation analyses in additional family members could not be
202 assessed. Given the uncertainty of the molecular consequence, the differences in phenotypic
203 features, and the lack of additional segregation data, we classified this as a VUS, with a case-
204 level designation of Uncertain.

205 Proband 4 has a complex tandem *de novo* duplication affecting 16p13.2
206 (NC_000016.9:g.(8742452_9220783)ins[(8742452_8879961)_(8879962_9000189)x2_(900019
207 0_9220783)], Table 2, Supplemental Figure S2). Overlapping duplications have been reported in
208 gnomAD but are rare (Lek et al. 2016). One individual in Decipher (Patient: 251349) has also
209 been reported with a very similar duplication of uncertain consequence (Deciphering
210 Developmental Disorders 2015). These duplications span seven genes, three of which are
211 associated with disease: *ABAT*, *PMM2*, and *USP7*. The first (*ABAT*) is intersected by a duplication
212 breakpoint, but all other breakpoints lie within intergenic regions. *USP7* is the only gene
213 associated with autosomal dominant disease (Hao-Fountain Syndrome, MIM:616863), but this
214 gene is expected to have a copy number of four in this proband, while LOF is generally the
215 mechanism associated with disease (Hao et al. 2015). Some general features of Hao-Fountain
216 Syndrome overlap this proband, but it is not a strong phenotypic fit. The proband is reported to
217 have moderate ID, seizures, microcephaly, and facial dysmorphisms. Proband 4 first had a trio
218 srES and no variants were returned, and this duplication was not called. Trio lrGS identified the
219 order, orientation, and copy number of segments of these tandem duplications. We have
220 classified this variant as a VUS, with a case-level designation of Uncertain. Note that this variant
221 is pending orthogonal validation. Proband 4 was one of the two individuals who only previously
222 had srES rather than srGS, and this duplication was not easily visualized in srES data.

223

224 **lrGS-informed SNV/Indels**

225 Improved variant calling in repeat regions is also a benefit of lrGS (Nurk et al. 2022). In addition
226 to analysis of SNVs and SVs in our standard pipeline, we assessed variant calls in 66 tandem
227 repeat expansion (TRE) regions, including both known disease-associated and candidate
228 disease-associated loci (Supplemental Table S4). We intersected TRE regions of interest with
229 pbsv-called variants in each individual and compared these calls to known pathogenic expansion
230 sizes from the literature. We identified several large heterozygous insertions in repeat regions

231 that, while longer than the expected pathogenicity threshold, were predicted after manual curation
232 to be benign based on their sequence content (Nakamura et al. 2020, Supplemental Figure S3A).
233 We also identified two large heterozygous insertions in *RFC1* (MIM: 102579), one of which is
234 benign based on sequence content and one of which is expected to be a pathogenic insertion.
235 However, *RFC1*-associated disease (CANVAS, MIM: 614575) is caused by biallelic expansions,
236 which we did not observe (Supplemental Figure S3B), suggesting this proband is merely a
237 heterozygous carrier.

238 In Proband 5, we observed a *de novo* 18-bp alanine tract expansion in *PHOX2B* (MIM:
239 603851, NM_003924.4:c.741_758dup, p.(Ala255_Ala260dup), Table 2, Supplemental Figure
240 S4), associated with Central Hypoventilation Syndrome, with or without Hirschsprung Disease
241 (MIM: 209880). This disorder was clinically suspected, but variation in *PHOX2B* was missed by
242 initial clinical genetic testing and srGS, which was performed on a PCR+ srGS library. This variant
243 is reported as pathogenic in literature (Amiel et al. 2003), is a strong match to the proband's
244 observed symptoms, and was confirmed independently by additional clinical testing. We have
245 classified this variant as pathogenic, with a case-level classification of Definitive Diagnostic.

246 In proband 6, we identified a 270 bp insertion in *AFF3* (MIM:601464, NM_001386135.1:c.-
247 64-281_-64-280insGGC[90], Table 2, Supplemental Figure S5). While missense variants in *AFF3*
248 have been associated with KINSSHIP Syndrome (MIM:619297), an expansion of a CGG-repeat
249 in the promoter of this gene and subsequent hypermethylation of the promoter has recently been
250 reported to be associated with NDDs (Jadhav et al. 2023). The majority of probands in our cohort
251 (77/96) had both alleles matching hg38 in this region. Among the 19 probands harboring non-
252 reference alleles, proband 6 had a 270 bp insertion, and the remaining 18 probands had insertions
253 ranging from 24-84 bp in length (8-28 triplet repeat units). Similarly, when comparing to the larger,
254 "set 2" database of HiFi genomes (n=266), only 53/266 individuals harbor a non-reference allele;
255 the longest insertions outside of Proband 6 are 105 bp (35 repeats) and 93 bp (31 repeats), each
256 in different individuals, while the median non-reference insertion is 36 bp (12 repeats). Thus, the

257 270 bp insertion (90 repeat units) in Proband 6 is a clear outlier, being nearly three-fold longer
258 than the second longest insertion (Supplemental Table S4, Supplemental Figure S6). Proband 6
259 was sequenced as a neonate and presented with intrauterine growth restriction (IUGR) and
260 hypoplastic left heart (HLH). Jadhav et al. suggest that normal variation ranges to up to ~38 repeat
261 units, with ≥ 61 repeats being a likely pathogenic threshold; our results are consistent with those
262 observations. Nevertheless, there is a need to further replicate and confirm the spectrum of
263 normal and pathogenic variation in *AFF3* repeat lengths. Given this uncertainty and the
264 uncertainty regarding the proband's cognitive development, we have classified this TRE as a
265 VUS, and the case-level classification is Uncertain. Also note that this variant is pending
266 orthogonal validation.

267 A *de novo* *SHANK3* SNV (NM_033517.1:c.3161delT, p.(Leu1054Argfs*10)) was identified
268 by IrGS in proband 7. While two reads in the srGS data support this 1 bp deletion, the variant was
269 not called by our srGS variant calling pipeline (Supplemental Figure S7). LOF variation in
270 *SHANK3* (MIM: 606230) is associated with Phelan-McDermid syndrome (MIM: 606232). Features
271 of this syndrome are consistent with the proband's features, and we classified this variant as
272 pathogenic (case-level Definitive Diagnostic). Coverage of this region in short read data does not
273 indicate a systematic coverage deficiency, as mean coverage within 50 bp of this variant in the
274 srGS data for the cohort is 26.1x (n=96), while it is 15.1x for proband 7. Further, this gene is not
275 present in the list of medically relevant genes that tend to be poorly covered by srGS (Wagner et
276 al. 2022), suggesting it may have resulted from a stochastic loss of alternative allele reads in the
277 srGS data. Nevertheless, IrGS has been shown to provide better overall sensitivity and specificity
278 to SNVs and indels in Genome-In-A-Bottle (GIAB) gold-standard datasets compared to srGS
279 (Logsdon et al. 2020; Hiatt et al. 2021) and thus detection failures to variants such as this *SHANK3*
280 event are more likely in srGS data in general.

