



Published in final edited form as:

Nat Genet. 2021 August ; 53(8): 1125–1134. doi:10.1038/s41588-021-00899-8.

Recent ultra-rare inherited variants implicate novel autism candidate risk genes

Amy B. Wilfert¹, Tychele N. Turner^{1,†}, Shwetha C. Murali^{1,2}, PingHsun Hsieh¹, Arvis Sulovari¹, Tianyun Wang¹, Bradley P. Coe¹, Hui Guo^{1,3}, Kendra Hoekzema¹, Trygve E. Bakken⁴, Lara H. Winterkorn⁵, Uday S. Evani⁵, Marta Byrska-Bishop⁵, Rachel K. Earl⁶, Raphael A. Bernier⁶, The SPARK Consortium⁷, Michael C. Zody⁵, Evan E. Eichler^{1,2,*}

¹Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA.

²Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA.

³Center for Medical Genetic & Hunan Key Laboratory of Medical Genetics, School of Life Sciences, Central South University, Changsha, Hunan, China.

⁴Allen Institute for Brain Science, Seattle, WA, USA.

⁵New York Genome Center, New York, NY, USA.

⁶Department of Psychiatry and Behavioral Sciences, University of Washington, Seattle, WA, USA.

⁷A list of consortium authors and affiliations appears at the end of the paper.

Abstract

Autism is a highly heritable complex disorder where *de novo* mutation (DNM) variation contributes significantly to risk. Using whole-genome sequencing data from 3,474 families, we investigate another source of large-effect risk variation, ultra-rare variants. We report and replicate a transmission disequilibrium of private, likely gene-disruptive (LGD) variants in probands but find that 95% of this burden resides outside of known DNM-enriched genes. This variant class more strongly affects multiplex family probands and supports a multi-hit model for autism. Candidate genes with private LGD variants preferentially transmitted to probands converge on the E3 ubiquitin-protein ligase complex, intracellular transport, and Erb signaling protein networks. We estimate these variants are ~2.5 generations old and significantly younger than other variants of similar type and frequency in siblings. Overall, private LGD variants are under strong purifying selection and appear to act on a distinct set of genes not yet associated with autism.

* eee@gs.washington.edu .

† Current address: Washington University School of Medicine, Department of Genetics, St. Louis, MO, USA.

AUTHOR CONTRIBUTIONS

A.B.W. and E.E.E. designed and conceived the study. L.H.W. and M.C.Z. coordinated samples and sequencing for CCDG cohorts. The SPARK Consortium coordinated samples and sequencing for the SPARK cohort. A.B.W., T.N.T., S.C.M., A.S., T.W., B.P.C., U.S.E., M.B.-B., and H.G. called variants and ran QC. K.H. performed Sanger validations. A.B.W., T.N.T., S.C.M., P.H., A.S., and T.W. conducted analyses and data interpretation. T.E.B. and A.B.W. performed gene expression analyses. A.B.W., R.A.B. and R.K.E. performed phenotypic analyses. A.B.W. and E.E.E. wrote the manuscript with input from co-authors.

COMPETING INTERESTS

All authors declare no competing interests.

Autism spectrum disorder (ASD) is a phenotypically heterogeneous disorder affecting about 1 in 59 children in the United States¹. Studies to date have primarily focused on high-impact, sporadic variants such as *de novo* copy number variants (CNVs) and single-nucleotide variants (SNVs). Despite their large effect sizes, *de novo* mutations (DNMs) account for ~3-25%²⁻⁵ of autism cases. Although this genetic model is highly relevant to simplex ASD, where only one child is affected in a family, it does not explain most cases and is less likely for multiplex families, where more than one child is affected⁶. This has led to the reassessment of various classes of inherited variation and their contribution to autism risk^{3,7-12}.

It is well established that large, sometimes inherited CNVs underlie a small percentage of sporadic and multiplex autism^{3,4}. Preferential transmission of likely gene-disruptive (LGD) variants have been observed in both simplex³ and multiplex autism⁸ in genes that converge on related functional networks^{3,8}. Genetic studies of ASD and developmental delay families have found affected children are enriched for multiple gene-disruptive variants (CNVs and SNVs)¹³⁻¹⁶. Recent analyses suggest that common inherited risk variants also contribute to ASD pathology^{5,17,18}.

In this study, we focus on private variants, or ultra-rare variants unique to a family. In contrast to other studies^{8,11}, we do not rely on support from DNMs or DNM rates and exclude genes known to be enriched with DNMs in ASD and neurodevelopmental disorder (NDD) cases to facilitate the discovery of novel genes. We expand the number of multiplex and simplex autism families sequenced, taking advantage of the increased sensitivity afforded by whole-genome sequencing (WGS) over whole-exome sequencing (WES) data¹⁹, to create a highly sensitive variant callset from WGS data from 3,474 autism families from the Centers for Common Disease Genomics (CCDG) (Table 1 and Supplementary Table 1). We assess transmission biases in probands and unaffected siblings after controlling for population structure^{8,20-23} and replicate these analyses in WES data from 5,879 families from SPARK (Simons Foundation Powering Autism Research for Knowledge)^{24,25}. Our results provide strong support for private LGD variants contributing to autism, particularly multiple hits in different genes. We show these variants have arisen more recently in autism families (2.5 generations) when compared to other classes of variants. Importantly, genes enriched for DNMs contribute little to this burden; rather, we suggest new gene candidates enriched in specific functional pathways.

RESULTS

Whole-genome sequencing.

We generated WGS data (30-fold sequence coverage) from 2,507 individual DNA samples from 394 multiplex and 251 simplex autism families (Table 1, Supplementary Table 1, and Online Methods). Combined with published WGS data^{8,21,22}, we created a standardized callset of SNVs and indels from 13,547 genome samples using two callers and made these publicly available (**Data and Code Availability**). The set consists of data from 4,364 probands and 2,235 siblings and includes parent-child SNV data from 774 multiplex and 2,700 simplex families. Focusing initially on DNMs, we employed two additional callers and performed 582 random Sanger sequencing validation experiments, combining these with

published validation experiments (Supplementary Table 2)⁴. We report an overall validation rate of 99.5% for our DNM callset and estimate a false negative rate of 3.5%. On average, we observe 65.14 DNMs per child and an increase of 1.11 and 0.37 mutations per year of paternal and maternal age, respectively (Supplementary Fig. 1). This estimate is lower than what was reported in Turner et al.⁴ because we required three out of four variant callers to agree to increase the specificity of the callset (Online Methods). As expected, *de novo* LGD and severe missense mutations are significantly enriched in probands when compared to siblings (LGD OR = 1.8, $P = 1.43 \times 10^{-23}$; MIS30 OR = 1.25, $P = 1.10 \times 10^{-4}$; Supplementary Table 3). Combining these *de novo* calls with published DNM callsets (Supplementary Table 4), we identify 100 genes after Benjamini-Hochberg correction (False Discovery Rate (FDR) < 5%, DNM count > 1) and 45 genes after Bonferroni correction ($P < 5.1 \times 10^{-7}$, DNM count > 1) with an excess of DNMs in autism probands (Online Methods and Supplementary Table 5). One gene, *MED13*, published as a case report²⁹, reached Bonferroni significance for an excess of DNMs in autism; four other genes (*CPA6*, *FRA10AC1*, *MPHOSPH10*, *RALGAPB*) reached FDR < 5% significance. These results replicate the reported excess of *de novo* LGD mutations in *RALGAPB*^{24,30} and identify three genes associated with other neurodevelopmental or neurodegenerative disorders with an excess of *de novo* missense mutations in ASD. *CPA6* is associated with epilepsy³¹; *FRA10AC1*, a gene of unknown function, is associated with Alzheimer's disease³²; and *MPHOSPH10* is associated with early-onset Parkinson's disease³³. These genes represent excellent candidates for future investigation.

Discovery and properties of private variants.

While previous studies focused on the contribution of DNMs or common variants underlying ASD^{2-4,6}, we focused on the contribution of transmitted variants^{3,8,11,12}. Because our previous study showed transmission disequilibrium signals increased with rarer inherited variants³, we focused on private inherited variants. We define these as heterozygous variants observed only once in the parent population and transmitted to at least one child, regardless of potential *de novo* status in unrelated children within the cohort. We note that 0.036% of our private variants overlap with mutations in our DNM callset. Based on our sample size, private variants are ultra-rare in nature and correspond to an approximate allele frequency 7×10^{-5} . We identified 26,606,722 unique private variants (35,871,117 total) in our discovery cohort of 6,599 children and detected no difference in the average number of private variants between probands and unaffected siblings genome-wide (Mann-Whitney test, autosomes, $P = 0.168$; female X chromosome, $P = 0.328$; male X chromosome, $P = 0.534$). We detect no difference in the number of private autosomal variants transmitted from fathers as compared to mothers genome-wide (Wilcoxon signed-rank test, $P = 0.1995$) but did detect a difference, as expected, if we consider the female X chromosome (Wilcoxon signed-rank test, $P = 0.0215$; mean paternally transmitted: 98.7, mean maternally transmitted: 102).

Since individuals of similar ancestry have an increased chance of allele sharing as compared to individuals from different populations³⁴, we considered private variants in the context of ancestry. We assigned individuals to one of six super populations (EUR, AFR, AMR, EAS, SAS, and OCN) based on maximum likelihood estimations of ancestry using a human

diversity panel (Online Methods and Supplementary Fig. 2). Consistent with previous studies^{35,36}, children of European ancestry carry the fewest private variants per genome (Figs. 1b,c and Supplementary Tables 6 and 7). This is because most individuals in the discovery cohort (85.6%) are of European ancestry (Supplementary Fig. 2). Private variants among the EUR subgroup will be of the lowest frequency, providing the greatest specificity, in principle, to detect pathogenic events³.

We tested whether filtering against a genetic database (dbSNP150) would be sufficient to eliminate this effect and improve our specificity for private events in other populations (Fig. 1c, Supplementary Fig. 3, and Supplementary Table 8). Although dbSNP filtering did reduce the average number of private variants per child, the magnitude of the effect varied by population. This treatment reduced the number of private variants by 69.2% among individuals of African ancestry but a reduction of only 43.6% and 44.9% among individuals of East and South Asian ancestry, respectively. Children of African and East Asian ancestry had, on average, similar variant counts prior to dbSNP filtering (mean: 11,630 AFR vs. 11,653 EAS; Supplementary Table 7). This suggests that the composition of the population genetic database may introduce additional biases because sampling across populations has been nonuniform, and allele frequency filtering alone is not sufficient to account for population stratification. These differences highlight the need to evaluate the impact of ancestry and increase underrepresented populations for gene discovery—even in rare variant analyses. We evaluated the impact of population stratification on our results by comparing burden estimates with and without ancestry as a covariate (Supplementary Table 8). We find the results are comparable and conclude that variation in the number of private variants between populations is generalizable and does not introduce biases into our analyses. Nonetheless, all analyses reported in this study have been replicated in the SPARK cohort and confirmed in the European subset of our discovery cohort (Supplementary Figs. 4-7).

Patterns of private, transmitted variants in protein-coding regions.

In this study, we restrict our analyses to autosomal, protein-coding regions of the genome where we expect to have the greatest power to detect enrichment of private, transmitted variants^{2,3,8}. Missense variants are the most abundant followed by synonymous and then LGD variants, defined here as stop-gain, stop-loss, splice-altering SNVs or frameshift indels (Fig. 1a and Supplementary Table 7). We observe no significant difference between the overall proportion of proband and unaffected sibling carriers for missense, synonymous, or LGD variants and detect no significant enrichment when considering all genes (Logistic regression, LGD: OR = 1.03, Bonferroni-corrected $P = 0.153$; MIS: OR = 1, Bonferroni-corrected $P = 1$; SYN: OR = 1, Bonferroni-corrected $P = 1$).

