

Forecasting autism gene discovery with machine learning and genome-scale data

Leo Brueggeman^{1,2,3}, Tanner Koomar^{1,3}, Jacob J Michaelson^{1,4,5}

¹Department of Psychiatry, Carver College of Medicine

²Medical Scientist Training Program, Carver College of Medicine

³Interdisciplinary Graduate Program in Genetics

⁴Department of Biomedical Engineering, College of Engineering

⁵Department of Communication Sciences and Disorders, College of Liberal Arts and Sciences

University of Iowa

Iowa City, IA 52242

USA

Abstract

Genes are one of the most powerful windows into the biology of autism, and it has been estimated that perhaps a thousand or more genes may confer risk. However, less than 100 genes are currently viewed as having robust enough evidence to be considered true "autism genes". Massive genetic studies are underway to produce data to implicate additional genes, but this approach, although necessary, is costly and slow-

moving. Here, we approach autism gene discovery as a machine learning problem, rather than a genetic association problem, and use genome-scale data as predictors for identifying further genes that have similar properties in the feature space compared to established autism risk genes. This approach, which we call forecASD, integrates spatiotemporal gene expression, heterogeneous network data, and previous gene-level predictors of autism association to yield a single score that represents each gene's likelihood of being involved in the etiology of autism. We demonstrate that forecASD has substantially increased sensitivity and specificity compared to previous gene-level predictors of autism association, including genetic-based measures such as TADA. On an independent test set, consisting of newly-released pilot data from the SPARK Genomics Consortium, we show that forecASD best predicts which genes will have an excess of likely gene disrupting (LGD) mutations. Using forecASD results, we show which molecular pathways are currently under-represented in the autism literature and likely represent under-appreciated biological mechanisms of autism. Finally, the larger importance of this work is that by enumerating the genes that are most likely involved in the pathogenesis of autism, we have an opportunity to consider what molecular research in autism might look like in a post-gene discovery era.

Introduction

Autism Spectrum Disorder (ASD) is a heterogeneous grouping of developmental disorders caused by a range of genetic and environmental factors. The core diagnostic features of ASD, which manifest at a young age, are impairments in social communication and restrictive and repetitive behaviors and interests. Evidence for the role of genetics in ASD is strong, with monozygotic twins having near 90% concordance of ASD diagnosis¹. Further population and twin studies have confirmed these findings², and further estimated the narrow-sense heritability of ASD to be in the range of 50-95%.

While there is an abundance of evidence for the role of genetics in autism, our understanding of the genetic etiology of the disorder is still limited. It is estimated that there may over 1000 genes which contribute to autism risk³. However, the current list of high-confidence autism genes stands at 84 genes⁴. This discrepancy is partly explained by the relatively limited number of genomic studies compared with the vast genetic heterogeneity underlying autism.

To close this gap between the number of anticipated and known autism genes, several network-biology approaches have been applied in the past decade. These studies leverage large, publicly-available datasets to add context and amplify the genetic signals observed through sequencing studies. These network-biology studies have predicted genes that then became bona fide autism genes⁵, but have fallen short of providing a useful genome-wide metric that indicates the evidence of autism

involvement for every gene. More recently, machine learning based methods have used gene interaction networks⁶ and cell-specific expression profiles⁷ to predict gene involvement in autism. Importantly, the results of these studies lead to a quantitative metric that scores every gene in the genome according to evidence of a role in autism. Despite the demonstrated effectiveness of these studies in prioritizing autism risk genes, our preliminary investigations suggested there was still room for appreciable improvement in the form of the classification algorithm, the training set, and the predictors used. In particular, these approaches do not incorporate indicators of autism involvement that are based on genetic association (e.g., TADA scores) into their predictive features.

We introduce a new score, forecASD, that integrates prior network-biology approaches, scores of genetic association, brain gene expression, and topological information from large gene interaction networks relevant to the brain into a single gene-level score for autism involvement. We show that forecASD successfully outperforms existing methods in a diverse range of gene and mutation prioritization tasks. Further, using the recent sequencing studies MSSNG⁸ and SPARK⁹, we show that forecASD generalizes to previously unseen data. Importantly, this generalization holds even when excluding genes with known links to autism, emphasizing forecASD's ability to identify novel ASD genes. Through comparing the top decile of forecASD identified genes (1787 genes; hereafter forecASD genes) with known autism genes, we identify numerous biological pathways that are currently underrepresented in our understanding of autism risk. By reanalyzing the results of autism brain differential gene expression studies, we show

that the current list of known autism genes is significantly depleted for upregulated biological pathways, whereas forecASD captures both up- and downregulated pathways. We show that direction of differential expression is related to haploinsufficiency status, with low pLI genes showing a trend towards upregulation. Importantly, this relationship between direction of differential expression and pLI is dependent on forecASD gene inclusion, signifying forecASD's ability to capture low- and high-pLI disease genes. Through these studies, we show evidence that current methods of autism gene discovery have biases, and that forecASD mitigates these biases through its integrative approach, thus providing a view of the full spectrum of genes and biological pathways underlying autism.

Methods

Overview:

The forecASD method relies upon stacked Random Forest models, organized in two levels (shown in Figure 1). In the first level, two models are trained using BrainSpan¹⁰ gene expression and the STRING¹¹ shortest paths network as features, respectively. Our training dataset consists of high-confidence genes scored in SFARI gene⁴ as either 1 or 2 (SFARI HC genes), and 1,000 random background genes not contained within SFARI gene. These two models produce genome-wide predictions for autism involvement. These scores are then used as features in the second level's Random Forest model, along with other genome-wide scores obtained from previous studies.

BrainSpan, STRING, and TADA data assembly

BrainSpan data was obtained from the Allen Institute, and brain regions containing fewer than 20 samples were excluded. This filtered BrainSpan dataset was loess-smoothed, with the purpose of reducing noise and imputing missing data points.

The STRING database¹¹ was thresholded at their recommended score of 0.4, and transformed into a gene by gene matrix with each cell representing the shortest path between two genes.

TADA summary statistics were downloaded from the largest meta-analysis for autism available at the time of publication¹². TADA summary statistics were also obtained from the secondary supplementary table of another comprehensive study of autism risk³. All available TADA summary statistics were used as features in the final model, with `tadaFdrAscSscExomeSscAgpSmallDel`¹² used as a representative comparator in the ROC curve displayed in Figure 3.

Model training and genome wide prediction

We used a stacked Random Forest classifier to generate genome-wide predictions of autism gene involvement. All models were trained using SFARI HC genes as positive examples (of which there are 76 common to both STRING & BrainSpan), and a randomly sampled set of 1,000 background genes (i.e., not listed in the SFARI Gene database) as negative examples.

The first level of our stacked model consists of two genome-wide scores based on data from BrainSpan or the STRING interaction network. The features used in training these two models include the loess-smoothed observations in the BrainSpan database, and the STRING shortest path matrix, respectively. The random forest models were trained with 1000 total trees constructed, and the strata option enabled to insure a balance of 70 positive and 70 negative training examples during the construction of each tree. Given the large number of features for the STRING-based random forest model, we performed feature selection, wherein each feature not used in any of the constructed trees was dropped. This variable selection step was repeated until the final model contained only features which were selected at least once during tree construction. With the STRING and BrainSpan models, we then predicted autism involvement scores for the remaining genes not included in our training set. These scores are in Supplementary Table 1 in the columns BrainSpan_score and STRING_score. Scores for training set genes are the out-of-bag estimates.

