1
2
3
4

# Supplementary Materials for

Paternally inherited cis-regulatory structural variants contribute to autism

William M Brandler, Danny Antaki, Madhusudan Gujral, Morgan L Kleiber, Joe Whitney, Michelle S Maile, Oanh Hong, Timothy R Chapman, Shirley Tan, Prateek Tandon, Timothy Pang, Shih C Tang, Keith K Vaux, Yan Yang, Eoghan Harrington, Sissel Juul, Daniel J Turner, Bhooma Thiruvahindrapuram, Gaganjot Kaur, Stephen F Kingsmore, Joseph G Gleeson, Denis Bisson, Boyko Kakaradov, Amalio Telenti, J Craig Venter, Roser Corominas, Claudio Toma, Bru Cormand, Isabel Rueda, Silvina Guijarro, Karen S Messer, Caroline M Nievergelt, Maria J Arranz, Eric Courchesne, Karen Pierce, Alysson R Muotri, Lilia M Iakoucheva, Amaia Hervas, Stephen W Scherer, Christina Corsello & Jonathan Sebat
correspondence to: jsebat@ucsd.edu

**This PDF file includes:**

     Materials and Methods
     Figs. S1 to S12
     Captions for Tables S1 to S11

**Other Supplementary Materials for this manuscript includes the following:**

     Tables S1 to S11

## Materials and Methods
### Study design

A key challenge in the analysis of large biological datasets is to account for the full array

of hypotheses that are tested during the process of data collection, data filtering,

annotation and statistical analysis. Analysis choices that are made throughout the process

are influenced by properties of the data. Thus, correcting for all of the formal statistical

tests that are performed in a study may not fully account for the "garden of forking paths"

that led to the formulation of these hypotheses (*31*). It is therefore difficult for the reader

to know for certain if the multiple test correction performed accounts for the full

hypothesis space that could potentially be explored.

This is particularly problematic when investigating genetic association in non-coding

regions of the genome. The analyst has almost infinite degrees of freedom in terms of the

selection of functional annotations and gene-sets that could potentially be tested for

association. The solution to this challenge, as proposed by Gelman and Loken (*31*) is to

pair an initial experiment with a "pre-registered" replication, in which the hypotheses and

all details of analysis are specified in advance. Our study has followed this design, and it

is structured in three stages:

1. Target functional elements were selected from a larger set of annotations based

   on evidence of SV intolerance from this study and from a SV call set from the

   1000 genomes project.

2. Target categories were tested for association in a primary sample of 829 ASD

   families. One category of non-coding annotation was significant after correction for

56    multiple testing. This association and all details of the analysis were then posted as a

57    manuscript to the preprint server bioRxiv (https://www.biorxiv.org/).

58    3. The primary hypothesis was subsequently replicated in an independent sample of

59    1,771 families

60    Stage 1 of this study provides an effective means for reducing the number of tests to a

61    limited set of functional annotations in which SVs are under strong natural selection. The

62    prepublication manuscript posted in stage 2 provides a transparent way to state our

63    primary hypothesis and describe all analysis methods prior to obtaining the replication

64    dataset. The addition of stage 3, prompted by peer review of the primary study, allows for

65    confirmation of the primary scientific claim.

66    **<u>Recruitment</u>**

67    Our discovery sample consisted of ASD families derived from two cohorts, which will be

68    referred to as 'REACH' or 'SSC1' in the following sections.

69    REACH cohort individuals were referred from clinical departments at Rady Children's

70    Hospital, including the Autism Discovery Institute, Psychiatry, Neurology, Speech and

71    Occupational Therapy and the Developmental Evaluation Clinic (DEC) as part of the

72    Relating genes to Adolescent and Child Health (REACH) study. Further referrals came

73    directly through the REACH project website (http://reachproject.ucsd.edu/). In total 612

74    individuals from 161 families came from the REACH project. The Autism Center of

75    Excellence at the University of California San Diego contributed 11 trios. A further 452

76    samples from 139 families were recruited at Hospital Universitari Mútua de Terrassa in

77    Barcelona. The REACH families combined consisted of 112 controls and 362 affected

78    individuals - 285 with ASD, 43 with pervasive developmental disorder – not otherwise

79    specified (PDD-NOS), 10 with attention deficit hyperactivity disorder (ADHD), and 24

80    with speech delay, epilepsy, anxiety, or other related developmental disorders that were

81    therefore classified as 'cases' for bioinformatics analyses.

82    The Simons Simplex Collection Phase 1 (SSC1) Whole Genome Sequencing dataset

83    (http://bit.ly/2jc34rU) consisted of 518 quad families with sibling pairs discordant for an

84    ASD diagnosis that were selected from the full cohort of 2,644 families after excluding

85    those where offspring carried any plausible contributory *de novo* or inherited SNVs,

86    indels, deletion or duplications identified from microarray or exome sequencing data. The

87    exclusion criteria for exome- or array-'positive' individuals are described below and were

88    applied to both ASD cases and sibling-controls:

89    1) *De novo* deletions and duplications (189 families): Any confirmed/published *de novo*

90    copy number variant (CNV) (*10*), Illumina SNP genotyping data, or exome CNV data

91    that is: Rare (≤0.1 population frequency based on parents and DGV) or genic (≥1 exon).

92    2) Inherited CNVs (92 families): Any CNV from Illumina genotyping data, or exome

93    CNV data that is: rare (≤0.1 population frequency based on parents and DGV), or

94    intersects ≥10 genes.

95    3) *De novo* LoF (564 families): Any *de novo* loss of function from published sequencing

96    data that is: rare (≤0.1 population frequency based on the exome variant server),

97    nonsense, canonical splice site, or frameshift (*2, 26*).

98

99

**Whole Genome Sequencing**

100

101    Our combined dataset consisted of WGS data collected for two cohorts and sequenced at

102    three sites (**table S1**). All WGS data were generated from whole blood DNA. All

103    members of individual families were sequenced within the same batch of samples.

104    **REACH cohort**

105    Whole genome sequencing was performed on blood-derived genomic DNA samples of

106    1,126 individuals from 319 families, including 893 individuals from 260 families.

107    Sequencing was performed at Human Longevity Inc. (HLI) on an Illumina HiSeq X10

108    system (150 bp paired ends at mean coverage of 50X) and an additional 204 individuals

109    from 59 families that were sequenced at the Illumina FastTrack service laboratory on the

110    Illumina HiSeq 2500 platform as described in our previous publication (*9*). We

111    performed initial quality control (QC) steps to ensure relatedness and gender matched the

112    sample sheets, excluding any mismatches or half-siblings. We also tested for an excess of

113    Mendelian errors in the children, and an excess of single nucleotide variants called in

114    either parent (>3 SD from the mean) indicative of low quality DNA. In total 29 samples

115    were removed, including eight complete families. Therefore, 1,097 individuals from 311

116    families were taken forward for structural variant calling and analysis.

117    **SSC1 Cohort**

118    Whole genome sequencing of the SSC phase 1 (SSC1) cohort of 540 families was

119    performed at the New York Genome Center on an Illumina HiSeq X10 (150 bp paired

120    ends at mean coverage of 40X). Of the 540 SSC families, 518 were complete quad

121    families. Incomplete families were excluded from the dataset. All 518 met the above QC

122    criteria for inclusion in the study. Mean coverage (39-50X) and insert sizes (348-420) and

123    were similar at all three sequencing sites (**table S1**).

124    Sequence alignment and variant calls for REACH samples were generated on families

125    using our WGS analysis pipeline implemented on the Comet compute cluster at the San

126    Diego Supercomputer Center (SDSC, https://goo.gl/C4bVoe). For SSC samples the same

127    pipeline was adapted for use on Amazon Web Services (AWS). In brief, short reads were

128    mapped to the hg19 reference genome by BWA-mem version 0.7.12. Subsequent

129    processing was carried out using SAMtools version 1.2, GATK version 3.3, and Picard

130    tools version 1.129, which consisted of the following steps: sorting and merging of the

131    BAM files, indel realignment, removal of duplicate reads, base quality score recalibration

132    for each individual.