When considering subsets of genes at increasing thresholds of gene constraint using the probability of loss-of-function intolerance (pLI), we replicate^{3,8,37,38} the relationship of increasing burden of LGD variants in probands with increasing gene constraint for the discovery, replication, and combined cohort (Fig. 2a, Supplementary Figs. 4 and 8-10, Supplementary Table 9, and Supplementary Note). A similar trend was reported in Satterstrom et al. when considering ultra-rare variants (cohort AC 5) in a case-control

cohort and in ~6,305 families. We note the larger effect size reported in Satterstrom et al. for the case–control cohort is likely due to the presence of DNMs in these samples, and the smaller effect size reported for families is likely due to the higher allele frequency threshold used (Supplementary Fig. 11).

We expect this increased burden to be the result of a transmission bias and used a rare variant transmission disequilibrium test to confirm an overtransmission of LGD variants to probands (pLI = 0.99, Bonferroni-corrected $P=0.0137$ in probands, Bonferroni-corrected $P=0.52$ in siblings; Supplementary Table 10). Surprisingly, we observe a significant undertransmission of LGD variants in multiplex families with two probands and find no significant increase in allele sharing among affected siblings, suggesting an LGD variant in one child with autism is not predictive of a second child with autism (Supplementary Fig. 12). We do not observe an increased burden of missense or synonymous variants using pLI gene constraint thresholds. However, we observe an increase in potentially pathogenic, private missense variants in genes with increasing intolerance to missense mutation (Supplementary Fig. 13).

We estimate that the effect size of private, transmitted variants is ~8x smaller than the effect size of DNMs (OR = 11.67 for LGD DNMs in genes enriched for DNMs in ASD patients vs. 1.43 for private, transmitted LGD variants in genes with pLI = 0.99). We specifically excluded DNM-enriched genes in ASD cases as part of this calculation to estimate the effect size excluding well-established genes with an excess of DNM. In contrast, we compare the burden of private LGD variants among various autism risk gene sets to examine whether the private inherited and DNM signals were exclusive. These included genes shown to be enriched for DNMs in ASD and NDD cases^{10,11,39} and 845 genes from the Simons Foundation Autism Research Initiative (SFARI) (Online Methods). All gene sets show a trend toward enrichment of private LGD variants among probands when compared to unaffected siblings but to varying degrees. The Coe et al. gene set³⁹ shows nominal significance for enrichment in our discovery (Fig. 2b, Supplementary Figs. 5 and 14, Supplementary Table 11, and Supplementary Note; OR = 1.36, nominal $P=0.040$) and combined cohorts (Fig. 2b, Supplementary Table 11, and Supplementary Note; OR = 1.29, nominal $P=0.018$). The trends are consistent between replication and discovery cohorts, suggesting larger sample sizes are required to achieve significance that survives multiple-test correction. In general, DNM-derived gene sets show greater enrichment than a more general set of autism risk genes (i.e., SFARI). DNM-enriched gene sets derived from ASD and NDD studies perform as well (if not better) than those derived strictly from autism cohorts. Importantly, all trends disappear if we consider variants at higher allele counts in the parent population (Supplementary Fig. 15), indicating that the signal is strongest for inherited private variants.

Based on the initial sequencing of the SPARK autism families, Feliciano et al. reported most of the rare LGD variant transmission bias could not be accounted for by known ASD/NDD genes²⁴. We reevaluated the burden of private, transmitted LGD variants at increasing thresholds of gene constraint, excluding genes enriched for DNMs to quantify this effect. We find that 95.4% of private, transmitted LGD variant burden in probands remains (Fig. 3a and Supplementary Table 12) at pLI = 0.99 in the discovery cohort. We estimate that

private LGD variants in these DNM-enriched genes account for 1.45% of ASD risk, whereas private transmitted LGD variants in the remaining genes at $pLI = 0.99$ account for 2.64% of ASD risk (Table 2). Unlike *de novo* LGD mutations associated with autism, we estimate that most of the attributable risk for private variants awaits discovery and this risk will be identified among genes not already associated with DNM burden. Taken together, these results confirm that DNM-enriched genes confer substantial risk for ASD; however, there is additional burden in the less penetrant set of constrained genes ($pLI = 0.99$) yet to be discovered.

Simplex versus multiplex and a multi-hit model for ASD.

Both our discovery and replication cohorts consist of simplex families and multiplex families. Simplex families have been shown^{2,26} to be enriched for sporadic or *de novo* genetic events²⁰, while multiplex families are more likely to inherit ASD-predisposing variants⁴⁰. We compared the proportion of probands versus siblings carrying at least one private LGD variant at increasing thresholds of gene constraint considering simplex and multiplex families independently ($n = 2,700$ simplex vs. 774 multiplex families; Table 1 and Supplementary Table 1). Despite having 3.5-fold fewer families, multiplex families show a 25.7% higher burden of private, transmitted LGD variants in probands when compared to simplex families, with the greatest effect in less constrained genes (Fig. 3b, Supplementary Figs. 6 and 17, and Supplementary Table 15; multiplex vs. simplex OR = 1.37 vs. 1.09, permuted $P = 0.004$ at $pLI = 0.1$). Among simplex families, significant burden is observed, in contrast, among genes intolerant to mutation ($pLI > 0.99$).

Previous CNV work and analysis of putative noncoding DNMs^{14,19} have shown enrichment of multiple deleterious mutations in autism probands, while other recent studies have reported an additive effect between common and rare risk variants¹⁸. If the signal we observed was relevant to the genetic etiology of autism, we hypothesized that affected children would be more likely to carry multiple private LGD variants, partially explaining why both parents are unaffected or less severely affected in multiplex families. We compared the transmission of two or more private LGD variants in probands and unaffected siblings conditioning on intolerance to mutation. We find that probands are significantly more likely to carry multiple inherited LGD variants in less constrained genes when compared to unaffected siblings (Fig. 3c, Supplementary Fig. 18, Supplementary Table 16, and Supplementary Note; OR = 1.29, Bonferroni-corrected $P = 0.026$ at $pLI = 0.1$). Under an additive model, which represents independent assortment and random segregation, we would expect the odds ratio for the two-hit model to equal the square of the odds ratio for the one-hit model. This is exactly what we observe, and the effect becomes stronger if we restrict the analysis to individuals of European ancestry (Supplementary Fig. 12), indicating that this signal is not an artifact of population stratification.

Novel candidate genes and interconnected functional networks.

We investigated whether highly constrained genes not enriched for DNMs showed enrichment for expression or protein-protein interaction (PPI) networks. Previous studies^{8,10,41} typically performed such analyses by integrating candidates with DNM-enriched genes as opposed to considering them separately. We focus on 163 highly

constrained genes ($pLI = 0.99$, Supplementary Table 17) where private LGD variants are exclusively transmitted to probands and have not been reported in SFARI or as DNM enriched in three ASD/NDD studies^{10,11,39}. Among these genes, there are a total of 276 LGD variants and 28 genes with independent LGD variants observed in two or more unrelated families.

Gene ontology (GO) analysis shows that the candidate gene set is highly enriched for encoded phosphoproteins (Supplementary Table 18; KW-0597, 129/163 genes, $q = 1.93 \times 10^{-20}$), and the genes are more likely to be interconnected as part of PPI networks (Fig. 4; 102 observed vs. 75 expected edges, $P = 0.00164$). A subset of the genes (74/163 genes), including half the genes with events in multiple families, converge on several functional pathways (Fig. 4 and Supplementary Table 18). This includes a small network of genes enriched for the E3 ubiquitin-protein ligase pathway by both the GO and Reactome databases, which are involved in proteasome degradation (HSA-98316) and regulation of protein modification by small protein conjugation or removal (GO:1903320). Similarly, there is a set of more than a dozen genes associated with internal cellular transport and specifically transport between the Golgi and endoplasmic reticulum. Other subnetworks are significantly enriched for nucleobase-containing compound metabolic process (GO:006139) and Erb signaling (hsa04012).

This proband candidate gene set is also enriched for cell-type-specific expression at the early and mid-fetal cortical stages of human brain development (Supplementary Fig. 19). We observe no enrichment in a set of 83 genes in siblings ascertained using the same criteria (not DNM enriched, $pLI = 0.99$, no private LGD variants in probands) (Supplementary Fig. 19). If we focus our expression analyses from brain regions to individual cell types in the adult human cortex, we find that our candidate genes are significantly enriched for expression in both excitatory and inhibitory neurons (Supplementary Fig. 20; excitatory $P = 4.7 \times 10^{-4}$, inhibitory $P = 5.0 \times 10^{-4}$) but not enriched for expression in non-neuronal cell types (Supplementary Fig. 20; $P = 0.24$) as compared to control sets. There is no difference between proband and sibling genes ascertained using the same criteria. It should be noted that these pathway enrichments are only observed when compared to the whole genome. If we compare to only genes intolerant to mutation ($pLI = 0.99$), no pathways remain significant.

Private LGD variants in children with autism are evolutionarily younger.

Classical population genetics predicts that deleterious variants, such as disease-associated alleles, should be, on average, younger than neutral alleles of the same allele frequency due to purifying selection⁴². Focusing on children of European ancestry, we applied a genome-wide genealogy method developed by Speidel and colleagues⁴³ that uses the local ancestry (i.e., linkage disequilibrium) surrounding a single-nucleotide polymorphism (SNP) of interest to construct a coalescent tree and estimate the generational age of the allele based on the coalescent branch length. We selected 101 private LGD variants transmitted only to children with autism where none of the 163 candidate genes were previously associated with ASD. We compare them to a random subset of ~500 private LGD variants in other genes obtained from both probands and siblings. We estimate the average age of disease-

associated LGD variants to be 2.5 generations and find these are significantly younger than other classes of private LGD variants (Fig. 5). We estimate that other proband-associated LGD variants in highly constrained genes ($pLI < 0.99$) outside the candidate gene set have a median age of 3.7 generations and are significantly older (Mann-Whitney U test, Bonferroni-corrected $P = 0.0133$). Sibling LGD variants in highly constrained genes ($pLI < 0.99$) are estimated to be almost two generations older (4.3 generations; Mann-Whitney U test, Bonferroni-corrected $P = 0.000255$ candidate vs. sibling). As a negative control, we do not observe any difference between the age of private LGD variants in genes outside of the candidate gene set between probands and siblings (Fig. 5; Mann-Whitney U test, $P = 0.139$) or for synonymous or private variants mapping to intergenic regions (Supplementary Figs. 21 and 22). Since alleles in these candidate genes are carried in unaffected parents, we hypothesize that these variants are under weaker selection than deleterious DNMs but under stronger selection than a neutral allele. Specifically, if we assume mutation-selection balance under an additive (e.g., two-hit; $h = 0.5$) model⁴², we can apply gene-specific mutation rates and allele frequencies within the cohort to estimate the median selection coefficient for the 101 private, transmitted LGD variants in probands of European ancestry. We estimate a rather strong selection coefficient of 0.27 (s.d. = 0.24) for private candidate LGD variants transmitted to only autism probands in this study.

Contribution of known and novel ASD-associated variation.

Since common variants are implicated in autism risk, we calculated the polygenic risk score (PRS) from our larger sample set and assessed transmission disequilibrium as recently described¹⁷. We find an even larger difference in the transmission of polygenic risk between probands and unaffected siblings when compared to Weiner and colleagues, observing the signal in both multiplex and simplex families (Supplementary Fig. 23). We quantified the relative increase in risk for ASD conferred by common variants, DNMs, private SNVs, and a set of CNVs previously implicated in autism and NDD. Using a multivariate logistic regression, we estimated the effect size of these four variant categories and computed the population attributable risk (PAR) (Supplementary Table 19). We restricted this analysis to the Simons Simplex Collection (SSC; $n = 1,765$ quads) where we had variant calls across all four mutation classes for all samples. We find that children with *de novo* LGD mutations in DNM-enriched genes are 11.7 times more likely to have autism, accounting for 4.4% of the PAR for ASD. Although CNVs associated with ASD and NDD are, collectively, the rarest events included here, children with such an event are 2.7 times more likely to have autism, accounting for 0.9% of the PAR. Carrying one or more private, LGD variants in highly constrained genes increases the likelihood of developing autism by 1.4-fold. We estimate that these events account for 3.3% of the PAR, which is comparable to the amount of risk associated with LGD DNMs. Lastly, we find a 1.1-fold increase in the likelihood of developing autism as polygenic risk increases, further supporting a polygenic transmission disequilibrium. While PRS accounts for a large fraction of ASD heritability, we estimate that having a PRS in the top 10% of all children accounts for 1.8% of the PAR for ASD. We note that the contribution of polygenic risk is likely an underestimate as ASD genome-wide association studies to date are underpowered. Only a small number of robust loci have been identified, so we are likely missing much of the common variant liability for ASD¹⁷. These

four categories of risk variants only account for 10.4% of the PAR for autism, suggesting many more risk factors for autism are yet to be discovered.