We used these scores, along with DAWN⁵, TADA^{12,3}, DAMAGES⁷, and the score from Krishnan *et al.*⁶ score, as predictive features in a final Random Forest, using the same training labels described previously. Genome-wide predictions were then obtained, again using out-of-bag estimates for training set genes. This final score is listed under forecASD_score in Supplementary Table 1.

SPARK and MSSNG data sources

De novo mutation (DNM) data from the MSSNG dataset was obtained through the *de novo db* database⁶. Mutations were filtered for LGD or missense status. De novo mutation data was obtained from the SPARK dataset from the consortium's recently released de novo mutation table. For both SPARK and MSSNG, only DNMs for probands were used.

Pathway enrichment and comparison with case/control brain gene expression data

We used Reactome annotations¹³, and unless otherwise noted, PantherDB¹⁴ to assess functional enrichment in both forecASD genes and SFARI HC genes using Fisher's method. Odds ratios and p-values were used to compare these two prioritization methods (Fig. 4) in terms of the pathways they implicate. The full list of results of these enrichment analyses are provided in Supplemental Table 2. Statistical analyses described in results and discussion were all performed in R¹⁵ using either `glm()` or `fisher.test()`. Pathway-summarized haploinsufficiency (pLI: probability of loss-of-function intolerance¹⁶) was calculated by counting the proportion of genes in a Reactome pathway satisfying $pLI > 0.9$. Gene-wise and pathway-level comparisons with ASD case/control brain gene expression data were performed using frontal cortex RNA-seq summary statistics from Gandal *et al.*¹⁷. Our preliminary tests showed that both SFARI HC and forecASD showed the highest agreement with expression data from the frontal cortex.

Class and functional enrichment of top forecASD genes

Data used for functional enrichment in Figure 2D was taken from PubMed, STRING¹¹, and BrainSpan¹⁰, using forecASD genes as the subject. PubMed literature enrichment scores were calculated by summing total mentions of the gene list in abstracts also containing the word autism. The network interaction scores were derived using the STRING database, accessed via the STRINGdb package¹⁸ in R¹⁵. Using a score threshold of 0.4, we keep all STRING interactions between top forecASD genes. The total number of interactions above this threshold is then summed. Fetal brain coexpression scores are based on average Pearson correlation between top scoring forecASD genes in early developmental timepoints in the BrainSpan dataset. Given these three functional enrichment scores, average background values were permuted by randomly drawing a set of 1787 genes 1000 times. P-values and enrichment were computed relative to the permuted samples. Datasets used in the class enrichment in Figure 2D were taken from Sugathan *et al.*¹⁹, Darnel *et al.*²⁰, and Abrahams *et al.*⁴. P-values were computed by the hypergeometric statistical test of overlap between forecASD genes and these three gene sets.

Cluster analysis of top scoring forecASD genes

Using the STRING database, interactions were obtained for forecASD genes and loaded into a network using the igraph package²¹ in R¹⁵. No filter for interaction strength was enforced. Hierarchical greedy clustering based on optimization of the modularity score²² was performed using the fastgreedy.community function in the igraph package. Clustering was performed iteratively, with clusters containing more than 200 genes being subject to further clustering. Clusters with fewer than 30 genes were

discarded. The annotated network of clusters was loaded into Cytoscape²³, using the STRING application¹⁸. Functional enrichment of clusters was assessed using the STRING application in Cytoscape, with the p-value threshold set to 0.05. For annotation of the network plot shown in figure 6A, either the top annotated term or commonality between several top terms was chosen as representative. The p-value of enrichment with pLI scores was performed using Fisher's exact test of genes within each cluster with a pLI score above 0.5. The p-value for overlap with SFARI HC genes was performed using the hypergeometric test, assuming a background of 18,000 total genes.

Results

forecASD model and performance

The goal of our approach was to create a gene-wise score that indexes the level of evidence for involvement in ASD using both systems biology (i.e., network and transcriptional data) and genetic features. An initial forecASD systems biology model was built (forecASD:sys) using only BrainSpan expression and the STRING database shortest paths matrices as features. This model was trained on the high confidence set of 76 SFARI genes scoring 1 or 2 (SFARI HC genes), with negative training labels assigned to 1,000 background genes that were not listed in the SFARI gene database.

As an initial test of performance, we scored genes hit by coding *de novo* mutations (DNMs) in the recently published MSSNG study. As shown in Figure 2A, there is a

significant enrichment of likely gene disrupting (LGD) DNMs in the 90th percentile of both the TADA p-value (OR = 3.76, P = 5.45 x 10⁻⁹) and the forecASD:sys scores (OR = 3.15, P = 7.33 x 10⁻⁸). However, by far the greatest enrichment (OR = 12.81, P < 2.2 x 10⁻¹⁶) is seen when restricting to DNMs passing both a TADA q-value and forecASD:sys 90th percentile threshold.

To leverage both the genetic signal and the systems biology signal, we next built the final forecASD model, which incorporates forecASD:sys, the Krishnan *et al.* score⁶, DAMAGES⁷, DAWN⁵, and several TADA genetic scores from two recent studies^{3,12}. After training the forecASD model, we visualized the variable importance in figure 2B by mean decrease in the Gini impurity measure. The most informative feature was the STRING score from the forecASD:sys model, followed closely by two TADA score variables.

To facilitate a comparison with manually curated gene prioritizations, we scored all genes in the SFARI gene database using forecASD, forecASD:sys, and the most comprehensive TADA feature in the forecASD model. Shown in figure 2C, the forecASD model ranks SFARI genes scoring 3, 4, 5 and syndromic-only as significantly more autism-related than TADA (P: 7.7x10⁻⁴, 4.7x10⁻¹¹, 2.3x10⁻⁴, 7.7x10⁻⁶). The forecASD model also significantly outperforms the limited forecASD:sys model in gene categories 2, 3, and 4 (P: 8.4x10⁻⁵, 2.15x10⁻⁷, 4.0x10⁻⁵, respectively). In all cases, forecASD prioritizes SFARI genes as well, or better than TADA and forecASD:sys.

As an initial validation of genes prioritized by forecASD, we tested for an enrichment of gene sets and characteristics well known to be overrepresented in autism genes (Fig. 2D). We first performed several overrepresentation tests and found that genes receiving forecASD scores in the top decile (1,787 genes, referred to as forecASD genes) had a significant overlap with known targets of CHD8 ($P < 1 \times 10^{-16}$), FMRP ($P < 1 \times 10^{-16}$), and the full SFARI gene database ($P < 1 \times 10^{-16}$). We next performed a series of functional enrichment tests, comparing forecASD genes to randomly sampled sets of background genes. Text mining in PubMed showed that forecASD genes were significantly overrepresented in abstracts which mention autism ($P < 0.001$). Given the established role of autism genes early in fetal development, we next tested and found that forecASD genes showed significantly higher rates of coexpression across all regions of the fetal brain ($P < 0.001$). Lastly, forecASD genes were shown to have significantly enriched rates of interaction in the STRING database ($P < 0.001$).

We next tested the ability of these scores to discriminate both high confidence (Fig. 3A) and trending (Fig. 3B) autism genes from negative background genes. High confidence autism genes (SFARI HC) are defined as scoring 1 or 2 in SFARI Gene, with trending autism genes scoring 3. Importantly, the negative set of non-autism genes was sampled to have the same background mutation rate as the autism genes ($P > 0.1$ by the Kolmogorov-Smirnov test). In both comparisons, forecASD showed the highest level of performance of all methods tested (AUC=0.97 for SFARI 1+2 and AUC=0.82 for SFARI Gene score 3; Fig. 3). Furthermore, while the SFARI HC genes were used to train the forecASD model, only “out of bag” predictions were used as the forecASD score for

those genes, i.e., only those trees where the gene was not included in the bootstrap sample voted for the class of the gene. None of the trending autism genes (Fig. 3B) were used to train forecASD, and consequently they provide an unbiased estimate of performance.

