

1 Contribution of *de novo* retroelements to birth defects and childhood cancers

2

3 Chong Chu¹, Viktor Ljungström¹, Antuan Tran¹, Hu Jin¹, Peter J. Park^{1*}

4

5 ¹ Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

6

7 * Correspondence should be addressed to P.J.P. (peter_park@hms.harvard.edu)

8

9 Abstract

10 Insertion of active retroelements—L1s, *Alus*, and SVAs—can disrupt proper genome function
11 and lead to various disorders including cancer. However, the role of *de novo* retroelements
12 (DNRTs) in birth defects and childhood cancers has not been well characterized due to the lack
13 of adequate data and efficient computational tools. Here, we examine whole-genome sequencing
14 data of 3,244 trios from 12 birth defect and childhood cancer cohorts in the Gabriella Miller Kids
15 First Pediatric Research Program. Using an improved version of our tool xTea (x-Transposable
16 element analyzer) that incorporates a deep-learning module, we identified 162 DNRTs, as well
17 as 2 pseudogene insertions. Several variants are likely to be causal, such as a *de novo Alu*
18 insertion that led to the ablation of a whole exon in the *NFI* gene in a proband with brain tumor.
19 We observe a high *de novo* SVA insertion burden in both high-intolerance loss-of-function genes
20 and exons as well as more frequent *de novo Alu* insertions of paternal origin. We also identify
21 potential mosaic DNRTs from embryonic stages. Our study reveals the important roles of
22 DNRTs in causing birth defects and predisposition to childhood cancers.

23

24

25 INTRODUCTION

26 Three types of retroelements are still active in the human genome: long interspersed element-1
27 (L1), *Alu*, and SINE-VNTR-*Alu* (SVA). These retroelements replicate through RNA
28 intermediates by a “copy and paste” mechanism mediated by the LINE-1-encoded proteins; the
29 L1 machinery can also mediate retroduplication of protein-coding genes to generate processed
30 pseudogenes (PPGs). Retroelement insertions into genes may disrupt the function of the gene,
31 potentially leading to a wide spectrum of diseases (Hancks and Kazazian 2016; Burns 2017;
32 Chuong et al. 2017). In particular, *de novo* retroelement (DNRT) insertions have been associated
33 with several developmental disorders and other genetic diseases (Gardner et al. 2019; Brandler et
34 al. 2016; Brandler et al. 2018; Werling et al. 2018; Borges-Monroy et al. 2021). Such DNRT
35 insertions may occur in the gonadal tissues, resulting in a heterozygous germline variant in the
36 proband, or during embryonic development, resulting in a mosaic variant in the proband.

37
38 Compared to other types of *de novo* mutations, DNRTs have been less well studied due to (i) a
39 lack of large whole-genome sequencing (WGS) datasets (especially trios, which are critical for
40 accurately finding *de novo* insertions), and (ii) a lack of reliable computational methods designed
41 specifically for DNRTs. Although several tools have been developed for identifying germline
42 (Gardner et al. 2017; Keane et al. 2013; Thung et al. 2014; Zhuang et al. 2014) and somatic (Lee
43 et al. 2012; Tubio et al. 2014) retroelement insertions, they typically give many high false
44 positive calls for DNRTs because of the low rate of DNRTs and the low variant allele frequency
45 (VAF) of mosaic events. For trio data, one could treat a trio as two pairs of case-control samples
46 and find events that exist in the proband but are absent in both parents; however, we find that a

47 more sensitive detection requires more sophisticated filtering steps specifically designed for the
48 trio design.

49

50 Here, we extend our xTea (Chu et al. 2021) pipeline for *de novo* retroelement insertion
51 identification by further integrating a newly developed machine learning based filtering module.

52 We apply our pipeline to 3,244 trios from the Gabriella Miller Kids First Pediatric Research
53 Program (GMKF) composed of 12 cohorts of different birth defects and childhood cancers, as
54 well as 596 trios from the 1000 Genomes Project (Byrska-Bishop et al. 2022) as reference (Tab.

55 S1). We identified 162 DNRTs from the GMKF cohorts, several of which mobilized to genes
56 where disruptive variants previously have been deemed causative for the diseases. Below, we

57 describe several analyses including detailed examination of likely pathogenic insertions, trio-

58 based phasing to determine whether paternal and maternal contributions are equal, identification

59 of mosaic insertions that occurred at early embryonic stages and PPG insertions, characterization

60 of genes with a higher burden of insertions, and the activity of different SVA subfamilies.

61

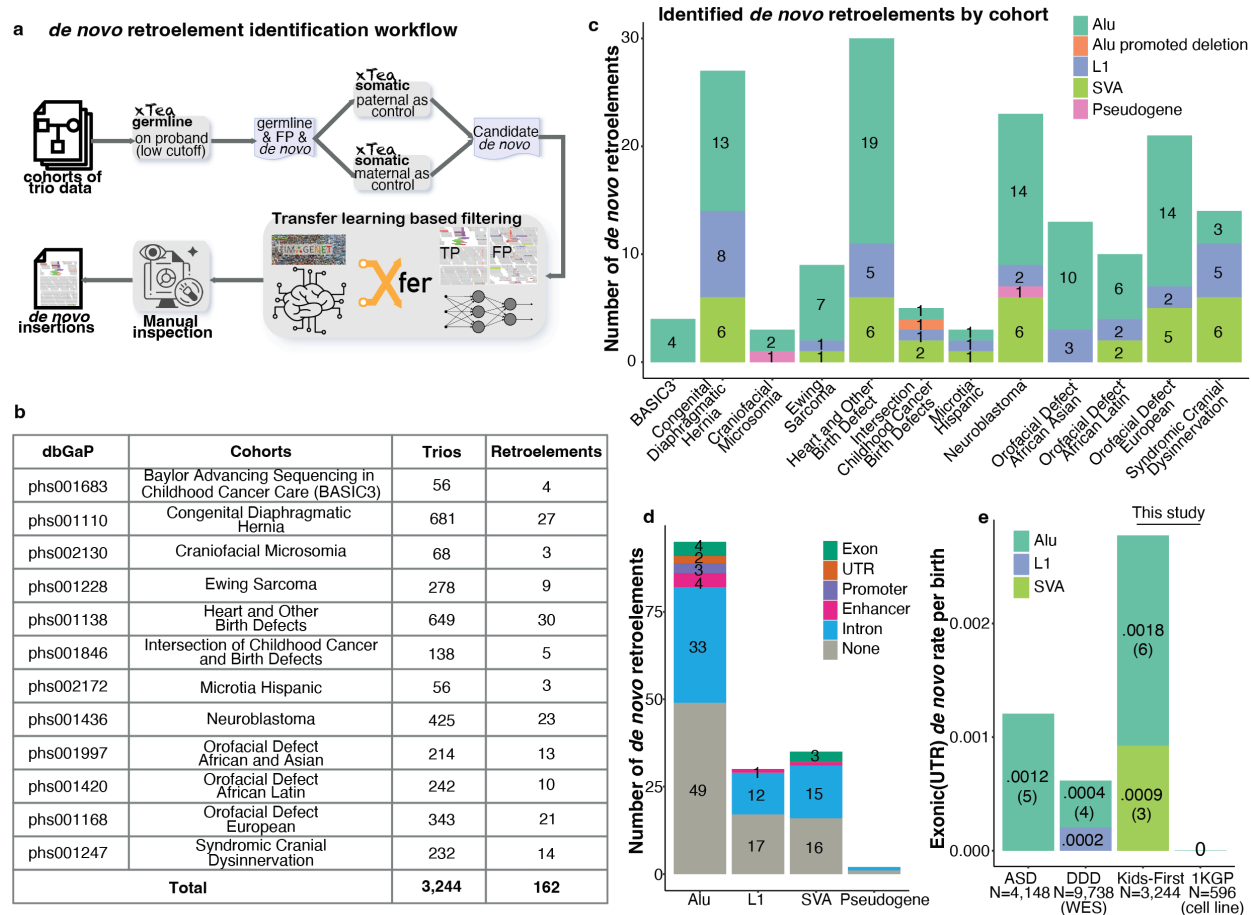
62 RESULTS

63 *De novo* retroelements identification in Kids First and 1KGP data

64 We built an efficient pipeline for DNRT identification (Fig. 1a). Given trio data, we first run the
65 xTea germline module on the proband to identify a set of initial candidates. Because many of the
66 *de novo* events are potentially mosaic events in the proband with low VAFs, we set lenient
67 criteria including on the VAF threshold to ensure high sensitivity. Next, we run the xTea somatic
68 module on the initial candidates with each of the parents as a control. The output from this step
69 still has a high fraction of false positives due to the low cutoff settings. Thus, we apply two
70 additional filtering steps. First, we convert the *de novo* insertion identification problem to an
71 image classification problem by training a machine learning model based on training data labeled
72 from both real and simulated data (see Methods for details). Second, we manually inspect the
73 candidates to curate the final set of DNRTs. Combining the low cutoffs with the filtering steps
74 allows us to design a pipeline with both high sensitivity and high specificity, while its efficient
75 implementation enables application on large cohorts.