281

282

283

284 **Reinterpretation of SNVs**

285 The remaining 9 cases that had relevant variation identified following IrGS illustrate the value of
286 reanalysis and/or reinterpretation (Table 2, Supplemental Case Reports). In three cases
287 (Probands 8-11, *HNRNPU*, *CSNK2B*, *GNB2*, *MCF2*) we identified variation in genes that had
288 additional published support for association of the gene or the variant with disease since the time
289 of the most recent analysis. In another four cases (Probands 12-15, *NOTCH3*, *AFF4*, *KCNT2*,
290 *KIF21A*, *NRXN1*) we identified variation in established disease genes that conflicted with the
291 published data regarding molecular mechanisms or expected mode of inheritance. For these
292 cases, the SNVs were also identified by srGS but they were prioritized more effectively following
293 IrGS, particularly after no new P/LP/VUS SVs of interest were observed. We note that in 7/9
294 cases, variants were identified in an unaffected or mildly affected parent, which was somewhat
295 unexpected in these cases due to suspicion of high penetrance (also see Supplemental Case
296 Reports). Lastly, we identified a variant of interest in *SCN1A* that resulted from a targeted analysis
297 of candidate “poison exon” variants (Proband 16, Felker et al. 2023).

298

299 Discussion

300 There remains a considerably large fraction of disease suspected to have genetic causes that is
301 refractory to genomic testing, a finding that has been repeatedly shown across many clinical and
302 research projects (e.g., Srivastava et al. 2019; Baxter et al. 2022). Several non-mutually exclusive
303 hypotheses exist to explain these observations. Environmental risk factors, such as teratogenic
304 exposure or infectious disease, may be relevant to some phenotypes. Multigenic contributors are
305 also likely to explain at least some cases. For example, a small but appreciable fraction of “double
306 diagnoses”, in which an affected individual is observed to harbor two distinct diseases resulting
307 from highly penetrant variation in two distinct genes, has been observed in clinical genomic
308 studies (e.g. Posey et al. 2017). Notably, such discoveries are typically only made when the
309 variation in both genes is independently amenable to a pathogenicity determination (e.g., would
310 be P/LP regardless of P/LP variation in the other gene), and it is possible if not probable that at
311 least some conditions result from combinations of variants in different genes that are not
312 pathogenic in isolation (Papadimitriou et al. 2019). At the further end of this spectrum is a
313 polygenic accumulation of many risk-factor alleles, which is known to be relevant to many
314 common, complex diseases and which may contribute to some rare conditions, as has been
315 suggested for at least a subset of neurodevelopmental disorders (Niemi et al. 2018).

316 We find it likely that a substantial fraction of unexplained rare disease arises from highly
317 penetrant, monogenic variation that we have not yet been able to precisely identify or confidently
318 interpret. The results from this study are consistent with that hypothesis, with over 15% of
319 probands with previously negative testing being now found to harbor a relevant or potentially
320 relevant genetic variant. Further, these observations are consistent with the general picture of
321 rare disease testing in recent years. With the advent of exome capture and sequencing ~15 years
322 ago (Ng et al. 2010, 2009) and subsequent improvements in cost and efficiency, genome-wide
323 detection of highly penetrant variation has greatly accelerated in recent years (Bamshad et al.
324 2019; Baxter et al. 2022; Hamosh et al. 2022; Boycott et al. 2022). Long-read genome sequencing

325 represents the next phase of that acceleration by facilitating a substantial increase in variant
326 comprehensiveness and accuracy. We note that this improvement derives from both sensitivity
327 and specificity improvements. For example, effectively all the variation we describe in this study
328 as being newly visible within lrGS is, in fact, visible in srGS data, at least at the breakpoint levels.
329 However, being retrospectively visible, once the location and structure of a variant is known to
330 exist, is a much lower bar than the ability to prospectively detect, define, filter, and curate such
331 variation. For example, we describe here a 4.9 kb insertion of SVA and L1 sequence into the 3'
332 UTR of *HCFC1*. In retrospective analysis, the Mobile-Element Locator Tool (MELT, Gardner et
333 al. 2017) detected a 1.2 kb SVA mobile element at this location in srGS. However, this call is
334 incorrect with respect to size and sequence composition and MELT generally produces too many
335 calls in our srGS data to allow for a sustainable level of manual curation (Hiatt et al. 2021). For
336 example, this proband had 1,311 total MELT calls, 203 of which are within 50 bp of a RefSeq
337 exon. Further, this individual call has an unbalanced number of reads supporting the left and right
338 breakpoints for this event (LP=17; RP=1). The low number of reads supporting the right
339 breakpoint and the incorrect length of the insertion appear to be due to the presence of L1
340 sequence at the 3' end of the insertion, compounded with a long polyA sequence (79 nt).

341 Our results are also consistent with other studies of the benefits of lrGS for discovering
342 genetic contributors to disease. In 2021, for example, we showed in a small pilot project that 2 of
343 6 previously srGS-negative probands harbored clinically relevant variation uniquely interpretable
344 by lrGS (Hiatt et al. 2021). Since that time, several other studies have also used lrGS for molecular
345 diagnosis of rare disease. For example, Cohen and colleagues showed increased yield of GS
346 (both lrGS and srGS) in exome-negative cases (Cohen et al. 2022). Unique discoveries for lrGS
347 included detection of novel repeat expansions of *STARD7* and a compound heterozygous
348 SNV/deletion that was easier to detect in lrGS. However, the increase in diagnostic yield (~13%)
349 was mainly from SVs detectable by either lrGS or srGS. Other studies have also shown the
350 diagnostic value of lrGS, especially in small, well-phenotyped cohorts or families (Sakamoto et al.

351 2024; Kilich et al. 2024; Audet et al. 2023; Del Gobbo et al. 2023; Fukuda et al. 2023; Miller et al.
352 2021). Our results are similar to these studies at a high-level. However, some important
353 differences are worth noting. One notable difference is that we have shown the value of IrGS in
354 singletons. While srGS data was available for most of the probands' parents, only 10/96 probands
355 had parent IrGS data. Our filtering strategies were sufficient to allow prioritization of proband
356 variants without parental data; further, once flagged for interest, inheritance data could often be
357 assessed by looking at parental srGS reads. Based on this experience, prioritization of proband
358 sequencing with targeted validation in parents is an effective way to increase diagnostic yield with
359 IrGS while reducing IrGS sequencing needs.

360 We have also provided a direct, systematic comparison of IrGS to contemporaneously
361 analyzed srGS in previously negative cases. Our results thus allow for the separation and
362 description of clinically relevant variants that are “new” by virtue of being truly unique to IrGS
363 versus those that are “new” by virtue of reanalysis, and which could be found via IrGS or srGS-
364 reanalysis. In that context, we note that the benefits of reanalysis of older srGS data remain
365 considerable. In the results described here, the “new” discoveries in 9 of 16 cases in IrGS were
366 called correctly and interpretable within srGS data. These observations reflect the rapid pace of
367 gene discovery in rare disease (Baxter et al. 2022; Boycott et al. 2022; Hamosh et al. 2022), and
368 are consistent with other studies. For example, we previously showed that the probability of a
369 negative srGS dataset harboring a clinically relevant variant increased from 1% within one year
370 of a previous analysis to ~22% if more than three years have passed since a previous analysis
371 (Hiatt et al. 2018). Several other studies found similar results, with many reanalysis findings being
372 due to recent publications of new gene-disease associations (Liu et al. 2019; Schobers et al.
373 2022; Hartley et al. 2023).