DISCUSSION

Despite the high heritability of autism, most gene discovery in autism research has been driven by studies of *de novo* variation^{2,8-10,26}. Our analysis shows that ultra-rare transmitted LGD variants are not only significantly enriched in children with autism but contribute to at least 4.5% of autism risk in the human population. This estimate is in line with other studies^{3,37} and suggests this understudied class of variation may confer almost as much risk as *de novo* SNVs and indels (6-9% of cases using the same PAR estimator)^{2,3}. While the burden of private LGD variants in affected children is higher in multiplex families, both simplex and multiplex families show evidence of biased transmission of private LGD variants. This effect is significant in simplex families only for genes intolerant to mutation, while in multiplex families the effect is larger and significant for genes more tolerant to mutation (Fig. 3b). This may explain why we observe a significant excess of multiple private LGD variants in probands as multiple gene disruptions may be required to reach the diagnostic threshold for ASD.

Some studies focused on identifying risk genes have combined *de novo* and ultra-rare variant risk burden to improve sensitivity, such as the Transmission And *De novo* Association (TADA)⁴¹ analysis employed by Ruzzo and colleagues⁸. Because a significant fraction of DNM-enriched genes have been discovered^{8,11,39}, we sought to tease apart these effects by excluding known DNM-enriched genes. We estimate that about half of private LGD risk is conferred from genes identified through DNM enrichment studies, and excluding known DNM-associated risk genes has a marginal effect on the burden we observe. To enrich for pathogenicity, we identified a set of 163 candidate genes according to gene constraint (pLI > 0.99) and the absence of private LGD variants in unaffected siblings. Although there has been no reported evidence of DNM enrichment in these genes, we find that several of our candidate genes and gene networks identify pathways previously implicated in autism.

For example, we identified three independent private LGD variants in *HDAC9* transmitted exclusively to probands. Pinto et al. identified a transmitted *HDAC9* deletion in a patient with ASD and five additional gene deletions in patients with intellectual disability and schizophrenia⁴⁴, supporting the role of private, transmitted LGD variants in *HDAC9* in ASD pathogenesis. Several other *HDAC* genes have also been implicated in ASD, including *HDAC8*⁴⁵ and *HDAC4*⁴⁶, and the chromatin-remodeling pathway is known to play a key role in autism^{9,26,47}. Another gene in our network, *TOP2A*, is part of the topoisomerase gene family thought to be critical in regulating the expression of ASD-related genes⁴⁸. Although this specific topoisomerase has not yet been reported as enriched for DNMs in autism cases, inhibitors of this gene alter the expression of imprinted genes, and the topoisomerase acts by resolving transcription-associated supercoiling of long genes, including ASD-related genes critical for synaptic function. Our findings suggest that private inherited variants may identify a subset of genes with variants of smaller effect sizes; however, we and others have shown³⁹ that more than half of all genes enriched for DNMs have yet to be discovered, and none of the inherited genes reach gene-level significance.

This is in large part due to sample size and the locus heterogeneity underlying autism. With greater sample sizes, there will likely be more extensive overlap between inherited risk and DNMs-risk genes. The case-control study design may be particularly well suited to validate individual candidate genes with an increased burden of private, transmitted variants in autism family studies.

Additionally, we identified a small network of seven genes in the E3 ubiquitin-protein ligase pathway, which has a well-characterized role in autism^{49,50}. There are several genes in this pathway enriched for DNMs in children with autism^{10,39}, indicating that DNMs and private, transmitted LGD variants converge on the same pathway but may be hitting distinct sets of genes. An interesting finding from this study is the discovery of a subnetwork of genes (e.g., dyneins, kinesins, and coatomer subunits) related to vesicular intracellular transport between the Golgi and endoplasmic reticulum (Fig. 4). This process is important in the transport of synaptic molecules, such as neuroligins and neuroligins, to the cell surface, endocytic cycling of receptors, and vesicular cargo transport along microtubules⁵¹⁻⁵³. Mutations in related genes in both autosomal recessive and dominant form have been implicated in autism, peripheral neuropathies, and NDD. Disruptions in gene function alter synaptic plasticity and morphology of neuronal dendrites and axons. While these associations are exciting, we caution that network and enrichment analyses are often biased toward the most well-studied genes and pathways⁵⁴, and thus, more than half the genes that failed to associate with a functional network likely await discovery.

Finally, we report evidence supporting a multi-hit model of autism. We find that private truncating variants in different genes are 50% more likely to occur in autism probands than siblings—a signal consistent with the pathogenicity of this class of variant (Figs. 2a and 3c, Supplementary Table 11, and Supplementary Note). There are other instances of such models reported in ASD ranging from a simple two-hit model⁵⁵ to an oligogenic model of disease^{4,13,14,19,55}. For example, the 16p12.1 deletion^{14,55} is often inherited but requires a secondary CNV to reach the genetic liability threshold for disease. Similarly, carrying three or more potentially deleterious DNMs (in the absence of an LGD DNM or large CNV) can be attributed to about 7.3% of autism cases⁴. A targeted study of seven genes identified a significant overrepresentation of probands with two or more nonsynonymous variants and suggests multiple moderate impact events in the same pathway are necessary to cause nonsyndromic forms of autism¹³. Efforts focusing on patient recontact, not only for the purpose of re-phenotyping families as diagnostic criteria evolve but also for providing additional counseling as novel genetic candidates are identified, will be critical in the task of understanding genotype-phenotype relationships and has already been proposed by others⁵⁶. Understanding the diversity of genetic etiologies underlying autism as well as their corresponding phenotypic outcomes will be critical for providing accurate risk assessments for family planning and genetic counseling.

These findings highlight some key considerations for future ASD studies. Specifically, family composition of the cohort will influence what types of and to what degree different variant classes contribute to ASD risk. This is important for replication of the findings reported here as well as findings from other groups⁸⁻¹⁰. Most autism families characterized by exomes and genomes are simplex in origin, and a greater effort must be taken to recruit

and characterize multiplex families as part of large-scale sequencing efforts. Additionally, these results highlight the weakness of assuming *de novo* and rare transmitted variants will impact genes in a similar manner (e.g., monogenic and highly penetrant mutations in constrained genes)^{38,41,57}. Although we find *de novo* and private variants converge on related pathways, our data suggest these two variant classes may act through different genetic mechanisms and modulate distinct sets of genes in ASD pathogenesis. Our allele age estimates are consistent with the action of strong selection operating on these variants. Our analysis suggests that the variants we identify in candidate genes persist for two to three generations before being removed from the gene pool by selection. In contrast, most of the LGD variants associated with *de novo*-enriched genes are removed from the gene pool almost immediately due to the action of stronger purifying selection.

ONLINE METHODS

Sequencing and quality control of cohorts.

Individuals enrolled in the Autism Genetic Resource Exchange (AGRE), SSC, The Autism Simplex Collection (TASC), and Study of Autism Genetics Exploration (SAGE) studies were whole-genome sequenced at the New York Genome Center (NYGC) as part of the CCDG (<http://ccdg.rutgers.edu/>) (Table 1 and Supplementary Table 1). This study was approved for sequencing by the local institutional review board (IRB) at NYGC (Biomedical Research Alliance of New York (BRANY) IRB File # 17-08-26-385). All participants provided informed consent prior to participation in the study (SSC: IRB STUDY00001619, SAGE: IRB protocol #44219, TASC: STUDY00002514 at the University of Washington). Sequencing was performed on an Illumina HiSeq X Ten platform using 1 µg of DNA and an Illumina PCR-free library protocol. Post-sequencing, the data were processed using the standard pipeline for the CCDG⁵⁸ and the GRCh38_full_analysis_set_plus_decoy_hla.fa reference genome. Briefly, raw reads were aligned to the GRCh38 reference genome (BWA-MEM, v0.7.15)⁵⁹, duplicate reads were marked (Picard v2.5.0, <http://broadinstitute.github.io/picard/>), base scores recalibrated (GATK v3.8.0)⁶⁰, and indels were realigned (GATK). CRAM quality control (QC) metrics for the SAGE cohort have been previously published²¹; SSC, TASC, and AGRE QC metrics were determined using Picard WGS metrics, Picard insert-size metrics, and SAMtools⁶¹ flagstat. The average sequence depths for SSC, TASC, and AGRE were 34.99 ± 4.09 -fold, 33.89 ± 5.46 -fold, and 33.03 ± 4.31 -fold, respectively. The average insert sizes for SSC, TASC, and AGRE were 444.9 ± 17.86 bp, 455.1 ± 6.40 bp, and 384.8 ± 26.39 bp, respectively.

Individuals enrolled in the SPARK study were whole-exome sequenced at Regeneron (unpublished) (Table 1). Exomes were sequenced to an average coverage of 61.84 ± 14.99 -fold. QC analysis included HybridizationMetrics (Picard) and SAMtools flagstat.

Variant calling.

We called SNVs and indels in families using four different callers: GATK HaplotypeCaller (v.3.5.0)⁶², FreeBayes (v1.1.0)⁶³, Platypus (v0.8.1)⁶⁴, and Strelka2 (Illumina, v2.9.2). In addition, multi-nucleotide variants were called using FreeBayes and Platypus. Post-calling,

BCFtools (v1.3.1)⁶⁵ norm was used to left-align and normalize indels. Following variant calling, we partitioned the genome into the high-quality regions, consisting of unique space as well as ancient repeats and the recent repeat regions, which consisted of repeats <10% diverged from the consensus in RepeatMasker. Variants in high-quality portions of the genome retained for analysis and recent repeat region variants were removed from the study.

Kinship and sample redundancy.

All samples from both the discovery and validation cohorts were merged together and kinship coefficients were calculated with KING (v1.4)⁶⁶. Samples with kinship coefficients that did not match their reported relationship were identified as potential sample swaps or contamination and were either removed or, when possible, their relationships were corrected. Samples with kinship coefficients greater than 0.35 were identified as potential sample duplicates (Supplementary Table 20). We first checked whether potential duplicates were known monozygotic twin pairs or known duplicates within a cohort (some individuals had both blood and cell line DNA sequenced for QC purposes). We retained one sample from each of the known duplicate pairs, preferentially retaining the sample generated from blood DNA when possible and randomly selecting the retained sample if there was no difference in DNA source. The remaining duplicates, which represented samples sequenced as part of multiple cohorts, were retained for one and only one of the cohorts according to the following prioritization scheme: (1) sample was sequenced as part of an SSC family; (2) the sample was from a complete family, their DNA was from blood, and was WGS; (3) the sample was from a complete family and sample was WGS; or (4) the sample was from a complete family and contained unaffected siblings. Families with twins were retained for private variant discovery but excluded from all statistical analyses.

Principal component analysis (PCA).

In addition to our discovery and validation cohorts, we included two reference cohorts, 1000 Genomes Project (1KG, 20140818 release)³⁵ and Simons Genome Diversity Project (SGDP, available at NCBI under BioProject ID PRJNA522307)⁶⁷, in our PCA. Each cohort was cleaned separately (described below), merged together, and then cleaned additionally. Reference data from the SGDP and 1KG were prepared for PCA by left normalizing variants with BCFtools (v1.9), followed by filtering for individual missingness (<10% missing genotypes within an individual), SNP missingness (<50% missing genotypes across a SNP), minor allele frequency (>5%), and linkage disequilibrium pruning with PLINK 1.90⁶⁸. Sites were then converted from hg19 to GRCh38 using UCSC liftOver. Since the 1KG data were generated with a lower density SNP array than SPARK, the remaining 1KG sites were the only sites considered in the remaining cohorts. The SSC, SAGE, TASC, and AGRE samples were prepared using GATK joint-genotype files provided by the NYGC and then iteratively merged together within the respective cohort (most cohorts had to be processed in multiple batches). Each joint-genotype file was prepared as described above. The SPARK samples were prepared using the joint-genotype generated by Regeneron using Illumina InfiniumCoreExome-24_v1.1 array data. Intensity data files were processed using Illumina Genome Studio Software. Since this data set was already in plink format, it did not undergo additional processing. Prior to merging all six cohorts together, the 1KG target sites were extracted from each cohort. After merging, the combined autism and reference cohorts were

filtered for genotype missingness within the individual and SNP (both < 5%). Finally, the data were input into EIGENSTRAT (v5.0.1)^{69,70} for PCA. The results of this analysis are summarized in Supplementary Figure 24.

