

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository:<https://orca.cardiff.ac.uk/id/eprint/101274/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Singh, Tarjinder, Walters, James T. R. , Johnstone, Mandy, Curtis, David, Suvisaari, Jaana, Torniaainen, Minna, Rees, Elliott , Iyegbe, Conrad, Blackwood, Douglas, McIntosh, Andrew M., Kirov, George , Geschwind, Daniel, Murray, Robin M, Di Forti, Marta, Bramon, Elvira, Gandal, Michael, Hultman, Christina M., Sklar, Pamela, Palotie, Aarno, Sullivan, Patrick F., O'Donovan, Michael C. , Owen, Michael J. and Barrett, Jeffrey C. 2017. The contribution of rare variants to risk of schizophrenia in individuals with and without intellectual disability. *Nature Genetics* 49 , pp. 1167-1173. 10.1038/ng.3903

Publishers page: <http://dx.doi.org/10.1038/ng.3903>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



# The contribution of rare variants to risk of schizophrenia in individuals with and without intellectual disability

Tarjinder Singh<sup>1</sup>, James T. R. Walters<sup>2</sup>, Mandy Johnstone<sup>3</sup>, David Curtis<sup>4,5</sup>, Jaana Suvisaari<sup>6</sup>, Minna Torniainen<sup>6</sup>, Elliott Rees<sup>2</sup>, Conrad Iyegbe<sup>7</sup>, Douglas Blackwood<sup>3</sup>, Andrew M. McIntosh<sup>8</sup>, Georg Kirov<sup>2</sup>, Daniel Geschwind<sup>9</sup>, Robin M. Murray<sup>7</sup>, Marta Di Forti<sup>7</sup>, Elvira Bramon<sup>10</sup>, Michael Gandal<sup>9</sup>, Christina M. Hultman<sup>11</sup>, Pamela Sklar<sup>12</sup>, INTERVAL Study<sup>13</sup>, UK10K Consortium<sup>13</sup>, Aarno Palotie<sup>14,15</sup>, Patrick F. Sullivan<sup>16,17</sup>, Michael C. O'Donovan<sup>2</sup>, Michael J. Owen<sup>2</sup>, Jeffrey C. Barrett<sup>1</sup>

<sup>1</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton CB10 1HH, Cambridge, UK. <sup>2</sup>MRC Centre for Neuropsychiatric Genetics and Genomics, Division of Psychological Medicine and Clinical Neurosciences, School of Medicine, Cardiff University, Cardiff CF24 4HQ, UK. <sup>3</sup>Division of Psychiatry, The University of Edinburgh, Royal Edinburgh Hospital, Edinburgh EH10 5HF, UK. <sup>4</sup>University College London (UCL) Genetics Institute, University College London, Darwin Building, Gower Street, London WC1E 6BT, UK. <sup>5</sup>Centre for Psychiatry, Barts and the London School of Medicine and Dentistry, London, UK. <sup>6</sup>National Institute for Health and Welfare (THL), Helsinki FI-00271, Finland. <sup>7</sup>Institute of Psychiatry, King's College London, 16 De Crespigny Park, London SE5 8AF, UK. <sup>8</sup>Centre for Cognitive Ageing and Cognitive Epidemiology, The University of Edinburgh, 7 George Square, Edinburgh EH8 9JZ, UK. <sup>9</sup>UCLA David Geffen School of Medicine, Los Angeles, California 90095, USA. <sup>10</sup>Division of Psychiatry, University College London, Charles Bell House, Riding House Street, London W1W 7EJ, UK. <sup>11</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, SE-17177 Stockholm, Sweden. <sup>12</sup>Division of Psychiatric Genomics, Department of Psychiatry, Icahn School of Medicine at Mount Sinai, New York, New York 10029, USA. <sup>13</sup>The members of this consortium are listed in the Supplementary Note. <sup>14</sup>Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki FI-00014 Finland. <sup>15</sup>Program in Medical and Population Genetics and Genetic Analysis Platform, The Broad Institute of MIT and Harvard, Cambridge MA 02132, USA. <sup>16</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, SE-17177 Stockholm, Sweden. <sup>17</sup>Departments of Genetics and Psychiatry, University of North Carolina, Chapel Hill, NC, 27599-7264, USA.

Correspondence should be addressed to Tarjinder Singh (ts14@sanger.ac.uk) and Jeffrey C. Barrett (barrett@sanger.ac.uk).

## Abstract

By meta-analyzing rare coding variants in whole-exome sequences of 4,133 schizophrenia cases and 9,274 controls, de novo mutations in 1,077 trios, and copy number variants from 6,882 cases and 11,255 controls, we show that individuals with schizophrenia carry a significant burden of rare damaging variants in 3,488 genes previously identified as having a near-complete

depletion of loss-of-function variants. In schizophrenia patients who also have intellectual disability, this burden is concentrated in risk genes associated with neurodevelopmental disorders. After excluding known neurodevelopmental disorder risk genes, a significant rare variant burden persists in other loss-of-function intolerant genes, and while this effect is notably stronger in schizophrenia patients with intellectual disability, it is also seen in patients who do not have intellectual disability. Together, our results show that rare damaging variants contribute to the risk of schizophrenia both with and without intellectual disability, and support an overlap of genetic risk between schizophrenia and other neurodevelopmental disorders.

## Introduction

Schizophrenia is a common and debilitating psychiatric illness characterized by positive symptoms (hallucinations, delusions, disorganized speech and behaviour), negative symptoms (social withdrawal and diminished emotional expression), and cognitive impairment that result in social and occupational dysfunction<sup>1,2</sup>. Operational diagnostic criteria for the disorder as described in the DSM-V require the presence of at least two of the core symptoms over a period of six months with at least one month of active symptoms<sup>3</sup>. It is increasingly recognized that current categorical psychiatric classifications have a number of shortcomings, in particular that they overlook the increasing evidence for etiological and mechanistic overlap between psychiatric disorders<sup>4</sup>.

A diverse range of pathophysiological processes may contribute to the clinical features of schizophrenia<sup>5</sup>. Indeed, previous studies have suggested a number of hypotheses about schizophrenia pathogenesis, including abnormal pre-synaptic dopaminergic activity<sup>6</sup>, postsynaptic mechanisms involved in synaptic plasticity<sup>7</sup>, dysregulation of synaptic pruning<sup>8</sup>, and disruption to early brain development<sup>9,10</sup>. This complexity is underpinned by the varied nature of genetic contributions to risk of schizophrenia. Genome-wide association studies have identified over 100 independent loci defined by common (minor allele frequency [MAF] > 1%) single nucleotide variants (SNVs)<sup>11</sup>, and a recent analysis determined that more than 71% of all one-megabase regions in the genome contain at least one common risk allele<sup>12</sup>. The modest effects of these variants (median odds ratio [OR] = 1.08) combine to produce a polygenic contribution that explains only a fraction ( $h_g^2 = 0.274$ ) of the overall liability<sup>12</sup>. In addition, a number of rare variants have been identified that have far larger effects on individual risk. These are best exemplified by eleven large, rare recurrent copy number variants (CNVs) but evidence from whole-exome sequencing studies implies that many other rare coding SNVs and *de novo* mutations also confer substantial individual risk<sup>13–17</sup>. There is growing evidence that some of the same genes and pathways are affected by both common and rare variants<sup>7,18</sup>. Pathway analyses of common variants and hypothesis-driven gene set analyses of rare variants have begun to enumerate some of these specific biological processes, including histone methylation, transmission at glutamatergic synapses, calcium channel signaling, synaptic plasticity, and translational regulation by the fragile X mental retardation protein (FMRP)<sup>11,13,14,19,20</sup>.

In addition to exploring the biological mechanisms underlying schizophrenia, genetic analyses can also be used to understand its relationship to other neuropsychiatric and neurodevelopmental disorders. For instance, schizophrenia, bipolar disorder, and autism (ASD) show substantial sharing of common risk variants<sup>21,22</sup>. Sequencing studies of neurodevelopmental disorders suggest that this sharing of genetic risk may extend to rare variants of large effect. In the largest sequencing study of ASD to date, 20 of the 46 genes and all six CNVs implicated (false discovery rate [FDR] < 5%) had been previously described as dominant causes of developmental disorders<sup>23</sup>. Furthermore, an analysis of 60,706 whole exomes led by the ExAC consortium identified 3,230 genes with near-complete depletion of protein-truncating variants, and *de novo* loss-of-function (LoF) mutations identified in individuals with ASD or developmental disorders were concentrated in this set of “LoF intolerant” genes<sup>23–25</sup>. Similarly, evidence from rare variants for a broader shared genetic etiology between schizophrenia and neurodevelopmental disorders has begun to emerge. Analyses of whole-exome data provided support for an enrichment of schizophrenia rare variants in intellectual disability genes, and schizophrenia cases were also found to have a higher concentration of ultra-rare disruptive SNVs in the ExAC LoF intolerant genes compared to controls<sup>13,17,26</sup>.

However, the contribution of these rare variants to risk in the wider population of individuals diagnosed with schizophrenia, including those without intellectual disability, remains unclear. Intriguingly, the 11 rare CNVs found to be highly penetrant for schizophrenia also increased risk for intellectual disability and other congenital defects<sup>16,27</sup>, and more recently, a meta-analysis of whole-exome sequence data showed that LoF variants in *SETD1A* conferred substantial risk for both schizophrenia and neurodevelopmental disorders<sup>18</sup>. Concurrent analyses of autism whole-exome data found that *de novo* loss-of-function (LoF) mutations identified in ASD probands, particularly those that disrupt genes associated with neurodevelopmental disorders, were disproportionately found in individuals with intellectual disability<sup>23,28</sup>. These emerging results raise the possibility that rare schizophrenia risk variants may be concentrated in a subset of schizophrenia patients with co-morbid intellectual disability. Here, we present the one of the largest accumulation of schizophrenia rare variant data to date, which we jointly analyze with phenotype data on cognitive function. Using this data set, we attempt to identify groups of genes disrupted by schizophrenia rare risk variants, and determine if a subset of patients disproportionately carry these damaging alleles.

