# Unified AI framework to uncover deep interrelationships between gene expression and Alzheimer's disease neuropathologies

Nicasia Beebe-Wang[1], Safiye Celik[2], Ethan Weinberger[1], Pascal Sturmfels[1], Philip L. De Jager[3],

Sara Mostafavi[1,4,*] and Su-In Lee[1,*]


[1] Paul G. Allen School of Computer Science and Engineering, University of Washington, WA, Seattle, USA.

[2] Benevolent Artificial Intelligence, NY, USA.

[3] Center for Translational and Computational Neuroimmunology, Department of Neurology, Columbia University Medical Center, New York, NY, USA.

[4] Department of Statistics, University of British Columbia, BC, Canada.

* These authors contributed equally to this work.

## ABSTRACT

Deep neural networks offer a promising approach for capturing complex, non-linear relationships among variables. Because they require immense sample sizes, their potential has yet to be fully tapped for understanding complex relationships between gene expression and human phenotypes. Encouragingly, a growing number of diseases are being studied through consortium efforts. Here we introduce a new analysis framework, namely MD-AD (**M**ulti-task **D**eep learning for **A**lzheimer's **D**isease neuropathology), which leverages an unexpected synergy between deep neural networks and multi-cohort settings. In these settings, true joint analysis can be stymied using conventional statistical methods, which (1) require "harmonized" phenotypes (i.e., measured in a highly consistent manner) and (2) tend to capture cohort-level variations, obscuring the subtler true disease signals. Instead, MD-AD incorporates multiple related phenotypes sparsely measured across cohorts, and learns complex, non-linear interactions between genes and phenotypes not discovered using conventional expression data analysis methods (e.g., component analysis and module detection), enabling the model to capture subtler signals than cohort-level variations. Applied to the largest available collection of brain samples (N=1,758), we demonstrate that MD-AD learns a truly generalizable relationship between gene expression program and AD-related neuropathology. The learned program generalizes in several important ways, including recapitulation of the disease progress in animal models and across tissue types, and we show that such generalizability is not achieved by previous statistical paradigms. Its ability to identify genes with high non-linear relevance to neuropathology enabled us to identify a sex-specific relationship between neuropathology and immune response across microglia, providing a nuanced context for association between inflammatory genes and AD.

37  **INTRODUCTION**

38  Alzheimer's disease (AD), the sixth leading cause of death in the United States, is a degenerative brain
39  condition with no known treatment to prevent, cure, or delay its progression. Primary challenges to
40  treating and preventing AD include *extensive heterogeneity* in the clinicopathologic state of older
41  individuals[1] and *limited knowledge* about genetic and molecular drivers and suppressors of AD-related
42  (amyloid and tau) proteinopathies and AD dementia[2]. Recent efforts to identify molecular mechanisms
43  underlying AD and its progression focus on two complimentary approaches. First, the assembly of large
44  genome-wide association studies (GWAS) (N>100K subjects) enabled case/control analyses of genetic
45  variants correlated with a clinical diagnosis of AD. Interestingly, some identified variants have implicated
46  tau protein binding, amyloid precursor protein (APP) metabolism or immune pathways that play a role in
47  their aggregation and/or uptake[3–5]. These results reinforce the need for detailed investigations of the
48  drivers of neuropathological variation across individuals. Second, moderate-scale post-mortem
49  transcriptomic studies have investigated molecular correlates of a richer set of phenotypic and
50  neuropathological outcomes[6–9]. Early work in this domain examined pairwise correlations among gene
51  expression levels and AD related traits[10] or a diagnosis of AD[11]. More recent attempts have focused on
52  learning statistical dependencies among gene expression using AD expression data collected from one
53  cohort, in order to infer gene regulatory networks[7] or co-expressed modules[6] associated with AD related
54  phenotypes (see Supplementary Methods for details). The relative scarcity of brain gene expression data
55  collected from each cohort has posed a challenge to the use of complex models, such as deep neural
56  networks.

57  The collection of postmortem brain RNA-sequencing datasets, assembled by the AMP-AD (**A**ccelerating
58  **M**edicines **P**artnership **A**lzheimer's **D**isease) consortium, provides a unique opportunity to combine
59  multiple data sets in an integrative analysis. Previous work has applied existing co-expression methods to
60  each dataset and used consensus methods to identify consistent gene expression modules across datasets[9].
61  To our knowledge, there has not yet been a *unified* approach to learn a single joint model that
62  incorporates multiple AMP-AD datasets, which would enable the use of all samples to capture intricate
63  interactions between gene expression levels and phenotypes. A unified approach has been hindered by:
64  (1) the need for "harmonized" phenotypes consistently measured across datasets, and (2) the limitation of
65  current analysis methods that focus on linear relationships between variables (e.g., module analysis[9])
66  which tend to capture broader patterns in gene expression that often correspond to cohort-level variations,
67  and to consequently obscure true disease signals.

68  Here, we develop MD-AD (**M**ulti-task **D**eep learning for **A**lzheimer's **D**isease neuropathology), a *unified*
69  *framework for analyzing heterogeneous AD datasets* to improve our understanding of expression basis for
70  AD neuropathology (**Figure 1a-d**). Unlike previous approaches, MD-AD learns a single neural network
71  by jointly modeling multiple neuropathological measures of AD (**Figure 1a**), and hence it incorporates a
72  large collection of postmortem brain RNA-sequencing datasets. The combined AMP-AD dataset contains
73  1,758 samples distributed across 9 brain regions which are labeled with up to six neuropathological
74  outcomes that are *sparsely* available across cohorts (**Figure 1e**). This *unified* framework has key
75  advantages over separately trained models. First, MD-AD can accommodate sparsely labeled data, which
76  is a natural characteristic of datasets aggregated through consortium efforts (**Figure 1e**). Even if different
77  phenotypes only partially overlap in the measured samples, each sample contributes to the training of both
78  phenotype-specific and shared layers (**Figure 1a**). Predicting multiple phenotypes at once biases shared
79  network layers to capture relevant features of these AD phenotypes at the same time. This is of critical

80    importance: each phenotype represents a *different* noisy measurement of the same underlying true
81    biological process, and, as we demonstrate, joint training allows MD-AD to average out the noise to
82    extract the true hidden signal. Additionally, the increased sample size enables MD-AD to capture
83    complex non-linear interactions between genes and phenotypes. Multi-layer perceptrons (MLPs) offer
84    another powerful approach for directly capturing complex relations between gene expression and a
85    phenotype. However, training separate MLPs for each phenotype (**Supplementary Figure 1a**) has limited
86    scope: it can utilize only the samples measured for a specific phenotype, and it cannot share information
87    across related phenotypes. We demonstrate that these advantages improve MD-AD prediction accuracy,
88    enabling its predictions to generalize across species and tissue types (**Figure 1b**).

89    MD-AD's ability to capture complex non-linear relationships provides an opportunity to gain new
90    insights into the expression basis of AD neuropathology, which were not identified by previous
91    approaches. However, an obvious drawback of deep neural networks is their black-box nature, making it
92    difficult to biologically interpret gene-phenotype associations. This paper presents two ways to address
93    this challenge. First, MD-AD adopts a well-known feature attribution method[12], which quantifies how
94    much each input variable (here, gene expression level) contributes to a prediction (here, a
95    neuropathological phenotype) to identify genes and pathways relevant to each neuropathological
96    phenotype (**Figure 1d**). Second, because MD-AD is a deep learning model, we can interpret its
97    intermediate layers as biologically relevant *high-level feature representation* of gene expression levels
98    and its predictions as the amalgamation of AD-specific molecular markers. The last shared layer of MD-
99    AD can be viewed as *a supervised embedding* influenced by each neuropathological phenotype used
100   during training. Thus, by interpreting this layer's embedding, we gain understanding of model
101   components and high-level dependencies between expression and neuropathology (**Figure 1c**). As the
102   first deep learning attempt to relate gene expression to multiple AD neuropathological phenotypes, we
103   identify globally important genes not previously implicated in linear methods and perform sex-specific
104   analyses to explore implicitly captured non-linear effects among genes and AD severity predictions.

105   In sum, our new MD-AD framework makes the following contributions: (1) It is able to *effectively impute*
106   *accurate AD neuropathological phenotype predictions* from broad compendia of heterogeneous brain
107   gene expression data; (2) it produces learned representations that are more robust than separately learned
108   models, improving *generalizability to other datasets, species, and even tissue types*; (3) it provides an
109   *improved understanding of inter-relationships* among molecular drivers of AD neuropathology that is
110   missed by linear methods; and (4) from a biological standpoint, MD-AD highlights a sex-specific
111   relationship between microglial immune activation and neuropathology.