ADMIXTURE and ancestry assignment.

The files we used for PCA input were split by reference cohort (SGDP and 1KG) and discovery cohort (SSC, SAGE, TASC, and AGRE) and filtered for sites present in the reference cohort and individual-level missingness a second time. We ran the software ADMIXTURE (v1.3.0)⁷¹ with 10-fold cross validation (CV) on our reference cohort of 1,964 unrelated individuals to determine the optimal value for the K parameter. We found that K = 10 resulted in the smallest CV error (Supplementary Fig. 2); however, there is little difference in CV error for values of K between 8 and 14, and we recognize that a lower value of K would result in similar population assignments. We assessed the quality of our inferences for our reference cohort by visualizing the proportion of ancestry from each cluster for a random subset of 15 individuals from each known population (250 individuals total).

Due to the underlying relationships between individuals in our autism cohort, we chose to use the allele frequencies learned by ADMIXTURE from our reference cohort to assess the ancestry of our discovery cohort in a supervised manner by using projection with ADMIXTURE. We assigned each individual to the cluster that contributed the largest proportion of ancestry and then grouped clusters into six super populations (EUR, European; AFR, African; EAS, East Asian; SAS, South Asian; AMR, Amerindian; OCN, Oceanian) according to membership of known populations from the reference cohort (Supplementary Fig. 2). We were unable to assign ancestry to 1.01% of our discovery cohort due to missing data in the joint-genotype files and find that the majority of our cohort (85.9%) is European ancestry (Supplementary Fig. 2 and Supplementary Table 6).

***De novo* mutation (DNM) calls.**

DNMs in the SSC, SAGE, and TASC cohorts were called using a custom pipeline. DNMs were not called in AGRE because DNA for most samples in this cohort was derived from cell lines, which are prone to introducing artifacts in DNM analyses. First, variants that were *de novo* based on genotype (father and mother genotypes were equal to 0/0 and the genotype in the child was 0/1 or 1/1) were retained for further assessment. Second, variants from Platypus with a filter of LowGQX or NoPassedVariantGTs were removed, and Strelka2 variants had to have the filter field equal to PASS. Third, variants needed to have the support of at least two of the four callers. Fourth, variants were resequenced with FreeBayes using default settings and needed to remain as *de novo*. Fifth, variants in a homopolymer A or T of length 10 or greater were removed. Sixth, we removed all variants in low-complexity regions, recent repeats, or centromeres. Finally, we applied the following sample-level filters: the father alternate allele count = 0, mother alternate allele count = 0, child allele balance > 0.25, father depth > 9, mother depth > 9, child depth > 9, and either child genotype quality (GQ) > 20 (GATK) or sum of quality of the alternate observations (QA) > 20 (FreeBayes). For variants on the X chromosome, we separately considered variants in the

pseudoautosomal regions (chrX:10000-2781479, chrX:155701382-156030895) and the X/Y duplicatively transposed region (chrX:89201803-93120510).

We performed random Sanger validation and combined these data with published validations to look at a total of 3,233 sites in a conditional inference analysis (Supplementary Table 2). The metrics we included in this analysis included: (1) the mer150 mappability, which we calculated on build 38 of the human genome using a workflow originally designed as part of the ENCODE project; (2) the average mapping quality of the read ± 100 base pairs (bp) around the variant in the child; (3) the average mismatch in the reads ± 25 bp around the variant in the child; and (4) the callers that supported the event as *de novo*. Based on this analysis, the final dataset for *de novo* SNVs and indels were sites that either had the support of all four callers or were supported by three callers and had an average mapping quality greater than 57 for the reads in the 100 bp region around the variant. For the multi-nucleotide variants, we also inspected all sites using SAMtools tview, and the sites had to have visual inspection support of *de novo* status and an average mapping quality greater than 57 for the reads ± 100 bp around the variant. We estimate our validation rate in this dataset at 99.5% and our false negative rate at 3.5%. In addition, we removed samples that were statistically defined as outliers, in terms of *de novo* counts, based on the boxplot function in R.

Private SNV calls.

Each cohort was assessed separately to identify ultra-rare, inherited variants using a custom pipeline. Briefly, SNVs and indels were called using FreeBayes (v1.1.0) and GATK on a per-family basis. Sites were left-aligned, normalized, and multiallelic sites were split into separate lines using BCFtools (v1.9). Sites from the two callers were merged using GATK CombineVariants. To ensure a high level of specificity, we counted all alleles in the parent population that were present in the union set of the two callers and passed the following QC filters: (1) site quality score (QUAL) > 50 , and (2) read depth (DP) ≥ 10 for genomes and DP ≥ 20 for exomes. We used slightly different DP filters for the exome and genome data to account for differences in sequencing depth between the two sequencing platforms. All sites that were heterozygous and observed only once in the parent population were designated as candidate private variants (cohort-level parental frequency $\approx 7 \times 10^{-5}$; approximate equivalent ExAC frequency $\approx 2.5 \times 10^{-5}$).

The set of private variants for each cohort was comprised of candidate private variants that were present in the intersection set of GATK and FreeBayes and did not violate the rules of Mendelian inheritance. We annotated variants using SnpEff (v4.3t)^{72,73} with gene and transcript information (GRCh38.86), predicted effect of the variant on the transcript, ExAC (r0.3, non-neuropsych subset) lifted over to GRCh38 using the UCSC liftOver tool, and dbSNP (v150). Finally, variants were filtered against recent repeats (see DNM methods for details), low-complexity regions, centromeres and gaps, and pseudoautosomal regions (hg38 chrY:10,000-2,781,479, chrY:56,887,902-57,217,415, chrX:10,000-2,781,479, chrX:155,701,382-156,030,895) using BEDTools (v2.24.0)^{74,75}.

The set of private variants from each cohort was compared to all variants observed in the other cohorts. Candidate private variants not observed in any other cohorts were retained

for our final set of private variants. For example, the discovery cohort private variants are comprised of sites unique to one parent across only the WGS cohorts, whereas the combined set private variants are comprised of sites unique to one parent across both the WES and WGS cohorts. When combining the WES and WGS cohorts, we only included regions with an average coverage of 20-fold in the exomes.

Private CNV calls.

WES CNVs were called from 21,442 individuals among 5,904 complete families (including 6,582 probands and 3,045 siblings) in SPARK_WES_1 release using CoNIFER (v0.2.2)⁷⁶ and XHMM (version statgen-xhmm-3c57d886bc96)⁷⁷, as described previously²⁴. An independent SNP microarray dataset from 99.2% of the samples (21,271/21,442) was generated using Illumina Infinium Global Screening Array-24v1 (GSA-24v1), SNP CNVs were detected using CNVPartition (Illumina, v3.2.0), PennCNV (v1.0.4)⁷⁸, and CRLMM (v1.38.0)⁷⁹ as described previously²⁴. We assessed inheritance using both SNP and WES data and filtered putative valid private CNVs based on the inheritance (paternal or maternal), CNV frequency in parents ($n = 1$), number of exome probes (> 4), percentage of overlap with segmental duplication ($< 75\%$), and microarray validation (support by at least one of CNVPartition, PennCNV, and CRLMM approaches). In addition, we required the CNV to only interrupt a single gene to detect “gene-killing” CNVs. A pLI score was assigned to each gene spanned by a private CNV, and CNVs were binned by pLI scores and copy number type (deletions and duplications). Fisher’s exact test (two-sided) was performed to compare the number of probands and siblings carrying “gene-killing” private CNVs. Whole-genome structural variant (SV) calling was conducted on short-read WGS data from the SSC cohort, comprised of 8,617 samples from 2,276 families, after SV calling and QC, and the SAGE cohort. Calling and merging of SVs was done as described in Turner et al.⁴, with the exception that this study used Delly²⁸⁰ instead of VariationHunter⁸¹. Each genome underwent SV calling by six different callers, and merging across callers was done in the following order, representing most to least accurate breakpoint callers: WhamG⁸², Lumpy⁸³, Delly²⁸⁰, GenomeSTRiP⁸⁴, dCGH⁸⁵, and CNVnator⁸⁶.

De novo burden and enrichment.

DNMs were integrated from CCDG genomes (SSC, SAGE, and TASC) in this study, unpublished SPARK exomes (SPARK_WES_1), and three other major published autism exome or genome studies. DNMs, if on hg38, were lifted over to hg19 to enable a merged DNM set. DNMs were restricted to a high-coverage (average $> 20x$) coding regions⁸⁷ to combine exome and genome datasets. We also removed mutations that fell within known segmental duplication regions and known recent repeat and low-complexity regions. Sample duplicates, in cohorts like CCDG genomes and SPARK where the underlying sequencing data are available, were identified using King software to estimate pairwise relatedness between samples. Any samples with a kinship value > 0.35 were considered identical and counted only once. Identical samples from the same cohort were also checked for reported monozygotic twin status. Note, samples in SPARK that overlapped with SSC samples were already removed in the final release by the SPARK Consortium. For other published cohorts where the underlying exome data are unavailable, we relied on the published studies to eliminate within-study overlap. For example, we excluded any potential sample overlap

across the CCDG genomes in this study with samples in published literature. This included all SSC samples in the ASC study¹⁵ that overlap with CCDG SSC genomes. We also excluded CCDG TASC genomes ($n = 246$) to avoid overlap with the TASC samples in the ASC study and only retained TASC samples in the ASC study ($n = 855$) as it has a larger sample size. We further excluded samples with ten or more coding DNMs, and removed DNMs seen in five or more different individuals after above filtering. These measures yielded a total of 15,182 unique ASD trios in the integrated *de novo* enrichment analysis (Supplementary Table 4). Annotation by VEP (Ensembl GRCh37 release 94)⁸⁸ and Combined Annotation Dependent Depletion (CADD) score (v1.3)⁸⁹ were applied to ensure uniformity, and the analysis was restricted to the canonical transcript with the most deleterious annotation. A chimpanzee–human divergence model (CH model)²⁶ and denovolyzeR^{27,28} were used to identify genes with an excess of DNMs. Genes were considered to have an excess of DNMs if both models were significant after multiple test correct (Benjamini-Hochberg FDR < 5% or Bonferroni).

Transmission bias and burden.

We partition variants from protein-encoding regions of the genome into three classes: (1) likely gene disrupting, which we define as any mutation that introduces a stop codon, ablates a stop, changes the frame of the open reading frame, or introduces a change at a predicted splice donor or splice acceptor site; (2) missense, which is any mutation that causes an amino acid change, or (3) synonymous, or any mutation that results in no amino acid change. We quantified the number of private, transmitted variants observed in probands and unaffected siblings by gene set and variant type and compared the proportion of carriers using both a Fisher's exact test and logistic regression (one model for each variant type and pLI threshold). For the DNM-enriched gene set analyses, we compared the proportion of carriers between probands and siblings using Fisher's exact test and logistic regression (Supplementary Table 22). Multi-hit analyses and simplex versus multiplex analyses were conducted using a Fisher's exact test to compare the proportion of individuals carrying two or more hits in probands and siblings. We applied Bonferroni and FDR corrections to all *P*-values using the R function `p.adjust` for each analysis.

Rare variant transmission disequilibrium test (rvTDT).

To perform the rvTDT, each child and their parents represented a separate trio. We used the following formula from He et al.⁹⁰ to compare the rate of transmitted variants within each set of genes to the expected transmission rate of 0.5. We applied the rvTDT separately to affected and unaffected children. Finally, we calculated an approximate odds ratio based on the fraction of observed transmissions over the expected number of 0.5 transmissions if there was no bias.