133    **Replication Cohorts**

134    Our hypothesis and all analytic details were pre-registered by posting a preprint

135    describing the results of our primarily analysis (*16*). We then carried out a replication of

136    our primary scientific claim in an independent sample.

137    The replication WGS dataset consists of data from two cohorts, the Autism Speaks'

138    MSSNG program (Principal investigator: S.S.) (*17*), and the SSC phase 2 (SSC2) sample.

139    The MSSNG sample consisted of 30X WGS of 3,769 individuals using Illumina HiSeq

140    X10 platform, including 1,395 ASD cases from 1,187 families (998 trios, 157 quads, 28

141    quintets, 3 sextets, and 1 septet concordant for ASD). A complete breakdown and list of

142    samples is provided in **table S1**. The SSC2 cohort consisted of 2,336 individuals from

143    584 quads discordant for ASD, sequenced and processed in the same way as the first

144      phase of SSC quads. In total the replication cohorts consisted of 6,105 individuals from

145      1,771 families, including 1,979 ASD cases and 584 sibling controls.

146      **SV detection, genotyping and filtering – discovery cohort**

147      We utilized four complementary algorithms to detect SVs: ForestSV, Lumpy, Manta, and

148      Mobster. ForestSV is designed to detect deletions and duplications based on a

149      combination of signatures including, coverage, discordant paired ends and other metrics

150      such as mapping quality (http://sebatlab.ucsd.edu/software-data). Lumpy

151      (https://github.com/arq5x/lumpy-sv) and Manta  (https://github.com/Illumina/manta;

152      workflow version 0.29.0), the latter being a new addition to the SV analysis pipeline

153      since our previous publication (*9*), both utilize a combination of discordant paired ends

154      and split reads and have greater sensitivity for small (<500 bp) deletions, duplications,

155      inversions and complex rearrangements. Finally, Mobster

156      (http://sourceforge.net/projects/mobster) uses discordant paired-end and split-read signal

157      in combination with consensus sequences of known active transposable elements to

158      identify mobile element insertions (MEIs). A consensus callset was generated by merging

159      calls from ForestSV, Lumpy, Manta and Mobster. SV calls from multiple methods were

160      combined, and overlapping variants detected in the same sample were collapsed as

161      described in our previous structural variant publication (*9*). The unfiltered consensus

162      callset consisted of the union of calls from the four methods. As a preliminary filtering

163      step, SVs were removed from the consensus callset if they overlapped by more than 66%

164      with centromeres, segmental duplications, regions with low mappability with 100bp

165      reads, regions subject to somatic V(D)J recombination (parts of anitbodies and T-cell

166      receptor genes). SVs called by Manta or Lumpy were filtered if they had one or both

167    breakpoints overlapping one of these regions. Regions used for filtering can be found in

168    our previous publication (*9*).

169    We generated a set of uniformly-called genotypes for the combined set of deletions and

170    duplications detected by three methods Lumpy, Manta, or ForestSV, using a single

171    genotyping algorithm $SV^2$ (https://github.com/dantaki/SV2). $SV^2$ (*11*) provides estimates

172    of genotype likelihoods for deletions and duplications across a broad size range (10bp-

173    10Mb), and this metric was used as our primary filtering criterion for these. Lumpy and

174    Manta also provide genotype likelihoods for the subset of calls that were generated by

175    these methods, which include SVs that are not genotyped by $SV^2$ such as inversions and

176    non-tandem duplications. These genotype likelihoods were also used as quality metrics

177    during the filtering of SV callset as described below.

178    $SV^2$ designates SV calls as "PASS" or "FAIL" at two levels of stringency: "standard" and

179    "de novo", which are described in detail in our companion paper (*11*).  Standard filters

180    were used to generate the main SV callset and these genotypes were used for family

181    based association testing. The more stringent "de novo" filters were used for de novo

182    mutation calling. In addition, we included in the consensus callset SVs, which passed

183    genotype likelihood thresholds for Lumpy and Manta. Manta and Lumpy genotype-

184    likelihood thresholds for SV filtering were determined based on estimates of FDR, which

185    were performed from Illumina 2.5M SNP array data on a subset of 205 genomes using

186    the Intensity Rank Sum test implemented using the Structural Variation Toolkit.

187    Thresholds were selected for SVs across a range of sizes and depending on sequence

188    context (short tandem repeats, segmental duplications, etc.). FDR estimates for SV calls

189   filtered at standard and de novo stringency and genotype likelihood thresholds for Lumpy

190   and Manta are provided in **table S3**.

191   Due to the requirements of this study for high genotyping accuracy, we have applied

192   additional filtering measures that were not used in a previous publication from our group

193   (*9*). The FDR of variants intersecting STRs was an order of magnitude higher than SVs

194   that did not; therefore more stringent genotype likelihood filters were applied to SVs

195   overlapping STRs (**table S3**). Furthermore because STRs were depleted in probes on the

196   Illumina 2.5M SNP array, only 7.2% of deletions and 12.9% of duplications overlapping

197   an STR had one or more probes, compared to 28.5% of deletions and 56.3% of

198   duplications that do not. FDR estimates for these variants could be less accurate.

199   Therefore, for all analyses in this study, we have excluded SVs with breakpoints

200   overlapping STRs. We have also annotated these in the callset VCF (which can be

201   downloaded from NDAR study #434), and we suggest that these SVs be treated with

202   caution. Hence, the number of deletions and duplications reported in the SV callset here

203   is lower than in our previous publication (*8, 9*).

204   In total we detected 11.87 million alleles from 89,123 distinct loci encompassing 19.4%

205   of the GRCh37 (hg19) release of the 'mappable' reference human genome

206   (0.497/2.57Gb, excluding SVs larger than 1Mb, which are likely to be pathogenic and

207   would contribute disproportionately to this estimate, **table S2**). 12.5% (320Mb) of the

208   reference genome was deleted and 7.3% (186Mb) duplicated in our cohort of 829

209   families.

210

211

212  **SV detection, genotyping and filtering – replication cohort**

213  **MSSNG**

214  Data processing was performed by Scherer laboratory, and functional annotation of SV

215  calls was performed using an annotation file that we provided.  Briefly, for 2,945

216  individuals alignment was performed using BWA version 0.7.10. SV calling was

217  performed on a per family basis using Manta and Lumpy, with genotyping using $SV^2$

218  following the pre-registered protocol (described above). For a subset of individuals (n =

219  824) sequence alignment was performed with the ISAAC aligner and SV calling was

220  performed by Manta on a per individual basis, but with genotyping of each SV call on a

221  per family basis using $SV^2$.

222  **SSC2**

223  The SSC phase 2 (SSC2) data was processed on the Amazon Web Services cloud in a

224  manner identical to that for SSC1. SV genotypes from the replication dataset were

225  intersected with our original SV callset based on the confidence intervals for SV

226  breakpoints given by Manta / Lumpy. We then identified SVs that had an allele

227  frequency <0.0003 (the allele frequency for private variants in our original study).

228  *De novo* **calling and phasing**

229  *De novo* SVs were called if they occurred in a child and were genotyped reference in both

230  parents and the parent allele frequency for the variant was less than 1%. We also applied

231  more stringent $SV^2$ genotype likelihood filters for *de novo* SVs and TDT analyses, which

232  are detailed in **table S3**. The average rate of Mendelian errors for deletions and

233  duplications in the callset as a whole was 0.99% (95% CI: 0.03) and 4.66% (95% CI:

10

234   0.15) respectively (**fig. S4**). *De novo* genotype likelihood filters applied to variants with

235   parent allele frequencies <1% reduced the rate to 0.21% (95% CI: 0.1) for deletions and

236   0.5% (95% CI: 0.2) for duplications.

