

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

3

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

5

6 ¹Department of Psychiatry, Carver College of Medicine

7 ²Medical Scientist Training Program, Carver College of Medicine

8 ³Interdisciplinary Graduate Program in Genetics

9 ⁴Department of Biomedical Engineering, College of Engineering

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

11 Sciences

12 *Corresponding Author

13 University of Iowa

14 Iowa City, IA 52242

15 USA

16

17

18

19

20

21

22

23

24 **Abstract**

25 **Background**

26 Genes are one of the most powerful windows into the biology of autism, and it has been
27 estimated that perhaps a thousand or more genes may confer risk. However, less than
28 100 genes are currently viewed as having robust enough evidence to be considered
29 true "autism genes". Massive genetic studies are underway to produce data to
30 implicate additional genes, but this approach, although necessary, is costly and slow-
31 moving.

32 **Methods**

33 We approach autism gene discovery as a machine learning problem, rather than a
34 genetic association problem, and use genome-scale data as predictors for identifying
35 further genes that have similar properties in the feature space compared to established
36 autism risk genes. This approach, which we call forecASD, integrates spatiotemporal
37 gene expression, heterogeneous network data, and previous gene-level predictors of
38 autism association into an ensemble classifier that yields a single score that indexes
39 each gene's evidence for being involved in the etiology of autism.

40 **Results**

41 We demonstrate that forecASD has substantially increased sensitivity and specificity
42 compared to previous gene-level predictors of autism association, including genetic
43 measures such as TADA. On an independent test set, consisting of newly-released
44 pilot data from the SPARK Genomics Consortium, we show that forecASD best predicts
45 which genes will have an excess of likely gene disrupting (LGD) *de novo* mutations. We
46 further use independent data from a recent post mortem study of case/control gene

47 expression to show that forecASD is also a significant predictor of genes implicated in
48 ASD through differential expression. Using forecASD results, we show which molecular
49 pathways are currently under-represented in the autism literature and likely represent
50 under-appreciated biological mechanisms of autism. Finally, forecASD correctly
51 predicted 12 of 16 genes implicated at FDR=0.2 by the latest ASD gene discovery
52 study, while also identifying the most likely false positives among the candidate genes.

53 **Conclusions**

54 These results demonstrate that forecASD bridges the gap between genetic- and
55 expression-based ASD gene discovery, and provides a data-driven replacement to
56 much of the manual filtering and curation that is a critical step in ensuring the
57 robustness of gene discovery studies.

58

59 **Keywords:**

60 Autism, genetics, machine learning, disease-gene discovery

61

62 **Background**

63

64 Autism Spectrum Disorder (ASD) is a heterogeneous grouping of developmental
65 disorders caused by a range of genetic and environmental factors. The core diagnostic
66 features of ASD, which manifest at a young age, are impairments in social
67 communication and restrictive and repetitive behaviors and interests. Evidence for the
68 role of genetics in ASD is strong, with monozygotic twins having near 90% concordance
69 of ASD diagnosis(1). Further population and twin studies have confirmed these
70 findings(2), and further estimated the narrow-sense heritability of ASD to be in the range
71 of 50-95%.

72

73 While there is an abundance of evidence for the role of genetics in autism, our
74 understanding of the genetic etiology of the disorder is still limited. It is estimated that
75 there may over 1000 genes which contribute to autism risk(3). However, the current list
76 of high-confidence autism genes stands at 84 genes(4). This discrepancy is partly
77 explained by the relatively limited number of genomic studies compared with the vast
78 genetic heterogeneity underlying autism.

79

80 To close this gap between the number of anticipated and known autism genes, several
81 network-biology approaches have been applied in the past decade. These studies
82 leverage large, publicly-available datasets to add context and amplify the genetic
83 signals observed through sequencing studies. These network-biology studies have
84 predicted genes that then became bona fide autism genes(5), but have fallen short of

85 providing a useful genome-wide metric that indicates the evidence of autism
86 involvement for every gene. More recently, machine learning based methods have used
87 gene interaction networks(6) and cell-specific expression profiles(7) to predict gene
88 involvement in autism. Importantly, the results of these studies lead to a quantitative
89 metric that scores every gene in the genome according to evidence of a role in autism.
90 Despite the demonstrated effectiveness of these studies in prioritizing autism risk
91 genes, our preliminary investigations suggested there was still room for appreciable
92 improvement in the form of the classification algorithm, the training set, and the
93 predictors used. In particular, these approaches do not incorporate indicators of autism
94 involvement that are based on genetic association (e.g., TADA scores) into their
95 predictive features.

96
97 We introduce a new score, forecASD, that integrates prior network-biology approaches,
98 scores of genetic association, brain gene expression, and topological information from
99 large gene interaction networks relevant to the brain into a single gene-level score for
100 autism involvement. We show that forecASD successfully outperforms existing methods
101 in a diverse range of gene and mutation prioritization tasks. Further, using the recent
102 sequencing studies MSSNG(8) and SPARK(9), we show that forecASD generalizes to
103 previously unseen data. Importantly, this generalization holds even when excluding
104 genes with known links to autism, emphasizing forecASD's ability to identify novel ASD
105 genes. We also demonstrate that forecASD correctly predicts 12 of 16 genes implicated
106 at FDR=0.2 by the latest ASD gene discovery study, while also identifying the most
107 likely false positives among the candidate genes. Through comparing the top decile of

108 forecASD identified genes (1787 genes; hereafter forecASD genes) with known autism
109 genes, we identify numerous biological pathways that are currently underrepresented in
110 our understanding of autism risk. By reanalyzing the results of autism brain differential
111 gene expression studies, we show that the current list of known autism genes is
112 significantly depleted for upregulated biological pathways, whereas forecASD captures
113 both up- and downregulated pathways. We show that direction of differential expression
114 is related to haploinsufficiency status, with low pLI genes showing a trend towards
115 upregulation. Importantly, this relationship between direction of differential expression
116 and pLI is dependent on forecASD gene inclusion, signifying forecASD's ability to
117 capture low- and high-pLI disease genes. Through these studies, we show evidence
118 that current methods of autism gene discovery have biases, and that forecASD
119 mitigates these biases through its integrative approach, thus providing a view of the full
120 spectrum of genes and biological pathways underlying autism.

121

122 **Methods**

123

124 Overview:

125 The forecASD method relies upon stacked Random Forest models, organized in two
126 levels (shown in Figure 1). In the first level, two models are trained using BrainSpan(10)
127 gene expression and the STRING(11) shortest paths network as features, respectively.
128 Our training dataset consists of high-confidence genes scored in SFARI gene(4) as
129 either 1 or 2 (SFARI HC genes), and 1,000 random background genes not contained
130 within SFARI gene. These two models produce genome-wide predictions for autism

131 involvement. These scores are then used as features in the second level's Random
132 Forest model, along with other genome-wide scores obtained from previous studies.

133

134 BrainSpan, STRING, and TADA data assembly

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

138

139 The STRING database(11) was thresholded at their recommended score of 0.4, and
140 transformed into a gene by gene matrix with each cell representing the shortest path
141 between two genes.

142

143 TADA summary statistics were downloaded from the largest meta-analysis for autism
144 available at the time of publication(12). TADA summary statistics were also obtained
145 from the secondary supplementary table of another comprehensive study of autism
146 risk(3). All available TADA summary statistics were used as features in the final model,
147 with tadaFdrAscSscExomeSscAgpSmallDel(12) used as a representative comparator in
148 the ROC curve displayed in Figure 3.

149

150 Model training and genome wide prediction

151 We used a stacked Random Forest classifier to generate genome-wide predictions of
152 autism gene involvement. All models were trained using SFARI HC genes as positive
153 examples (of which there are 76 common to both STRING & BrainSpan), and a

154 randomly sampled set of 1,000 background genes (i.e., not listed in the SFARI Gene
155 database) as negative examples.

156

157 The first level of our stacked model consists of two genome-wide scores based on data
158 from BrainSpan or the STRING interaction network. The features used in training these
159 two models include the loess-smoothed observations in the BrainSpan database, and
160 the STRING shortest path matrix, respectively. The random forest models were trained
161 with 1000 total trees constructed, and the strata option enabled to insure a balance of
162 70 positive and 70 negative training examples during the construction of each tree.

163 Given the large number of features for the STRING-based random forest model, we
164 performed feature selection, wherein each feature not used in any of the constructed
165 trees was dropped. This variable selection step was repeated until the final model
166 contained only features which were selected at least once during tree construction. With
167 the STRING and BrainSpan models, we then predicted autism involvement scores for
168 the remaining genes not included in our training set. These scores are in
169 Supplementary Table 1 in the columns BrainSpan_score and STRING_score. Scores
170 for training set genes are the out-of-bag estimates.

171

172 We used these scores, along with DAWN(5), TADA(12)⁽³⁾, DAMAGES(7), and the score
173 from Krishnan *et al.*(6) score, as predictive features in a final Random Forest, using the
174 same training labels described previously. Genome-wide predictions were then
175 obtained, again using out-of-bag estimates for training set genes. This final score is
176 listed under forecASD_score in Supplementary Table 1.