112   **RESULTS**

113   **MD-AD provides a unified framework to learn a single model of multiple neuropathological**
114   **phenotypes across multiple cohort datasets**

115   The MD-AD model takes as input brain gene expression profiles and simultaneously predicts several AD-
116   related neuropathological phenotypes (**Figure 1a**). In particular, the model is trained on expression data
117   from the ROSMAP[6,13,14], ACT[15] and MSBB[16] cohort studies, which together have 1,758 gene expression
118   profiles for 925 distinct individuals. These data are normalized for study batch (Supplementary Methods,
119   **Supplementary Figure 1b-c**) [17]. As shown in **Figure 1a**, the MD-AD model simultaneously predicts six
120   AD-related neuropathological phenotypes: three related to amyloid plaques and three to tau tangles. The
121   former include: (1) **Aβ IHC**: amyloid-β protein density via immunohistochemistry, (2) **NPs:** neuritic

122 amyloid plaque counts from stained slides, and (3) **CERAD score:** a semi-quantitative measure of
123 neuritic plaque severity[18]. The latter include: (4) **τ IHC:** abnormally phosphorylated τ protein density via
124 immunohistochemistry, (5) **tangles:** neurofibrillary tangle counts from silver stained slides, and (6)
125 **Braak stage:** a semi-quantitative measure of neurofibrillary tangle pathology [19]. Thus, MD-AD generates
126 six highly related predictions simultaneously and covers each of the two main hallmarks of AD
127 neuropathology (plaques and tangles) at three levels of granularity. The three studies measure partially
128 overlapping subsets of the six phenotypes described above (**Figure 1e** and **Table 1**), so across our
129 combined dataset some variables are sparsely labeled, although Braak and CERAD are each measured in
130 all studies (**Figure 1e**). During training, the MD-AD model continually updates model parameters via
131 backpropagation, but only for labeled phenotypes from a given sample. Thus, for each phenotype for a
132 given sample, MD-AD updates parameters from associated separate layers along with all shared layers.
133 This allows us to train a unified model from all available samples despite having many missing labels.

134

135 **MD-AD accurately predicts neuropathology from gene expression, and its predictions are**
136 **generalizable to external datasets.**

137 In the first pass at model evaluation, we assessed MD-AD using standard five-fold cross-validation (CV),
138 quantifying the average mean squared error (MSE) on the test samples (**Figure 2a and Supplementary**
139 **Figure 1d**). We compared MD-AD to two simpler baseline models: a regularized linear model (ridge
140 regression) and a single output deep neural network (MLP). These alternative results helped us assess two
141 significant components of the MD-AD model: (1) its non-linear modeling of the relation between gene
142 expression and neuropathological phenotypes, and (2) its joint modeling of multiple related
143 neuropathological phenotypes. In general, MLP models outperformed linear models, highlighting a
144 general advantage of deep learning over a linear approach. Furthermore, compared to the MLP models,
145 MD-AD showed MSE reductions of 7% for CERAD score, 13% for Braak stage, 7% for NPs, 25% for
146 tangles, 10% for Aβ immunohistochemistry (IHC), and 14% for τ IHC (**Figure 2a**). Interestingly, MD-
147 AD showed its largest performance gain for the tangles variable, which also had the most missing labels
148 (**Figure 1e**), highlighting a specific advantage of joint learning for sparsely labeled data.

149 Because our model was trained and evaluated on ACT, MSBB, and ROSMAP datasets, we assessed
150 whether residual (uncorrected) batch effects affected performance. To do so, we performed additional
151 validation experiments by leaving out specific datasets during training and then evaluating their
152 performance for MD-AD trained on the other datasets (**Figure 2b, Supplementary Figure 2a**). We
153 evaluated MSE performance for ROSMAP alone since it was the only dataset with all six phenotype
154 labels; further, by evaluating a single dataset's performance, we can identify the influence of adding
155 "external" data. We make several observations from this analysis. First, as one may expect, larger training
156 samples always helped reducing prediction error on test samples from the unseen study (ROSMAP), and
157 especially so when datasets from multiple cohorts were included in the training (i.e., ACT and MSBB)
158 (circular markers in **Figure 2b**). Second, when considering the effects of augmenting ROSMAP data with
159 other datasets during training (diamond markers in **Figure 2b**), we observed that errors initially increased
160 when adding a new dataset but tended to decline as more datasets were included in training. This may
161 result from small differences in labeling conventions across studies, or batch effects in gene expression
162 data. However, we find that the benefits of additional heterogeneous samples ultimately outweigh
163 potential batch effects in prediction performance. Third, interestingly, we observed that adding new
164 samples improved performance for a phenotype even when the phenotype in question was not measured

165      in the new samples (see gray footprints around markers in **Figure 2b**). This suggests that the shared
166      representation learned by MD-AD captures the underlying biological signal common across noisy
167      neuropathological phenotype measurements.

168      Next, as the ultimate test of MD-AD out-of-sample predictions, we assessed performance on three
169      independent studies never seen by the model: Mount Sinai Brain Bank Microarray (MSBB; N=1,053),
170      Harvard Brain Tissue Resource Center (HBTRC; N=460), and Mayo Clinic Brain Bank (N=323).
171      Because these datasets provide a sparse set of neuropathological labels, we evaluated whether MD-AD
172      predictions were consistent with the (binary) neuropathological diagnosis of AD by calculating "MD-AD
173      neuropathology scores" for each sample (by averaging ranked predictions across the six phenotypes). For
174      comparison with other methods, we also generated "neuropathology score" predictions for our baseline
175      models.

176      As shown in **Figure 2c**, we observed a highly significant difference in predicted neuropathology scores
177      between AD cases and controls (two-sided t-test: $t=22.98$, $p<0.001$), and these differences were more
178      pronounced for MD-AD compared to the other baseline models (results split by dataset are shown in
179      **Supplementary Figure 3a**). More convincingly, when split by age group (**Figure 2c** right panel), we
180      consistently observed a significant increase in predicted neuropathology for AD vs control samples, but
181      the difference was largest in individuals under 75 (between-groups $p$-values are shown in **Supplementary**
182      **Figure 3b**). This is consistent with the observation that aging individuals who are cognitively non-
183      impaired often have substantial neuropathology[15]. Together, these results indicate that MD-AD can
184      identify generalizable gene expression patterns that are predictive of AD-related neuropathology across
185      varied age ranges, and thus it is unlikely that these patterns merely capture normal aging.

186

### Complex transcriptomic predictors of neuropathology are conserved across species.

188      We next evaluated how well MD-AD's learned expression patterns predictive of neuropathology
189      recapitulated neuropathology in mouse models. We applied MD-AD trained on human datasets to make
190      predictions based on brain (hippocampal and cortical) gene expression data from 30 TASTPM mice that
191      harbored double transgenic mutation in APP and PSEN1 and compared the predictions to those from 76
192      wild type mice[20]. We focused on TASTPM mice because they were found to robustly exhibit early signs
193      of amyloid aggregation and plaque formation. As above, to simplify MD-AD predictions, we then
194      predicted all six neuropathological phenotypes via MD-AD and generated an aggregate "neuropathology
195      score" per mouse (as described in Supplementary Methods).

196      As shown in **Figure 2d**, MD-AD predicted significantly higher neuropathology scores for the
197      homozygous cross TASTPM than wild type mice (two-sided t-test: $t=3.45$, $p<.001$). The MLP baseline
198      method also produced significant differences between homozygous and wild type mice, but less
199      effectively ($t=3.01$, $p<.01$). Furthermore, there was a stronger trend for higher predictions in the
200      heterozygous TASTPM cross (N=32) than wild type mice for MD-AD ($t=1.38$, $p=.17$) compared to MLP
201      baselines ($p=.38$). Interestingly, our linear baseline tended to predict lower average neuropathology levels
202      for these AD strains than wild type, suggesting that a linear approach may fail to effectively model cross-
203      species AD signal. None of the models produced significantly different neuropathology scores between
204      other strains (i.e., TPM, TAS10, Tau) and wild type mice, consistent with lower neuropathological burden
205      in these models (data not shown). Notably, when we stratified the samples by age, we found that MD-AD
206      tended to predict higher neuropathology in older mice regardless of strain), but in particular it made

207  higher neuropathology predictions for homozygous than heterozygous crosses followed by wild type mice
208  (many of these groups differed significantly from one another, as shown in **Supplementary Figure 3c**).
209  Overall, these results indicate that MD-AD learns a generalizable expression pattern associated with
210  neuropathology that is conserved across species.