374 One possible optimal path to maximizing overall yield in previously srGS-negative
375 individuals is to include reanalysis prior to IrGS. However, this reflects a cost/benefit ratio that
376 depends on the cost of the analysis step in relation to the costs of sequencing. As IrGS costs

377 decrease there may reach a point where the cost of variant analysis alone (which requires both
378 compute resources and manual curation time and thus represents a non-trivial cost) is substantial
379 relative to lrGS costs per se and which might favor a process of simply performing lrGS. Further,
380 lrGS allows analysis of a more complete genomic picture, including both SNVs and larger variant
381 types that are truly not called in srGS data. This can allow more accurate and confident variant
382 prioritization (including SNVs) as evidenced by several of our cases.

383 We note that the results in this study were generated over a period of time with
384 considerable change in lrGS protocols. For example, most probands in this study were sequenced
385 on Sequel IIE machines (n=64) before the Revio (n=32) became available. Given the costs of
386 data production, the 64 Sequel IIE samples were covered at lower-depth (median coverage
387 24.65X) than the 32 Revio-sequenced samples (median coverage 30.09X, Supplemental Figure
388 S8), which may have reduced our sensitivity to variants in some loci in the earlier samples.
389 Further, methylation data were not available in the early period of this study, although methylation
390 calls are now routine and available for the most recent 44 probands, which may also impact
391 diagnostic yield. For example, evaluation of the, at present VUS, *AFF3* expansion (Supplemental
392 Figure S5, S6) would benefit from an assessment of methylation levels at this locus, as
393 hypermethylation is likely associated with disease risk (Jadhav et al. 2023). Additionally, some
394 probands harbor methylation alterations that are clinically relevant even in the absence of a
395 pathogenic genetic variant (Aref-Eshghi et al. 2021).

396 In addition to changes in sequencing, informatic changes have also been considerable
397 over the course of the data generation for this study. While we present and describe a uniform
398 set of variant-calls, assemblies, and annotations (see Methods), analysis of individual samples
399 took place simultaneously with the optimization of variant-calling and annotation pipelines. One
400 particularly relevant change is variant-frequency annotations. While the key strength of lrGS is its
401 ability to see variants that are invisible or poorly detected in srGS, the ability to discriminate
402 genuine highly penetrant variation from the background of benign alleles depends on the ability

403 to annotate and remove alleles that are common in the general population. As the main allele
404 frequency resources are built from srGS data (e.g., Lek et al. 2016), we have limited ability to filter
405 away likely benign alleles among the variants uniquely detected by lrGS. This ability, however,
406 improved as more lrGS datasets were produced over the course of this study (Ebert et al. 2021;
407 Nurk et al. 2022). Projects like the COLORs consortium (<https://colordb.org/>) are likely to improve
408 frequency annotation in the future and continue to improve variant curation efficacy. Accumulating
409 lrGS data from as many samples and studies as possible is critical for the long-term maximization
410 of lrGS benefits.

411 In sum, rare disease genetics continues to be a rapidly advancing field. With data-sharing
412 (Sobreira et al. 2015; Philippakis et al. 2015) and technology improvement (Wenger et al. 2019),
413 the overall diagnostic yield for individuals with rare disease is increasing at a considerable pace
414 each year. In that context, lrGS is a decisive improvement over srGS, providing substantial gains
415 to variant specificity and sensitivity, especially for complex and repeat-associated variants. We
416 anticipate that the degree of improvement will widen over time, as sequencing and analysis
417 pipelines mature and as lrGS datasets grow.

418

419 **METHODS**

420 *Short-read sequencing and variant calling*

421 Probands, their parents, and, when appropriate, affected siblings were enrolled in one of four
422 research studies aimed at identifying genetic causes of rare disease (Bowling et al. 2017; East et
423 al. 2021; Bowling et al. 2022 and Pediatric Genomics (PGEN), unpublished). These studies were
424 monitored by Western IRB and UAB IRB (WIRB 0071, UAB IRB protocols 170303004,
425 300000328, and 130201001). A parent or legal guardian gave consent for the proband, and
426 assent was also obtained from those probands who were capable. In some cases, adult probands
427 who were capable consented to participation in the study. All enrolled individuals consented to
428 publication of de-identified data. Short read exome (srES) or short-read genome sequencing
429 (srGS) was performed as described (Bowling et al. 2022; East et al. 2021; Hiatt et al. 2021;
430 Bowling et al. 2017). Briefly, whole blood genomic DNA was isolated using the QIAasympyphony
431 (Qiagen), and sequencing libraries were constructed by the HudsonAlpha Genomic Services Lab
432 or the Clinical Services Laboratory, LLC, using a standard protocol that generally included PCR
433 amplification (86/96). Genomes were sequenced at an approximate mean depth of 30X, with at
434 least 80% of base positions reaching 20X coverage. Exomes were sequenced to a mean depth
435 of 71X. For short read reanalysis, srES and srGS reads were aligned and small variants and
436 CNVs were called with a uniform pipeline to hg38, and SNVs/Indels and CNVs were curated using
437 an in-house software tool, as previously described (Hiatt et al. 2021). Expected sample
438 relatedness was confirmed with Somalier (v. 0.2.10, (Pedersen et al. 2020) and predicted major
439 genetic ancestries were calculated with peddy (v. 0.4.1, Pedersen and Quinlan 2017).

440

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

442 Long-read sequencing was performed using HiFi (CCS) mode on either a PacBio Sequel II or
443 Revo instrument (Pacific Biosciences of California, Inc.). Libraries were constructed using a
444 SMRTbell Template Prep Kit (V1.0, 2.0 or 3.0) and tightly sized on a SageELF or BluePippin

445 instrument (Sage Science, Beverly, MA, USA). Sequencing was performed using a 2 hour pre-
446 extension with either 24 or 30 hour movie times. The resulting raw data was processed using
447 either the CCS3.4 or CCS4 algorithm, as the latter was released during the course of the study.
448 Comparison of the number of high-quality indel events in a read versus the number of passes
449 confirmed that these algorithms produced comparable results. Probands were sequenced on 2-3
450 Sequel II or one Revio SMRT cell. Top off sequencing was performed if the sequencing did not
451 meet the desired coverage (>20X). This resulted in an average estimated HiFi read depth of 26.1X
452 (range 16.7-41) of raw, unaligned sequence data for probands. For 10 families, parents were also
453 sequenced on 2-3 Sequel II SMRT cells, with an average estimated depth of 21.6x (range 14-28)
454 of raw, unaligned sequence data for parents. Aligned sequencing metrics are shown in
455 Supplemental Tables S1 and S2. HiFi reads were aligned to the were aligned to the hg38 no-alt
456 analysis set
457 (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)
458 using pbmm2 (v. 1.10.0). SNVs and small indels, were called with DeepVariant (v. 1.5.0) and used
459 to haplotag the aligned reads with whatshap (v. 1.7). Structural variants were called using pbsv
460 v2.9.0 (<https://github.com/PacificBiosciences/pbsv>). For the 10 cases with parent lrGS, candidate
461 de novo SVs required a proband genotype of 0/1 and parent genotypes of 0/0, with ≥ 6 alternate
462 reads in the proband and 0 alternate reads, and ≥ 5 reference reads in the parents.