177

178 SPARK and MSSNG data sources

179 De novo mutation (DNM) data from the MSSNG dataset was obtained through the *de*
180 *novo db* database(6). Mutations were filtered for LGD or missense status. De novo
181 mutation data was obtained from the SPARK dataset from the consortium's recently
182 released de novo mutation table. For both SPARK and MSSNG, only DNMs for
183 probands were used.

184

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

186 We used Reactome annotations(13), and unless otherwise noted, PantherDB(14) to
187 assess functional enrichment in both forecASD genes and SFARI HC genes using
188 Fisher's method. Odds ratios and p-values were used to compare these two
189 prioritization methods (Fig. 4) in terms of the pathways they implicate. The full list of
190 results of these enrichment analyses are provided in Supplemental Table 2. Statistical
191 analyses described in results and discussion were all performed in R(15) using either
192 `glm()` or `fisher.test()`. Pathway-summarized haploinsufficiency (pLI: probability of loss-of-
193 function intolerance(16)) was calculated by counting the proportion of genes in a
194 Reactome pathway satisfying $pLI > 0.9$. Gene-wise and pathway-level comparisons with
195 ASD case/control brain gene expression data were performed using frontal cortex RNA-
196 seq summary statistics from Gandal *et al.*(17). Our preliminary tests showed that both
197 SFARI HC and forecASD showed the highest agreement with expression data from the
198 frontal cortex.

199

200 Class and functional enrichment of top forecASD genes

201 Data used for functional enrichment in Figure 2D was taken from PubMed, STRING(11),
202 and BrainSpan(10), using forecASD genes as the subject. PubMed literature
203 enrichment scores were calculated by summing total mentions of the gene list in
204 abstracts also containing the word autism. The network interaction scores were derived
205 using the STRING database, accessed via the STRINGdb package(18) in R(15). Using
206 a score threshold of 0.4, we keep all STRING interactions between top forecASD
207 genes. The total number of interactions above this threshold is then summed. Fetal
208 brain coexpression scores are based on average Pearson correlation between top
209 scoring forecASD genes in early developmental timepoints in the BrainSpan dataset.
210 Given these three functional enrichment scores, average background values were
211 permuted by randomly drawing a set of 1787 genes 1000 times. P-values and
212 enrichment were computed relative to the permuted samples. Datasets used in the
213 class enrichment in Figure 2D were taken from Sugathan *et al.*(19), Darnel *et al.*(20),
214 and Abrahams *et al.*(4). P-values were computed by the hypergeometric statistical test
215 of overlap between forecASD genes and these three gene sets.

216

217 Cluster analysis of top scoring forecASD genes

218 Using the STRING database, interactions were obtained for forecASD genes and
219 loaded into a network using the igraph package(21) in R(15). No filter for interaction
220 strength was enforced. Hierarchical greedy clustering based on optimization of the
221 modularity score(22) was performed using the fastgreedy.community function in the
222 igraph package. Clustering was performed iteratively, with clusters containing more than

223 200 genes being subject to further clustering. Clusters with fewer than 30 genes were
224 discarded. The annotated network of clusters was loaded into Cytoscape(23), using the
225 STRING application(18). Functional enrichment of clusters was assessed using the
226 STRING application in Cytoscape, with the p-value threshold set to 0.05. For annotation
227 of the network plot shown in figure 6A, either the top annotated term or commonality
228 between several top terms was chosen as representative. The p-value of enrichment
229 with pLI scores was performed using Fisher's exact test of genes within each cluster
230 with a pLI score above 0.5. The p-value for overlap with SFARI HC genes was
231 performed using the hypergeometric test, assuming a background of 18,000 total
232 genes.

233

234 **Results**

235

236 forecASD model and performance

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

244

245 As an initial test of performance, we scored genes hit by coding *de novo* mutations
246 (DNMs) in the recently published MSSNG study. As shown in Figure 2A, there is a
247 significant enrichment of likely gene disrupting (LGD) DNMs in the 90th percentile of
248 both the TADA p-value (OR = 3.76, P = 5.45 x 10⁻⁹) and the forecASD:sys scores (OR
249 = 3.15, P = 7.33 x 10⁻⁸). However, by far the greatest enrichment (OR = 12.81, P < 2.2
250 x 10⁻¹⁶) is seen when restricting to DNMs passing both a TADA q-value and
251 forecASD:sys 90th percentile threshold.

252
253 To leverage both the genetic signal and the systems biology signal, we next built the
254 final forecASD model, which incorporates forecASD:sys, the Krishnan *et al.* score(6),
255 DAMAGES(7), DAWN(5), and several TADA genetic scores from two recent
256 studies(3)⁽¹²⁾. After training the forecASD model, we visualized the variable importance
257 in figure 2B by mean decrease in the Gini impurity measure. The most informative
258 feature was the STRING score from the forecASD:sys model, followed closely by two
259 TADA score variables.

260
261 To facilitate a comparison with manually curated gene prioritizations, we scored all
262 genes in the SFARI gene database using forecASD, forecASD:sys, and the most
263 comprehensive TADA feature in the forecASD model. Shown in figure 2C, the forecASD
264 model ranks SFARI genes scoring 3, 4, 5 and syndromic-only as significantly more
265 autism-related than TADA (P: 7.7x10⁻⁴, 4.7x10⁻¹¹, 2.3x10⁻⁴, 7.7x10⁻⁶). The forecASD
266 model also significantly outperforms the limited forecASD:sys model in gene categories

267 2, 3, and 4 (P : 8.4×10^{-5} , 2.15×10^{-7} , 4.0×10^{-5} , respectively). In all cases, forecASD
268 prioritizes SFARI genes as well, or better than TADA and forecASD:sys.

269
270 As an initial validation of genes prioritized by forecASD, we tested for an enrichment of
271 gene sets and characteristics well known to be overrepresented in autism genes (Fig.
272 2D). We first performed several overrepresentation tests and found that genes receiving
273 forecASD scores in the top decile (1,787 genes, referred to as forecASD genes) had a
274 significant overlap with known targets of CHD8 ($P < 1 \times 10^{-16}$), FMRP ($P < 1 \times 10^{-16}$),
275 and the full SFARI gene database ($P < 1 \times 10^{-16}$). We next performed a series of
276 functional enrichment tests, comparing forecASD genes to randomly sampled sets of
277 background genes. Text mining in PubMed showed that forecASD genes were
278 significantly overrepresented in abstracts which mention autism ($P < 0.001$). Given the
279 established role of autism genes early in fetal development, we next tested and found
280 that forecASD genes showed significantly higher rates of coexpression across all
281 regions of the fetal brain ($P < 0.001$). Lastly, forecASD genes were shown to have
282 significantly enriched rates of interaction in the STRING database ($P < 0.001$).

283
284 We next tested the ability of these scores to discriminate both high confidence (Fig. 3A)
285 and trending (Fig. 3B) autism genes from negative background genes. High confidence
286 autism genes (SFARI HC) are defined as scoring 1 or 2 in SFARI Gene, with trending
287 autism genes scoring 3. Importantly, the negative set of non-autism genes was sampled
288 to have the same background mutation rate as the autism genes ($P > 0.1$ by the
289 Kolmogorov-Smirnov test). In both comparisons, forecASD showed the highest level of

290 performance of all methods tested (AUC=0.97 for SFARI 1+2 and AUC=0.82 for SFARI
291 Gene score 3; Fig. 3). Furthermore, while the SFARI HC genes were used to train the
292 forecASD model, only “out of bag” predictions were used as the forecASD score for
293 those genes, i.e., only those trees where the gene was not included in the bootstrap
294 sample voted for the class of the gene. None of the trending autism genes (Fig. 3B)
295 were used to train forecASD, and consequently they provide an unbiased estimate of
296 performance.

297

298 Generalization to new data: *de novo* mutation enrichment

299 To compare forecASD and prior methods’ ability to generalize to new data, we
300 combined two recently released autism genetics resources. Specifically, we used *de*
301 *novo* mutations in gene regions from the SPARK(9) and MSSNG(8) cohorts.
302 Importantly, none of our model training used information from these studies, thus any
303 subsequent validation is unbiased.

304

305 We first compared forecASD and competing ASD gene scores with respect to
306 enrichment of genes with recurrent *de novo* loss of function and damaging missense
307 mutations in probands. forecASD significantly outperformed all prior approaches
308 (OR=26.8, $P=3.1 \times 10^{-24}$; Fig. 3C). We next tested whether forecASD continued to show
309 significant enrichment when known autism genes (here, any gene listed in the SFARI
310 gene database, regardless of score) were excluded (Fig. 3D), since the ideal method
311 should detect both known and potentially novel autism genes. forecASD had superior

312 performance in this test as well (OR=6.7, P=0.0004), with most of the other external
313 methods lacking a statistically significant enrichment over baseline.