## Results

### Study design

To maximize our power to detect enrichment of damaging variants in schizophrenia cases in groups of genes, we performed a meta-analysis of three different types of rare coding variant studies: (1) high-quality SNV calls from whole-exome sequences of 4,133 schizophrenia cases and 9,274 matched controls, (2) *de novo* mutations identified in 1,077 schizophrenia parent-proband

trios (Figure 1), and (3) CNV calls from genotyping array data of 6,882 cases and 11,255 controls. The ascertainment of these samples, data production, and quality control were described previously<sup>18,29</sup>. All *de novo* mutations included in our analysis had been validated through Sanger sequencing, and stringent quality control steps were performed on the case-control data to ensure that sample ancestry and batch were closely matched between cases and controls (Online Methods).

For each data type, we used appropriate methods to test for an excess of rare variants (Figure 1, Online Methods). In analyses of case-control SNV data, we applied an extension of the variant threshold burden test that corrected for exome-wide differences between cases and controls<sup>30</sup>. We tested all allele frequency thresholds below 0.1% observed in our data, and assessed statistical significance by permutation testing. In analyses of *de novo* SNV data, we compared the observed number of *de novo* mutations to random samples from an expected distribution based on a gene-specific mutation rate model to calculate an empirical *P*-value. For both types of whole-exome sequencing data, we restricted our analyses to loss-of-function variants. Finally, in analyses of case-control CNV data, we used a logistic regression framework that compares the rate of CNVs overlapping a specific gene set while correcting for differences in CNV size and number of genes disrupted<sup>7,19,31</sup>. To ensure our model was well calibrated, we restricted our analyses to small deletions and duplications overlapping fewer than seven genes with MAF < 0.1% (Supplementary Figure 1, Online Methods).

We tested for an excess of rare damaging variants in schizophrenia patients in 1,766 gene sets (Online Methods, Supplementary Table 1, and detailed results below). Gene set *P*-values were computed using the three methods and variant definitions described above, and then meta-analyzed using Fisher's Method to provide a single *P*-value for each gene set. Because we gave each data type equal weight, gene sets achieving significance typically show at least some signal in all three types of data. We observed a marked inflation in the quantile-quantile (Q-Q) plot of gene set *P*-values (Supplementary Figure 2), so we conducted two analyses to ensure our results were robust and not biased due to methodological or technical artifacts. First, we observed no inflation of *P*-values when testing for enrichment of synonymous variants in our case-control and *de novo* analyses (Supplementary Figure 2). Second, we created random gene sets by sampling uniformly across the genome, and observed null distributions in Q-Q plots regardless of variant class and analytical method (Supplementary Figure 3). These findings suggested that our methods sufficiently corrected for known genome-wide differences in LoF and CNV burden between cases and controls, and other technical confounders like batch and ancestry.

### **Rare, damaging schizophrenia variants are concentrated in LoF intolerant genes**

We first tested whether rare schizophrenia risk variants were consistently concentrated in genes defined loss-of-function intolerant across study design and variant type. Because some of our schizophrenia exome data

was included in the ExAC database, we focused on the subset of 45,376 ExAC exomes without a known psychiatric diagnosis and that were not present in our study. From this subset, 3,488 genes were found to have near-complete depletion of such variants, which we defined as the LoF intolerant gene set. We found that rare damaging variants in schizophrenia cases were enriched in LoF intolerant genes ( $P < 3.6 \times 10^{-10}$ , Table 1, Figure 2), with support in case-control SNVs ( $P < 5 \times 10^{-7}$ ; OR 1.24, 1.16-1.31, 95% CI), case-control CNVs ( $P = 2.6 \times 10^{-4}$ ; OR 1.21, 1.15 – 1.28, 95% CI), and *de novo* mutations ( $P = 6.7 \times 10^{-3}$ ; OR 1.36, 1.1 – 1.68, 95% CI). While this result was consistent with observations in intellectual disability and ASD<sup>24,32</sup> the absolute effect size is smaller (e.g. *de novos*, Supplementary Figure 4 and 5). We observed no excess burden of rare damaging variants in the remaining 14,753 genes (Figure 2, Supplementary Figure 5). Furthermore, this signal was spread among many different LoF intolerant genes: if we rank genes by decreasing significance, the enrichment disappears in the case-control SNV analysis ( $P > 0.05$ ) only after the exclusion of the top 50 genes. This suggests that the contribution of damaging rare variants in schizophrenia is not concentrated in just a handful of genes, but instead spread across many genes.

### Schizophrenia risk genes are shared with other neurodevelopmental disorders

Given the significant enrichment of rare damaging variants in LoF intolerant genes in developmental disorders, autism and schizophrenia, we next asked whether these variants affected the same genes. We found that autism risk genes identified from exome sequencing meta-analyses<sup>23</sup> and genes in which LoF variants are known causes of severe developmental disorders as defined by the DDD study<sup>33,34</sup> were significantly enriched for rare variants in individuals with schizophrenia ( $P_{ASD} = 9.5 \times 10^{-6}$ ;  $P_{DD} = 2.3 \times 10^{-6}$ ; Table 1, Online Methods). Previous analyses have shown an enrichment of rare damaging variants in genes whose mRNA are bound by FMRP in both schizophrenia and autism<sup>35,13,32</sup>, so we sought to identify further shared biology by testing targets of neural regulatory genes previously implicated in autism<sup>32,36</sup>. We observed enrichment of both such sets: promoter targets of *CHD8* ( $P = 1.1 \times 10^{-6}$ ) and splice targets of *RBFox* ( $P = 1.3 \times 10^{-5}$ ) (Table 1). We noted that some published gene lists attributed to same biological process differed due to choices of assay, cell type, method of sample extraction, and threshold of statistical significance, leading to distinct results in our gene set analyses. For example, we observed a significant enrichment in the published FMRP binding gene set based on mouse brain data<sup>37</sup>, but with no signal in one based on a human kidney cell line<sup>38</sup>.

We also tested an additional 1,759 gene sets from databases of biological pathways with at least 100 genes, as we lacked power to detect weak enrichments in smaller sets (Online Methods). We observed enrichment of damaging rare variants in schizophrenia cases at FDR  $q < 0.05$  in 35 of these gene sets (Supplementary Table 1, 2). These included previously implicated gene sets, like the NMDA receptor and ARC complexes<sup>13,14,35,37</sup>, as well as novel gene sets, such as genes involved in cytoskeleton (GO: 0007010), chromatin modification (GO:0016568), and chromatin organization (GO: 0006325). Furthermore, the gene sets most significantly enriched (FDR  $q < 0.01$ ) for

schizophrenia rare variants (Table 1) had all been previously linked to autism, intellectual disability, and severe developmental disorders<sup>23,32,33</sup>. Our enrichment results matched some of the findings from a pathway analysis of common risk variants in psychiatric disorders, which also implicated neuronal and chromatin gene sets<sup>20</sup>. However, unlike that study, we found no enrichment of rare variants in immune-related gene sets.

We noticed that the 1,759 gene sets we tested were collectively enriched with LoF intolerant genes when compared to a random sampling of genes from the genome (Supplementary Figure 6 and 7). For some of the gene sets associated with schizophrenia, this over-representation was quite substantial: 67% of the gene targets of FMRP and 74% of the genes associated with severe neurodevelopmental disorders are LoF intolerant. To better understand the consequences of this overlap on our results, we extended the gene set enrichment methods (Online Methods) to condition on LoF intolerance and brain-expression for the 35 gene sets with FDR  $q < 0.05$  in the previous analysis (Supplementary Table 2). We first observed that 22 of the 35 gene sets remained significant even after conditioning on brain expression (Supplementary Tables 3, Online Methods), suggesting they represent more specific biological processes involved in schizophrenia. However, only known autism risk genes ( $P = 4.4 \times 10^{-4}$ ) and neurodevelopmental disorder genes ( $P = 3 \times 10^{-5}$ ) had an excess of rare coding variants above the enrichment already observed in LoF intolerant genes (Supplementary Table 3). Thus, in addition to biological pathways implicated specifically in schizophrenia, at least a portion of the schizophrenia risk conferred by rare variants of large effect is shared with childhood onset disorders of neurodevelopment.

### **Schizophrenia patients with intellectual disability have a greater burden of rare damaging variants**

In autism spectrum disorders, the observed excess of rare damaging variants has been shown to be greater in individuals with intellectual disability than those with normal levels of cognitive function<sup>28</sup>. We observed a similar phenomenon in schizophrenia cases carrying *SETD1A* LoF variants<sup>18</sup>, so next sought to explore whether this pattern is consistent in gene sets implicated in schizophrenia. We acquired relevant cognitive phenotype data for 2,971 of the 4,131 schizophrenia patients with whole-exome sequencing data (Supplementary Figure 8). Of these individuals, 279 were clinically diagnosed with intellectual disability in addition to fulfilling the full diagnostic criteria for schizophrenia (SCZ-ID subgroup, Online Methods). We also identified 1,165 individuals for whom we could rule out cognitive impairment (by excluding pre-morbid IQ  $< 85$ , fewer than 12 years of schooling or lowest decile of composite cognitive measures, depending on available data, Online Methods). Finally, we identified 1,527 individuals who were not diagnosed with intellectual disability, but in whom some cognitive impairment could not be excluded.