$b = \#$ children with reference allele transmitted

$c = \#$ children with alternate allele transmitted

$$rvTDT: \chi^2 = (b - c)^2 / (b + c), df = 1$$

Population attributable risk (PAR).

PARs were calculated using the following formula published by Cole and MacMahon⁹¹. Our calculations assume that siblings are representative of the general population. Even if siblings are sub-threshold for ASD, these estimates would serve as a lower bound for PAR.

P_e = proportion of population (controls) exposed

RR = relative risk or ratio of the risk between the exposed and unexposed

$$PAR(\%) = \frac{P_e \times (RR - 1)}{(1 + P_e \times (RR - 1))}$$

Polygenic risk scores (PRS).

PRS were calculated using the additive model implemented in PLINK (v1.9)⁶⁸. Briefly, genome-wide association study summary statistics from Grove et al.¹⁷ were linkage disequilibrium pruned, and variants with a study $P < 0.01$ were retained for scoring. Each genotype in a sample was weighted with the variant's odds ratio, and all of the weighted variants were summed together into a PRS (Supplementary Table 21). To ensure risk scores are comparable across studies, each cohort was quantile normalized before combining across cohorts.

Polygenic transmission disequilibrium test (pTDT).

To estimate the burden of common variation in probands and siblings, we performed a pTDT as described in Weiner et al.¹⁸. We used the following formulas to calculate the pTDT:

n = number of trios/cases

$$\begin{aligned} PRS_{MP} &= PRS_{mother} + PRS_{father} / 2 \\ pTDT \text{ deviation} &= PRS_c - PRS_{MP} / SD(PRS_{MP}) \\ t_{pTDT} &= \text{mean}(pTDT \text{ deviation}) / SD(pTDT \text{ deviation}) / \sqrt{n} \end{aligned}$$

P -values were calculated using a two-sided, one-sample t -test with the null hypothesis defined as the mean pTDT deviation equal to zero. All P -values were Bonferroni corrected for 14 tests and two conditions (probands and siblings).

Allele age estimation.

We estimate the age of a private, transmitted variant at a site of interest using the software Relate⁴³. In short, Relate reconstructs the local genealogy of the region of interest using a scalable computation, which guarantees the inferred genealogy exactly producing the observed data under the infinite-sites model and, thus, can effectively apply to datasets with thousands of individuals while taking into account recombination. Mutations are then mapped onto the branches of the resulting local tree to estimate mutation age. Each private, transmitted variant is annotated as LGD, synonymous, or intergenic as described above and then classified into one of the three datasets for candidate gene, affected proband, and unaffected sibling depending on the carrier of the private variant. To reduce the

computational burden for the inferences for synonymous and intergenic sets, we divided the genome into 100,000-bp windows and randomly selected a locus of interest from up to 550 windows (up to 25 windows per chromosome). We removed sites where the derived allele could not be determined and only included alleles carried by individuals of European descent. For the 163 candidate genes, we included sites within genes where there is a transmitted event within only one family, as well as genes where there are transmitted events of different variants across multiple families. Additionally, we required the genes to carry private LGD variants to have been transmitted only to autism children and previously associated with ASD. This left us with 101 variants for this gene set.

For a site of interest, we first generated phased haplotypes for the 100,000-bp region surrounding the site of interest for all samples using BEAGLE (v5.1)⁹² without imputation. For our analysis, we included only one individual from each family, where the family was of European descent, and removed samples that are likely related (see Methods for kinship above). To model recombination wherever applicable, we used the HapMap genetic map. Ancestral and derived states for individual sites are based on the sequences downloaded from Ensembl. To ensure the quality of the genotypes, we masked sites that are within segmental duplications, low-complexity, and repeat masker sequences (see Methods for private SNV calls above). The software Relate outputs two age estimates: the lower and upper ages represent the ages of the coalescence events below and above the mutation of interest, respectively. We determine the age of a given variant by taking the average of the two estimates. Note that we recognize the presence of natural selection at sites where deleterious mutations occur would affect the inference of allele age, and thus, the age estimates for deleterious alleles inferred in this study are overestimated and can be deemed as an upper bound for their ages. We compare the distributions of allele age among different datasets using the Mann-Whitney U test, and *P*-values are all Bonferroni corrected using *p.adjust*. Potential caveats, such as phasing errors and cryptic relatedness, might affect the individual age estimates but are expected to have limited impacts on the observation of differences among the variant sets because the same procedure was applied to individual sets.

Selection coefficient estimation.

We apply the classic mutation-selection balance model to estimate selection coefficients for the 101 private, transmitted variants in the candidate gene set. Our rationale is that because these variants are predominantly recent (on the order of $1/s$)⁴², selection must act relatively strong on these variants, and thus, the presence of these variants are primarily due to mutation. Assuming the multi-hit, additive model ($h = 0.5$), the individual selection coefficient (s) can be approximated as μ/qh , where μ is gene-specific mutation rate, h is the dominance coefficient, and q is the observed allele frequency in the entire population sample.

Gene expression analyses.

Cell-type specific expression analyses (CSEA) were conducted using the CSEA tool⁹³. Candidate genes from probands and siblings were uploaded to the available web server for CSEA across brain regions and development in humans. Gene expression was identified in a

published set of transcriptomically defined cell types in the human temporal cortex⁹⁴. Gene sets were tested for enriched expression in three broad cell classes—inhibitory neurons, excitatory neurons, and non-neuronal cells—by counting the number of cell types within each class that expressed each gene with average counts per million greater than 1. For each gene set, the number of cell types with expression were calculated, and gene sets were visually compared by plotting cumulative distributions. For each cell class, Wilcoxon rank-sum tests were used to identify statistically significant differences in the number of cell types with expression for ASD and control gene sets. *P*-values were adjusted for multiple comparisons using Bonferroni correction.

PPI network analysis.

We used the STRING database (STRING, v11)⁹⁵ to perform PPI network analyses via Cytoscape (v3.7.2)^{96,97}. We used the multiple protein input option with all default settings except we required interactions to be limited to those which were high confidence (0.700). Disconnected nodes were hidden from the network output (but not from enrichment analyses). In addition, we also used STRING to calculate network statistics and run functional enrichment analyses from the Gene Ontology resource, KEGG, and the Reactome pathway database to identify shared functions across the full set of genes and three subnetworks. A subnetwork was identified as any group of genes that contained at least five genes. The genes *KIAA0430* and *ATP5B* are also known as *MARF1* and *ATP5F1B*, respectively.

Statistics.

All statistics were calculated using R versions 3.5.1 and 3.5.2.

DATA AVAILABILITY

The WGS data used in this study are available from the following resources. The AGRE study is available at the Database of Genotypes and Phenotypes (dbGaP) under accession phs001766. Access to the AGRE WGS data is subject to approval by Autism Speaks and AGRE. All sequencing and phenotype data for the SSC are available through the Simons Foundation for Autism Research Initiative (SFARI) and are available to approved researchers at SFARI Base (<http://base.sfari.org>, accession IDs: SFARI_SSC_WGS_p, SFARI_SSC_WGS_1, and SFARI_SSC_WGS_2). The genomic and phenotypic data for the SPARK study are available by request from SFARI Base (<http://base.sfari.org>, accession ID: SFARI_SPARK_WES_1). Data from the SAGE study are available at dbGaP under accession phs001740.v1.p1. Data from the TASC study are available at dbGaP under accession phs001741. Family-level FreeBayes and GATK VCF files for SAGE, SSC, and TASC are available at dbGaP accession phs001874.v1.p1 and also at SFARI Base under accession SFARI_SSC_WGS_2a

CODE AVAILABILITY

All software used in this study is publicly available. Code for the ultra-rare transmitted variant pipeline can be found here: https://github.com/EichlerLab/ultra_rare_transmitted.git. Code for figures and analyses are available upon request.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

ACKNOWLEDGEMENTS

We thank Tonia Brown for assistance in editing this manuscript and Sunday Stray, Mary Eng, James Moore, Hannah Kortbawi, and Anne Thornton from the laboratory of Mary-Claire King for isolation of DNA from whole blood. We thank Tom Maniatis and the New York Genome Center for conducting sequencing and initial QC. We are grateful to all of the families at the participating Simons Simplex Collection (SSC) sites, as well as the principal investigators (A. Beaudet, R. Bernier, J. Constantino, E. Cook, E. Fombonne, D. Geschwind, R. Goin-Kochel, E. Hanson, D. Grice, A. Klin, D. Ledbetter, C. Lord, C. Martin, D. Martin, R. Maxim, J. Miles, O. Ousley, K. Pelphrey, B. Peterson, J. Piggot, C. Saulnier, M. State, W. Stone, J. Sutcliffe, C. Walsh, Z. Warren, E. Wijsman). We are grateful to all of the families in SPARK, the SPARK clinical sites, and SPARK staff. We appreciate obtaining access to phenotypic and genetic data on SFARI Base. Approved researchers can obtain the SSC population dataset described in this study (<https://www.sfari.org/resource/simons-simplex-collection/>) and SPARK population dataset described in this study (<https://www.sfari.org/resource/spark/>) by applying at <https://base.sfari.org>. We gratefully acknowledge the resources provided by the Autism Genetic Resource Exchange (AGRE) Consortium and the participating AGRE families. Genomic data for the AGRE cohort was provided by iHART, an initiative led by Hartwell Foundation and Directed by D. Wall and D. Geschwind. This work was supported, in part, by grants from the US National Institutes of Health (NIH R01 MH101221 to E.E.E.; R01 MH100047 to R.A.B.; K99 MH117165 to T.N.T.; K99 HG011041 to P.H.; UM1 HG008901 to M.C.Z.) and the Simons Foundation (SFARI 608045 to E.E.E.). The CCDG is funded by the National Human Genome Research Institute and the National Heart, Lung, and Blood Institute. The GSP Coordinating Center (U24 HG008956) contributed to cross-program scientific initiatives and provided logistical and general study coordination. AGRE is a program of Autism Speaks and is supported in part by grant 1U24MH081810 from the National Institute of Mental Health to C.M. Lajonchere. E.E.E. is an investigator of the Howard Hughes Medical Institute.

Appendix

The SPARK Consortium

Xueya Zhou^{8,9}, Tianyun Wang¹, Pamela Feliciano¹⁰, Jacob Hall¹⁰, Irina Astrovskaya¹⁰, Shwetha C. Murali^{1,2}, Simon Xu¹⁰, Chang Shu^{8,9}, Joseph Obiajulu^{8,9}, Leo Brueggeman¹¹, Jessica Wright¹⁰, Olena Marchenko¹⁰, Chris Fleisch¹⁰, Timothy S. Chang^{12,13}, LeeAnne Green Snyder¹⁰, Sarah D. Barns¹⁰, Tychele N. Turner^{1,†}, Bing Han¹⁰, William Harvey¹, Andrew Nishida¹⁴, Ryan Doan^{15,16}, Aubrey Soucy^{15,16}, Brian J. O’Roak¹⁴, Timothy W. Yu^{15,16,17}, Daniel Geschwind^{12,13,18}, Jacob Michaelson¹¹, Natalia Volfovsky¹⁰, Evan E. Eichler^{1,2}, Yufeng Shen⁹, and Wendy K. Chung^{8,10}

⁸Department of Pediatrics, Columbia Medical Center, New York, NY, USA.

⁹Department of Systems Biology, Columbia University, New York, NY, USA.

¹⁰Simons Foundation, New York, NY, USA.

¹¹Department of Psychiatry, University of Iowa Carver College of Medicine, Iowa City, IA, USA.

- ¹²Program in Neurogenetics, David Geffen School of Medicine, Los Angeles, CA, USA.
- ¹³Department of Neurology, University of California Los Angeles, Los Angeles, CA, USA.
- ¹⁴Department of Molecular & Medical Genetics, Oregon Health & Science University, Portland, OR, USA.
- ¹⁵Division of Genetics and Genomics, Boston Children's Hospital, Boston, MA, USA.
- ¹⁶Department of Pediatrics, Harvard Medical School, Boston, MA, USA.
- ¹⁷The Broad Institute, Cambridge, MA, USA.
- ¹⁸Department of Human Genetics, University of California Los Angeles, Los Angeles, CA, USA.