314

315 Functional enrichment and clustering of forecASD genes

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

328

329 While SFARI HC genes show a strong bias toward genes with high pLI (P<0.001,
330 Fisher's exact test; Fig. 5A), forecASD is significantly less biased (P<0.001, Fisher's
331 exact test). We also discovered a significant relationship between pLI and differential
332 expression (DE) t-statistics in case/control brain gene expression studies(17) (beta=-
333 0.13, t-statistic-4.3, P=1.9x10⁻⁵, Fig. 5B), potentially exposing a form of bias in current
334 gene discovery approaches that leads to under-ascertainment of ASD risk genes with

335 low pLI and upregulation in ASD cases. We also found a significant interaction between
336 forecASD and pLI ($F=54.1$, $P=3.9 \times 10^{-24}$) such that the pLI-expression relationship exists
337 among forecASD genes ($\beta=-0.24$, $t=-2.6$, $P=0.009$; Fig. 5D) but is absent in non-
338 forecASD genes ($\beta=0.004$, $t=0.1$, $P=0.91$; Fig. 5C).

339
340 Lastly, forecASD genes were loaded into the STRING network and clustered using a
341 greedy hierarchical approach which maximizes the modularity score. The resulting
342 networks consisted of 17 clusters composed of 1452 genes. All clusters were found to
343 be significantly enriched with numerous GO and KEGG pathways (Supplemental Table
344 3). Similarly, all clusters contained a significant enrichment of haploinsufficiency genes
345 ($pLI > 0.5$), except for the small cluster of 31 genes with functions related to the
346 mediator complex. Clusters were also tested for overlap with SFARI HC genes, of which
347 8 clusters failed to reach significance, suggesting groupings of genes currently missing
348 from the known list of autism genes. Clusters lacking significant overlap includes those
349 with functions: signal transduction, cytoskeleton, cell migration, neuron projection,
350 steroid signaling, neuron differentiation, potassium signaling, development and
351 morphogenesis. A marginally significant correlation was seen between a clusters
352 enrichment for high pLI genes and its overlap with SFARI HC genes (Spearman's $r =$
353 0.48 , $p\text{-value} = 0.053$), further suggesting a bias in SFARI HC genes towards
354 haploinsufficiency status.

355

356 **Discussion**

357

358 We present forecASD, a machine learning approach that combines systems biology and
359 genetic models into a single score that indexes the strength of evidence for a gene's
360 involvement in autism. This genome-wide score can be a useful prior, filter, or positive
361 control in molecular studies of autism. It can also be used as a starting point to generate
362 new hypotheses to investigate currently under-appreciated aspects of the molecular
363 etiology of autism. In our tests of predictive performance and generalization, forecASD
364 outperformed other systems biology and genetic approaches for autism gene
365 prioritization.

366 Because it draws upon multiple approaches for identifying autism genes, forecASD is
367 less biased than gene discovery based only on one form of data (e.g., genetic data).
368 This is particularly important because current SFARI HC genes, which rely heavily on
369 studies of *de novo* mutation, are strongly biased towards genes that are loss-of-function
370 intolerant (Fig. 5A). While these haploinsufficient genes represent a sizable and
371 important component of genetic risk for autism, this ascertainment bias has led to
372 molecular “blind spots” that will not be resolved simply by sequencing more probands
373 and identifying additional *de novo* mutations. For instance, pathways implicated
374 preferentially by SFARI HC genes had significantly higher pLI, whereas pathways with
375 lower pLI were under-represented (compared to forecASD-implicated pathways;
376 $OR=0.38$, $P=5.6 \times 10^{-7}$). Furthermore, while SFARI HC pathways significantly predicted
377 case/control expression-implicated pathways ($Z=4.5$, $P=7.9 \times 10^{-6}$, binomial model), only
378 3% of the deviance could be explained. In contrast, forecASD pathways explained an
379 order of magnitude more deviance (31%, $Z=12.1$, $P=1.5 \times 10^{-33}$) when predicting
380 expression-implicated pathways. When a model of dichotomous DE significance was fit

381 that included terms from both SFARI HC and forecASD pathways, the SFARI HC term
382 became redundant, and the forecASD-only model yielded a superior Bayesian
383 information criterion (BIC; 518 for forecASD-only vs. 521 for full model and 720 for the
384 SFARI HC-only model). When considering directionality, SFARI HC gene pathways
385 were significantly depleted for ASD-upregulated pathways (OR=0.48, $P=1.3\times 10^{-6}$),
386 further illustrating the bias in SFARI HC genes. These results demonstrate that
387 forecASD showed greater representation of low pLI and ASD-upregulated pathways,
388 without sacrificing sensitivity to well-known ASD risk pathways where haploinsufficiency
389 plays a dominant role.

390 In our analyses, we noted a trend that is a potential bridge between gene discovery
391 studies based on DNA sequence variants and those based on differential expression.
392 Specifically, pLI is significantly and negatively correlated with previously published
393 frontal cortex differential expression (case/control) t-statistics(17) ($\beta=-0.13$, t-statistic-
394 4.3, $P=1.9\times 10^{-5}$). This suggests that low-pLI genes are more likely to be up-regulated
395 and high-pLI genes are more likely to be down-regulated in ASD cases (Fig. 5B). This is
396 consistent with our observation of SFARI HC gene pathways (which have an
397 ascertainment bias in favor of haploinsufficiency) being significantly under-represented
398 in both low pLI and ASD-upregulated pathways. We further observed a significant
399 interaction ($F=54.1$, $P=3.9\times 10^{-24}$) between forecASD and pLI when explaining variation
400 in ASD brain gene expression: forecASD genes (i.e., top decile) show the significant
401 negative relationship between pLI and t-statistic ($\beta=-0.24$, $t=-2.6$, $P=0.009$; Fig. 5D),
402 while non-forecASD genes show no relationship ($\beta=0.004$, $t=0.1$, $P=0.91$; Fig. 5C).
403 Consequently, we propose that this pLI-expression relationship is a hallmark of robust

404 ASD risk genes, and may be used as a criterion when identifying optimal thresholds in
405 genome-wide scores like forecASD. Indeed, although initially chosen as a convenient
406 but arbitrary threshold for identifying a discrete set of ASD candidate genes, the top
407 decile proved to be the optimal split point for forecASD, maximizing the significance of
408 the pLI/t-statistic relationship among candidate genes, while minimizing the same
409 relationship in the remaining, non-candidate genes. Interestingly, when applying this
410 approach to TADA FDR values, although TADA-implicated genes showed the expected
411 pLI-expression relationship, no TADA threshold was able to eliminate the trend from
412 non-candidate genes, suggesting lower sensitivity in identifying ASD risk genes
413 compared to forecASD. Taken together, these analyses demonstrate that the reduced
414 bias in forecASD contributes to increased sensitivity to autism risk pathways identified in
415 gene expression studies (Fig. 4C,4D) as well as those implicated by genetic studies
416 (Fig. 3).

417 Some pathways, although represented (but not necessarily enriched) in the current
418 SFARI HC list, showed a substantial relative increase in enrichment when considering
419 forecASD (Fig. 4A, Supplemental Table 2). This suggests that these pathways
420 represent noted and plausible but still under-appreciated molecular themes in our
421 understanding of autism. The pathway that underwent the largest relative increase in
422 enrichment from SFARI HC to forecASD is Rho GTPase signaling ($OR=2.2$, $P=4.8 \times 10^{-5}$),
423 which plays a critical role in cytoskeletal dynamics in neurodevelopment(24),
424 including interactions with SHANK proteins and the formation and maturation of
425 dendritic spines(25). As another example, although chromatin modification in general is
426 a well-established theme in autism genetic risk, histone acetyltransferases showed

427 relatively little representation in the SFARI HC list, but were significantly enriched in
428 forecASD genes ($OR=4.1$, $P=3.5 \times 10^{-9}$). Histone acetylation was recently shown to be a
429 pervasive genomic predictor of affected status in a large autism case/control
430 postmortem brain study(26), underscoring the importance of this mechanism that is
431 under-represented in established risk genes but that forecASD was sensitive to. As a
432 final example of these under-appreciated molecular mechanisms, the circadian clock
433 pathway was implicated by forecASD as an important source of risk for autism ($OR=6.5$,
434 $P=7.6 \times 10^{-13}$). Sleep disturbances are a well-known and problematic comorbidity in
435 autism, and molecular deficits in circadian regulation related to autism have been
436 documented(27),(28),(29). Although literature support is available for these processes
437 playing a role in autism, our results indicate that their current sparse representation in
438 lists of accepted genetic risk factors is not representative of their importance in the
439 disorder.