When stratifying into these three groups (intellectual disability, no intellectual disability but cognitive impairment not excluded, no cognitive impairment), we observed that the burden of rare damaging variants in LoF

intolerant genes was significantly greater in the SCZ-ID subgroup than in the remaining schizophrenia cases ( $P = 2.6 \times 10^{-4}$ ; OR 1.3, 1.12– 1.51, 95% CI) or controls ( $P < 5 \times 10^{-7}$ ; OR 1.61, 1.37 – 1.89, 95% CI; Figure 3). In the LoF intolerant gene set, 0.27 (0.2 – 0.35, 95% CI) extra singleton (defined as having an allele count of one in our data set) LoF variants were observed per exome in SCZ-ID cases compared to controls, while 0.10 (0.065 – 0.13, 95% CI) extra singleton LoF variants per exome were observed in the remaining schizophrenia cases compared to controls (Online Methods). Furthermore, SCZ-ID individuals had significant enrichment of rare LoF variants in developmental disorder genes compared to the other cases ( $P = 9 \times 10^{-4}$ ; OR 2.36, 1.41– 3.92, 95% CI) or to controls ( $P = 9.5 \times 10^{-6}$ ; OR 3.43, 2.01– 5.86, 95% CI; Figure 4). Compared to controls, the SCZ-ID individuals carried 0.045 (0.03 – 0.06, 95% CI) extra singleton LoF variants in developmental disorder genes per exome, suggesting that around 4% of these cases had a LoF variant that is relevant to their clinical presentation. No enrichment in neurodevelopmental disorder genes was observed in schizophrenia patients without intellectual disability, suggesting that these genes were relevant only for that subset of schizophrenia patients (Figure 4, Supplementary Table 4). Notably, even after excluding known developmental disorder genes from the set of LoF intolerant genes, we still observed an enrichment of rare variants in SCZ-ID patients compared to the remaining cases ( $P = 1 \times 10^{-3}$ ; 1.26, 1.08 – 1.47, 95% CI) or to controls ( $P < 5 \times 10^{-7}$ ; OR 1.54, 1.31– 1.81, 95% CI; Supplementary Figure 9). Rare variation in these genes contributes more to disease risk in the subset of patients with both schizophrenia and intellectual disability.

### **Rare variants confer risk for schizophrenia in individuals without intellectual disability**

While rare damaging variants in LoF intolerant genes were most enriched in the subset of schizophrenia patients with intellectual disability, we still observed a weaker but significant enrichment in individuals with schizophrenia for whom we could confirm do not have intellectual disability ( $P = 5.5 \times 10^{-4}$ ; 1.16, 1.05 – 1.27, 95% CI; Figure 3). Therefore, rare risk variants for schizophrenia follow the pattern previously described in autism: concentrated in individuals with intellectual disability, but not exclusive to that group. To produce a more accurate estimate of the effect of damaging rare variants on schizophrenia conditional on their effects on overall cognition, we recalculated the enrichment of rare variants in LoF intolerant genes in a subset of 2,161 schizophrenia cases and 2,398 controls for which data on years of education was available and for whom intellectual disability could be excluded (Supplementary Figure 8). After controlling for differences in educational attainment (Online Methods), individuals with schizophrenia have a 1.26-fold excess of rare variants in LoF intolerant genes ( $P = 2 \times 10^{-6}$ ; 1.14 – 1.38, 95% CI). This increase in our observed odds ratio is consistent with previous accounts that rare damaging variants also affect educational attainment in controls<sup>39</sup>, thus biasing our unconditional estimate.

### **Discussion**

Our integrated analysis of thousands of whole-exome sequences demonstrates that rare damaging variants increase risk of schizophrenia both with and without co-morbid intellectual disability. While the identification of individual genes remains difficult at current sample sizes, we show that the burden of damaging *de novo* mutations, rare SNVs and CNVs in schizophrenia is not scattered across the genome but is primarily concentrated in 3,488 genes intolerant of loss-of-function variants. This observation is shared with autism, intellectual disability, and severe neurodevelopmental disorders<sup>32,40</sup>. We recapitulate enrichment in previously published gene sets, including transmission at glutamatergic synapses and translational regulation by FMRP, and implicate other gene sets previously linked to autism, intellectual disability, and severe developmental disorders. However, we find that all of these gene sets share a large number of underlying genes, and are especially enriched with the 3,488 genes intolerant of LoF variants. These overlaps among gene sets originating from very different analyses, as well as the subtleties of how they are defined, suggest caution in interpreting biological explanations from observed enrichments.

We jointly analyzed the case-control SNV data with information on cognitive function for 2,971 patients, and find that LoF variants disrupting genes associated with severe developmental disorders are disproportionately found in individuals with schizophrenia with co-morbid intellectual disability, with 4% of these cases having a single LoF variant that is relevant to their clinical presentation. Even after excluding variants in known developmental disorder genes, rare variants contribute a greater degree to schizophrenia risk in the SCZ-ID subgroup of patients than the remaining schizophrenia population. These results show that some of these genetic perturbations have clear manifestations in childhood, and that rare risk variants in schizophrenia are particularly associated with co-morbid intellectual disability. Our observations are consistent with results in autism in which rare risk variants are associated with intellectual disability<sup>22,23,28</sup>. Notably, a weaker but still significant rare variant burden was observed in schizophrenia patients without cognitive impairment, and this signal persists even after controlling for educational attainment. Together, these results demonstrate that rare variants have different contributions to schizophrenia risk depending on the degree of cognitive impairment. Importantly, they do not simply confer risk for a small subset of patients but contribute to disease pathogenesis more broadly.

Our study supports the observation that genetic risk factors for psychiatric and neurodevelopmental disorders do not follow clear diagnostic boundaries. Coding variants disrupting the same genes, and quite possibly, the same biological processes, increase risk for a range of phenotypic manifestation. This clinically variable presentation is reminiscent of LoF variants in *SETD1A* and 11 large copy number variant syndromes, previously shown to confer risk for schizophrenia in addition to other prominent developmental defects<sup>16,18</sup>. It is possible that these genes contain an allelic series of variants conferring gradations of risk. A recent schizophrenia GWAS meta-analysis demonstrated that the common variant association signal was similarly enriched in LoF intolerant genes<sup>41</sup>, suggesting that schizophrenia risk genes may be perturbed by

common variants of subtle effects and disrupted by rare variants of high penetrance in the population. This possibility is also supported by the overlap in at least some of the pathways affected by both rare and common variation, such as chromatin remodeling. However, the most common deletion in the 22q11.2 locus and a recurrent two base deletion in *SETD1A* are associated with both schizophrenia and more severe neurodevelopmental disorders, suggesting the same variants can also confer risk for a range of clinical features<sup>18,42,43</sup>. Ultimately, it may prove difficult to clearly partition patients genetically into subtypes with similar clinical features, especially if genes and variants previously thought to cause well-characterized Mendelian disorders can have such varied outcomes. This pattern is consistent with the hypothesis that LoF variants in genes under genic constraint result in a spectrum of neurodevelopmental outcomes with the burden of mutations highest in intellectual disability and least in schizophrenia, corresponding to a gradient of neurodevelopmental pathology indexed by the degree of cognitive impairment, age of onset, and severity<sup>4</sup>.

Despite the complex nature of genetic contributions to risk of schizophrenia, it is notable that across study design (trio or case-control) and variant class (SNVs or CNVs), risk loci of large effect are concentrated in a small subset of genes. Previous rare variant analyses in other neurodevelopmental disorders, such as autism, have successfully integrated information across *de novo* SNVs and CNVs to identify novel risk loci<sup>23</sup>. As sample sizes increase, meta-analyses leveraging the shared genetic risk across study designs and variant types, including those we did not consider here, such as classical recessive inheritance, will be similarly well powered to identify additional risk genes in schizophrenia.

## Acknowledgements

We gratefully thank all participants in these studies. We thank Timi Touloupoulou, Marco Picchioni, Chiara Nosarti, Fiona Gaughran, and Oliver Howes for contributing clinical data used in this study. The UK10K project was funded by Wellcome Trust grant WT091310. The INTERVAL sequencing studies are funded by Wellcome Trust grant WT098051. T.S. is supported by the Williams College Dr. Herchel Smith Fellowship. A.P. is supported by Academy of Finland grants 251704 and 286500, NIMH U01MH105666 and the Sigrid Juselius Foundation. The work at Cardiff University was funded by Medical Research Council (MRC) Centre (G0801418) and Program Grants (G0800509). P.F.S. gratefully acknowledges support from the Swedish Research Council (Vetenskapsrådet, award D0886501). Creation of the Sweden schizophrenia study data was supported by NIMH R01 MH077139 and the Stanley Center of the Broad Institute. Participants in INTERVAL were recruited with the active collaboration of NHS Blood and Transplant England, which has supported fieldwork and other elements of the trial. DNA extraction and genotyping was funded by the National Institute of Health Research (NIHR), the NIHR BioResource and the NIHR Cambridge Biomedical Research Centre. The academic coordinating centre for INTERVAL was supported by core funding from: NIHR Blood and Transplant Research Unit in Donor Health and Genomics,