211

212  **Deep transcriptomic signatures of neuropathology are predictive of AD dementia**

213  Hidden layers of a deep neural network capture the embedding of input examples in the derived feature
214  space, yielding a "hidden" representation that is predictive of the outcome(s) of interest. In this case, the
215  last shared layer of MD-AD (**Figure 1a, c**) captures a latent (lower) dimensional representation of gene
216  expression that is predictive of multiple types of neuropathology related to AD. To derive the biological
217  basis of MD-AD predictions, we first visualized this embedding space in 2D using the t-SNE algorithm
218  (**Figure 3a**) [21] (to improve stability, we used a consensus approach over many re-trainings of the MD-AD
219  model, **Supplementary Figure 4a**). We observed that the representation in this space was impressively
220  coherent with respect to all six neuropathological variables: individuals with similar overall
221  neuropathology severities had similar MD-AD consensus representations for their gene expression
222  profiles, and this observation was true for external test samples not used for model training (**Figure 3d-e,
223  Supplementary Figure 3d**). This was remarkable because representations derived by unsupervised
224  dimensionality reduction (e.g., K-means or PCA) failed to capture the components of gene expression
225  relevant to neuropathology, and mainly captured batch effects, while those derived by standard single
226  output MLP tended to overfit to each neuropathology variable and were incoherent *across*
227  neuropathological measurements (**Figure 3c** and **Supplementary Figure 5**).

228  Next, we evaluated whether the MD-AD embedding can go beyond neuropathology to also capture the
229  molecular manifestation of AD dementia. In particular, we considered three "higher-level" phenotype
230  variables: AD dementia (a *clinical* diagnosis of AD), assessment of cognitive function, and assessment of
231  AD duration. We then correlated the latent representation captured by the hidden nodes in the last shared
232  layer with each of these three higher-level phenotypes. As shown in **Figure 3b**, we found that MD-AD
233  consistently produced nodes that were significantly correlated with high-level AD phenotypes; using
234  paired *t*-tests, these correlations often outperformed nodes from our MLPs and always outperformed
235  unsupervised methods and module-based approaches ($p<.05$ after FDR correction over nodes). This
236  indicates that MD-AD creates embeddings that most consistently capture the relationship between gene
237  expression and general AD severity. Together, these results show that by jointly predicting several
238  neuropathological phenotypes, the MD-AD framework produces a low dimensional representation of
239  gene expression data, in the form of embedding nodes, that robustly captures a generalizable signature of
240  AD beyond individual neuropathological phenotypes alone. Detailed annotations for MD-AD embedding
241  nodes are provided in **Supplementary Table 2** and **Supplementary Figure 4b-d**.

242

243  **MD-AD reveals an interrelationship between sex and immune genes predictive of AD
244  neuropathology**

245  We next sought to interpret MD-AD's learned parameters to identify the set of genes (and their
246  relationships) that underlie its impressive predictive performance. Here, we applied the Integrated
247  Gradient (IG) algorithm[12] on the fully trained model in an ensemble fashion to ensure robustness
248  (Supplementary Methods, **Supplementary Figure 6a-b**), producing an "importance score" for each gene.

249 For a global view, we first performed functional enrichment analysis (GSEA[22,23]) using these importance
250 scores, and found that relevant genes for the MD-AD model were enriched for several pathways,
251 including metabolism of RNA and proteins, immune system, cell-to-cell communication, and signal
252 transduction (**Figure 4b**). **Figure 4a** shows the top 50 genes and their pathway annotations where the
253 particular relevance of immune function is even more prominent.

254 We next assessed to what extent the learned gene importance varied between a linear model and a non-
255 linear model like MD-AD. With a simple linear correlation-based gene ranking, we found that the top 50
256 genes had a much lower prevalence of REACTOME pathways (**Supplementary Figure 7a**). When we
257 directly compared the top 1% of genes from MD-AD versus a correlation-based approach in **Figure 4c**,
258 we observed that many genes belonging to metabolism, immune system, and signal transduction
259 pathways were highly ranked for MD-AD but not for correlation-ranking. In contrast, transcription-
260 related genes were more frequently highly ranked for correlation-based rankings compared to MD-AD's
261 rankings. Overall, gene importance scores generated via correlations alone were enriched for a much
262 larger set of REACTOME categories (**Supplementary Figure 7b**), whereas MD-AD pathways tended to
263 be more specific (**Figure 5b**). We saw similar results when performing the same analyses with KEGG
264 pathways (**Supplementary Figure 8**) [24].

265 The nonlinear relationships identified by MD-AD can implicitly capture interaction effects with other
266 covariates observable from expression data (e.g., sex, age, medication intake). Leveraging the fact that, if
267 our model captures a nonlinear effect, then two samples with the same expression level for a single gene
268 could receive different IG ("importance") scores by MD-AD (e.g., **Figure 5d**; in contrast, a linear model
269 would have no vertical dispersion), we assessed whether a covariate like sex could explain discrepancy
270 between expression levels and IG scores. (Sex is a major risk factor in AD and has prominent gene
271 expression signatures[25]). Thus, we modeled each gene's IG score as a linear combination of the gene's
272 expression, the individual's sex, and the interaction between them to identify sex-interacting genes
273 relevant to AD. Of the 14,591 genes in our dataset, 6,465 showed differential MD-AD importance
274 between sexes ($p<0.05$ after FDR), demonstrating that sex-specific expression effects in AD may be
275 widespread. To confirm that genes are not sex-differential by chance, we show the distribution of sex-
276 differential genes compared with the same analysis conducted with shuffled sex labels (**Supplementary
277 Figure 9a**). However, we were particularly interested in genes with high overall MD-AD importance.
278 When focusing on the top 100 genes with the highest MD-AD scores, we consistently observed high
279 degrees of interaction between sex and immune system genes (as well as reproduction and hemostasis-
280 related genes) (**Figure 5a-b;** we saw similar patterns for KEGG pathways in **Supplementary Figure 9b-
281 c**).

282 We next explored specific examples of genes with high MD-AD rankings and strong interactions with sex
283 (i.e., the six genes from the top 100 MD-AD list with the strongest interaction *p*-values; **Figure 5c-d**):
284 *KNSTRN, C4B, CMTM4, TREM2, P2RY11*, and *SERPINA3*. In particular, for each of these genes, we
285 observed high expression values associated with higher neuropathology predictions but some
286 stratification across sexes: high expression in females led to especially high neuropathology predictions
287 for *KNSTRN* and *P2RY11*, while the opposite was true for the other four genes. More broadly, our finding
288 immune genes display sex-differential contributions to MD-AD scores appears to be consistent with
289 conclusions from recent studies about sex differences in neuroinflammatory activity and the role these
290 differences may play in neurodegenerative disorders[26].

We note that some of our top sex-interacting genes may play important roles in immune response, particularly in microglia. *TREM2*, which is genetically implicated in AD, interacts with *CD33* (another AD susceptibility gene) [27], is an important contributor in the clearance of toxic Amyloid-β by microglia in mice [28], and is correlated with Aβ deposition in the human brain [27]. Similarly, *KNSTRN* is known to be upregulated in mouse microglial cells' early response to neurodegeneration[29]. These findings indicate that MD-AD may capture patterns related to sex-differential microglia activity. To explore this idea further, we obtain lists of upregulated genes from nine clusters of single cell microglial transcriptomes[30], and compare them to our MD-AD gene rankings. As expected, many top MD-AD genes are upregulated in multiple microglial clusters (**Figure 6a**); correlation-based methods ranked these microglial genes less highly (**Supplementary Figure 9d**). Furthermore, genes upregulated in clusters related to stress, immune function and proliferation tended to be sex-differential in their gene importance (**Figure 6b**), further strengthening the finding that sex differences in immune response and inflammation may be an important factor in the molecular basis of age-related neuropathology.

To more broadly identify possible cell-type specific effects of MD-AD's important genes, we tested for the enrichment of 41 different cell type clusters (across six cell types) found by a single cell transcriptomic analysis of AD[8]. Here, we found an enrichment of 2 different microglia clusters, as well as astrocytes and inhibitory neuron clusters (**Figure 6c**). Hence, MD-AD's predictions of neuropathology rely on broader transcriptomic events that goes beyond microglia genes, suggesting a heterogeneity in the underlying molecular biology that is predictive of accumulation of AD-related neuropathology.

**Complex transcriptomic predictors learned by MD-AD are conserved across tissues.**

Although MD-AD was developed for brain gene expression data, we next asked whether the learned transcriptomic signatures generalize to blood. To this end, we applied our brain-trained MD-AD model to gene expression datasets from two batches of the AddNeuroMed cohort, which we called Blood1 and Blood2 (NCBI GEO database accessions GSE63060 and GSE63061, respectively; summarized in **Supplementary Table 3**)[31]. As shown in **Figure 7a**, MD-AD predicted significantly higher neuropathology scores for individuals with both mild cognitive impairment (MCI) (two-sided t-test: $t$=7.34, $p$ <.001) and AD dementia (two-sided t-test: $t$=5.87, $p$ <.001) compared to cognitively normal controls (CTL). Consistent with external brain samples shown in **Figure 2d** and **2f**, MD-AD predictions tended to increase with age for cognitively normal individuals, while they were consistently significantly higher for MCI and AD individuals compared to controls for individuals under 80 years old (**Figure 7b, Supplementary Figure 10b**). Importantly, we noted that a linear model failed to make meaningful predictions (**Figure 7a** and **Supplementary Figure 10a**), suggesting that complex models like MD-AD have better performance in extracting the true underlying signal transferrable between tissues than linear models.