Generalization to new data: *de novo* mutation enrichment

To compare forecASD and prior methods' ability to generalize to new data, we combined two recently released autism genetics resources. Specifically, we used *de novo* mutations in gene regions from the SPARK⁹ and MSSNG⁸ cohorts. Importantly, none of our model training used information from these studies, thus any subsequent validation is unbiased.

We first compared forecASD and competing ASD gene scores with respect to enrichment of genes with recurrent *de novo* loss of function and damaging missense mutations in probands. forecASD significantly outperformed all prior approaches (OR=26.8, $P=3.1 \times 10^{-24}$; Fig. 3C). We next tested whether forecASD continued to show significant enrichment when known autism genes (here, any gene listed in the SFARI gene database, regardless of score) were excluded (Fig. 3D), since the ideal method should detect both known and potentially novel autism genes. forecASD had superior performance in this test as well (OR=6.7, $P=0.0004$), with most of the other external methods lacking a statistically significant enrichment over baseline.

Functional enrichment and clustering of forecASD genes

Having demonstrated the predictive performance characteristics of forecASD, we next turned to practical applications that could further illuminate the underlying biological mechanisms at play in autism. Functional enrichment using Reactome annotations showed that forecASD genes are highly enriched for pathways known to play an important role in autism etiology, including chromatin modification, synaptic transmission, and developmental biology (full list in Supplemental Table 2). To highlight new biological themes that forecASD detects but that are not clear from the list of SFARI HC genes, we prioritized pathways based on differential enrichment (Fig. 4). Figure 4A highlights pathways that were represented in SFARI HC genes, but that showed significantly greater enrichment in forecASD genes. Figure 4B shows a sampling of the most significant forecASD pathways not represented by any SFARI HC gene, thus highlighting under-appreciated mechanisms in autism.

While SFARI HC genes show a strong bias toward genes with high pLI ($P < 0.001$, Fisher's exact test; Fig. 5A), forecASD is significantly less biased ($P < 0.001$, Fisher's exact test). We also discovered a significant relationship between pLI and differential expression (DE) t-statistics in case/control brain gene expression studies¹⁷ ($\beta = -0.13$, $t\text{-statistic} = -4.3$, $P = 1.9 \times 10^{-5}$, Fig. 5B), potentially exposing a form of bias in current gene discovery approaches that leads to under-ascertainment of ASD risk genes with low pLI and upregulation in ASD cases. We also found a significant interaction between forecASD and pLI ($F = 54.1$, $P = 3.9 \times 10^{-24}$) such that the pLI-expression relationship exists among forecASD genes ($\beta = -0.24$, $t = -2.6$, $P = 0.009$; Fig. 5D) but is absent in non-forecASD genes ($\beta = 0.004$, $t = 0.1$, $P = 0.91$; Fig. 5C).

Lastly, forecASD genes were loaded into the STRING network and clustered using a greedy hierarchical approach which maximizes the modularity score. The resulting networks consisted of 17 clusters composed of 1452 genes. All clusters were found to be significantly enriched with numerous GO and KEGG pathways. Similarly, all clusters contained a significant enrichment of haploinsufficiency genes ($pLI > 0.5$), except for the small cluster of 31 genes with functions related to the mediator complex. Clusters were also tested for overlap with SFARI HC genes, of which 8 clusters failed to reach significance, suggesting groupings of genes currently missing from the known list of autism genes. Clusters lacking significant overlap includes those with functions: signal transduction, cytoskeleton, cell migration, neuron projection, steroid signaling, neuron differentiation, potassium signaling, development and morphogenesis. A marginally significant correlation was seen between a clusters enrichment for high pLI genes and its overlap with SFARI HC genes (Spearman's $r = 0.48$, $p\text{-value} = 0.053$), further suggesting a bias in SFARI HC genes towards haploinsufficiency status.

Discussion

We present forecASD, a machine learning approach that combines systems biology and genetic models into a single score that indexes the strength of evidence for a gene's involvement in autism. This genome-wide score can be a useful prior, filter, or positive control in molecular studies of autism. It can also be used as a starting point to generate new hypotheses to investigate currently under-appreciated aspects of the molecular

etiology of autism. In our tests of predictive performance and generalization, forecASD outperformed other systems biology and genetic approaches for autism gene prioritization.

Because it draws upon multiple approaches for identifying autism genes, forecASD is less biased than gene discovery based only on one form of data (e.g., genetic data). This is particularly important because current SFARI HC genes, which rely heavily on studies of *de novo* mutation, are strongly biased towards genes that are loss-of-function intolerant (Fig. 5A). While these haploinsufficient genes represent a sizable and important component of genetic risk for autism, this ascertainment bias has led to molecular “blind spots” that will not be resolved simply by sequencing more probands. For instance, pathways implicated preferentially by SFARI HC genes had significantly higher pLI, whereas pathways with lower pLI were under-represented (compared to forecASD-implicated pathways; OR=0.38, $P=5.6 \times 10^{-7}$). Furthermore, while SFARI HC pathways significantly predicted case/control expression-implicated pathways ($Z=4.5$, $P=7.9 \times 10^{-6}$, binomial model), only 3% of the deviance could be explained. In contrast, forecASD pathways explained an order of magnitude more deviance (31%, $Z=12.1$, $P=1.5 \times 10^{-33}$) when predicting expression-implicated pathways. When a model of dichotomous DE significance was fit that included terms from both SFARI HC and forecASD pathways, the SFARI HC term became redundant, and the forecASD-only model yielded a superior Bayesian information criterion (BIC; 518 for forecASD-only vs. 521 for full model and 720 for the SFARI HC-only model). When considering directionality, SFARI HC gene pathways were significantly depleted for ASD-upregulated pathways (OR=0.48, $P=1.3 \times 10^{-6}$), further illustrating the bias in SFARI HC

genes. These results demonstrate that forecASD showed greater representation of low pLI and ASD-upregulated pathways, without sacrificing sensitivity to well-known ASD risk pathways where haploinsufficiency plays a dominant role.

Notably, pLI is significantly negatively correlated with previously published frontal cortex differential expression t-statistics¹⁷ ($\beta = -0.13$, $t = -4.3$, $P = 1.9 \times 10^{-5}$), suggesting that low-pLI genes are more likely to be up-regulated and high-pLI genes are more likely to be down-regulated in ASD cases (Fig. 5B). This is consistent with our observation of SFARI HC gene pathways (which have an ascertainment bias in favor of haploinsufficiency) being significantly under-represented in both low pLI and ASD-upregulated pathways. We further observed a significant interaction ($F = 54.1$, $P = 3.9 \times 10^{-24}$) between forecASD and pLI when explaining variation in ASD brain gene expression: forecASD genes show the significant negative relationship between pLI and t-statistic ($\beta = -0.24$, $t = -2.6$, $P = 0.009$; Fig. 5D), while non-forecASD genes show no relationship ($\beta = 0.004$, $t = 0.1$, $P = 0.91$; Fig. 5C). Consequently, we propose that this pLI-expression relationship is a hallmark of robust ASD risk genes, and may be used as a criterion when identifying optimal thresholds in genome-wide scores like forecASD.