UK Medical Research Council (G0800270), and British Heart Foundation (SP/09/002). We would like to acknowledge the contribution of data from outside sources: (i) Genetic Architecture of Smoking and Smoking Cessation accessed through dbGAP: Study Accession: phs000404.v1.p1. Funding support for genotyping, which was performed at the Center for Inherited Disease Research (CIDR), was provided by 1 X01 HG005274-01. CIDR is fully funded through a federal contract from the National Institutes of Health to The Johns Hopkins University, contract number HHSN268200782096C. Assistance with genotype cleaning, as well as with general study coordination, was provided by the Gene Environment Association Studies (GENEVA) Coordinating Center (U01 HG004446). Funding support for collection of datasets and samples was provided by the Collaborative Genetic Study of Nicotine Dependence (COGEND; P01 CA089392) and the University of Wisconsin Transdisciplinary Tobacco Use Research Center (P50 DA019706, P50 CA084724). (ii). High-Density SNP Association Analysis of Melanoma: Case-Control and Outcomes Investigation, dbGaP Study Accession: phs000187.v1.p1. Research support to collect data and develop an application to support this project was provided by 3P50CA093459, 5P50CA097007, 5R01ES011740 and 5R01CA133996. (iii) Genetic Epidemiology of Refractive Error in the KORA Study, dbGaP Study Accession: phs000303.v1.p1. Principal investigators: Dwight Stambolian, University of Pennsylvania, Philadelphia, PA, USA; H. Erich Wichmann, Institut für Humangenetik, Helmholtz-Zentrum München, Germany, National Eye Institute, National Institutes of Health, Bethesda, MD, USA. Funded by R01 EY020483, National Institutes of Health, Bethesda, MD, USA. (iv) WTCCC2 study: Samples were downloaded from <https://www.ebi.ac.uk/ega/> and include samples from the National Blood Donors Cohort, EGAD00000000024 and samples from the 1958 British Birth Cohort, EGAD00000000022. Funding for these projects was provided by the Wellcome Trust Case Control Consortium 2 project (085475/B/08/Z and 085475/Z/08/Z), the Wellcome Trust (072894/Z/03/Z, 090532/Z/09/Z and 075491/Z/04/B) and NIMH grants (MH 41953 and MH083094).

467

## 468 **Author contributions**

469

470 T.S., J.C.B conceived and designed the experiments.

471 T.S performed the statistical analysis.

472 T.S., J.T.R.W., M.J., D.C., J.S., M.T., E.R., P.F.S analysed the data.

473 T.S., J.T.R.W., M.J., J.S., M.T., E.R., C.I., D.B., A.M.M., G.K., D.G., R.M.M., M.D.F., E.B.,

474 M.G., C.M.H., P.S., A.P., M.C.O., M.J.O., J.C.B contributed

475 reagents/materials/analysis tools.

476 T.S., D.C., M.J.O., J.C.B wrote the paper

477

## 478 **Competing financial interests statement**

479

480 We have no competing financial interests to declare.

## 481 **References**

482

483 1. van Os, J. & Kapur, S. Schizophrenia. *Lancet* **374**, 635–45 (2009).

- 484 2. American Psychiatric Association. *Diagnostic and statistical manual of*  
485 *mental disorders (DSM-5®)*. (American Psychiatric Publishing, 2013).
- 486 3. Tandon, R. *et al.* Definition and description of schizophrenia in the DSM-5.  
487 *Schizophr. Res.* **150**, 3–10 (2013).
- 488 4. Owen, M. J. New approaches to psychiatric diagnostic classification.  
489 *Neuron* **84**, 564–571 (2014).
- 490 5. Owen, M. J., Sawa, A. & Mortensen, P. B. Schizophrenia. *Lancet* **6736**, 1–12  
491 (2016).
- 492 6. Howes, O. D. & Kapur, S. The dopamine hypothesis of schizophrenia:  
493 version III--the final common pathway. *Schizophr. Bull.* **35**, 549–62 (2009).
- 494 7. Pocklington, A. J. *et al.* Novel Findings from CNVs Implicate Inhibitory and  
495 Excitatory Signaling Complexes in Schizophrenia. *Neuron* **86**, 1203–1214  
496 (2015).
- 497 8. Sekar, A. *et al.* Schizophrenia risk from complex variation of complement  
498 component 4. *Nature* **530**, 177–183 (2016).
- 499 9. Owen, M. J., O'Donovan, M. C., Thapar, A. & Craddock, N.  
500 Neurodevelopmental hypothesis of schizophrenia. *Br. J. Psychiatry* **198**,  
501 173–5 (2011).
- 502 10. Rapoport, J. L., Giedd, J. N. & Gogtay, N. Neurodevelopmental model of  
503 schizophrenia: update 2012. *Mol. Psychiatry* **17**, 1228–38 (2012).
- 504 11. Schizophrenia Working Group of the Psychiatric Genomics Consortium.  
505 Biological insights from 108 schizophrenia-associated genetic loci. *Nature*  
506 **511**, 421–7 (2014).
- 507 12. Loh, P.-R. *et al.* Contrasting genetic architectures of schizophrenia and  
508 other complex diseases using fast variance-components analysis. *Nat.*  
509 *Genet.* **47**, 1385–1392 (2015).
- 510 13. Fromer, M. *et al.* De novo mutations in schizophrenia implicate synaptic  
511 networks. *Nature* **506**, 179–184 (2014).
- 512 14. Kirov, G. *et al.* De novo CNV analysis implicates specific abnormalities of  
513 postsynaptic signalling complexes in the pathogenesis of schizophrenia.  
514 *Mol. Psychiatry* **17**, 142–53 (2012).
- 515 15. The International Schizophrenia Consortium. Rare chromosomal deletions  
516 and duplications increase risk of schizophrenia. *Nature* **455**, 237–41  
517 (2008).
- 518 16. Rees, E. *et al.* Analysis of copy number variations at 15 schizophrenia-  
519 associated loci. *Br. J. Psychiatry* **204**, 108–14 (2014).
- 520 17. Zhu, X., Need, A. C., Petrovski, S. & Goldstein, D. B. One gene, many  
521 neuropsychiatric disorders: lessons from Mendelian diseases. *Nat.*  
522 *Neurosci.* **17**, 773–781 (2014).
- 523 18. Singh, T. *et al.* Rare loss-of-function variants in SETD1A are associated  
524 with schizophrenia and developmental disorders. *Nat. Neurosci.* **19**, 571–  
525 577 (2016).
- 526 19. Szatkiewicz, J. P. *et al.* Copy number variation in schizophrenia in Sweden.  
527 *Mol. Psychiatry* **19**, 762–773 (2014).
- 528 20. Psychiatric Genetics Consortium. Psychiatric genome-wide association  
529 study analyses implicate neuronal, immune and histone pathways. *Nat.*  
530 *Neurosci.* (2015). doi:10.1038/nn.3922
- 531 21. Lee, S. H. *et al.* Genetic relationship between five psychiatric disorders  
532 estimated from genome-wide SNPs. *Nat. Genet.* **45**, 984–94 (2013).

- 533 22. Robinson, E. B. *et al.* Genetic risk for autism spectrum disorders and  
534 neuropsychiatric variation in the general population. *Nat. Genet.* **48**, 552–  
535 555 (2016).
- 536 23. Sanders, S. J. *et al.* Insights into Autism Spectrum Disorder Genomic  
537 Architecture and Biology from 71 Risk Loci. *Neuron* **87**, 1215–1233  
538 (2015).
- 539 24. Samocha, K. E. *et al.* A framework for the interpretation of de novo  
540 mutation in human disease. *Nat. Genet.* **46**, 944–950 (2014).
- 541 25. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706  
542 humans. *Nature* **536**, 285–291 (2016).
- 543 26. Genovese, G. *et al.* Increased burden of ultra-rare protein-altering variants  
544 among 4,877 individuals with schizophrenia. *Nat. Neurosci.* (2016).  
545 doi:10.1038/nn.4402
- 546 27. Kirov, G. *et al.* The penetrance of copy number variations for schizophrenia  
547 and developmental delay. *Biol. Psychiatry* **75**, 378–85 (2014).
- 548 28. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism  
549 spectrum disorder. *Nature* **515**, 216–21 (2014).
- 550 29. Rees, E. *et al.* CNV analysis in a large schizophrenia sample implicates  
551 deletions at 16p12.1 and SLC1A1 and duplications at 1p36.33 and CGNL1.  
552 *Hum. Mol. Genet.* **23**, 1669–76 (2014).
- 553 30. Price, A. L. *et al.* Pooled Association Tests for Rare Variants in Exon-  
554 Resequencing Studies. *Am. J. Hum. Genet.* **86**, 832–838 (2010).
- 555 31. Raychaudhuri, S. *et al.* Accurately assessing the risk of schizophrenia  
556 conferred by rare copy-number variation affecting genes with brain  
557 function. *PLoS Genet.* **6**, (2010).
- 558 32. De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted  
559 in autism. *Nature* **515**, 209–15 (2014).
- 560 33. Firth, H. V *et al.* DECIPHER: Database of Chromosomal Imbalance and  
561 Phenotype in Humans Using Ensembl Resources. *Am. J. Hum. Genet.* **84**,  
562 524–33 (2009).
- 563 34. Deciphering Developmental Disorders Study. Prevalence and architecture  
564 of de novo mutations in developmental disorders. *Nature* **542**, 433–438  
565 (2017).
- 566 35. Purcell, S. M. *et al.* A polygenic burden of rare disruptive mutations in  
567 schizophrenia. *Nature* **506**, 185–90 (2014).
- 568 36. Cotney, J. *et al.* The autism-associated chromatin modifier CHD8 regulates  
569 other autism risk genes during human neurodevelopment. *Nat. Commun.*  
570 **6**, 6404 (2015).
- 571 37. Darnell, J. C. *et al.* FMRP stalls ribosomal translocation on mRNAs linked to  
572 synaptic function and autism. *Cell* **146**, 247–61 (2011).
- 573 38. Ascano, M. *et al.* FMRP targets distinct mRNA sequence elements to  
574 regulate protein expression. *Nature* **492**, 382–386 (2012).
- 575 39. Ganna, A. *et al.* Ultra-rare disruptive and damaging mutations influence  
576 educational attainment in the general population. *Nat. Neurosci.* **19**, 1563–  
577 1565 (2016).
- 578 40. The Deciphering Developmental Disorders Study. Large-scale discovery of  
579 novel genetic causes of developmental disorders. *Nature* **519**, 223–8  
580 (2015).
- 581 41. Pardiñas, A. F. *et al.* Common schizophrenia alleles are enriched in

- 582 mutation-intolerant genes and maintained by background selection.  
583 *bioRxiv* 68593 (2016). doi:10.1101/068593  
584 42. Ben-Shachar, S. *et al.* 22q11.2 Distal Deletion: A Recurrent Genomic  
585 Disorder Distinct from DiGeorge Syndrome and Velocardiofacial  
586 Syndrome. *Am. J. Hum. Genet.* **82**, 214–221 (2008).  
587 43. Michaelovsky, E. *et al.* Genotype-phenotype correlation in 22q11.2  
588 deletion syndrome. *BMC Med. Genet.* **13**, 122 (2012).