237   ### **SV validation**

238   We validated large putative *de novo* deletions and duplications using an in silico SNP-

239   based approach that utilizes read depth from the VCF files from GATK Haplotype Caller.

240   For each SNP we normalized allelic read depth relative to the genome average for

241   reference / alternate alleles, and calculated a z-score for each SNP. We also calculated the

242   B allele frequency (BAF) by taking the average of the allele (reference or alternate) with

243   the fewest number of supporting reads across the locus. Since deletions are hemizygous,

244   the expected BAF is 0. For duplications we calculated the BAF only for heterozygote

245   SNPs, which have an expected BAF of 0.33 for autosomal variants. If the child showed

246   an average elevated or depleted SNP read depth more than one standard deviation from

247   both parents, and a BAF consistent with the SV type, and/or the variant could be phased,

248   then the SV was designated as valid. Furthermore this SNP data was used to determine

249   the parent of origin, by performing a paired t-test on phased SNP allelic depth within the

250   locus. We plotted the validation results for each member of the trio using the R package

251   CNVplot, which was developed in house (https://github.com/dantaki/CNVplot). This

252   approach is orthogonal to the SV calling steps above, which do not phase variants,

253   calculate their BAF, or estimate coverage using SNP data.

254   Small deletions, duplications, inversions, complex SVs, and MEIs were validated using

255   PCR. Both *de novo* inversion calls were validated. We attempted PCR validation on 13

256   *de novo Alu* elements, all of which validated as *de novo*. *Alu* insertions have poly-A tails;

257    we therefore used a lower extension temperature (65ºC), because A/T rich sequences

258    have a low melting temperature. We also used longer extension times (90 seconds) to an

259    otherwise standard PCR protocol.

260    **<u>Validation and FDR estimation by Nanopore sequencing</u>**

261    We validated deletions and duplications predicted in Illumina short read paired-end

262    genomes in three unrelated individuals with Oxford Nanopore (ONP) long read

263    sequencing. ONP reads were aligned to the human hg19 reference with bwa-mem

264    (version 0.7.10-r789) and ngmlr (https://github.com/philres/ngmlr, version 0.2.3) with the

265    "-x ont2d" and "-x ont" options.  The average coverage was 7.4X and average read length

266    was 2,574bp for bwa-mem alignments and 7.3X and 2,525bp for ngmlr alignments. We

267    restricted validation to variants with less than 50% overlap to elements in our genome

268    mask. Additionally, we ensured that the median base-pair depth of coverage was greater

269    than 0X in 1000bp regions flanking the breakpoints, totaling 3,252 deletion and 62

270    duplication candidates for validation. We then searched for supporting reads in bwa-mem

271    and ngmlr alignments, defined as supplementary alignments or CIGAR string deletions

272    and insertions with breakpoints that overlap at least 50% reciprocally to the SV in

273    question. Short-read SV predictions were considered validated if at least 1 supporting

274    read was detected in either bwa-mem or ngmlr alignments. We then calculated the false

275    discovery rate (FDR) specifying false positives as SVs without supporting reads while

276    binning on allele frequency and SV length. The overall FDR was 10.4% for deletions and

277    30.6% for duplications; for private variants of SV length 100bp-1000bp the FDR was 0%

278    for deletions.

279

**Oxford Nanopore Targeted Validation of *LEO1***

281     Recurrent deletions of the *LEO1* locus were validated and fine mapped by single

282     molecule sequencing. Deletion and reference haplotype sequences were amplified by

283     long range PCR (LongAmp® Taq 2X Master Mix, New England BioLabs, M0287L) in

284     carriers of *LEO1* deletions from three families (14-59, F0182, and F0208). Target

285     sequences were amplified from each sample using one set of primers that span the

286     deletion breakpoint and another set that specifically amplifies the reference (non-

287     deletion) haplotype, and PCR amplicons were gel purified. Samples were barcoded using

288     Oxford Nanopore Technologies' (ONT) Native Barcoding Kit 1D (EXP-NBD103) and

289     sequencing adapters were added using Ligation Sequencing Kit 1D (SQK-LSK108).

290     Libraries were sequenced for 48 hours on ONT's MinION Mk1B, using the SpotON

291     Flow Cell Mk I (R9.4, FLO-SPOTR9) and MinKNOW software (v.1.3.30). In total

292     approximately 2.3 Gb of fasta data was generated.

293     Quality and length filters were applied to the unaligned reads. Reads with a mean quality

294     score of 8.5 or less or which differed from the expected amplicon length by 2kb or more

295     were removed. Reads were aligned to the human genome (hg19) using BWA-mem

296     (v.0.7.15-r1140) with the '-x ont2d -M' flags and filtered to keep only those that

297     overlapped 95% of the target region. Consensus sequences were generated from the

298     alignment of multiple reads using Mummer (http://mummer.sourceforge.net/), and

299     deletion breakpoints were identified by aligning the consensus sequence to the reference

300     genome using BLAT. The consensus fasta sequences can be downloaded from NDAR.

301

302

**Evaluation of SV calling across data from multiple sequencing centers**

The average SV numbers for each class of SV were similar between cohorts sequenced at

different sequencing centers (**table S1**). Modest differences in SV calling were observed

between sequencing centers. We compared SV calls for one individual (REACH000236)

who was sequenced twice, on the Illumina HiSeq 2500 with 100bp reads (at 43X

coverage) and on the Illumina HiSeq X with 150bp reads (also at 43X coverage). Since

the coverage is the same between the two samples but the read length is 50% longer on

the HiSeq X, this sample has only 2/3 as many reads when sequenced on the HiSeq X.

This affects SV calling for two reasons, there will be on average more split reads

supporting each call on the HiSeq X, but fewer discordant paired-end reads. The overlap

between the SVs called on each platform in this sample ranged from 66-96% for each SV

type (**fig. S12**).

**Selection of target functional elements based on SV intolerance**

We investigated the enrichment/depletion of private deletions, duplications, and mobile

element insertions within specific genomic features compared to a random distribution of

SVs. Random distributions of SVs were simulated using two different models of random

mutation: (1) a uniform random model (UM) in which we shuffled the position of sites

that were private to families (i.e. observed in only one parent) across the genome using

BedTools  and (2) a non-uniform random model (NUM) based on a concept used

previously (*15*), in which the correlation of SVs to genome features was modeled by

fitting a linear regression to the observed rate of SV breakpoints to GC content, coverage,

low-complexity repetitive elements, and segmental duplications. A probability density

function derived from the linear model was then used to simulate random SVs. In both

326      cases we excluded regions of the genome that cannot be sequenced with short reads. We

327      counted the number of times a shuffled SV overlapped (at least 1bp) the following

328      genomic features: protein coding exons, transcription start sites (TSS), 5'UTRs, 3'UTRs,

329      promoters, noncoding RNAs, enhancers, conserved noncoding regions, human

330      accelerated regions, CTCF binding sites, exon flanking (one breakpoint within 100bp of

331      an exon), 1kb upstream, 1kb downstream, and introns. Events that overlapped multiple

332      features were prioritized in the order above, so for example if a variant overlapped a

333      protein coding exon, a 3'UTR and an intron, it is counted as protein coding but not

334      3'UTR or intronic. Each feature is explained in detail below and we've summarized each

335      in a table included as part of **table S5**. We performed 10,000 permutations and compared

336      the observed overlap to the expected overlap. P values were corrected using a Benjamini–

337      Hochberg false-discovery rate adjustment, and Q values are reported in **table S5**.

338      Categories that were depleted among variant-intolerant genes (ExAC pLI>90$^{th}$ percentile)

339      were selected as targets in our primary analysis. Significant depletion was defined as

340      OR<1 and FDR adjusted Q<0.01.

341      **Generating a random distribution of SVs using a linear regression model**

342      Structural mutation rates vary across the genome, and regional differences in the rate of

343      SVs introduce biases in the distribution of SVs that could confound our estimates of SV

344      intolerance. To address this concern, we adapted a model from Ruderfer et al. (*15*) to

345      estimate SV mutation rate and to simulate variants according to a non-uniform,

346      empirically-derived random distribution that is more reflective of true genomic

347      background than assuming uniform random mutation. Our NUM model fits a linear

348      regression for the observed rate of SV breakpoints (in 1000 bp windows) in relation to

349    GC content, coverage, overlap with low complexity repetitive elements, and intersection

350    with segmental duplication regions.