76
77 We ran our pipeline on 12 GMKF cohorts (Fig. 1b), totaling 3,244 WGS trios (Tab. S1). Across
78 all cohorts, we identified 162 DNRTs, including 95 *Alu*, 30 L1, 35 SVA, and 2 PPG insertions
79 (Tab. S2). Fig. 1c shows the number of DNRTs identified from each cohort by repeat type.
80 Besides the classic *Alu*, L1, and SVA insertions, we also identified 1 *Alu*-promoted deletion and
81 2 PPG insertions from 3 different cohorts. Of the 162 *de novo* insertions, 7 (4 *Alu* and 3 SVA)
82 are exonic, 2 *Alu* affect UTRs, 3 *Alu* are within promoter regions, and 6 (4 *Alu*, 1 L1, and 1
83 SVA) fall in enhancer regions. Among the others, 33 *Alu*, 12 L1 and 15 SVA are intronic and the
84 rest are intergenic (Fig. 1d). We also analyzed the recently-released trio data from the 1000

85 Genomes Project (1KGP) consisting of 603 trios of which 596 were successfully processed,
86 leading to identification of 26 *Alu*, 12 L1, and 8 SVA DNRTs (Tab. S3).
87
88 We identified 9 and 0 exonic/UTR DNRTs from the GMKF (3,244 births) and 1KGP (596
89 births) cohorts, respectively. In comparison, the Deciphering Developmental Disorders (DDD)
90 study (Gardner et al. 2019) revealed 6 exonic/UTR variants in 9,738 births analyzed with whole-
91 exome sequencing (WES); another study (Borges-Monroy et al. 2021) of autism spectrum
92 disorder (ASD) showed 5 exonic/UTR variants in 4,184 births (Tab. S4). Thus, the exonic
93 (including UTR) *de novo* rate from the GMKF cohorts in this study is $4.5\times$ (0.0027 vs 0.0006)
94 and $2.25\times$ (0.0027 vs 0.0012) higher than in the DDD and ASD studies, respectively. An earlier
95 systematical study on DDD (Deciphering Developmental Disorders Study, 2017) has shown that
96 developmental disorders have higher *de novo* mutation rate than autism, which is concordant
97 with the much higher rate in the GMKF study than in the ASD study. The much lower rate of the
98 DDD study may be due to the difference in the cohorts or possibly due to the lower sensitivity of
99 the method they used.



100

101 **Fig. 1: Overview of the pipeline and identified de novo retroelements.** **a** Schematic outline of
 102 the upgraded xTea workflow. The proband sample from trio data is analyzed using the germline
 103 module with low cutoffs and subsequently filtered using the somatic module and incorporating
 104 parental data to identify candidate de novo events. A transfer learning based filtering step
 105 following a manual inspection step is applied to filter out the false positives. **b** We ran our
 106 pipeline on 3,244 trios from 12 disease cohorts released by the GMKF studies and identified 162
 107 de novo retroelements. **c** The number of identified de novo retroelements by cohorts and repeat
 108 types. Besides the classical TE insertions, we also identified 1 Alu promoted deletion and 2
 109 pseudogene insertions from 3 different cohorts. **d** 7 (4 Alu and 3 SVA), 2 (Alu), 3 (Alu), and 6 (4
 110 Alu, 1 L1 and 1 SVA) de novo retroelements fell in exons, UTRs, promoters, and enhancers
 111 respectively. Out of the rest, 33 Alu, 12 L1, 15 SVA and 1 pseudogene are intronic insertions,
 112 and the remaining are all intergenic ones. **e** The exonic/UTR de novo rate from the GMKF and
 113 1KGP cohorts (this study) compared to the Deciphering Developmental Disorders (DDD) study
 114 and a study of autism spectrum disorders (ASD).

115 **Potentially pathogenic *de novo* retroelement insertions**

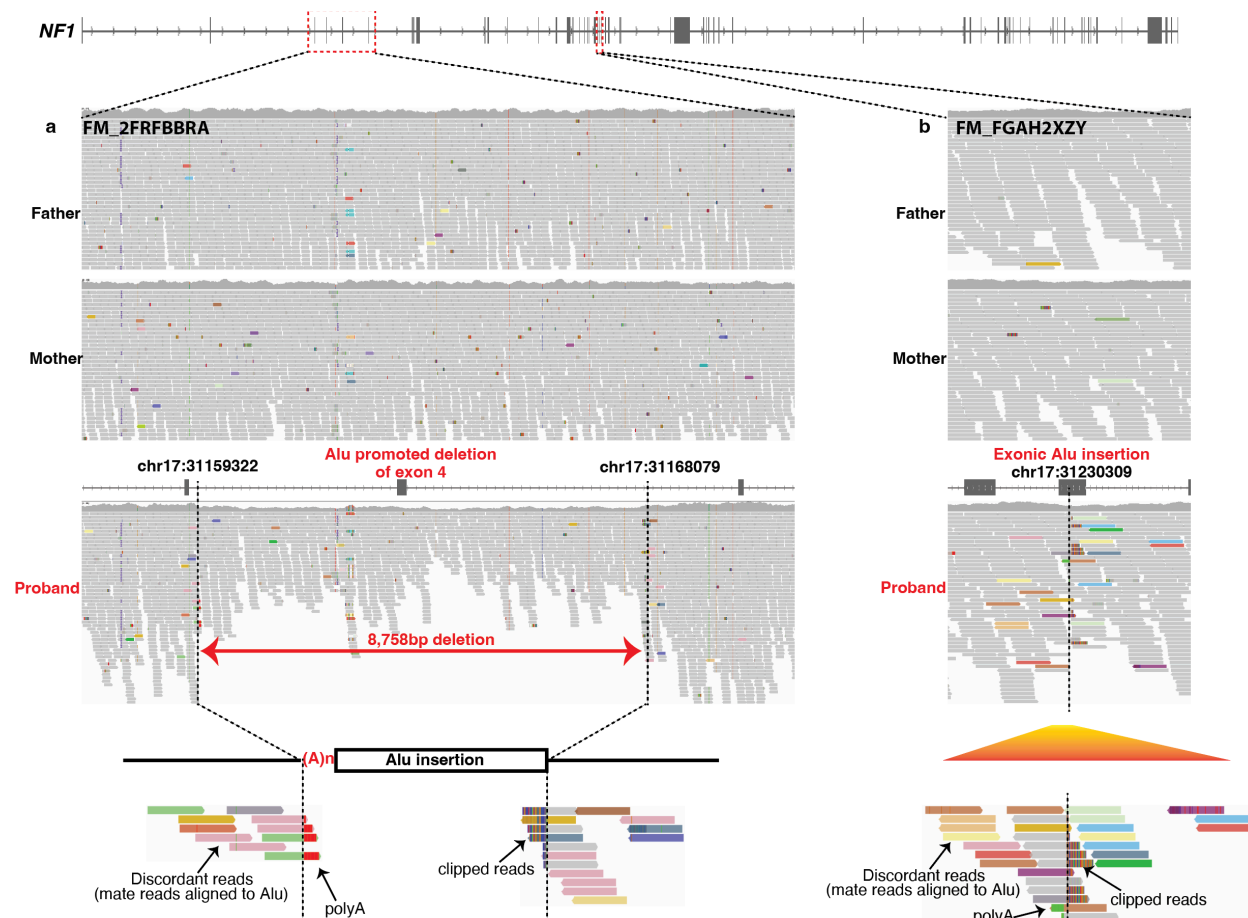
116 WGS data on trios provide an opportunity to identify pathogenic mutations and discover novel
117 disease-associated genes. We prioritized the identified 162 DNRTs by checking whether they
118 occur in (i) exonic or UTR regions; (ii) promoter or enhancer regions, as annotated in
119 ANNOVAR (Yang and Wang 2015); (iii) genes that are associated with the disease. In addition,
120 to account for potential compound heterozygosity, we also checked whether there are detectable
121 second hits for those DNRTs that fall in autosomal recessive genes. In Table 1, we show the
122 selected exonic, UTR, and other potentially pathogenic DNRTs and their annotations.

123
124 We identified two *Alu* insertions that fall in the exonic regions of *NFI*. One insertion promoted
125 an 8,758 bp deletion (chr17:31159322-31168079 on hg38) spanning all of exon 4 (Fig. 2a). The
126 polyA reads at the breakpoints, discordant pairs, and the copy number change strongly support
127 the existence of this complex event. Trio-based phasing and the high VAF indicates this complex
128 event is inherited from the father (a mosaic mutation in the father), but we cannot rule out the
129 possibility that this is an early embryonic event (a mosaic mutation in the child). The other *Alu*
130 insertion was mobilized to the 24th exon of *NFI*. Fig. 2b shows the clipped reads, polyA tail, and
131 the discordant reads present in the proband but absent in the parents. Both probands presented
132 with brain tumor and, since neurofibromatosis type 1 caused by a mutation in *NFI* is an
133 autosomal dominant disorder associated with brain tumor, we propose these two exonic *de novo*
134 *Alu* insertions to be causal mutations. Besides the *NFI* cases, we also identified one SVA
135 insertion mobilized into an exon of the *TRRAP* gene in one proband who was diagnosed with
136 Ptosis (HP:0000508) within the Syndromic Cranial Dysinnervation (SCD) cohort. There have
137 been several reports demonstrating that variants in *TRRAP* are causative for several
138 developmental diseases including autism and syndromic intellectual disability (Cogné et al.

139 2019; Xia et al. 2019; Mavros et al. 2018), although no directly related cases have been reported
140 for SCD. In addition, we also identified one *de novo Alu* insertion in the 3' UTR region of
141 *RREB1* in one proband with Ewing Sarcoma. The *RREB1* gene has been associated with Ewing
142 Sarcoma in several of the earlier studies (Shi et al. 2020; Machiela et al. 2018).