464
465 De novo assemblies were generated for all probands using hifiasm (v. 0.15.2 (r334 or v.0.16.1-
466 r375, Cheng et al. 2021). Hifiasm was used to create two assemblies. First, the default parameters
467 were used, followed by two rounds of Racon (v1.4.10) polishing of contigs. For cases with parent
468 GS data, trio-binned assemblies were built using kmers (srGS). The kmers were generated using
469 yak (v0.1) using the suggested parameters for running a hifiasm trio assembly (kmer size=31 and
470 Bloom filter size of 2^{*37}). Maternal and paternal contigs went through two rounds of Racon

471 (v1.4.10) polishing. Individual parent assemblies were also built with hifiasm (v0.15.2) using
472 default parameters. The resulting contigs went through two rounds of Racon (v1.4.10) polishing.
473
474 Coordinates of breakpoints were defined by a combination of assembly-assembly alignments
475 using minimap2 (Li 2018) (followed by use of bedtools bamToBed), visual inspection of CCS read
476 alignments, and BLAT. Dot plots illustrating sequence differences were created using Gepard
477 (Krumstiek et al. 2007).

478

479 *Structural Variant Merging and Filtering*

480 Overall structural variant counts in Supplemental Table S1 were generated from the proband pbsv
481 VCFs using bcftools (v1.15.1) and awk (e.g., bcftools filter -i 'SVTYPE=="DEL"' <input.vcf> |
482 {a[i++]=\$1; sum+=\$1} END{asort(a); min=a[1]; max=a[i]; if(i%2==1) median=a[int(i/2)+1]; else
483 median=(a[i/2]+a[i/2+1])/2; mean=sum/i; print mean, median, min, max;}).

484

485 An internal allele frequency callset was constructed and periodically updated using all available
486 internally sequenced HiFi pbsv variant calls and pbsv calls generated from public HiFi sequencing
487 data (n=266 at the end of this analysis). Samples included a majority of participants from this
488 cohort and our previous NDD pilot cohort (Hiatt et al. 2021, 120/266); samples from HudsonAlpha
489 non-NDD projects (35/266); HudsonAlpha-sequenced data from HG001, HG003, HG004, HG006,
490 and HG007 (5/266); HG00514, HG00731, HG00732, NA19240 from HGSC2 (Ebert et al. 2021)
491 (4/266); CHM13 (1/266) (Nurk et al. 2022),
492 (https://www.ncbi.nlm.nih.gov/sra/?term=SRX789768*+CHM13) the HPRC Year 1 and
493 HPRC_PLUS Year 2 releases (101/266)(Ebert et al. 2021). Pbsv calls were merged naively using
494 bcftools merge (v. 1.15.1), i.e., merging only variants that were identical in terms of location,
495 reference sequence, and alternate sequence. The internal pbsv allele frequency callset included
496 affected probands as well as parent/child trios.

497
498 A more robustly merged cohort callset (n=96) was constructed using Jasmine, which allows for
499 merging of similar structural variants that may have non-identical representation in terms of
500 genomic position or variant sequence via spanning a structural variant proximity graph. Jasmine
501 was run with options `--centroid_merging --min_overlap=0.65 --min_sequence_id=0.75 --`
502 `output_genotypes` to create the merged callset. This cohort callset was then annotated with five
503 sets of structural variant frequency annotation. Cohort allele frequencies were generated from the
504 merged set with `bcftools +fill-tags` (set 1, n=96). A second Jasmine merge was used to combine
505 the cohort callset allele count with additional allele frequency resources: the HudsonAlpha internal
506 pbsv callset (described above) (set 2, n=266); gnomAD structural variants (Collins et al. 2020,
507 v4.0, <https://gnomad.broadinstitute.org/news/2023-11-v4-structural-variants>) (set 3, n=63,046);
508 HGVC2 structural variants (Ebert et al. 2021)
509 [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_insdels.vcf.gz)
510 [llset/variants_freeze4_sv_insdels.vcf.gz](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGVC2/release/v2.0/integrated_callset/variants_freeze4_sv_insdels.vcf.gz)) (set 4, n=18); pbsv calls for individuals from the
511 Human Pangenome Reference Consortium and Genome in a Bottle (PacBio;
512 [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)
513 [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 -
514 `-output_genotypes` option and records without genotypes were discarded. A custom python script
515 was used to transfer allele frequency annotations from the annotation source VCFs to the merged
516 VCF based on unique variant identifiers (which were created for each allele frequency resource
517 as needed) and the IDLIST record from Jasmine. The annotated callset was split into individual
518 VCFs with `bcftools (bcftools view -I -s <id>)` which were filtered with `bcftools` to generate the
519 counts in Supplemental Table S2. The “rare” filter described in Supplemental Table S3 was
520 defined as excluding variants with any of the following: within-cohort allele count > 3, in-house
521 allele count > 3, gnomAD maximum population allele frequency (POP_MAX_AF) >= 1%,
522 HPRC_GIAB allele frequency >=1%, or HGVC2 allele frequency >= 1%. The “proband-

523 exclusive” filter described in Supplemental Table S3 was defined as excluding variants with any
524 of the following: within-cohort allele count > 1, in-house allele count > 2, gnomAD allele count >
525 0, HPRC_GIAB allele count > 0, or HGSC2 allele count > 0. The in-house allele count cutoff
526 was set at 2 in order to account for the fact that some parental samples were included in the in-
527 house frequency database. Exon regions were defined as RefSeq exons +/- 50bp (calculated
528 from bedtools slop). Variant filtering and counting was performed with bcftools view and
529 command-line tools, e.g., `bcftools view -R refseq_exons_plus50bp.bed.gz -e 'AC>3 |`
530 `inhouse_pbsv_AC >3 | gnomad4_POPMAX_AF >= 0.01 | hgsvc2_AF >= 0.01 |`
531 `hprc_giab_pbsv_AF >= 0.01' -H <input_file> | wc -l` for the rare filter and `bcftools view -R`
532 `refseq_exons_plus50bp.bed.gz -e 'AC>1 | inhouse_pbsv_AC>2 | gnomad4_AC>0 |`
533 `hgsvc2_AC>0 | hprc_giab_pbsv_AC>0' -H <input_file> | wc -l` for the exclusive filter.

534

535 *Structural Variant Annotation and Curation*

536 For individual case structural variant analysis, a frequency-filtered subset of the proband’s pbsv
537 calls was generated using bcftools annotate and bcftools filter, requiring that calls be located
538 within +/- 50bp of a RefSeq exon and have an allele count of < 4 in the HudsonAlpha internal
539 allele frequency dataset (described above). Each call was visualized using a custom pipeline to
540 automatically generate IGV screenshots (Robinson et al. 2011). Additionally, these filtered pbsv
541 variants were annotated, prioritized, and visualized with SvAnna (v1.0.4, annotations v.2204 or
542 v.230, Danis et al. 2022) based on manually curated HPO terms for each case.