351    Coverage tracks were generated from SSC, 1000 genomes, and REACH samples. At least

352    30 samples were used to generate each track. Fine-grain GC content, repetitive element,

353    and segmental duplication overlap tracks were generated from raw data available on

354    UCSC genome browser. These tracks were then used to fit three linear regression models

355    to predict the empirical density of SV breakpoints in each data set for benign variants.

356    These linear regression models were then converted into probability density functions

357    (pdfs) that could be used to simulate new background variants.

358    10,000 simulations were performed to shuffle the variants in each data set. Our empirical

359    pathogenic predictions were compared against the generated null distributions. The

360    results did not differ significantly after correcting for SV background mutation rates

361    according to the Ruderfer model.

362    **<u>Definitions of gene disrupting SVs versus noncoding</u>**

363    Gene disrupting deletions were defined as those that directly disrupted at least one

364    protein coding exon from one transcript of a gene (transcripts were extracted from hg19

365    RefSeq). Noncoding deletions could delete UTRs, introns, enhancers, or promoters of

366    genes, but not protein coding exonic sequence or the start position of the first exon of a

367    transcript. Protein coding duplications were divided into four categories. Whole gene

368    duplications encompassed at least one full-length transcript of a gene. Internal exon

369    duplications intersected at least one protein coding exon internal to a transcript, but not

370    the UTRs. Duplications that intersected at least one exon and with one breakpoint outside

371    of the gene and the other internal to the gene were divided into two categories, those that

372     encompassed the 5'UTR (and promoter), and those that encompassed the 3'UTR. Gene

373     disrupting inversions were classified as variants that either had one or both breakpoints

374     inside a protein coding exon of a gene, or that had one breakpoint in an intron of a gene

375     and the other breakpoint either outside of that gene or in another intron. Inversions that

376     inverted an entire gene or genes but had intergenic breakpoints were considered

377     noncoding.

378     **<u>Definition and selection of noncoding elements</u>**

379     Transcription start sites, 3'UTRs, and 5'UTRs were defined using full-length protein-

380     coding transcripts from RefSeq. Two types of noncoding RNAs, micro-RNAs and natural

381     antisense transcripts were defined. Human micro-RNAs were downloaded from miRBase

382     (mirbase.org, v21), lifted over to hg19 annotated to genes if they were intronic in a sense

383     orientation and therefore transcribed with the gene itself. Exons of natural antisense

384     transcripts (NATs) were assigned to genes if they were transcribed in an antisense

385     direction and overlapped with a gene. NAT data was downloaded from GENCODE v25

386     (only including transcripts with support level of 1, 2 or 3).

387     Conserved noncoding regions were defined from two studies; one that defined

388     ultraconserved elements > 100bp conserved in human, mouse and rat genomes (*32*), and

389     the other that defined ultrasensitive noncoding regions with almost as much selective

390     constraint as coding genes (*33*).

391     Promoters and enhancers were defined using fetal brain data Epigenomics Roadmap

392     Project and data from ENCODE. The Epigenomics Roadmap Project integrated

393     combinatorial interactions between five different chromatin marks to define 15 chromatin

394    states using a Hidden Markov Model

395    (http://egg2.wustl.edu/roadmap/web_portal/chr_state_learning.html).

396    Four states were used to define promoters, active transcription start site (1_TssA), TSS

397    flank (2_TssAFlnk), bivalent TSS (10_TssBiv), and bivalent TSS flank (11_BivFlnk).

398    Three states were used to define fetal brain enhancers, genic enhancer (6_EnhG),

399    enhancer (7_Enh), and bivalent enhancer (12_EnhBiv).

400    For the Epigenomics Roadmap Project data, fetal brain promoters/enhancers were

401    defined using the intersection of male and female fetal brain tissue (epigenomes: E081 &

402    E082). Adult brain promoters/enhancers were defined using the intersection of

403    epigenomes from eight brain regions (E067 (Angular gyrus), E068 (Anterior Caudate),

404    E069 (Cingulate Gyrus), E070 (Germinal Matrix), E071 (Hippocampus), E071 (Inferior

405    Temporal Lobe), E073 (Dorsolateral Prefrontal Cortex), & E074 (Substantia Nigra)),

406    excluding any elements that intersected with those in fetal brain.

407    ENCODE enhancers and promoters were defined based on chromatin state segmentations

408    from six human cell lines (GM12878, K562, H1-hESC, HeLa-S3, HepG2, and HUVEC),

409    which integrated ENCODE ChIP-seq, DNase-seq, and FAIRE-seq data from two

410    algorithms (chromHMM and Segway) to segment the genome into seven states. Data for

411    all six cell types was downloaded from UCSC genome browser, two states were used to

412    defined ENCODE promoters, predicted promoter or transcription start site (state: TSS),

413    predicted promoter flanking region (state: PF). One state was used to define ENCODE

414    enhancers, predicted strong enhancer (State: E). ENCODE CTCF enriched elements were

415    used to define CTCF binding sites (State: CTCF). Promoters and Enhancers were

416    assigned to genes based on proximity, if they intersected or were within 10kb of the

417    transcription start site of an isoform of the gene.

418    Assigning enhancers to genes based purely on proximity is not the most effective

419    approach, as the majority of annotated enhancers do not interact with the nearest gene.

420    We therefore implemented TargetFinder (https://github.com/shwhalen/targetfinder), a

421    machine-learning algorithm that annotates to genes with an FDR <15% by integrating

422    features such as DNA methylation, histone marks, and cap analysis of gene expression

423    (CAGE) data to predict distal enhancers (distance 10kb-2Mb) that interact with

424    promoters. We extracted all enhancers predicted to directly activate genes in six cell

425    types from ENCODE (GM12878, HeLa-S3, HUVEC, IMR90, K562, & NHEK). We also

426    attempted to assign enhancers to genes using the correlation of expression between

427    enhancers and promoters within 500kb of each other using data from FANTOM5

428    (http://fantom.gsc.riken.jp/data/).

429    We downloaded chromatin interaction analysis by paired-end tag (ChIA-PET) data

430    detailing the interactome map between noncoding elements and transcription start sites

431    through CTCF or RNA polymerase II interactions (*21, 22*). For each interacting pair of

432    elements if one member of the pair overlapped a promoter of a gene (within 10kb) we

433    assigned its pair to the target gene as a putative noncoding interacting element.

434    Finally fetal central nervous system DNase hypersensitivity data (*6*) and 'human

435    accelerated regions' that have undergone rapid evolution since the split from

436    chimpanzees (*5*) were also tested. Both these features were assigned to genes based on

437    proximity as for enhancers and promoters.

438

### Defining variant-intolerant genes and annotating known ASD genes

Genes were categorized based on their probability of being loss-of-function (LoF) intolerant (pLI) as assessed by large-scale exome sequencing of populations by the Exome Aggregation consortium (ExAC) (*12*). The EXAC release 0.3.1 dataset (January 2016) was downloaded, and we used the published pLI scores that were calculated on the subset of the cohort after excluding individuals with schizophrenia. The pLI score ranges from 0-1 for 18,421 genes, with higher scores indicating that a gene is more intolerant to inactivating mutations.

Our set of known autism genes were taken from the integration of ASD array data and exome sequencing of the SSC cohort (*10*), and genes with an FDR < 0.1 from another large scale whole exome sequencing study (*18*). In total there are 71 ASD associated genes.

### Transmission Disequilibrium Test

Family-based association tests were performed using $SV^2$ genotype calls for SVs filtered at standard stringency. We tested whether variants private to families in our callset were transmitted to affected children or controls more or less than expected by chance, using a two-tailed haplotype-based group-wise transmission disequilibrium test (gTDT) (*34*), assuming a dominant model. Variants smaller than 100bp or overlapping STRs (>50%) were excluded as it is challenging to validate them or estimate their FDR. We further excluded two families from this analysis, one family where the parents DNA was cell line derived (MT_121), and one family where the mother and child had an excess of coverage based calls from ForestSV (F0226).

461    Our analysis focused on genes with pLI scores $>= 90^{th}$ percentile, which we determined

462    are enriched for genes associated with autism from published exome studies. We also

463    only tested features that were SV intolerant from the callset permutation analyses above

464    as we hypothesize that these features will be enriched for variants associated with autism.

465    P values were corrected for multiple testing using a Benjamini–Hochberg false-discovery

466    rate adjustment.

467    To compare paternal and maternal transmission rates to cases we performed a binomial

468    test under the assumption that 50% of transmitted variants should derive from each

469    parent.

470    **Considering potential biases or technical artifacts in the TDT**

471    The transmission disequilibrium test requires accurate genotyping of variants.

472    Genotyping error can result in the apparent biased transmission of parental variants to

473    offspring. For example false-positive SV calls in parents or false negative genotype calls

474    in children can lead to an apparent under-transmission bias. For instance, given an FDR

475    of 2% for SV calls in parents, and no transmission of the false calls, a rate of 48%

476    transmission would be consistent with random segregation. This modest under-

477    transmission bias, is not specific to SVs, and is also apparent for single nucleotide

478    variants genotyped using GATK (*34*).

479    We have therefore evaluated the potential for genotyping error to lead to spurious results

480    in the TDT as part of a companion study (*11*) and in this study, we further examined the

481    rates of Mendelian error and transmission to offspring for private SVs across a broad size

482    range (**fig. S4**). Our results suggest that private >100 bp deletions and duplications

483    respectively have low FDR (2.3% and 1.7%) and Mendelian error rates (2.0% and 0.6%).

484    Since only a small fraction (2.7%) of SVs <100bp in length overlapped with probes on

485    the Illumina 2.5M SNP microarray we could not accurately estimate the FDR for these;

486    therefore SVs <100bp in size were not included in our analysis.

487    As an additional control in the TDT we also demonstrate that there is no transmission

488    bias for SVs in a non-depleted control category (intronic), which has a similar length

489    distribution (mean = 1,988 bp) to the cis-regulatory category (mean = 2,920 bp). We also

490    observe 50% transmission in tolerant genes for all functional categories of private SVs

491    that were tested (**table S6**). We are therefore able to rule out a systematic transmission

492    bias as an explanation for our results. Lastly, over-transmission of private coding and

493    noncoding SVs was specific to cases, not observed in controls, and the association was

494    replicated in an independent cohort.

495    **Test for enrichment of recurrent SVs in cases**

496    To permute the relative enrichment / depletion of SVs overlapping the same functional

497    elements (e.g. exons) in different families, we permuted these variants across the genome

498    ensuring that permuted variants intersected at least one functional element of a gene with

499    a pLI score >= 90$^{th}$ percentile using bedtools shuffle (by implementing the –incl

500    command). Variants could overlap because of an elevated mutation rate We excluded

501    variants that overlapped a functional element that was also overlapped by a variant from

502    the 1000 Genomes phase 3 SV callset, or that overlapped ≥50% with a 1000 Genomes

503    variant, to exclude variants that may reside in hotspots for structural mutation. We

504    repeated the analysis for controls and for genes with pLI scores <90$^{th}$ percentile. For

505    analysis of coding variants we required that observed / permuted variants impacted any

506    exon of the same gene to be considered recurrent. For noncoding analysis we required

507  that variants impacted the same element (e.g. a 5'UTR from the same transcript) to be

508  considered recurrent. We counted the number of times we observed a gene or functional

509  element was intersected by more than one distinct SV and compared this to 10,000

510  permutations.

511  **Testing the association of *LEO1 de novo* mutations with ASD and DD**

512  A series of 20 different studies have been published that reported all *de novo* mutations

513  detected across the exome in cases. For a specific candidate locus in this study we have

514  investigated the potential association with developmental disorders base on tests of *de*

515  *novo* mutation burden in a large combined sample of 13,391 subjects.

516  **SV Burden**

517  The burden of *de novo* structural variants between individuals with ASD in this study and

518  the controls from this study was assessed using a case-control permutation test

519  implemented in PLINK.

520  **Mutational Clustering**

521  To assess whether *de novo* SVs cluster with *de novo* nucleotide substitutions or indels,

522  we used a window based permutation approach. We took windows of 100bp, 1kb, 10kb,

523  100kb, 1Mb, and 10Mb around the breakpoints of *de novo* SVs and intersected the

524  windows with *de novo* SNVs and indels in the same individuals (*de novo* detection of

525  SNVs and indels was performed as described in our previous publication (*9*). We then

526  shuffled the position of these windows in the genome either randomly (excluding regions

527  that were filtered during SV calling) or across detected inherited SV breakpoints using

528    BedTools and calculated the expected number of window overlapping DNMs using

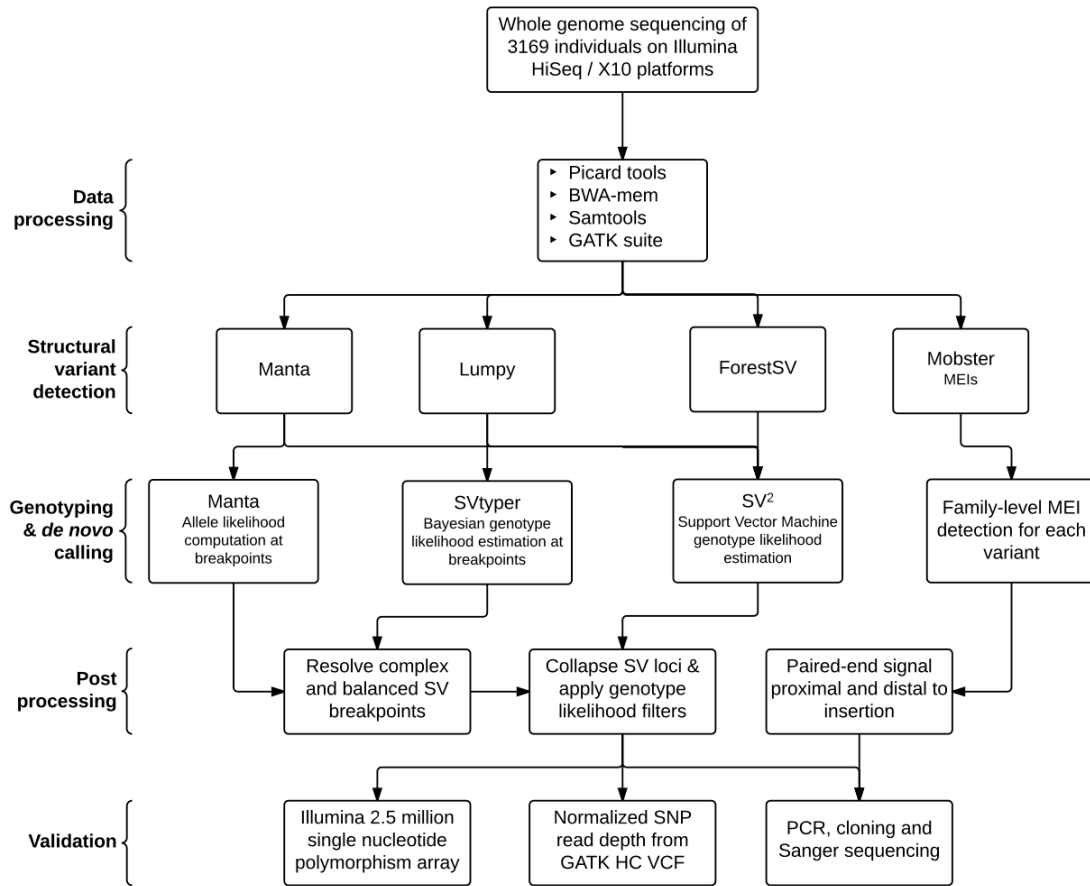529    100,000 permutations.