440 Other pathways were identified by forecASD as significantly enriched for autism risk, but
441 were not represented at all among SFARI HC genes (Supplemental Table 2, Figure 4B).
442 Consequently, we expect that new insights into the molecular basis of autism will come
443 disproportionately from these pathways as their constituent genes are associated with
444 autism. One gene set in particular, potassium channels, showed highly significant
445 enrichment in forecASD genes ($OR=4.1$, $P=7.2 \times 10^{-9}$, $N=35$ genes) despite the absence
446 of potassium channel genes among currently accepted autism risk genes. However, the
447 literature shows support for a role for potassium channels in ASD risk(30),(31),(32),(33), and
448 the pathway was enriched for differential regulation in a recently published brain gene
449 expression study of autism ($P=0.001$, downregulated)(17). Notably, this pathway has a

450 lower proportion of genes with pLI>0.9 (0.22) compared to SFARI HC gene-implicated
451 pathways (median=0.47), potentially explaining its absence due to ascertainment bias.
452 Overall, pathways that demonstrated forecASD-specific excess enrichment showed a
453 significant agreement with pathway enrichment from independent case/control brain
454 gene expression studies (OR=28.8; P=2.9x10⁻⁴⁸), and were more likely to support
455 pathways that were up-regulated in the gene expression data (OR=2.1, P=1.27x10⁻⁶,
456 Fig. 4C) compared to pathways implicated by SFARI HC.

457 To group forecASD genes into distinct functional categories, we performed iterative
458 clustering and identified a total of 17 clusters enriched for specific functional
459 annotations. While nearly all clusters showed significant enrichment for
460 haploinsufficiency genes, many lacked a significant overlap with SFARI HC genes, after
461 Bonferroni correction. Similar to conclusions reached above, we found an entire cluster
462 enriched for Potassium signaling (P=1.8x10⁻⁴⁹) which lacked significant overlap with
463 SFARI HC genes. In addition to this cluster, there were also seven others lacking
464 significant overlap with SFARI HC genes. Notable examples include clusters related to
465 cell migration (P=6.0x10⁻¹¹) and endocytosis (P=2.7x10⁻²¹). These pathways have more
466 recently been explored in their ability to regulate brain connectivity(34) and postsynaptic
467 organization(35), respectively. In agreement with the proposed haploinsufficiency bias
468 of autism gene discovery, we observed a marginally significant relationship between
469 cluster pLI enrichment and SFARI HC gene overlap (Spearman's rho: 0.48; P=0.053).

470 During the development of forecASD, another ASD gene prediction method was
471 published(36), ASD-FRN, which utilizes a brain-specific functional network. We
472 evaluated this method using the performance benchmarks presented in Figure 3, and

473 found it to have performance comparable to DAMAGES and Krishnan, et al.
474 (Supplemental Figure 1 and Supplemental Figure 2), but in none of these tests did it
475 surpass the performance observed using forecASD as an ASD gene predictor. We also
476 added ASD-FRN to the ensemble that comprises forecASD, but we did not observe a
477 significant increase in predictive performance, suggesting that the latent information in
478 ASD-FRN is already accounted for in the forecASD ensemble as presented here.

479 Our use of TADA to impart genetic association information to the forecASD ensemble is
480 unique among the ASD gene prediction approaches we used as benchmarks. However,
481 this raises concerns about the potential for circularity: TADA is emerging as the most
482 popular way to compute and update gene-wise genetic association statistics for ASD
483 studies, and previous TADA scores are strongly correlated with updated TADA scores.
484 Furthermore, TADA scores are among the most important predictive features in the
485 forecASD ensemble (Fig. 2). Consequently, it would be concerning if forecASD's ability
486 to predict new ASD genes was due entirely to the inclusion of previously published
487 TADA scores. To examine this possibility, we fit bivariate logistic regression models,
488 always including TADA rankings as a covariate, with Krishnan, DAMAGES, ASD-FRN,
489 or forecASD rankings as the predictor of membership among SFARI Gene score 3
490 (Supplemental Figure 2; score 3 genes were not used in training forecASD, but are
491 nevertheless assumed to be enriched for true ASD genes). All other genes listed in the
492 SFARI Gene database were removed from the analysis, and the remainder of the
493 genome was considered the negative class. Strong association with score 3
494 membership (as measured by the regression coefficient and Z-score) is a desirable trait
495 for an ASD gene prediction method, and forecASD proved to be the most strongly

496 associated of the tested methods, even after correcting for the information imparted by
497 TADA rankings ($Z=11.9$, $P=8.9\times 10^{-33}$).

498 Rather than simply parroting gene prioritization according to TADA, forecASD purifies
499 TADA's signal by imposing biological plausibility through the inclusion of other predictive
500 features, such as network and expression data, as well as previously published
501 machine learning approaches to ASD gene prediction. To illustrate this filtering effect
502 that forecASD has on TADA rankings, we performed a GSEA on the differences in gene
503 ranks between forecASD and TADA. This analysis highlights pathways and processes
504 favored by forecASD, vs. those favored by TADA (Supplemental Figure 3). Pathways
505 where there is consensus fall in the middle of the difference-in-ranks distribution, and
506 consequently are not highlighted as enriched or depleted in this analysis. In this
507 analysis, forecASD prioritized genes involved in well-established neurodevelopmental
508 pathways, such as signaling by receptor tyrosine kinases, axon guidance, Rho GTPase
509 signaling, and chromatin modification. Conversely, TADA prioritized genes belonging to
510 large gene families, such as olfactory receptors, ribosomal proteins, and defensins.
511 These kinds of genes are routinely manually removed from consideration even in the
512 presence of statistical support, because they often lack biological plausibility. forecASD
513 accomplishes this same task, but in a more automated, data-driven, and objective way.

514 One bias in forecASD uncovered by the above analysis, and that is intuitive based on
515 forecASD's reliance on functional data, is that forecASD tends to de-prioritize genes
516 that are poorly annotated or lack strong signals in expression or other functional data
517 sets (in the above GSEA, forecASD-preferred genes were depleted for "Uncategorized"
518 at $P=6.8\times 10^{-9}$). On one hand, this limits forecASD's ability to make entirely unexpected

519 discoveries among poorly characterized genes. On the other hand, the probability of
520 highly important genes completely evading decades of neuroscience investigation is
521 diminishing, and penalizing these “unassuming” genes may be justifiable. Nevertheless,
522 the value of new discovery should not be discounted, and this will be an area of active
523 investigation in the ongoing development of forecASD.

524 As a final proof of principle for one of forecASD’s intended uses, we investigated genes
525 nominated as new ASD genes by a currently unpublished large-scale gene discovery
526 study(37). The data underlying this study was not available to us during the
527 development of forecASD, so this represents a true test of forecASD’s generalization
528 and interpretive value. This study nominated *BTRC*, *C16orf13*, *CCSER1*, *CMPK2*,
529 *DDX3X*, *FAM98C*, *GRIA1*, *MLANA*, *MYO5A*, *PCM1*, *PRKAR1B*, *RAPGEF4*, *SMURF1*,
530 *TMEM39B*, *TSPAN4*, and *UIMC1* as newly discovered ASD genes at a TADA FDR of
531 0.2, meaning that of these 16 genes, 2-3 are expected to be false positives. Overall,
532 the nominated genes are strongly enriched for elevated forecASD scores ($P=4.8 \times 10^{-8}$,
533 Wilcoxon test), and 12 of the 16 genes scored in the top decile of forecASD. This result
534 strengthens both the case for forecASD as an effective predictor of ASD risk genes, and
535 the findings of this gene discovery study. However, *FAM98C* and *CCSER1* have very
536 low forecASD scores (0.022 and 0.024, respectively), indicating little support in network
537 and functional data. *UIMC1* and *C16orf13* have slightly higher scores (0.210 and 0.346,
538 respectively), but are still below the top-decile threshold we have used in this study
539 (0.374). Consequently, because of weak support beyond the genetic association, these
540 four genes are candidates for being the false positives, and should require additional

541 experimental data supporting a role in autism before they are considered true ASD risk
542 genes.

543 For the foreseeable future, traditional gene discovery studies will continue to add to the
544 list of bona fide ASD risk genes. Eventually, as sample sizes saturate and gene
545 discovery decelerates, the field will be faced with the challenge of developing new and
546 useful applications of this acquired knowledge. By combining new and previously
547 published predictors into a high-performance ensemble classifier, forecASD provides a
548 glimpse of that future and gives an opportunity right now to begin thinking about what
549 we would do with a definitive list of autism genes.

550 **Conclusions**

551 We introduce a new model, forecASD, for prioritizing autism-associated risk genes. We
552 show that forecASD has significantly improved sensitivity and specificity compared to
553 previous gene-level predictors, including the purely genetic-based approach TADA. We
554 demonstrate forecASD's usefulness as an autism risk-gene discovery post-hoc filter
555 through its ability to prioritize 12 of 16 genes implicated at FDR=0.2 by the latest ASD
556 gene discovery study. Using forecASD prioritized risk-genes, we highlight which
557 molecular pathways are currently under-represented in the autism literature and likely
558 represent under-appreciated biological mechanisms of autism.

559 **List of abbreviations**

560 LGD- likely gene disrupting, DNM- de novo mutation, TADA- Transmission And De novo
561 Association, DE- differential expression, GSEA- gene set enrichment analysis

562 **Declarations**

563 **Ethics approval and consent to participate:** All human genetic data used in this study
564 was accessed in a deidentified manner from associated sequencing consortia with their
565 approval. Due to the method of access, this study is not formally considered human
566 subjects research and therefore not subject to associated restrictions, as outlined by the
567 National Institutes of Health.