143
144 Although it is more challenging to functionally annotate the intronic DNRTs, we identified two
145 cases with variants that have a high chance of explaining the phenotype observed in the proband.
146 One case is a full length *de novo* L1 insertion mobilized in the sense orientation to the intron of
147 *LTBP1*. The patient was diagnosed with congenital diaphragmatic hernia (CDH)
148 (MONDO:0005711). The *LTBP* gene family has been demonstrated to be highly associated with
149 cutis laxa (a connective tissue disorder; Latin for loose skin) in several studies and the disorder is
150 linked to an increased risk of CDH (Bultmann-Mellin et al. 2015; Zhang et al. 2020; Pottie et al.
151 2021; Urban and Davis 2014). Since most types of cutis laxa are autosomal recessive, we further
152 screened for other types of mutations in the gene and detected one exonic deleterious SNP (2-
153 33567971-C-T; hg38). The SNP is a rare mutation with a population allelic frequency of ~0.78%
154 in the gnomAD database (Chen et al. 2022). We therefore infer that the *de novo* L1 insertion and
155 the SNV lead to a compound heterozygous state resulting in CDH, although functional studies
156 are needed for validation. The other case is a *de novo Alu* insertion identified in a patient
157 diagnosed with right aortic arch (MONDO:0020417). The insertion mobilized to a strong
158 intronic enhancer of the gene *NEK1*. An on-going study (manuscript in preparation) shows that
159 variants in *NEK1* are associated with congenital heart defect (CHD) in an autosomal dominant
160 pattern.

161



162

163 **Fig. 2: Two exonic *de novo* Alu insertions identified on *NF1* gene.** a An *Alu* insertion
164 promoted deletion was identified in one trio. The deletion spans 8,758bp and the two breakpoints
165 fell in intron 3 and 4 leading to deletion of the exon of number 4 of *NF1*. The copy number
166 change in the proband shows the existence of the deletion. The clipped reads, polyA reads, and
167 the discordant reads indicate the existence of an *Alu* insertion. Together, all these features
168 demonstrate the existence of the *Alu* insertion promoted deletion. b In another trio, one *NF1*
169 exonic *Alu* insertion was identified in proband but absent from parents. Similarly, the clipped
170 reads, polyA reads and discordant pairs strongly support the presence of an *Alu* insertion.

171 Table 1. Potential pathogenic *de novo* retroelement insertions identified from the GMKF cohorts.

Cohort	Proband	Repeat	Position	Region	Gene	pLI
CDH	PT 0MDDJ6T0	SVA	chr19:47152964	Exon	<i>SAE1</i>	0.989
CDH	PT 0XTN2CCZ	SVA	chr15:75648240	Exon	<i>IMP3</i>	5.333E-05
ODE	PT PNVAXZ9A	<i>Alu</i>	chr2:19931364	Exon	<i>WDR35</i>	2.648E-16
ODE	PT ZV2GG6F2	<i>Alu</i>	chr1:78643160	Exon	<i>IFI44L</i>	2.218E-10
SCD	PT BGKBJ0JS	SVA	chr7:98892523	Exon	<i>TRRAP</i>	1.0
CCBD	PT MYK8V1XH	<i>Alu</i> promoted deletion	chr17:31159322- 31168079	Exon	<i>NFI</i>	1.0
CCBD	PT 8B6VTS55	<i>Alu</i>	chr17:31230309	Exon	<i>NFI</i>	1.0
ES	PT 469265KR	<i>Alu</i>	chr6:7250370	UTR	<i>RREB1</i>	1.0
ODAA	PT W8FX4W4P	<i>Alu</i>	chr21:46567806	UTR	<i>DIP2A</i>	0.726
ODE	PT GMC33B00	<i>Alu</i>	chr7:75740317	Promoter	<i>HIP1</i>	1.0
HBD	PT 6FZ9C7MC	<i>Alu</i>	chr4:169411142	Enhancer	<i>NEK1</i>	4.219E-12
Neuroblastoma	PT 58J0PB4V	SVA	chr2:227461602	Enhancer	<i>AGFG1</i>	NA
ODAL	PT BPY27QQT	<i>Alu</i>	chr6:26028622	Enhancer	<i>H4C2</i>	NA
ODE	PT KA81JM7G	<i>Alu</i>	chr10:89350311	Enhancer	<i>LIPA</i>	0.0113
SCD	PT AXN0W87J	<i>Alu</i>	chr3:42059106	Enhancer	<i>TRAK1</i>	9.709E-05
Microtia Hispanic	PT P6REK66H	L1	chr2:142867366	Enhancer	<i>KYNU</i>	NA
ES	PT CXSMCH24	<i>Alu</i>	chr13:99200420	Promoter	<i>UBAC2</i>	0.111
Neuroblastoma	PT 7X4TQ0S0	<i>Alu</i>	chr2:30447127	Promoter	<i>LCLAT1</i>	NA
CDH	PT W4Z36EKV	L1	chr2:33137094	Intron	<i>LTBP1</i>	0.526
ODAL	PT 1DAAEZYX	SVA	chr2:229907097	Intron	<i>TRIP12</i>	1.0
Neuroblastoma	PT GA9F5STK	<i>Alu</i>	chr10:91988963	Intron	<i>BTAF1</i>	1.0
SCD	PT QA254K2D	L1	chr5:39017595	Intron	<i>RICTOR</i>	1.0
Neuroblastoma	PT ZPA02FW4	SVA	chr2:121483826	Intron	<i>CLASP1</i>	1.0
ES	PT TGW274S6	SVA	chr2:197421627	Intron	<i>SF3B1</i>	1.0
SCD	PT QA254K2D	SVA	chr3:51564902	Intron	<i>RAD54L2</i>	1.0
Neuroblastoma	PT CBXEYWC5	<i>Alu</i>	chr4:92442143	Intron	<i>GRID2</i>	1.0
ODE	PT 9GBCW4SS	SVA	chr20:47589939	Intron	<i>NCOA3</i>	1.0
SCD	PT EZ1E9P9V	SVA	chr4:13582834	Intron	<i>BOD1L1</i>	1.0
CM	PT V6GS089W	<i>Alu</i>	chr1:232535752	Intron	<i>SIPA1L2</i>	1.0
SCD	PT 24A962K2	L1	chr19:46381076	Intron	<i>PPP5C</i>	1.0

HBD	PT_0M50Q933	<i>Alu</i>	chr1:176182399	Intron	<i>RFWD2</i>	0.999
ES	PT_SNKAQV35	<i>Alu</i>	chr6:1991026	Intron	<i>GMDS</i>	0.9983
HBD	PT_XP2CHGBB	<i>Alu</i>	chr12:8938247	Intron	<i>PHC1</i>	0.997
CDH	PT_0PN34B34	<i>Alu</i>	chr10:1069389	Intron	<i>WDR37</i>	0.997
HBD	PT_DJBYTWQ3	<i>Alu</i>	chr2:127881330	Intron	<i>AMMECR1L</i>	0.996
HBD	PT_X7GE7E9N	<i>Alu</i>	chr12:99293035	Intron	<i>ANKS1B</i>	0.993
HBD	PT_WXTBZHQG	<i>Alu</i>	chr6:110147613	Intron	<i>WASF1</i>	0.914

- 172
- 173 pLI: The probability of loss-of-function intolerance
- 174 CDH: Congenital Diaphragmatic Hernia
- 175 ODE: Orofacial Defect European
- 176 SCD: Syndromic Cranial Dysinnervation
- 177 CCBD: Childhood Cancer Birth Defects
- 178 ES: Ewing Sarcoma
- 179 ODAA: Orofacial Defect African and Asian
- 180 ODAL: Orofacial Defect African Latin
- 181 CM: Craniofacial Microsomia
- 182 HBD: Heart and Other Birth Defects

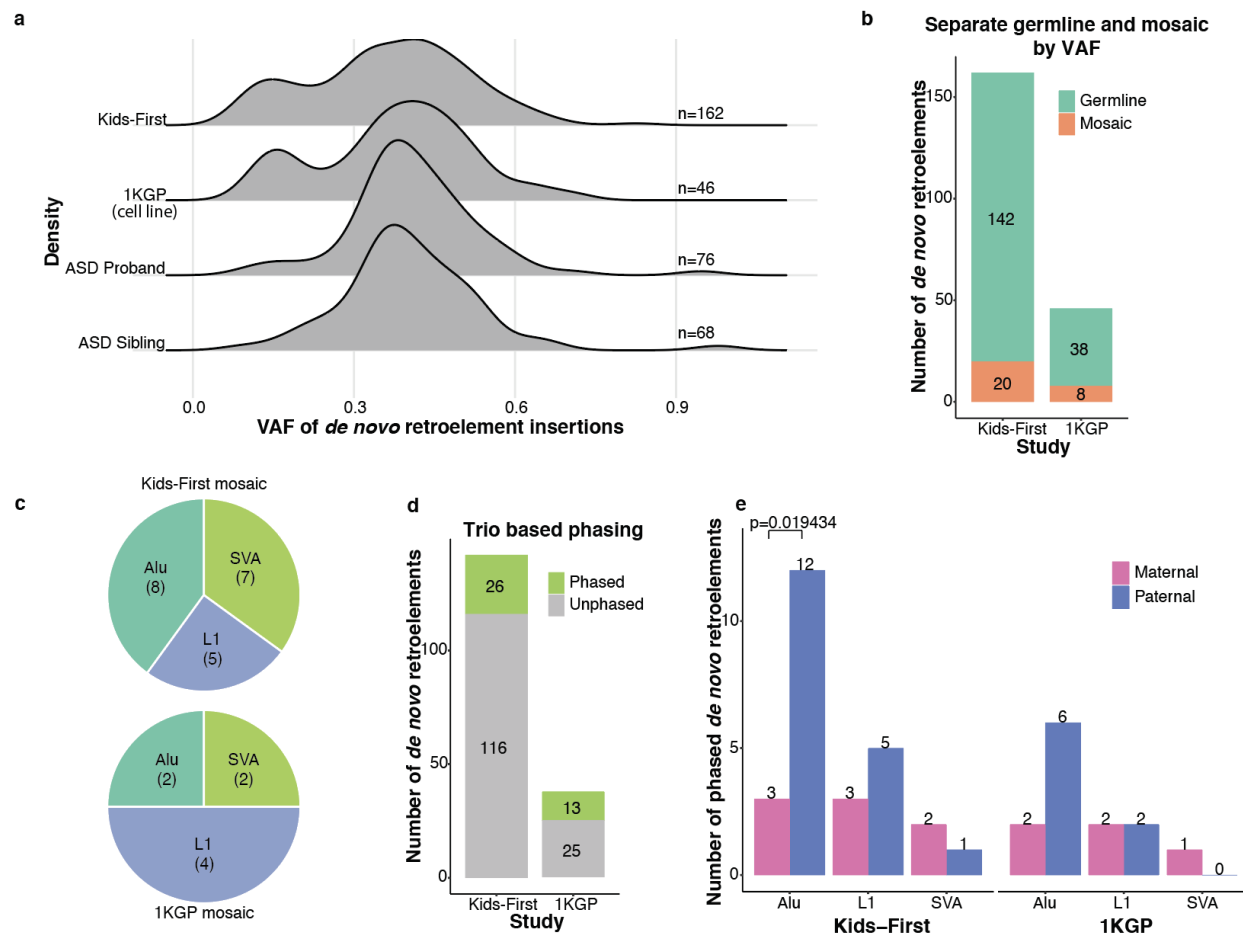
183 **Mosaic and parental origin of *de novo* retroelements**