Indeed, although initially chosen as a convenient but arbitrary threshold for identifying a discrete set of ASD candidate genes, the top decile proved to be the optimal split point for forecASD, maximizing the significance of the pLI/t-statistic relationship among candidate genes, while minimizing the same relationship in the remaining, non-candidate genes. Interestingly, when applying this approach to TADA FDR values, although TADA-implicated genes showed the expected pLI-expression relationship, no TADA threshold was able to eliminate the trend from non-candidate genes, suggesting

lower sensitivity in identifying ASD risk genes compared to forecASD. Taken together, these analyses demonstrate that the reduced bias in forecASD contributes to increased sensitivity to autism risk pathways identified in gene expression studies (Fig. 4C,4D) as well as those implicated by genetic studies (Fig. 3).

Some pathways, although represented (but not necessarily enriched) in the current SFARI HC list, showed a substantial relative increase in enrichment when considering forecASD (Fig. 4A, Supplemental Table 2), suggesting that these pathways represent noted and plausible but still under-appreciated molecular themes in our understanding of autism. The pathway that underwent the largest relative increase in enrichment from SFARI HC to forecASD is Rho GTPase signaling ($OR=2.2$, $P=4.8 \times 10^{-5}$), which plays a critical role in cytoskeletal dynamics in neurodevelopment²⁴, including interactions with SHANK proteins and the formation and maturation of dendritic spines²⁵. As another example, although chromatin modification in general is a well-established theme in autism genetic risk, histone acetyltransferases showed relatively little representation in the SFARI HC list, but were significantly enriched in forecASD genes ($OR=4.1$, $P=3.5 \times 10^{-9}$). Histone acetylation was recently shown to be a pervasive genomic predictor of affected status in a large autism case/control postmortem brain study²⁶, underscoring the importance of this mechanism that is under-represented in established risk genes but that forecASD was sensitive to. As a final example of these under-appreciated molecular mechanisms, the circadian clock pathway was implicated by forecASD as an important source of risk for autism ($OR=6.5$, $P=7.6 \times 10^{-13}$). Sleep disturbances are a well-known and problematic comorbidity in autism, and molecular deficits in circadian regulation related to autism have been documented^{27,28,29}. Although

literature support is available for these processes playing a role in autism, our results indicate that their current sparse representation in lists of accepted genetic risk factors is not representative of their importance in the disorder.

Other pathways were identified by forecASD as significantly enriched for autism risk, but were not represented at all among SFARI HC genes (Supplemental Table 3, Figure 4B). Consequently, new insights into the molecular basis of autism will come disproportionately from these pathways as their constituent genes are associated with autism. One gene set in particular, potassium channels, showed highly significant enrichment in forecASD genes ($OR=4.1$, $P=7.2 \times 10^{-9}$, $N=35$ genes) despite the absence of potassium channel genes among currently accepted autism risk genes. However, the literature shows support for a role for potassium channels in ASD risk^{30,31,32,33}, and the pathway was enriched for differential regulation in a recently published brain gene expression study of autism ($P=0.001$, downregulated)¹⁷. Notably, this pathway has a lower proportion of genes with $pLI>0.9$ (0.22) compared to SFARI HC gene-implicated pathways (median=0.47), potentially explaining its absence due to ascertainment bias. Overall, pathways that demonstrated forecASD-specific excess enrichment showed a significant agreement with pathway enrichment from case/control brain gene expression studies ($OR=28.8$; $P=2.9 \times 10^{-48}$), and were more likely to support pathways that were up-regulated in the gene expression data ($OR=2.1$, $P=1.27 \times 10^{-6}$, Fig. 4C) compared to pathways implicated by SFARI HC.

To group forecASD genes into distinct functional categories, we performed iterative clustering and identified a total of 17 clusters enriched for specific functional annotations. While nearly all clusters showed significant enrichment for

haploinsufficiency genes, many lacked a significant overlap with SFARI HC genes, after Bonferroni correction. Similar to conclusions reached above, we found an entire cluster enriched for Potassium signaling ($P=1.8 \times 10^{-49}$) which lacked significant overlap with SFARI HC genes. In addition to this cluster, there were also seven others lacking significant overlap with SFARI HC genes. Notable examples include clusters related to cell migration ($P=6.0 \times 10^{-11}$) and endocytosis ($P=2.7 \times 10^{-21}$). These pathways have more recently been explored in their ability to regulate brain connectivity³⁴ and postsynaptic organization³⁵, respectively. In agreement with the proposed haploinsufficiency bias of autism gene discovery, we observed a marginally significant relationship between cluster pLI enrichment and SFARI HC gene overlap (Spearman's rho: 0.48; $P=0.053$).

For the foreseeable future, traditional gene discovery studies will continue to add to the list of bona fide ASD risk genes. Eventually, as sample sizes saturate and gene discovery decelerates, the field will be faced with the challenge of developing new and useful applications of this acquired knowledge. By providing a glimpse of that future, forecASD gives an opportunity right now to begin thinking about what we would do with a definitive list of autism genes.