## 589 Figure captions

590  
591 **Figure 1:** Analysis workflow. Data sets are shown in blue, statistical methods  
592 and analysis steps are shown in green, and results (figures and tables) from the  
593 analysis are shown in orange. **A:** Enrichment analyses in 1,766 gene sets using  
594 the entire rare variant data set. **B:** Enrichment analyses in LoF intolerant and  
595 developmental disorder genes in the subset of cases with information on  
596 cognitive function. ID: intellectual disability; SCZ: schizophrenia; SCZ-ID:  
597 schizophrenia patients with intellectual disability.

598 **Figure 2:** Enrichment of schizophrenia rare variants in genes intolerant of loss-  
599 of-function variants. **A:** Schizophrenia cases compared to controls for rare SNVs  
600 and indels; **B:** Rates of *de novo* mutations in schizophrenia probands compared  
601 to control probands; **C:** Case-control CNVs. *P*-values shown were from the test of  
602 LoF enrichment in **A**, LoF enrichment in **B**, and all CNVs enrichment in **C**. Error  
603 bars represent the 95% CI of the point estimate. LoF intolerant: 3,448 genes with  
604 near-complete depletion of truncating variants in the ExAC database; Rest: the  
605 remaining genes in the genome with pLI < 0.9; Damaging missense: missense  
606 variants with CADD phred > 15. Asterisk:  $P < 1 \times 10^{-3}$ .

607  
608 **Figure 3:** Enrichment of rare loss-of-function variants in LoF intolerant genes in  
609 schizophrenia cases stratified by information on cognitive function compared to  
610 controls. The *P*-values shown were calculated using the variant threshold  
611 method comparing LoF burden between the corresponding cases and controls.  
612 Error bars represent the 95% CI of the point estimate. Damaging missense:  
613 missense variants with CADD phred > 15.

614  
615 **Figure 4:** Enrichment of rare loss-of-function variants in known severe  
616 developmental disorder genes in schizophrenia cases stratified by information  
617 on cognitive function compared to controls. The *P*-values shown were calculated  
618 using the variant threshold method comparing LoF burden between the  
619 corresponding cases and controls. Error bars represent the 95% CI of the point  
620 estimate. Damaging missense: missense variants with CADD phred > 15.  
621

Name	N <sub>genes</sub>	Est <sub>SNV</sub>	95% CI of Est <sub>SNV</sub>	P <sub>SNV</sub>	Est <sub>DNM</sub>	95% CI of Est <sub>DNM</sub>	P <sub>DNM</sub>	Est <sub>CNV</sub>	95% CI of Est <sub>CNV</sub>	P <sub>CNV</sub>	P <sub>meta</sub>	Q <sub>meta</sub>
ExAC LoF intolerant genes (pLI > 0.9)	3488	1.24	1.16-1.31	< 5.0 x 10 <sup>-7</sup>	1.36	1.1-1.68	0.0067	1.21	1.15-1.28	0.00026	< 3.60 x 10 <sup>-10</sup>	4.30 x 10 <sup>-7</sup>
Dominant, diagnostic DDG2P genes, in which LoF variants result in developmental disorders with brain abnormalities	156	1.42	1.07-1.88	0.011	4.18	2.21-8.03	0.00073	1.92	1.54-2.39	0.0016	2.30 x 10 <sup>-6</sup>	0.00067
Sanders <i>et al.</i> autism risk genes (FDR < 10%)	66	1.28	0.97-1.69	0.0095	3.96	1.65-9.94	0.019	2.21	1.75-2.79	0.00033	9.50 x 10 <sup>-6</sup>	0.0017
Darnell <i>et al.</i> targets of FMRP	790	1.24	1.13-1.36	8.5 x 10 <sup>-6</sup>	1.31	0.83-2.09	0.17	1.32	1.2-1.47	0.0032	9.30 x 10 <sup>-7</sup>	0.00038
Cotney <i>et al.</i> CHD8-targeted promoters (hNSC and human brain tissue)	2920	1.09	1.02-1.16	0.0008	1.77	1.36-2.31	0.00025	1.11	1.05-1.18	0.027	1.10 x 10 <sup>-6</sup>	0.00038
G2CDB: mouse cortex post-synaptic density consensus	1527	1.20	1.11-1.3	2.5 x 10 <sup>-6</sup>	1.57	1.06-2.33	0.028	1.04	0.96-1.11	0.32	3.90 x 10 <sup>-6</sup>	0.00097
Weynvanhentenryck <i>et al.</i> CLIP targets of RBFOX	967	1.21	1.11-1.33	4.8 x 10 <sup>-5</sup>	1.84	1.21-2.8	0.0085	1.07	0.98-1.17	0.2	1.30 x 10 <sup>-5</sup>	0.002
NMDAR network (defined in Purcell <i>et al.</i> )	61	1.66	1.09-2.54	0.0061	5.60	2.06-16.09	0.017	2.46	1.78-3.4	0.0028	3.70 x 10 <sup>-5</sup>	0.0044
GOBP: chromatin modification (GO:0016568)	519	1.29	1.13-1.49	0.00018	2.26	1.32-3.94	0.0099	1.12	0.99-1.28	0.18	4.20 x 10 <sup>-5</sup>	0.0046

**Table 1:** Gene sets enriched for rare coding variants conferring risk for schizophrenia at FDR < 1%. The effect sizes and corresponding P-values from enrichment tests of each variant type (case-control SNVs, DNM, and case-control CNVs) are shown for each gene set, along with the Fisher's combined P-value (P<sub>meta</sub>) and the FDR-corrected Q-value (Q<sub>meta</sub>). We only show the most significant gene set if there are multiple ones from the same data set or biological process (see Supplementary Table 1 for all 1,766 gene sets). N<sub>genes</sub>: number of genes in the gene set; Est: effect size estimate and its lower and upper bound assuming a 95% CI; DNM: *de novo* mutation.

## 627 **Supplementary Table captions**

628

629 **Supplementary Table 1:** Full results from enrichment analyses of 1,766 gene  
630 sets. The  $P$ -values from enrichment tests of each variant type (case-control SNVs,  
631 DNM, and case-control CNVs) are shown for each gene set, along with the  
632 Fisher's combined  $P$ -value ( $P_{\text{meta}}$ ) and the FDR-corrected  $Q$ -value ( $Q_{\text{meta}}$ ).  $N_{\text{genes}}$ :  
633 number of genes in the gene set; SNV: single nucleotide variants from whole-  
634 exome data; DNM: *de novo* mutations.

635

636 **Supplementary Table 2:** Gene sets enriched for rare coding variants conferring  
637 risk for schizophrenia at  $\text{FDR} < 5\%$ . The effect sizes and corresponding  $P$ -values  
638 from enrichment tests of each variant type (case-control SNVs, DNM, and case-  
639 control CNVs) are shown for each gene set, along with the Fisher's combined  $P$ -  
640 value ( $P_{\text{meta}}$ ) and the FDR-corrected  $Q$ -value ( $Q_{\text{meta}}$ ).  $N_{\text{genes}}$ : number of genes in  
641 the gene set; Est: effect size estimate and its lower and upper bound assuming a  
642 95% CI; SNV: single nucleotide variants from whole-exome data; DNM: *de novo*  
643 mutations.

644

645 **Supplementary Table 3:** Results from enrichment analyses of  $\text{FDR} < 5\%$  gene  
646 sets, conditional on brain-expressed and ExAC LoF intolerant genes. We restrict  
647 enrichment analyses to genes that reside in two different background gene sets,  
648 one defined on brain-enriched expression in GTEx, and the second on genic  
649 constraint (ExAC LoF intolerant genes), and determined if gene sets with  $\text{FDR} <$   
650  $5\%$  in the meta-analysis still had significance above the specific background. The  
651  $P$ -values from enrichment tests of each variant type (case-control SNVs, DNM,  
652 and case-control CNVs) are shown for each gene set, along with the Fisher's  
653 combined  $P$ -value ( $P_{\text{meta}}$ ). SNV: single nucleotide variants from whole-exome  
654 data; DNM: *de novo* mutations

655

656 **Supplementary Table 4:** Results from enrichment analyses of rare loss-of-  
657 function variants in LoF intolerant genes and developmental disorder genes  
658 comparing schizophrenia cases stratified by information on cognitive function  
659 and matched controls. Each comparison is defined in the Table, and the  $P$ -values  
660 shown were calculated using the variant threshold method comparing LoF  
661 burden between the corresponding case and baseline samples.  $N_{\text{case}}$ : number of  
662 case samples;  $N_{\text{comparison}}$ : number of comparison samples; Estimates: effect size  
663 estimate and its lower and upper bound assuming a 95% CI.