568 **Consent for publication:** Not applicable.

569 **Availability of data and material:** All data and code used to generate the forecASD
570 model is available for download from the repository associated with this project,
571 <https://github.com/LeoBman/forecASD>.

572 **Competing interests:** The authors declare that they have no competing interests.

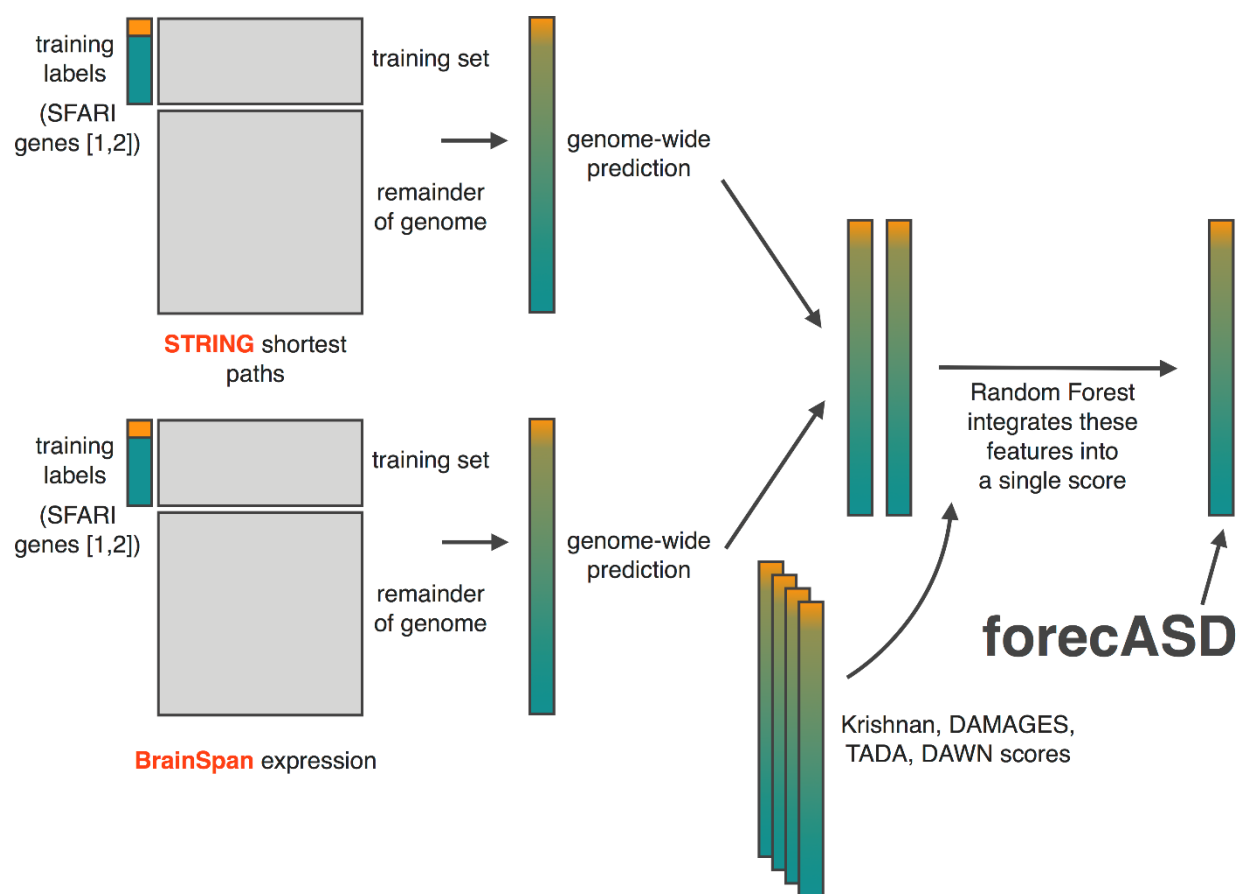
573 **Funding:** This work was supported by the National Institutes of Health [MH105527 and
574 DC014489 to JJM]. This work was supported by a grant from the Simons Foundation
575 (SFARI # 516716, [JJM]).

576 **Authors' contributions:** LB contributed to design, testing, and implementation of
577 forecASD model and was a major contributor in writing the manuscript. TK contributed
578 to analysis of forecASD model results and contributed to github implementation of
579 forecASD. JM contributed to design, testing, and implementation of forecASD model
580 and was a major contributor in writing the manuscript. All authors read and approved
581 the final manuscript.

582 **Acknowledgements**

583 We are grateful to all of the families in SPARK, the SPARK clinical sites and SPARK
584 staff. We appreciate obtaining access to exome sequencing and phenotypic data on
585 SFARI Base. Approved researchers can obtain the SPARK population dataset
586 described in this study by applying at <https://base.sfari.org>.

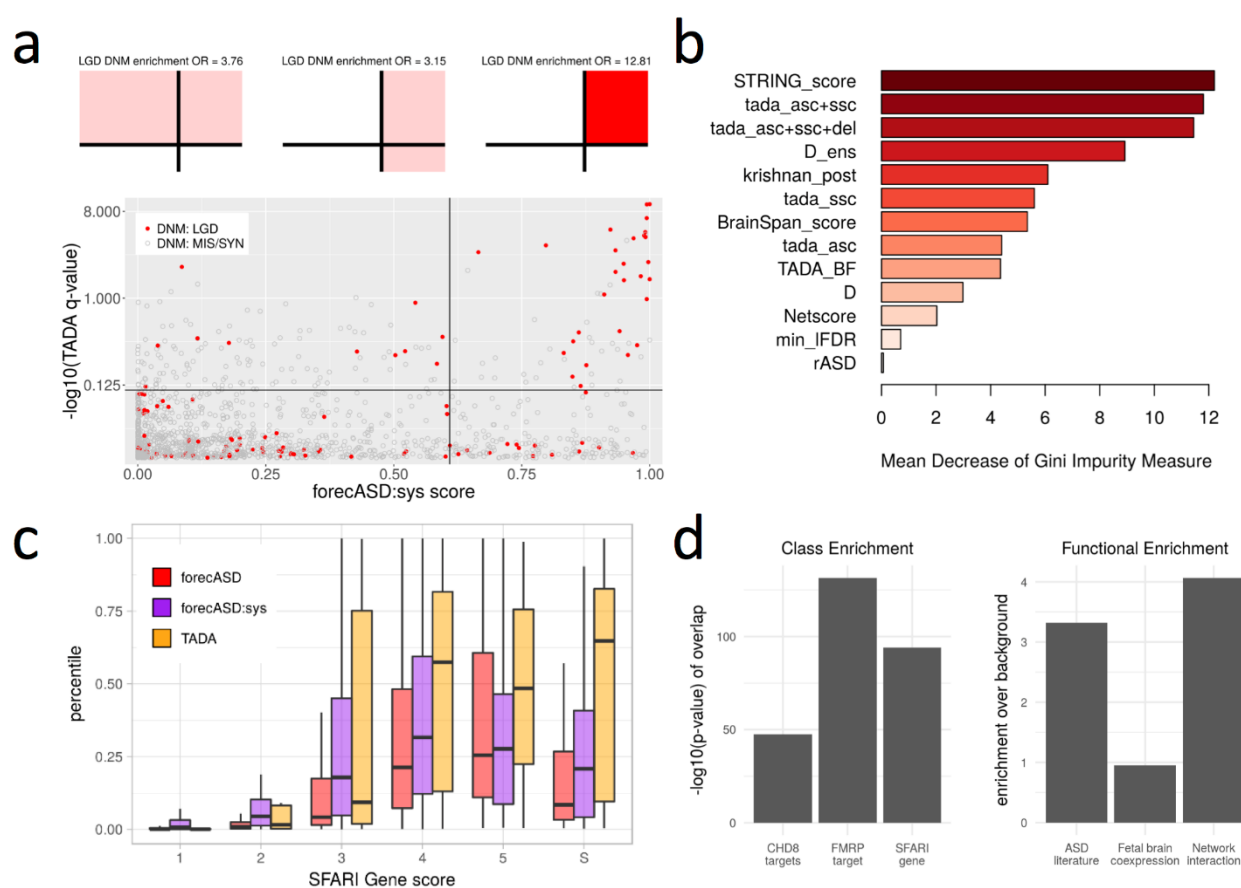
587



588

589 **Figure 1 - Overview of forecASD.** Two Random Forest classifiers, one using
590 BrainSpan gene expression and the other using the STRING network as predictors, are
591 trained to discriminate high confidence autism genes (SFARI HC, scores 1 and 2) from
592 a set of 1,000 genes drawn randomly from those not listed at all in the SFARI Gene
593 database. Predictions are then made on the remainder of the genome, and these are

594 combined with the out-of-bag (OOB) estimates from the training process to yield a
 595 prediction for each gene in the genome. A subsequent classifier is then trained using
 596 the output of these two RFs and previously published autism gene scores as predictive
 597 features, and again predictions are made on the remainder of the genome, with OOB
 598 predictions being used for those genes in the training set. The RF vote proportion for
 599 class “autism gene” is then the final forecASD score.
 600