REFERENCES

1. Baio J et al. Prevalence of autism spectrum disorder among children aged 8 years – Autism and Developmental Disabilities Monitoring Network, 11 Sites, United States, 2014. *MMWR Surveill. Summ*67, 1–23 (2018).
2. Iossifov J et al. The contribution of de novo coding mutations to autism spectrum disorder. *Nature*515, 216–221 (2014). [PubMed: 25363768]
3. Krumm N et al. Excess of rare, inherited truncating mutations in autism. *Nat. Genet*47, 582–588 (2015). [PubMed: 25961944]
4. Turner T N et al. Genomic patterns of de novo mutation in simplex autism. *Cell*171, 710–722 e712 (2017). [PubMed: 28965761]
5. Gaugler T et al. Most genetic risk for autism resides with common variation. *Nat. Genet*46, 881–885 (2014). [PubMed: 25038753]
6. Constantino J N et al. Autism recurrence in half siblings: strong support for genetic mechanisms of transmission in ASD. *Mol. Psychiatry*18, 137–138 (2013). [PubMed: 22371046]
7. Ganna A et al. Quantifying the impact of rare and ultra-rare coding variation across the phenotypic spectrum. *Am. J. Hum. Genet*102, 1204–1211 (2018). [PubMed: 29861106]
8. Ruzzo E K et al. Inherited and de novo genetic risk for autism impacts shared networks. *Cell*178, 850–866 e826 (2019). [PubMed: 31398340]
9. De Rubeis S et al. Synaptic, transcriptional and chromatin genes disrupted in autism. *Nature*515, 209–215 (2014). [PubMed: 25363760]
10. Sanders S J et al. Insights into autism spectrum disorder genomic architecture and biology from 71 risk loci. *Neuron*87, 1215–1233 (2015). [PubMed: 26402605]
11. Satterstrom F K et al. Large-scale exome sequencing study implicates both developmental and functional changes in the neurobiology of autism. *Cell*180, 568–584 e523 (2020). [PubMed: 31981491]
12. Satterstrom F K et al. Autism spectrum disorder and attention deficit hyperactivity disorder have a similar burden of rare protein-truncating variants. *Nat. Neurosci*22, 1961–1965 (2019). [PubMed: 31768057]
13. Schaaf C P et al. Oligogenic heterozygosity in individuals with high-functioning autism spectrum disorders. *Hum. Mol. Genet*20, 3366–3375 (2011). [PubMed: 21624971]
14. Girirajan S et al. A recurrent 16p12.1 microdeletion supports a two-hit model for severe developmental delay. *Nat. Genet*42, 203–209 (2010). [PubMed: 20154674]
15. Du Y et al. Nonrandom occurrence of multiple de novo coding variants in a proband indicates the existence of an oligogenic model in autism. *Genet. Med*22, 170–180 (2019). [PubMed: 31332282]

16. Jiang YHet al. A mixed epigenetic/genetic model for oligogenic inheritance of autism with a limited role for UBE3A. *Am. J. Med. Genet. A*131, 1–10 (2004). [PubMed: 15389703]
17. Grove Jet al. Identification of common genetic risk variants for autism spectrum disorder. *Nat. Genet*51, 431–444 (2019). [PubMed: 30804558]
18. Weiner DJet al. Polygenic transmission disequilibrium confirms that common and rare variation act additively to create risk for autism spectrum disorders. *Nat. Genet*49, 978–985 (2017). [PubMed: 28504703]
19. Turner TNet al. Genome sequencing of autism-affected families reveals disruption of putative noncoding regulatory DNA. *Am. J. Hum. Genet*98, 58–74 (2016). [PubMed: 26749308]
20. Fischbach GD & Lord C The Simons Simplex Collection: a resource for identification of autism genetic risk factors. *Neuron* 68, 192–195 (2010). [PubMed: 20955926]
21. Guo Het al. Genome sequencing identifies multiple deleterious variants in autism patients with more severe phenotypes. *Genet. Med*21, 1611–1620 (2019). [PubMed: 30504930]
22. An JY et al. Genome-wide de novo risk score implicates promoter variation in autism spectrum disorder. *Science*362, eaat6576 (2018). [PubMed: 30545852]
23. Buxbaum JDet al. The Autism Simplex Collection: an international, expertly phenotyped autism sample for genetic and phenotypic analyses. *Mol. Autism*5, 34 (2014). [PubMed: 25392729]
24. Feliciano Pet al. Exome sequencing of 457 autism families recruited online provides evidence for autism risk genes. *NPJ Genom. Med*4, 19 (2019). [PubMed: 31452935]
25. SPARK Consortium. SPARK: a US cohort of 50,000 families to accelerate autism research. *Neuron*97, 488–493 (2018). [PubMed: 29420931]
26. O'Roak BJet al. Sporadic autism exomes reveal a highly interconnected protein network of de novo mutations. *Nature*485, 246–250 (2012). [PubMed: 22495309]
27. Samocha KE et al. A framework for the interpretation of de novo mutation in human disease. *Nat. Genet*46, 944–950 (2014). [PubMed: 25086666]
28. Ware JS, Samocha KE, Homsy J & Daly MJ Interpreting de novo variation in human disease using denovolyzeR. *Curr. Protoc. Hum. Genet* 87, 7.25.1–7.25.15 (2015). [PubMed: 26439716]
29. Snijders Blok Let al. De novo mutations in MED13, a component of the Mediator complex, are associated with a novel neurodevelopmental disorder. *Hum. Genet*137, 375–388 (2018). [PubMed: 29740699]
30. Shah AA et al. Excess of RALGAPB de novo variants in neurodevelopmental disorders. *Eur. J. Med. Genet*63, 104041 (2020). [PubMed: 32853829]
31. Sapio MR et al. Novel carboxypeptidase A6 (CPA6) mutations identified in patients with juvenile myoclonic and generalized epilepsy. *PLoS One*10, e0123180 (2015). [PubMed: 25875328]
32. Li QS, Parrado AR, Samtani MN, Narayan VA & Alzheimer's Disease Neuroimaging Initiative. Variations in the FRA10AC1 fragile site and 15q21 are associated with cerebrospinal fluid Aβ1–42 level. *PLoS One* 10, e0134000 (2015). [PubMed: 26252872]
33. Siitonen A et al. Genetics of early-onset Parkinson's disease in Finland: exome sequencing and genome-wide association study. *Neurobiol. Aging*53, 195.e7–195.e110 (2017).
34. Gravel Set al. Demographic history and rare allele sharing among human populations. *Proc. Natl. Acad. Sci. USA*108, 11983–11988 (2011). [PubMed: 21730125]
35. Genomes Project Consortium et al. A global reference for human genetic variation. *Nature*526, 68–74 (2015). [PubMed: 26432245]
36. Tennesen JA et al. Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *Science*337, 64–69 (2012). [PubMed: 22604720]
37. Iossifov I et al. Low load for disruptive mutations in autism genes and their biased transmission. *Proc. Natl. Acad. Sci. USA*112, E5600–5607 (2015). [PubMed: 26401017]
38. Epi25 Collaborative. Ultra-rare genetic variation in the epilepsies: a whole-exome sequencing study of 17,606 individuals. *Am. J. Hum. Genet*105, 267–282 (2019). [PubMed: 31327507]
39. Coe BP et al. Neurodevelopmental disease genes implicated by de novo mutation and copy number variation morbidity. *Nat. Genet*51, 106–116 (2019). [PubMed: 30559488]
40. Sebat Jet al. Strong association of de novo copy number mutations with autism. *Science*316, 445–449 (2007). [PubMed: 17363630]

41. He X et al. Integrated model of de novo and inherited genetic variants yields greater power to identify risk genes. *PLoS Genet.* 9, e1003671 (2013). [PubMed: 23966865]
42. Maruyama T The age of a rare mutant gene in a large population. *Am. J. Hum. Genet* 26, 669–673 (1974). [PubMed: 4440678]
43. Speidel L, Forest M, Shi S & Myers SR A method for genome-wide genealogy estimation for thousands of samples. *Nat. Genet* 51, 1321–1329 (2019). [PubMed: 31477933]
44. Pinto D et al. Convergence of genes and cellular pathways dysregulated in autism spectrum disorders. *Am. J. Hum. Genet* 94, 677–694 (2014). [PubMed: 24768552]
45. Deardorff M A et al. HDAC8 mutations in Cornelia de Lange syndrome affect the cohesin acetylation cycle. *Nature* 489, 313–317 (2012). [PubMed: 22885700]
46. Williams S R et al. Haploinsufficiency of HDAC4 causes brachydactyly mental retardation syndrome, with brachydactyly type E, developmental delays, and behavioral problems. *Am. J. Hum. Genet* 87, 219–228 (2010). [PubMed: 20691407]
47. Bernier R et al. Disruptive CHD8 mutations define a subtype of autism early in development. *Cell* 158, 263–276 (2014). [PubMed: 24998929]
48. King I F et al. Topoisomerases facilitate transcription of long genes linked to autism. *Nature* 501, 58–62 (2013). [PubMed: 23995680]
49. Sanders S J et al. Multiple recurrent de novo CNVs, including duplications of the 7q11.23 Williams syndrome region, are strongly associated with autism. *Neuron* 70, 863–885 (2011). [PubMed: 21658581]
50. Glessner J T et al. Autism genome-wide copy number variation reveals ubiquitin and neuronal genes. *Nature* 459, 569–573 (2009). [PubMed: 19404257]
51. Fairless R et al. Polarized targeting of neurexins to synapses is regulated by their C-terminal sequences. *J. Neurosci* 28, 12969–12981 (2008). [PubMed: 19036990]
52. Gromova K V et al. Neurobeachin and the kinesin KIF21B are critical for endocytic recycling of NMDA receptors and regulate social behavior. *Cell Rep.* 23, 2705–2717 (2018). [PubMed: 29847800]
53. Tomaselli P J et al. A de novo dominant mutation in KIF1A associated with axonal neuropathy, spasticity and autism spectrum disorder. *J. Peripher. Nerv. Syst* 22, 460–463 (2017). [PubMed: 28834584]
54. The Gene Ontology Consortium. The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic Acids Res.* 47, D330–D338 (2019). [PubMed: 30395331]
55. Girirajan S et al. Phenotypic heterogeneity of genomic disorders and rare copy-number variants. *N. Engl. J. Med* 367, 1321–1331 (2012). [PubMed: 22970919]
56. Stessman H A, Bernier R & Eichler E E A genotype-first approach to defining the subtypes of a complex disease. *Cell* 156, 872–877 (2014). [PubMed: 24581488]
57. Epi4K Consortium & Epilepsy Phenome/Genome Project. Ultra-rare genetic variation in common epilepsies: a case-control sequencing study. *Lancet Neurol.* 16, 135–143 (2017). [PubMed: 28102150]

METHODS-ONLY REFERENCES

58. Regier A A et al. Functional equivalence of genome sequencing analysis pipelines enables harmonized variant calling across human genetics projects. *Nat. Commun* 9, 4038 (2018). [PubMed: 30279509]
59. Li H & Durbin R Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25, 1754–1760 (2009). [PubMed: 19451168]
60. McKenna A et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20, 1297–1303 (2010). [PubMed: 20644199]
61. Li H et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078–2079 (2009). [PubMed: 19505943]
62. Poplin R et al. Scaling accurate genetic variant discovery to tens of thousands of samples. *bioRxiv*, 10.1101/201178 (2018).