184 Although *de novo* retroelement insertions refer to those identified in the proband but not in the
185 parents, two types of retroelement insertions are identified as *de novo* in practice. One type is a
186 mosaic insertion in a parent that appears in the child as a germline insertion. The other type is a
187 mosaic insertion that occurs in the early development of a proband, with VAF that depends on
188 the timing of the event. As large cohort studies usually have blood or saliva samples with
189 standard depth (~30X) WGS, only very early embryonic mosaic mutations whose prevalence
190 across tissues is high enough are identified. Earlier trio-based studies on SNVs with ~30-40X
191 WGS have shown that besides the ~100 germline *de novo* small mutations (Jónsson et al. 2017;
192 Kong et al. 2012), a small number of proband mosaic mutations were identified (Byrska-Bishop
193 et al. 2022; Ng et al. 2021). Similarly for *de novo* retroelements, we identified variants of both
194 germline and proband mosaic origin.

195
196 VAF calculation for retroelement insertions is more challenging compared to SNVs/indels
197 because bias will be introduced when aligning the reads containing repetitive sequences (Fig.
198 S1). We optimized our procedure for retroelements VAF estimation (see Method for details) and
199 calculated the VAF for each of the identified 162, 46, and 144 DNRTs in the GMKF, 1KGP, and
200 ASD cohorts, respectively (Fig. 3a). Note that the VAFs previously reported for the ASD
201 DNRTs (separated by probands and siblings) (Borges-Monroy et al. 2021) used an earlier
202 version of xTea, and they have now been recalculated with the latest version of xTea.

203
204 In the VAF density plot (Fig. 3a), we observed two peaks for GMKF and 1KGP, indicating the
205 existence of germline and mosaic DNRTs in the proband. Using a gaussian mixture model to
206 separate the two event classes, we identified 20 (out of 162) and 8 (out of 46) mosaic DNRTs for

207 the GMKF and 1KGP cohorts, respectively (Fig. 3b). We found mosaic DNRTs for all three
208 types of TE insertions (8 Alu, 5 L1, and 7 SVA for GMKF and 2, 4, 2 for 1KGP, respectively;
209 Fig. 1c). The prevalence of mosaic DNRTs in 1KGP (17%) could be explained by somatic
210 DNTRs occurring in cell culture. Indeed, a higher mosaic rate for SNP/Indel has been reported
211 on the same data, attributed to ongoing mutation processes in cell culture (Byrska-Bishop et al.
212 2022; Ng et al. 2021). As a comparison, the mosaic DNRTs in GMKF was 12%, likely due to
213 early embryonic events. This result is consistent with an independent study on normal colon that
214 identified one mosaic L1 insertion occurring at the fourth cell division (Nam et al. 2022).
215
216 For the 142 (out of 162) and 38 (out of 46) germline DNRTs in GMKF and 1KGP respectively,
217 we ran trio-based phasing using nearby heterozygous SNPs (see Method for details) to determine
218 their parental origins. Limited by the read length and insert size of the sequencing libraries, we
219 were only able to phase only 26 out of 142 GMKF DNRTs and 13 out of 38 1KGP DNRTs (Fig.
220 3d). Nonetheless, our inspection of the phased DNRTs by repeat type revealed a significant
221 enrichment ($p=0.019$; binomial test) for *Alu* DNRTs of paternal origin in the GMKF cohort while
222 no statistical significance was reached for L1 and SVA (Fig. 3e). A similar trend was also
223 observed in the 1KGP cohort, although it was not statistically significant due to small sample size
224 (Fig. 3e). A similar pattern was also reported in an earlier study on large pedigree data of three
225 generations (Feusier et al. 2019). Based on these results, we hypothesize that the *Alu* insertion
226 rate is higher in the cell types involved in spermatogenesis compared to oocytes.



227

228 **Fig. 3: VAF characterization and trio-based phasing of *de novo* retroelements.** **a** The
 229 variation allele frequency (VAF) for the *de novo* retroelements identified in this study (162 from
 230 GMKF and 46 from 1KGP), as well as the 144 (76 from proband and 68 from sibling) reported
 231 *de novo* retroelements from the ASD study (Borges-Monroy et al. 2021). The two peaks in
 232 density plots of GMKF and 1KGP cohorts indicate the presence of mosaic retroelements in
 233 proband. **b** Gaussian mixture model separation of germline and mosaic *de novo* retroelements
 234 based on the density of VAF. Out of the 162 and 46 *de novo* retroelements, 20 and 8 mosaic ones
 235 were identified for the GMKF and 1KGP respectively. **c** We checked the repeat types of the
 236 mosaic events. 8 *Alu*, 5 L1, and 7 SVA were annotated for the 20 GMKF mosaic retroelements,
 237 and 2 *Alu*, 4 L1, and 2 SVA were annotated for the 1KGP mosaic ones. **d** For the germline *de*
 238 *de novo* retroelements, we applied a trio-based phasing step and identified 26 (out of 142) and 13
 239 (out of 38) phasable ones for the GMKF and 1KGP respectively. **e** Further checking the phased
 240 ones from the GMKF cohorts, we found more *de novo Alu* insertions were transmitted from
 241 father than from mother (p -value=0.019, exact binomial test), while for L1 and SVA it is unclear.
 242 The same pattern was observed from the 1KGP data, although not statistically significant due to
 243 the sample size being small.

244

245

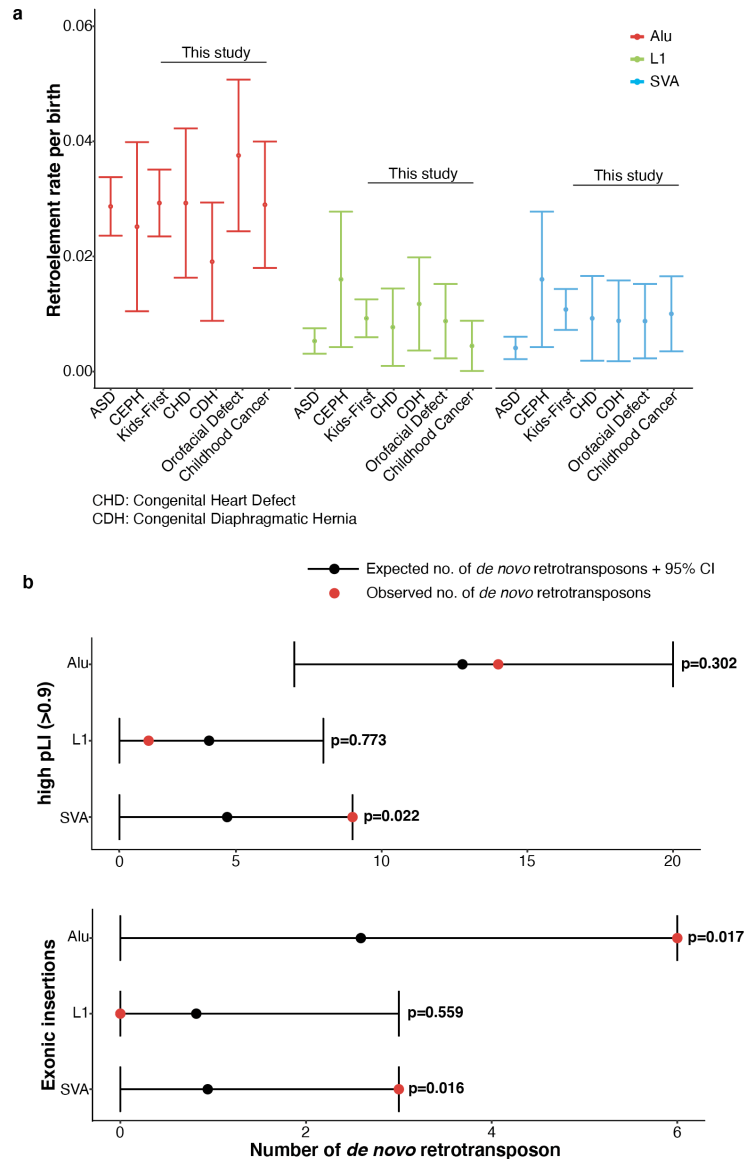
246 ***De novo* rate estimation and enrichment of deleterious retroelements**

247 The large number of trios allowed us to estimate the *de novo* rate of the retroelements in the birth
248 defect cohorts. Previous estimates were highly variable (Borges-Monroy et al. 2021), with our
249 recent study on the Simons Simplex Collection (SSC) of ASD cases (2,288 families with
250 proband and siblings) using an earlier version of xTea giving 1/26 birth (adjusted to 1/21 after
251 considering detection sensitivity benchmarked from long reads). Below, we computed the exact
252 binomial confidence intervals on all the GMKF cohorts combined as well as for four disease
253 groups that have sufficient sample sizes: Congenital Heart Defect (CHD), Congenital
254 Diaphragmatic Hernia (CDH), Orofacial Defect, and Pediatric Tumor. The Orofacial Defect
255 group consists of 3 cohorts (phs001997, phs001420, and phs001168 in Fig. 1b) of the same
256 defect but from different populations. The Pediatric Tumor group consists of 4 cohorts
257 (phs001683, phs001228, phs001846, and phs001436 in Fig. 1b) of different tumor types. For
258 reference, we list the *de novo* insertions for two earlier studies on the SSC ASD (Tab. S3), 1KG
259 (Tab. S4), and the Utah Centre d'Etude du Polymorphisme Humain (CEPH) cohorts (Tab. S5).