543

544 *Variant interpretation and orthogonal confirmation*

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

549

550 *Sequencing Metrics*

551 Sequencing metrics were generated from the aligned BAMs using the Sentieon implementations
552 of Picard sequencing metrics (Kendig et al. 2019). Sentieon algorithms QualityYield, GCBias,
553 AlignmentStat were run with default settings. Sentieon WgsMetricsAlgo was run with settings --
554 include_unpaired true true --min_map_qual 0 --min_base_qual 0 --coverage_cap 5000. Further
555 coverage metrics were generated with cramino using the --phased option. Supplemental Table
556 S5 maps each tool or command and output field name to the corresponding entry in Supplemental
557 Tables S1 and S2. SNV and indel counts and ratios were calculated with rtg vcfstats
558 (<https://github.com/RealTimeGenomics/rtg-tools>). Summary statistics and graphs were calculated
559 with R (v4.3.1), RStudio (v2023.9.1 build 494), and ggplot2 (v3.4.3). Coverage across regions of
560 interest, such as *SHANK3* was calculated using samtools bedcov with default parameters.

561

562 *Repeat region detection and analysis*

563 We curated a BED file of disease-associated low-complexity repeat regions in 66 genes from
564 previous studies (Hiatt et al. 2021; Cohen et al. 2022 and references therein). Variant calls from
565 pbsv were extracted from these regions +/- 30bp (Supplemental Table S4). Reads were also
566 visualized using the Integrated Genomics Viewer (IGV). Coverage across low complexity repeat
567 regions was calculated using samtools bedcov with default parameters. A coverage of at least 8x
568 across the low-complexity region was required for inclusion in Supplemental Table S4 and
569 Supplemental Figure S6. We also used TRGT (Dolzhenko et al. 2024) and companion tool TRVZ
570 for visualization of a subset of calls including those for display in Figure S3. TRGT was fed an
571 hg38 reference genome FASTA, BED catalog of tandem repeats, and a sample's BAM to
572 generate a VCF containing genotypes for each tandem repeat from the provided catalog in the
573 given sample and a BAM containing only reads that span the repeat sequences. Output VCFs

574 and BAMs were sorted, indexed, and fed into TRVZ with the same hg38 reference FASTA and
575 BED catalog of tandem repeats to generate pileup plots for any desired variants.

576

577 *3' mRNA-seq*

578 Total RNA was isolated from blood samples in PAXgene RNA tubes (PreAnalytiX #762165)
579 according to the manufacturer's instructions and stored short-term at -20°C. RNA was isolated
580 using the PAX gene Blood RNA Kit (Qiagen #762164) according to the manufacturer's
581 instructions. Isolated RNA was quantified by the Qubit RNA HS Assay Kit (Thermo Q32855). 425
582 ng of RNA was used as input for the QuantSeq 3' mRNA-Seq Library Prep Kit FWD for Illumina
583 and UMI Second Strand Synthesis Module for QuantSeq FWD (Illumina, Read 1) from Lexogen
584 (015.96 and 081.96, respectively). Libraries were quantified using the Qubit DNA HS Assay Kit
585 (Thermo Q32854) and visualized with the BioAnalyzer High Sensitivity DNA Analysis kit (Agilent
586 5067–4626) and 2100 BioAnalyzer Instrument (Agilent). Sequencing was carried out using
587 Illumina NextSeq 75 bp single-end. UMIs were first extracted from the reads with UMI-tools extract
588 with regex. Reads were then trimmed with bbdduk and aligned to hg38-GENCODEv42 using STAR
589 with the Lexogen recommended parameters for QuantSeq. Bams were deduplicated by UMI and
590 mapping coordinates using UMI-tools dedup. Count tables were generated using htseq-count with
591 the intersection-nonempty method.

592

593 **DATA ACCESS**

594 For participants who consented to controlled-access sharing, the IrGS data generated in this study
595 will be submitted to dbGAP and/or AnVIL under accession number phs003537.v1
596 (https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs003537.v1).

597 srGS data for participants who consented to controlled-access sharing in NIH-funded studies are
598 available via dbGAP and/or AnVIL (CSER1: <https://www.ncbi.nlm.nih.gov/projects/gap/cgi->

599 bin/study.cgi?study_id=phs001089.v3.p1, dbGaP accession phs001089; SouthSeq:
600 <https://anvilproject.org/data/studies/phs002307>, dbGaP accession phs002307).

601

602

603 **COMPETING INTEREST STATEMENT**

604 The authors declare no competing interests.

605

606 **ACKNOWLEDGMENTS**

607 We thank the families who have contributed to our studies and our collaborating physicians and
608 clinical staff for recruitment and enrollment for these research studies. Some reagents were
609 provided by PacBio as part of an early-access testing program. The CSER1 project was
610 supported by a grant from the US National Human Genome Research Institute (NHGRI;
611 UM1HG007301). The SouthSeq project (U01HG007301) was supported by the Clinical
612 Sequencing Evidence-Generating Research (CSER2) consortium, which is funded by the
613 National Human Genome Research Institute with co-funding from the National Institute on
614 Minority Health and Health Disparities and the National Cancer Institute. The Alabama Genomic
615 Health Initiative is an Alabama-State earmarked project (F170303004) through the University of
616 Alabama in Birmingham. The PGEN cohort was funded by the Alabama Pediatric Genomics
617 Initiative. IrGS of some samples was supported by a Research Grant from the Muscular Dystrophy
618 Association (MDA 963255).

