

# 1 **Strength of functional signature correlates with effect size in autism**

2

## 3 **Authors:**

4 Sara Ballouz,<sup>1</sup>

5 Jesse Gillis\*,<sup>1</sup>

## 6 **Affiliations:**

7 1 The Stanley Institute for Cognitive Genomics, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY,

8 11724, USA

9 \* Corresponding author: Dr J Gillis, The Stanley Institute for Cognitive Genomics, Cold Spring Harbor

10 Laboratory, Cold Spring Harbor 11724, NY, USA

11 JG: [jgillis@cshl.edu](mailto:jgillis@cshl.edu)

12 SB: [sballouz@cshl.edu](mailto:sballouz@cshl.edu)

## 1 **Abstract**

### 2 **Background**

3 Disagreements over genetic signatures associated with disease have been particularly prominent in the  
4 field of psychiatric genetics, creating a sharp divide between disease burdens attributed to common and  
5 rare variation, with study designs independently targeting each. Meta-analysis within each of these  
6 study designs is routine, whether using raw data or summary statistics, but combining results across  
7 study designs is atypical. However, tests of functional convergence are used across all study designs,  
8 where candidate gene sets are assessed for overlaps with previously known properties. This suggests  
9 one possible avenue for combining not study data, but the functional conclusions that they reach.

### 10 **Method**

11 In this work, we test for functional convergence in autism spectrum disorder (ASD) across different  
12 study types, and specifically whether the degree to which a gene is implicated in autism is correlated  
13 with the degree to which it drives functional convergence. Because different study designs are  
14 distinguishable by their differences in effect size, this also provides a unified means of incorporating the  
15 impact of study design into the analysis of convergence.

### 16 **Results**

17 We detected remarkably significant positive trends in aggregate ( $p < 2.2e-16$ ) with 14 individually  
18 significant properties ( $FDR < 0.01$ ), many in areas researchers have targeted based on different reasoning,  
19 such as the fragile X mental retardation protein (*FMRP*) interactor enrichment ( $FDR 0.003$ ).

20 We are also able to detect novel technical effects and we see that network enrichment from protein-  
21 protein interaction data is heavily confounded with study design, arising readily in control data.

### 22 **Conclusions**

1 We see a convergent functional signal for a subset of known and novel functions in ASD from all sources  
2 of genetic variation. Meta-analytic approaches explicitly accounting for different study designs can be  
3 adapted to other diseases to discover novel functional associations and increase statistical power.

#### 4 **Keywords**

5 autism spectrum disorder; rare variation; common variation; loss-of-function; recurrence; effect sizes;  
6 functional enrichment; gene candidate score; meta-analysis

7

8

## 1 Background

2 Over the last decade, enormous progress has been made in characterizing sources of DNA variation  
3 contributing to disease. Most of this progress has been enabled by study designs which are carefully  
4 tailored to exploit technologies targeting particular classes of variation. Researchers have used  
5 chromosomal analysis arrays [1-4], genotyping arrays [5-8], whole-exome sequencing (WES)[9-14], and  
6 whole genome sequencing (WGS)[15, 16], to identify risk loci and alleles. The results from these studies  
7 cannot be naively compared; common variants are limited to regions of the genome with known  
8 variation (a SNP is known) but only reach significance with large numbers, while rare or ultra-rare  
9 variants are conditioned on not being in this list of common variants. Trio and quad studies are used  
10 mainly in WES and WGS study designs, while large case and control cohorts are required for signals in  
11 genome-wide association studies (GWAS). Thus for each study design, we are asking distinct questions  
12 that relate to the population prevalence, disease mechanism, burden and risk.

13 Within each study, however, it is commonplace to look to overlapping functional properties of candidate  
14 disease genes to find the biologically meaningful signal among the positive results. Candidate genes are  
15 prioritized based on enrichment analyses in pathways related to the phenotype (e.g., neuronal activity  
16 regulation) or some other disease feature shared by the genes (e.g., expression in the brain). If these  
17 methods return no significant results, more complex methods are performed to extract common  
18 features from the disease gene set [17], such as co-regulatory module detection from co-expression  
19 networks [18] or binding from protein-protein interaction (PPI) networks [19]. Regardless of the study  
20 design, the analysis with respect to functional convergence follows a similar (and largely separable)  
21 design: genes selected as hits are tested for the presence of some joint signature with the null provided  
22 by genes which are not hits. By the same logic that suggests testing hits for functional convergence  
23 relative to the background, we hypothesize that sets of genes which are “strong” hits will show more  
24 functional convergence than those which are “weak” hits.

1 We suggest that the degree of functional convergence may be hypothesized to vary (monotonically)  
2 with the degree to which genes are causal for the disease. Genes only weakly causal, whether due to  
3 high false positive rates in the study design or low effect sizes, are not strongly implicated as sharing a  
4 joint role by their co-occurrence as disease-related. For instance, disease candidates from GWAS have  
5 low relative risks (and therefore low effect sizes) as they are inherited common variation in the  
6 population. On the other hand, *de novo* mutations are a form of genetic variation which evolutionary  
7 forces have had little time to act upon [20] (e.g., unless embryonically lethal), and are of high risk (and  
8 high effect sizes). Studies also suffer from type I errors (false positives), and this too should be reflected  
9 in an aggregate disease signal of the candidate genes, as quantified by their common functional  
10 properties. A set of genes with *de novo* mutations will show a strong aggregate disease signal, while we  
11 might expect a weaker signal from the gene candidates from GWAS [21]. Measuring their ‘functional  
12 convergence’, as determined by a gene set enrichment test or network analysis, we can thus exploit our  
13 knowledge of gene candidates’ effect sizes and false positive rates. For a true disease property, we  
14 expect the correlation between gene set effect size and functional convergence to be strong, and for a  
15 weak or artifactual property, we expect no significant correlation.

16 We propose to test this hypothesis by running a meta-analytic study on autism spectrum disorder (ASD  
17 [MIM 209850]) candidates across numerous genetic studies and over a wide range of gene properties  
18 and functions. ASD is a neurodevelopmental disease commonly characterized by behavioral traits such  
19 as poor social and communication skills [22]. In more severe cases, ASD is comorbid with mild to severe  
20 intellectual disability, facial and cranial dysmorphology and gastrointestinal disorders. Perhaps because  
21 of grouping these multiple and sometimes distinct phenotypes into one disorder, and the complexity of  
22 behavior as a trait, understanding the genetic architecture of this cognitive disease has been non-trivial  
23 [23]. The genetic component of ASD is estimated to be 50-60% [24], however there are still a substantial  
24 number of cases where the underlying genetic factors of the disease are unknown. Due to these levels

1 of heterogeneity, multiple studies and study designs have been used to determine the underlying  
2 genetics which we make use of here. Taking these different studies, we construct several disease gene  
3 candidate collections, each containing genes of similar levels of risk, as determined by their odds ratios  
4 and relative risks. On every gene collection, we run a number of analyses, calculating the functional  
5 convergence using standard enrichment methods, and more complex network analysis enrichments. By  
6 exploiting trends in targeted genetic variation and their known effect sizes, we demonstrate it is  
7 possible to discriminate biologically convergent signals from likely technical artifacts at a very fine  
8 resolution. The disease properties with strong trend signals are largely consistent with the known  
9 literature on ASD (e.g., FMRP interactor enrichment) but we also see a few otherwise interesting  
10 properties as unlikely to be disease specific. Particularly protein-protein interaction networks and some  
11 co-expression networks, which extract artifactual signals from the study design, show signals in control  
12 data using that study design. Our focus here is on autism due to our interest in the disorder, its well-  
13 powered data, and also its phenotypic and genetic heterogeneity.

## 14 **Methods**

### 15 **Study design**

16 An overview of our study design and method is shown in **Fig 1**. Briefly, we start by characterizing the  
17 ASD gene sets collected for this analysis. Each study's results were collapsed individually into a set of  
18 genes, with an estimated average effect size for that candidate set (**Fig 1A**). We calculate a functional  
19 effect (e.g., statistical overlaps with known functions, **Fig 1B**) for disease-specific and more general gene  
20 properties. We then calculate the correlation of these functional convergences with the estimated effect  
21 size of that variant class (**Fig 1C**). More specifically, we test to see if the set of genes with high effect  
22 sizes have strong relative functional convergences as measured by a functional enrichment of some  
23 disease property across them, and those with low effect sizes, have weaker functional signals. We apply

1 this test to numerous functional properties on candidate gene sets from a variety of study designs.  
2 Functions with positive correlations (positive trends) we believe will show signatures that are likely  
3 associated with autism and can be used for further functional characterization of the disease.  
4 Throughout our work we refer to the “effect size” as the disease burden or risk of a gene candidate (or  
5 the average of such values within a gene set), and the “functional convergence” as the significance of a  
6 functional test for a disease gene set after controlling for the set size.

## 7 **Study data**

### 8 *Disease gene candidate sets*

9 We first collected candidate disease gene sets from available autism studies. We selected the largest  
10 study of whole-exome sequencing (WES) of families from the Simons Simplex Collection (SSC)[25]. We  
11 defined different sets of genes from over 2000 gene candidates, splitting into recurrent (at least 2  
12 probands having the mutation) and non-recurrent mutations, according to mutation type (loss-of-  
13 function, missense and silent mutations). We selected copy number variant (CNV) data also from the  
14 individuals in the SSC [26], and parsed it into similar sets. We then used the CNVs as parsed by Gilman et  
15 al.[3], which prioritized genes with their NETBAG algorithm. For GWAS gene sets, we generated two lists  
16 from the Psychiatric Genomics Consortium (PGC) study on autism and related psychiatric disorders [23]:  
17 one from the reported gene list and a second list of all adjacent genes as listed in the GWAS NHGRI-EBI  
18 catalog [6]. For our control and test gene sets, we took all the GWAs data in the GWAs catalog [6],  
19 totaling over 1,396 traits across 2,066 studies. For each trait, we created gene lists with the reported  
20 genes. We conditioned on traits with at least 27 genes, which left us with approximately 200 traits. Our  
21 negative control sets included using the genes with mutations in the unaffected siblings of the probands  
22 from the SSC studies. Overall, we had 11 gene sets for the main autism analysis, and 148 trait gene sets  
23 from GWAS.