Next, we evaluated whether the patterns captured by the MD-AD model were consistent across training brain gene expression samples and blood. To this end, we again visualized MD-AD's learned embedding using the t-SNE algorithm (**Figure 7c**). We noted a clear difference in expression patterns between blood and brain samples (as seen by the clustering of blood samples in **Figure 7c**); however, MD-AD nevertheless produced an embedding for blood data that stratified blood samples along predicted neuropathological phenotypes in a manner highly consistent with the blood donor's cognitive status (**Figure 7c; Supplementary Figure 10c**). Together, these analyses indicate that jointly learning the

333 relationship among brain gene expression and several neuropathological phenotypes may allow for
334 learned representations that span tissues. This in turn can open up new avenues for early identification of
335 individuals at risk, and provide new clues into tissue-agnostic molecular mechanisms underlying AD
336 dementia.

337

338 **DISCUSSION**

339 We introduce MD-AD, a deep neural network approach for jointly modeling the relationship between
340 brain gene expression data and multiple sparsely labeled neuropathological phenotypes in a multi-cohort
341 setting. By exploiting the synergy between deep learning and a multi-cohort, multi-task setting, we
342 demonstrated that MD-AD can capture complex, non-linear feature representations that are not learned
343 using conventional expression data analysis methods. Specifically, we observed that multi-task learning
344 improves prediction performance over singly trained models. Adding data from different cohorts
345 improves performance for various phenotypes, even those that lacked labels. When we extended our
346 method to other datasets, it captured AD-related biological signals, showing that MD-AD can transfer
347 effectively to out-of-sample, out-of-species (mouse), and even out-of-tissue (blood) datasets.

348 As a neural network framework, MD-AD's last shared layer embedding reveals high-level features of
349 gene expression that are predictive of neuropathology according to the intermediate components of the
350 model. As expected, due to multi-task supervision, our embedding nodes tend to relate to AD-associated
351 neuropathology far more effectively than do standard unsupervised approaches and earlier reported
352 (unsupervised) module-based approaches. Compared to singly task-supervised neural networks, the joint
353 training MD-AD performs consistently provided a more stable and coherent AD-related embedding. By
354 exploring the molecular pathways relevant to each node, we identified relevant gene sets contributing to
355 these high-level AD-related features of gene expression.

356 Finally, we leveraged the complex relationships learned by MD-AD to refine our understanding of the
357 molecular drivers of AD neuropathology. By interpreting genes relevant to our model's predictions, we
358 uncovered that MD-AD relied on many genes not found in earlier linear-based methods, including several
359 immune system genes. These findings expand the general narrative established by human genetic studies
360 of AD and now a proteomic study of AD[32]; in particular, we see enrichment for complement pathway
361 genes (**Figure 4**) which likely connect with the role of the complement receptor 1 (*CR1*) gene which
362 harbors an AD susceptibility variant whose functional consequences remain poorly understood but do
363 include an influence on the accumulation of neuritic plaque pathology[33–36]. Thus, MD-AD results
364 converge with human genetic results to emphasize the role of complement in AD; interestingly
365 complement protein C4B emerges as one of the top pathology-related genes that display a strong
366 interaction with sex, with men showing a much stronger association than women (**Figure 5c**). This is
367 similar to the behavior of TREM2, another well-validated AD susceptibility gene (**Figure 5c**); however,
368 its relation to amyloid pathology in ROSMAP data was previously reported as being modest[27]. MD-AD
369 was able to uncover its more prominent role in transcriptional data, which is obscured by its sex-
370 dependent nature. Likewise, women reported to have higher expression of a signature of aged microglia
371 in these data[26], and two modules of co-expressed cortical genes enriched for microglial genes and
372 associated with amyloid (module m114) or tau (module m5) pathology are also influenced by sex[37].
373 However, the role of neither group of genes is explained by sex; this indicates that the role of sex in the
374 impact of the immune system in AD is complex. MD-AD was able to uncover this complexity more

375 effectively, as is illustrated in **Figure 5c** where some genes have greater effects in men and other in
376 women. Thus, it is not the case that role of the immune system is polarized in one of the two sexes; rather,
377 some pathways and perhaps certain cell subsets may have a larger role in women while others are
378 dysfunctional in men. This could explain why the role of immune genes is more prominent in our
379 analyses: reports from simpler linear models often included immune pathways[6] but other pathways
380 usually figured more prominently in these earlier RNA-based network models. A meta-analysis of RNA
381 studies (which include the ROSMAP data) highlighted the larger number of sex-influenced genes among
382 the AD-associated gene modules and noted that microglial cells appear to be enriched for both male and
383 female-specific expression effects. With our list of results and our careful evaluation of sex effects we
384 now have an important new road map with which to guide our exploration of the role of microglia in AD
385 in a sex-informed manner. This perspective will be critical not only for mechanistic studies whose results
386 could be obscured by sex effects but also, more importantly, by guiding the study design of clinical trials
387 as highly targeted therapeutic agents emerge to modulate the immune system in AD.

388 This is but one of the narratives that has emerged from our initial deployment of the MD-AD approach in
389 the aging brain. As new cohorts are characterized, sample sizes expand and new data such as single
390 nucleus RNA sequencing profiles emerge, our approach will help to facilitate data integration and to
391 uncover insights that would not otherwise emerge. Beyond enabling good predictions, our report may
392 actually highlight a more important contribution of MD-AD in resolving key elements of the data
393 structure in the nodes that we defined: these are more than simple aggregates of factors with predictive
394 power. They are beginning to uncover complex interactions, such as the impact of sex which is involved
395 in both men and women, but in different ways, making it difficult to appreciate the role of certain immune
396 pathways in simpler statistical models.

397

398

399 **FIGURE LEGENDS**

400 **Figure 1**. Overview of the MD-AD method and analyses. (**a**) Overview of the MD-AD framework: MD-
401 AD is trained to predict six neuropathology phenotypes simultaneously from brain gene expression
402 samples. During model training, samples do not need to have all available phenotypes; they influence
403 only the layers for which they have labels (including shared layers). (**b**) Illustrates out-of-sample datasets
404 we used to validate MD-AD's predictions (**c**) Illustrates analyses used to validate the last shared layer of
405 MD-AD. (**d**) By using model interpretability methods, we highlight genes relevant to MD-AD's
406 predictions. Further analyses reveal non-linear effects among genes and their relationship with AD
407 severity prediction.

408 **Figure 2**. MD-AD prediction performance for within-sample and out-of-sample data. (**a**) Average test set
409 mean squared error (MSE) for phenotype predictions across 5 test splits. MLP: Multiple Layer
410 Perceptron. Linear: linear model using L2 regularization. (**b**) Average MSE for ROSMAP test set samples
411 when training on subsets of the available data sets in the training set. (**c**) For samples from three external
412 validation data sets, we obtain neuropathology scores for each sample by averaging the percentiles of
413 predictions across all six neuropathology variables. *Left*: t-test statistics measuring the difference between
414 each model's predicted neuropathology scores for AD-diagnosed vs. control individuals. All tests results
415 were statistically significant ($p<.001$). *Right*: Box plots displaying the distribution of MD-AD's predicted
416 neuropathology scores split by age group and diagnosis (see **Supplementary Figure 3b** for sample sizes
417 and significance of pair-wise differences). (**d**) *Left*: t-test statistics measuring the difference between each
418 model's predicted neuropathology score for heterozygous TASTPM vs. wild type mice. *Middle*: t-test
419 statistics measuring the difference between each model's predicted neuropathology score for homozygous
420 TASTPM vs. wild type mice (*: $p<.05$, **: $p<.01$, ***:$p<.001$). *Right*: Box plots displaying the
421 distribution of MD-AD's predicted neuropathology scores for mice split by age and strain (See
422 **Supplementary Figure 3c** for sample sizes and significance of pair-wise differences).