63. Garrison E & Marth G Haplotype-based variant detection from short-read sequencing. arXiv preprint arXiv:1207.3907 (2012).
64. Rimmer A et al. Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nat. Genet* 46, 912–918 (2014). [PubMed: 25017105]
65. Li H A statistical framework for SNP calling, mutation discovery, association mapping and population genetic parameter estimation from sequencing data. *Bioinformatics* 27, 2987–2993 (2011). [PubMed: 21903627]
66. Manichaikul A et al. Robust relationship inference in genome-wide association studies. *Bioinformatics* 26, 2867–2873 (2010). [PubMed: 20926424]
67. Hsieh P et al. Adaptive archaic introgression of copy number variants and the discovery of previously unknown human genes. *Science* 366, eaax2083 (2019). [PubMed: 31624180]
68. Chang C C et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 4, 7 (2015). [PubMed: 25722852]
69. Patterson N, Price A L & Reich D Population structure and eigenanalysis. *PLoS Genet.* 2, e190 (2006). [PubMed: 17194218]
70. Price A L et al. Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet* 38, 904–909 (2006). [PubMed: 16862161]
71. Alexander D H, Novembre J & Lange K Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* 19, 1655–1664 (2009). [PubMed: 19648217]
72. Cingolani P et al. Using *Drosophila melanogaster* as a model for genotoxic chemical mutational studies with a new program, SnpSift. *Front. Genet* 3, 35 (2012). [PubMed: 22435069]
73. Cingolani P et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* 6, 80–92 (2012). [PubMed: 22728672]
74. Quinlan A R B E D Tools: the Swiss-Army tool for genome feature analysis. *Curr. Protoc. Bioinformatics* 47, 11.12.1–34 (2014).
75. Quinlan A R & Hall I M B E D Tools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26, 841–842 (2010). [PubMed: 20110278]
76. Krumm N et al. Copy number variation detection and genotyping from exome sequence data. *Genome Res.* 22, 1525–1532 (2012). [PubMed: 22585873]
77. Fromer M et al. Discovery and statistical genotyping of copy-number variation from whole-exome sequencing depth. *Am. J. Hum. Genet* 91, 597–607 (2012). [PubMed: 23040492]
78. Wang K et al. PennCNV: an integrated hidden Markov model designed for high-resolution copy number variation detection in whole-genome SNP genotyping data. *Genome Res.* 17, 1665–1674 (2007). [PubMed: 17921354]
79. Scharpf R B, Irizarry R A, Ritchie M E, Carvalho B & Ruczinski I Using the R package crrmm for genotyping and copy number estimation. *J. Stat. Softw.* 40, 1–32 (2011).
80. Rausch T et al. DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics* 28, i333–i339 (2012). [PubMed: 22962449]
81. Hormozdiari F et al. Next-generation VariationHunter: combinatorial algorithms for transposon insertion discovery. *Bioinformatics* 26, i350–357 (2010). [PubMed: 20529927]
82. Kronenberg Z N et al. Wham: identifying structural variants of biological consequence. *PLoS Comput. Biol.* 11, e1004572 (2015). [PubMed: 26625158]
83. Layer R M, Chiang C, Quinlan A R & Hall I M LUMPY: a probabilistic framework for structural variant discovery. *Genome Biol.* 15, R84 (2014). [PubMed: 24970577]
84. Handsaker R E et al. Large multi-allelic copy number variations in humans. *Nat. Genet* 47, 296–303 (2015). [PubMed: 25621458]
85. Sudmant P H et al. Global diversity, population stratification, and selection of human copy-number variation. *Science* 349, aab3761 (2015). [PubMed: 26249230]
86. Abyzov A, Urban A E, Snyder M & Gerstein M CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs from family and population genome sequencing. *Genome Res.* 21, 974–984 (2011). [PubMed: 21324876]

87. Turner TNet al. Sex-based analysis of de novo variants in neurodevelopmental disorders. *Am. J. Hum. Genet* 105, 1274–1285 (2019). [PubMed: 31785789]
88. McLaren Wet al. The Ensembl Variant Effect Predictor. *Genome Biol.* 17, 122 (2016). [PubMed: 27268795]
89. Kircher Met al. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet* 46, 310–315 (2014). [PubMed: 24487276]
90. He Zet al. Rare-variant extensions of the transmission disequilibrium test: application to autism exome sequence data. *Am. J. Hum. Genet* 94, 33–46 (2014). [PubMed: 24360806]
91. Cole P & MacMahon B. Attributable risk percent in case-control studies. *Br. J. Prev. Soc. Med* 25, 242–244 (1971). [PubMed: 5160433]
92. Browning SR & Browning BL. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am. J. Hum. Genet* 81, 1084–1097 (2007). [PubMed: 17924348]
93. Dougherty JD, Schmidt EF, Nakajima M & Heintz N. Analytical approaches to RNA profiling data for the identification of genes enriched in specific cells. *Nucleic Acids Res.* 38, 4218–4230 (2010). [PubMed: 20308160]
94. Hodge RDet al. Conserved cell types with divergent features in human versus mouse cortex. *Nature* 573, 61–68 (2019). [PubMed: 31435019]
95. Szklarczyk Det al. STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res* 47, D607–D613 (2019). [PubMed: 30476243]
96. Ono K, Muetze T, Kolishovski G, Shannon P & Demchak B. CyREST: turbocharging Cytoscape access for external tools via a RESTful API. *F1000Res.* 4, 478 (2015). [PubMed: 26672762]
97. Shannon Pet al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* 13, 2498–2504 (2003). [PubMed: 14597658]

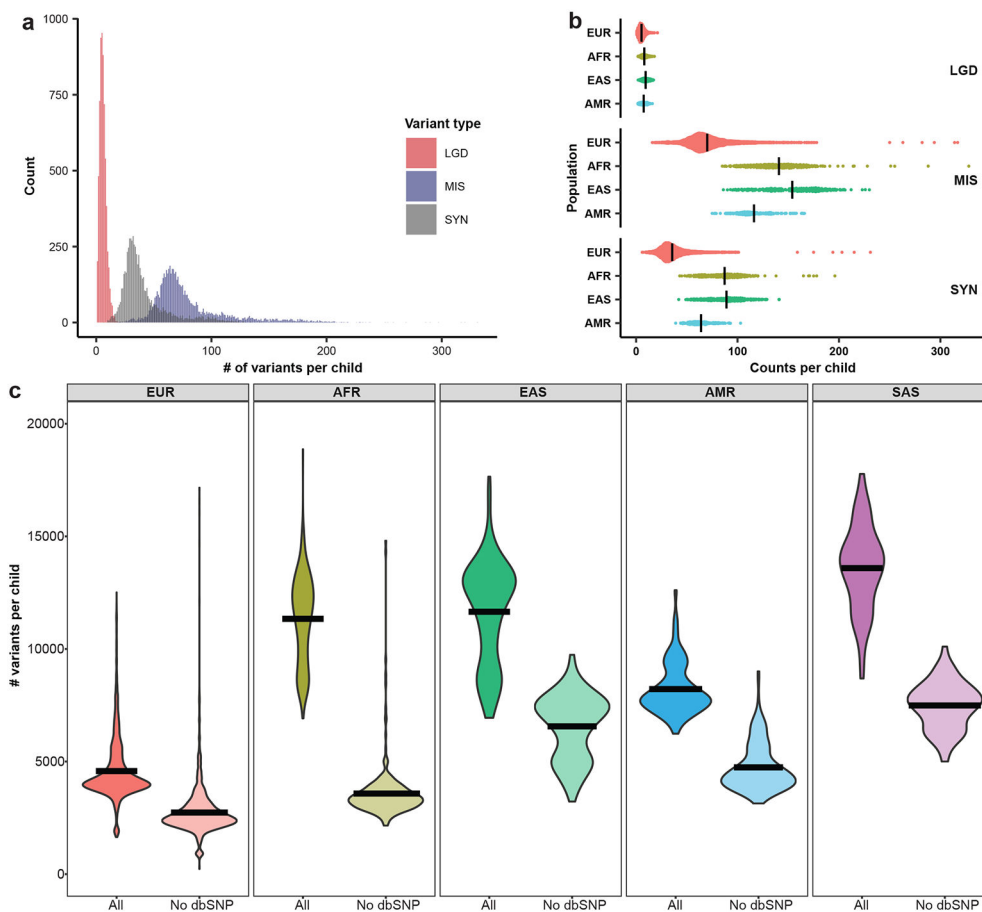


Figure 1 | Overview of private variants in discovery cohort.

Private variants are defined as variants observed in one and only one parent in the cohort.

a, Distribution of likely gene-disruptive (LGD), missense (MIS), and synonymous (SYN) private variants per child (probands and unaffected siblings). **b**, The cumulative number of each variant class by assigned population group (EUR, European ($n = 5,685$); AFR, African ($n = 290$); EAS, East Asian ($n = 252$); AMR, Amerindian ($n = 193$); SAS, South Asian ($n = 103$)), excluding SAS. **c**, Private, transmitted variant counts per child grouped by ancestry before (All) and after (No dbSNP) filtering with dbSNPv150. Excess of private variants is partially but not fully resolved after excluding sites observed in dbSNP. We were unable to assign ancestry to one of these five population groups for 74 of the children in this study. The y-axis was truncated at 20,000 variants per child; however, both the AFR and EUR populations had a small number of children with variant counts above this threshold (see Supplementary Tables 6 and 7 for details). Black lines indicate the average variant count per population in **b** and **c**.

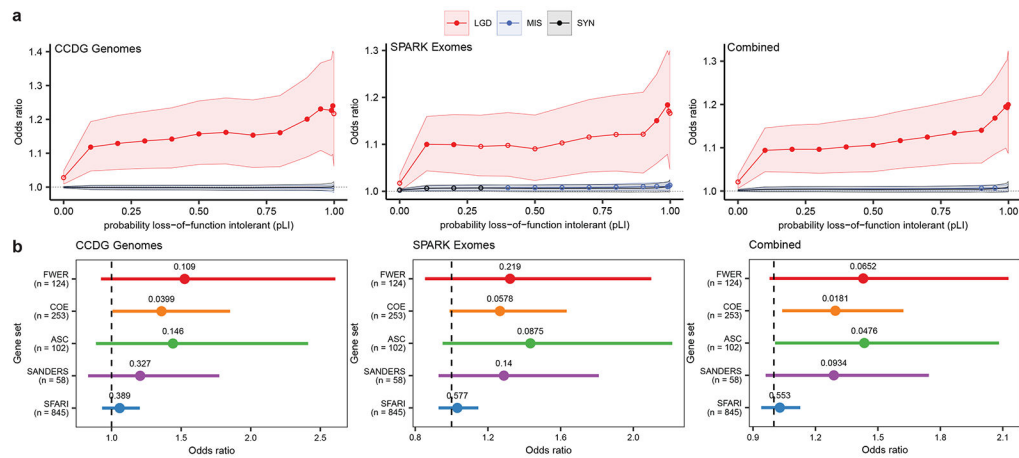


Figure 2 | Burden of private LGD variants in affected children.

a, Burden of private LGD variants in probands as compared to siblings was quantified (odds ratio (OR)) at increasing thresholds of gene constraint (pLI) in our discovery ($n = 4,201$ affected and 2,191 unaffected children), replication ($n = 6,453$ affected and 3,007 unaffected children), and combined discovery and replication ($n = 10,657$ affected and 5,199 unaffected children) cohorts. Filled circles indicate Bonferroni-corrected $P < 0.05$ (42 tests per cohort), unfilled circles indicate nominal $P < 0.05$, and shaded areas indicate 95% confidence intervals around the OR estimate. OR and confidence intervals were calculated using logistic regression (see Supplementary Table 11 for details). **b**, Enrichment of private, LGD variant transmission to probands for five autism risk gene sets (FWER, COE, ASC, SANDERS, SFARI). With the exception of SFARI, most gene sets were identified based on an excess of *de novo* mutations (DNMs) in parent-child trios (see Online Methods). OR was based on a comparison of the proportion of carriers between probands and siblings in our discovery ($n = 4,201$ affected and 2,191 unaffected children), replication ($n = 6,453$ affected and 3,007 unaffected children), and combined ($n = 10,657$ affected and 5,199 unaffected children) cohorts using a two-sided Fisher’s exact test (see Supplementary Table 5 for details). Dashed black line indicates OR = 1, which represents no difference between probands and siblings. Families with monozygotic twins ($n = 75$ in discovery, $n = 63$ in replication, and $n = 138$ in combined) were removed from analysis. For the combined set, variants were restricted to regions with at least 20x average coverage in the exomes. Reported P -values are nominal, points indicate the OR estimate, and error bars indicate 95% confidence intervals around the OR estimate.