260
261 The *de novo* retroelement rate by repeat type are shown in Fig. 4a for the different cohorts. The
262 rate of all the GMKF cohorts is around 1/34 (1 per 34 births) for *Alu*, 1/108 for L1, and 1/93 for
263 SVA, and 1/20 for all 3 TE types combined. Compared to the ASD cohort, we observed a similar
264 rate (1/34 vs. 1/35) for *Alu* insertions, but much higher rate (1/108 vs. 1/189) for L1 and (1/93 vs.
265 1/244) SVA insertions. Our results are consistent with a previous study (Deciphering
266 Developmental Disorders Study 2017) that showed a higher rate for *de novo* SNVs/indels in the
267 DDD cohort compared to the ASD cohort.

268

269 Given the higher rate of insertions in the GMKF cohorts, we examined whether the DNRTs are
270 enriched in genes likely to be intolerant of loss of function (LoF), as determined by the pLI score
271 (Lek et al. 2016) and whether they enriched in exonic regions. We first simulated the distribution
272 of DNRTs across the genome with the following idea. When L1 insertions (also for *Alu* and SVA
273 as they rely on the L1 protein) are integrated into the genome, they prefer specific motifs
274 (consensus TTTTT/AA). Inspired by the landmark experimental study on L1 endonuclease
275 activity (Flasch et al. 2019) that engineered L1 insertions in cultured cell lines to characterize the
276 pattern of the cleavage sites, we simulated the insertion sites based on the frequency of cleavage
277 motif sequences. A key step here is to construct the weight matrix based on the frequency of the
278 cleavage motifs. We calculated the frequency from germline insertions, different from the
279 experimental study (Flasch et al. 2019) that inferred the matrix from engineered somatic
280 insertions, but the logo plots of cleavage site sequences were almost identical (Fig. S2; see
281 Method for details on the simulation procedure). Compared to the background distribution
282 generated from 10,000 simulations, we observed an enrichment ($p=0.022$) for *de novo* SVA
283 insertions in the high pLI (>0.9) genes, but not for *Alu* and L1. We also find an enrichment of *de*
284 *novo Alu* ($p=0.017$) and SVA ($p=0.016$) insertions in the exons, but no statistical significance for
285 *de novo* L1 insertions (Fig. 4b). Both analyses suggest that *de novo* retrotransposons—especially
286 *de novo* SVA insertions—found in the birth defect cohorts are more likely to be deleterious.
287



288

289 **Fig. 4: *De novo* rate and enrichment analysis of retroelements.** **a** We calculated the *de novo*
 290 rate and 95% confidence interval using an exact binomial confidence interval estimate with
 291 x =number of retroelements and N =number of births. ASD and CEPH are from two previous
 292 studies and the rest are from the GMKF cohorts in this study. GMKF is for all the 12 cohorts,
 293 “Orofacial Defect” is combined from 3 orofacial defect cohorts (phs001997, phs001420, and
 294 phs001168), and “Childhood Cancer” is combined from 4 cancer related cohorts (phs001683,
 295 phs001228, phs001846, and phs001436). The *de novo* rate for L1 and SVA in this study is
 296 clearly higher than ASD, with a similar rate for *Alu*. **b** We checked whether the identified *de*
 297 *de novo* retroelements from the GMKF cohorts were enriched in genes whose pLI>0.9 (top) or
 298 enriched in exonic regions (bottom). We ran 10,000 simulations (details in Method) for *Alu*, L1
 299 and SVA, and compared with the number of observed ones. There is a statistically significant
 300 enrichment of SVA insertions fallen in genes with pLI>0.9. Both *Alu* and SVA insertions were
 301 also found enriched in exonic regions.

302

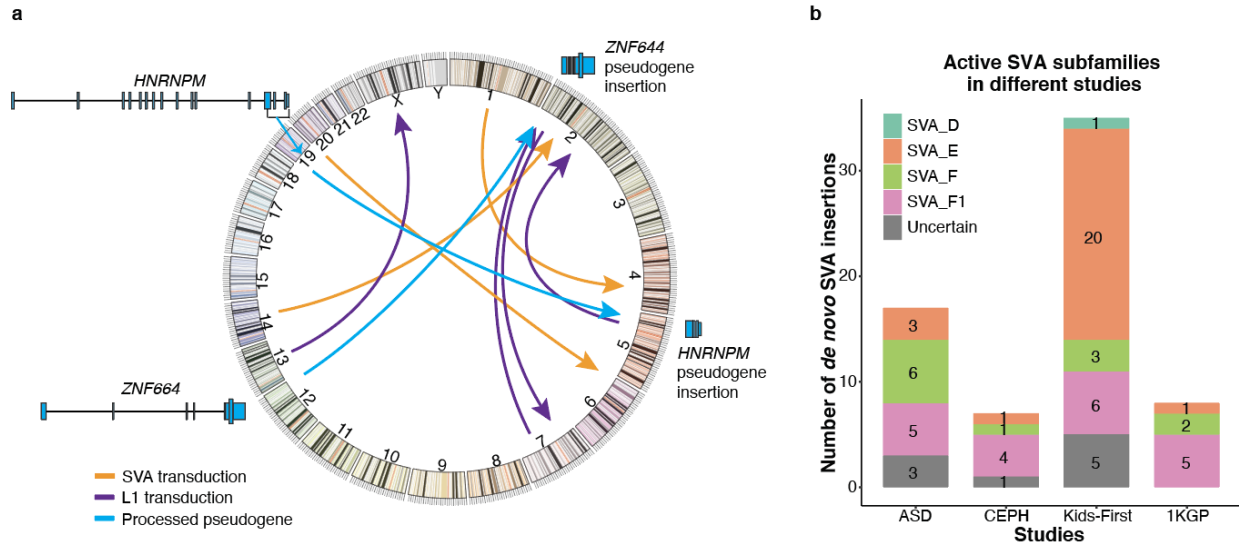
303 ***De novo* processed pseudogene insertion and *de novo* TE insertion activity**

304 Although germline and somatic PPG insertions in human have been reported (Esnault et al.
305 2000; Feng and Li 2021; Schrider et al. 2013; Ewing et al. 2013; Cooke et al. 2014), *de novo*
306 PPG insertions have not been well characterized due to their low *de novo* rate and the lack of
307 detection tools designed for PPGs. The study on the DDD cohort reported 2 *de novo* PPG
308 insertions from 9,738 WES trios (Gardner et al. 2019). Here, we extended our xTea method for
309 *de novo* PPG insertion detection and identified 2 *de novo* PPG insertions from 3,244 trios (Fig.
310 S4), suggestive of a 3-fold increase in insertion rates. One insertion originated from gene
311 *HNRNPM* and reverse-transcribed into a truncated PPG insertion within the intronic region of
312 *FBXL7*. The other originated from gene *ZNF664* and mobilized to an intergenic region on
313 chromosome 2 (Fig. 5a; S4).

314
315 For L1 and SVA, the flanking regions (mainly 3' for L1, and 5' and 3' for SVA) sometimes
316 transpose together with the retrotransposon to form transductions (TDs). Because most of the TD
317 sequences are unique in the genome, they could be used to trace the source elements of the
318 insertions. We ran xTea on all the proband cases having *de novo* L1 and SVA insertions and
319 identified 4, 1, 1, and 1 transduction from GMKF, 1KGP, ASD and CEPH cohort, respectively.
320 In the original CEPH paper (Feusier et al. 2019), the authors reported 3 L1 TDs, but after manual
321 inspection, we could only confirm one of them. Our results indicate a small number of active *de*
322 *nov*o transduction events.

323
324 Finally, we investigated the activity of *de novo* insertions for SVA, one of the youngest
325 retrotransposons in the human genome. SVA_E and SVA_F subfamilies are the known major

326 active subfamilies (Wang et al. 2005; Hancks and Kazazian 2010); SVA_F has two subfamilies,
327 SVA_F1 and CH10_SVA_F, that fused with part of *MAST2* to form new structure and transpose.
328 CH10_SVA_F has two *Alu* copies at the two ends, making the annotation from short reads
329 difficult; thus, the activity of CH10_SVA_F is not explored here. SVA_F1 insertions can have
330 reads on one side that cover the *MAST2* region and reads on the other sides aligned to the SVA_F
331 region, thus allowing for annotation. Using an SVA annotation module we recently developed
332 (Chu et al. 2023), we examined each of the *de novo* SVA insertions from the ASD (n=17), CEPH
333 (n=7), GMKF (n=35), and 1KGP (n=8) study (Fig. 5b) and found that SVA_E (25/58; 9 could
334 not be classified) and SVA_F1 (20/58) are the two most active SVA subfamilies in human
335 genome.