448

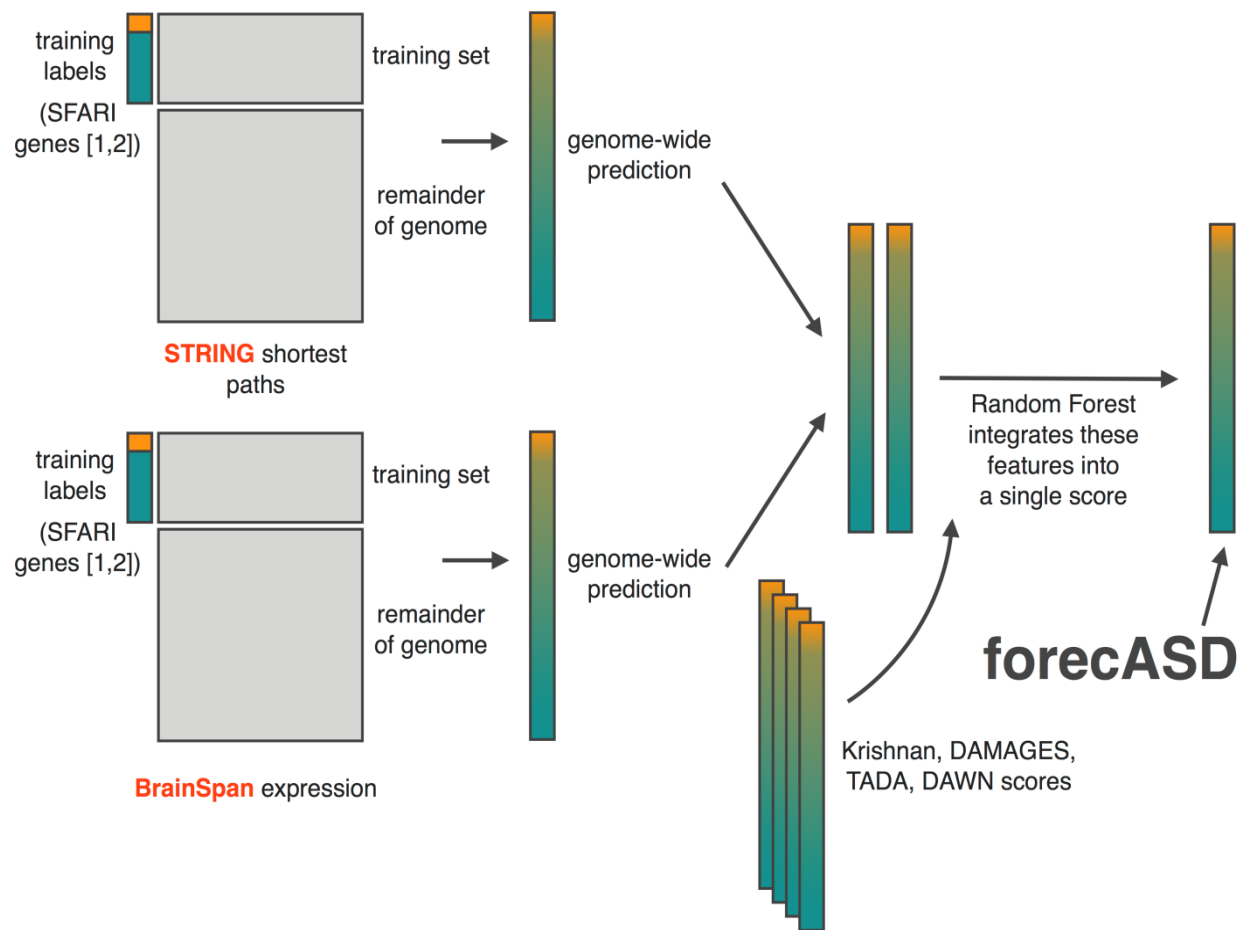


Figure 1 - Overview of forecASD. Two Random Forest classifiers, one using BrainSpan gene expression and the other using the STRING network as predictors, are trained to discriminate high confidence autism genes (SFARI HC, scores 1 and 2) from a set of 1,000 genes drawn randomly from those not listed at all in the SFARI Gene database. Predictions are then made on the remainder of the genome, and these are combined with the out-of-bag (OOB) estimates from the training process to yield a prediction for each gene in the genome. A subsequent classifier is then trained using the output of these two RFs and previously published autism gene scores as predictive features, and again predictions are made on the remainder of the genome, with OOB predictions being used for those genes in the training set. The RF vote proportion for class “autism gene” is then the final forecASD score.

449

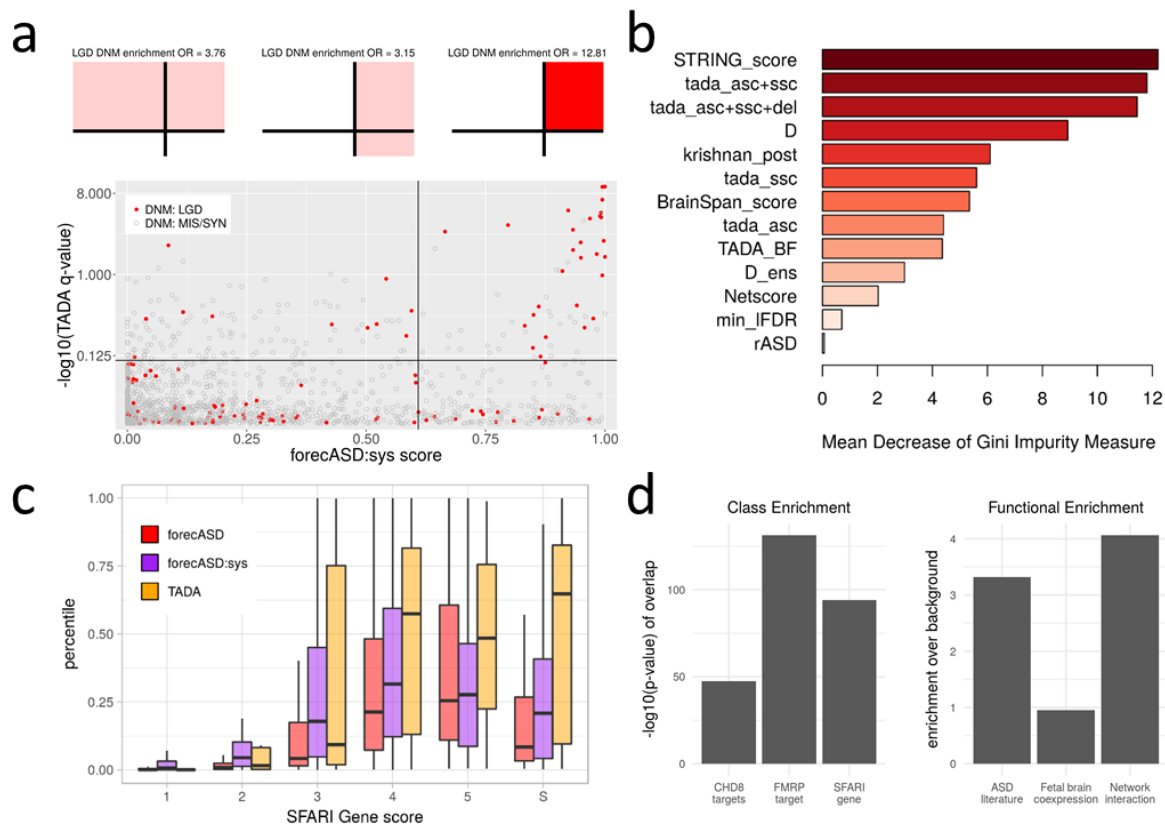


Figure 2 - Prioritization of de novo likely gene-disrupting mutations and enrichment of gene sets in forecASD. Training a limited model, forecASD:sys, using brain gene expression and interaction data shows optimal prioritization of de novo LGDs when combined with a genetic measure of autism association (a). Building the full forecASD model, we test all features for their informativeness, finding that the STRING score is primary (b). Using the three mentioned scores, we assess their genome-wide ranking of SFARI genes at all levels, and find that the full forecASD model at least ties, and often significantly outperforms TADA and forecASD:sys in the prioritization of SFARI genes (c). As an initial assessment of forecASD prioritized genes, we find the top decile of genes ranked by forecASD (1787 genes) shows enrichment typical of classical autism genes (d).

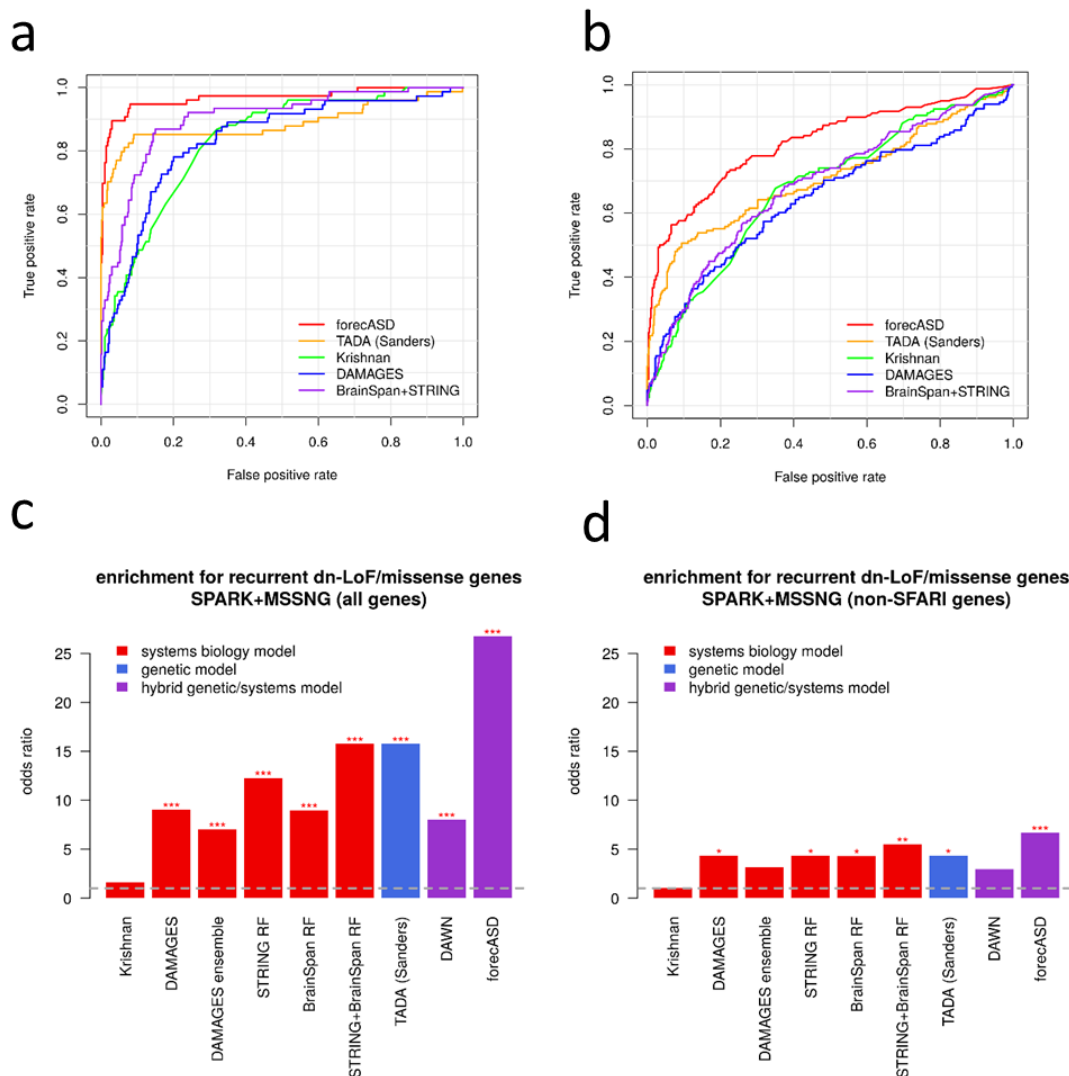


Figure 3 - Comparison of forecASD with prior models of autism gene prioritization. To compare forecASD with competitors, we evaluate performance by each methods' ability to prioritize SFARI genes and genes which were subject to recurrent de novo loss-of-function or missense mutations. Starting with SFARI genes scoring 1 or 2 as a positive set and size-matched random background genes as the negative set, forecASD out-of-bag estimates showed superior classification over all methods (a). In a fully unbiased test, forecASD estimates also showed superior classification of trending SFARI genes (score: 3) over all other methods (b). Using two sequencing cohorts which no methods draw information from, the top decile of forecASD genes (1787 genes) shows the greatest overlap with genes containing recurrent de novo loss-of-function and missense mutations (c). When excluding genes in the SFARI gene database, forecASD still shows superior prioritization of genes accumulating de novo mutations (d).

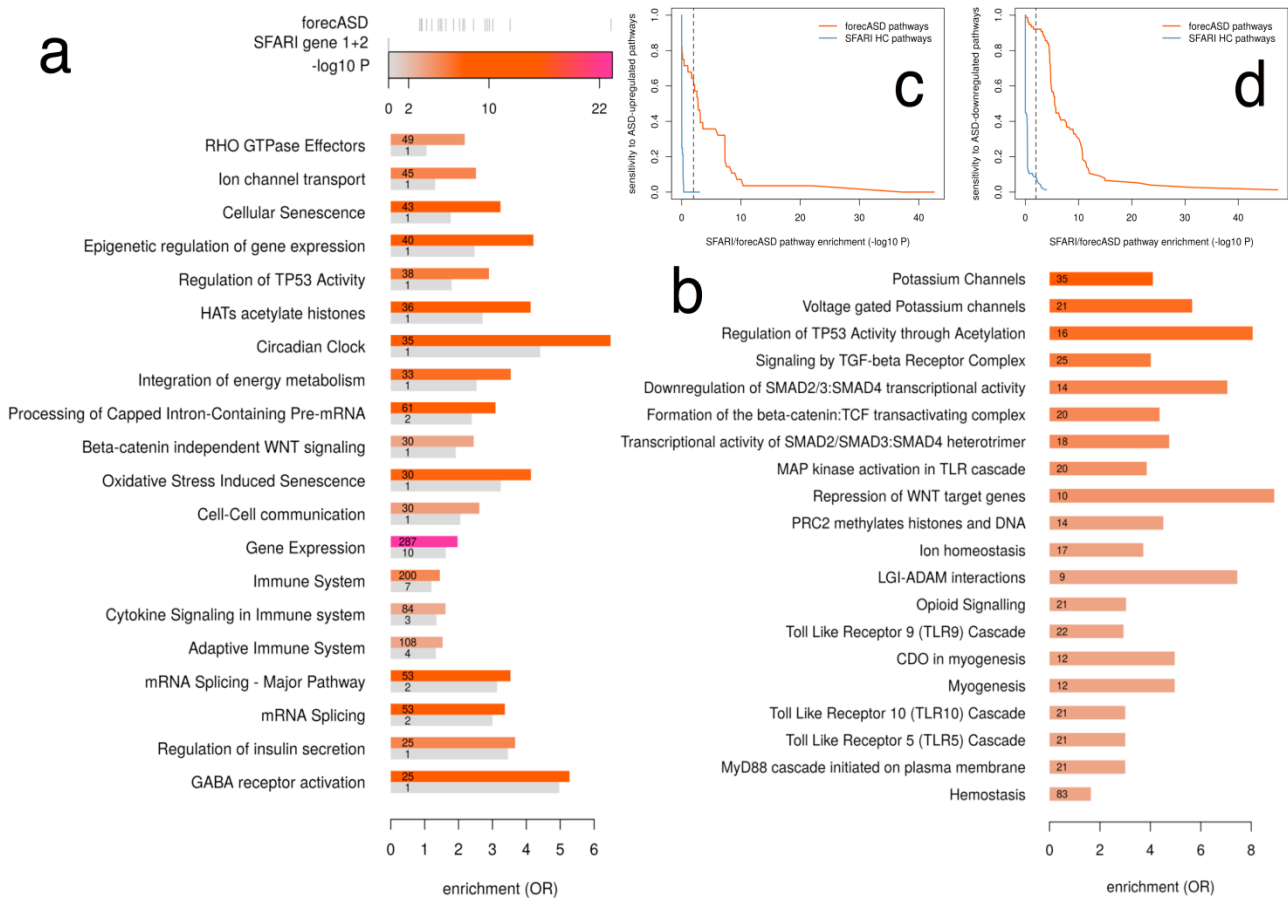


Figure 4 - forecASD-specific pathway enrichment and sensitivity to gene expression-implicated pathways. When testing the top-decile genes according to forecASD for Reactome pathway enrichment, pathways emerged that were represented, but not enriched in the SFARI HC list (a). Other pathways were highly enriched in forecASD genes that were not represented at all in the SFARI HC list, even though they have associated literature suggesting a role in autism (b). forecASD is more sensitive than SFARI HC to pathways that are differentially regulated in the brains of individuals with autism, particularly in ASD-upregulated pathways (c), but also in downregulated pathways (d). Using the top decile of TADA -log10 FDR genes showed similar sensitivity to SFARI HC (not shown), suggesting that rare variant approaches may be less sensitive in implicating genes found through gene expression studies.

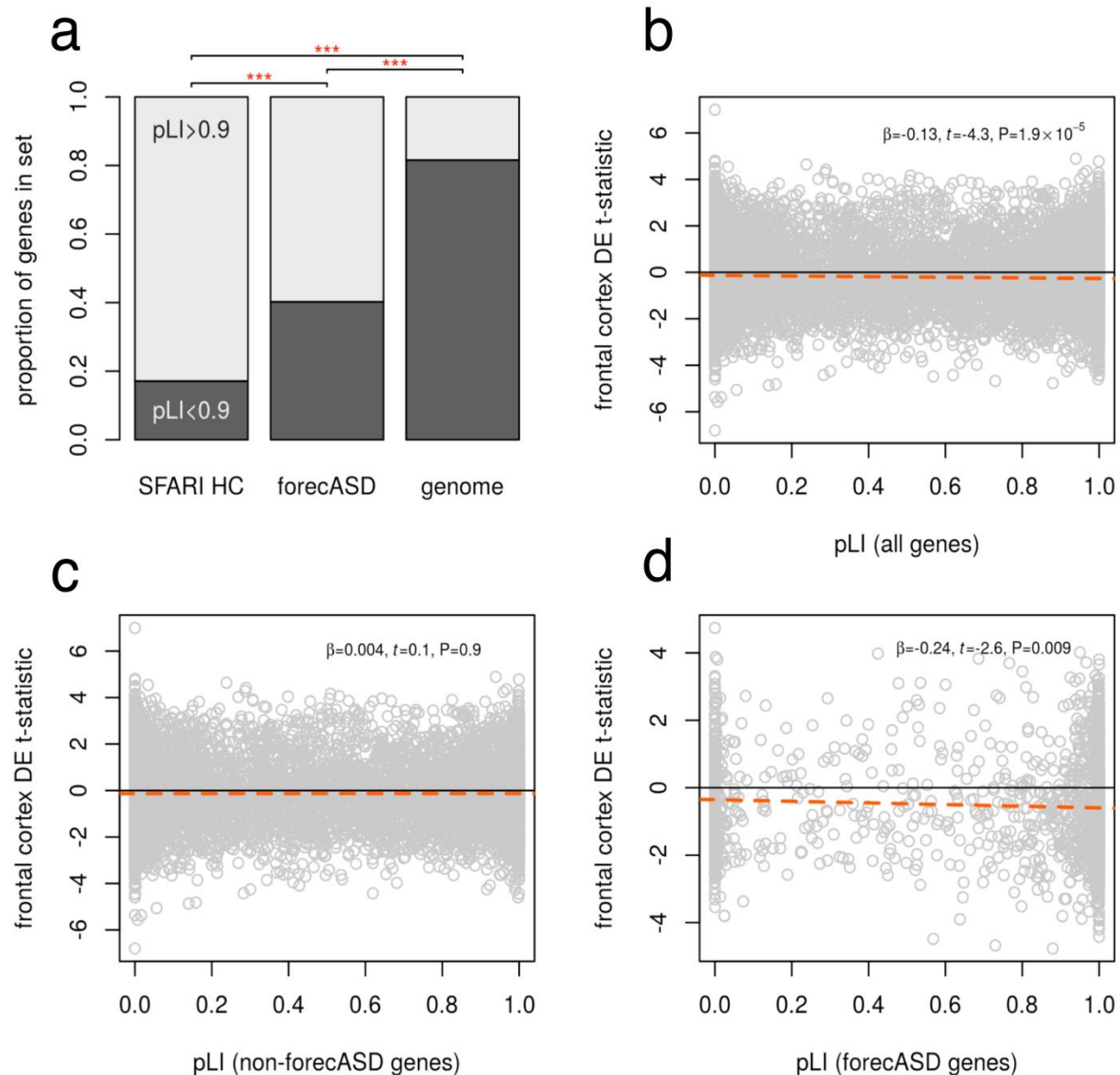


Figure 5 - Relationship between pLI and ASD-specific up- and down-regulation of brain gene expression. SFARI HC is strongly biased toward genes with high pLI (a), while forecASD is significantly less biased. We found a significant relationship between pLI and differential expression (DE) in the brains of autism cases (b), such that low pLI genes tend toward upregulation in cases, while high pLI genes tend to be downregulated. We also observed a significant interaction between forecASD and pLI such that the observed pLI-DE trend (b) is absent in non-forecASD genes (c), and present and significant among forecASD genes (d). We propose that the presence of the pLI-DE trend is a hallmark of ASD risk genes, and an optimal ASD gene prioritization method will concentrate the trend among risk genes and remove it from non-risk genes. Notably, no threshold of TADA (tested to the 50th percentile) was able to remove the trend from the non-prioritized genes, suggesting the persistence of residual risk genes that were not selected.

459

460

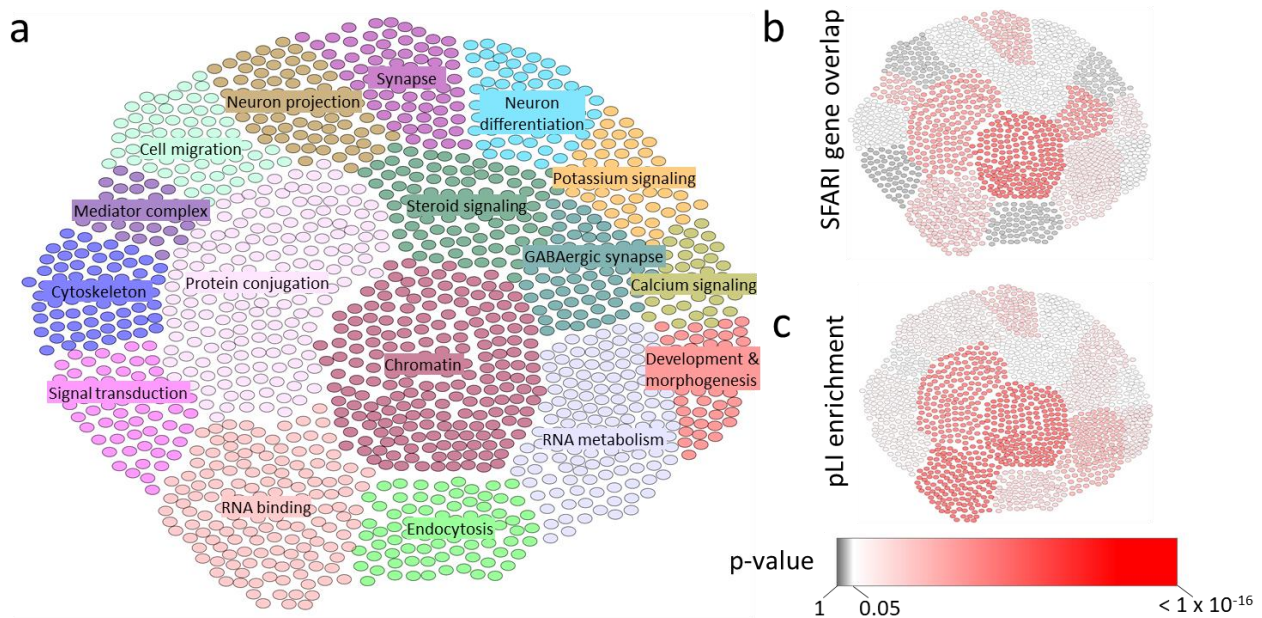


Figure 6 - Clustering of top foreASD genes with enrichment analysis of ExAC pLI and SFARI high-confidence gene overlap. Greedy hierarchical optimization of the modularity score yielded 17 clusters consisting of 1452 foreASD genes (a). All clusters have several significantly enriched biological pathways, of which the top terms were overlaid in figure 6A. Clusters were tested for significance of overlap with the list of SFARI HC genes (b), and enrichment of haploinsufficiency genes ($pLI > 0.5$; c).

463 1 Rosenberg, R. E. *et al.* Characteristics and concordance of autism spectrum disorders among 277
464 twin pairs. *Arch Pediatr Adolesc Med* **163**, 907-914, doi:10.1001/archpediatrics.2009.98 (2009).

465 2 Colvert, E. *et al.* Heritability of Autism Spectrum Disorder in a UK Population-Based Twin Sample.
466 *JAMA Psychiatry* **72**, 415-423, doi:10.1001/jamapsychiatry.2014.3028 (2015).

467 3 De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in autism. *Nature*
468 **515**, 209-215, doi:10.1038/nature13772 (2014).

469 4 Abrahams, B. S. *et al.* SFARI Gene 2.0: a community-driven knowledgebase for the autism
470 spectrum disorders (ASDs). *Mol Autism* **4**, 36, doi:10.1186/2040-2392-4-36 (2013).

471 5 Liu, L. *et al.* DAWN: a framework to identify autism genes and subnetworks using gene
472 expression and genetics. *Mol Autism* **5**, 22, doi:10.1186/2040-2392-5-22 (2014).

473 6 Krishnan, A. *et al.* Genome-wide prediction and functional characterization of the genetic basis
474 of autism spectrum disorder. *Nat Neurosci* **19**, 1454-1462, doi:10.1038/nn.4353 (2016).

475 7 Zhang, C. & Shen, Y. A Cell Type-Specific Expression Signature Predicts Haploinsufficient Autism-
476 Susceptibility Genes. *Hum Mutat* **38**, 204-215, doi:10.1002/humu.23147 (2017).

477 8 RK, C. Y. *et al.* Whole genome sequencing resource identifies 18 new candidate genes for autism
478 spectrum disorder. *Nat Neurosci* **20**, 602-611, doi:10.1038/nn.4524 (2017).

479 9 pfeliciano@simonsfoundation.org, S. C. E. a. & Consortium, S. SPARK: A US Cohort of 50,000
480 Families to Accelerate Autism Research. *Neuron* **97**, 488-493, doi:10.1016/j.neuron.2018.01.015
481 (2018).

482 10 Sunkin, S. M. *et al.* Allen Brain Atlas: an integrated spatio-temporal portal for exploring the
483 central nervous system. *Nucleic Acids Res* **41**, D996-D1008, doi:10.1093/nar/gks1042 (2013).

484 11 von Mering, C. *et al.* STRING: a database of predicted functional associations between proteins.
485 *Nucleic Acids Res* **31**, 258-261 (2003).

486 12 Sanders, S. J. *et al.* Insights into Autism Spectrum Disorder Genomic Architecture and Biology
487 from 71 Risk Loci. *Neuron* **87**, 1215-1233, doi:10.1016/j.neuron.2015.09.016 (2015).

488 13 Fabregat, A. *et al.* The Reactome Pathway Knowledgebase. *Nucleic Acids Res* **46**, D649-D655,
489 doi:10.1093/nar/gkx1132 (2018).

490 14 Mi, H. *et al.* The PANTHER database of protein families, subfamilies, functions and pathways.
491 *Nucleic Acids Res* **33**, D284-288, doi:10.1093/nar/gki078 (2005).

492 15 R: A language and environment for statistical computing (R Foundation for Statistical
493 Computing, Vienna, Austria, 2008).

494 16 Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285-
495 291, doi:10.1038/nature19057 (2016).

496 17 Gandal, M. J. *et al.* Shared molecular neuropathology across major psychiatric disorders parallels
497 polygenic overlap. *Science* **359**, 693-697, doi:10.1126/science.aad6469 (2018).

498 18 Szklarczyk, D. *et al.* The STRING database in 2017: quality-controlled protein-protein association
499 networks, made broadly accessible. *Nucleic Acids Res* **45**, D362-D368, doi:10.1093/nar/gkw937
500 (2017).

501 19 Sugathan, A. *et al.* CHD8 regulates neurodevelopmental pathways associated with autism
502 spectrum disorder in neural progenitors. *Proc Natl Acad Sci U S A* **111**, E4468-4477,
503 doi:10.1073/pnas.1405266111 (2014).

504 20 Darnell, J. C. *et al.* FMRP stalls ribosomal translocation on mRNAs linked to synaptic function and
505 autism. *Cell* **146**, 247-261, doi:10.1016/j.cell.2011.06.013 (2011).

506 21 The igraph software package for complex network research (InterJournal, 2006).

507 22 Newman, M. E. Modularity and community structure in networks. *Proc Natl Acad Sci U S A* **103**,
508 8577-8582, doi:10.1073/pnas.0601602103 (2006).

509 23 Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular
510 interaction networks. *Genome Res* **13**, 2498-2504, doi:10.1101/gr.1239303 (2003).

511 24 Reichova, A., Zatkova, M., Bacova, Z. & Bakos, J. Abnormalities in interactions of Rho GTPases
512 with scaffolding proteins contribute to neurodevelopmental disorders. *J Neurosci Res* **96**, 781-
513 788, doi:10.1002/jnr.24200 (2018).

514 25 Martin-Vilchez, S. *et al.* RhoGTPase Regulators Orchestrate Distinct Stages of Synaptic
515 Development. *PLoS One* **12**, e0170464, doi:10.1371/journal.pone.0170464 (2017).

516 26 Sun, W. *et al.* Histone Acetylome-wide Association Study of Autism Spectrum Disorder. *Cell* **167**,
517 1385-1397 e1311, doi:10.1016/j.cell.2016.10.031 (2016).

518 27 Lipton, J. O. *et al.* Aberrant Proteostasis of BMAL1 Underlies Circadian Abnormalities in a
519 Paradigmatic mTOR-opathy. *Cell Rep* **20**, 868-880, doi:10.1016/j.celrep.2017.07.008 (2017).

520 28 Monyak, R. E. *et al.* Insulin signaling misregulation underlies circadian and cognitive deficits in a
521 *Drosophila* fragile X model. *Mol Psychiatry* **22**, 1140-1148, doi:10.1038/mp.2016.51 (2017).

522 29 Kozlov, S. V. *et al.* The imprinted gene Magel2 regulates normal circadian output. *Nat Genet* **39**,
523 1266-1272, doi:10.1038/ng2114 (2007).

524 30 Guglielmi, L. *et al.* Update on the implication of potassium channels in autism: K(+)
525 channelautism spectrum disorder. *Front Cell Neurosci* **9**, 34, doi:10.3389/fncel.2015.00034
526 (2015).

527 31 Deng, P. Y. & Klyachko, V. A. Genetic upregulation of BK channel activity normalizes multiple
528 synaptic and circuit defects in a mouse model of fragile X syndrome. *J Physiol* **594**, 83-97,
529 doi:10.1113/JP271031 (2016).

530 32 Lee, H., Lin, M. C., Kornblum, H. I., Papazian, D. M. & Nelson, S. F. Exome sequencing identifies
531 de novo gain of function missense mutation in KCND2 in identical twins with autism and seizures

532 that slows potassium channel inactivation. *Hum Mol Genet* **23**, 3481-3489,
533 doi:10.1093/hmg/ddu056 (2014).

534 33 Sicca, F. *et al.* Gain-of-function defects of astrocytic Kir4.1 channels in children with autism
535 spectrum disorders and epilepsy. *Sci Rep* **6**, 34325, doi:10.1038/srep34325 (2016).

536 34 Reiner, O., Karzbrun, E., Kshirsagar, A. & Kaibuchi, K. Regulation of neuronal migration, an
537 emerging topic in autism spectrum disorders. *J Neurochem* **136**, 440-456, doi:10.1111/jnc.13403
538 (2016).

539 35 Loebrich, S. The role of F-actin in modulating Clathrin-mediated endocytosis: Lessons from
540 neurons in health and neuropsychiatric disorder. *Commun Integr Biol* **7**, e28740,
541 doi:10.4161/cib.28740 (2014).

542