423 **Figure 3.** Comparing MD-AD's supervised embedding to other embedding methods. (**a**) For each colored
424 box, *Left*: 2-dimensional t-SNE embedding of MD-AD's last shared layer colored by neuropathological
425 phenotype indicated in the title of the box, *Right*: -$\log_{10}(p$-value) of correlations between "best" node
426 from each embedding method and the neuropathological phenotype across 5 test folds. The "best" node
427 was identified as the most significantly correlated in the training set, but the figure reports correlation -
428 $\log_{10}(p$-value)'s in their corresponding test sets. Bar graph columns (left to right): two unsupervised
429 embeddings (green; K-Means and PCA), three module-based embeddings (orange; Modules #1 [7] and
430 Modules #2 [6], and Modules #3 [9]), six singly-trained MLPs (blue), and MD-AD (red). (**b**) Highest
431 correlation -$\log_{10}(p$-values) (averaged across 5 training folds) found between each embedding method and
432 high-level AD phenotypes: dementia (diagnosis prior to death), dementia duration (approximate time
433 between dementia diagnosis and death; available for ACT and ROSMAP), and last available cognition
434 score (controlling for age, sex and education; available for ROSMAP only). All $p$-values listed are shown
435 after FDR correction over the nodes within each method. (**c**) 2-dimensional t-SNE embedding of
436 alternative embedding methods (described in **a**). (**d**) 2-dimensional t-SNE embeddings of MD-AD
437 embeddings for training and external data sets. Each point represents a sample colored by dataset (Left),
438 AD status for external samples (Middle) and MD-AD's predicted neuropathology score (Right). (**e**) 2-
439 dimensional t-SNE embeddings of MD-AD embeddings for external human and mouse samples.

440 **Figure 4.** Top predictive genes for the consensus MD-AD model. (**a**) Top 50 MD-AD genes and whether
441 they are negatively (-) or positively (+) associated with high neuropathology. Colored squares indicate

442     that the gene belongs at least one pathway in the column-labeled REACTOME category. **(b)** Gene set
443     enrichment -log10(*p*-value) across the final MD-AD gene ranking for REACTOME pathways. Bars are
444     colored by the pathway's REACTOME category. We show all pathways with significant enrichment
445     (*p*<.01).  REACTOME pathways with long names are indicated by their REACTOME stable IDs. **(c)**
446     Comparison of top genes from MD-AD vs a linear correlation-based approach. For each ranking method,
447     we identify the top 1% of all genes and check their membership in REACTOME categories. For each
448     REACTOME category with at least 15 genes in the top 1% of MD-AD and/or correlation rankings, we
449     generate the following plot: each line represents a gene, with left endpoint at the percentile rank for MD-
450     AD and right endpoint at percentile rank for correlations. For clarity, we color the line purple if the gene
451     falls in the top 1% of both MD-AD and correlations, red if it is only in the top 1% of MD-AD, and blue if
452     it is only in the top 1% of correlations. Finally, the title indicates the ratio of MD-AD to correlation-based
453     top genes for the given REACTOME category.

454     **Figure 5.** MD-AD's top genes and their interactions with sex. **(a)** For the top 100 MD-AD genes, we
455     compute the significance of the interaction between expression and sex for its MD-AD score. The bars
456     indicate the gene's –log10(p-value) of the interaction term with sex (after FDR correction), and pathway
457     categories each gene belongs to are indicated below. A filled square indicates that the gene significantly
458     interacts with sex (*p<.05* after FDR correction), and an "x" marker indicates that it does not. **(b)** For
459     genes with significant sex interactions, we compute the significance of the overlap between REACTOME
460     category genes and sex-differential genes among:  Left:  all genes, and Right: the top 100 MD-AD genes
461     only. **(c)** For the top 100 MD-AD genes, we identify the genes with the most significant sex interaction
462     for MD-AD scores. We show the significance of the interaction (Top) and the interaction coefficients
463     (Bottom) for the top 6 most sex-differential genes. Each gene's MD-AD rank is indicated in their x-axis
464     labels **(d)** For the top 6 most-sex differential top 100 MD-AD genes, we display scatter plots of
465     expression by MD-AD score, coloring each sample by sex of the donor.

466     **Figure 6.** MD-AD's reliance on microglial cluster genes and gene set signatures. **(a)** Bars indicate the
467     gene's –log10(p-value) of the interaction term with sex (after FDR correction), and gene membership in
468     microglial cluster gene sets from Olah et al.[30] is indicated below. A filled square indicates that the gene
469     significantly interacts with sex (*p<.05* after FDR correction),  and an "x" marker indicates that it does not.
470     **(b)** For genes with significant sex interactions, we compute the significance of the overlap between
471     microglial cluster genes and sex-differential genes among:  Left:  all genes, and Right: the top 100 MD-
472     AD genes only. **(c)** Gene set enrichment -log10(*p*-value) across the final MD-AD gene ranking for cell
473     type signatures.[8]

474     **Figure 7.** MD-AD's transfer performance for blood gene expression data sets. **(a)** Shows t-test statistics
475     comparing average predicted neuropathology between individuals with mild cognitive impairment (MCI;
476     Left) and Alzheimer's dementia (AD; Right) vs. cognitively normal (CTL) individuals. **(b)** Box plots
477     show the differences in predicted neuropathology for blood samples from individuals stratified by age
478     group and cognitive status. Significant differences are shown in **Supplementary Figure 10b**. **(c)** t-SNE
479     embedding of last shared layer from MD-AD models trained for Blood1 and Blood2 datasets. Samples are
480     colored by their dataset (Left), cognitive status (while brain samples are shown in grey; Middle), and
481     predicted neuropathology score (Right).

482

483  **SUPPLEMENTARY METHODS**

484  **1. DATA PROCESSING**

485  For developing the MD-AD model, we used data from the following RNA-Seq and neuropathology
486  datasets available through the AMP-AD Knowledge Portal: (1) Adult Changes in Thought (ACT) [15],
487  (2) Mount Sinai Brain Bank (MSBB) [16], and (3) Religious Orders Study/Memory and Aging Project
488  (ROSMAP)[6,13,14]. Details of sample collection and sequencing methods are described in previously
489  published work [6,13–16]. We pooled together brain gene expression data from the temporal cortex, parietal
490  cortex, hippocampus, and forebrain white matter from ACT, Brodmann areas 10, 22, 36, and 40 from
491  MSBB, and the dorsolateral prefrontal cortex from ROSMAP. To avoid confounding conditions, we
492  excluded samples from individuals who had neuropathological diagnoses other than AD. Taken together,
493  the studies provide 1,758 gene expression samples.

494  In order to compile gene expression samples across the three cohorts, we retain expression levels for
495  genes which are present in all datasets. Within each dataset, we exclude genes with null values for over
496  two-thirds of samples. Before combining datasets, we log-transformed the expression values and then
497  normalized them for each gene to vary between 0 and 1. We then combined the gene expression datasets
498  and performed batch effect correction with ComBat[17] to reduce systematic differences across studies
499  (**Supplementary Figure 1b-c**)[17]. The resulting dataset contains 1,758 gene expression samples, each with
500  14,591 genes measured.

501  Next, for each gene expression sample, we incorporated the available corresponding neuropathology
502  labels: (1) **Aβ IHC**: amyloid-β protein density via immunohistochemistry, (2) **plaques:** neuritic amyloid
503  plaque counts from stained slides, and (3) **CERAD score:** a semi-quantitative measure of neuritic plaque
504  severity[38], (4) **τ IHC:** abnormally phosphorylated τ protein density via immunohistochemistry, (5)
505  **tangles:** neurofibrillary tangle counts from silver stained slides, and (6) **Braak stage:** a semi-quantitative
506  measure of neurofibrillary tangle pathology [19]. Detailed descriptions for each phenotype within each
507  dataset are provided in **Supplementary Table 1**. Because Braak stage and CERAD score are global
508  measurements of neuropathological damage, if an individual had multiple available gene expression
509  measurements from different regions, they each sample was labeled with the same Braak and CERAD
510  values. However, Aβ-IHC and τ-IHC were provided for several brain regions for both ROSMAP and
511  ACT studies. Therefore, each expression sample was labeled with the Aβ-IHC and τ-IHC measurements
512  for the same or nearest region. Because the available plaques label provided by MSBB was averaged over
513  several brain regions, we similarly used ROSMAP's average plaques and tangles labels (aggregated from
514  several regions) for consistency with MSBB's metrics (see **Supplementary Table 1**). Finally, for
515  consistency across datasets, we first normalized all neuropathological variables to vary between 0 and 1
516  before combing datasets.

517

518  **2. COMPUTATIONAL METHODS**
519      **A. Review of previous approaches**

520  Post-mortem transcriptomic studies have investigated molecular phenotypic and neuropathological
521  outcomes in AD. Early work in this domain examined simple correlations among gene expression and
522  AD symptoms[10] or compared gene expression levels across AD-patients versus controls[11]. More recently,
523  more systematic network-based analyses have contributed to the understanding of AD biology. In

13

particular, Zhang et al. [7] constructed molecular networks based on bulk gene expression data separately for individuals with and without AD, and identified modules with remodeling effects in the AD network. More recently, Mostafavi et al. [6] used co-expressed genes in the aging human frontal cortex to build a single molecular network and identified modules related to AD neuropathological and cognitive endophenotypes. Using single-cell RNA sequencing data, Mathys et al. [8] clustered cells within brain cell-types to identify and characterize AD-related cellular sub-populations. Each of these approaches have been applied to single cohorts. Until recently, a unified and robust modeling of AD neuropathology based on brain gene expression has been hindered by relative scarcity and regional heterogeneity of brain gene expression datasets. One possible solution is to combine multiple data sets to gain statistical power. The collection of postmortem brain RNA-sequencing datasets, assembled by the AMP-AD (**A**ccelerating **M**edicines **P**artnership **A**lzheimer's **D**isease) consortium, provides new opportunities to combine multiple data sets. However, such heterogeneous datasets pose challenges to many methods, which must account for inter-study differences. In a recent attempt, Logsdon et al.[9] used a meta-analysis approach to identify co-expressed modules separately for 7 brain regions across 3 datasets, then subsequently applied consensus methods to identify modules that were conserved across multiple regions and studies. As of now, we're not aware of any methods that directly model all data in a unified way.

## B. The MD-AD Model

MD-AD (**M**ulti-task **D**eep learning for **A**lzheimer's **D**isease neuropathology), is a *unified framework for analyzing heterogeneous AD datasets* to improve our understanding of expression basis for AD neuropathology (**Figure 1**). Unlike previous approaches, MD-AD learns a single neural network by jointly modeling multiple neuropathological measures of AD severity phenotypes, and hence can incorporate data collected from multiple datasets. This *unified* framework has key advantages over separately trained models. First, MD-AD allows sparsely labeled data, which is a natural characteristic of datasets aggregated through consortium efforts (**Figure 1e**). Even if different phenotypes only partially overlap in the measured samples, each sample contributes to the training of both phenotype-specific and shared layers. Predicting multiple phenotypes at once biases shared network layers to capture relevant features of these AD phenotypes at the same time. This is of critical importance: each phenotype represents a *different type* of noisy measurement of the same underlying true biological process, and as we demonstrate by joint training MD-AD is able to average out the noise to extract the true hidden signal. Additionally, the increased sample size enables MD-AD to capture complex non-linear interactions between genes and phenotypes. In contrast, Multi-layer perceptrons (MLPs) offer another powerful approach for directly capturing complex relations between gene expression and a phenotype. However, training separate MLPs for each phenotype (**Supplementary Figure 1a**) has limited scope: it can utilize only the samples measured for a specific phenotype, and it cannot share information across related phenotypes. We demonstrate that these advantages improve MD-AD prediction accuracy, enabling it predictions to generalize across species and tissue types (**Figure 1b**). As illustrated in **Figure 1a**, the MD-AD network jointly predicts six neuropathological phenotypes from gene expression input data via shared hidden layers followed by task-specific hidden layers.

## 3. TRAINING & EVALUATING MD-AD

As described above, we build the MD-AD model in Python using the TensorFlow and Keras packages. In order to have efficient and robust training and to reduce overfitting, we apply a principal component

566     analysis (PCA) transformation to the data and use resulting top 500 principal components – a 500-
567     dimensional representation of our 14,591 gene expression values – as the input to the MD-AD and all
568     baseline models. For comparison to MD-AD, we generate six analogous MLP networks with un-shared
569     representations, and six linear models containing no hidden layers, to serve as baseline models (see
570     **Supplementary Figure 1a**).

571     In order to robustly evaluate the performance of the models, we segment the dataset into five parts, and
572     each part is treated as a test set once. Within each of the five training and test split splits, each model
573     architecture was trained and hyperparameter-tuned using five-fold cross validation within the training set.
574     We then train each model with the best hyperparamters found by cross validation using the full training
575     set before performance was evaluated on the corresponding test set (see **Supplementary Figure 1d**).
576     Thus, prediction performance reported in the results section are the average of these five test performance
577     values. For training the models, we use a mean squared error (MSE) loss function applied to each
578     phenotype prediction. For the MLP and linear baselines, parameters of the networks are updated via back-
579     propagation for 200 epochs from the mean-squared error (MSE) of the network's prediction on the given
580     variable's label among training batches. Similarly, MD-AD's parameters are also updated via back-
581     propagation, with the loss function calculated as the sum over MSEs across all six prediction tasks
582     (masking losses for missing phenotypes). For MD-AD, we explored several different options for
583     architectures with different amounts of shared and task-specific layers (**Supplementary Figure 2b-c**).
584     We selected the final architecture (shown in **Figure 1a**) because we wanted to have multiple hidden
585     layers in both the shared portion and task-specific portion of the network to allow for non-linear
586     interactions to be learned in both the shared representation and in the task-specific branches, and
587     **Supplementary Figure 2b-c** shows that alternatives to this approach tended to perform similarly or
588     worse.

### A. Internal test-set validation

590     As described above, for each training and test split, we use five-fold cross-validation to make modeling
591     choices for the MD-AD model and baselines before training each model with the full training set and
592     reporting and reporting test MSEs (averaged over all five test splits). We evaluate model performance in
593     two ways: (**1**) standard train and test sets, and (**2**) ROSMAP test performance for different subsets of the
594     available datasets.

595     First, separately for each of our five cross validation training sets, we calculate the final test MSE on the
596     corresponding hold-out set. To test whether these effects are significant, for each baseline method, we
597     performed one-sided paired t-tests to determine whether there is a significant difference between the
598     baseline method's error and MD-AD's across the five test folds (**Figure 2a**). Next, in order to evaluate
599     the contributions of each dataset to prediction performance, we performed the above procedure with
600     different subsets of available datasets. Because ROSMAP is the only dataset with all available
601     phenotypes, we evaluate performance specifically on ROSMAP. In **Figure 2b**, we show ROSMAP test
602     samples' MSE performance when trained on all subsets of ACT, MSBB, and ROSMAP training samples
603     (following the same cross-validation procedure described above).

### B. External dataset validation (Human)

605     In order to evaluate MD-AD's ability to generalize to out of sample data, we assessed performance on
606     three datasets: Mount Sinai Brain Bank Microarray (MSBB-M; N=1,053), Harvard Brain Tissue

607  Resource Center (HBTRC; N=460), and Mayo Clinic Brain Bank (N=323). These datasets were collected
608  from AMP-AD, but were left out of the original MD-AD training because they were microarray samples
609  or lacked many neuropathology labels.

610  After normalizing gene expression samples from external data sets in the same way as described for the
611  ACT, MSBB RNA Seq, and ROSMAP datasets, we then adjust the expression values to have similar
612  distributions to our batch corrected training data sets. We evaluated the MD-AD model on our new
613  processed data to obtain predictions for all six phenotypes. Because these three external datasets provide a
614  sparse set of neuropathological labels, we do not have access to labels for many of the six MD-AD labels.
615  Instead, we evaluated whether MD-AD's predictions were consistent with the (binary) neuropathological
616  diagnosis of AD, by aggregating MD-AD's various neuropathology predictions into one "neuropathology
617  score". The "neuropathology score" was produced by first calculating percentiles across samples (within
618  each dataset) for each neuropathological phenotype, then averaging over the six phenotypes.

619  **Figure 2c** shows that MD-AD provides the largest differences in neuropathology scores between
620  individuals with and without neuropathological diagnoses of AD. We further compared neuropathology
621  scores between AD and non-AD individuals split by age group (significance between groups shown in
622  **Supplementary Figure 3b**)

### C. Cross-species validation (Mouse)

624  To evaluate how well expression patterns predictive of neuropathology learned by MD-AD recapitulates
625  neuropathology in mouse models. To that end, we obtained gene expression data from Matarin et al. [20]
626  for 30 TASTPM mice which harbor double transgenic mutation in APP and PSEN1, as well as 76 wild
627  type mice. Data were quantile-normalized and log transformed. For this experiment, we mapped mouse to
628  human genes (via gene symbols) for a total of 7,057 intersecting genes between our training dataset and
629  the mouse expression data, which were again normalized to follow the same distributions as our MD-AD
630  training data. We retrained our MD-AD model on only these 7057 genes for all MD-AD samples and then
631  generated "neuropathology scores" for the mouse samples exactly as described in the previous section. As
632  with out-of-sample experiments described above, we compare MD-AD to MLPs and linear models in
633  separating neuropathology scores between TASTPM and wild type mice (Figure 2E). We also show
634  differences in neuropathology scores between different age groups (**Figure 2d, Supplementary Figure
635  3c**).

### D. Supervised embedding validation

637  The output of an intermediate layer of a neural network can be viewed as lower dimensional embedding
638  of the input features. In this paper, we focus on the last shared layer of the MD-AD network because it is
639  a supervised embedding of gene expression data which is influenced by all six training phenotypes. We
640  evaluate the embedding compared with those generated by both singly-trained MLPs as well as
641  unsupervised methods (i.e., K-Means and principal components analysis (PCA)) in two ways: (1) high
642  level visualization with t-SNE, and (2) evaluating the correspondence between individual nodes and AD-
643  related features.