## 664 **Online Methods**

### 665 **Sample collections**

666

667 The ascertainment, data production, and quality control of the  
668 schizophrenia case-control whole-exome sequencing data set had been  
669 described in detail in an earlier publication<sup>18</sup>. Briefly, the data set was composed  
670 of schizophrenia cases recruited as part of eight collections in the UK10K  
671 sequencing project, and matched population controls from non-psychiatric arms  
672 of the UK10K project, healthy blood donors from the INTERVAL project, and five

673 Finnish population studies. The UK10K data set was combined and analyzed  
674 with published data from a Swedish schizophrenia case-control study<sup>35</sup>. The data  
675 production, quality control, and analysis of the case-control CNV data set was  
676 described in an earlier publication<sup>29</sup>. The schizophrenia cases were recruited as  
677 part of the CLOZUK and CardiffCOGS studies, which consisted of both  
678 schizophrenia individuals taking the antipsychotic clozapine and a general  
679 sample of cases from the UK. Matched controls were selected from four publicly  
680 available non-psychiatric data sets. All samples were genotyped using Illumina  
681 arrays, and processed and called under the same protocol. Sanger-validated *de*  
682 *nov* mutations identified through whole exome-sequencing in seven published  
683 studies of schizophrenia parent-proband trios were aggregated and re-annotated  
684 for enrichment analyses<sup>13,44–49</sup>. A full description of each trio study, including  
685 sequencing and capture technology and sample recruitment was previously  
686 described<sup>18</sup>.

## 687 **Sample and variant quality control**

688  
689 We jointly called each case data set with its nationality-matched controls,  
690 and excluded samples based on contamination, coverage, non-European  
691 ancestry, and excess relatedness<sup>18</sup>. A number of empirically derived filters were  
692 applied at the variant and genotype level, including filters on GATK VQSR,  
693 genotype quality, read depth, allele balance, missingness, and Hardy-Weinberg  
694 disequilibrium<sup>18</sup>. After variant filtering, the per-sample transition-to-  
695 transversion ratio was ~3.2 across the entire data set, as expected for  
696 populations of European ancestry<sup>50</sup>. For the case-control CNV analysis, we  
697 similarly excluded samples based on excess relatedness, and only CNVs  
698 supported by more than 10 probes and greater than 10 kilobases in size were  
699 retained to ensure high quality calls. All *de novo* mutations in our study had been  
700 validated using Sanger sequencing.

701  
702 We used the Ensembl Variant Effect Predictor (VEP) version 75 to  
703 annotate all variants (SNVs and CNVs) according to Gencode v.19 coding  
704 transcripts. We defined frameshift, stop gained, splice acceptor, and donor  
705 variants as loss-of-function (LoF), and missense or initiator codon variants with  
706 the recommended CADD Phred score cut-off of greater than 15 as damaging  
707 missense<sup>51</sup>. A gene was annotated as disrupted by a deletion if part of its coding  
708 sequence overlapped the copy number event. We more conservatively defined  
709 genes as duplicated only if the entire canonical transcript of the gene overlapped  
710 with the duplication event.

711  
712 Statistical tests of the case-control exome data used case-control  
713 permutations within each population (UK, Finnish, Swedish) to generate  
714 empirical *P*-values to test hypotheses. No genome-wide inflation was observed in  
715 burden tests of individual genes<sup>18</sup>. In the curated set of *de novo* mutations, we  
716 observed the expected exome-wide number of synonymous mutations given  
717 gene mutation rates from previously validated models<sup>24</sup>, suggesting variant  
718 calling was generally unbiased across Gencode v.19 coding genes. Lastly, the  
719 case-control CNV data set had been previously analyzed for burden of CNVs  
720 affecting individual genes, and enrichment analyses in targeted gene sets<sup>7,29</sup>.

## 721 Rare variant gene set enrichment analyses

722 **Case-control enrichment burden tests** For the case-control SNV data set, we  
 723 performed permutation-based gene set enrichment tests using an extension of  
 724 the variant threshold method<sup>30</sup>. This method assumed that variants with a MAF  
 725 below an unknown threshold  $T$  were more likely to be damaging than variants  
 726 with a MAF above  $T$ , and this threshold was allowed to differ for every gene or  
 727 pathway tested. To consider different possible values for threshold  $T$ , a gene or  
 728 gene set test statistic  $t(T)$  was calculated for every allowable  $T$ , and the  
 729 maximum test-statistic, or  $t_{\max}$ , was selected. The statistical significance of  $t_{\max}$   
 730 was evaluated by permuting phenotypic labels, and calculating  $t_{\max}$  from the  
 731 permuted data such that different values of  $T$  could be selected following each  
 732 permutation. In Price *et al.*,  $t(T)$  was defined as the z-score calculated from  
 733 regressing the phenotype on the sum of the allele counts of variants in a gene  
 734 with  $\text{MAF} < T$ . We extended this method to test for enrichment in gene sets by  
 735 regressing schizophrenia status on the total number of damaging alleles in the  
 736 gene set of interest with  $\text{MAF} < T$  ( $X_{in,T}$ ) while correcting for the total number of  
 737 damaging alleles genome-wide with  $\text{MAF} < T$  ( $X_{all,T}$ ).  $X_{all,T}$  controlled for  
 738 exome-wide differences between schizophrenia cases and controls, ensuring any  
 739 significant gene set result was significant beyond baseline differences.  $t(T)$  was  
 740 defined as the  $t$ -statistic testing if the regression coefficient of  $X_{in,T}$  deviated  
 741 from 0. We then calculated  $t(T)$  for all observed thresholds below a minor allele  
 742 frequency of 0.1%, and selected the maximum value for the  $t_{\max}$  based on the  
 743 observed data. To calculate a null distribution for  $t_{\max}$ , we performed two  
 744 million case-control permutations within each population (UK, Finnish, and  
 745 Swedish) to control for batch and ancestry, and calculated  $t_{\max}$  for each  
 746 permuted sample while allowing  $T$  to vary. The  $P$ -value for each gene set was  
 747 calculated as the fraction of the two million permuted samples that had a greater  
 748  $t_{\max}$  than what was observed in the unpermuted data. The odds ratio and 95%  
 749 confidence interval of each gene set was calculated using a logistic regression  
 750 model, regressing schizophrenia status on  $X_{in}$  while controlling for total number  
 751 of variants genome-wide ( $X_{all}$ ) and population (UK, Finnish, and Swedish).  
 752 Unlike gene set  $P$ -values which were calculated using permutation across  
 753 multiple frequency thresholds, the odds ratios and 95% CI were calculated using  
 754 only variants observed once in our data set (allele count of 1) to ensure they  
 755 were comparable between tested gene sets.

756 **CNV logistic regression** We adapted a logistic regression framework described in  
 757 Raychaudhuri *et al.* and implemented in PLINK to compare the case-control  
 758 differences in the rate of CNVs overlapping a specific gene set while correcting  
 759 for differences in CNV size and total genes disrupted<sup>7,19,31</sup>. We first restricted our  
 760 analyses to coding deletions and duplications, and tested for enrichment using  
 761 the following model:

$$762 \quad \log\left(\frac{p_{i,\text{case}}}{1-p_{i,\text{case}}}\right) = \beta_0 + \beta_1 s_i + \beta_2 g_{\text{all}} + \beta_3 g_{\text{in}} + \epsilon,$$

763 where for individual  $i$ ,  $p_i$  is the probability they have schizophrenia,  $s_i$  is the  
 764 total length of CNVs,  $g_{\text{all}}$  is the total number of genes overlapping CNVs, and  $g_{\text{in}}$   
 765 is the number of genes within the gene set of interest overlapping CNVs. It has been  
 766 shown that  $\beta_1$  and  $\beta_2$  sufficiently controlled for the genome-wide differences in

767 the rate and size of CNVs between cases and control, while  $\beta_3$  captured the true  
768 gene set enrichment above this background rate<sup>7,19,31</sup>. For each gene set, we  
769 reported the one-sided *P*-value, odds ratio, and 95% confidence interval of  $\beta_3$ .

770 **Weighted permutation-based sampling of *de novo* mutations** For each variant  
771 class of interest, we first determined the total number of *de novo* mutations  
772 observed in the 1,077 schizophrenia trios. We then generated 2 million random  
773 samples with the same number of *de novo* mutations, weighting the probability  
774 of observing a mutation in a gene by its estimated mutation rate. The baseline  
775 gene-specific mutation rates were obtained using the method described in  
776 Samocha *et al.* and adapted to produce LoF and damaging missense rates for  
777 each Gencode v.19 gene. These mutation rates adjusted for both sequence  
778 context and gene length, and were successfully applied in the primary analyses  
779 of large-scale exome sequencing of autism and severe developmental disorders  
780 with replicable results<sup>23,32,40</sup>. For each gene set, one-sided enrichment *P*-values  
781 were calculated as the fraction of two million random samples that had a greater  
782 or equal number of *de novo* mutations in the gene set of interest than what is  
783 observed in the 1,077 trios. The effect size of the enrichment was calculated as  
784 the ratio between the number of observed mutations in the gene set of interest  
785 and the average number of mutations in the gene set across the two million  
786 random samples. We adapted a method in Fromer *et al.* to calculate 95% credible  
787 intervals for the enrichment statistic<sup>13</sup>. We first generated a list of one thousand  
788 evenly spaced values between 0 and ten times the point estimate of the  
789 enrichment. For each value, the mutation rates of genes in the gene set of  
790 interest were multiplied by that amount, and 50,000 random samples of *de novo*  
791 mutations were generated using these weighted rates. The probability of  
792 observing the number of mutations in the gene set of interest given each effect  
793 size multiplier was calculated as the fraction of samples in which the number of  
794 mutations in the gene set is the same as the observed number in the 1,077 trios.  
795 We normalized the probabilities across the 1,000 values to generate a posterior  
796 distribution of the effect size, and calculated the 95% credible interval using this  
797 empirical distribution.