## 1 **Gene functional annotation data**

### 2 *Co-expression networks*

3 The majority of recent studies used co-expression networks from BrainSpan to illustrate network  
4 convergence among disease genes of ASD (e.g., see [27-29]). In a similar fashion, we generated a brain  
5 specific network from the BrainSpan RNA-seq data (578 samples). In addition to this, we generated an  
6 aggregate co-expression network from 28 brain tissue and cell specific microarray experiments (3,362  
7 samples). For more general networks, we used our aggregate RNA-seq and microarray co-expression  
8 networks as previously described in Ballouz et al, [30]. In brief, these are the aggregates of 50 networks  
9 (1,970 samples) and 43 networks (5,134) samples respectively, across various tissues, cell types and  
10 conditions. As a comparison to the aggregate networks we recommend, we constructed and tested  
11 individual networks from single experiments that are more commonly used. This includes tissue-specific  
12 co-expression networks from the GTEx data [31] (29 tissues), and age specific co-expression networks (5  
13 age groups). As additional tests, we took a further 227 RNA-seq expression datasets with at least 20  
14 samples within each experiment from GEMMA [8], and have generated a further 454 individual human  
15 co-expression networks, using all annotated transcripts (30K, GENCODE [32]), and then only protein-  
16 coding genes (18K).

### 17 *Protein-protein interaction networks*

18 We used the human physical protein-protein interactions from BIOGRID (version 3.2.121)[33] and  
19 created a binary protein-protein interaction network, where each protein was a node and each protein-  
20 protein interaction is an edge. Because of the sparseness of the network, we extended the network by  
21 modelling indirect connections [34], taking the inverse minimum path length between two proteins as  
22 the weighted edge, with a maximum distance of 6 jumps roughly as described in Gillis et al,[35]. We  
23 repeated this for alternate PPI datasets including: I2D [36] (v 2.9), HPRD [37] (Release 9), HIPPIE [38]



1 (v1.8), IntAct [39], the CCSB interactome database [40](HI-III v2.2), STRING [41](v 10), and PIPs [42]. A  
2 non-interacting protein-protein network was created from data from the negatome [43](v2).

### 3 *Gene sets and collections*

4 We considered common functional gene sets and neurological specific sets, as used in numerous  
5 studies, as gene sets to test for ASD candidate enrichment. These included the post synaptic density  
6 (HPSD) gene set [44], synapse sets [45], the synaptosome [46], chromatin remodelling set [47], fragile X  
7 mental retardation protein (FMRP) set [48], and gene essentiality [49]. For more standard sets, we also  
8 took the Gene Ontology [50] (GO) terms (April 2015) and KEGG pathways [51]. For each GO term, we  
9 only used evidence codes that were not inferred electronically, propagated annotations through the  
10 ontology (parent node terms inherited the genes of their leaf node terms). To minimize redundancy  
11 from GO, we restricted our enrichment analyses to GO terms groups with sizes between 20 to 1000  
12 genes. While these GO terms and KEGG groups are used in the enrichment analyses (with the full  
13 multiple hypothesis test correction penalty). As an extension to the original study, we collected  
14 alternate gene property sets for more functional enrichment tests. For this we used all the collections  
15 from MSigDB [7] (gene sets H, C1-C7). We calculated the multifunctionality of a gene based on the  
16 number of times a gene is seen as being annotated to a function (using GO, see [52]).

### 17 *Disease gene score sets*

18 We used disease gene scoring methods that rank genes according to likely having damaging effects if  
19 they are mutated. This included the Residual Variation Intolerance Score (RVIS)[53], haploinsufficiency  
20 (HI) scores [54], mutational rates and constrained gene scores and probabilities (pLI) from ExAC [55].

### 21 *Expression data*

22 To obtain brain specific expression and differential expression information, we used three common and  
23 large sample size brain-specific transcriptomic sets. These included the Human Brain Transcriptome

1 (GSE25159)[56], BrainSpan [39] and the Human Prefrontal Cortex transcriptome (GSE30272)[57]. We  
2 divided the samples into fetal (post-conception week – PCW) and post-birth stages, and performed a  
3 straightforward differential expression (DE) fold change analysis (averaging across these stages)[58].

#### 4 **Calculating average disease effect sizes**

5 For the 11 candidate disease and control gene sets (**Table 1, Fig 2A**), we ranked the set according to the  
6 overall or average “effect size” of the genes within it. For the *de novo* mutation candidates, we took the  
7 ratio of observed counts of mutations to silent mutations within the study for that class of mutations,  
8 and then the ratio of those odds between siblings to probands (as calculated in Sanders et al,[10]). To  
9 calculate this effect size for the GWAS results, we took the average odds ratios from the individual  
10 studies of each the SNP, which ranged between 1.01-1.1. For the control sets (siblings and the silent  
11 mutations), we took the effect size to be null. We then ranked the sets based on these overall effect  
12 sizes. After these calculations, we end up with three general classes: null effects (as controls), weak  
13 effects (missense and common variants) and strong effects (rarer loss-of-function and copy number  
14 variants).

#### 15 **Calculating functional convergences**

16 Our functional tests, described below, return  $p$ -values which are dependent on the size of the gene set  
17 being considered. The statistical tests differ depending on the mode of analysis (e.g., enrichment or  
18 network), but by ‘functional convergence’ we simply mean significance ( $p$ -value) after correcting for the  
19 set size, typically by downsampling. For the downsampling, we took a subset of genes, recalculated the  
20  $p$ -value and then took geometric means of the adjusted  $p$ -values. Throughout, where we write  
21 ‘functional convergence’ it is possible to read ‘ $p$ -value after correcting for set size’.

## 1 *Network connectivity*

2 We measure the clustering of sets of genes within networks through the use of a network modularity  
3 calculation. We compare the degree of connections a gene has to all the genes in the network (global  
4 node degree), and to those of interest within the sub-network they form (local node degree). The null  
5 expectation is that genes will be connected equally well to genes within the sub-network as to those  
6 outside. Genes with large positive residuals have more weighted internal connections than external  
7 connections, implying a well inter-connected module. We test the significance of this distribution of  
8 residuals to a null set (random similarly sized set of genes, Mann-Whitney-Wilcoxon test, `wilcox.test` in  
9 R) to determine our test statistic.

## 10 *Gene set enrichment testing*

11 As a way to determine the level of enrichment of the candidate gene sets within other functional sets,  
12 we used a hypergeometric test with multiple test correction (`phyper` in R). The downsampled  $p$ -value  
13 was used as the functional convergence measure.

## 14 *Disease gene property testing*

15 For the disease gene scoring properties, we tested the significance of the scores of the candidate genes  
16 using the Mann-Whitney-Wilcoxon test (`wilcox.test` in R). The functional convergence was the  $p$ -value of  
17 this test.

## 18 *Measuring functional convergence trends*

19 For each gene property tested, we then measured the “trend” by calculating the correlation of the  
20 ranked functional effect sizes of our gene sets, whereby the gene sets are ordered according to their  
21 effect size ranks. A positive correlation is one where the function tested is correlated with our ordering.  
22 We computed this using Spearman’s rank coefficient to capture the degree of variation, but the  
23 significant subsets identified are generally robust to choice of measurement metric such as the

1 Pearson's coefficient. We limited our functional convergence tests to the subset of functions where at  
2 least one gene set of the 11 showed a significant functional convergence signal ( $p < 0.05$ ). In essence this  
3 filtering removes gene sets where there are, for example, no overlaps with any disease sets and should  
4 not affect our analysis.

## 5 **Determining significance of the functional convergence trends**

6 To calculate a null, we permute the labels of the gene sets, and calculate the functional convergence  
7 trends. Note that in the ranked case, this is simply the null distribution of a spearman correlation, with  
8 similarly associated significances. We first filter for functional tests where any one of the disease and  
9 control gene sets have a functional convergence of 0.05, but report both pre- and post- filtering results.  
10 Because our hypothesis (and test) are concerned with the ordering of functional effect sizes, filtering so  
11 that the data has at least one significant value changes the null distribution only slightly (e.g., probability  
12 of ties). We calculate the number of significant correlations based on the false discovery rate (FDR) at  
13 0.01 and 0.05. Known confounds of disease gene sets are gene length [59] and gene multifunctionality,  
14 and to test this we generated matched gene set controls by sampling genes with similar gene lengths,  
15 GO multifunctionality and disease multifunctionality measures. Using the ranked CDS (coding DNA  
16 sequence) region of the genes, we generated sets of genes of similar ranked length distributions to the  
17 11 real gene sets in the analysis. Downsampled, we then ran the analyses on these gene sets that are  
18 specifically not involved in the phenotype. This was repeated for multifunctionality as calculated using  
19 GO and then disease (using Phenocarta [60]).

## 20 **Results**

### 21 **Little overlap of the autism candidate genes across gene sets of different effect sizes**

22 We find genes with loss-of-function *de novo* mutations to be little implicated in GWA studies, with only  
23 4 candidate genes overlapping those two sets (**Fig 2B**, hypergeometric test  $p = 0.76$ ). Interestingly, the

1 more recurrent genes in the loss-of-function *de novo* set, the more unlikely they are to be found in other  
2 gene sets. For gene sets with the lower average effect sizes (e.g., the genes with missense mutations),  
3 their overlap with other gene sets is greater, in particular with the control sets (**Fig 2B**, hypergeometric  
4 tests  $p \sim 4.4e-3$  to  $2.4e-6$ ). The *de novo* variants are conditioned on being rare (low frequency) and novel  
5 by not appearing in the parents. The SNPs used in GWAS are generally conditioned on being common by  
6 having minor allele frequencies greater than 0.05 [61]. Even if this filtering is done on the variant level,  
7 and not on the gene level, it still creates selection trends within our observations of variants and thus  
8 genes. This is possibly a version of Berkson's effect[62] – where selecting for an outcome generates  
9 negative correlations between potential causes for it. An additional cause is largely technical; since  
10 we've conditioned on frequency, genes with higher mutability are depleted in our rare lists, and  
11 enriched in our common lists. Thus the lack of overlap is at least potentially not largely reflective of  
12 underlying genetics or biology, but likely due to the selection bias in obtaining them. There is also poor  
13 overlap within the rarer variation itself, for instance of genes within CNVs and those with loss-of-  
14 function SNVs (3 genes,  $p \sim 0.37$ ); there is generally a discrepancy between study designs focused on  
15 (different) sources of rare variation, and not just rare versus common. It should be noted that whether  
16 biological or technical, the lack of overlap does nothing to discredit either common or rare variation as a  
17 contributor to the disease – but highlights the need for a framework to combine and analyze the results  
18 of these studies that is aware of these biases and can distinguish biology from technical effects.

### 19 **Functional convergence trends as shown through enrichment and connectivity tests**

20 While enrichment analysis is comparatively straightforward, we demonstrate an example in **Fig 3A** using  
21 the genes with *de novo* loss-of-function mutations from lossifov et al,[1] (341 genes) and their overlap  
22 with essential genes (see Methods). In **Fig 3A**, we represent this enrichment test as a Venn diagram of  
23 the overlap of the candidate disease gene set with the essential gene set, and calculate the significance  
24 of the overlap with a hypergeometric test ( $n=82$ ,  $p \sim 9.8e-9$ ). We continue this analysis on the other

1 candidate disease gene sets from recent ASD studies, varying across study designs and technologies  
2 (WES, GWAS and arrays). Splitting each gene set by mutational class, recurrence and gender, we  
3 perform the same hypergeometric tests. To make comparable assessments between studies and gene  
4 sets, we calculate the functional convergence by downsampling – selecting a subset of genes within that  
5 set and averaging the results over a 1000 permutations (schematic in **Fig 3B**). Taking a representative set  
6 of studies (**Table 1**), we use the degree of disease effect to rank these sets, noting that recurrence leads  
7 to a higher effect size even for variation and study designs of the same class by reducing the number of  
8 false positives. Placing the controls sets on the far left, and the highest disease rank set (recurrent *de*  
9 *novo* loss-of-function genes) on the far right, and plotting their functional effect values, we observe an  
10 upward trend (**Fig 3C** Spearman's  $r_s=0.95$ , Fisher's transformation  $p<8.24e-06$ ). The slope (i.e., the  
11 correlation) of this trend line represents the “functional convergence trend”, with higher correlations  
12 indicating higher functional effects.

13 A less common (likely due to complexity) yet important functional test is network connectivity. Genes  
14 that are co-regulated or form parts of a functional unit, protein complex or pathway, are preferentially  
15 co-expressed, and this information is captured in co-expression networks. We next demonstrate how  
16 network-style effect sizes can be similarly calculated through a modularity analysis. In **Fig 3D**, we plot  
17 the global node degrees (x-axis) against their connectivity to the remainder of genes in the set (y-axis).  
18 In the null (grey line), the genes would be connected to other autism genes in proportion to the  
19 incidence of those genes within the genome. Deviations from this null across all genes generate excess  
20 modularity within this set (studentized residuals shown in inset **Fig 3D**) and determine the statistical  
21 results reported for the set overall (Wilcoxon test). A large number of genes are highly interconnected in  
22 this set, as shown by the number of points above the line (Wilcoxon test on the studentized residuals,  
23  $p\sim 7.83e-41$ ). It is important to note that this network analysis is calculated against the empirical null for  
24 each gene individually (x-axis) and so is unaffected by any gene-specific bias (such as length). Only

1 higher-order topological properties across gene-gene relationships for a given gene can produce a  
2 signal. Even assortativity, the tendency for genes of high node degree to preferentially interact, is quite  
3 low within this data ( $r=0.064$ ). As in the previous steps, we repeat the network connectivity tests across  
4 all gene sets (**Fig 3E**), also downsampling to calculate the functional convergence. Once again, gene sets  
5 with higher proportion of burden genes correlate with functional convergence tests (**Fig 3E**, Spearman's  
6  $r_s=0.69$ , Fisher's transformation  $p<0.02$ ).

### 7 **A subset of functional properties are correlated with disease effect sizes**

8 We extend our analysis to other disease gene property tests, and calculate their effect size correlations,  
9 plotting the distribution of correlations in **Fig 4A** (4210 functional tests performed, 4164 with calculable  
10 correlations). We then calculated the null distribution for the variation across effect sizes by permuting  
11 the estimated effect size for each real set and rerunning our analysis. Only limiting our functional tests  
12 to those where we had at least one gene set returning a significant enrichment signal, we observe a  
13 strong signal (61 tests, **Fig 4A**, 14 functions  $FDR<0.01$  **Table 2**). Reducing the stringency of the  
14 underlying enrichment (383 tests, **Fig 4B**), we observe a weaker signal (10 functions  $FDR<0.01$ ).  
15 Removing the underlying enrichment constraint, we observe that most functional tests are ordered  
16 consistent with the null, with a few highly correlated functions (**Fig 4C** enrichment at positive end, 3  
17 functions  $FDR<0.01$ ). The results are broadly reassuring that some weak artifact is not driving the  
18 tendency of the functional convergence and effect size to be correlated because that correlation occurs  
19 almost exclusively where the underlying tests themselves are detecting significance. In other words, the  
20 ordering of significances is only non-random where the underlying values are also non-random. We  
21 focus on the 14 functional properties identified in the first filtered assessment (**Table 2**).  
22 Each property can be defined by its vector of effect sizes across gene sets and so we can cluster the  
23 properties by their Euclidean distance in this space. Taking the 61 properties and highlighting the

1 properties that are significant (FDR 0.01), they split into approximately 7 clusters and a singleton (**Fig**  
2 **4E**). The interesting clusters are 1 and 7 as they have the highest correlations (as depicted by the dark  
3 purple scale), and a stronger significant signal from the *de novo* set (white/yellow in heatmap). Cluster 1,  
4 specifically, has the most consistent trends and contains the expression analyses (overexpression and  
5 fold change), the gene essentiality scores and some of the neural gene sets. Cluster 3 has the co-  
6 expression networks clustered, and the mutational probabilities, but is slightly weaker as the control  
7 sets also show some enrichment. Cluster 5 contains most of the GO groups. Cluster 6 has some tests  
8 which are functionally enriched in the CNV and missense gene sets but are not significantly enriched for  
9 any of the genes in the *de novo* recurrent gene set and are thus not showing a substantially positive  
10 functional convergence trend. The clustering speaks to the similarity of some of the tests (i.e., GO  
11 groups clustering), but also to a likely neuronal signature across the disease gene sets.

## 12 **Significant functional properties are consistent with the autism literature**

13 One of the properties with the highest correlation was network connectivity in the BrainSpan co-  
14 expression network; however, all disease gene sets had a significant functional convergence with  
15 Brainspan, indicating that in addition to the real signal, there is a background signal affecting even  
16 control data. In particular, the signal from the silent recurrent mutations in the probands (functional  
17 convergence  $p=7.5e-7$ ) shows that control data subject to only one study design may select genes in a  
18 highly non-random pattern. Most top scoring disease properties are consistent with the literature on  
19 autism candidates such as average RVIS and haploinsufficiency scores, [63] along with gene length and  
20 enrichment for FMRP interactors. RVIS scores are highly enriched in the loss-of-function recurrent set  
21 and the CNVs, but not significant in any of the other sets (**Fig 5A**); as with any meta-analysis significance  
22 in any one set is not necessary for aggregate significance. Genes with high haploinsufficiency scores –  
23 those that cannot maintain normal function with a single copy - are overrepresented in the loss-of-  
24 function recurrent genes, and there is also a significant effect in the GWAS results. Many interaction



1 networks and traditional functional categories appear to be poor candidates to determine convergence  
2 in disease genes, as they cluster control gene sets and sets of low effects as well as those of disease  
3 genes. For instance, the extended PPI network has a high effect in the sibling controls sets (e.g., silent  
4 functional convergence  $p \sim 1.3e-5$ , **Fig 5B**). GO terms and KEGG pathways typically do not survive  
5 correcting for multiple testing, although there is a general deviation from the null and the extremal GO  
6 functions are concordant with the known literature (e.g., GO: 0016568 chromatin modification  
7 hypergeometric test  $p \sim 1e-3$  for the *de novo* recurrent set or GO: 0048667 cell morphogenesis involved  
8 in neuron differentiation, hypergeometric test  $p \sim 0.04$  for the CNV set). So although functional  
9 convergence trends are concentrated in more clearly disease related -properties such as RVIS,  
10 traditional functional categories from, e.g., GO remain of modest use.

### 11 **Robustness and relative contributions of study designs and variants**

12 In order to determine whether the functional convergence trend rose preferentially from a subset of  
13 studies, we conducted a series of robustness analyses (**Additional file 1: Fig S1**). Ideally, the significant  
14 functional convergence trend we see is due only to effect size estimates across studies which are  
15 themselves robust. Nor do we want the trends to be strongly affected by ordering of the gene sets with  
16 similar effect sizes. Even though the average effect sizes for the GWAS sets were the same, the number  
17 of false positives within these sets varies, and this was incorporated into the ranking scheme. It is also  
18 arguable that the silent mutations in the probands may have some regulatory effect, or are false  
19 negatives. As a more stringent test, we removed whole classes of variants from the analyses (e.g., all  
20 the controls or all the common/weaker gene sets) and calculated the trends once again (**Additional file**  
21 **1: Fig S1D-F**). This is a negative control experiment in the sense that if the functional convergence trend  
22 arises meta-analytically, it should be largely robust to changing things we are not certain about (e.g., as  
23 above, whether effect sizes are 1.1 or 1.09) and not robust to changing things we are certain about  
24 (common variants play some role in autism). Removing either controls, genes sets with the highest

1 effects or the common variants from the trend analyses removes all the number of significant  
2 correlations although some deviation from the null remained (**Additional file 1: Fig S1**). When rare  
3 variants are excluded, the distribution of correlations is most similar to the null, but still significantly  
4 different (Student's paired T-test  $p \sim 0.03$ ), while the total significance of the test is closest to the full  
5 version when common variation is excluded (Student's paired T-test  $p \sim 8.2e-7$ ). Since our common data  
6 is likely the weakest due to the tremendous focus of autism data collection toward rare variation in the  
7 SSC, this makes sense, but common variation still contributes substantial joint signal. These tests  
8 confirm that the approximate order of gene sets by effect sizes correctly drives the results and that we  
9 are robust to minor variation in the exact effect sizes listed, but do rely on the joint use of the extremely  
10 divergent study results (rare and common) within the meta-analysis to attain significant results.

11 To control for the impact of gene length and multifunctionality (number of functions a gene is listed as  
12 possessing), we repeated a control version of our analyses. In this case, the real disease gene sets were  
13 swapped out with gene sets matched with respect to multifunctionality or length. We then reran the  
14 evaluation of functional convergence trends to determine if any previously identified properties arise as  
15 correlated with these control sets (ordered by their match to a specific disease sets, e.g., gene length  
16 distribution). Repeating the analysis in this control case, we find the derived correlations are for the  
17 most part extremely similar to the null (reference). We can additionally use these controls versions as a  
18 slightly more stringent null distribution for expected correlations when we evaluate the real disease  
19 sets. In the analysis where we do not condition on the underlying tests having reached some level of  
20 significance (as in **Fig 4A**), we see even more correlations passing significance (**Additional file 1: Fig S2**),  
21 indicating the multifunctionality or gene length do little to explain the general trends we see.

22 Promiscuous or absent enrichment have both historically been problematic within disease gene data;  
23 both diminish the specificity of functional results. When too many functions are returned from an

1 analysis, we need to cherry pick and with too few, we have no “leads” and are left in the dark. We  
2 suggest that the strong aggregate effect we see and small number of significant functions is likely near  
3 to a useful and biologically plausible type of specificity for downstream analysis, as suggested by the fact  
4 that ad hoc filtering (i.e., top ten lists) usually are at about this level when not constrained by  
5 significance. Our set of functional tests and results are shown in **Additional file 2: Table S1** and the full  
6 data set is available online.

7 One potential failure mode of this analysis comes from the GWAs we have used. Because the number of  
8 autism GWAs available and well-powered for analysis was relatively small, we used a combined  
9 psychiatric genomics dataset, which included bipolar and schizophrenia. We now wish to test how  
10 specific our results were to our disease and not a signal of GWAs in general. We repeated our analysis  
11 using each of the 148 GWAS traits in the GWAS catalog that had enough genes to be included in our  
12 tests. We did not recalculate the effect sizes specifically for each, but used the mean estimates from the  
13 autism set. Using the number of correlations calculated as significant to rank the 148 traits, the top ten  
14 traits include the autism and schizophrenia GWAs, and a few larger studies such as “Body mass index”  
15 (**Additional file 1: Fig S3**). This is a fairly striking confirmation of our original hypothesis: the degree of  
16 correlation between functional convergence is so specific that it correctly distinguishes particular  
17 disease sets as belonging to the same trend (as defined by a particular disease). The larger GWAs, also  
18 found in the top 10, are not related to psychiatric disorders show a signal in very broad disease  
19 properties, such as the gene mutability scores.

## 20 **Expanding the functional gene tests show no further significant properties**

21 We wished to see if we could find other significant associations if we expanded our repertoire of  
22 functions within each type of test. Our first set of network analyses focused on general aggregate co-  
23 expression networks and brain sample only aggregates. In most analyses, researchers use individual

1 datasets to build their networks and we wished to compare our results to these. Thus we expanded our  
2 tests to a total of 540 networks. We repeated the same analysis, using an additional 6 PPI networks, 76  
3 condition specific networks (tissues, sex, and age), and a further 454 RNAseq co-expression networks  
4 (227 across 18K protein coding genes and 227 across 30K transcripts). Once again, we see functional  
5 convergence across almost all the gene sets with little ordered trend by effect size. The network  
6 convergence exists in even the control data and is therefore likely due to study selection biases alone;  
7 none pass an FDR of 0.01 (**Additional file 1: Fig S4A**).

8 Initially, we focused on expression data for the brain, but were curious about how tissue-specific these  
9 patterns were, or whether the genes were generally highly expressed. To this end, we repeated the  
10 expression analyses using tissue specific expression datasets from GTEx data [31]. We were also curious  
11 to determine if we could see sex specific differences, and used additional data from the GEUVADIS  
12 project [64]. Repeating the functional convergence trends on all these expression datasets shows little  
13 to no significant expression in the individual gene sets, and no significant functional correlations  
14 (**Additional file 1: Fig S4B**). The functional test with the greatest correlation were also from brain  
15 specific expression datasets ( $r_s=0.78$ ).

16 One last set of gene properties typically used by researchers in their analyses are the curated gene sets  
17 from MSigDB. We repeated our analyses on all 8 collections, and calculated the functional convergence  
18 trends, using the hypergeometric test as in the case of calculating enrichment in GO. The gene sets  
19 range from curated data sets from the known literature, to computationally derived gene sets from  
20 cancer microarrays. Perhaps unsurprisingly, we see no enrichment in these gene sets (**Additional file 1:**  
21 **Fig S5**), as most are inflammatory or oncogenic collections, or versions of GO terms and KEGG pathways  
22 which we had already found to have no enrichment.

## 1 Discussion

2 Our contribution in this work has been to establish that there is a significant correlation between the  
3 effect size of candidate autism genes and the degree to which tests assessing their functional  
4 convergence find a signal. As with any meta-analysis, the hope is that by incorporating multiple data, the  
5 aggregate signal may be stronger. While our work suggests an approach to do this and shows strong  
6 statistical trends, we anticipate that divisions in the field of genetics may play some role in the  
7 interpretation of this work [65, 66]. GWAS researchers may question the power of our GWAS analysis  
8 and assume that if it just had a large enough 'N', it alone would be the dominant player in understanding  
9 psychiatric disease. Similarly, they may suspect rare variants of reflecting 'anomalous' versions of the  
10 disorder and thus think they are less likely to be specifically linked to autism. Rare variant researchers  
11 may question the precision of the data underlying our rare variant analysis and assume that if we just  
12 had large enough 'N' to remove false positives, it would be the dominant player in understanding  
13 psychiatric disease. Similarly, they may suspect that GWAS data is affected by confounds such as  
14 population stratification. Both reactions are perfectly reasonable. Our analysis does not establish that  
15 all functional properties are distributed across all classes of autism, but rather, for a subset, there is a  
16 very significant trend. This is further supported by the functions that arise being either of specific  
17 relevance to autism or of well-known importance to disease in general; however, even this division  
18 implies that we are capturing multiple factors affecting the genetic architecture of disease.

19 Likewise, our specific experimental design for assessing a relationship between functional convergence  
20 and effect size may well be open to elaboration and emendation. Our principle interest was in ensuring  
21 that any observed trends would be reflective of tests and data used within the literature and not choices  
22 of our own. Thus, we aimed for the simplest and most conventional means of assessing functional  
23 convergence and focused on being exhaustive (in terms of properties) within this domain, rather than  
24 optimizing our design for observing functional convergence trends. We also developed a framework

1 which is readily extensible to new tests, regardless of their form or complexity. In each case, the test  
2 can be equivalently applied to the downsampled disease gene sets and the significance of any  
3 correlation simply calculated as would be conventional via permutation test. This is both simple,  
4 general, and allows easy comparison across studies using the equivalent approach; however, more  
5 theoretically grounded or alternative means of calculating functional convergence trends are surely  
6 possible. Particular weaknesses in our design are our use of non-parametric tests and downsampling to  
7 control for set size. These are, we think, natural choices for robustness but more finely tuned  
8 alternatives are likely to exist and could easily be a target of research since our results suggests the  
9 observation of key functional convergence trends is highly robust and salient within the data.

10 As the number of disease gene sets expands, and further refinement of risk assessment is achieved, the  
11 resolution of functional convergence trends should grow. Indeed, incorporating effect size as a meta-  
12 analytic constraint offers a diverse range of novel applications. That integration may be across study  
13 designs and classes of variation, as we have done, or may involve phenotype or other properties. So, for  
14 example, one could determine functional convergence trends that grow or shrink depending on how  
15 patients were classified, or even broken down in a sex-specific manner for interpreting protective  
16 effects. More broadly, as data and the means for obtaining it grows, techniques to statistically assess its  
17 structured dependencies will grow more useful and important. Our robustness analysis speaks to this in  
18 that while we are robust to modest losses of data, it is clear that more data will only improve the signals  
19 of the individual classes. More finely-tuned effect size estimates and better separations of the gene sets  
20 and variant classifications will also help refine the distinction between biological and artifactual signals,  
21 ideally allowing us conduct yet more focused study designs in a productive feedback loop.

## 1 **Conclusions**

2 In this work we have found that the stronger the effect size of autism candidate genes, the more likely  
3 they are to exhibit a joint functional signal. The functional properties identified exhibit some specificity  
4 to autism and neuropsychiatric disease (e.g. FMRP interactors), but also some more general links to  
5 disease (e.g., RVIS). While there remains substantial heterogeneity between study designs and the  
6 genetic architectures of disease which they may uncover, we have shown that there is some  
7 commonality across study designs. The commonality across study designs is not a literal overlap in risk  
8 genes, or even functional effect, but that functions weakly identified in GWA studies are likely to be  
9 more strongly identified in rare variation studies. As evidence for autism and other disorders continues  
10 to develop and continues to be heterogeneous with respect to ascertainment biases and study designs,  
11 we suspect approaches related to the one we describe will be of increasing importance.

## 12 **Declarations**

## 13 **Abbreviations**

14 ASD: autism spectrum disorder  
15 CNV: copy number variant  
16 FDR: false discovery rate  
17 FMRP: fragile X mental retardation protein  
18 FWER: family-wise error rate  
19 GO: Gene Ontology  
20 GWAS: genome-wide association studies  
21 HI: haploinsufficiency  
22 LoF: Loss-of-function mutations  
23 PGC: Psychiatric Genomics Consortium  
24 PPI: protein-protein interactions  
25 PPIN: protein-protein interaction network  
26 RVIS: Residual Variation Intolerance Score  
27 SNP: single nucleotide polymorphism  
28 SNV: single nucleotide variant  
29 SSC: Simons Simplex Collection  
30 WES: whole exome sequencing  
31 WGS: whole genome sequencing

32  
33

1 **Availability of data and materials**

2 The datasets supporting the conclusions of this article are included within the article and in the  
3 additional files.

4 ***Additional files***

5 **Additional file 1: Figures**

6 **Additional file 1: Fig S1 Trend line robustness analysis.**

7 **Additional file 1: Fig S2 Functional convergences null for matched length and multifunctionality**  
8 **controls**

9 **Additional file 1: Fig S3 Functional convergence correlation/trend distributions for GWA studies.**

10 **Additional file 1: Fig S4 Functional convergence correlation/trend distributions for all network connectivity**  
11 **tests and all gene expression tests**

12 **Additional file 1: Fig S5 Functional convergence correlation/trend distributions for all MSigDB collections.**

13 **Additional file 2: Tables**

14 **Additional file 2: Table S1 Functional convergence correlation/trend distributions**

15 **Additional file 2: Table S2 Additional expression functional convergences and correlations/trends**

16 **Additional file 2: Table S3 Additional gene properties functional convergences and**  
17 **correlations/trends**

18 **Additional file 2: Table S4 Additional MSigDB functional convergences and correlations/trends**

19 **Additional file 2: Table S5 Additional network connectivity functional convergences and**  
20 **correlations/trends**



1

## 2 **Competing interests**

3 The authors declare that they have no competing interests.

## 4 **Funding**

5 This work was supported by a grant from T. and V. Stanley.

## 6 **Authors' contributions**

7 SB wrote the manuscript, and conducted the experiments. JG wrote the manuscript and designed the  
8 experiments. All authors read and approved the final manuscript.

## 9 **Acknowledgments**

10 We would like to thank members of the CSHL Wigler lab for access to their data. We thank Paul Pavlidis  
11 for helpful comments on a draft of the manuscript.

## 12 **References**

- 13 1. Sebat J, Lakshmi B, Malhotra D, Troge J, Lese-Martin C, Walsh T, Yamrom B, Yoon S, Krasnitz A,  
14 Kendall J, et al: **Strong Association of De Novo Copy Number Mutations with Autism.** *Science*  
15 2007, **316**:445-449.
- 16 2. Christian SL, Brune CW, Sudi J, Kumar RA, Liu S, KaraMohamed S, Badner JA, Matsui S, Conroy J,  
17 McQuaid D, et al: **Novel Submicroscopic chromosomal abnormalities detected in Autism**  
18 **Spectrum Disorder.** *Biological psychiatry* 2008, **63**:1111-1117.
- 19 3. Marshall CR, Noor A, Vincent JB, Lionel AC, Feuk L, Skaug J, Shago M, Moessner R, Pinto D, Ren  
20 Y, et al: **Structural Variation of Chromosomes in Autism Spectrum Disorder.** *American Journal*  
21 *of Human Genetics* 2008, **82**:477-488.
- 22 4. Glessner JT, Wang K, Cai G, Korvatska O, Kim CE, Wood S, Zhang H, Estes A, Brune CW, Bradfield  
23 JP, et al: **Autism genome-wide copy number variation reveals ubiquitin and neuronal genes.**  
24 *Nature* 2009, **459**:569-573.
- 25 5. Ma DQ, Salyakina D, Jaworski JM, Konidari I, Whitehead PL, Andersen AN, Hoffman JD, Slifer SH,  
26 Hedges DJ, Cukier HN, et al: **A genome-wide association study of autism reveals a common**  
27 **novel risk locus at 5p14.1.** *Annals of human genetics* 2009, **73**:263-273.
- 28 6. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P, Manolio T,  
29 Hindorff L, Parkinson H: **The NHGRI GWAS Catalog, a curated resource of SNP-trait**  
30 **associations.** *Nucleic Acids Research* 2014, **42**:D1001-D1006.
- 31 7. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A,  
32 Pomeroy SL, Golub TR, Lander ES, Mesirov JP: **Gene set enrichment analysis: A knowledge-**

- 1           **based approach for interpreting genome-wide expression profiles.** *Proceedings of the National*  
2           *Academy of Sciences* 2005, **102**:15545-15550.
- 3    8.    Zoubarev A, Hamer KM, Keshav KD, McCarthy EL, Santos JRC, Van Rossum T, McDonald C, Hall A,  
4    Wan X, Lim R: **Gemma: a resource for the reuse, sharing and meta-analysis of expression**  
5    **profiling data.** *Bioinformatics* 2012, **28**:2272-2273.
- 6    9.    O'Roak BJ, Deriziotis P, Lee C, Vives L, Schwartz JJ, Girirajan S, Karakoc E, MacKenzie AP, Ng SB,  
7    Baker C: **Exome sequencing in sporadic autism spectrum disorders identifies severe de novo**  
8    **mutations.** *Nature genetics* 2011, **43**:585-589.
- 9    10.    Sanders SJ, Murtha MT, Gupta AR, Murdoch JD, Raubeson MJ, Willsey AJ, Ercan-Sencicek AG,  
10    DiLullo NM, Parikshak NN, Stein JL: **De novo mutations revealed by whole-exome sequencing**  
11    **are strongly associated with autism.** *Nature* 2012, **485**:237-241.
- 12   11.    Iossifov I, Ronemus M, Levy D, Wang Z, Hakker I, Rosenbaum J, Yamrom B, Lee Y-h, Narzisi G,  
13    Leotta A, et al: **De Novo Gene Disruptions in Children on the Autistic Spectrum.** *Neuron* 2012,  
14    **74**:285-299.
- 15   12.    Neale BM, Kou Y, Liu L, Ma'ayan A, Samocha KE, Sabo A, Lin CF, Stevens C, Wang LS, Makarov V,  
16    et al: **Patterns and rates of exonic de novo mutations in autism spectrum disorders.** *Nature*  
17    2012, **485**:242-245.
- 18   13.    Lohmueller KE, Pearce CL, Pike M, Lander ES, Hirschhorn JN: **Meta-analysis of genetic**  
19    **association studies supports a contribution of common variants to susceptibility to common**  
20    **disease.** *Nat Genet* 2003, **33**:177-182.
- 21   14.    Wang K, Zhang H, Ma D, Bucan M, Glessner JT, Abrahams BS, Salyakina D, Imielinski M, Bradfield  
22    JP, Sleiman PM, et al: **Common genetic variants on 5p14.1 associate with autism spectrum**  
23    **disorders.** *Nature* 2009, **459**:528-533.
- 24   15.    Cobb JP, Mindrinos MN, Miller-Graziano C, Calvano SE, Baker HV, Xiao W, Laudanski K,  
25    Brownstein BH, Elson CM, Hayden DL, et al: **Application of genome-wide expression analysis to**  
26    **human health and disease.** *Proceedings of the National Academy of Sciences of the United*  
27    *States of America* 2005, **102**:4801-4806.
- 28   16.    Deng Q, Ramsköld D, Reinius B, Sandberg R: **Single-Cell RNA-Seq Reveals Dynamic, Random**  
29    **Monoallelic Gene Expression in Mammalian Cells.** *Science* 2014, **343**:193-196.
- 30   17.    Bauer-Mehren A, Bundschuh M, Rautschka M, Mayer MA, Sanz F, Furlong LI: **Gene-Disease**  
31    **Network Analysis Reveals Functional Modules in Mendelian, Complex and Environmental**  
32    **Diseases.** *PLoS ONE* 2011, **6**:e20284.
- 33   18.    Ben-David E, Shifman S: **Networks of Neuronal Genes Affected by Common and Rare Variants**  
34    **in Autism Spectrum Disorders.** *PLoS Genet* 2012, **8**:e1002556.
- 35   19.    Sakai Y, Shaw CA, Dawson BC, Dugas DV, Al-Mohtaseb Z, Hill DE, Zoghbi HY: **Protein Interactome**  
36    **Reveals Converging Molecular Pathways Among Autism Disorders.** *Science Translational*  
37    *Medicine* 2011, **3**:86ra49.
- 38   20.    Sullivan PF, Daly MJ, O'Donovan M: **Genetic architectures of psychiatric disorders: the**  
39    **emerging picture and its implications.** *Nat Rev Genet* 2012, **13**:537-551.
- 40   21.    Schizophrenia Working Group of the Psychiatric Genomics C: **Biological insights from 108**  
41    **schizophrenia-associated genetic loci.** *Nature* 2014, **511**:421-427.
- 42   22.    Miles JH: **Autism spectrum disorders--a genetics review.** *Genetics in Medicine* 2011, **13**:278-  
43    294.
- 44   23.    Cross-Disorder Group of the Psychiatric Genomics C: **Identification of risk loci with shared**  
45    **effects on five major psychiatric disorders: a genome-wide analysis.** *Lancet* 2013, **381**:1371-  
46    1379.
- 47   24.    Gibson G: **Rare and common variants: twenty arguments.** *Nat Rev Genet* 2012, **13**:135-145.

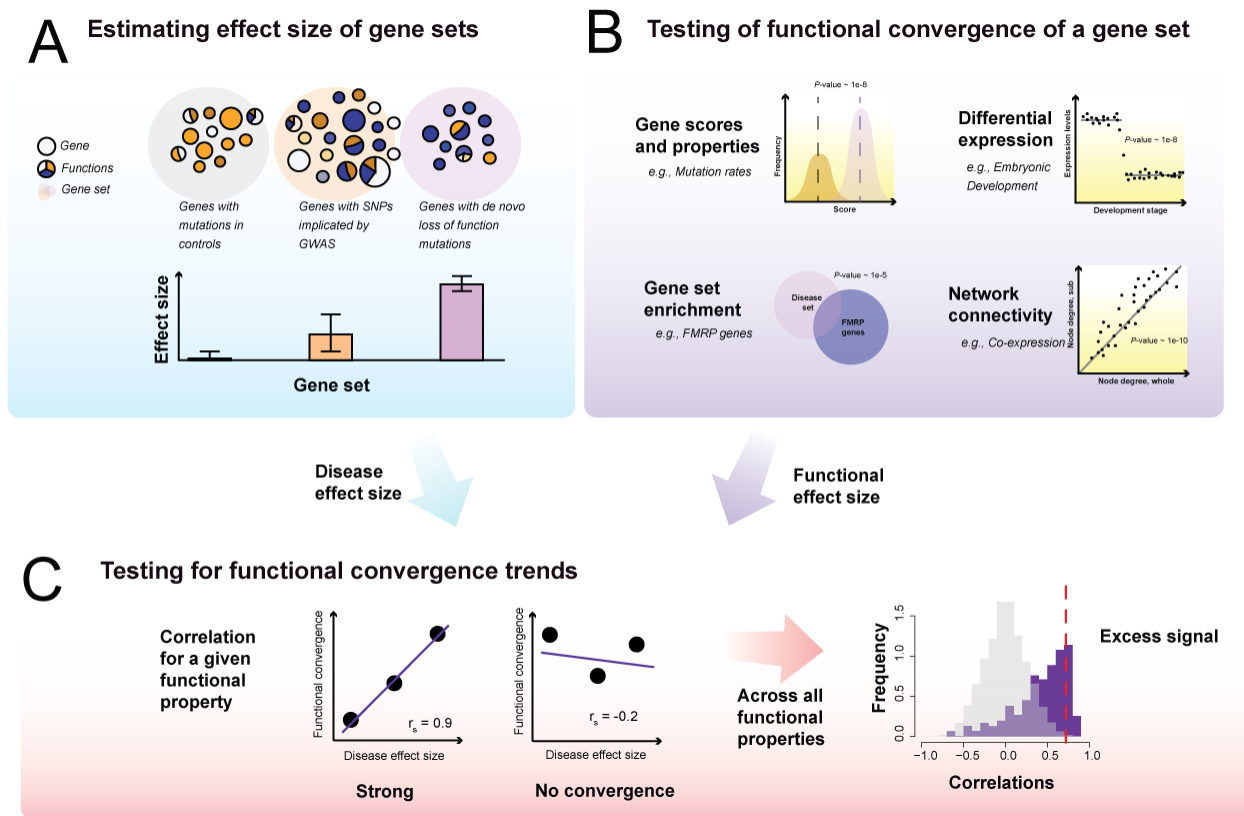
- 1 25. Iossifov I, O’Roak BJ, Sanders SJ, Ronemus M, Krumm N, Levy D, Stessman HA, Witherspoon K,  
2 Vives L, Patterson KE, et al: **The contribution of de novo coding mutations to autism spectrum  
3 disorder.** *Nature* 2014, **515**:216-221.
- 4 26. Levy D, Ronemus M, Yamrom B, Lee YH, Leotta A, Kendall J, Marks S, Lakshmi B, Pai D, Ye K, et  
5 al: **Rare de novo and transmitted copy-number variation in autistic spectrum disorders.**  
6 *Neuron* 2011, **70**:886-897.
- 7 27. Willsey AJ, Sanders Stephan J, Li M, Dong S, Tebbenkamp Andrew T, Muhle Rebecca A, Reilly  
8 Steven K, Lin L, Fertuzinhos S, Miller Jeremy A, et al: **Coexpression Networks Implicate Human  
9 Midfetal Deep Cortical Projection Neurons in the Pathogenesis of Autism.** *Cell* 2013, **155**:997-  
10 1007.
- 11 28. Sugathan A, Biagioli M, Golzio C, Erdin S, Blumenthal I, Manavalan P, Ragavendran A, Brand H,  
12 Lucente D, Miles J, et al: **CHD8 regulates neurodevelopmental pathways associated with  
13 autism spectrum disorder in neural progenitors.** *Proceedings of the National Academy of  
14 Sciences* 2014, **111**:E4468-E4477.
- 15 29. Hormozdiari F, Penn O, Borenstein E, Eichler EE: **The discovery of integrated gene networks for  
16 autism and related disorders.** *Genome Research* 2015, **25**:142-154.
- 17 30. Verleyen W, Ballouz S, Gillis J: **Positive and negative forms of replicability in gene network  
18 analysis.** *Bioinformatics* 2015.
- 19 31. Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R, Walters G, Garcia F, Young N,  
20 et al: **The Genotype-Tissue Expression (GTEx) project.** *Nature Genetics* 2013, **45**:580-585.
- 21 32. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D,  
22 Zadissa A, Searle S, et al: **GENCODE: the reference human genome annotation for The ENCODE  
23 Project.** *Genome Research* 2012, **22**:1760-1774.
- 24 33. Stark C, Breitkreutz B-J, Reguly T, Boucher L, Breitkreutz A, Tyers M: **BioGRID: a general  
25 repository for interaction datasets.** *Nucleic Acids Research* 2006, **34**:D535-D539.
- 26 34. Chua HN, Sung W-K, Wong L: **Exploiting indirect neighbours and topological weight to predict  
27 protein function from protein–protein interactions.** *Bioinformatics* 2006, **22**:1623-1630.
- 28 35. Gillis J, Pavlidis P: **The role of indirect connections in gene networks in predicting function.**  
29 *Bioinformatics* 2011, **27**:1860-1866.
- 30 36. Brown KR, Jurisica I: **Unequal evolutionary conservation of human protein interactions in  
31 interologous networks.** *Genome Biology* 2007, **8**:1-11.
- 32 37. Peri S, Navarro JD, Kristiansen TZ, Amanchy R, Surendranath V, Muthusamy B, Gandhi TKB,  
33 Chandrika KN, Deshpande N, Suresh S, et al: **Human protein reference database as a discovery  
34 resource for proteomics.** *Nucleic Acids Research* 2004, **32**:D497-D501.
- 35 38. Schaefer MH, Fontaine J-F, Vinayagam A, Porras P, Wanker EE, Andrade-Navarro MA: **HIPPIE:  
36 Integrating Protein Interaction Networks with Experiment Based Quality Scores.** *PLoS ONE*  
37 2012, **7**:e31826.
- 38 39. Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, Campbell NH, Chavali G,  
39 Chen C, del-Toro N, et al: **The MIntAct project—IntAct as a common curation platform for 11  
40 molecular interaction databases.** *Nucleic Acids Research* 2014, **42**:D358-D363.
- 41 40. Rolland T, Taşan M, Charlotheaux B, Pevzner Samuel J, Zhong Q, Sahni N, Yi S, Lemmens I,  
42 Fontanillo C, Mosca R, et al: **A Proteome-Scale Map of the Human Interactome Network.** *Cell*,  
43 **159**:1212-1226.
- 44 41. Szklarczyk D, Franceschini A, Wyder S, Forslund K, Heller D, Huerta-Cepas J, Simonovic M, Roth  
45 A, Santos A, Tsafou KP, et al: **STRING v10: protein–protein interaction networks, integrated  
46 over the tree of life.** *Nucleic Acids Research* 2015, **43**:D447-D452.
- 47 42. McDowall MD, Scott MS, Barton GJ: **PIPs: human protein–protein interaction prediction  
48 database.** *Nucleic Acids Research* 2009, **37**:D651-D656.

- 1 43. Blohm P, Frishman G, Smialowski P, Goebels F, Wachinger B, Ruepp A, Frishman D: **Negatome**  
2 **2.0: a database of non-interacting proteins derived by literature mining, manual annotation**  
3 **and protein structure analysis.** *Nucleic Acids Research* 2014, **42**:D396-D400.
- 4 44. Ronan JL, Wu W, Crabtree GR: **From neural development to cognition: unexpected roles for**  
5 **chromatin.** *Nature reviews Genetics* 2013, **14**:347-359.
- 6 45. Bayés À, Collins MO, Croning MDR, van de Lagemaat LN, Choudhary JS, Grant SGN: **Comparative**  
7 **Study of Human and Mouse Postsynaptic Proteomes Finds High Compositional Conservation**  
8 **and Abundance Differences for Key Synaptic Proteins.** *PLoS ONE* 2012, **7**:e46683.
- 9 46. Collins MO, Husi H, Yu L, Brandon JM, Anderson CNG, Blackstock WP, Choudhary JS, Grant SGN:  
10 **Molecular characterization and comparison of the components and multiprotein complexes in**  
11 **the postsynaptic proteome.** *Journal of neurochemistry* 2006, **97 Suppl 1**:16-23.
- 12 47. Petrovski S, Wang Q, Heinzen EL, Allen AS, Goldstein DB: **Genic Intolerance to Functional**  
13 **Variation and the Interpretation of Personal Genomes.** *PLoS Genet* 2013, **9**:e1003709.
- 14 48. Georgi B, Voight BF, Bućan M: **From Mouse to Human: Evolutionary Genomics Analysis of**  
15 **Human Orthologs of Essential Genes.** *PLoS Genet* 2013, **9**:e1003484.
- 16 49. Huang N, Lee I, Marcotte EM, Hurles ME: **Characterising and Predicting Haploinsufficiency in**  
17 **the Human Genome.** *PLoS Genet* 2010, **6**:e1001154.
- 18 50. Verleyen W, Ballouz S, Gillis J: **Measuring the wisdom of the crowds in network-based gene**  
19 **function inference.** *Bioinformatics* 2015, **31**:745-752.
- 20 51. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes.** *Nucleic Acids Research*  
21 2000, **28**:27-30.
- 22 52. Gillis J, Pavlidis P: **The Impact of Multifunctional Genes on "Guilt by Association" Analysis.** *PLoS*  
23 *ONE* 2011, **6**:e17258.
- 24 53. Darnell JC, Van Driesche SJ, Zhang C, Hung KYS, Mele A, Fraser CE, Stone EF, Chen C, Fak JJ, Chi  
25 SW, et al: **FMRP stalls ribosomal translocation on mRNAs linked to synaptic function and**  
26 **autism.** *Cell* 2011, **146**:247-261.
- 27 54. Lips ES, Cornelisse LN, Toonen RF, Min JL, Hultman CM, Holmans PA, O'Donovan MC, Purcell SM,  
28 Smit AB, Verhage M, et al: **Functional gene group analysis identifies synaptic gene groups as**  
29 **risk factor for schizophrenia.** *Mol Psychiatry* 2012, **17**:996-1006.
- 30 55. Lek M, Karczewski K, Minikel E, Samocha K, Banks E, Fennell T, O'Donnell-Luria A, Ware J, Hill A,  
31 Cummings B: **Analysis of protein-coding genetic variation in 60,706 humans.** *BioRxiv*  
32 2016:030338.
- 33 56. Levy D, Ronemus M, Yamrom B, Lee Y-h, Leotta A, Kendall J, Marks S, Lakshmi B, Pai D, Ye K, et  
34 al: **Rare De Novo and Transmitted Copy-Number Variation in Autistic Spectrum Disorders.**  
35 *Neuron* 2011, **70**:886-897.
- 36 57. Colantuoni C, Lipska BK, Ye T, Hyde TM, Tao R, Leek JT, Colantuoni EA, Elkahouloun AG, Herman  
37 MM, Weinberger DR, Kleinman JE: **Temporal dynamics and genetic control of transcription in**  
38 **the human prefrontal cortex.** *Nature* 2011, **478**:519-523.
- 39 58. Chang J, Gilman SR, Chiang AH, Sanders SJ, Vitkup D: **Genotype to phenotype relationships in**  
40 **autism spectrum disorders.** *Nat Neurosci* 2015, **18**:191-198.
- 41 59. Ouwenga RL, Dougherty J: **Fmrp targets or not: long, highly brain-expressed genes tend to be**  
42 **implicated in autism and brain disorders.** *Molecular Autism* 2015, **6**:16.
- 43 60. Portales-Casamar E, Ch'ng C, Lui F, St-Georges N, Zoubarev A, Lai A, Lee M, Kwok C, Kwok W,  
44 Tseng L, Pavlidis P: **Neurocarta: aggregating and sharing disease-gene relations for the**  
45 **neurosciences.** *BMC Genomics* 2013, **14**:129.
- 46 61. Bush WS, Moore JH: **Chapter 11: Genome-Wide Association Studies.** *PLoS Comput Biol* 2012,  
47 **8**:e1002822.

- 1 62. Westreich D: **Berkson’s bias, selection bias, and missing data.** *Epidemiology (Cambridge, Mass)*  
2 2012, **23**:159-164.
- 3 63. Ji X, Kember RL, Brown CD, Bućan M: **Increased burden of deleterious variants in essential**  
4 **genes in autism spectrum disorder.** *Proceedings of the National Academy of Sciences* 2016,  
5 **113**:15054-15059.
- 6 64. Lappalainen T, Sammeth M, Friedlander MR, t Hoen PAC, Monlong J, Rivas MA, Gonzalez-Porta  
7 M, Kurbatova N, Griebel T, Ferreira PG, et al: **Transcriptome and genome sequencing uncovers**  
8 **functional variation in humans.** *Nature* 2013, **501**:506-511.
- 9 65. Krumm N, Turner TN, Baker C, Vives L, Mohajeri K, Witherspoon K, Raja A, Coe BP, Stessman HA,  
10 He Z-X: **Excess of rare, inherited truncating mutations in autism.** *Nature genetics* 2015, **47**:582-  
11 588.
- 12 66. Pinto D, Delaby E, Merico D, Barbosa M, Merikangas A, Klei L, Thiruvahindrapuram B, Xu X,  
13 Ziman R, Wang Z, et al: **Convergence of Genes and Cellular Pathways Dysregulated in Autism**  
14 **Spectrum Disorders.** *The American Journal of Human Genetics* 2014, **94**:677-694.  
15  
16  
17

1

2 **Figures**



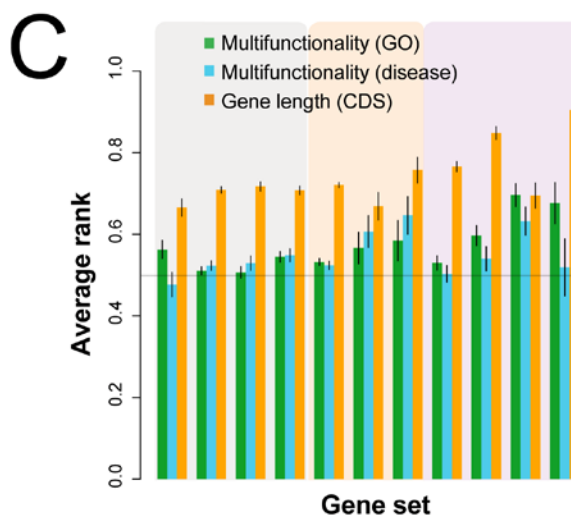
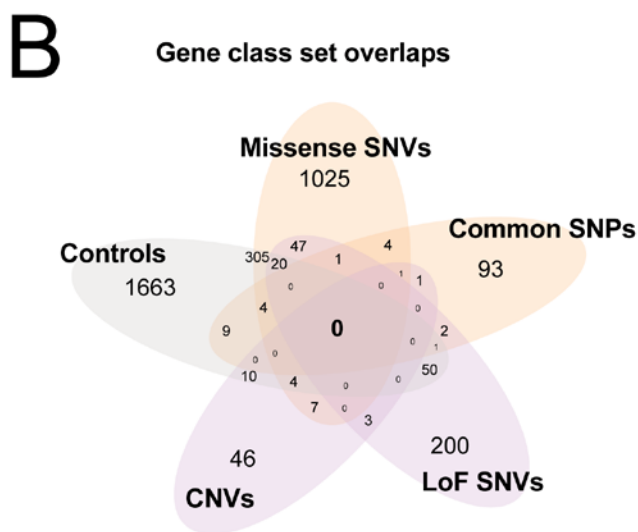
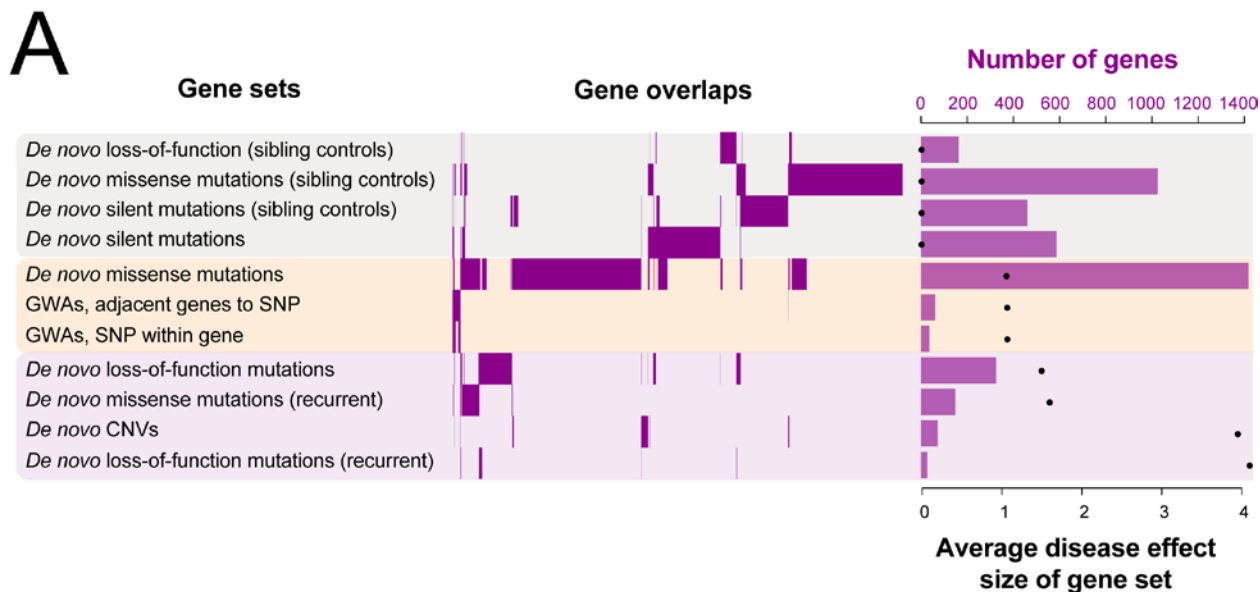
3

4 **Fig 1 Schematic of functional convergence trend calculation.**

5 (A) Starting with disease gene set collections, we rank each by the average effect size of the genes within that set.  
 6 (B) We then run ‘functional tests’ on these genes sets and calculate a functional convergence for each. (C) Then,  
 7 using the ranking of the disease gene sets, we measure the functional convergence signature – the correlation of  
 8 the trend line of the functional convergences versus the rank.

9





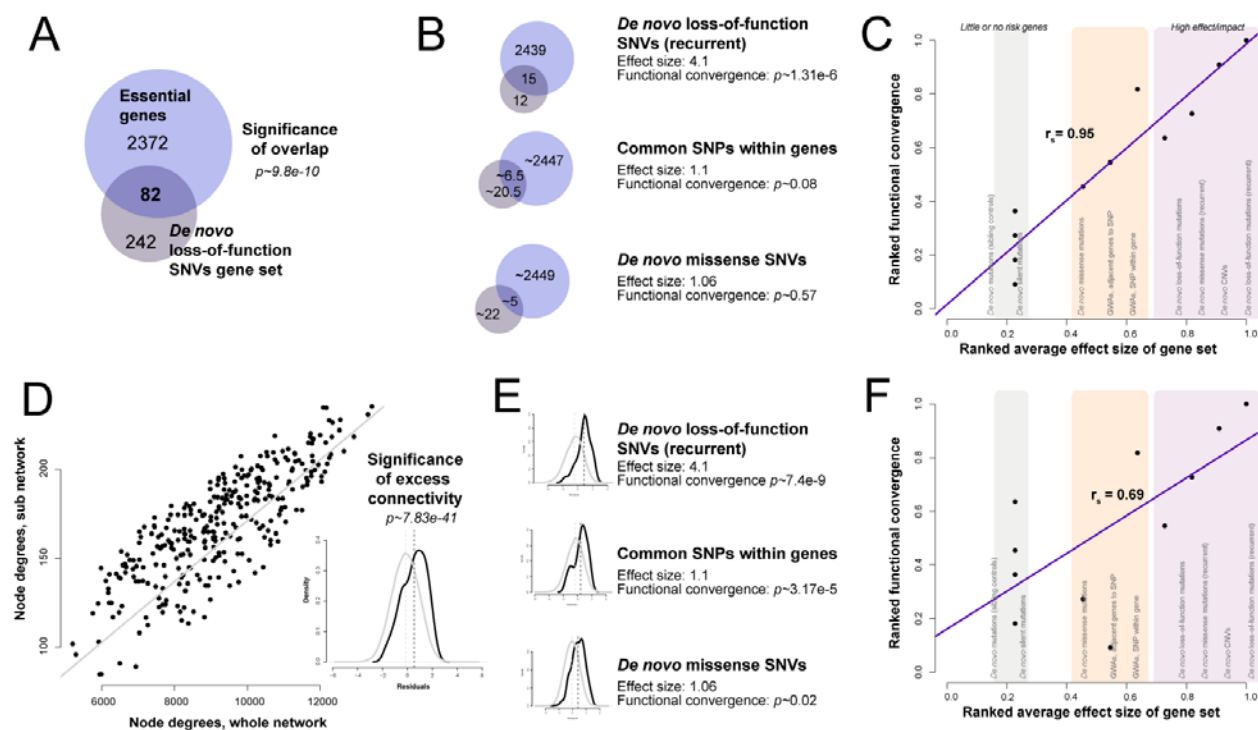
1

2 **Fig 2 Characterization of autism candidate gene sets.**

3 (A) We first classified the 11 gene sets used in the study into three larger groups: no effects, weak effects and  
 4 strong effects. We see little overlap in the individual gene sets themselves (mid panel). The total number of genes  
 5 in each set also varies (right-most panel), and is negatively correlated with the average effect size ( $r_s=-0.69$ ). (B)  
 6 Control gene sets overlap significantly with missense genes (333 genes hypergeometric test  $p=2.54e-6$ ), common  
 7 SNPs gene sets (11 genes,  $p=4.5e-3$ ), and the loss-of-function (LoF) SNV gene sets (71 genes,  $p=3.2e-3$ ) but not the  
 8 CNV gene sets (14 genes,  $p=0.03$ ). Missense and common SNPs overlap significantly (4 genes,  $p=2.4e-4$ ). However,  
 9 loss-of-function SNVs do not overlap significantly with either common (4 genes,  $p=0.62$ ), missense (68 genes,  
 10  $p=0.75$ ), or CNVs (3 genes,  $p=0.37$ ). (C) Common biases that affect studies are gene length and number of  
 11 functional annotations. The average standardized rank (+/- SE) of genes with respect to these properties shows  
 12 that the "rare" disease sets contain longer genes but are not much more multifunctional than random.

13

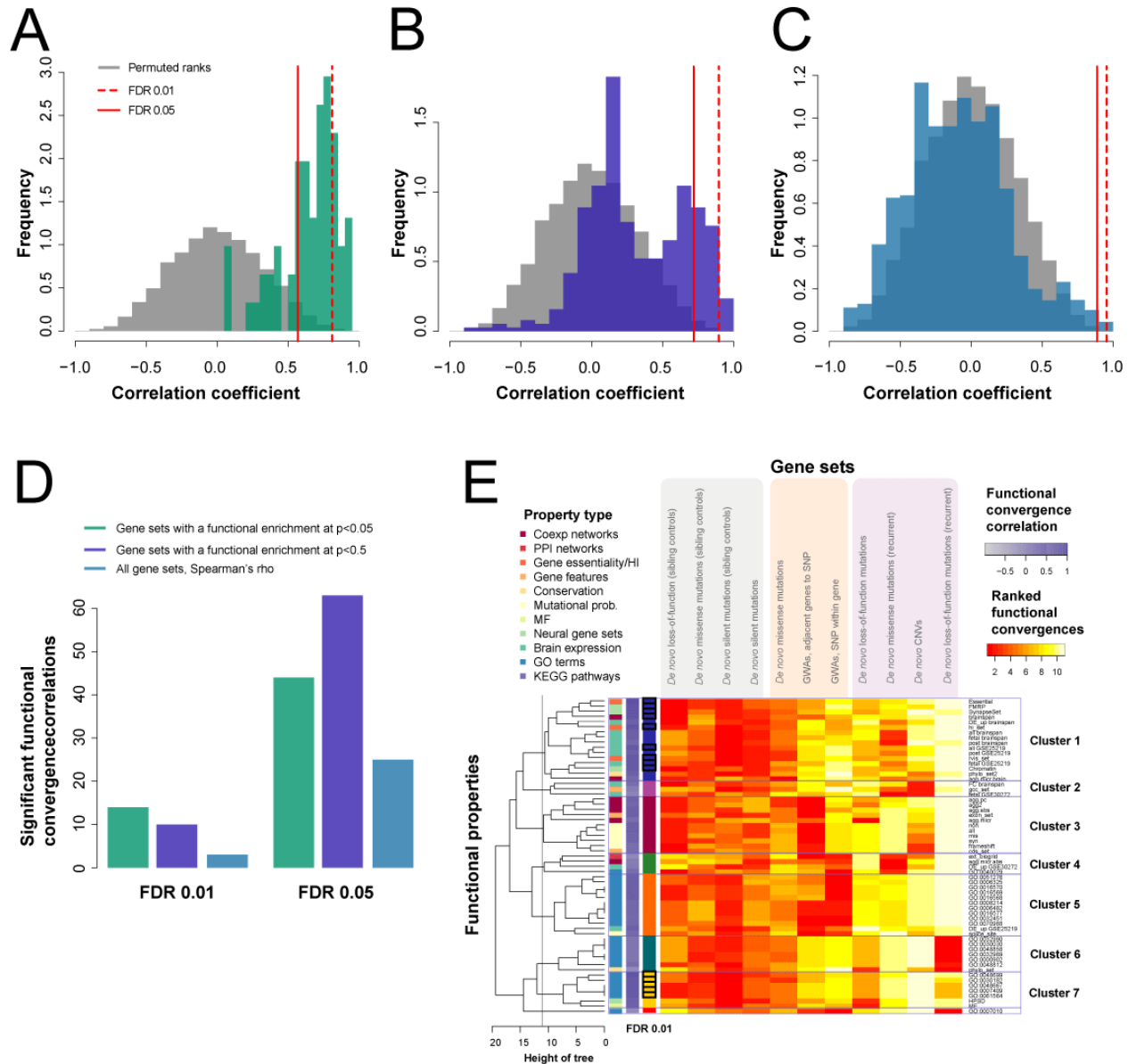
14



**Fig 3 Functional properties of disease gene sets are tested using gene set enrichment (top panels) and co-expression network connectivity (bottom panels).**

(A) Gene set enrichment is calculated with a hypergeometric test. A large number (34%) of the genes in the *de novo* loss-of-function set overlap with essential genes (hypergeometric test  $p \sim 9.8e-10$ ). (B) This is repeated across all disease gene sets, (a subset shown here). Sample size is controlled through downsampling. Gene sets with the higher effect sizes also have the higher functional convergences. (C) We can now demonstrate how to calculate the functional convergence trend for the “essential genes” test. The disease gene sets are ordered by an estimate of the average effect size of genes within the set (from low to high on the x-axis) and the functional convergence is of that disease gene set is plotted (y-axis). A trend between the effect size of the candidate genes and their essentiality can be clearly observed. The network connectivity functional test (D) consists of calculating the ratio of disease genes’ total connectivity (node degrees calculated from the whole network; sum of their connections) to their internal connectivity (node degrees of their subnetwork; sum of their connections to one another). The line (in grey) reflects the expected values if there is no preferential connectivity within the set. We see that a large number (72%) of the genes lie above the identity line. The Wilcoxon p-value of the mean residuals is shown in the inset ( $p \sim 7.83e-41$ ) (E) Once again, controlling for sample size through downsampling, the functional convergence of each gene set is calculated (subset shown). (F) A weak trend between the effect size of the candidate genes and their degree of co-expression is visible. Empirical nulls are calculated by permuting disease gene sets and FDRs through a Benjamini-Hochberg correction against the resultant functional convergence trends.



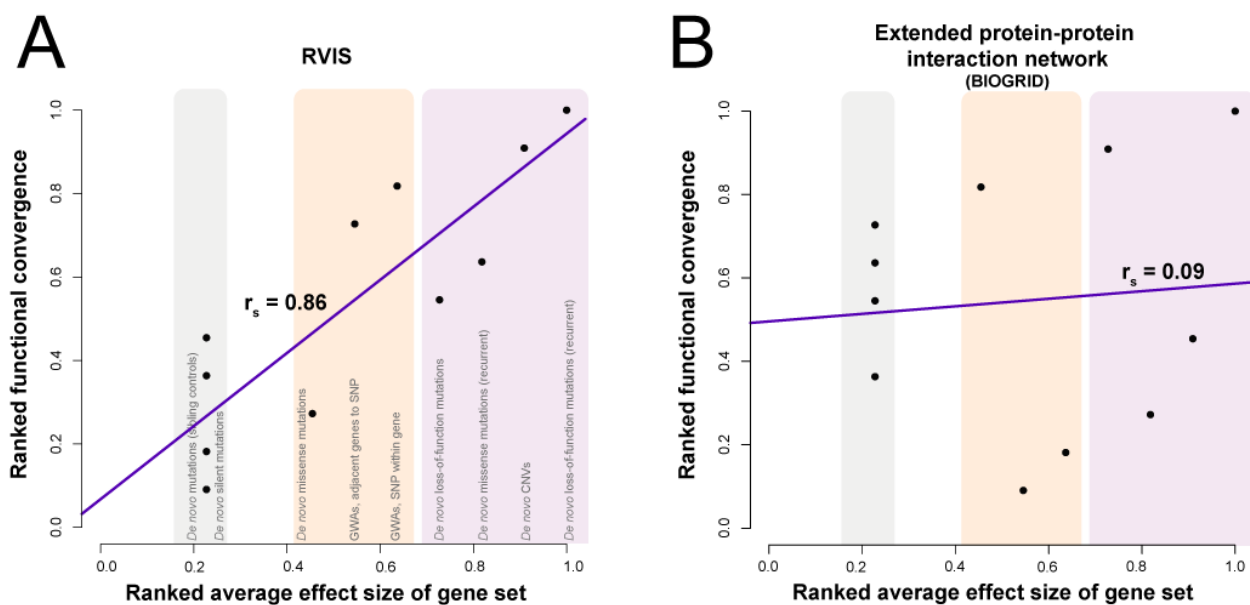


1

2 **Fig 4 Clustering disease property tests by their functional convergence trends.**

3 (A) The correlations of the ranked effect size trends are significantly different from null distributions for the 61  
4 functional properties (effect size permutation null, Student's paired T-test  $p < 2.2e-16$ ). We've drawn the red line to  
5 indicate the false discovery rates (FDRs) of 0.01, where 14 functions are significant, and 0.05 where 44 are  
6 significant. (B) If we filter for the functions with some weak signal in the underlying functional tests ( $p < 0.5$ ), 383  
7 correlations are considered with 10 functions as significant (dark blue, Student's paired T-test  $p < 2.2e-16$ ). Note  
8 that we are not filtering with respect to our own functional effect size test, which assesses variation in the  
9 underlying functional tests, merely that the underlying tests do return some values. (C) And when we have no  
10 constraints (all tests, 4164 shown), 3 pass an FDR of 0.01. (D) We enumerate these in a barplot. (E) A heatmap of  
11 all the ranked scores of the test gene sets (columns) for the subset of significant effect 61 properties (rows). The  
12 properties clustered into 6 groups when we cut the dendrogram at a height of  $\sim 12$ . Their functional convergence  
13 correlations (ranked) show that most high correlations cluster (in clusters 1 and 3). White/yellow is a high rank, red  
14 is low. The property type is color coded as described in the figure key. A high correlation is shaded purple,  
15 low/negative correlations are grey. Clusters are labeled and colored, with functions FDR < 0.01 outlined in black.

1



2

3 **Fig 5 Example correlations and slopes of functional convergence.**

4 (A) RVIS enrichment test has a very high correlation across the gene sets (B) while the extended PPIN has mostly  
5 artifactual signal as demonstrated by the flatness of the line, highly elevated from the null but in no effect-  
6 correlated way and suggestive of consistent biases.

7

## 1 Tables and captions

### 2 Table 1 Disease gene sets used in the study

Gene set	Set size (genes)	Odds Ratios/ effect sizes	Rank
<b>WES results.</b>			
<b><i>Varying effect sizes depending on mutation and recurrence.</i></b>			
<b><i>Odds ratio is calculated as in Sanders et al,[10] (ratio of observed counts of mutations to silent mutations, and then ratio of those odds between siblings to probands).</i></b>			
<b><i>De novo</i></b> loss-of-function, recurrent	27	4.1	11
<b><i>De novo</i></b> CNVs	72	3.95	10
<b><i>De novo</i></b> missense, recurrent	153	1.6	9
<b><i>De novo</i></b> loss-of-function	341	1.5	8
<b><i>De novo</i></b> missense	1339	1.06	5
<b><i>GWAS. Multiple hit models, low effect size per gene. Odds ratios between 1 and 1.1</i></b>			
GWAS, reported genes	49	1.08	7
GWAS, adjacent genes to SNP	116	1.08	6
<b><i>Control sets, no effect</i></b>			
<b><i>De novo</i></b> silent	590	NA	2.5
<b><i>Control groups, no disease</i></b>			
<b><i>De novo</i></b> loss-of-function (sibling controls)	174	NA	2.5
<b><i>De novo</i></b> missense (sibling controls)	1066	NA	2.5
<b><i>De novo</i></b> silent (sibling controls)	468	NA	2.5

3

4

1

2 **Table 2 Functional properties with significant functional convergence trends**

<b>Functional property</b>	<b>Spearman's correlation</b>	<b>FDR</b>
Gene essentiality (Georgi et al.)	0.95	0.001
GO:0048699 (generation of neurons)	0.92	0.003
FMRP interactors (Darnell et al.)	0.91	0.003
Brainspan co-expression network	0.90	0.003
SynapseSet (Lips et al.)	0.90	0.003
Haploinsufficiency scores	0.89	0.003
RVIS	0.86	0.007
GSE25219: Overall expression	0.84	0.008
GO:0007409 (axonogenesis)	0.84	0.006
GO:0030182 (neuron differentiation)	0.84	0.006
GO:0048667 (cell morphogenesis involved in neuron differentiation)	0.84	0.006
GO:0061564 (axon development)	0.84	0.006
Chromatin remodeling gene set (Ronan et al.)	0.83	0.006
GSE25219: Fetal expression	0.83	0.006

3

4