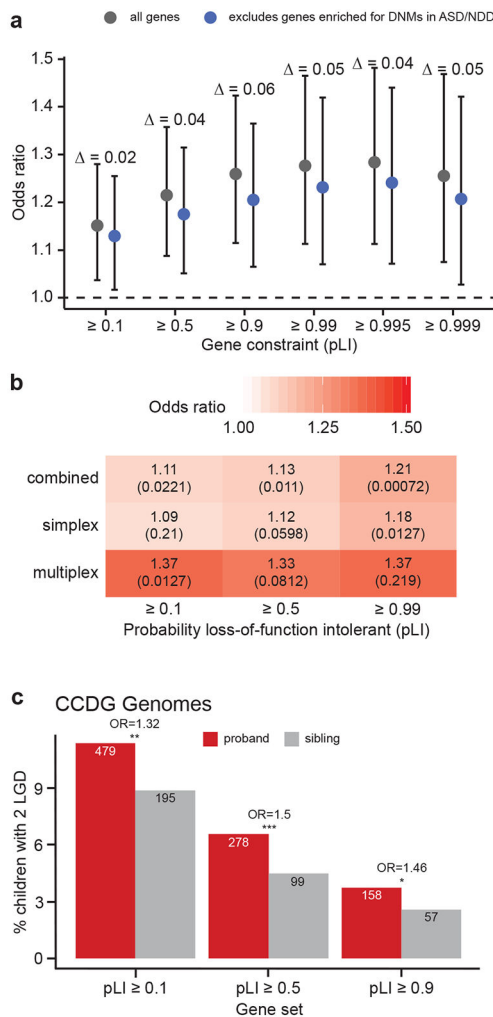


Figure 3 |. Genetic properties of inherited LGD variant burden.

a, At least 95.4% of private, transmitted LGD variant burden resides outside of genes identified with an excess of DNMs in ASD/NDD cases (321 genes considered and 154 genes with transmissions) based on analysis of CCDG autism genomes ($n = 4,201$ affected and 2,191 unaffected children). We observe 141 DNM-enriched genes with transmissions to probands and 85 genes with transmissions to siblings (Supplementary Table 12). OR for five cumulative pLI bins were compared before and after excluding DNM-enriched genes in ASD/NDD cases. The percentage of remaining burden is calculated as quotient of the OR for the pLI bin after removing genes enriched for DNMs in ASD/NDD cases and the OR for all genes in that pLI bin. Families with monozygotic twins ($n = 75$) were excluded from this analysis. OR and associated P -values were calculated using a two-sided Fisher’s exact test. Points indicate the OR estimate, and error bars indicate the 95% confidence interval around the OR estimate. **b**, Multiplex families ($n = 1,268$ families, 2,691 probands, 533 siblings) show a higher burden of private, transmitted LGD variants in probands as compared to siblings across three pLI thresholds compared to simplex families ($n = 7,962$ families, 7,962 probands, 4,666 siblings). **c**, We observe a significant enrichment of probands carrying two private, transmitted LGD variants (2 LGD) when compared to unaffected siblings

at various levels of gene constraint (3 cumulative pLI bins considered) based on CCDG genomes sequenced from autism families ($n = 4,201$ probands, 2,191 siblings). Families with monozygotic twins ($n = 75$) were excluded from this analysis. OR was calculated using a two-sided Fisher's exact test, and reported P-values are Bonferroni corrected for nine (**b**) and three (**c**) tests (see Supplementary Tables 7 and 8 for details).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

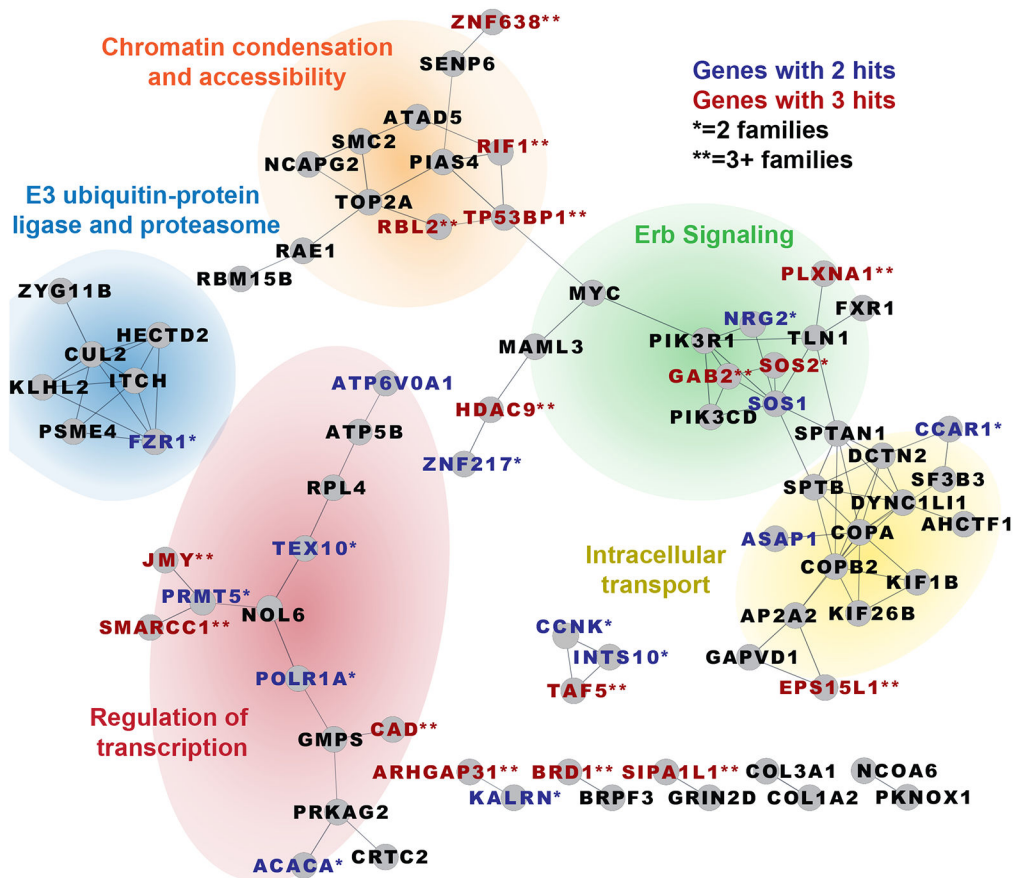


Figure 4 |. PPI network for autism candidate genes.

We identified 163 constrained genes (pLI = 0.99) carrying private LGD variants transmitted only to autism probands based on combined dataset and not previously identified as a DNM-enriched ASD gene (Supplementary Table 9). STRING network shows a significant excess of PPI ($P = 0.00164$). Gene names are colored if observed in two (blue) or three or more (red) probands and labeled if observed in two independent families (*) or more (**). Families with monozygotic twins ($n = 138$) were removed from analysis. Analyses were restricted to regions with at least 20x average coverage in the exomes.

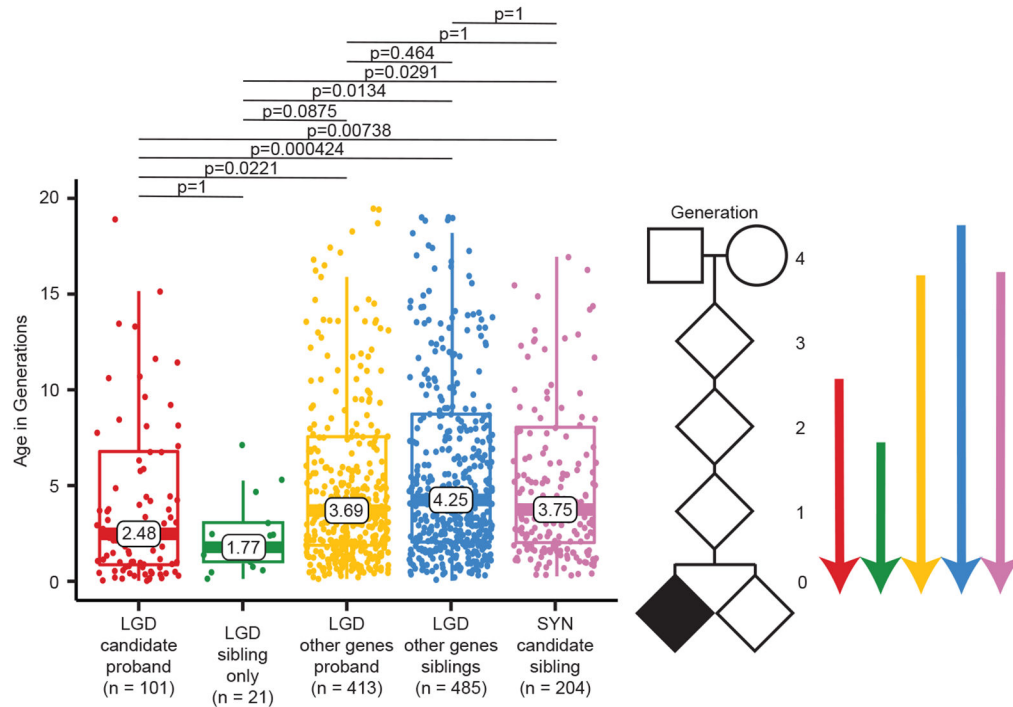


Figure 5 I. Estimate of allele age.

The software Relate was used to estimate the coalescent age (in generations) for private LGD (red) and SYN (blue) variants in 163 candidate genes, private LGD variants (green) in 83 sibling-only genes, and ~500 sites from all remaining genes for European probands (yellow, $n = 3,776$) and siblings (pink, $n = 1,909$). P -values were calculated using a two-sided t -test and Bonferroni corrected for six tests. Plot was truncated at 20 generations. Data points older than this are included in calculating represented statistics (e.g., boxplots, medians, P -values) but are not visualized. To view all data points, see Supplementary Figure 14. Boxplot whiskers represent 1.5 times the upper and lower interquartile ranges. Upper and lower hinges correspond to the 25th and 75th percentiles, and the middle line represents the median. Mean values are noted on the plot.

Table 1 |

Summary of whole-genome and whole-exome sequencing

Cohort	Individuals					Families					Assay
	Genomes	Probands	Siblings	Twin pairs	Simplex	Multiplex	Total	Publication			
SSC	8,757	2,299	1,860	1	2,299	0	2,299	An et al. ²²	WGS		
SAGE	547	202	5	4	144	26	170	Guo et al. ²¹	WGS		
TASC	750	250	0	1	250	0	250	Unpublished	WGS		
AGRE/iHart Phase I *	1,736	822	194	69	6	354	360	Ruzzo et al. ⁸	WGS		
AGRE/iHart Phase II	1,757	791	176	0	1	394	395	Unpublished	WGS		
Discovery set	13,547	4,364	2,235	75	2,700	774	3,474	-	WGS		
SPARK	21,331	6,539	3,034	63	5,278	601	5,879	Unpublished	WES		
Combined	34,880	10,905	5,269	164	7,978	1,375	9,353	-	WGS & WES		

Summary information of the five cohorts included in our study. WGS samples have been sequenced to an average depth of 34x.

* 10 families had additional members sequenced in AGRE/iHart Phase II. These families were still included in the number of simplex/multiplex families for AGRE/iHart phase I; however, the family was identified as simplex or multiplex based on the full family data, which includes the additional members from AGRE/iHart Phase II.

Table 2 |Population attributable risk for *de novo* and private LGD variants

Variant class	Genes enriched for DNMs in ASD/NDD	Remaining genes, pLI 0.99
<i>De novo</i> *	4.39%	1.45%
Private	1.45%	2.64%

Population attributable risk (PAR) percentages were calculated in our discovery cohort for *de novo* and private LGD variants in children (Methods). DNM calculations do not include the AGRE study. We defined the DNM-enriched ASD/NDD gene set as the genes reported in Coe et al.³⁹, Sanders et al.¹⁰, and Satterstrom et al.¹¹

* Does not include AGRE cohort.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript