# Identification of moderate effect size genes in autism spectrum disorder through a novel gene pairing approach

Madison Caballero[1], F Kyle Satterstrom[3,4,5], Joseph D. Buxbaum[1,2,6,7,8,9], and Behrang Mahjani[1,2,6,7,10,11,12,*]

[1]Seaver Autism Center for Research and Treatment, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
[2]Department of Psychiatry, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
[3]Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA
[4]Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA
[5]Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, MA, USA
[6]Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
[7]The Mindich Child Health and Development Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
[8]Friedman Brain Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
[9]Department of Neuroscience, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
[10]Department of Artificial Intelligence and Human Health, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
[11]Department of Molecular Medicine and Surgery, Karolinska Institutet, Stockholm, Sweden.
[12]Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden.

* Corresponding author: behrang.mahjani@mssm.edu

## Abstract

Autism Spectrum Disorder (ASD) arises from complex genetic and environmental factors, with inherited genetic variation playing a substantial role. This study introduces a novel approach to uncover moderate effect size (MES) genes in ASD, which individually do not meet the ASD liability threshold but collectively contribute when paired with specific other MES genes. Analyzing 10,795 families from the SPARK dataset, we identified 97 MES genes forming 50 significant gene pairs, demonstrating a substantial association with ASD when considered in tandem, but not individually. Our method leverages familial inheritance patterns and statistical analyses, refined by comparisons against control cohorts, to elucidate these gene pairs' contribution to ASD liability. Furthermore, expression profile analyses of these genes in brain tissues underscore their relevance to ASD pathology. This study underscores the complexity of ASD's genetic landscape, suggesting that gene combinations, beyond high impact single-gene mutations, significantly contribute to the disorder's etiology and heterogeneity. Our findings pave the way for new avenues in understanding ASD's genetic underpinnings and developing targeted therapeutic strategies.

## Introduction

Autism Spectrum Disorder (ASD) is a complex neurodevelopmental disorder characterized by challenges in social interaction and communication, as well as a tendency toward restricted and

50    repetitive behaviors. Extensive research has underscored the complex interplay of genetic and
51    environmental factors contributing to its risk [1–6]. Diverse environmental influences, including
52    pregnancy complications, parental age, pollutants, and socioeconomic status, contribute to the
53    overall risk of ASD [7]. Nonetheless, there is consensus that genetic factors exert greater
54    influence on ASD's etiology, underscored by its substantial heritability exceeding 50% [2,8]. Large-
55    scale genomic studies have further revealed that ASD is a polygenic disorder influenced by a
56    broad spectrum of genetic variants rather than a single gene mutation [9]. These variants, ranging
57    in effect sizes, collectively shape the likelihood and clinical manifestations of ASD. The
58    threshold liability model for ASD, defined by an individual's unique combination of genetic and
59    environmental factors, has become the predominant way to understand and evaluate the
60    complex interplay of influences contributing to ASD risk and manifestation. Thus, the likelihood
61    of developing ASD is dependent on whether the summation of these factors crosses the
62    threshold for an ASD phenotype.

64    Genetic factors contributing to ASD liability range in rarity and effect size. Inherited common
65    variants (those with a population frequency of at least 1%) explain approximately 50% of
66    individual ASD liability [1,2,6]. Though common variants play a substantial role when considered
67    together, an individual common variant provides little contribution to individual ASD liability and
68    often offers little biological insight. The pressure of negative selection forces variants that
69    individually confer substantial heritability for ASD to be rare in the population; thus, studies of
70    rare variants – ranging from single nucleotide variants to large structural changes – have been
71    undertaken to help identify risk genes and developmental pathways. In recent years, research
72    efforts have focused on identifying *de novo* mutations in particular. *De novo* loss-of-function
73    variants are rare in the general population, particularly in genes which are haploinsufficient or
74    which exhibit evolutionary constraint. Genes in which such variants are overrepresented in
75    individuals with ASD can be identified as ASD risk genes, and the variants themselves can
76    confer extreme phenotypes that provide meaningful biological insight [3,4]. However, *de novo*
77    mutations only account for approximately 5% of ASD cases [2], and their sporadic origin renders
78    them less representative of ASD's overall heritability.

80    While there are an estimated 1100 genes implicated in ASD [10], only a few hundred have been
81    confidently identified [3,4,6,11]. Rare variant analyses aggregate variants across a gene and employ
82    advanced statistical models such as the Transmission and *De Novo* Association (TADA) test [12].
83    TADA increases statistical power by combining *de novo* mutations and case-control inherited
84    variants while incorporating per-gene mutation rates and gene constraint scores like the loss-of-
85    function tolerance [13,14]. Driven by variants of large effect, TADA has significantly improved the
86    identification of genes associated with ASD [4]. A notable set of genes, namely *SCN2A, SHANK3,*
87    *CHD8, ADNP*, and *SYNGAP1*, demonstrate especially compelling evidence, as deleterious
88    variants within these genes are effectively absent in control groups. However, this approach
89    loses power when considering genes where deleterious variants have a smaller effect size, as
90    these variants are more frequently found in non-ASD controls. As a result, there is a gap in our
91    understanding of genes where variants make moderate contributions to ASD liability.

93    Here we propose and implement a novel method for identifying moderate effect size (MES)
94    genes. We define MES genes as genes that individually do not surpass the ASD liability
95    threshold when disrupted and therefore have no significant association when tested in isolation
96    (**Fig 1A**). However, when deleterious variants in one gene co-occur with deleterious variants in
97    a specific second MES gene (always in concert with common genetic variation and
98    environmental factors), their cumulative contributions can cross the liability threshold. This does
99    not imply that gene pairs are epistatic or form dependent gene-gene interactions. Instead, we
100   find paired genes are predominantly independent with co-occurrence being strongly associated

101  with ASD. Thus, our strategy targets genes with meaningful but not singular impacts on ASD
102  liability.
103
104  The process of testing gene combinations poses a statistical challenge, particularly given the
105  vast number of potential gene pairs. For example, with 20,000 protein-coding genes, the
106  number of two-gene combinations is approximately 200 million, necessitating an extremely
107  stringent significance threshold. To overcome this challenge, we selected candidate pairs based
108  on familial inheritance patterns of ultra-rare deleterious variants, narrowing the focus to a
109  manageable subset of potential risk genes. Subsequent statistical analysis on this subset aimed
110  to detect gene pairs more prevalent in individuals with ASD compared to non-ASD controls.
111  After establishing statistical significance and validating findings, we then characterized 97
112  predicted MES genes to further support their roles in ASD etiology and heterogeneity.
113
114
115  **Results**
116
117  ***Identification of candidate MES genes inherited as pairs***
118
119  Our search for MES genes began by identifying pairs of candidate genes through family-based
120  inheritance patterns. This approach circumvents conducting approximately 200 million pairwise
121  comparisons among all protein-coding genes which would impose a harsh significance
122  threshold set by multiple testing corrections. Our primary objective was to identify instances
123  where inherited variants in two distinct genes were overrepresented in offspring with ASD, with
124  non-ASD offspring potentially inheriting either variant, but not both. This streamlined approach
125  enabled us to focus on candidate gene pairs exhibiting the most promising associations with
126  ASD.
127
128  The first and second steps of our candidate screen focused on the identification of familial
129  structures wherein ultra-rare (frequency <0.1%) deleterious variants in two genes were inherited
130  by all offspring with ASD within a given family (**Fig 2; Fig 1B**). In this context, the term 'A*'
131  designates an ultra-rare deleterious genetic alteration in a gene present in one parent, with the
132  potential for transmission to the offspring. The second step searched for mirrored inheritance
133  patterns involving an ultra-rare deleterious variant in a different gene, denoted as 'B*'. In such
134  instances, 'B*' originates from the parent who did not transmit 'A*', thereby maintaining genetic
135  independence.  We then filtered to instances where all ASD offspring that share the same
136  parents exhibited 'A*B*' inheritance, while none of their non-ASD siblings carried both alleles. In
137  different families, the specific ultra-rare deleterious variants in genes were predominantly
138  distinct. Additionally, as there were instances of multiple 'A*' or 'B*' variants within the same
139  gene and family, we used the count of unique families rather than unique sites for identifying
140  'A*B*' pairs. We initially analyzed 10,795 families from the SPARK[15] iWES v1 dataset ("SPARK
141  v1", comprising waves WES1-WES4), each with both parents and at least one ASD offspring
142  genotyped with whole exome sequencing data (WES; see **Methods**). For deleterious ultra-rare
143  variants (gnomAD[16] v3 non-neuro allele frequency of <0.1%), families on average exhibited an
144  'A*' inheritance pattern for 94 different genes.
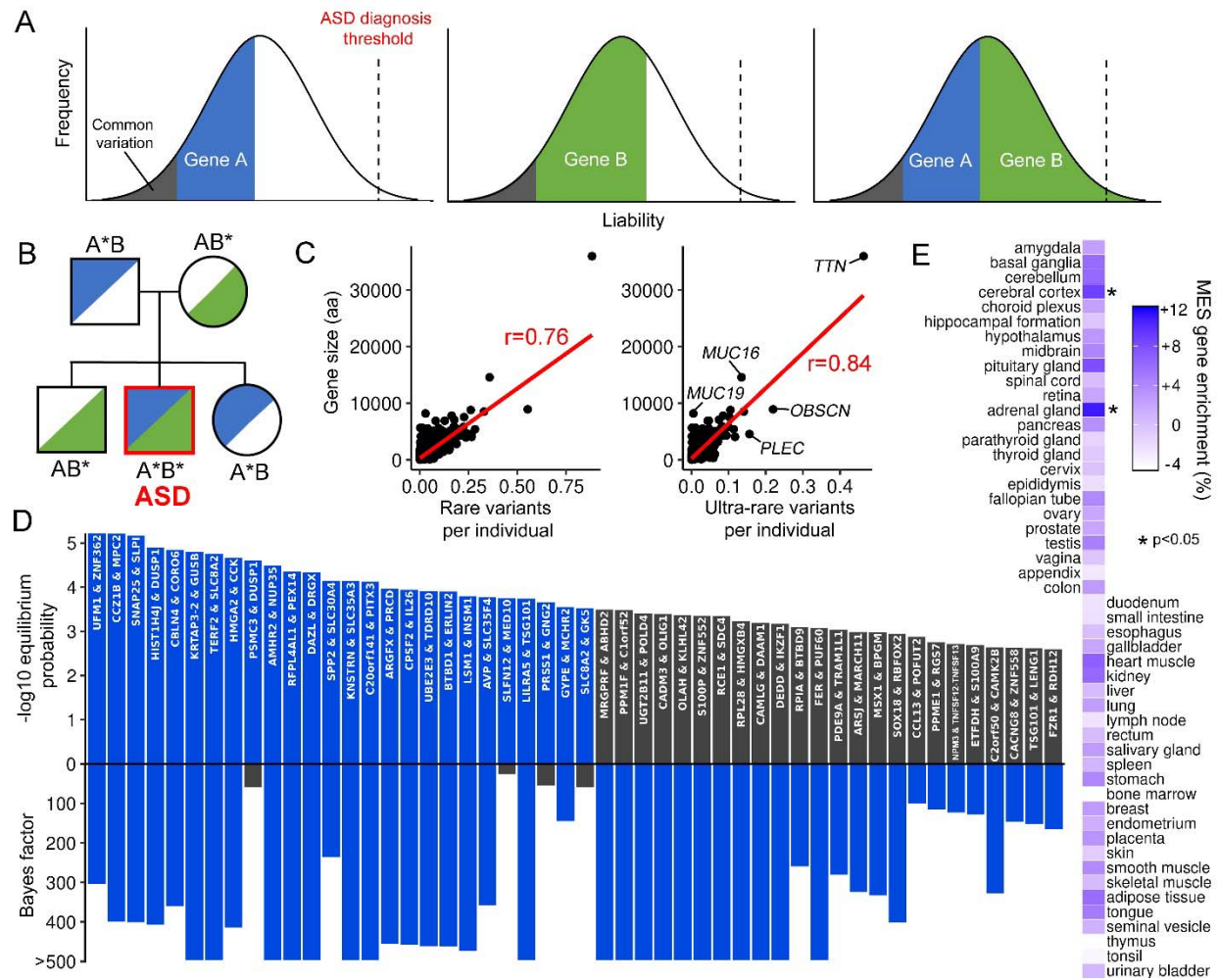145

146



147
148
149 **Figure 1**. **Discovery and enrichment of MES genes.** (A) Hypothesized role of individual MES genes
150 and their cumulative influence on the genetic liability of ASD in an individual. (B) Example inheritance
151 pattern illustrating ultra-rare deleterious variants in two genes within the parents (depicted in blue or
152 green) that are exclusively jointly inherited by offspring with ASD. This exemplifies a potential genetic
153 mechanism contributing to ASD susceptibility. (C) Mean number of rare (<1%) or ultra-rare (<0.1%)
154 deleterious variants in SPARK parents versus gene size. The red line represents linear regression.
155 Without *TTN*, correlation is 0.71 and 0.80, respectively. (D) Equilibrium probability and Bayes factors for
156 the predicted MES gene pairs in ASD. Bars in blue denote an equilibrium probability <3.03x10$^{-4}$ or a
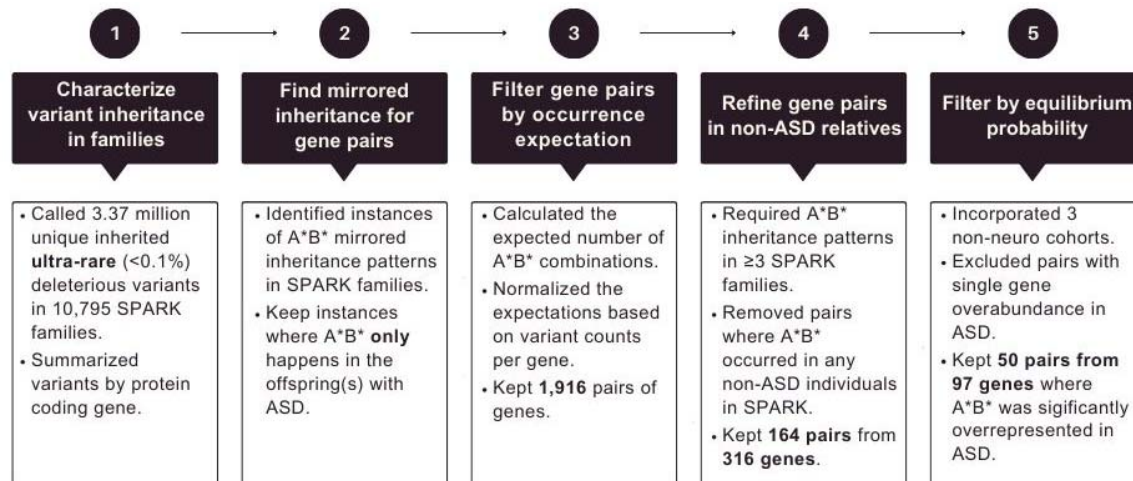157 Bayes factor >100. (E) Enrichment of predicted MES genes per tissue relative to all protein-coding genes.
158
159
160

**Figure 2**. **Flowchart of steps involved in MES gene discovery.**

In the third step, to identify promising gene pairs for further analysis, we assessed the likelihood of observing specific patterns of mirrored 'B*' inheritance within families displaying 'A*' inheritance. This entailed calculating the likelihood of ultra-rare deleterious 'B*' variants manifesting in a second gene, contingent upon our knowledge of the prevalence of such variants in the population. For instance, let us consider 'GeneA' which features 'A*' inheritance in seven families, and 'GeneB,' another gene with 'B*' inheritance in 20 families. Among these, four families exhibited the aforementioned mirrored 'A*B*' inheritance pattern. Utilizing the binomial distribution and adjusting by the frequency of ultra-rare deleterious variants per gene in the SPARK parent population, we determined (1) the expected frequency of observing this four-in-seven and, conversely, four-in-twenty arrangement, and (2) the relative likelihood that the genes observed were 'GeneA' and 'GeneB' (see **Methods**). This latter step is essential as genes with longer coding sequences generally have a greater frequency of ultra-rare or rarer deleterious variants (Pearson's $r$: 0.84; 0.76 for rare or rarer variants; **Fig1C**). To illustrate, 'A*' inheritance patterns within the large gene *TTN* were observed in 32.0% of all SPARK families. Consequently, the occurrence of a mirrored 'A*B*' inheritance pattern with *TTN* is considerably more anticipated than with *SNAP25*, a small gene where 'A*' inheritance patterns were present in only 0.048% of the SPARK families.

We retained 1,916 pairs of genes where the normalized expected frequency was less than one in both 'A*' and 'B*' configurations (**Table S1**). Additionally, in the fourth step, for each candidate pair, we imposed the condition that at least three families must exhibit the 'A*B*' inheritance pattern. To refine candidate pairs, we required that there were no occurrences of 'A*B*' in the remaining non-ASD individuals within the broader SPARK dataset, which comprised 36,323 individuals sequenced with WES. This filtering process resulted in a final count of 164 candidate pairs from 316 unique protein coding genes (**Table S1**).

Of note, we tested additional variant types and frequencies, including missense variants and PTVs of rare (<1%) frequency and analyses restricted to only PTVs (see **Methods**). We selected PTVs and missense variants of ultra-rare or rarer frequency for the following analyses as it produced the greatest number of candidate pairs after SPARK family filtering.

198    ***Statistical assessment of MES genes and incorporating non-neuro cohorts***
199
200    Though SPARK includes individuals without ASD, all are within families where at least one
201    individual has ASD. This makes these individuals an inaccurate representation of the general
202    population. Therefore, to perform statistical analyses below that incorporate the population
203    frequency of ultra-rare deleterious variants by gene, we included individuals without ASD or
204    other known mental, behavioral, or neurodevelopmental disorders. This included 3,202
205    individuals from the 1000 genomes project (1kGP)[17,18], 156,550 individuals from the All of Us
206    (v7) consortium, and 16,586 individuals from the BioMe Biobank. As with SPARK, we identified
207    ultra-rare deleterious variants in the 316 candidate MES genes in these cohorts. Lacking familial
208    information for most individuals, we could not determine with certainty whether any variant was
209    inherited or *de novo*. However, given that there are approximately 70 *de novo* mutations per
210    genome per generation[19], we can assume that variants observed in these cohorts are almost
211    always inherited.
212
213    To ensure comparability of detected variants in the 1kGP, All of Us, and BioMe cohorts as
214    controls, we examined the frequency of ultra-rare deleterious variants within the candidate MES
215    genes across these cohorts. Positive correlations were observed among all cohorts, with a more
216    pronounced correlation in cohorts characterized by larger sample sizes (**Fig S1**). Notably, non-
217    ASD individuals in the SPARK cohort displayed a slightly higher burden of ultra-rare deleterious
218    variants per gene compared to individuals in the other cohorts. On average, candidate genes in
219    the SPARK cohort contained 0.02%, 0.12%, and 0.18% more variants than in the All of Us,
220    BioMe, and 1kGP cohorts, respectively. Across all cohorts, SPARK showed a notable and
221    consistent enrichment for ultra-rare deleterious variants in *UGT2B11*, *AFAP1L2*, *C20orf141,* and
222    *OLAH* (**Fig S1**). However, these genes also contained relatively more synonymous variants in
223    the SPARK cohort, suggesting their enrichment for deleterious variants is not biologically
224    meaningful. Of note, variants in Z*NF717* were highly enriched in 1kGP compared to the other
225    cohorts (5.74-8.78% more). However, this gene's role in cellular proliferation and the
226    lymphoblastoid cell line source of 1kGP suggest these variants may be somatic and under
227    positive selection.
228
229    Among the 212,661 individuals without ASD, 'A*B*' was conspicuously absent in 26 out of the
230    164 candidate pairs and occurred in only two or fewer individuals in 83 candidate pairs. Across
231    all candidate pairs, the incidence of 'A*B*' was consistently higher among individuals with ASD
232    than their non-ASD counterparts. Notably, in 71 pairs, the frequency of 'A*B*' was more than 10
233    times greater in individuals with ASD. Together, this illustrates a pronounced disequilibrium
234    between individuals with ASD and those without ASD.
235
236    Our above search (encompassing steps one through four in **Fig 2**) aimed to identify potential
237    MES genes for in-depth analysis, mitigating the influence of large genes with a high frequency
238    of ultra-rare deleterious variants. In the fifth and final step, our objective was to assess the
239    statistical significance of the skew in MES gene pairs in relation to ASD. We hypothesized that
240    the concurrent inheritance of ultra-rare deleterious variants in two paired MES genes was
241    strongly associated with ASD development. We anticipated that the prevalence of 'A*B*' in the
242    population would adhere to an expected equilibrium, being the product of 'A*' and 'B*'
243    frequencies in the population. Their significant enrichment in individuals with ASD is presumably
244    the result of their involvement in ASD etiology. In contrast, our null hypothesis is that the
245    overrepresentation of 'A*B*' in individuals with ASD, contingent upon the individual frequencies
246    of 'A*' and 'B*' in the population, occurred by chance. Furthermore, when examining the
247    frequency of 'A*' or 'B*' in isolation, we required that neither gene may be significantly
248    overrepresented (p<0.05) in the individuals with ASD. Therefore, with respect to ASD, we

249    required only 'A*B*' to be in strong disequilibrium.
250
251    To illustrate, consider ultra-rare deleterious variants in the genes *SNAP25* (as 'A*')  and *SLPI* (as
252    'B*'), which are present in 0.056% and 0.196% of individuals, respectively. Notably, these
253    frequencies do not show significant differences when queried in individuals with or without ASD
254    (0.053% versus 0.057% for *SNAP25* and 0.220% versus 0.192% for *SLPI,* respectively). At
255    equilibrium, 'A*B*' should occur at a frequency determined by the product of their individual
256    frequencies. For *SNAP25* and *SLPI,* 'A*B*' occurs in three of 34,164 individuals with ASD and in
257    none of the 212,661 individuals without ASD. Using the cumulative density function of the
258    binomial distribution, we calculated the probability of observing at least three individuals with
259    ASD having 'A*B*' and none without ASD to be $6.88 \times 10^{-6}$. We calculated this equilibrium
260    probability for all 164 paired candidate MES genes. For 33 pairs, 'A*B*' was significantly
261    overrepresented in individuals with ASD ($p < 3.03 \times 10^{-4}$, corrected for multiple testing). We
262    discarded seven candidate pairs where an individual gene in isolation was significantly
263    overrepresented in individuals with ASD.
264
265    We additionally calculated the Bayes factor for each pair. As described above, our null
266    hypothesis posits that the increased occurrence of 'A*B*' in individuals with ASD is a random
267    event occurring at the equilibrium probability observed above. Our alternate hypothesis is that
268    'A*B*' is enriched among individuals with ASD precisely because it results in an ASD
269    phenotype. Therefore 'A*B*' would be more at equilibrium when not considering ASD status.
270    Using the example of *SNAP25* and *SLPI*, the Bayes factor was 400, strongly supporting the
271    alternate hypothesis. Interestingly, one of the largest Bayes factors was 5,807 for the genes
272    *PITX3,* a gene with a known role in dopaminergic neuron differentiation, and *C20orf141*, a gene
273    of unknown function. Across the 164 paired candidate MES genes, 46 had a Bayes factor
274    greater than 100.
275
276    We retained 50 candidate gene pairs (97 unique genes) where the equilibrium probability was
277    $< 3.03 \times 10^{-4}$ or the Bayes factor was >100 (**Fig 1D; Table S1**). Importantly, the genes in pairs
278    showed statistically significant associations to ASD but not when tested in isolation. From here,
279    we will refer to these as 'predicted MES genes'.
280
281    According to the recent and most comprehensive screens for ASD genes[4], the mean ASD FDR
282    of the predicted MES genes was 0.804, demonstrating that single MES genes by themselves
283    are not generally associated with ASD. Only one gene, *DEDD*, had an ASD FDR <0.05. Six
284    genes (*PUF60, SNAP25, PSMC3, CAMK2B,* and *HIST1H4J*) were associated with
285    neurodevelopmental disorders (FDR <0.05) though not for ASD. Similarly, the 97 MES genes
286    were also not associated with ASD in other major studies [3,10,11,20] and are therefore not
287    considered consensus ASD-associated genes.
288
289    The median probability of being loss-of-function intolerant (pLI) score of the predicted MES
290    genes was 0.029 and only 15 genes were >0.9, the common threshold for extreme
291    haploinsufficiency. When breaking up genes by probability of pLI decile, MES genes inhabited
292    more moderate deciles than large effect size genes characterized in the most recent
293    comprehensive assessment of ASD genes [4] (**Fig S2**). This was also the same result for loss-of-
294    function observed over expected upper bound fraction (LOEUF) scores (**Fig S2**). This suggests
295    the predicted MES genes are not strongly constrained and are more haplosufficient. This is
296    again consistent with our model of MES genes where each gene alone is not evidently
297    associated with ASD. Importantly, deleterious variants in both copies of predicted MES genes
298    are not present nor assessed in this study. Loss-of-function in both copies of a MES gene could
299    produce a separate pathological phenotype. Nonetheless, we present 97 predicted MES genes

300    that when disrupted in concert with a second MES gene are strongly associated with ASD.
301
302
303    ***Validations of MES genes***
304
305    To mitigate the challenge of a highly stringent multiple testing correction threshold inherent in
306    comparing every protein-coding gene, we streamlined our analysis to focus on candidate gene
307    pairs. This approach, while efficient, raises concerns of bias due to selective inference, as it
308    prioritizes pairs of genes with a pre-established likelihood of association. To address potential
309    bias, we implemented two validation strategies: (1) conducting a replication study using
310    synonymous variants as a neutral benchmark, and (2) incorporating additional ASD samples.
311
312    Our initial validation utilized ultra-rare synonymous variants, which, due to their non-functional
313    impact, were anticipated to yield fewer or less significant associations with ASD. This analysis
314    was conducted identically to that for deleterious variants. We found a similarly strong correlation
315    between the quantity of ultra-rare synonymous variants and CDS length (Pearson r: 0.83), as
316    well as between the counts of ultra-rare synonymous and deleterious variants (r: 0.82; **Fig S3 A,**
317    **B**). Moreover, the variability in the frequency of ultra-rare synonymous variants across different
318    cohorts paralleled that observed for missense variants (**Fig S3C**), affirming that synonymous
319    variants replicate our analyses under the same conditions. Utilizing ultra-rare synonymous
320    variants, we identified merely five gene pairs (10 unique genes; **Table S1**), in contrast to the 50
321    pairs identified with deleterious variants. Together, this suggests a potential false-positive rate of
322    approximately 10%.
323
324    Further validation involved a new cohort of 36,997 samples from SPARK iWES v2 ("SPARK v2",
325    adding WES5), including 10,128 ASD individuals. Applying the same methodological
326    framework, we first reassessed the equilibrium tests using only the new SPARK samples.
327    However, only eight of the original 50 gene pairs could be validated as we required A*B* in at
328    least one individual with ASD. Among the eight, seven pairs met the criteria of an equilibrium
329    probability p<0.05 or a Bayes Factor >10, adjusted for the smaller independent dataset (**Fig**
330    **S4A**). None of the five pairs identified with ultra-rare synonymous variants could be validated. In
331    addition, by combining data from both the SPARK v1 and v2 cohorts, we found 38 of the original
332    50 pairs maintained significance (equilibrium probability <$3.03\times10^{-4}$ or a Bayes factor >100;
333    **Table S1; Fig S4B**). By contrast, only two synonymous variant pairs retained significance post-
334    combination, none surpassing a significance level of $1.2\times10^{-4}$, a stark contrast to the maximum
335    significance of $5.9\times10^{-8}$ observed in the primary analysis with ultra-rare deleterious variants
336    (**Table S1**).
337
338    To summarize, our methodology efficiently identifies MES genes while maintaining a low false-
339    positive rate. Of the initial 50 gene pairs, 38 (74 unique genes) were corroborated by integrating
340    the SPARK v2 dataset, with seven out of eight subjected to independent validation using
341    SPARK v2. Subsequent analyses will delve into the 97 unique genes within these pairs,
342    focusing on their expression patterns and phenotypes among carriers as further evidence of
343    their involvement in ASD.
344
345
346    ***Enrichment of MES gene expression in the brain***
347
348    To further our identification of MES genes implicated in ASD, we conducted an analysis of their
349    tissue expression profiles. While ASD etiology encompasses factors beyond neural influences,

350   such as those associated with the gut microbiome [21], it is imperative that genes contributing to
351   ASD development exhibit expression in the brain. Accordingly, we hypothesized that our
352   predicted MES genes would exhibit an enrichment for expression in brain tissues. Utilizing
353   consensus normalized expression values (nTPM) from the Human Protein Atlas (v23)[22], and
354   encompassing 50 tissues, we observed that 77.0% of the 20,151 protein coding genes
355   demonstrated detectable expression (nTPM > 1) in at least one brain tissue, with 70.5%
356   exhibiting expression in the cerebral cortex. Strikingly, among the 97 predicted MES genes,
357   85.6% displayed detectable expression in at least one brain tissue, and 79.4% specifically in the
358   cerebral cortex. To assess the statistical significance of these findings, we performed random
359   gene sampling and quantified per-tissue expression, revealing a significant overrepresentation
360   of predicted MES genes in the cerebral cortex (p=0.022) and the adrenal glands (p=0.009) (**Fig
361   1E; Fig S5A**). The predicted MES genes did not exhibit significant underrepresentation in any
362   tissues. Among the 74 unique genes that maintained significance with the addition of SPARK
363   v2, we observed the same cerebral cortex and adrenal gland enrichment. Furthermore, we did
364   not observe that the predicted MES genes were expressed at a significantly higher level than
365   other cerebral cortex or adrenal expressed genes. Nonetheless, the observed enrichment of
366   MES gene expression in the cerebral cortex aligns with our initial predictions. The enrichment in
367   the adrenal glands was unexpected but not unprecedented as increases in androgen and
368   cortisol levels, which are both secreted by the adrenal glands, are associated with ASD [23,24].