644  *Visualizations with t-SNE*: For each of the MD-AD, MLP, and unsupervised models, we train the models
645  on the full combined dataset. For the deep learning models, we then generate "supervised" embeddings by
646  obtaining the output of the last shared layer (or analogous layer of the MLP model). For the unsupervised

647 methods, K-Means and PCA, we generate an embedding of 100 dimensions to be consistent with the MD-
648 AD and MLP models. After generating these embeddings for all samples, we then compress them to 2
649 dimensions via the t-SNE algorithm [21]. T-SNE Visualizations of MD-AD's supervised embedding are
650 shown in **Figure 3a** (left side for each phenotype), and the figure is replicated with six times, with each
651 plot showing samples colored by neuropathological phenotype severity for each of the six phenotypes.
652 For comparison, t-SNE visualizations for the singly-trained MLPs and unsupervised methods are shown
653 in **Figure 3c** (colored by CERAD Score only) and colored by other phenotypes and covariates of interest
654 in **Supplementary Figure 5**.

655 *Node-phenotype correlations*: To test whether MD-AD's embedding generalizes more to AD phenotypes
656 than the alternative methods, we compare the nodes that best capture each phenotype among MD-AD,
657 MLPs, and unsupervised methods. We perform the following analysis with the same five training and test
658 splits described earlier: for each of the six phenotypes used in MD-AD's training, we identify the node in
659 MD-AD's last shared layer whose output is most significantly correlated with that phenotype in the
660 training set. We then report the –log10($p$-value) (after FDR correction over nodes) for the correlation
661 between that node's output and the training phenotype in the test set, averaged across the train/test splits.
662 (**Figure 3a**, right side for each phenotype).

663 We also perform a similar analysis with higher-level AD phenotypes not used during model training:
664 dementia diagnosis (binary variable available in all datasets), last available cognition score (controlling
665 for age, sex, and education; only available for the ROSMAP dataset), and AD duration (i.e., time between
666 dementia diagnosis and death; available for the ACT and ROSMAP datasets). For this analysis, we report
667 the highest –log10(p-value) after FDR correction between nodes and the high-level phenotypes, average
668 over the five test sets (**Figure 3b**).

669

## 4. MODEL INTERPRETATION
### A. Constructing and annotating MD-AD consensus nodes (Figure S7)

672 Because deep neural networks have non-convex loss functions, randomness in our training procedure
673 produces networks with different weights from run to run. In order to capture robust nodes and highly
674 relevant genes, we repeat our training procedure 100 times, in order to simulate a "consensus network".
675 As shown in **Supplementary Figure 6a**, we construct "MD-AD consensus nodes" by clustering nodes
676 from many runs: (1) we train 100 MD-AD networks, (2) we obtain last shared layer node outputs for all
677 samples and normalize them (0-mean, unit variance), (3) we combine all nodes across all runs and then
678 cluster them using k-means (where the dimensions used to calculate similarity are samples) with k=50,
679 (4) we summarize each cluster of nodes by their medoid. Thus, for each sample, the MD-AD consensus
680 embedding is made up of 50 nodes which are medoids of clusters generated from 100 re-trainings.

681 In **Supplementary Figure 4b**, we provide a visual overview of the MD-AD consensus embedding
682 generated as described above. To provide a simple view of clusters, we select a subset of samples for
683 which we have clear high or low pathology, excluding ambiguous cases. We include (1) individuals with
684 Braak stage of at least 5 and CERAD scores at least 3 (i.e., "moderate"), or (2) individuals with Braak
685 stage of 3 or lower and a CERAD score of 1 (i.e., "absent") who are at least 85 years old and have no
686 dementia. Case 1 captures all individuals with pathologic AD diagnoses (with and without dementia),
687 whereas case 2 captures all individuals considered "resistant" to AD due to their old age but lack of

688  cognitive or neurological decline (consistent with previous literature, e.g. Latimer et al. (2019)).    To
689  annotate each node in the consensus embedding, we display their correlations with various phenotypes
690  and covariates, as well as their enrichment for REACTOME pathways.

691  *Correlations:* For each variable (neuropathological phenotypes, high-level AD phenotypes, and
692  covariates), we compute the correlation –log10(p-value) between the variable and each consensus node
693  output. In **Supplementary Figure 4c**, a high –log10(p-value) indicates that a node captures (or is highly
694  linearly related to) a variable.

695  *Pathway enrichment:* Beyond relationships between nodes and phenotypes, we annotated nodes with
696  which gene sets are relevant to their outputs.  For each of the fully trained MD-AD model, we use
697  integrated gradients (IG) [12] to obtain sample-level gene importances for each consensus node. Note that
698  each consensus node (as medoid within a cluster) is some node in one of the 100 re-training runs of MD-
699  AD, thus we perform integrated gradients for the specific node in that network. By generating sample-
700  level gene attributions for each sample, we are able to aggregate the absolute IG values across samples to
701  obtain average gene attributions for each gene on each node. For each MD-AD consensus node, this
702  method therefore provides us with a ranking over all genes by their importance. We then test for
703  enrichment of REACTOME pathways [40] in these gene rankings via gene set enrichment analysis (GSEA)
704  [22,23] to identify whether certain pathways seem to be involved in the activation these nodes. Enriched
705  pathways for the MD-AD consensus nodes are shown in **Supplementary Figure 4d**. **Supplementary
706  Table 2** provides detailed annotations for each node.

707      B.   **Identifying MD-AD's top genes**

708  In order to identify genes that drive MD-AD predictions, we used integrated gradients (IG) [12] to provide
709  importance estimates of each gene on the predicted outcomes. Again, in order to improve model stability,
710  we calculate gene rankings based on 100 re-trainings. After each run of training, we take our trained
711  model and apply IG for each sample to get the importance of each gene on each phenotype prediction. We
712  then calculate a weighted average by sample (weighted by relative pathology) to compute a global
713  importance value for each gene on each phenotype, where positive values indicate that high expression of
714  the gene relates to more severe AD phenotypes. Finally, by averaging over all phenotypes, we obtain our
715  final "IG score" the given round of MD-AD training. By averaging these score across 100 re-trainings, we
716  arrive at our "consensus IG score" for MD-AD. Negative scores imply that higher expression is
717  associated with less pathology, while positive scores imply that higher expression is associated with more
718  pathology, according to MD-AD.  We note that 100 re-trainings are more than enough to converge to a
719  stable gene ranking (**Supplementary Figure 6c**). The top genes for MD-AD are shown in **Figure 4a**, and
720  enriched REACTOME pathways in the top ranked MD-AD genes (via GSEA) are shown in **Figure 4b**.
721  The full gene ranking, generated separately for each phenotype, is provided in **Supplementary Table 4**.

722  For comparison with a linear gene ranking method, we also calculate the correlations between each gene
723  with each neuropathological phenotype (across all samples in our dataset), and then rank the genes by
724  their average correlation coefficients across all six phenotypes. Comparisons between REACTOME
725  categories represented in the top MD-AD vs correlation-based rankings are shown in **Figure 4c**.

726      C.   **Calculating nonlinear effects for MD-AD genes**

18

727    As a deep learning method MD-AD has the capacity to identify non-linear relationships among genes'
728    expression levels and neuropathological phenotypes. These non-linear relationships may reveal an
729    implicit capture of interaction effects with other covariates observable from expression data. Thus, we
730    sought to investigate the presence of interactions between sample-level covariates and specific genes in
731    their contributions to the MD-AD predictions. To monitor the presence of these interaction effects, we
732    modeled the consensus IG scores as a linear combination of a gene's expression level, a covariate of
733    interest, and the interaction of the two. Specifically, $score_{g,i} = a\ expr_{g,i} + b\ feat_i + c\ expr_{g,i}\ feat_i +$
734    $d$, where $score_{g,i}$ is the consensus IG value for gene $g$ and sample $i$, $expr_{g,i}$ is the sample $i$'s expression
735    level for gene $g$, and $feat_i$ is sample $i$'s value for the covariate. Based on this representation, we consider
736    there to be an interaction effect between a gene and feature on its importance in the MD-AD model if the
737    learned $c$ coefficient is statistically significant ($p<.05$, after FDR correction over all genes). We primarily
738    focus on identifying an interaction effects with sex ($feat_i = 1$ if sample $i$ comes from a male), and rank
739    interactions between genes and sex for MD-AD based on the –log10(p-value) of the interaction term.