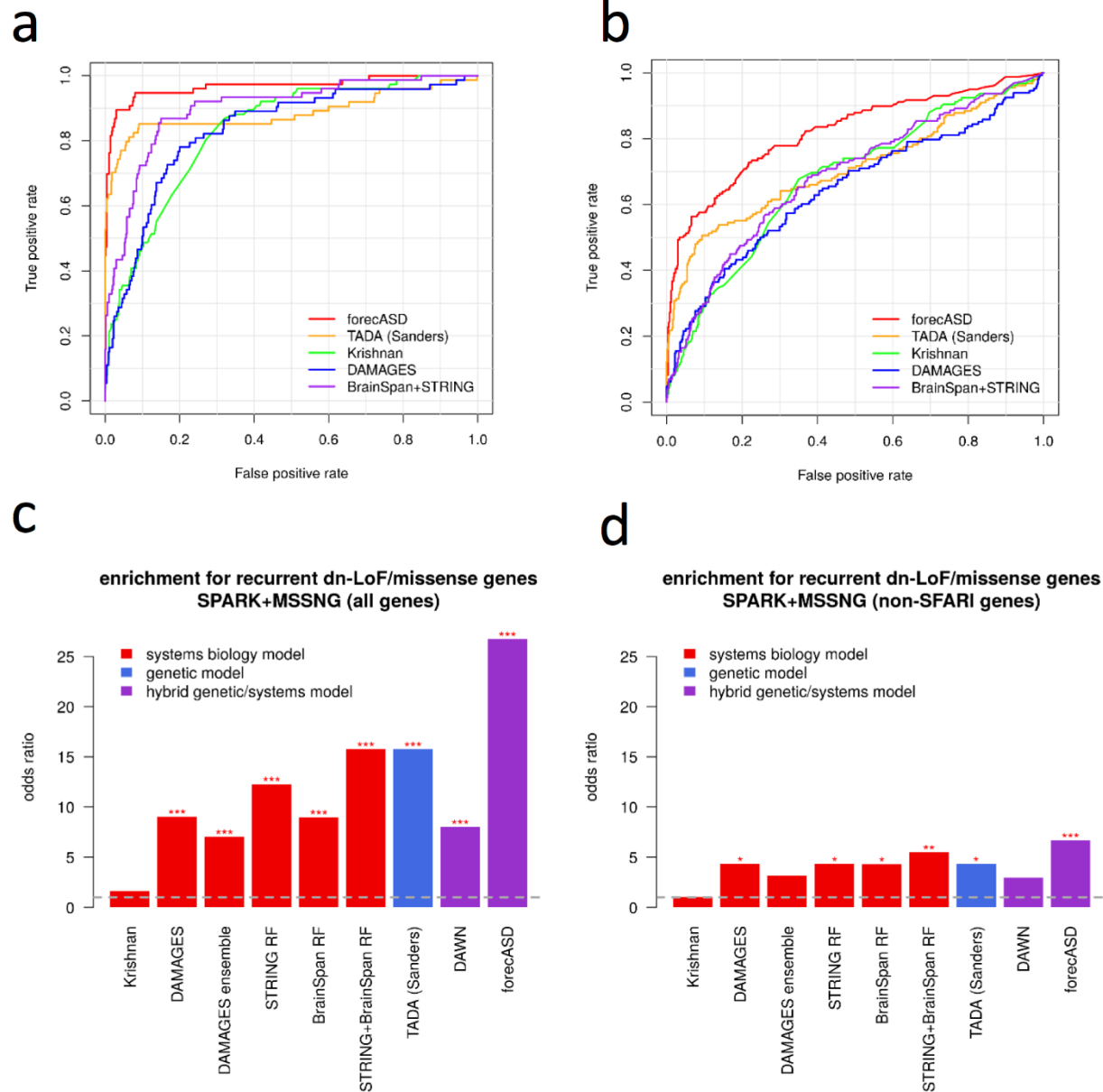


601

602 **Figure 2 - Prioritization of de novo likely gene-disrupting mutations and**
 603 **enrichment of gene sets in forecASD.** Training a limited model, forecASD:sys, using
 604 brain gene expression and interaction data shows optimal prioritization of de novo LGDs
 605 when combined with a genetic measure of autism association (a). Building the full

606 forecASD model, we test all features for their informativeness, finding that the STRING
607 score is primary (b). Using the three mentioned scores, we assess their genome-wide
608 ranking of SFARI genes at all levels, and find that the full forecASD model at least ties,
609 and often significantly outperforms TADA and forecASD:sys in the prioritization of
610 SFARI genes (c). As an initial assessment of forecASD prioritized genes, we find the
611 top decile of genes ranked by forecASD (1787 genes) shows enrichment typical of
612 classical autism genes (d).

613



614

615 **Figure 3 - Comparison of forecASD with prior models of autism gene**

616 **prioritization.** To compare forecASD with competitors, we evaluate performance by

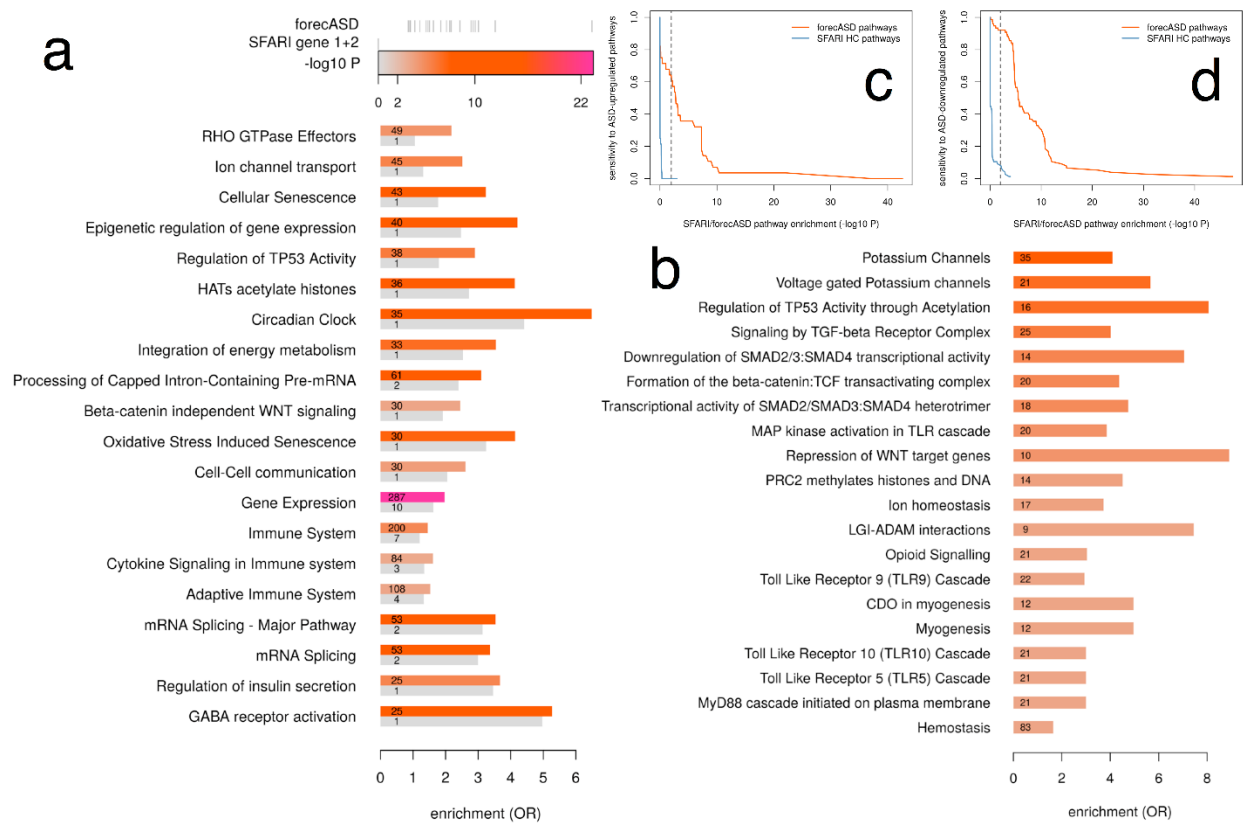
617 each methods' ability to prioritize SFARI genes and genes which were subject to

618 recurrent de novo loss-of-function or missense mutations. Starting with SFARI genes

619 scoring 1 or 2 as a positive set and size-matched random background genes as the

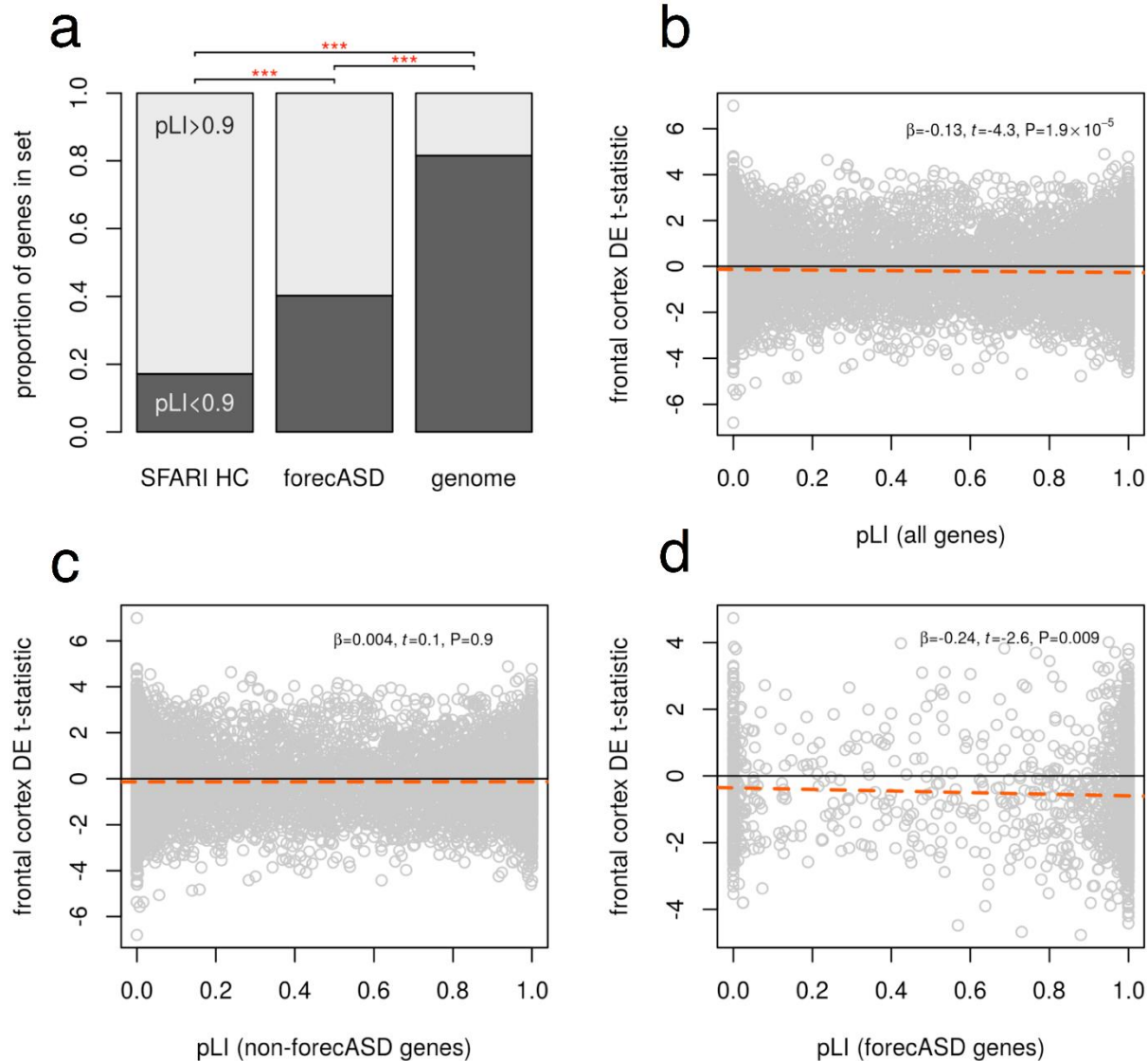
620 negative set, forecASD out-of-bag estimates showed superior classification over all

621 methods (a). In a fully unbiased test, forecASD estimates also showed superior
 622 classification of trending SFARI genes (score: 3) over all other methods (b). Using two
 623 sequencing cohorts which no methods draw information from, the top decile of
 624 forecASD genes (1787 genes) shows the greatest overlap with genes containing
 625 recurrent de novo loss-of-function and missense mutations (c). When excluding genes
 626 in the SFARI gene database, forecASD still shows superior prioritization of genes
 627 accumulating de novo mutations (d).
 628



629
 630 **Figure 4 - forecASD-specific pathway enrichment and sensitivity to gene**
 631 **expression-implicated pathways.** When testing the top-decile genes according to
 632 forecASD for Reactome pathway enrichment, pathways emerged that were
 633 represented, but not enriched in the SFARI HC list (a). Other pathways were highly

634 enriched in forecASD genes that were not represented at all in the SFARI HC list, even
635 though they have associated literature suggesting a role in autism (b). forecASD is
636 more sensitive than SFARI HC to pathways that are differentially regulated in the brains
637 of individuals with autism, particularly in ASD-upregulated pathways (c), but also in
638 downregulated pathways (d). Using the top decile of TADA $-\log_{10}$ FDR genes showed
639 similar sensitivity to SFARI HC (not shown), suggesting that rare variant approaches
640 may be less sensitive in implicating genes found through gene expression studies.
641



642

643 **Figure 5 - Relationship between pLI and ASD-specific up- and down-regulation of**

644 **brain gene expression.** SFARI HC is strongly biased toward genes with high pLI (a),

645 while forecASD is significantly less biased. We found a significant relationship between

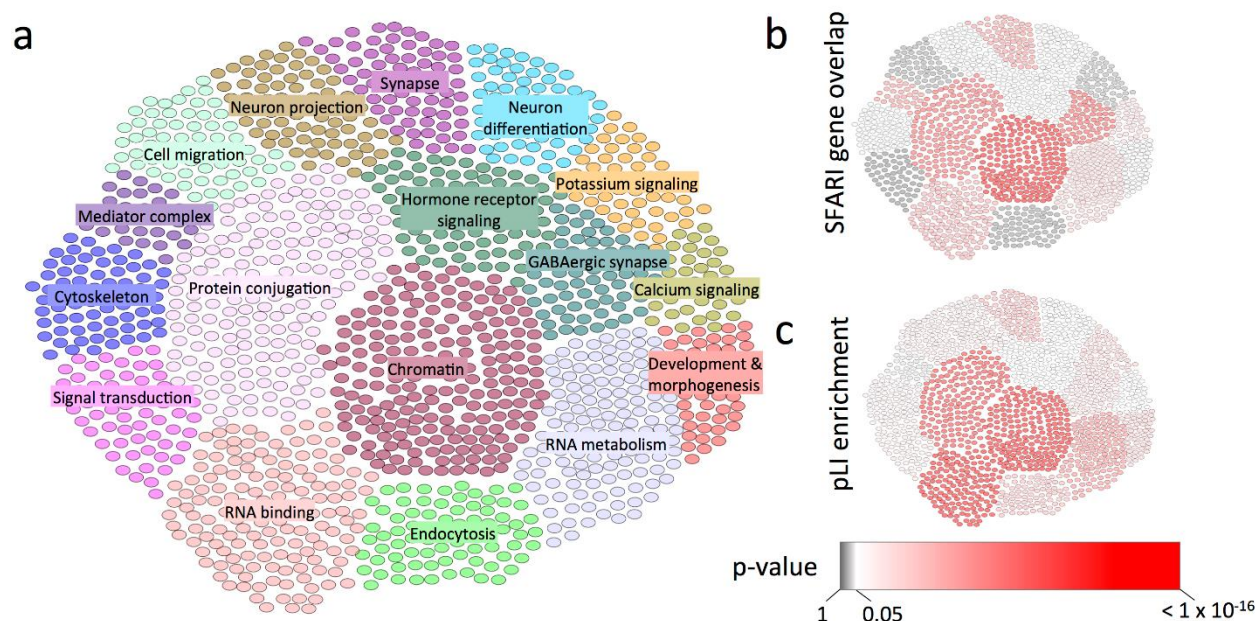
646 pLI and differential expression (DE) in the brains of autism cases (b), such that low pLI

647 genes tend toward upregulation in cases, while high pLI genes tend to be

648 downregulated. We also observed a significant interaction between forecASD and pLI

649 such that the observed pLI-DE trend (b) is absent in non-forecASD genes (c), and

650 present and significant among forecASD genes (d). We propose that the presence of
651 the pLI-DE trend is a hallmark of ASD risk genes, and an optimal ASD gene
652 prioritization method will concentrate the trend among risk genes and remove it from
653 non-risk genes. Notably, no threshold of TADA (tested to the 50th percentile) was able to
654 remove the trend from the non-prioritized genes, suggesting the persistence of residual
655 risk genes that were not selected.
656



657
658 **Figure 6 - Clustering of top forecASD genes with ExAC pLI and SFARI high-**
659 **confidence gene overlap enrichment analysis..** Greedy hierarchical optimization of
660 the modularity score yielded 17 clusters consisting of 1452 forecASD genes (a). All
661 clusters have several significantly enriched biological pathways, of which the top terms
662 were overlaid in figure 6A. Clusters were tested for significance of overlap with the list of
663 SFARI HC genes (b), and enrichment of haploinsufficiency genes (pLI > 0.5; c).
664