370   To conduct a more in-depth investigation, we queried gene expression within 193 subregions of
371   the brain (**Fig S5B**). Notably, we identified the most pronounced enrichment of the 97 predicted
372   MES genes in the ventral tegmental area and central amygdala, although the statistical
373   significance fell just below the conventional threshold (p=0.08 and p=0.07, respectively). Again,
374   the same pattern was observed with the 74 unique genes that maintained significance with
375   SPARK v2. It is noteworthy that these particular brain regions exhibited the highest proportion of
376   MES genes with detectable expression (83.5% for both regions), surpassing the expression
377   levels observed in other brain regions or any alternative tissues that we analyzed. The 97
378   predicted MES genes are not canonically associated with ASD because they do not show
379   strong associations when tested alone. Yet, their expression is most pronounced in a region of
380   the brain linked to psychiatric disorders including ASD [25,26]. Together, we showed the predicted
381   MES genes are characterized by expression in neural tissues, suggesting their disruptions may
382   result in psychiatric phenotypes.

386   ***Phenotypes of non-ASD carriers of MES genes***

387   While MES genes were identified in pairs, the nature of their contribution to ASD remained
388   ambiguous – whether each gene within a pair independently influenced ASD or, conversely, if
389   their impact stemmed from interdependent gene-gene (GxG) interactions. In the scenario of
390   independent contributions, the disruption of a single MES gene may manifest mild ASD-like
391   phenotypes, while the disruption of two MES genes could result in an additive or synergistic
392   effect. In contrast, within the context of dependent GxG interactions, the disruption of a single
393   MES gene yields no discernible phenotype. To observe if any predicted MES pairs show
394   independent contributions, we queried the phenotypes and questionnaire responses of non-
395   ASD individuals within the SPARK v1 cohort. Per our MES gene discovery methodology, non-
396   ASD individuals may have an ultra-rare deleterious variant in either gene (herein simply referred
397   to as 'carriers') of an MES gene pair – but never both.

398   We first examined responses from the Social Communication Questionnaire (SCQ), completed

399  for adolescent offspring with and without ASD. The SCQ, comprising 40 yes-or-no questions
400  designed for ASD screening, was available for 3,831 non-ASD offspring in our study. We
401  focused on 68 out of the 97 predicted MES genes that exhibited a minimum of five non-ASD
402  MES gene carriers with SCQ responses. For each gene and each SCQ question, we employed
403  a $\chi^2$ test to compare responses between carriers and non-carriers among the non-ASD
404  offspring. We found 31 out of the 68 tested genes exhibited a significant ($p<0.05$) skew for at
405  least one SCQ question. For 15 genes, more than one question exhibited significantly increased
406  ASD-like responses, with a maximum of six questions per gene. In total, 65 questions
407  demonstrated significant skew, with 63 revealing a higher frequency of ASD-like responses in
408  carriers compared to non-carriers (**Fig S6**). Our results remained consistent when comparing
409  carriers to the 3,007 non-ASD offspring devoid of ultra-rare deleterious variants in any of the 97
410  predicted MES genes (65 of 70 significant questions indicated carriers with heightened ASD-like
411  responses).

412  In addition to examining individual SCQ questions, we assessed the total SCQ score. No
413  discernible elevation in total scores was observed among carriers of any MES genes. The most
414  substantial mean score difference between non-ASD carriers and non-carriers was less than
415  two points. Nonetheless, our study demonstrated that non-ASD carriers of MES genes may
416  exhibit significant differences in responses to specific SCQ questions, but not in overall scores.
417  Nevertheless, our results suggest that individual MES genes can contribute to a mild social
418  deficiency phenotype.

419  Though ASD-related phenotype measurements were less abundant for non-ASD individuals, we
420  were able to evaluate the presence of other psychiatric or developmental conditions through the
421  basic medical questionnaire. Across 28,883 non-ASD individuals and 29 questions on the
422  presence of developmental or psychiatric conditions (e.g., schizophrenia, obsessive-compulsive
423  disorder, sleep disorders, motor delays), we found 47 of 97 tested MES genes exhibited a
424  significant ($p<0.05$) skew for at least one condition (**Fig S7**). In 83.6% of cases, carriers of the
425  MES gene showed increased incidence of the condition. For example, 15.5% of non-ASD
426  carriers of ultra-rare deleterious variants in *PITX3* reported obsessive-compulsive disorder
427  (OCD) compared to 4.3% of non-carriers ($p=0.004$). Of note, the minority of cases where MES
428  gene carriers had significantly lower incidence occurred exclusively for conditions that were
429  relatively common (>5% of non-ASD respondents) such as depression and social anxiety.
430  Again, these results were consistent when comparing to non-ASD offspring devoid of ultra-rare
431  deleterious variants in any of the 97 predicted MES genes. In summary, we found evidence that
432  many individual MES genes may contribute independently to specific ASD-related traits and
433  conditions in non-ASD carriers.

434

435  ***Phenotypes of ASD carriers of MES genes***

436  The heightened prevalence of psychiatric conditions in non-ASD carriers of individual MES
437  genes suggests a potential influence on the diverse presentation of ASD. It is conceivable that
438  MES genes contribute to both the overall liability of ASD and the liability of comorbid symptoms
439  associated with ASD. In the SPARK v1 cohort, 99.95% of individuals with ASD exhibit at least
440  one of the 29 comorbid disorders assessed. Notably, language delays are the most common,
441  affecting 51% of individuals, followed by attention-deficit/hyperactivity disorder (ADHD) at 38%.

442  Among 58 of the 97 MES genes analyzed, carriers with ASD displayed significantly altered
443  frequencies of comorbid disorders compared to ASD individuals without any MES genes (**Fig**

444   **3A**). For 31 MES genes, comorbid disorders were always more prevalent. For instance,
445   individuals with ASD carrying an ultra-rare deleterious variant in the gene *NPM3* demonstrated
446   a significantly higher incidence of motor delay (p=0.002), diagnosed sleep disorders (p=0.007),
447   separation anxiety (p=0.02), social communication disorder (p= 0.02), eating disorders (p=0.02),
448   and encopresis (p=0.03). Importantly, carriers of *NPM3* among non-ASD individuals also
449   exhibited a significantly higher incidence of learning disabilities (p=0.01) and separation anxiety
450   (p=0.04). In total, among the 31 MES genes associated with increased comorbidity frequency in
451   individuals with ASD, 14 were also linked to increased comorbidity frequency in non-ASD
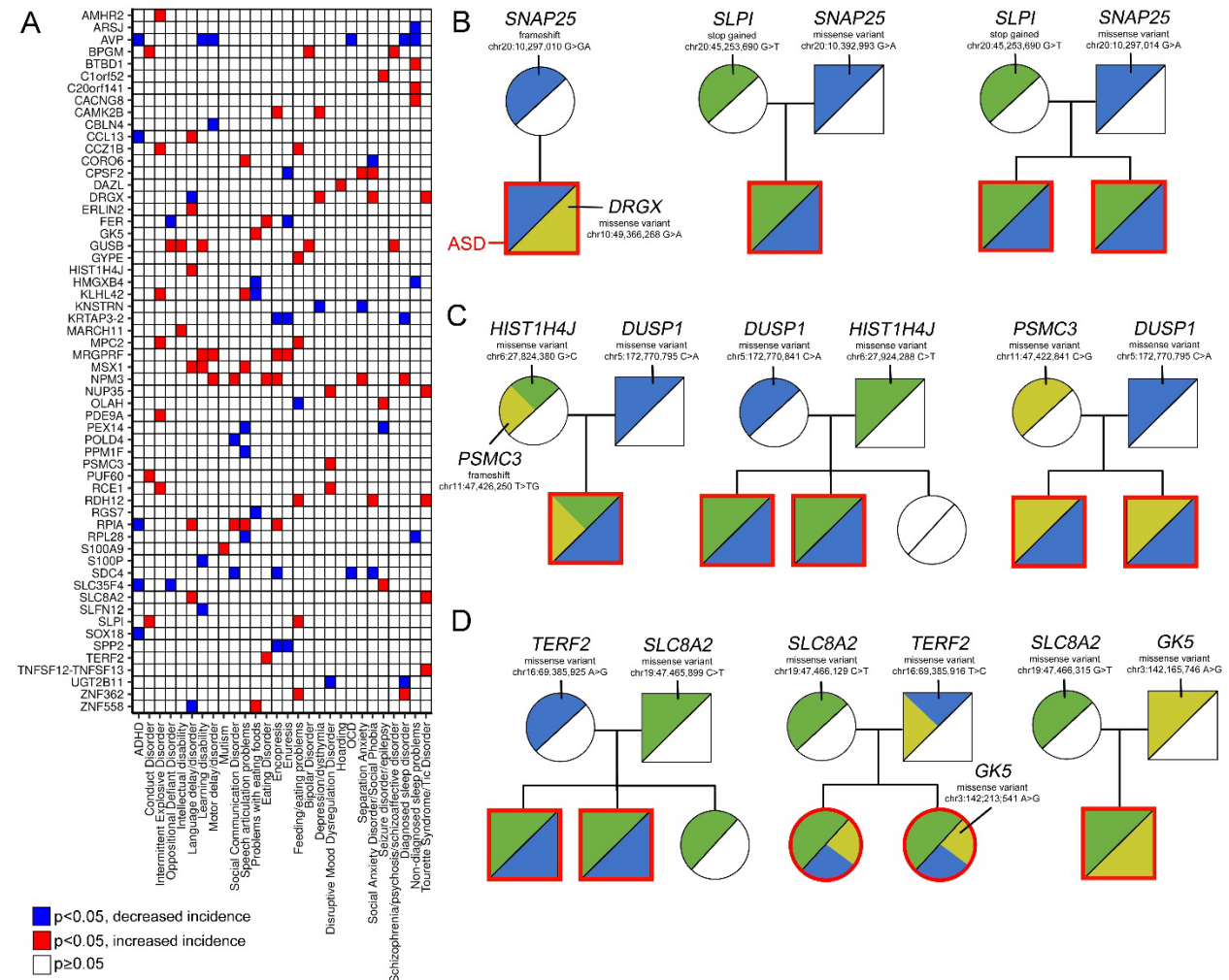452   individuals, though not always for the same comorbidity.

453



454

455   **Figure 3**. **Phenotypes and familial inheritance of MES genes.** (A) Differences in comorbidity frequency
456   of carriers of MES genes with ASD. Significant differences in ASD carrier versus non-carrier condition
457   frequencies are highlighted. Only genes with at least one significant comorbidity bias are shown. (B)
458   Example families from SPARK v1 that carry mutations in *SNAP25.* (C) As in panel B for *DUSP1.* (E) As in
459   panel B for *SLC8A2.*

460
461
462

463   Interestingly, for 17 MES genes, individuals with ASD exclusively exhibited significantly lower
464   incidence of comorbidities (**Fig 3A**). Carriers of ultra-rare deleterious variants in *AVP*, the gene
465   responsible for the production of the neuropeptide hormone arginine vasopressin, demonstrated
466   a reduced incidence of ADHD (p=5x10$^{-4}$), sleep disorders (p=0.004), OCD (p=0.03), motor
467   delays (p=0.03), and learning disabilities (p=0.048). Notably, in non-ASD individuals, carriers of
468   *AVP* exhibited a decreased incidence of depression (p=0.02) and a noteworthy, albeit not
469   statistically significant, decrease in ADHD and sleep disorders. As *AVP* is anxiogenic and
470   promotes the stress response [27], we tangentially hypothesize its disruption could reduce these
471   effects. In addition, a further 10 MES genes showed increased frequencies for one comorbidity
472   while decreased frequency for another. Carriers of *DRGX*, a gene involved in nervous system
473   development, showed elevated incidence of social anxiety disorder (p=0.03), depression
474   (p=0.049), and Tourette Syndrome (p=0.04) but a decreased incidence of language delay
475   (p=0.01). Thus, we find evidence MES genes can both contribute to the genetic liability of ASD
476   and its heterogeneity.
477
478
479   ***Co-expression and interactions of MES gene pairs***
480
481   Above, we demonstrated that 76 of the 97 identified MES genes exhibited discernible
482   phenotypic associations, suggesting most MES genes do not form dependent interactions when
483   contributing to ASD risk. To further test dependence, we inquired into the co-expression and
484   protein interaction patterns of these MES gene pairs compared to random MES gene pairings.
485
486   Employing the tissue-based expression data used above, we found that for 48 of 50 gene pairs,
487   the correlation in tissue co-expression did not significantly surpass that observed between any
488   two MES genes (**Fig S8A**). Notably, only the paired genes *CADM3* and *OLIG1* demonstrated a
489   significantly elevated co-expression correlation (p=0.028; Pearson's r= 0.46). However, carriers
490   of *CADM3* demonstrated altered SCQ responses, suggesting the pair's high co-expression is
491   not compelling evidence of dependence. Interestingly, the other gene pair with significant co-
492   expression, *CACNG8* and *ZNF558,* was anti-correlated (p=0.032; r= -0.38). Using the
493   expression profiles of 193 brain subregions, we also found three pairs demonstrated
494   significantly elevated co-expression and three pairs with significant anti-correlation (**Fig S9A**).
495   Interestingly, one of these significantly anti-correlated pairs was *CADM3* and *OLIG1* (p= 0.012*,*
496   r= -0.72) indicating their cross-tissue expression pattern was alike while their brain subregion
497   pattern was opposing. Together, we do not find evidence that paired MES genes are more co-
498   expressed than any two MES genes, suggesting the pairs predominantly contribute
499   independently to ASD risk.
500
501   We additionally clustered MES co-expression patterns and compared them to 185 high
502   confidence ASD genes[4]. We found that MES genes did not cluster separately from these
503   established ASD genes, suggesting that MES genes integrate into similar tissue-expression
504   pathways (**Fig S8B**). This same finding was also observed using the expression from 193 brain
505   subregions (**Fig S9B**).
506
507   We additionally investigated protein-protein interactions between paired MES genes using
508   STRING (v12)[28]. Of the 97 MES genes, 40 interacted with some second MES gene with
509   medium confidence (interaction score ≥0.4). However, only one MES pair, *HMGA2* and *CCK*,
510   showed evidence of interaction (score=0.44). This again supports our hypothesis that MES pair
511   genes do not form dependent interactions.
512
513

514
### Illustrative examples of MES gene pairs

516 Within this investigation, we have delineated the identification of 97 MES genes in 50 pairs,
517 exploring their connections to brain expression patterns and ASD-like phenotypes. To conclude,
518 we will highlight specific MES genes in these pairs, with particular attention given to familial
519 case studies. These cases not only underscore the significance of inheriting multiple MES
520 genes but also highlight the disparate contributions of MES genes to the genetic susceptibility of
521 ASD.

522 The concurrent inheritance of ultra-rare deleterious variants in *SNAP25* (Synaptosome
523 Associated Protein 25) and *SLPI* (Secretory Leukocyte Peptidase Inhibitor) emerged as a robust
524 association with ASD (equilibrium probability of $6.9x10^{-6}$; Bayes factor of 400). *SNAP25,* a
525 synaptic protein, is one of the highest expressed genes in the brain. While variants in *SNAP25*
526 have been linked to ADHD, schizophrenia, bipolar disorder, and neurodevelopmental disorders,
527 it often fails to be linked to ASD [3,4,29,30]. Our findings support this inconsistency as loss-of-
528 function disruptions of one copy of *SNAP25* did not always coincide with ASD, best exemplified
529 by an inherited frameshift variant in a non-ASD mother transmitted to her child with ASD (**Fig
530 3B**). This suggests that meeting the genetic liability for ASD among *SNAP25* carriers
531 necessitates additional contributions from common variations or another MES gene. In the case
532 of the mother-to-son transmitted frameshift in *SNAP25*, the ASD-afflicted child, but not the
533 mother, harbored an ultra-rare deleterious variant in another MES gene, *DRGX*, implicated in
534 nervous system development. *SLPI*, another candidate paired with *SNAP25* (**Fig 3B**), is integral
535 to the innate immune system and serves a protective role against inflammation in tissues [31].
536 Notably, this candidate assumes relevance given the well-established association between
537 inflammatory conditions and ASD [32,33].

539 Another noteworthy candidate MES gene is *DUSP1* (Dual Specificity Phosphatase 1)*,* a gene
540 implicated in the regulation of inflammation and nervous system development [34]. Our
541 investigation revealed a strong association of *DUSP1* with ASD when co-occurring with
542 *HIST1H4J* (H4 Clustered Histone 11; *H4C11*) or *PSMC3* (Proteasome 26S Subunit, ATPase 3),
543 with both genes linked to intellectual disability and neurodevelopmental delay [35,36]. Interestingly,
544 the concurrent inheritance of ultra-rare deleterious variants in *PSMC3* and *HISTH4J* was
545 observed in only two individuals—a non-ASD mother and her son with ASD (**Fig 3C**). Though a
546 limited occurrence, this raises the possibility that *PSMC3* and *HISTH4J* may have comparatively
547 lesser contributions to the genetic liability of ASD when compared to *DUSP1*. This inference is
548 reinforced by the familial context, as the son with ASD inherited a *DUSP1* variant from the father
549 in addition to the *PSMC3* and *HISTH4J* variants from the mother. This unique familial pattern
550 suggests a more prominent role for *DUSP1* in the manifestation of ASD within this specific
551 genetic context.

553 The last candidate MES genes we will illustrate are *SLC8A2* (Solute Carrier Family 8 Member
554 A2) which associates with ASD when paired with either *TERF2* (Telomeric Repeat Binding
555 Factor 2) or *GK5* (Glycerol Kinase 5) (**Fig 3D**). *SLC8A2*, a sodium-calcium exchanger, is highly
556 expressed in neurons where it plays a role in sodium and calcium ion homeostasis [37]. *TERF2*
557 plays an important role in telomere maintenance and neuronal differentiation [38] whereas *GK5* is
558 a broadly expressed gene that is part of glycerol metabolism. The concurrent inheritance of
559 ultra-rare deleterious variants in *TERF2* and *GK5* was observed in four individuals – two with
560 and two without ASD. Three of these individuals were in the same family where variants in both
561 genes were transmitted from a non-ASD father to his two daughters with ASD (**Fig 3D**).
562 However, as was observed above with *DUSP1, PSMC3,* and *HISTH4J,* the daughters with ASD

563    also inherited variants in *SLC8A2* from the mother. This case study again suggests that
564    *SLC8A2* has a greater contribution to the genetic liability of ASD than *TERF2* and *GK5.* In
565    summation, the interplay of our predicted MES genes underscores the polygenic landscape of
566    genetic factors contributing to ASD susceptibility and heterogeneity.
567
568
569

570    **Discussion**
571
572    This study presents a novel approach to understanding the complex genetic landscape of ASD
573    by identifying MES genes – genes that individually, when harboring a deleterious variant, fall
574    below the ASD liability threshold but, when paired with another MES gene, collectively
575    contribute to ASD risk. By analyzing familial inheritance patterns of ultra-rare deleterious
576    variants, we identified candidate pairs of genes further assessed in 255,883 individuals of which
577    31,183 have ASD. We revealed 97 predicted MES genes forming 50 pairs, with significant
578    associations to ASD when considered together but not individually. This study emphasized the
579    importance of gene combinations, shedding light on the nuanced interplay of genetic factors in
580    ASD susceptibility and providing insights into the disorder's heterogeneity. Additionally, we
581    explored the expression profiles of MES genes in the brain and examined the phenotypic effects
582    of carriers, offering a comprehensive perspective on the multifaceted genetic contributions to
583    ASD.
584
585    Our study focused on identifying a novel but limited set of paired MES genes with the most
586    pronounced associations to ASD. A primary limitation of this study is that we did not consider
587    common genetic variation or environmental conditions, assuming that the co-inheritance of
588    paired MES genes alone was sufficient to result in ASD. However, it is well-established that
589    common genetic variation and environmental factors substantially contribute to ASD liability.
590    Further analyses of MES or even large effect size genes may involve comparing carriers with
591    high or low polygenic risk scores, potentially revealing that single MES genes and a high burden
592    of common variation may be adequate to cross the diagnostic threshold. The incorporation of
593    common variation may also aid in uncovering even smaller effect size genes. Detecting
594    combinations of three or more genes using our methodology would necessitate a larger number
595    of families. However, identifying smaller effect size genes conditioned on a high or low
596    polygenic risk score remains an avenue for future exploration.
597
598    Our investigation sought evidence for whether identified MES genes act independently or
599    participate in dependent GxG interactions to influence ASD susceptibility. SCQ responses
600    among non-ASD carriers revealed a significant skew for specific questions, suggesting that
601    many individual MES genes may independently contribute to subtle social deficiencies.
602    Furthermore, altered frequencies of developmental or psychiatric conditions in both non-ASD
603    and ASD carriers supported the notion that MES genes can independently influence diverse
604    phenotypic outcomes. In total, 76 out of the 97 predicted MES genes exhibited skewed
605    frequencies in at least one SCQ response or the incidence of a developmental or psychiatric
606    condition. It is plausible that some MES gene pairs form dependent GxG interactions, implying
607    carriers of one gene exhibit no discernable differences from non-carriers. Notably, there were
608    instances where non-ASD individuals harbored ultra-rare deleterious variants in up to five MES
609    genes. MES gene burden in non-ASD individuals was not associated with sex at birth,
610    suggesting this is not a result of the female protective effect. Instead, in cases where non-ASD
611    individuals possessed more than two MES genes, the 21 genes without a skewed frequency in
612    SCQ or developmental/psychiatric conditions constituted an average of 90% of the multiple
613    MES genes per individual. If dependent GxG interactions describe this minority of genes, it

614 would allow non-ASD individuals to carry multiple unpaired MES genes without affecting their
615 phenotype. This nuanced understanding of the independent and collective actions of MES
616 genes provides insights into the intricate nature of their contributions to ASD susceptibility,
617 unraveling the complex interplay between genetics and phenotypic expression.
618
619 We hypothesized that MES genes would exhibit less conservation than larger effect size genes
620 which was indeed reflected by their moderate pLI deciles. In the context of ASD, genes with
621 large effect sizes typically manifest autosomal-dominant effects, leading to profound phenotypes
622 such as epilepsy, developmental delay, and intellectual disability, alongside ASD [3].
623 Consequently, *de novo* mutations in these genes are seldom transmitted to subsequent
624 generations. In contrast, loss-of-function mutations in MES genes tend to result in less severe
625 phenotypes if any, possibly due to redundant pathways or functional paralogs. As a result,
626 disruptions in these genes can be transmitted and still contribute to the heritability of ASD. The
627 high heritability of ASD implies that its key genetic factors may generally be subject to more
628 lenient evolutionary constraints. Variants in genes that tolerate altered function and expression
629 may follow similar trends. The ability of MES genes to tolerate mutations and still contribute to
630 ASD's heritability underscores the importance of understanding the interplay between genetic
631 variation, phenotypic outcomes, and evolutionary pressures in the context of ASD susceptibility.
632
633 The most challenging aspect of MES genes is that the disruption of a single gene is generally
634 inadequate to induce ASD, yet still contributes to its etiology and heterogeneity. This
635 necessitates a nuanced approach in determining how terms such as "causal genes," "risk
636 genes," or "pathogenic variants" are to be applied. In the context of monogenic diseases, the
637 causality of a gene or mutation is often straightforward and direct. However, the causality
638 associated with complex diseases like ASD entails a combination of multiple genetic variants
639 and environmental factors contributing to the progression of the condition, rather than merely
640 elevating the risk. Although a MES gene may be considered causal as it influences key
641 pathways in ASD development, it does not singularly determine the condition. Guidelines, such
642 as those established by the American College of Medical Genetics and Genomics (ACMG), play
643 a pivotal role in interpreting and categorizing genetic variants, ranging from pathogenic to
644 benign. ACMG incorporates an array of factors, including functional predictions, inheritance
645 patterns, and prior reports, to assign pathogenicity to genes or variants [39]. This approach is
646 highly effective at identifying strong gene drivers of human diseases but less applicable to
647 complex conditions where MES genes may play crucial roles. Variants in MES genes may be
648 assigned as benign because they are not associated with a disorder, yet carrier status could be
649 clinically crucial such as predicting ASD recurrence or the likelihood of comorbidities. On the
650 other hand, assigning pathogenicity to these variants or genes is also an inaccurate term. This
651 complexity highlights a critical gap in the current genetic interpretation frameworks. As
652 researchers further study complex disorders (e.g., ADHD, OCD, bipolar disorder), there will be a
653 growing need to adapt or develop new guidelines to more accurately address the nuances of
654 polygenic diseases.
655
656
657
658 **Methods**
659
660 ***Family-based search for MES genes in SPARK***
661
662 For our initial family-based screening of MES genes, we harnessed a dataset comprising 10,908
663 families sourced from SPARK [15] iWES v1 (including sequencing waves WES1-WES4), in which
664 both parents and at least one ASD offspring underwent sequencing via WES. All genotypes

665 were annotated with Ensembl Variant Effect Predictor (VEP) [40]. We excluded genotype calls
666 that did not adhere to Mendelian inheritance patterns within the family (accounting for an
667 average of 0.72% of variants), effectively eliminating *de novo* mutations. Furthermore, during
668 the initial screening phase, we omitted 113 families in which a proband exhibited a *de novo* rare
669 deleterious variant (a PTV or a missense variant with a PolyPhen-2 [41] score ≥0.5 or SIFT [42]
670 score ≤ 0.05; gnomAD v3 [16] non-neuro allele frequency of <1%) in the genes *CHD8, SCN2A,*
671 *SYNGAP1, SHANK3,* or *ADNP*. These particular genes have a reported zero FDR for
672 association with ASD from TADA [4], signifying a strong and monogenic-like association with
673 ASD. Consequently, they were deemed unsuitable for the identification of MES genes. *PTEN*
674 was not removed for the screen as deleterious variants occurred in controls.

676 We isolated ultra-rare deleterious variants in the families (a PTV or a missense variant with a
677 PolyPhen-2 score ≥0.5 or SIFT score ≤ 0.05; gnomAD v3 non-neuro allele frequency of <0.1%).
678 Though many more missense pathogenicity estimators are available [43], we opted for ones that
679 scored for protein function impact without emphasis on species conservation or gene constraint
680 (e.g., MPC[44]). Per our hypothesis, deleterious variants in single MES genes do not result in a
681 strong phenotype and may not show as robust evidence of negative selection. For the same
682 reason, we did not exclude genes based on LOEUF or pLI score. Additionally, we only included
683 heterozygous genotypes on the reasoning that (1) homozygous combinations of ultra-rare
684 deleterious variants are effectively unobserved, and (2) individuals homozygous for deleterious
685 variants in a MES gene may have a separate non-ASD phenotype.

687 We first identified familial structures characterized by the inheritance of a specific ultra-rare
688 deleterious variant (denoted as 'A*') by all offspring with ASD within the family, with 'A*' also
689 being present in one non-ASD parent. Subsequently, we identified mirrored patterns of
690 inheritance, focusing on a second variant in a distinct gene ('B*'). In such scenarios, 'B*' must
691 derive from the parent who did not transmit 'A*'. 'A*B*' must be present in all offspring with ASD
692 and never in non-ASD siblings. Importantly, our methodology accommodated the potential
693 existence of multiple 'A*' or 'B*' variants within the same gene in a single family. We quantified
694 these patterns based on the number of unique families, rather than unique sites.

696 We next assessed the likelihood of observing 'B*' inheritance patterns within families exhibiting
697 'A*' inheritance. For this, we conducted binomial samplings to estimate the number of times we
698 should observe any number of secondary independent gene hits:

$$N = \sum_{gene=1}^{n\_genes} B \sim Binom(A, M_{gene})$$

701 where B is the number of families with 'B*' inheritance, A is the number of families with 'A*'
702 inheritance, and $M_{gene}$ is the frequency of ultra-rare deleterious variants in a gene among the
703 21,570 non-ASD parents. This outputs N, the expected number of genes with B secondary hits
704 given A.

706 Drawing from the frequency of ultra-rare deleterious variants per gene, calculated using parental
707 data, our aim was to standardize N, representing the expected number of genes with secondary
708 B hits, given A. Notably, larger genes, containing a higher occurrence of ultra-rare deleterious
709 variants, might bear reduced significance when compared to smaller genes with fewer such
710 variants. In our approach, for each pair of candidate MES genes, we multiplied N by the relative
711 frequency of ultra-rare deleterious variants in the respective genes, designated as 'N_norm'.
712 The gene possessing the highest number of ultra-rare deleterious variants, *TTN*, would have an

713  adjustment of one with all other genes adjusted by factors less than one, contingent on their
714  relative frequency of ultra-rare deleterious variants to *TTN*. To illustrate, consider the case of
715  *SNAP25* and *SLPI*. Among the families analyzed, five exhibited the 'A*' inheritance pattern for
716  *SNAP25*, while 21 displayed the same for *SLPI*. Two families featured the mirrored 'A*B*'
717  inheritance pattern for both genes. According to binomial sampling, the occurrence of '5-given-2'
718  is expected approximately 15 times, whereas '21-given-2' should transpire approximately 196
719  times (before normalization). Due to their smaller gene sizes, *SNAP25* and *SLPI*
720  correspondingly exhibited 'N_norm' values of 0.061 and 0.174, indicating that the co-occurrence
721  of these two genes inherited together is relatively unexpected and thus included as a candidate
722  pair for further investigation. We kept pairs where 'N_norm' was less than one for both
723  configurations of gene pairs.
724  
725  

726  ### Additional variant types and frequencies
727  
728  In addition to PTVs and missense variants of ultra-rare or rarer frequency, we tested three other
729  stringencies: PTVs and missense variants of rare (1%) or rarer frequency, only PTVs of ultra-
730  rare (0.1%) and rarer frequency, and only PTVs of rare or rarer frequency. For PTVs and
731  missense variants of rare or rarer frequency, we retained 2,507 pairs of genes (2,146 total
732  unique genes) where the normalized expected frequency was less than one in both 'A*' and 'B*'
733  configurations. After filtering for non-ASD 'A*B*' in the SPARK v1 dataset, only 179 candidate
734  pairs remained (339 total unique genes). For PTVs of rare or rarer frequency, we retained 834
735  (1,336 total unique genes) candidate pairs of which only 83 (161 unique genes) remained after
736  SPARK filtering. For PTVs of ultra-rare or rarer frequency, we retained 692 candidate pairs
737  (1,177 total unique genes) of which only 63 (122 unique genes) remained after SPARK filtering.
738  
739  

740  ### Non-neuro cohorts
741  
742  We included variants from the 3,202 individuals from the 1000 genomes project (1kGP)[17,18].
743  SIFT scores, PolyPhen-2 scores, and allele frequencies were extracted from GnomAD
744  (v3.1.2)[45].
745  
746  Analyses for All of Us were performed in the workspace titled "Detecting the prevalence of ultra-
747  rare gene mutations." We selected samples with short read whole genome sequencing and
748  excluded individuals with a survey response diagnosis of ASD (diagnosed with, receiving
749  treatment for, or seeking treatment for ASD) or any disorder characterized under 'mental
750  disorder.' The final cohort included 156,550 individuals of which 50.5% were of a self-reported
751  non-white race and 57.0% were assigned female sex at birth. Due to the security constraints of
752  the All of Us workbench, VEP was performed with SNPeff [42] (4.3t) as opposed to Ensembl.
753  GnomAD v3 frequencies, Polyphen-2 scores, and SIFT scores were acquired from
754  dbNSFP(v4)[47]. 98.8% of missense variants had an available Polyphen-2 and SIFT score. For
755  high impact insertions and deletions that were not present in dbNSFP, we used the allele
756  frequency in the 156,550 individuals as a proxy. To avoid disseminating individual-level data, we
757  did not report the precise number of individuals with 'A*B*' for any gene pair (unless zero or
758  greater than 20). For any single candidate gene, ultra-rare deleterious variants were always
759  present in at least 20 individuals.
760  
761  We additionally included individuals from the BioMe Biobank maintained by the Mount Sinai
762  Health System. We isolated 16,586 individuals with genotype information derived from WES
763  and excluded any individuals with a diagnosed or self-reported mental, behavioral, or

764   neurodevelopmental disorders (ICD codes F01-F99). Of the final cohort, the mean age was 61.7
765   years, 56.2% were of self-reported male sex, and 56.9% were of self-reported non-white race.
766   For consistency, VEP was performed in an identical manner as for the All of Us cohort.
767
768

769   ### *Equilibrium test and Bayes factors*
770

771   We determined the frequency of ultra-rare deleterious variants per gene across all samples,
772   encompassing both ASD and non-ASD individuals. The probability of 'A*B*' was computed as
773   the product of the individual gene frequencies. Using the R function 'pbinom,' we designated this
774   product as the probability of success for 'A*B*.' Subsequently, we calculated (1) the probability
775   of observing the given number or a greater number of ASD individuals with 'A*B*', considering
776   the total number of individuals with ASD, and (2) the probability of observing the given number
777   or a lesser number of non-ASD individuals with 'A*B*' mutations, considering the total number of
778   individuals without ASD. The equilibrium probability was then determined as the product of
779   these two probabilities.
780

781   For the Bayes factor, the null hypothesis was the equilibrium probability, and the alternate
782   hypothesis was the probability of observing the given number or a greater number of individuals
783   regardless of ASD status with 'A*B*', considering the total number of individuals. The Bayes
784   factor was calculated as the ratio of alternate and null equilibrium probabilities.
785
786

787   ### *Validation analyses*
788

789   We procured ultra-rare synonymous variants across all cohorts, excluding those with a
790   functional impact (e.g., missense, or more significant effects) on another gene. Ultra-rare
791   synonymous variants were, on average, 1.83 times more frequent per gene than ultra-rare
792   deleterious variants, which underwent further filtration for deleteriousness via VEP. We
793   determined the ratio of ultra-rare synonymous to ultra-rare deleterious variants across all
794   protein-coding genes. As some genes presented markedly divergent ratios, we excluded those
795   beyond the >95th or <5th percentile thresholds. This adjustment predominantly impacted
796   smaller genes where variants were only observed in a few samples.
797

798   Using ultra-rare synonymous variants, the initial screening (steps 1-3 of **Fig 2**) identified 690
799   candidate gene pairs, significantly fewer than the 1,916 pairs identified with deleterious variants.
800   Subsequent analysis incorporating the complete SPARK v1 cohort (step 4 of **Fig 2**) further
801   narrowed the candidate pairs to 129 (258 unique genes).
802

803   In the SPARK v1 cohort, we noted outlier individuals with high counts of ultra-rare synonymous
804   variants. On average, the occurrence of ultra-rare synonymous variants was less than one per
805   individual per gene. Among 258 unique genes from 129 candidate pairs identified with
806   synonymous variants, four samples exhibited over 100 variants in single genes, 50-fold higher
807   than the mean (SP0122762, SP0379928, SP0149880, SP0064380; **Fig S3D**). The excess
808   synonymous variants in the four samples were unique and not cataloged in gnomAD v4, casting
809   doubts on their quality. Contrastingly, such an excess was not observed with ultra-rare
810   deleterious variants (**Fig S3D**). We subsequently removed these four samples from analyses.
811

812   For additional validation, we incorporated novel samples from SPARK iWES v2 (WES5). Quality
813   control was identical to that of SPARK v1.
814

815
### *Gene expression*
817

818 We used nTPM values derived from RNAseq from the Human Protein Atlas (v23)[22]. The 50
819 tissue consensus used transcriptomics from both the Human Protein Atlas and the GTEx
820 consortium [48]. Expression for the 193 brain subregions was also from the Human Protein Atlas.
821 To test if the 97 predicted MES genes were significantly enriched in a tissue or brain subregion,
822 we sampled genes from all protein coding genes without replacement in 1000 iterations. For
823 each iteration, we measured the fraction of sampled genes expressed (nTPM >1) in each tissue
824 or brain subregion. We used the Shapiro–Wilk test[49] (R command 'shapiro.test') to confirm
825 sampled statistics were normally distributed (p>0.05). We then used the cumulative density
826 function of the normal distribution to determine the significance of predicted MES gene
827 enrichment per tissue or brain subregion.
828

829 Using the same two RNAseq datasets described above, we calculated the correlation of nTMP
830 values across tissues and brain subregions. To assess the significance of MES pair co-
831 expression, we calculated the co-expression correlation of all possible MES gene pairs and then
832 calculated the p-value of real MES pairs using the R function 'pnorm'. Hierarchical clustering of
833 the co-expression of MES and other ASD genes was performed with the R function 'hclust'
834 using the 'ward.D2' method and plotted with 'ComplexHeatmap'[50].
835
836

837 ### *Phenotypes of MES carriers*
838

839 We utilized the Social Communication Questionnaire (SCQ) and basic medical questionnaire
840 administered in the SPARK v1 cohort. For the SCQ total score, we removed 167 non-ASD
841 siblings with a score at or above 13 (two standard deviations above the mean among non-ASD
842 respondents). For the basic medical screening, we excluded questions related to birth defects
843 (as they were uncommon) and environmental exposures (e.g., lead poisoning, fetal alcohol
844 syndrome, traumatic brain injury).
845
846

## Acknowledgements
848

865  research cohorts used in this research.
866
870
871
872  **Author contributions**
873
874  M.C. led the study design, conducted analyses, and authored the manuscript. B.M. contributed
875  to the study design and manuscript editing and, along with J.D.M., oversaw project
876  management. F.K.S. performed data generation and quality control related to the SPARK
877  cohort.
878
879
880  **Competing interests**
881
882  The authors declare no competing interests.
883
884
885  **Data availability**
886
887  Whole genome sequencing from 1kGP is freely available from The International Genome
888  Sample Resource (https://www.internationalgenome.org/). All other individual-level data used in
889  this study are not publicly available with access permission subject to review. The genetic and
890  phenotypic data for SPARK can be requested at SFARI Base
891  (https://www.sfari.org/resource/spark/). Access to the Sinai BioMe data can be requested
892  through The Charles Bronfman Institute for Personalized Medicine
893  (https://icahn.mssm.edu/research/ipm/programs/biome-biobank) and is subject to compliance
894  with the Mount Sinai Health System's data use agreements. All of Us controlled tier data is
895  available for registered institutions and cohort selection is detailed in the workbench titled
896  "Detecting the prevalence of ultra-rare gene mutations."
897
898
899  **Code availability**
900
901  Code and resources used in this study is available on GitHub at
902  (https://github.com/MahjaniLab/MES_Code). For any further inquiries or requests for code not
903  available in the repository, please contact the corresponding author.
904
905
906  **References**
907
908  1.  Klei, L. *et al.* Common genetic variants, acting additively, are a major source of risk for

909      autism. *Mol. Autism* **3**, 9 (2012).

910    2. Gaugler, T. *et al.* Most genetic risk for autism resides with common variation. *Nat. Genet.*

911       **46**, 881–885 (2014).

912    3. Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional

913       Changes in the Neurobiology of Autism. *Cell* **180**, 568-584.e23 (2020).

914    4. Fu, J. M. *et al.* Rare coding variation provides insight into the genetic architecture and

915       phenotypic context of autism. *Nat. Genet.* **54**, 1320–1331 (2022).

916    5. Bai, D. *et al.* Association of Genetic and Environmental Factors With Autism in a 5-Country

917       Cohort. *JAMA Psychiatry* **76**, 1035–1043 (2019).

918    6. Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum disorder.

919       *Nat. Genet.* **51**, 431–444 (2019).

920    7. Modabbernia, A., Velthorst, E. & Reichenberg, A. Environmental risk factors for autism: an

921       evidence-based review of systematic reviews and meta-analyses. *Mol. Autism* **8**, 13 (2017).

922    8. Sandin, S. *et al.* The Heritability of Autism Spectrum Disorder. *JAMA* **318**, 1182–1184

923       (2017).

924    9. Rylaarsdam, L. & Guemez-Gamboa, A. Genetic Causes and Modifiers of Autism Spectrum

925       Disorder. *Front. Cell. Neurosci.* **13**, 385 (2019).

926    10. De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in autism.

927       *Nature* **515**, 209–215 (2014).

928    11. Sanders, S. J. *et al.* Insights into Autism Spectrum Disorder Genomic Architecture and

929       Biology from 71 Risk Loci. *Neuron* **87**, 1215–1233 (2015).

930    12. He, X. *et al.* Integrated Model of De Novo and Inherited Genetic Variants Yields Greater

931       Power to Identify Risk Genes. *PLOS Genet.* **9**, e1003671 (2013).

932    13. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**,

933       285–291 (2016).

934    14. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in

935       141,456 humans. *Nature* **581**, 434–443 (2020).

936    15. SPARK: A US Cohort of 50,000 Families to Accelerate Autism Research. *Neuron* **97**, 488–

937        493 (2018).

938    16. Siwei Chen *et al.* A genome-wide mutational constraint map quantified from variation in

939        76,156 human genomes. *bioRxiv* 2022.03.20.485034 (2022)

940        doi:10.1101/2022.03.20.485034.

941    17. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).

942    18. Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the expanded 1000

943        Genomes Project cohort including 602 trios. *Cell* **185**, 3426-3440.e19 (2022).

944    19. Sasani, T. A. *et al.* Large, three-generation human families reveal post-zygotic mosaicism

945        and variability in germline mutation accumulation. *eLife* **8**, e46922 (2019).

946    20. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum disorder.

947        *Nature* **515**, 216–221 (2014).

948    21. Taniya, M. A. *et al.* Role of Gut Microbiome in Autism Spectrum Disorder and Its

949        Therapeutic Regulation. *Front. Cell. Infect. Microbiol.* **12**, 915701 (2022).

950    22. Uhlén, M. *et al.* Tissue-based map of the human proteome. *Science* **347**, 1260419 (2015).

951    23. Gao, J. *et al.* Alteration of peripheral cortisol and autism spectrum disorder: A meta-

952        analysis. *Front. Psychiatry* **13**, 928188 (2022).

953    24. Al-Zaid, F. S., Alhader, A. F. A. & Al-Ayadhi, L. Y. A potential role for the adrenal gland in

954        autism. *Sci. Rep.* **11**, 17743 (2021).

955    25. Cai, J. & Tong, Q. Anatomy and Function of Ventral Tegmental Area Glutamate Neurons.

956        *Front. Neural Circuits* **16**, 867053 (2022).

957    26. Seguin, D. *et al.* Amygdala subnuclei volumes and anxiety behaviors in children and

958        adolescents with autism spectrum disorder, attention deficit hyperactivity disorder, and

959        obsessive–compulsive disorder. *Hum. Brain Mapp.* **43**, 4805–4816 (2022).

960    27. Beurel, E. & Nemeroff, C. B. Interaction of stress, corticotropin-releasing factor, arginine

961        vasopressin and behaviour. *Curr. Top. Behav. Neurosci.* **18**, 67–80 (2014).

962   28. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased

963   coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic*

964   *Acids Res.* **47**, D607–D613 (2019).

965   29. Antonucci, F. *et al.* SNAP-25, a Known Presynaptic Protein with Emerging Postsynaptic

966   Functions. *Front. Synaptic Neurosci.* **8**, 7 (2016).

967   30. Corradini, I., Verderio, C., Sala, M., Wilson, M. C. & Matteoli, M. SNAP-25 IN

968   NEUROPSYCHIATRIC DISORDERS. *Ann. N. Y. Acad. Sci.* **1152**, 93–99 (2009).

969   31. Doumas, S., Kolokotronis, A. & Stefanopoulos, P. Anti-Inflammatory and Antimicrobial Roles

970   of Secretory Leukocyte Protease Inhibitor. *Infect. Immun.* **73**, 1271–1274 (2005).

971   32. Siniscalco, D., Schultz, S., Brigida, A. L. & Antonucci, N. Inflammation and Neuro-Immune

972   Dysregulations in Autism Spectrum Disorders. *Pharmaceuticals* **11**, 56 (2018).

973   33. Arteaga-Henríquez, G., Gisbert, L. & Ramos-Quiroga, J. A. Immunoregulatory and/or Anti-

974   inflammatory Agents for the Management of Core and Associated Symptoms in Individuals

975   with Autism Spectrum Disorder: A Narrative Review of Randomized, Placebo-Controlled

976   Trials. *CNS Drugs* **37**, 215–229 (2023).

977   34. Pérez-Sen, R. *et al.* Dual-Specificity Phosphatase Regulation in Neurons and Glial Cells.

978   *Int. J. Mol. Sci.* **20**, 1999 (2019).

979   35. Ebstein, F. *et al.* PSMC3 proteasome subunit variants are associated with

980   neurodevelopmental delay and type I interferon production. *Sci. Transl. Med.* **15**, eabo3189

981   (2023).

982   36. Tessadori, F. *et al.* Recurrent de novo missense variants across multiple histone H4 genes

983   underlie a neurodevelopmental syndrome. *Am. J. Hum. Genet.* **109**, 750–758 (2022).

984   37. Spencer, S. A., Suárez-Pozos, E., Escalante, M., Myo, Y. P. & Fuss, B. Sodium-calcium

985   exchangers of the SLC8 family in oligodendrocytes: Functional properties in health and

986   disease. *Neurochem. Res.* **45**, 1287–1297 (2020).

987    38. Grammatikakis, I., Zhang, P., Mattson, M. P. & Gorospe, M. The long and the short of TRF2

988        in neurogenesis. *Cell Cycle* **15**, 3026–3032 (2016).

989    39. Richards, S. *et al.* Standards and Guidelines for the Interpretation of Sequence Variants: A

990        Joint Consensus Recommendation of the American College of Medical Genetics and

991        Genomics and the Association for Molecular Pathology. *Genet. Med. Off. J. Am. Coll. Med.*

992        *Genet.* **17**, 405–424 (2015).

993    40. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).

994    41. Adzhubei, I. A. *et al.* A method and server for predicting damaging missense mutations. *Nat.*

995        *Methods* **7**, 248–249 (2010).

996    42. Sim, N.-L. *et al.* SIFT web server: predicting effects of amino acid substitutions on proteins.

997        *Nucleic Acids Res.* **40**, W452–W457 (2012).

998    43. Dong, C. *et al.* Comparison and integration of deleteriousness prediction methods for

999        nonsynonymous SNVs in whole exome sequencing studies. *Hum. Mol. Genet.* **24**, 2125–

1000        2137 (2015).

1001   44. Kaitlin E. Samocha *et al.* Regional missense constraint improves variant deleteriousness

1002        prediction. *bioRxiv* 148353 (2017) doi:10.1101/148353.

1003   45. Chen, S. *et al.* A genome-wide mutational constraint map quantified from variation in 76,156

1004        human genomes. 2022.03.20.485034 Preprint at https://doi.org/10.1101/2022.03.20.485034

1005        (2022).

1006   46. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide

1007        polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w1118; iso-

1008        2; iso-3. *Fly (Austin)* **6**, 80–92 (2012).

1009   47. Liu, X., Li, C., Mou, C., Dong, Y. & Tu, Y. dbNSFP v4: a comprehensive database of

1010        transcript-specific functional predictions and annotations for human nonsynonymous and

1011        splice-site SNVs. *Genome Med.* **12**, 103 (2020).

1012    48. THE GTEX CONSORTIUM. The GTEx Consortium atlas of genetic regulatory effects across

1013        human tissues. *Science* **369**, 1318–1330 (2020).

1014    49. SHAPIRO, S. S. & WILK, M. B. An analysis of variance test for normality (complete

1015        samples)†. *Biometrika* **52**, 591–611 (1965).

1016    50. Gu, Z. Complex heatmap visualization. *iMeta* **1**, e43 (2022).

1017