530    **Fibroblast cell culture and quantitative RT-PCR**

531    Dermal fibroblasts were obtained from two carriers of LEO1 deletions (a father and son)

532    identified in our study and additional unrelated control subjects by punch biopsy.

533    Fibroblast cell lines were then derived by Cellular Dynamics international

534    (https://cellulardynamics.com/) as part of the California Institute for Regenerative

535    Medicine Tissue Collection for Neurodevelopmental Disabilities (http://bit.ly/2mKUhB2)

536    and then provided to our lab for further study. Samples used for analysis included

537    fibroblasts from F0182|REACH000322 (ASD proband and deletion heterozygote),

538    F0182|REACH000321 (father, deletion heterozygote), and three unrelated control

539    samples: CW60038, CW60044, and JS034. Cells were recovered from cryogenic storage

540    as per CIRM's protocol and cultured in Dulbecco's modified eagle medium (DMEM)

541    supplemented with 10% fetal bovine serum, 2 mM L-glutamine, 100µg/ml penicillin and

542    100µg/ml streptomycin (Thermo Fisher Scientific, Waltham, MA, USA). Cells were

543    maintained in an incubator at 37°C at 5% $CO_2$ and harvested for RNA isolation at

544    passage three.

545    Total RNA was isolated using the Quick-RNA Microprep kit (Zymo Research, Irvine,

546    CA, USA) protocol for adherent cells with in-column DNAse treatment. cDNA was

547    synthesized from 100 ng of RNA using random oligo primers as part of the High

548    Capacity cDNA Reverse Transcription kit (Applied Biosystems, Foster City, CA, USA)

549    according to the manufacturer's protocol. Multiplexed qPCR reactions were conducted in

550    triplicate for each sample using gene-specific predesigned PrimeTime® qPCR assays for

551 *LEO1* (Hs.PT.58.448164, FAM-labeled) and the housekeeping gene *HPRT1*

552 (Hs.PT.58v.45621572, HEX-labeled) (Integrated DNA Technologies, Coralville, IA,

553 USA) on a CFX Connect Real-Time PCR System (Bio-Rad, Hercules, CA, USA) along

554 with no-template and no-reverse-transcription controls. Changes in gene expression were

555 calculated using the comparative $C_T$ method  and the null hypothesis was assessed using
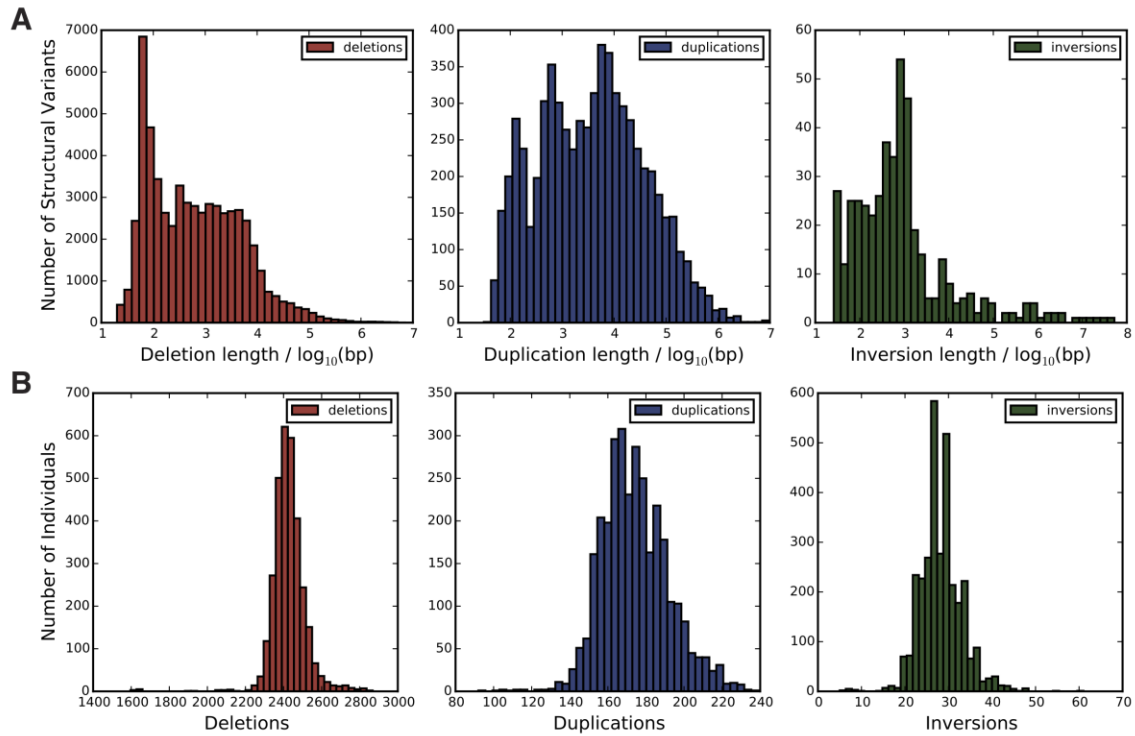
556 a Student's two-tailed unpaired T-test.

557

**Fig. S1.**

Flowchart detailing our custom pipeline for the discovery, genotyping, and validation of structural variants and *de novo* mutations. SV = Structural Variant; MEI = Mobile Element Insertion; PCR = Polymerase Chain Reaction.

564

**Fig. S2**

A) Histogram of the size distribution of deletions, duplications, and inversions per individual (log10 scale). B) Histogram of the number of deletions, duplications, and inversions per individual.

569

**Fig. S3**

Comparison of the SV call set from the discovery sample with the 1000 Genomes Phase 3 SV call set. A) Frequency of deletions, duplications, and inversions across parent allele frequency bins, stratified on known variants (from 1000 Genomes), and novel variants (detected only in this study). B) Venn diagrams of overlap of deletions, duplications, and inversions from our cohort with the 1000 Genomes.

576 **Fig. S4**

577 Metrics of genotyping accuracy for deletions and duplications by size. Bar charts
578 illustrating A) FDR based on intensity rank sum test from microarray, B) Mendelian error
579 rates, and C) variant transmission rates stratified on SV type (deletion and duplication)
580 and SV length bins for private variants. Quality metrics are reported for all private SVs in
581 the callset filtered based on $SV^2$ genotype likelihood at two levels of stringency
582 ("standard" and "*de novo*"). Whiskers represent 95% confidence intervals.

583

1

584

585  **Fig. S5**

586  Known autism genes are concentrated among genes that are most intolerant to loss-of-
587  function variants (pLI > 90th percentile).