336

337 **Fig. 5: Activity of *de novo* retroelements.** a 4 L1 transductions, 3 SVA transductions, and 2

338 pseudogene insertions were identified. Each arrow points from the source elements to the

339 insertion site. Specifically, the 2 pseudogene insertions were originated from gene *HNRNPM* and

340 gene *ZNF664*, with *HNRNPM* insertion is truncated while the *ZNF664* one is of full length. **b**

341 For all the *de novo* SVA retrotransposon insertions identified from the four studies (17 ASD, 7

342 CEPH, 35 GMKF and 8 1KGP), we further checked the subfamilies of each insertion. 3, 1, 5

343 were not well annotated from ASD, CEPH, and GMKF respectively. Out of those well annotated

344 ones, SVA_F1 (20/58) and SVA_E (25/58) are the most active subfamilies.

345

346 **DISCUSSION**

347 With a systematic analysis of the 12 GMKF WGS cohorts, our results revealed the different
348 pathogenic roles DNRTs may have in causing birth defects and childhood cancers. Our study
349 highlights the importance of WGS in identifying causal mutations beyond the standard variant
350 types and the effectiveness of the improved xTea pipeline we developed for characterizing
351 DNRTs in large cohort analyses and disease diagnosis.

352
353 Compared to germline DNRTs, mosaic DNRTs occurring in the proband are substantially more
354 difficult to identify due to their low VAFs. Our analysis found a set of DNRTs with low VAFs,
355 with a peak around 25% in the VAF density plot for GMKF. Although they could be variations
356 in read sampling or germline DNRTs in aneuploid regions, we suspect that at least some of them
357 are mosaic variants that occurred very early in embryonic development. A peak at the similar
358 VAF (~ 20%) was observed for 1KGP but those are likely to be somatic variants that arose in
359 cell culture; for ASD cohorts, no such peaks were found, increasing the likelihood that the peak
360 in the GMKF cohort may be due to mosaic variants. A further analysis on other types of
361 mutations (especially *de novo* SNP/Indels) or sequencing data with higher coverage may provide
362 additional information for the prevalence of mosaic mutations.

363
364 Strong support for the functional importance of *de novo* SVA insertions came from their
365 enrichment in genes that are estimated to be intolerant to loss of function ($pLI > 0.9$) as well as
366 exonic enrichment for *de novo* SVA and *Alu* insertions. We also identified several cases in which
367 DNRTs disrupted genes that are associated with the disease in the literature. Most notable are the
368 *Alu*-promoted deletion and an exonic insertion on the *NFI* gene. Together with the earlier studies
369 (Wallace et al. 1991; Vogt et al. 2014; Wimmer et al. 2011) on retroelements on *NFI* gene, our

370 study indicates a hot spot for RTs on the *NFI* gene. For other cases, experimental validation is
371 needed to ascertain whether those candidate DNRTs are causative of the observed phenotypes.
372
373 Most studies on retroelement insertions have focused on germline L1 and *Alu* insertions, as they
374 are the easiest to identify due to their relatively simple structure. SVA are substantially more
375 complicated in its structure, especially with the VNTR (variable number of tandem repeats)
376 region in the middle often causing mis-annotation, as revealed in our recent work on SVA
377 detection and annotation (Chu et al. 2023). The increasing adaption of long-read sequencing
378 platforms will greatly enhance detection sensitivity and specificity, as well as expanding the
379 regions of the genome that can be interrogated. The GMKF cohort continues to increase in size,
380 as do numerous other WGS studies for various genetic diseases. The resulting datasets will
381 enable discovery of additional disease-associated genes, allow for more accurate inferences on
382 the rates of mosaic insertion events, help pinpoint active source elements through transduction
383 events, and shed light on the expanding role that retroelements play in disease initiation.
384

385 **METHOD**

386 ***De novo* retroelements identification from trio data**

387 Compared to germline retroelement insertions, DNRTs are rare, thus sensitivity is important for
388 identification. Here, we optimized our xTea method for DNRT insertion identification. Fig. 1a
389 shows the major steps: for one given trio data, (i) we ran xTea germline module on the proband
390 sample (by default with parameters “--nclip 2 --cr 0 --nd 3 --nfclip 2 --nfdisc 3”); (ii) we ran
391 xTea somatic module on the candidates generated from step (i) with alignments from both
392 parents as controls; (iii) we further developed a machine learning based filtering module
393 (manuscript in preparation) to filter out the false positives. For each candidate insertion, we
394 converted the alignments to images using BamSnap (Kwon et al. 2021), where each candidate is
395 composed of three images from the trio on the same location. We first prepared a positive
396 training set from semi-simulated data, where we selected germline heterozygous retroelements
397 from one sample and viewed it as the “proband”, and then we selected two unrelated samples
398 that do not have retroelement insertions on this location as the “parents”. From these three
399 “combined” samples on this location we generated one positive image. In this way, we prepared
400 6952 positive training images. Then we prepared the same number of negative training images
401 from xTea output on two cohorts (phs001228 and phs001168) that we had manually inspected
402 for the true positives. Next, we trained a model from the positive and negative training sets.
403 Then, for each candidate image, we predict it to true positive or false positive; (iv) lastly, we ran
404 manual inspection on each of the candidates to select the true positive variants.

405

406 ***De novo* retroelements annotation**

407 For each identified DNRT, we first annotated it as exonic, 5' UTR, 3' UTR, intronic, or
408 intergenic based on the GENCODE (v28 on GRCh38) gene annotation file. Then, we ran

409 ANNOVAR (version downloaded on May 2022) to annotate DNRTs fall in promoter and
410 enhancer regions. We used the pLI score from the ExAC study (Lek et al. 2016) to annotate the
411 estimated intolerance of each gene to mutations, and $pLI > 0.9$ are annotated as “high intolerant”.

412
413 In addition, we also ran subfamily annotation specifically for *de novo* SVA insertions, because
414 the active subfamilies are less well characterized in large cohorts. To annotate the SVA insertion
415 subfamilies, we first collected all the discordant and clipped reads originated from the insertion.
416 Then we ran local assembly on the collected reads using the xTea assembly module. Next, we
417 ran the SVA annotation module (initially developed for a different study whose manuscript in
418 review) on the assembled sequences to get the subfamily information of each *de novo* SVA
419 insertion. Insertions annotated to more than one subfamily or failed to be assembled were
420 annotated as “uncertain”.

421

422 **Variation allele frequency (VAF) calculation**

423 Different from SNV VAF estimated by calculating the ratio between the number of reads
424 containing the mutation and the total number of reads at the site, calculating VAF for DNRTs is
425 more challenging. We illustrate the major biases introduced due to reads alignments in Fig. S1.
426 Generally, for one clipped read with part of the read from the retroelements and the other part
427 from the flanking regions, if the length of the retroelement part is short (by default, BWA mem
428 has a minimum kmer length of 19), then the read will be aligned elsewhere as the unique part
429 (part from flanking region) is short, as a result, these reads will not be counted when calculating
430 the VAF. To correct the introduced bias, when counting the number of full mapped reads
431 covering the breakpoint, we skipped reads having short overlap (by default < 19) with the
432 flanking region. In addition to calculate VAF from reads covering the breakpoint, discordant

433 pairs can also be used for calculating the VAF. Basically, within the given range (by default,
434 insert-size) we count the number of discordant pairs and concordant pairs, and then calculate the
435 VAF. In Fig. S1, we show the calculated clip- and discordant-based VAFs for 115,115 germline
436 TE insertions. In practice, we took the average of the clip- and discordant-based VAF as the
437 VAF of each DNRT.

438

439 **Trio based phasing**

440 For some of the identified DNRTs, we can phase them to derive paternal or maternal origin
441 based on the nearby heterozygous SNVs. For example, if we find heterozygous SNVs in the
442 father and identify the same SNVs within the discordant pairs of the DNRT in the proband, then
443 we can infer the DNRT to be of paternal origin (Fig. S3c,d). However, in practice, we can only
444 find nearby heterozygous SNVs for a small fraction of DNRTs. To broaden the range of phasable
445 DNRTs, for the germline DNRTs in proband, if we observe heterozygous SNVs in one of the
446 parents, but the same SNVs are found in the non-DNRT reads, then we can infer that the DNRT
447 is inherited from the other parent. To achieve this, we first adopt a Gaussian Mixture Model (2
448 mixture components) to classify the DNRTs as germline or mosaic, and then phase the set of
449 germline variants. To call heterozygous SNVs from the parents, we ran samtools “mpileup” and
450 “call” on the local region of each DNRT. To call SNVs for each DNRT from the proband, we
451 first separated the reads aligned to the local region to two groups: “DNRT reads” and “non-
452 DNRT reads”. For each group, we ran samtools “mpileup” and “call” to identify the SNVs. In
453 addition, we also ran manual inspection for each phased DNRTs from BamSnap screenshots to
454 further validate the phased the DNRTs. The whole pipeline is shown in Fig. S3.

455

456 ***De novo* rate estimation**

457 To estimate the *de novo* rate of retroelements, we adopted the exact binomial confidence interval
458 estimate, where X is the number of retroelements, and N is the number of births. To compare the
459 *de novo* rate among different disease cohorts, we show the results of “Syndromic cranial
460 dysinnervation”, “Heart birth defects”, “Orofacial birth defects”, “childhood cancer”, where
461 “Orofacial birth defects” results are merged from three orofacial birth defects of different
462 populations (phs001168, phs001420, and phs001997), and the “childhood cancer” results are
463 merged from the four tumor cohorts (phs001436, phs001228, phs001683, and phs001846). In
464 addition, we also compared the overall *de novo* rate in the birth defect disorders with two earlier
465 studies on autism (Borges-Monroy et al. 2021) and large pedigree of normal samples (Feusier et
466 al. 2019), where we used the number of *de novo* retroelements and the number of births in their
467 released results.