798  
799 **Combined joint analysis** Gene set *P*-values calculated using the case-control SNV,  
800 case-control CNV, and *de novo* data were meta-analyzed using Fisher's combined  
801 probability method with *df* = 6 to provide a single test statistic for each gene set.  
802 We corrected for the number of gene sets tested in the discovery analysis (*n* =  
803 1,776) by controlling the false discovery rate (FDR) using the Benjamini-  
804 Hochberg approach, and reported only results with a *q*-value of less than 5%.

## 805 806 **Description of gene sets**

807  
808 The full list of tested gene sets is found in Supplementary Table 1, and a  
809 detailed description is provided in the Supplementary Note. Briefly, we tested all  
810 gene sets with more than 100 genes from five public pathway databases. We  
811 additionally tested additional gene sets selected based on biological hypotheses  
812 about schizophrenia risk, and genome-wide screens investigating rare variants  
813 in intellectual disability, autism spectrum disorders, and other  
814 neurodevelopmental disorders. All gene identifiers were mapped to the

815 GENCODE v.19 release, and all non-coding genes were excluded. A total of 1,766  
816 gene sets were included in our analysis.

### 817 **Selection of allele frequency thresholds and consequence severity**

818  
819 For the case-control whole-exome data, we applied an extension of the  
820 variant threshold model (described above). With this method, we tested  
821 damaging variants at a number of frequency thresholds without specifying an *a*  
822 *priori* MAF cut-off. All thresholds below a MAF of 0.1% observed in our data  
823 were tested, and we assessed statistical significance by permutation testing. For  
824 all the whole-exome data (case-control and trio data), we restricted our analyses  
825 to loss-of-function variants. These variants have a clear and severe predicted  
826 functional consequence in that they putatively cause a single-copy loss of a gene.  
827 Furthermore, this class of variants had been demonstrated to have the strongest  
828 genome-wide enrichment between cases and controls across  
829 neurodevelopmental and psychiatric disorders<sup>18,32,40</sup>. When selecting MAF cut-  
830 offs for case-control CNVs, we found that while the bulk of the test statistics were  
831 not inflated, the tail of gene set *P*-values were dramatically inflated even when  
832 testing for enrichment in the random gene sets (Supplementary Figure 1). This  
833 inflation in the tail of the Q-Q plot was driven in part by very large (overlapping  
834 more than 10 genes), more common (MAF between 0.1% and 1%) CNVs  
835 observed mainly in cases or controls. Some of these, such as the known  
836 syndromic CNVs, likely harbored true risk genes. However, because these CNVs  
837 were highly recurrent in cases and depleted in controls, and disrupted a large  
838 number of genes, any gene set that included even a single gene within these  
839 CNVs would appear to be significant, even after controlling for total CNV length  
840 and genes overlapped. To ensure our model was well calibrated and its *P*-values  
841 followed a null distribution for random gene sets, we explored different  
842 frequency and size thresholds, and conservatively restricted our analysis to copy  
843 number events overlapping less than seven genes (excluding the largest 10% of  
844 CNVs) with MAF < 0.1% (Supplementary Figure 1). Our main conclusions  
845 remained unchanged even if we selected a more stringent (excluding the largest  
846 15% of CNVs) or less stringent (excluding the largest 5% of CNVs) size threshold.

847

### 848 **Robustness of enrichment analyses**

849

850 We uniformly sampled genes from the genome (as defined by Gencode  
851 v.19) to generate random gene sets with the same size distribution as the 1,776  
852 gene sets in our discovery analysis. For each random set, we calculated gene set  
853 *P*-values for the case-control SNV data, case-control CNV data, and *de novo* data  
854 using the appropriate method and frequency cut-offs across all variant classes. A  
855 Q-Q plot was generated using *P*-values from enrichment tests of each data set  
856 and variant type. Reassuringly, we observed null distributions in all such Q-Q  
857 plots (Supplementary Figure 3).

858

### 859 **Comparison of *de novo* enrichment with broader neurodevelopmental** 860 **disorders**

861

We aggregated and re-annotated *de novo* mutations from four studies: 1,113 severe DD probands<sup>40</sup>, 4,038 ASD probands<sup>23,32</sup>, and 2,134 control probands<sup>28,32</sup>. We used the Poisson exact test to calculate differences in *de novo* rates in constrained genes between schizophrenia, ASD, and DD and controls. Counts in each functional class (synonymous, missense, damaging missense, and LoF) were tested separately, and the one-sided *P*-value, rate ratio, and 95% CI of each comparison were reported and plotted in Figure 2, Supplementary Figure 4 and 5.

## Conditional analyses

In each of the three methods we used for gene set enrichment, we restricted all variants analyzed to those that reside in the background gene list, and tested for an excess of rare variants in genes shared between the gene set of interest (*K*) and the background list (*B*). Brain-enriched genes from GTEx, and the ExAC LoF intolerant genes (pLI > 0.9) were used as backgrounds (see above). For the case-control SNV data, we modified the variant threshold method to regress schizophrenia status on the total number of damaging alleles in genes present in both the gene set of interest and the background gene set ( $K \cap B$ ), while correcting for the total number of damaging alleles in the set of all background genes (*B*). The logistic regression model for the case-control CNV data was modified to:

$$\log\left(\frac{P_{i,\text{case}}}{1-P_{i,\text{case}}}\right) = \beta_0 + \beta_1 s_i + \beta_2 g_B + \beta_3 g_{K \cap B} + \epsilon,$$

where  $g_B$  is the total number of background genes overlapping a CNV, and  $g_{K \cap B}$  is the number of genes in the intersection of the gene set of interest and the background list overlapping a CNV. Finally, we determined the total number of *de novo* mutations within the background gene list observed in the 1,077 schizophrenia trios, and generated 2 million random samples with the same number of *de novo* mutations. For each gene set, one-sided enrichment *P*-values were calculated as the fraction of two million random samples that had a greater or equal number of *de novo* mutations in genes in  $K \cap B$  than what is observed in the 1,077 trios. Gene set *P*-values were combined using Fisher's method. We restricted our conditional enrichment analysis to gene sets with *q*-value < 5% in the discovery analysis, and adjusted for multiple testing using Bonferroni correction ( $P = 0.00071$ , or  $0.05/67$  tests; see Supplementary Table 3).

## Rare variants and cognition in schizophrenia

Within the UK10K study, 97 individuals from the MUIR collection were given discharge diagnoses of mild learning disability and schizophrenia (ICD-8 and -9). The recruitment guidelines of the MUIR collection were described in detail in a previous publication<sup>52</sup>. In brief, evidence of remedial education was a prerequisite to inclusion, and individuals with pre-morbid IQs below 50 or above 70, severe learning disabilities, or were unable to give consent were excluded. The Schizophrenia and Affective Disorders Schedule-Lifetime version (SADS-L) in people with mild learning disability, PANSS, RDC, and DSM-III-R, and St. Louis Criterion were applied to all individuals to ensure that any diagnosis of

908 schizophrenia was robust. Using the clinical information provided alongside the  
909 Swedish and Finnish case-control data sets, we identified additional 182  
910 schizophrenia individuals who were similarly diagnosed with intellectual  
911 disability, for a total of 279 individuals.

912 Cognitive testing and educational attainment data available for a subset of  
913 samples were used identify schizophrenia individuals without cognitive  
914 impairment. For 502 individuals from the Cardiff collection in the UK10K study,  
915 we acquired their pre-morbid IQ as extrapolated from National Adult Reading  
916 Test (NART), and identified 412 individuals for analysis after excluding all  
917 individuals with predicted pre-morbid IQ of less than 85 (or below one standard  
918 deviation of the population distribution for IQ). We additionally acquired  
919 information on educational attainment in 54 schizophrenia individuals in the  
920 UK10K London collection, and retained 27 individuals without intellectual  
921 disability and who completed at least 12 years of schooling. Lastly, the California  
922 Verbal Learning Test was conducted on 124 Finnish schizophrenia individuals  
923 sequenced as part of UK10K, and a composite score was generated from  
924 measures of verbal and visual working memory, verbal abilities,  
925 visuoconstructive abilities, and processing speed. All individuals with intellectual  
926 disability had been excluded from cognitive testing. Within this set of samples,  
927 we additionally excluded any individuals who ranked in the lowest decile in  
928 CVLT composite score, and retained 92 individuals for analysis. According to  
929 these criteria, we identified 531 of 697 schizophrenia individuals from the UK  
930 and Finnish data sets with cognitive data as not having intellectual disability. We  
931 additionally acquired data on educational attainment for the Swedish  
932 schizophrenia cases and controls from the Swedish National Registry. After  
933 excluding individuals with intellectual disability, we identified 1,527  
934 schizophrenia individuals who did not complete secondary school (less than 12  
935 years of schooling), and 634 schizophrenia individuals who completed at least  
936 compulsory and upper secondary schooling (at least 12 years of schooling). The  
937 last group with the greatest educational attainment and without intellectual  
938 disability was defined as cases without cognitive impairment. In the Swedish  
939 sample, 49.4% of control samples had lower educational attainment than the  
940 634 individuals with schizophrenia defined as having no cognitive impairment,  
941 suggesting that our definition was sufficiently strict. In total, combining the UK,  
942 Finnish, and Swedish data, we identified 1,165 schizophrenia individuals without  
943 cognitive impairment.