588
589 **Fig. S6**
590 Patterns of deletion intolerance in the 1000 genomes phase 3 SV call set were very similar to those observed in this study (see Fig. 1).
591 (A) Depletion of deletions within exons correlated with a SNP-based measure of gene loss-of-function intolerance (pLI) from the
592 Exome Aggregation Consortium. (B) Promoters, Transcription Start Sites and UTRs showed the strongest deletion depletion for
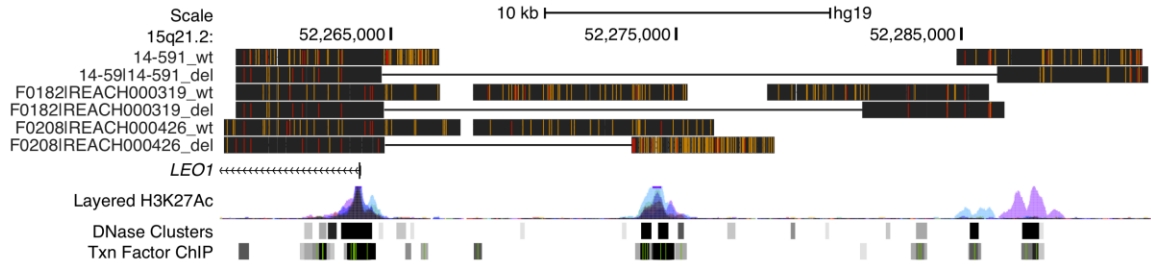593 variant intolerant genes (pLI >90th percentile).

1

| Cohort | % transmitted | p-value |
|---|---|---|
| REACH | 69.0 | 0.028 |
| SSC1 | 73.1 | 0.019 |
| MSSNG | 61.4 | 0.021 |
| SSC2 | 52.4 | 0.50 |
| **Combined Sample** | **63.4** | **0.00037** |

594

595 **Fig. S7**

596 Forest plot displaying the effect size (% transmitted) and 95% confidence intervals for
597 each of the four cohorts that were included in the study, including the two discovery
598 sample cohorts (REACH and SSC1), the two replication sample cohorts (MSSNG and
599 SSC2) and combined sample (discovery + replication). For detailed information see **table**
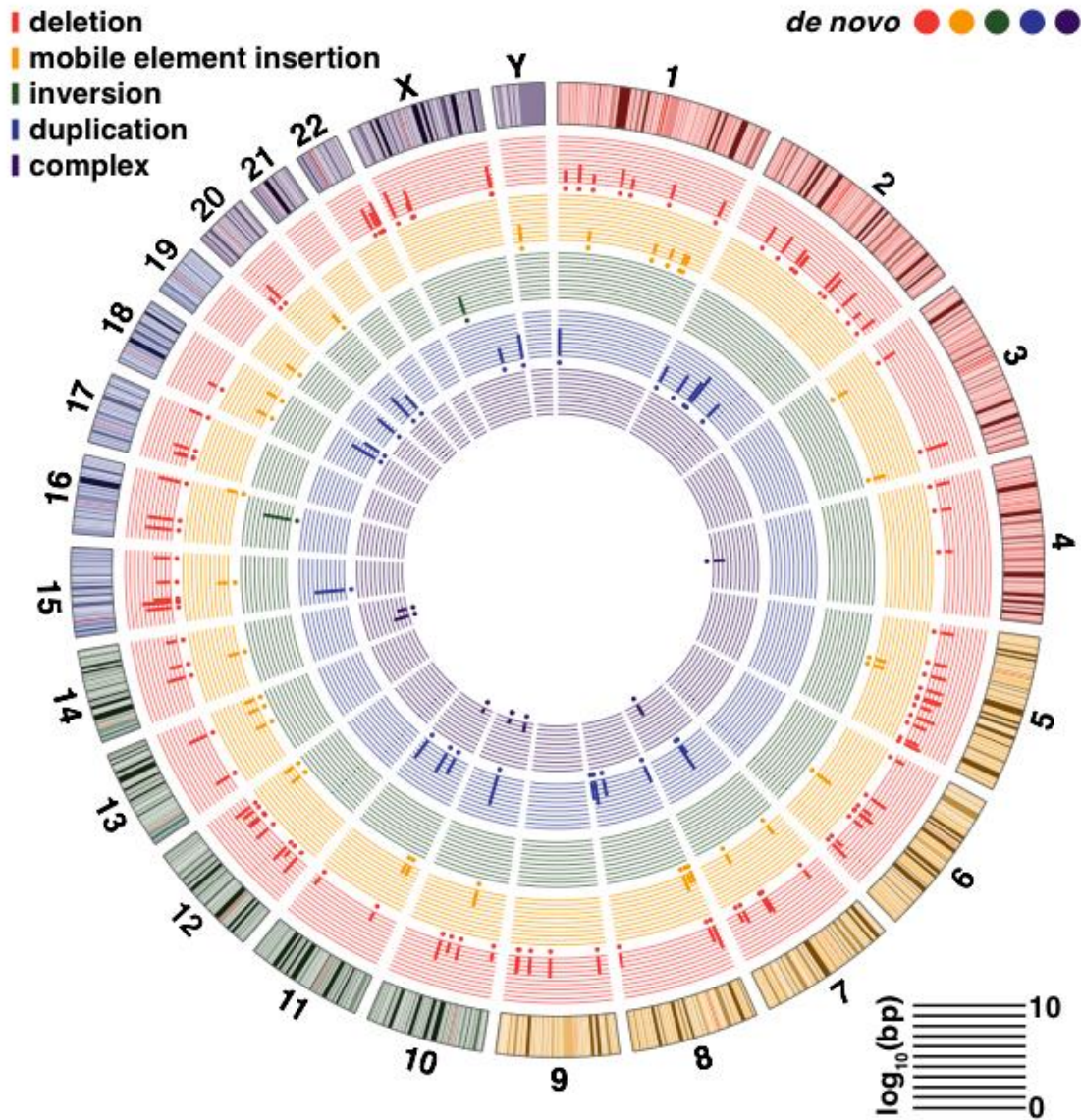600 **S6**.

601

602

1

603

**Fig. S8**

UCSC genome browser image showing BLAT alignments of Oxford Nanopore long read
sequences for three heterozygote deletions with corresponding wild type sequences. The
first two deletions are private to families 14-59 and F0182, and the third deletion is a
common polymorphism present in multiple families (an individual from F0208 was
selected for sequencing). Black bars show alignments with yellow lines indicating indels
and red lines SNPs. Wild type (wt) consensus contigs are shown within the breakpoint of
the deletion. Deletion (del) contigs mapping either side of the breakpoints are linked with
horizontal lines. Layered H3K27Ac = Histone 3 lysine 27 acetylation (an active promoter
associated mark) in seven cell types from ENCODE (GM12878, H1-hESC, HSMM,
HUVEC, K562, NHEK, and NHLF). DNase clusters = DNaseI Hypersensitivity Clusters
in 125 cell types from ENCODE (V3). Txn Factor ChIP = Transcription Factor ChIP-seq
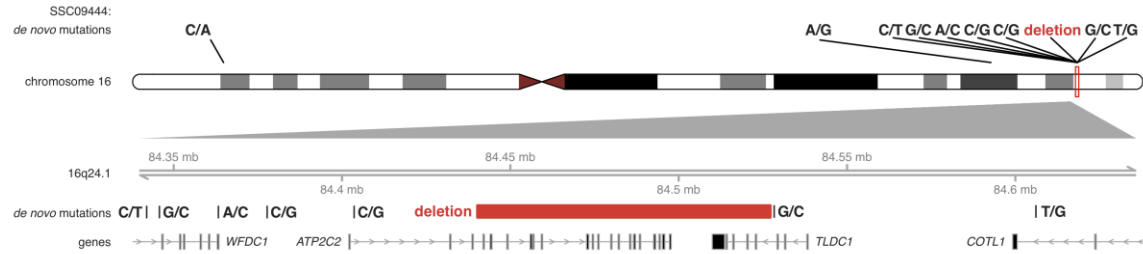(161 factors) from ENCODE with Factorbook Motifs (green).