468

469 ***De novo* retroelements enrichment analysis**

470 To test whether the identified DNRTs are enriched in the GMKF data or not, we need a control
471 model that simulates the random hits of DNRTs. However, it is known that RTs are not purely
472 randomly happened on the genome. A recent study based on engineered L1s in cell lines inferred
473 that there are specific endonuclease (EN) cleavage motifs (Flasch et al. 2019). Here, we adopt
474 the similar approach to build the control model as described in the endonuclease activity study
475 from engineered L1s (Flasch et al. 2019). As shown in Fig. S2, we first gather all the possible
476 EN cleavage motifs, then for each motif we estimate its frequency, which later will be used as
477 the probability of a simulated insertion occurring with the motif. Differently, here we use
478 germline insertion rather than engineered *de novo* L1s to gather all the possible motifs. To
479 achieve this, we first ran xTea on the 1KGP high depth WGSs and collect the high-quality TE
480 insertions (labeled with “tprt_both” in xTea output indicating exist of both the TSD and polyA

481 tail; and require the population $AF > 0.01$). For each TE insertion, we collect the first left 4 bases
482 and the right 3 bases at the breakpoint (adjusted accordingly for antisense cases). Then, we put
483 all the collected motifs in a table and calculate the frequency for each one. We also adopt the
484 same approach as described in the mentioned study (Flasch et al. 2019) to include potential motif
485 not recruited in the motif table: We split the 7-base motif to 3 independent segments: first 2
486 bases, middle 4 bases and the last base. For each segment, we calculate the frequency of each
487 sub-motif based on the frequency of the 7-base motifs. In this way, we have 3 tables of sub-
488 motifs whose frequency have been calculated. Thus, to generate a 7-base motif, we generate the
489 3 segments separately and for each segment we select a sub-motif based on the frequency table.
490 Then, we merge the 3 segments to one 7-base motif. Now, given one motif, we need to find out
491 where on the genome this motif can be generated. To achieve this, we build another 7-base motif
492 table for the whole genome, where we save all the positions of each motif. In this way, once
493 given a motif, we randomly select one from all the positions.

494

495 To construct the control model, for each round we simulated the same 163 DNRTs with our
496 pipeline, and we repeated the experiments for 10,000 times. We did enrichment analysis for
497 insertions fallen in exon regions and in pLI high (> 0.9) genes separately.

498

499 **Data availability**

500 The 12 cohorts of pediatric whole genome sequencing (WGS) data were accessed through the
501 portal of Gabriella Miller Kids First Pediatric Research Program <https://portal.kidsfirstdrc.org/>.

502 The high depth trio based WGS data from the 1000 Genomes Project were downloaded from the
503 International Genome Sample Resource (IGSR) at <https://www.internationalgenome.org/data/>.

504

505 **Code availability**

506 Source code for the de novo retroelements identification is available at
507 <https://github.com/parklab/xTea> (doi:10.5281/zenodo.4743788).

508

509 REFERENCES

- 510 Borges-Monroy, Rebeca, Chong Chu, Caroline Dias, Jaejoon Choi, Soohyun Lee, Yue Gao,
511 Taehwan Shin, Peter J. Park, Christopher A. Walsh, and Eunjung Alice Lee. 2021. “Whole-
512 Genome Analysis Reveals the Contribution of Non-Coding de Novo Transposon Insertions
513 to Autism Spectrum Disorder.” *Mobile DNA* 12 (1): 28.
- 514 Brandler, William M., Danny Antaki, Madhusudan Gujral, Morgan L. Kleiber, Joe Whitney,
515 Michelle S. Maile, Oanh Hong, et al. 2018. “Paternally Inherited Cis-Regulatory Structural
516 Variants Are Associated with Autism.” *Science* 360 (6386): 327–31.
- 517 Brandler, William M., Danny Antaki, Madhusudan Gujral, Amina Noor, Gabriel Rosanio,
518 Timothy R. Chapman, Daniel J. Barrera, et al. 2016. “Frequency and Complexity of De
519 Novo Structural Mutation in Autism.” *American Journal of Human Genetics* 98 (4): 667–
520 79.
- 521 Bultmann-Mellin, Insa, Anne Conradi, Alexandra C. Maul, Katharina Dinger, Frank Wempe,
522 Alexander P. Wohl, Thomas Imhof, et al. 2015. “Modeling Autosomal Recessive Cutis
523 Laxa Type 1C in Mice Reveals Distinct Functions for *Ltbp-4* Isoforms.” *Disease Models &*
524 *Mechanisms* 8 (4): 403–15.
- 525 Burns, Kathleen H. 2017. “Transposable Elements in Cancer.” *Nature Reviews. Cancer* 17 (7):
526 415–24.
- 527 Byrska-Bishop, Marta, Uday S. Evani, Xuefang Zhao, Anna O. Basile, Haley J. Abel, Allison A.
528 Regier, André Corvelo, et al. 2022. “High-Coverage Whole-Genome Sequencing of the
529 Expanded 1000 Genomes Project Cohort Including 602 Trios.” *Cell* 185 (18): 3426–40.e19.
- 530 Chen, Siwei, Laurent C. Francioli, Julia K. Goodrich, Ryan L. Collins, Qingbo Wang, Jessica
531 Alföldi, Nicholas A. Watts, et al. 2022. “A Genome-Wide Mutational Constraint Map
532 Quantified from Variation in 76,156 Human Genomes.” *bioRxiv*.
533 <https://doi.org/10.1101/2022.03.20.485034>.
- 534 Chu, Chong, Rebeca Borges-Monroy, Vinayak V. Viswanadham, Soohyun Lee, Heng Li,
535 Eunjung Alice Lee, and Peter J. Park. 2021. “Comprehensive Identification of Transposable
536 Element Insertions Using Multiple Sequencing Technologies.” *Nature Communications* 12
537 (1): 3836.
- 538 Chuong, Edward B., Nels C. Elde, and Cédric Feschotte. 2017. “Regulatory Activities of
539 Transposable Elements: From Conflicts to Benefits.” *Nature Reviews. Genetics* 18 (2): 71–
540 86.
- 541 Cogné, Benjamin, Sophie Ehresmann, Eliane Beauregard-Lacroix, Justine Rousseau, Thomas
542 Besnard, Thomas Garcia, Slavé Petrovski, et al. 2019. “Missense Variants in the Histone
543 Acetyltransferase Complex Component Gene *TRRAP* Cause Autism and Syndromic
544 Intellectual Disability.” *American Journal of Human Genetics* 104 (3): 530–41.
- 545 Cooke, Susanna L., Adam Shlien, John Marshall, Christodoulos P. Pipinikas, Inigo
546 Martincorena, Jose M. C. Tubio, Yilong Li, et al. 2014. “Processed Pseudogenes Acquired
547 Somatically during Cancer Development.” *Nature Communications* 5 (1): 1–9.
- 548 Esnault, Cécile, Joël Maestre, and Thierry Heidmann. 2000. “Human LINE Retrotransposons
549 Generate Processed Pseudogenes.” *Nature Genetics* 24 (4): 363–67.
- 550 Ewing, Adam D., Tracy J. Ballinger, Dent Earl, Broad Institute Genome Sequencing and
551 Analysis Program and Platform, Christopher C. Harris, Li Ding, Richard K. Wilson, and
552 David Haussler. 2013. “Retrotransposition of Gene Transcripts Leads to Structural
553 Variation in Mammalian Genomes.” *Genome Biology* 14 (3): R22.