619

620

621 REFERENCES

- 622 Amiel J, Laudier B, Attié-Bitach T, Trang H, De Pontual L, Gener B, Trochet D, Etchevers H,
623 Ray P, Simonneau M, et al. 2003. Polyalanine expansion and frameshift mutations of the
624 paired-like homeobox gene PHOX2B in congenital central hypoventilation syndrome. *Nat*
625 *Genet* **33**: 459–461. <https://pubmed.ncbi.nlm.nih.gov/12640453/> (Accessed February 5,
626 2024).
- 627 Aref-Eshghi E, Kerkhof J, Pedro VP, France G DI, Barat-Houari M, Ruiz-Pallares N, Andrau JC,
628 Lacombe D, Van-Gils J, Fergelot P, et al. 2021. Evaluation of DNA Methylation
629 Episignatures for Diagnosis and Phenotype Correlations in 42 Mendelian
630 Neurodevelopmental Disorders. *Am J Hum Genet* **108**: 1161–1163.
631 <https://pubmed.ncbi.nlm.nih.gov/34087165/> (Accessed February 27, 2024).
- 632 Audet S, Triassi V, Gelinat M, Legault-Cadieux N, Ferraro V, Duquette A, Tetreault M. 2023.
633 Integration of multi-omics technologies for molecular diagnosis in ataxia patients. *Front*
634 *Genet* **14**: 1304711. <http://www.ncbi.nlm.nih.gov/pubmed/38239855> (Accessed February
635 13, 2024).
- 636 Bamshad MJ, Nickerson DA, Chong JX. 2019. Mendelian Gene Discovery: Fast and Furious
637 with No End in Sight. *Am J Hum Genet* **105**: 448–455.
638 <https://pubmed.ncbi.nlm.nih.gov/31491408/> (Accessed July 19, 2022).
- 639 Baxter SM, Posey JE, Lake NJ, Sobreira N, Chong JX, Buyske S, Blue EE, Chadwick LH,
640 Coban-Akdemir ZH, Doheny KF, et al. 2022. Centers for Mendelian Genomics: A decade
641 of facilitating gene discovery. *Genet Med* **24**: 784–797.
642 <https://pubmed.ncbi.nlm.nih.gov/35148959/> (Accessed April 16, 2023).
- 643 Bowling KM, Thompson ML, Amaral MD, Finnila CR, Hiatt SM, Engel KL, Cochran JN, Brothers
644 KB, East KM, Gray DE, et al. 2017. Genomic diagnosis for children with intellectual
645 disability and/or developmental delay. *Genome Med* **9**: 43.
646 <http://www.ncbi.nlm.nih.gov/pubmed/28554332>.
- 647 Bowling KM, Thompson ML, Finnila CR, Hiatt SM, Latner DR, Amaral MD, Lawlor JM, East
648 KM, Cochran ME, Greve V, et al. 2022. Genome sequencing as a first-line diagnostic test
649 for hospitalized infants. *Genet Med* **24**: 851–861.
650 <https://pubmed.ncbi.nlm.nih.gov/34930662/> (Accessed May 29, 2023).
- 651 Boycott KM, Azzariti DR, Hamosh A, Rehm HL. 2022. Seven years since the launch of the
652 Matchmaker Exchange: The evolution of genomic matchmaking. *Hum Mutat* **43**: 659–667.
653 <https://pubmed.ncbi.nlm.nih.gov/35537081/> (Accessed July 19, 2022).
- 654 Cheng H, Concepcion GT, Feng X, Zhang H, Li H. 2021. Haplotype-resolved de novo assembly
655 using phased assembly graphs with hifiasm. *Nat Methods* **18**: 170–175.
656 <https://pubmed.ncbi.nlm.nih.gov/33526886/> (Accessed February 6, 2024).
- 657 Cohen ASA, Farrow EG, Abdelmoity AT, Alaimo JT, Amudhavalli SM, Anderson JT, Bansal L,
658 Bartik L, Baybayan P, Belden B, et al. 2022. Genomic answers for children: Dynamic
659 analyses of >1000 pediatric rare disease genomes. *Genetics in Medicine* **24**: 1336–1348.
660 <https://pubmed.ncbi.nlm.nih.gov/35305867/> (Accessed August 24, 2022).
- 661 Collins RL, Brand H, Karczewski KJ, Zhao X, Alföldi J, Francioli LC, Khera A V., Lowther C,
662 Gauthier LD, Wang H, et al. 2020. A structural variation reference for medical and

- 663 population genetics. *Nature* **581**: 444–451. <https://pubmed.ncbi.nlm.nih.gov/32461652/>
664 (Accessed February 15, 2024).
- 665 Danis D, Jacobsen JOB, Balachandran P, Zhu Q, Yilmaz F, Reese J, Haimel M, Lyon GJ,
666 Helbig I, Mungall CJ, et al. 2022. SvAnna: efficient and accurate pathogenicity prediction of
667 coding and regulatory structural variants in long-read genome sequencing. *Genome Med*
668 **14**. <https://pubmed.ncbi.nlm.nih.gov/35484572/> (Accessed February 12, 2024).
- 669 Deciphering Developmental Disorders S. 2015. Large-scale discovery of novel genetic causes
670 of developmental disorders. *Nature* **519**: 223–228.
671 <https://www.ncbi.nlm.nih.gov/pubmed/25533962>.
- 672 Del Gobbo GF, Wang X, Couse M, Mackay L, Goldsmith C, Marshall AE, Liang Y, Lambert C,
673 Zhang S, Dhillon H, et al. 2023. Long-read genome sequencing reveals a novel intronic
674 retroelement insertion in NR5A1 associated with 46,XY differences of sexual development.
675 *Am J Med Genet A*. <http://www.ncbi.nlm.nih.gov/pubmed/38131126> (Accessed February
676 13, 2024).
- 677 Dolzhenko E, English A, Dashnow H, De Sena Brandine G, Mokveld T, Rowell WJ, Karniski C,
678 Kronenberg Z, Danzi MC, Cheung WA, et al. 2024. Characterization and visualization of
679 tandem repeats at genome scale. *Nature Biotechnology* 2024 1–9.
680 <https://www.nature.com/articles/s41587-023-02057-3> (Accessed February 15, 2024).
- 681 East KM, Kelley W V., Cannon A, Cochran ME, Moss IP, May T, Nakano-Okuno M, Sodeke SO,
682 Edberg JC, Cimino JJ, et al. 2021. A state-based approach to genomics for rare disease
683 and population screening. *Genet Med* **23**: 777–781.
684 <https://pubmed.ncbi.nlm.nih.gov/33244164/> (Accessed July 26, 2022).
- 685 Ebert P, Audano PA, Zhu Q, Rodriguez-Martin B, Porubsky D, Bonder MJ, Sulovari A, Ebler J,
686 Zhou W, Mari RS, et al. 2021. Haplotype-resolved diverse human genomes and integrated
687 analysis of structural variation. *Science (1979)* **372**.
- 688 Felker SA, Lawlor JMJ, Hiatt SM, Thompson ML, Latner DR, Finnila CR, Bowling KM,
689 Bonnstetter ZT, Bonini KE, Kelly NR, et al. 2023. Poison exon annotations improve the
690 yield of clinically relevant variants in genomic diagnostic testing. *Genet Med* 100884.
691 <https://pubmed.ncbi.nlm.nih.gov/37161864/> (Accessed May 29, 2023).
- 692 Fukuda H, Mizuguchi T, Doi H, Kameyama S, Kunii M, Joki H, Takahashi T, Komiya H, Sasaki
693 M, Miyaji Y, et al. 2023. Long-read sequencing revealing intragenic deletions in exome-
694 negative spastic paraplegias. *J Hum Genet* **68**: 689–697.
- 695 Gardner EJ, Lam VK, Harris DN, Chuang NT, Scott EC, Stephen Pittard W, Mills RE, Devine
696 SE. 2017. The mobile element locator tool (MELT): Population-scale mobile element
697 discovery and biology. *Genome Res* **27**: 1916–1929.
- 698 Hamosh A, Wohler E, Martin R, Griffith S, Rodrigues E da S, Antonescu C, Doheny KF, Valle D,
699 Sobreira N. 2022. The impact of GeneMatcher on international data sharing and
700 collaboration. *Hum Mutat* **43**: 668–673. <https://pubmed.ncbi.nlm.nih.gov/35170833/>
701 (Accessed July 19, 2022).
- 702 Hao YH, Fountain MD, Fon Tacer K, Xia F, Bi W, Kang SHL, Patel A, Rosenfeld JA, Le Caignec
703 C, Isidor B, et al. 2015. USP7 Acts as a Molecular Rheostat to Promote WASH-Dependent
704 Endosomal Protein Recycling and Is Mutated in a Human Neurodevelopmental Disorder.
705 *Mol Cell* **59**: 956–969. <https://pubmed.ncbi.nlm.nih.gov/26365382/> (Accessed February 11,
706 2024).