665 References

- 666 1. Rosenberg RE, Law JK, Yenokyan G, McGready J, Kaufmann WE, Law PA. Characteristics and
667 concordance of autism spectrum disorders among 277 twin pairs. *Arch Pediatr Adolesc Med*.
668 2009;163(10):907-14.
- 669 2. Colvert E, Tick B, McEwen F, Stewart C, Curran SR, Woodhouse E, et al. Heritability of Autism
670 Spectrum Disorder in a UK Population-Based Twin Sample. *JAMA Psychiatry*. 2015;72(5):415-23.
- 671 3. De Rubeis S, He X, Goldberg AP, Poultney CS, Samocha K, Cicek AE, et al. Synaptic,
672 transcriptional and chromatin genes disrupted in autism. *Nature*. 2014;515(7526):209-15.
- 673 4. Abrahams BS, Arking DE, Campbell DB, Mefford HC, Morrow EM, Weiss LA, et al. SFARI Gene
674 2.0: a community-driven knowledgebase for the autism spectrum disorders (ASDs). *Mol Autism*.
675 2013;4(1):36.
- 676 5. Liu L, Lei J, Sanders SJ, Willsey AJ, Kou Y, Cicek AE, et al. DAWN: a framework to identify autism
677 genes and subnetworks using gene expression and genetics. *Mol Autism*. 2014;5(1):22.
- 678 6. Krishnan A, Zhang R, Yao V, Theesfeld CL, Wong AK, Tadych A, et al. Genome-wide prediction
679 and functional characterization of the genetic basis of autism spectrum disorder. *Nat Neurosci*.
680 2016;19(11):1454-62.
- 681 7. Zhang C, Shen Y. A Cell Type-Specific Expression Signature Predicts Haploinsufficient Autism-
682 Susceptibility Genes. *Hum Mutat*. 2017;38(2):204-15.
- 683 8. RK CY, Merico D, Bookman M, J LH, Thiruvahindrapuram B, Patel RV, et al. Whole genome
684 sequencing resource identifies 18 new candidate genes for autism spectrum disorder. *Nat Neurosci*.
685 2017;20(4):602-11.
- 686 9. pfeliciano@simonsfoundation.org SCEa, Consortium S. SPARK: A US Cohort of 50,000 Families to
687 Accelerate Autism Research. *Neuron*. 2018;97(3):488-93.
- 688 10. Sunkin SM, Ng L, Lau C, Dolbeare T, Gilbert TL, Thompson CL, et al. Allen Brain Atlas: an
689 integrated spatio-temporal portal for exploring the central nervous system. *Nucleic Acids Res*.
690 2013;41(Database issue):D996-D1008.
- 691 11. von Mering C, Huynen M, Jaeggi D, Schmidt S, Bork P, Snel B. STRING: a database of predicted
692 functional associations between proteins. *Nucleic Acids Res*. 2003;31(1):258-61.
- 693 12. Sanders SJ, He X, Willsey AJ, Ercan-Sencicek AG, Samocha KE, Cicek AE, et al. Insights into Autism
694 Spectrum Disorder Genomic Architecture and Biology from 71 Risk Loci. *Neuron*. 2015;87(6):1215-33.
- 695 13. Fabregat A, Jupe S, Matthews L, Sidiropoulos K, Gillespie M, Garapati P, et al. The Reactome
696 Pathway Knowledgebase. *Nucleic Acids Res*. 2018;46(D1):D649-D55.
- 697 14. Mi H, Lazareva-Ulitsky B, Loo R, Kejariwal A, Vandergriff J, Rabkin S, et al. The PANTHER
698 database of protein families, subfamilies, functions and pathways. *Nucleic Acids Res*. 2005;33(Database
699 issue):D284-8.
- 700 15. R Development Core Team. R: A language and environment for statistical computing. Vienna,
701 Austria: R Foundation for Statistical Computing; 2008.
- 702 16. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-
703 coding genetic variation in 60,706 humans. *Nature*. 2016;536(7616):285-91.
- 704 17. Gandal MJ, Haney JR, Parikshak NN, Leppa V, Ramaswami G, Hartl C, et al. Shared molecular
705 neuropathology across major psychiatric disorders parallels polygenic overlap. *Science*.
706 2018;359(6376):693-7.
- 707 18. Szklarczyk D, Morris JH, Cook H, Kuhn M, Wyder S, Simonovic M, et al. The STRING database in
708 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids
709 Res*. 2017;45(D1):D362-D8.

- 710 19. Sugathan A, Biagioli M, Golzio C, Erdin S, Blumenthal I, Manavalan P, et al. CHD8 regulates
711 neurodevelopmental pathways associated with autism spectrum disorder in neural progenitors. *Proc*
712 *Natl Acad Sci U S A*. 2014;111(42):E4468-77.
- 713 20. Darnell JC, Van Driesche SJ, Zhang C, Hung KY, Mele A, Fraser CE, et al. FMRP stalls ribosomal
714 translocation on mRNAs linked to synaptic function and autism. *Cell*. 2011;146(2):247-61.
- 715 21. Csardi G NT. The igraph software package for complex network research. *InterJournal*2006.
- 716 22. Newman ME. Modularity and community structure in networks. *Proc Natl Acad Sci U S A*.
717 2006;103(23):8577-82.
- 718 23. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: a software
719 environment for integrated models of biomolecular interaction networks. *Genome Res*.
720 2003;13(11):2498-504.
- 721 24. Reichova A, Zatkova M, Bacova Z, Bakos J. Abnormalities in interactions of Rho GTPases with
722 scaffolding proteins contribute to neurodevelopmental disorders. *J Neurosci Res*. 2018;96(5):781-8.
- 723 25. Martin-Vilchez S, Whitmore L, Asmussen H, Zareno J, Horwitz R, Newell-Litwa K. RhoGTPase
724 Regulators Orchestrate Distinct Stages of Synaptic Development. *PLoS One*. 2017;12(1):e0170464.
- 725 26. Sun W, Poschmann J, Cruz-Herrera Del Rosario R, Parikshak NN, Hajan HS, Kumar V, et al.
726 Histone Acetylome-wide Association Study of Autism Spectrum Disorder. *Cell*. 2016;167(5):1385-97 e11.
- 727 27. Lipton JO, Boyle LM, Yuan ED, Hochstrasser KJ, Chifamba FF, Nathan A, et al. Aberrant
728 Proteostasis of BMAL1 Underlies Circadian Abnormalities in a Paradigmatic mTOR-opathy. *Cell Rep*.
729 2017;20(4):868-80.
- 730 28. Monyak RE, Emerson D, Schoenfeld BP, Zheng X, Chambers DB, Rosenfelt C, et al. Insulin
731 signaling misregulation underlies circadian and cognitive deficits in a *Drosophila* fragile X model. *Mol*
732 *Psychiatry*. 2017;22(8):1140-8.
- 733 29. Kozlov SV, Bogenpohl JW, Howell MP, Wevrick R, Panda S, Hogenesch JB, et al. The imprinted
734 gene *Magel2* regulates normal circadian output. *Nat Genet*. 2007;39(10):1266-72.
- 735 30. Guglielmi L, Servettini I, Caramia M, Catacuzzeno L, Franciolini F, D'Adamo MC, et al. Update on
736 the implication of potassium channels in autism: K(+) channelautism spectrum disorder. *Front Cell*
737 *Neurosci*. 2015;9:34.
- 738 31. Deng PY, Klyachko VA. Genetic upregulation of BK channel activity normalizes multiple synaptic
739 and circuit defects in a mouse model of fragile X syndrome. *J Physiol*. 2016;594(1):83-97.
- 740 32. Lee H, Lin MC, Kornblum HI, Papazian DM, Nelson SF. Exome sequencing identifies de novo gain
741 of function missense mutation in *KCND2* in identical twins with autism and seizures that slows
742 potassium channel inactivation. *Hum Mol Genet*. 2014;23(13):3481-9.
- 743 33. Sicca F, Ambrosini E, Marchese M, Sforna L, Servettini I, Valvo G, et al. Gain-of-function defects
744 of astrocytic Kir4.1 channels in children with autism spectrum disorders and epilepsy. *Sci Rep*.
745 2016;6:34325.
- 746 34. Reiner O, Karzbrun E, Kshirsagar A, Kaibuchi K. Regulation of neuronal migration, an emerging
747 topic in autism spectrum disorders. *J Neurochem*. 2016;136(3):440-56.
- 748 35. Loebrich S. The role of F-actin in modulating Clathrin-mediated endocytosis: Lessons from
749 neurons in health and neuropsychiatric disorder. *Commun Integr Biol*. 2014;7:e28740.
- 750 36. Duda M, Zhang H, Li HD, Wall DP, Burmeister M, Guan Y. Brain-specific functional relationship
751 networks inform autism spectrum disorder gene prediction. *Transl Psychiatry*. 2018;8(1):56.
- 752 37. Ruzzo EK, Perez-Cano L, Jung J-Y, Wang L-k, Kashef-Haghighi D, Hartl C, et al. Whole genome
753 sequencing in multiplex families reveals novel inherited and de novo genetic risk in autism. *bioRxiv*.
754 2018.

755