- 554 Feng, Xiaowen, and Heng Li. 2021. “Higher Rates of Processed Pseudogene Acquisition in
555 Humans and Three Great Apes Revealed by Long-Read Assemblies.” *Molecular Biology*
556 *and Evolution* 38 (7): 2958–66.
- 557 Feusier, Julie, W. Scott Watkins, Jainy Thomas, Andrew Farrell, David J. Witherspoon, Lisa
558 Baird, Hongseok Ha, Jinchuan Xing, and Lynn B. Jorde. 2019. “Pedigree-Based Estimation
559 of Human Mobile Element Retrotransposition Rates.” *Genome Research* 29 (10): 1567–77.
- 560 Flasch, Diane A., Ángela Macia, Laura Sánchez, Mats Ljungman, Sara R. Heras, José L. García-
561 Pérez, Thomas E. Wilson, and John V. Moran. 2019. “Genome-Wide de Novo L1
562 Retrotransposition Connects Endonuclease Activity with Replication.” *Cell* 177 (4): 837–
563 51.e28.
- 564 Gardner, Eugene J., Vincent K. Lam, Daniel N. Harris, Nelson T. Chuang, Emma C. Scott, W.
565 Stephen Pittard, Ryan E. Mills, 1000 Genomes Project Consortium, and Scott E. Devine.
566 2017. “The Mobile Element Locator Tool (MELT): Population-Scale Mobile Element
567 Discovery and Biology.” *Genome Research* 27 (11): 1916–29.
- 568 Gardner, Eugene J., Elena Prigmore, Giuseppe Gallone, Petr Danecek, Kaitlin E. Samocha, Juliet
569 Handsaker, Sebastian S. Gerety, et al. 2019. “Contribution of Retrotransposition to
570 Developmental Disorders.” *Nature Communications* 10 (1): 4630.
- 571 Hancks, Dustin C., and Haig H. Kazazian Jr. 2010. “SVA Retrotransposons: Evolution and
572 Genetic Instability.” *Seminars in Cancer Biology* 20 (4): 234–45.
- 573 Hancks, Dustin C., and Haig H. Kazazian Jr. 2016. “Roles for Retrotransposon Insertions in
574 Human Disease.” *Mobile DNA* 7 (May): 9.
- 575 Jónsson, Hákon, Patrick Sulem, Birte Kehr, Snaedis Kristmundsdottir, Florian Zink, Eirikur
576 Hjartarson, Marteinn T. Hardarson, et al. 2017. “Parental Influence on Human Germline de
577 Novo Mutations in 1,548 Trios from Iceland.” *Nature* 549 (7673): 519–22.
- 578 Keane, Thomas M., Kim Wong, and David J. Adams. 2013. “RetroSeq: Transposable Element
579 Discovery from next-Generation Sequencing Data.” *Bioinformatics* 29 (3): 389–90.
- 580 Kong, Augustine, Michael L. Frigge, Gisli Masson, Soren Besenbacher, Patrick Sulem, Gisli
581 Magnusson, Sigurjon A. Gudjonsson, et al. 2012. “Rate of de Novo Mutations and the
582 Importance of Father’s Age to Disease Risk.” *Nature* 488 (7412): 471–75.
- 583 Chu, Chong, Eric W. Lin, Antuan Tran, Hu Jin, Natalie I. Ho, Alexander Veit, Isidro Cortes-
584 Ciriano, Kathleen H. Burns, David T. Ting, and Peter J. Park. 2023. “The Landscape of
585 Human SVA Retrotransposons.” *Nucleic Acids Research* 51 (21): 11453–65.
- 586 Kwon, Minseok, Soohyun Lee, Michele Berselli, Chong Chu, and Peter J. Park. 2021.
587 “BamSnap: A Lightweight Viewer for Sequencing Reads in BAM Files.” *Bioinformatics* ,
588 January. <https://doi.org/10.1093/bioinformatics/btaa1101>.
- 589 Lee, Eunjung, Rebecca Iskow, Lixing Yang, Omer Gokcumen, Psalm Haseley, Lovelace J.
590 Luquette 3rd, Jens G. Lohr, et al. 2012. “Landscape of Somatic Retrotransposition in
591 Human Cancers.” *Science* 337 (6097): 967–71.
- 592 Lek, Monkol, Konrad J. Karczewski, Eric V. Minikel, Kaitlin E. Samocha, Eric Banks, Timothy
593 Fennell, Anne H. O’Donnell-Luria, et al. 2016. “Analysis of Protein-Coding Genetic
594 Variation in 60,706 Humans.” *Nature* 536 (7616): 285–91.
- 595 Machiela, Mitchell J., Thomas G. P. Grünwald, Didier Surdez, Stephanie Reynaud, Olivier
596 Mirabeau, Eric Karlins, Rebeca Alba Rubio, et al. 2018. “Genome-Wide Association Study
597 Identifies Multiple New Loci Associated with Ewing Sarcoma Susceptibility.” *Nature*
598 *Communications* 9 (1): 3184.

- 599 Mavros, Chrystal F., Catherine A. Brownstein, Roshni Thyagrajan, Casie A. Genetti, Sahil
600 Tembulkar, Kelsey Graber, Quinn Murphy, et al. 2018. “De Novo Variant of TRRAP in a
601 Patient with Very Early Onset Psychosis in the Context of Non-Verbal Learning Disability
602 and Obsessive-Compulsive Disorder: A Case Report.” *BMC Medical Genetics* 19 (1): 197.
- 603 Nam, Chang Hyun, Jeonghwan Youk, Jeong Yeon Kim, Joonoh Lim, Jung Woo Park, Soo A.
604 Oh, Hyun Jung Lee, et al. 2022. “Extensive Mosaicism by Somatic L1 Retrotransposition in
605 Normal Human Cells.” *bioRxiv*. <https://doi.org/10.1101/2022.05.18.492429>.
- 606 Ng, Jeffrey K., Pankaj Vats, Elyn Fritz-Waters, Stephanie Sarkar, Eleanor I. Sams, Evin M.
607 Padhi, Zachary L. Payne, et al. 2021. “De Novo Variant Calling Identifies Cancer Mutation
608 Profiles in the 1000 Genomes Project.” *bioRxiv*. <https://doi.org/10.1101/2021.05.27.445979>.
- 609 Pottie, Lore, Christin S. Adamo, Aude Beyens, Steffen Lütke, Piyanoot Tapaneeyaphan,
610 Adelbert De Clercq, Phil L. Salmon, et al. 2021. “Bi-Allelic Premature Truncating Variants
611 in LTBP1 Cause Cutis Laxa Syndrome.” *American Journal of Human Genetics* 108 (6):
612 1095–1114.
- 613 Schrider, Daniel R., Fabio C. P. Navarro, Pedro A. F. Galante, Raphael B. Parmigiani, Anamaria
614 A. Camargo, Matthew W. Hahn, and Sandro J. de Souza. 2013. “Gene Copy-Number
615 Polymorphism Caused by Retrotransposition in Humans.” *PLoS Genetics* 9 (1): e1003242.
- 616 Shi, Xianping, Yueyuan Zheng, Liling Jiang, Bo Zhou, Wei Yang, Liyan Li, Lingwen Ding, et
617 al. 2020. “EWS-FLI1 Regulates and Cooperates with Core Regulatory Circuitry in Ewing
618 Sarcoma.” *Nucleic Acids Research* 48 (20): 11434–51.
- 619 Thung, Djie Tjwan, Joep de Ligt, Lisenka E. M. Vissers, Marloes Steehouwer, Mark Kroon,
620 Petra de Vries, Eline P. Slagboom, Kai Ye, Joris A. Veltman, and Jayne Y. Hehir-Kwa.
621 2014. “Mobster: Accurate Detection of Mobile Element Insertions in next Generation
622 Sequencing Data.” *Genome Biology* 15 (10): 488.
- 623 Tubio, Jose M. C., Yilong Li, Young Seok Ju, Inigo Martincorena, Susanna L. Cooke, Marta
624 Tojo, Gunes Gundem, et al. 2014. “Extensive Transduction of Nonrepetitive DNA Mediated
625 by L1 Retrotransposition in Cancer Genomes.” *Science* 345 (6196): 1251343.
- 626 Urban, Zsolt, and Elaine C. Davis. 2014. “Cutis Laxa: Intersection of Elastic Fiber Biogenesis,
627 TGF β Signaling, the Secretory Pathway and Metabolism.” *Matrix Biology: Journal of the*
628 *International Society for Matrix Biology* 33 (January): 16–22.
- 629 Vogt, Julia, Kathrin Bengesser, Kathleen B. M. Claes, Katharina Wimmer, Victor-Felix
630 Mautner, Rick van Minkelen, Eric Legius, et al. 2014. “SVA Retrotransposon Insertion-
631 Associated Deletion Represents a Novel Mutational Mechanism Underlying Large Genomic
632 Copy Number Changes with Non-Recurrent Breakpoints.” *Genome Biology* 15 (6): R80.
- 633 Wang, Hui, Jinchuan Xing, Deepak Grover, Dale J. Hedges, Kyudong Han, Jerilyn A. Walker,
634 and Mark A. Batzer. 2005. “SVA Elements: A Hominid-Specific Retroposon Family.”
635 *Journal of Molecular Biology* 354 (4): 994–1007.
- 636 Werling, Donna M., Harrison Brand, Joon-Yong An, Matthew R. Stone, Lingxue Zhu, Joseph T.
637 Glessner, Ryan L. Collins, et al. 2018. “An Analytical Framework for Whole-Genome
638 Sequence Association Studies and Its Implications for Autism Spectrum Disorder.” *Nature*
639 *Genetics* 50 (5): 727–36.
- 640 Wimmer, Katharina, Tom Callens, Annekatrin Wernstedt, and Ludwine Messiaen. 2011. “The
641 NF1 Gene Contains Hotspots for L1 Endonuclease-Dependent de Novo Insertion.” *PLoS*
642 *Genetics* 7 (11): e1002371.

- 643 Xia, Wenjun, Jiongiong Hu, Jing Ma, Jianbo Huang, Xu Wang, Nan Jiang, Jin Zhang, Zhaoxin
644 Ma, and Duan Ma. 2019. “Novel TRRAP Mutation Causes Autosomal Dominant Non-
645 Syndromic Hearing Loss.” *Clinical Genetics* 96 (4): 300–308.
- 646 Zhang, Qiang, Zailong Qin, Shang Yi, Hao Wei, Xun Zhao Zhou, and Jiasun Su. 2020. “Two
647 Novel Compound Heterozygous Variants of LTBP4 in a Chinese Infant with Cutis Laxa
648 Type IC and a Review of the Related Literature.” *BMC Medical Genomics* 13 (1): 183.
- 649 Zhuang, Jiali, Jie Wang, William Theurkauf, and Zhiping Weng. 2014. “TEMP: A
650 Computational Method for Analyzing Transposable Element Polymorphism in
651 Populations.” *Nucleic Acids Research* 42 (11): 6826–38.
- 652 Deciphering Developmental Disorders Study. 2017. “Prevalence and Architecture of de Novo
653 Mutations in Developmental Disorders.” *Nature* 542 (7642): 433–38.
- 654 Wallace, M. R., L. B. Andersen, A. M. Saulino, P. E. Gregory, T. W. Glover, and F. S. Collins.
655 1991. “A de Novo Alu Insertion Results in Neurofibromatosis Type 1.” *Nature* 353 (6347):
656 864–66.
657

658 **Acknowledgements**

659 This work was supported by R03CA249364 from the National Cancer Institute. VL is supported
660 by the Swedish Research Council (2020-00583).

661

662 **Competing interests**

663 The authors declare no competing interests.

664

665