617

618

**Fig. S9. *De novo* structural variation in 1,510 children**

Circos plot of *de novo* variants with concentric circles representing (from outermost to
inner): ideogram of the human genome with colored karyotype bands (hg19), deletions,
mobile element insertions, balanced inversions, tandem duplications, complex structural
variants. Circles indicate the location of *de novo* SVs, and their colors match the five SV
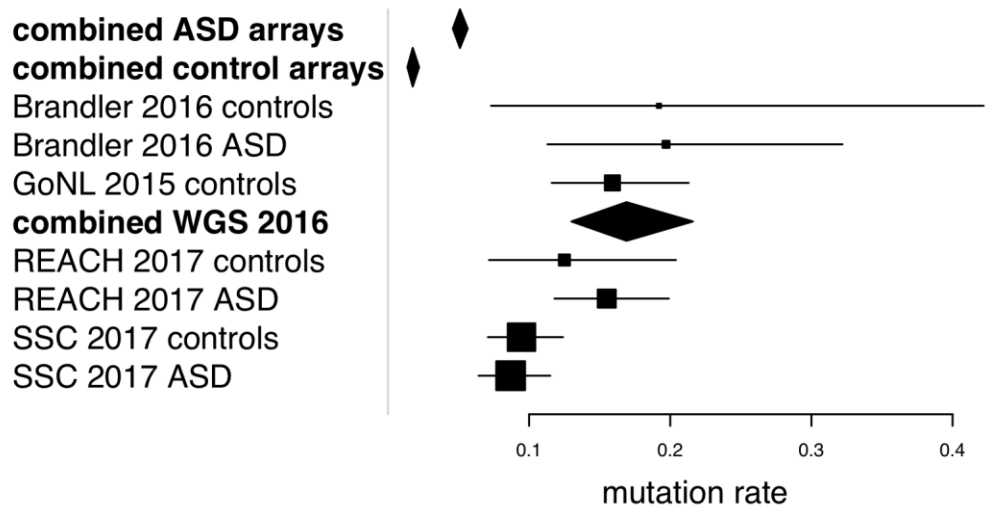types. Bars represent the $\log_{10}$ SV length of the *de novo* variants.

625

3

626
627
**Fig. S10**
One example of a complex mutation cluster are shown in the control individual from the
SSC, SSC09444 (alternate ID: 13874.s1). The 300kb zoomed in locus below the
ideogram shows the positions of *de novo* mutations relative to each other, an 82.3kb
deletion is clustered with six SNVs upstream and two downstream of it. Gene tracks
below the mutation show the longest transcript of each gene within the locus, with arrows
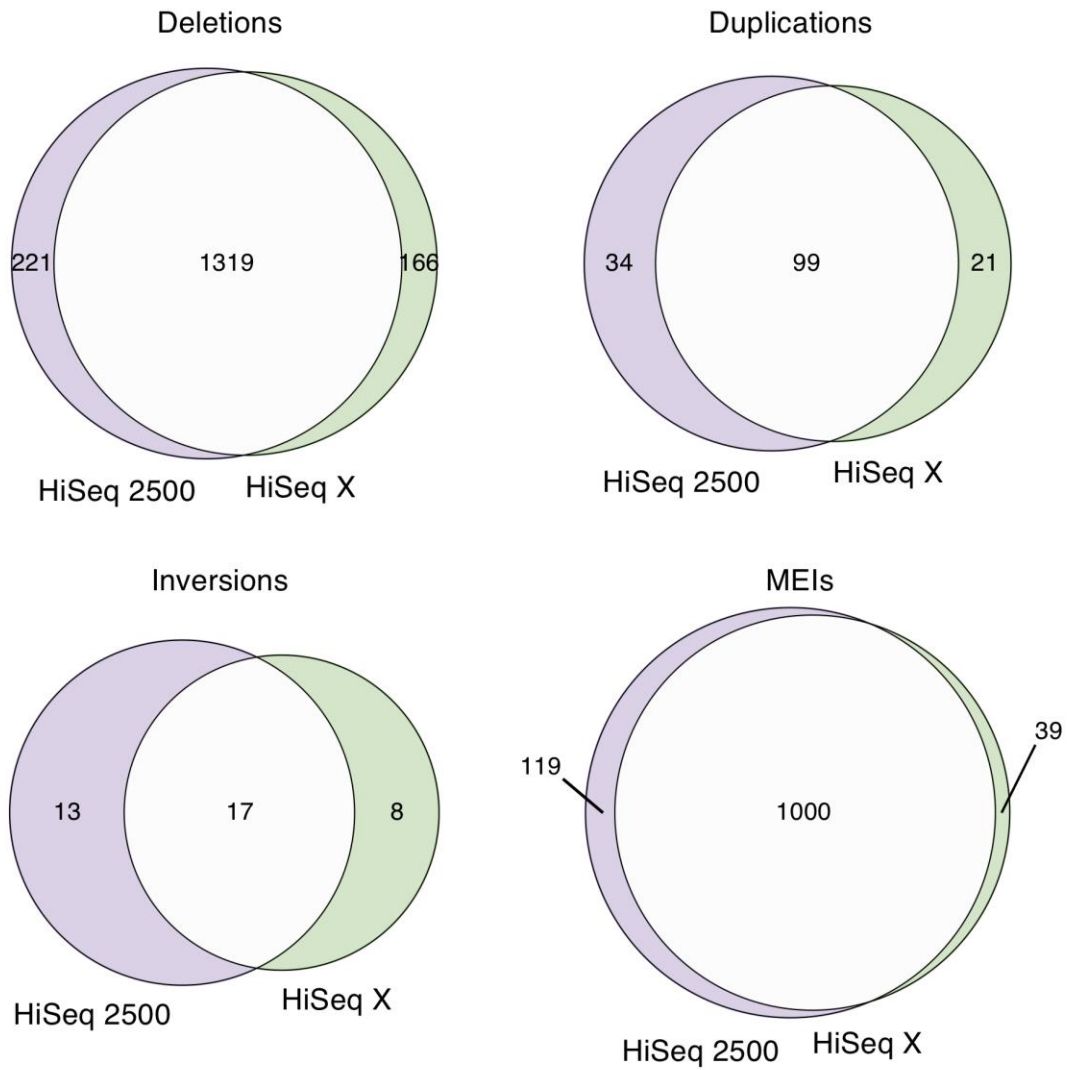indicating the strand and bars indicating the exons of genes.
635

mutation rate

636

637 **Fig. S11**

638 Forest plot of the *de novo* mutation rate in the two cohorts from the present study
639 (REACH 2017 and SSC1 2017) compared to previous whole genome sequencing and
640 microarray studies.

**REACH000236**

Deletions

Duplications

Inversions

MEIs

641
642    **Fig. S12**
643    Overlap between SV calls made from one sample sequenced on two platforms
644    Sample REACH000236 was sequenced at 43X coverage on both the Illumina HiSeq
645    2500 with 100bp reads and on the Illumina HiSeq X with 150bp reads. Venn diagrams
646    highlight the overlap for each SV type.
647

648     **table S1 (separate file)**
649     Information on samples used in this study
650
651     **table S2 (separate file)**
652     Descriptive statistics of the SV callset
653
654     **table S3 (separate file)**
655     False Discovery rate of copy number variants across size ranges and filters
656
657     **table S4 (separate file)**
658     Enrichment of known autism genes across pLI bins
659
660     **table S5 (separate file)**
661     Selection of target functional categories based on SV intolerance
662
663     **table S6 (separate file)**
664     Group-wise Transmission/Disequilibrium Test (TDT) results
665
666     **table S7 (separate file)**
667     SVs detected in the target functional categories in this study
668
669     **table S8 (separate file)**
670     Expression of *LEO1* and *MAPK6* in fibroblast cell lines from CRE-SV carriers and
671     controls
672
673     **table S9 (separate file)**
674     *De novo* SVs detected in the discovery sample
675
676     **table S10 (separate file)**
677     Complex Mutation Clusters
678
679     **table S11 (separate file)**
680     Known pathogenic SVs that were detected in the discovery sample