944 Using the variant threshold method, we tested for differences in rare LoF  
945 burden between the three case groups (intellectual disability, did not complete  
946 secondary school, no cognitive impairment) against controls. We restricted these  
947 analyses to three gene sets (LoF intolerant genes, genes in which LoF variants  
948 are diagnostic for severe developmental disorders, and LoF intolerant genes  
949 after excluding severe developmental disorders genes), and adjusted for multiple  
950 testing using Bonferroni correction ( $P = 0.0038$ , or  $0.05/13$  tests).  
951 Supplementary Table 4 enumerated all the statistical tests performed. To  
952 estimate the per-exome excess of rare singleton (defined as having an allele  
953 count of one in our data set) LoF variants in cases compared to controls, we  
954 regressed  $X_{in}$  (the number of LoF variants in the gene set of interest) on case  
955 status (0 or 1) while controlling for  $X_{all}$  (the total number of LoF variants

956 genome-wide) and population (UK, Finnish, and Swedish). The effect size and  
957 95% CI of the regression coefficient of case status predictor were reported.

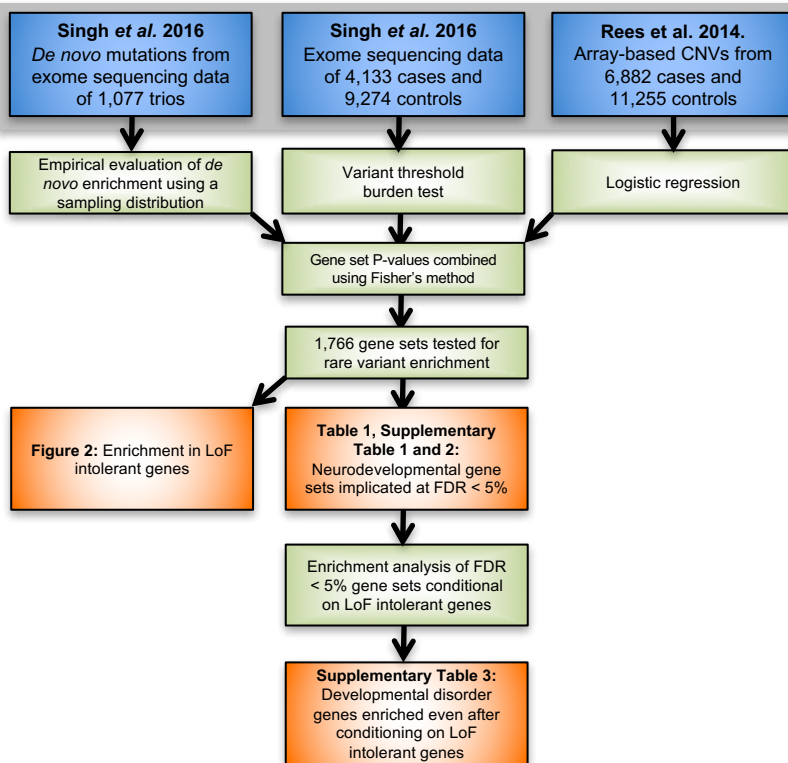
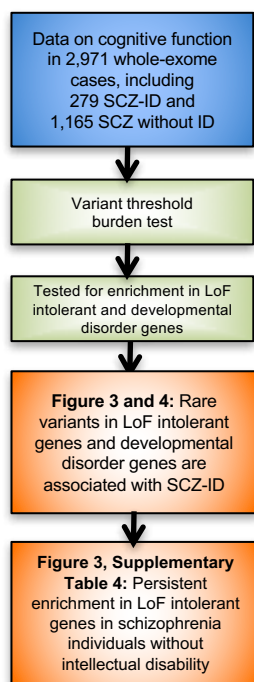
## 958 **Data Availability**

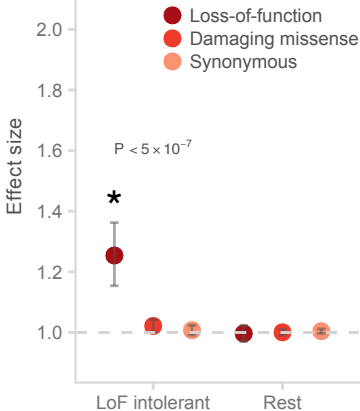
959  
960 Sequence data and processed VCFs for the UK10K project were deposited into  
961 the European Genome-phenome Archive (EGA) under study accession code  
962 EGAO00000000079. The processed VCFs from the Swedish case-control study  
963 were deposited in dbGAP under accession code (phs000473.v1.p1). Rare variant  
964 counts, and gene-level association results from combining the whole-exome  
965 sequencing data sets were described in a previous publication<sup>18</sup> and were made  
966 available on the PGC results and download page  
967 (<https://www.med.unc.edu/pgc/results-and-downloads>).  
968

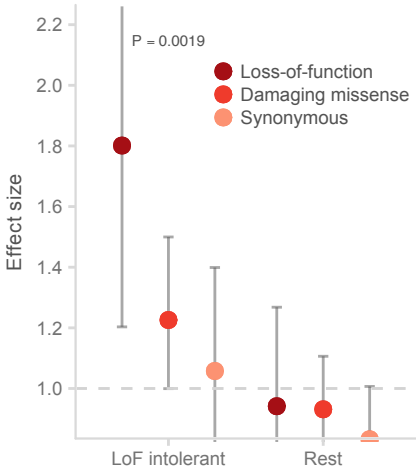
## 969 **References for Online Methods**

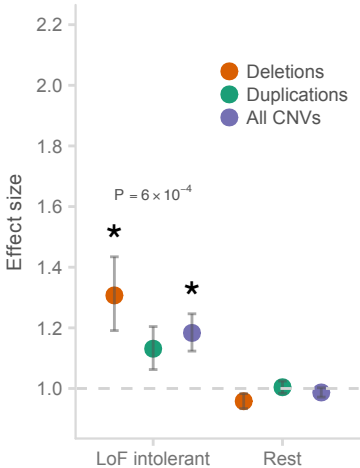
- 970  
971 44. Guipponi, M. *et al.* Exome sequencing in 53 sporadic cases of schizophrenia  
972 identifies 18 putative candidate genes. *PLoS One* **9**, e112745 (2014).  
973 45. Girard, S. L. *et al.* Increased exonic de novo mutation rate in individuals  
974 with schizophrenia. *Nat. Genet.* **43**, 860–3 (2011).  
975 46. McCarthy, S. E. *et al.* De novo mutations in schizophrenia implicate  
976 chromatin remodeling and support a genetic overlap with autism and  
977 intellectual disability. *Mol. Psychiatry* **19**, 652–8 (2014).  
978 47. Takata, A. *et al.* Loss-of-function variants in schizophrenia risk and  
979 SETD1A as a candidate susceptibility gene. *Neuron* **82**, 773–80 (2014).  
980 48. Xu, B. *et al.* Exome sequencing supports a de novo mutational paradigm for  
981 schizophrenia. *Nat. Genet.* **43**, 864–8 (2011).  
982 49. Xu, B. *et al.* De novo gene mutations highlight patterns of genetic and  
983 neural complexity in schizophrenia. *Nat. Genet.* **44**, 1365–9 (2012).  
984 50. Do, R. *et al.* Exome sequencing identifies rare LDLR and APOA5 alleles  
985 conferring risk for myocardial infarction. *Nature* **518**, 102–106 (2014).  
986 51. Kircher, M. *et al.* A general framework for estimating the relative  
987 pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–5 (2014).  
988 52. Doody, G. A., Johnstone, E. C., Sanderson, T. L., Owens, D. G. & Muir, W. J.  
989 ‘Pffropfschizophrenie’ revisited. Schizophrenia in people with mild learning  
990 disability. *Br. J. Psychiatry* **173**, 145–153 (1998).

991  
992  
993

**A****B**



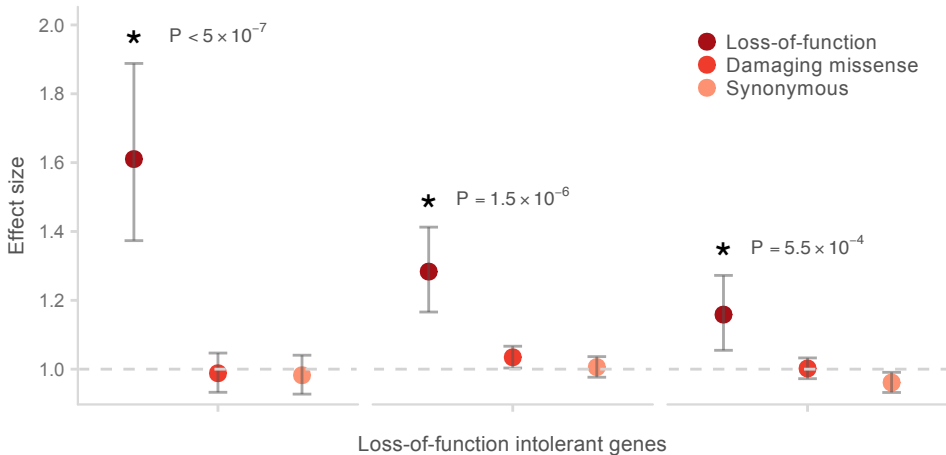




Schizophrenia individuals  
with intellectual disability  
v. controls

Schizophrenia individuals  
who did not complete  
secondary school  
v. controls

Schizophrenia individuals  
without intellectual disability  
v. controls



Schizophrenia individuals  
with intellectual disability  
v. controls

Schizophrenia individuals  
who did not complete  
secondary school  
v. controls

Schizophrenia individuals  
without intellectual disability  
v. controls