- 707 Hartley T, Soubry É, Acker M, Osmond M, Couse M, Gillespie MK, Ito Y, Marshall AE, Lemire
708 G, Huang L, et al. 2023. Bridging clinical care and research in Ontario, Canada:
709 Maximizing diagnoses from reanalysis of clinical exome sequencing data. *Clin Genet* **103**:
710 288–300.
- 711 Hiatt SM, Amaral MD, Bowling KM, Finnila CR, Thompson ML, Gray DE, Lawlor JM, Cochran
712 JN, Bebin EM, Brothers KB, et al. 2018. Systematic reanalysis of genomic data improves
713 quality of variant interpretation. *Clin Genet* **94**.
- 714 Hiatt SM, Lawlor JM, Handley LH, Ramaker RC, Rogers BB, Partridge EC, Boston LB,
715 Williams M, Plott CB, Jenkins J, et al. 2021. Long-read genome sequencing for the
716 molecular diagnosis of neurodevelopmental disorders. *HGG Adv* **2**: 100023.
717 <http://www.ncbi.nlm.nih.gov/pubmed/33937879> (Accessed June 9, 2021).
- 718 Jadhav B, Garg P, van Vugt JJ, Ibanez K, Gagliardi D, Lee W, Shadrina M, Mokveld T,
719 Dolzhenko E, Martin-Trujillo A, et al. 2023. A phenome-wide association study of
720 methylated GC-rich repeats identifies a GCC repeat expansion in *AFF3* as a significant
721 cause of intellectual disability. *medRxiv*. <https://pubmed.ncbi.nlm.nih.gov/37205357/>
722 (Accessed February 6, 2024).
- 723 Juven A, Nambot S, Piton A, Jean-Marçais N, Masurel A, Callier P, Marle N, Mosca-Boidron AL,
724 Kuentz P, Philippe C, et al. 2020. Primrose syndrome: a phenotypic comparison of patients
725 with a *ZBTB20* missense variant versus a 3q13.31 microdeletion including *ZBTB20*. *Eur J*
726 *Hum Genet* **28**: 1044–1055. <https://pubmed.ncbi.nlm.nih.gov/32071410/> (Accessed
727 February 15, 2024).
- 728 Kendig KI, Baheti S, Bockol MA, Drucker TM, Hart SN, Heldenbrand JR, Hernaez M, Hudson
729 ME, Kalmbach MT, Klee EW, et al. 2019. Sentieon DNaseq Variant Calling Workflow
730 Demonstrates Strong Computational Performance and Accuracy. *Front Genet* **10**.
731 <https://pubmed.ncbi.nlm.nih.gov/31481971/> (Accessed February 15, 2024).
- 732 Kilich G, Hassey K, Behrens EM, Falk M, Vanderver A, Rader DJ, Cahill PJ, Raper A, Zhang Z,
733 Westerfer D, et al. 2024. Kagami Ogata syndrome: a small deletion refines critical region
734 for imprinting. *NPJ Genom Med* **9**: 5. <http://www.ncbi.nlm.nih.gov/pubmed/38212313>
735 (Accessed February 13, 2024).
- 736 Kirsche M, Prabhu G, Sherman R, Ni B, Battle A, Aganezov S, Schatz MC. 2023. Jasmine and
737 Iris: population-scale structural variant comparison and analysis. *Nat Methods* **20**: 408–
738 417. <https://pubmed.ncbi.nlm.nih.gov/36658279/> (Accessed February 15, 2024).
- 739 Krumsiek J, Arnold R, Rattei T. 2007. Gepard: a rapid and sensitive tool for creating dotplots on
740 genome scale. **23**: 1026–1028. <http://mips.gsf.de/services/analysis/gepard> (Accessed May
741 18, 2020).
- 742 Lek M, Karczewski KJ, Minikel E V, Samocha KE, Banks E, Fennell T, O'Donnell-Luria AH,
743 Ware JS, Hill AJ, Cummings BB, et al. 2016. Analysis of protein-coding genetic variation in
744 60,706 humans. *Nature* **536**: 285–291. <https://www.ncbi.nlm.nih.gov/pubmed/27535533>.
- 745 Li H. 2018. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094–
746 3100. <https://github.com/ruanjue/smarddenovo>; (Accessed September 7, 2020).
- 747 Liu P, Meng L, Normand EA, Xia F, Song X, Ghazi A, Rosenfeld J, Magoulas PL, Braxton A,
748 Ward P, et al. 2019. Reanalysis of Clinical Exome Sequencing Data. *New England Journal*
749 *of Medicine* **380**: 2478–2480.

- 750 Logsdon GA, Vollger MR, Eichler EE. 2020. Long-read human genome sequencing and its
751 applications. *Nat Rev Genet* **21**: 597–614. [https://www.nature.com/articles/s41576-020-](https://www.nature.com/articles/s41576-020-0236-x)
752 [0236-x](https://www.nature.com/articles/s41576-020-0236-x) (Accessed November 4, 2020).
- 753 Mahmoud M, Huang Y, Garimella K, Audano PA, Wan W, Prasad N, Handsaker RE, Hall S,
754 Pionzio A, Schatz MC, et al. 2024. Utility of long-read sequencing for All of Us. *Nat*
755 *Commun* **15**: 837. <https://pubmed.ncbi.nlm.nih.gov/38281971/> (Accessed February 5,
756 2024).
- 757 Miller DE, Sulovari A, Wang T, Loucks H, Hoekzema K, Munson KM, Lewis AP, Fuerte EPA,
758 Paschal CR, Walsh T, et al. 2021. Targeted long-read sequencing identifies missing
759 disease-causing variation. *Am J Hum Genet* **108**: 1436–1449.
760 <https://pubmed.ncbi.nlm.nih.gov/34216551/> (Accessed February 13, 2024).
- 761 Nakamura H, Doi H, Mitsuhashi S, Miyatake S, Katoh K, Frith MC, Asano T, Kudo Y, Ikeda T,
762 Kubota S, et al. 2020. Long-read sequencing identifies the pathogenic nucleotide repeat
763 expansion in RFC1 in a Japanese case of CANVAS. *J Hum Genet* **65**: 475–480.
764 <https://www.nature.com/articles/s10038-020-0733-y> (Accessed June 27, 2021).
- 765 Ng SB, Buckingham KJ, Lee C, Bigham AW, Tabor HK, Dent KM, Huff CD, Shannon PT, Jabs
766 EW, Nickerson DA, et al. 2010. Exome sequencing identifies the cause of a mendelian
767 disorder. *Nat Genet* **42**: 30–35. <https://pubmed.ncbi.nlm.nih.gov/19915526/> (Accessed
768 February 27, 2024).
- 769 Ng SB, Turner EH, Robertson PD, Flygare SD, Bigham AW, Lee C, Shaffer T, Wong M,
770 Bhattacharjee A, Eichler EE, et al. 2009. Targeted capture and massively parallel
771 sequencing of 12 human exomes. *Nature* **461**: 272–276.
772 <https://pubmed.ncbi.nlm.nih.gov/19684571/> (Accessed February 13, 2024).
- 773 Niemi MEK, Martin HC, Rice DL, Gallone G, Gordon S, Kelemen M, McAloney K, McRae J,
774 Radford EJ, Yu S, et al. 2018. Common genetic variants contribute to risk of rare severe
775 neurodevelopmental disorders. *Nature* **562**: 268–271.
776 <https://pubmed.ncbi.nlm.nih.gov/30258228/> (Accessed February 11, 2024).
- 777 Nurk S, Koren S, Rhie A, Rautiainen M, Bzikadze A V., Mikheenko A, Vollger MR, Altemose N,
778 Uralsky L, Gershman A, et al. 2022. The complete sequence of a human genome. *Science*
779 **376**: 44–53. <https://pubmed.ncbi.nlm.nih.gov/35357919/> (Accessed January 28, 2024).
- 780 Papadimitriou S, Gazzo A, Versbraegen N, Nachtegael C, Aerts J, Moreau Y, Van Dooren S,
781 Nowé A, Smits G, Lenaerts T. 2019. Predicting disease-causing variant combinations. *Proc*
782 *Natl Acad Sci U S A* **116**: 11878–11887. <https://pubmed.ncbi.nlm.nih.gov/31127050/>
783 (Accessed February 14, 2024).
- 784 Pedersen BS, Bhetariya PJ, Brown J, Kravitz SN, Marth G, Jensen RL, Bronner MP, Underhill
785 HR, Quinlan AR. 2020. Somalier: rapid relatedness estimation for cancer and germline
786 studies using efficient genome sketches. *Genome Med* **12**.
787 <https://pubmed.ncbi.nlm.nih.gov/32664994/> (Accessed February 15, 2024).
- 788 Pedersen BS, Quinlan AR. 2017. Who’s Who? Detecting and Resolving Sample Anomalies in
789 Human DNA Sequencing Studies with Peddy. *Am J Hum Genet* **100**: 406–413.
790 <https://pubmed.ncbi.nlm.nih.gov/28190455/> (Accessed February 15, 2024).
- 791 Philippakis AA, Azzariti DR, Beltran S, Brookes AJ, Brownstein CA, Brudno M, Brunner HG,
792 Buske OJ, Carey K, Doll C, et al. 2015. The Matchmaker Exchange: A Platform for Rare

- 793 Disease Gene Discovery. *Hum Mutat* **36**: 915–921.
794 <http://doi.wiley.com/10.1002/humu.22858> (Accessed November 25, 2018).
- 795 Poplin R, Chang PC, Alexander D, Schwartz S, Colthurst T, Ku A, Newburger D, Dijamco J,
796 Nguyen N, Afshar PT, et al. 2018. A universal snp and small-indel variant caller using deep
797 neural networks. *Nat Biotechnol* **36**: 983.
- 798 Posey JE, Harel T, Liu P, Rosenfeld JA, James RA, Coban Akdemir ZH, Walkiewicz M, Bi W,
799 Xiao R, Ding Y, et al. 2017. Resolution of Disease Phenotypes Resulting from Multilocus
800 Genomic Variation. *N Engl J Med* **376**: 21–31. <https://pubmed.ncbi.nlm.nih.gov/27959697/>
801 (Accessed February 15, 2024).
- 802 Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, Grody WW, Hegde M, Lyon E,
803 Spector E, et al. 2015. Standards and guidelines for the interpretation of sequence
804 variants: a joint consensus recommendation of the American College of Medical Genetics
805 and Genomics and the Association for Molecular Pathology. *Genet Med* **17**: 405–424.
806 <https://www.ncbi.nlm.nih.gov/pubmed/25741868>.
- 807 Riggs ER, Andersen EF, Cherry AM, Kantarci S, Kearney H, Patel A, Raca G, Ritter DI, South
808 ST, Thorland EC, et al. 2020. Technical standards for the interpretation and reporting of
809 constitutional copy-number variants: a joint consensus recommendation of the American
810 College of Medical Genetics and Genomics (ACMG) and the Clinical Genome Resource
811 (ClinGen). *Genetics in Medicine* **22**: 245–257. <https://pubmed.ncbi.nlm.nih.gov/31690835/>
812 (Accessed September 8, 2020).
- 813 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP. 2011.
814 Integrative genomics viewer. *Nat Biotechnol* **29**: 24–26.
815 <https://pubmed.ncbi.nlm.nih.gov/21221095/> (Accessed February 15, 2024).
- 816 Sakamoto M, Kurosawa K, Tanoue K, Iwama K, Ishida F, Watanabe Y, Okamoto N, Tsuchida
817 N, Uchiyama Y, Koshimizu E, et al. 2024. A heterozygous germline deletion within USP8
818 causes severe neurodevelopmental delay with multiorgan abnormalities. *J Hum Genet* **69**:
819 85–90. <http://www.ncbi.nlm.nih.gov/pubmed/38030753> (Accessed February 13, 2024).
- 820 Sanghvi R V., Buhay CJ, Powell BC, Tsai EA, Dorschner MO, Hong CS, Lebo MS, Sasson A,
821 Hanna DS, McGee S, et al. 2018. Characterizing reduced coverage regions through
822 comparison of exome and genome sequencing data across 10 centers. *Genet Med* **20**:
823 855–866. <https://pubmed.ncbi.nlm.nih.gov/29144510/> (Accessed February 5, 2024).
- 824 Schobers G, Schieving JH, Yntema HG, Pennings M, Pfundt R, Derks R, Hofste T, de Wijs I,
825 Wieskamp N, van den Heuvel S, et al. 2022. Reanalysis of exome negative patients with
826 rare disease: a pragmatic workflow for diagnostic applications. *Genome Med* **14**.
- 827 Sobreira N, Schiettecatte F, Valle D, Hamosh A. 2015. GeneMatcher: a matching tool for
828 connecting investigators with an interest in the same gene. *Hum Mutat* **36**: 928–30.
829 <http://doi.wiley.com/10.1002/humu.22844> (Accessed November 25, 2018).
- 830 Sobreira NLM, Arachchi H, Buske OJ, Chong JX, Hutton B, Foreman J, Schiettecatte F, Groza
831 T, Jacobsen JOB, Haendel MA, et al. 2017. Matchmaker Exchange. *Curr Protoc Hum*
832 *Genet* **95**: 9.31.1-9.31.15. <http://doi.wiley.com/10.1002/cphg.50> (Accessed November 25,
833 2018).
- 834 Srivastava S, Love-Nichols JA, Dies KA, Ledbetter DH, Martin CL, Chung WK, Firth H V.,
835 Frazier T, Hansen RL, Prock L, et al. 2019. Meta-analysis and multidisciplinary consensus
836 statement: exome sequencing is a first-tier clinical diagnostic test for individuals with

837 neurodevelopmental disorders. *Genet Med* **21**: 2413–2421.
838 <https://pubmed.ncbi.nlm.nih.gov/31182824/> (Accessed December 6, 2022).
839 Wagner J, Olson ND, Harris L, McDaniel J, Cheng H, Fungtammasan A, Hwang YC, Gupta R,
840 Wenger AM, Rowell WJ, et al. 2022. Curated variation benchmarks for challenging
841 medically relevant autosomal genes. *Nat Biotechnol* **40**: 672–680.
842 <https://pubmed.ncbi.nlm.nih.gov/35132260/> (Accessed February 15, 2024).
843 Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, Ebler J,
844 Fungtammasan A, Kolesnikov A, Olson ND, et al. 2019. Accurate circular consensus long-
845 read sequencing improves variant detection and assembly of a human genome. *Nat*
846 *Biotechnol* **37**: 1155–1162.
847