740    *Gene set enrichment:* We evaluated whether sex-differential genes were enriched for the following gene
741    sets: (1) REACTOME pathways[40] and (2) microglial cluster gene signatures from a recent single cell
742    RNA Seq analysis of microglial cells from autopsied aging brains[30]. To evaluate whether the list of sex-
743    differential MD-AD genes are enriched for gene sets of interest, we use Fisher's exact tests to evaluate the
744    significance of overlap between all sex-differential genes and members of each gene set. Next, to evaluate
745    whether the top MD-AD sex-differential genes are enriched for the same gene sets, we perform Fisher's
746    exact tests again, but this time only consider the top 100 MD-AD genes in the calculations.

### 5.  BLOOD GENE EXPRESSION VALIDATION

748    To evaluate the ability of MD-AD to transfer to blood gene expression data, we downloaded publically
749    available AddNeuroMed cohort data from GEO (GSE63060 and GSE63061, which we refer to as Blood1
750    and Blood2, respectively). Details about the AddNeuroMed samples are provided in **Supplementary**
751    **Table 3**. As with the other validation datasets, each blood dataset was normalized such that each gene's
752    expression values have the same mean and variance as the processed MD-AD expression data. Because
753    each blood dataset had a different set of available genes, for each dataset, we re-trained MD-AD
754    consensus models for brain samples with only the genes available between them and blood samples
755    (12,104 and 11,392 genes for Blood1 and Blood2 respectively). Because these blood samples came from
756    living participants, we do not have access to the many neuropathology variables available across the brain
757    samples. Instead, we assess whether MD-AD's predictions align with individuals' cognitive diagnosis of
758    cognitively normal (CTL), mild cognitive impairment (MCI), or dementia.

759    We evaluate the effectiveness of the MD-AD model by comparing predicted MD-AD pathology scores
760    between CTL and MCI individuals, and between CTL individuals and individuals with dementia via two-
761    sided t-tests (together, and split by age). To evaluate the MD-AD embedding for blood samples,
762    separately for each blood dataset, we obtain the last shared layer embeddings of both the MD-AD brain
763    expression samples and blood samples from the first round of training.
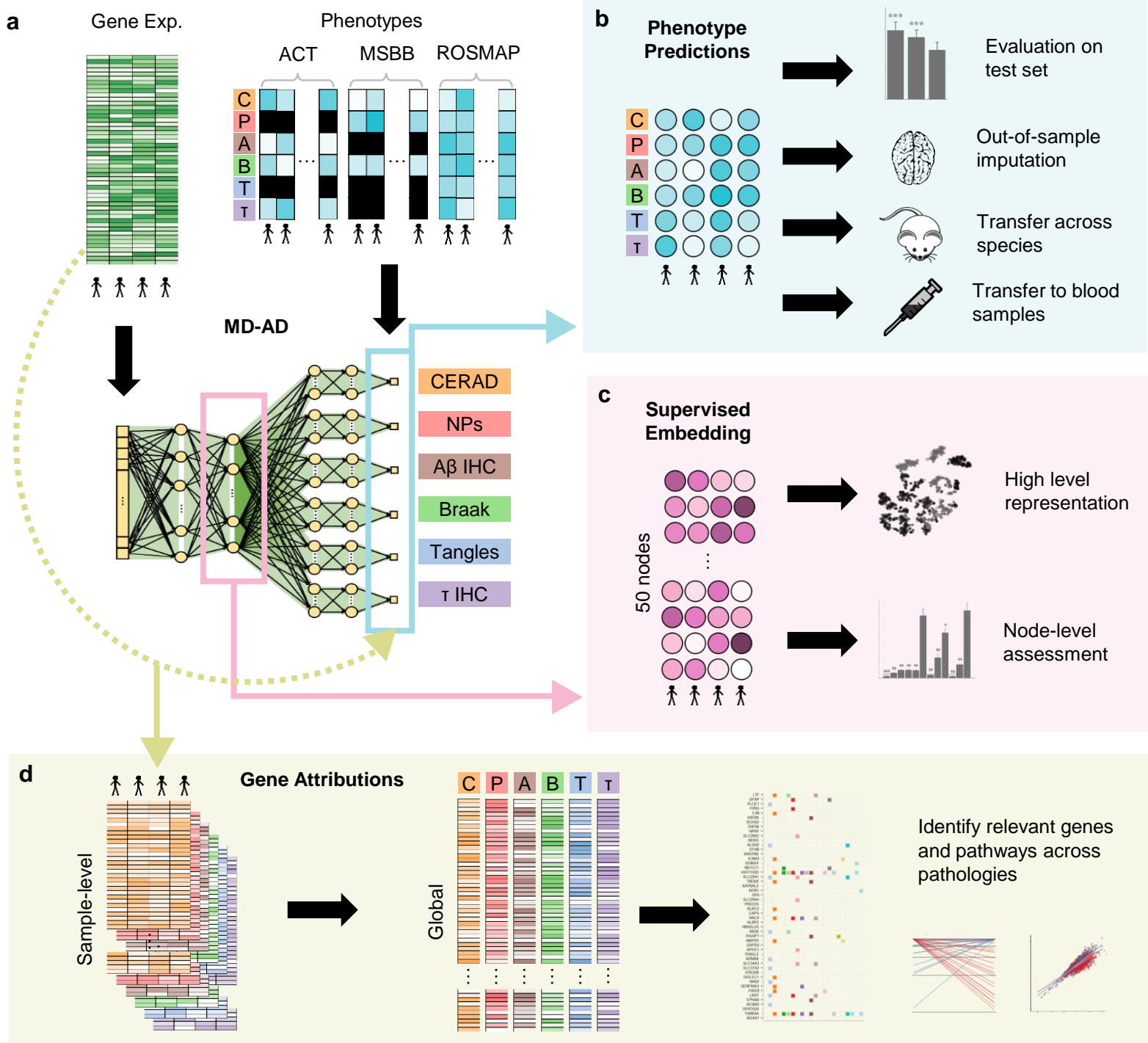
764

**REFERENCES**

1.    De Jager, P. L., Yang, H. S. & Bennett, D. A. Deconstructing and targeting the genomic architecture of human neurodegeneration. *Nat. Neurosci.* **21**, 1310–1317 (2018).

2.    Gaiteri, C., Mostafavi, S., Honey, C. J., De Jager, P. L. & Bennett, D. A. Genetic variants in Alzheimer disease-molecular and brain network approaches. *Nat. Rev. Neurol.* **12**, 413–427 (2016).

3.    Marioni, R. E. *et al.* GWAS on family history of Alzheimer's disease. *Transl. Psychiatry* **8**, 0–6 (2018).

4.    Jansen, I. E. *et al.* Genome-wide meta-analysis identifies new loci and functional pathways influencing Alzheimer's disease risk. *Nat. Genet.* **51**, 404–413 (2019).

5.    Kunkle, B. W. *et al.* Genetic meta-analysis of diagnosed Alzheimer's disease identifies new risk loci and implicates Aβ, tau, immunity and lipid processing. *Nat. Genet.* **51**, 414–430 (2019).

6.    Mostafavi, S. *et al.* A molecular network of the aging human brain provides insights into the pathology and cognitive decline of Alzheimer's disease. *Nat. Neurosci.* **21**, 811–819 (2018).

7.    Zhang, B. *et al.* Integrated systems approach identifies genetic nodes and networks in late-onset Alzheimer's disease. *Cell* **153**, 707–720 (2013).

8.    Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer ' s disease. *Nature* **570**, 332–337 (2019).

9.    Logsdon, B. A. *et al.* Meta-analysis of the human brain transcriptome identifies heterogeneity across human AD coexpression modules robust to sample collection and methodological approach. (2019). doi:10.7303/syn17114455

10.   Blalock, E. M. *et al.* Incipient Alzheimer's disease: Microarray correlation analyses reveal major transcriptional and tumor suppressor responses. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 2173–2178 (2004).

11.   Katsel, P., Li, C. & Haroutunian, V. Gene expression alterations in the sphingolipid metabolism pathways during progression of dementia and Alzheimer's disease: A shift toward ceramide accumulation at the earliest recognizable stages of Alzheimer's disease? *Neurochem. Res.* **32**, 845–856 (2007).

12.   Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. *34th Int. Conf. Mach. Learn. ICML 2017* **7**, 5109–5118 (2017).

13.   Bennett, D. A., Schneider, J. A., Arvanitakis, Z. & Wilson, R. S. Overview and findings from the religious orders study. *Curr. Alzheimer Res.* **9**, (2012).

14.   Bennett, D. A. *et al.* Overview and Findings from the Rush Memory and Aging Project. *Curr. Alzheimer Res.* **9**, 646–663 (2012).

15.   Miller, J. A. *et al.* Neuropathological and transcriptomic characteristics of the aged brain. *Elife* **6**, 1–26 (2017).

16.   Wang, M. *et al.* The Mount Sinai cohort of large-scale genomic, transcriptomic and proteomic data in Alzheimer's disease. *Sci. Data* **5**, 1–16 (2018).
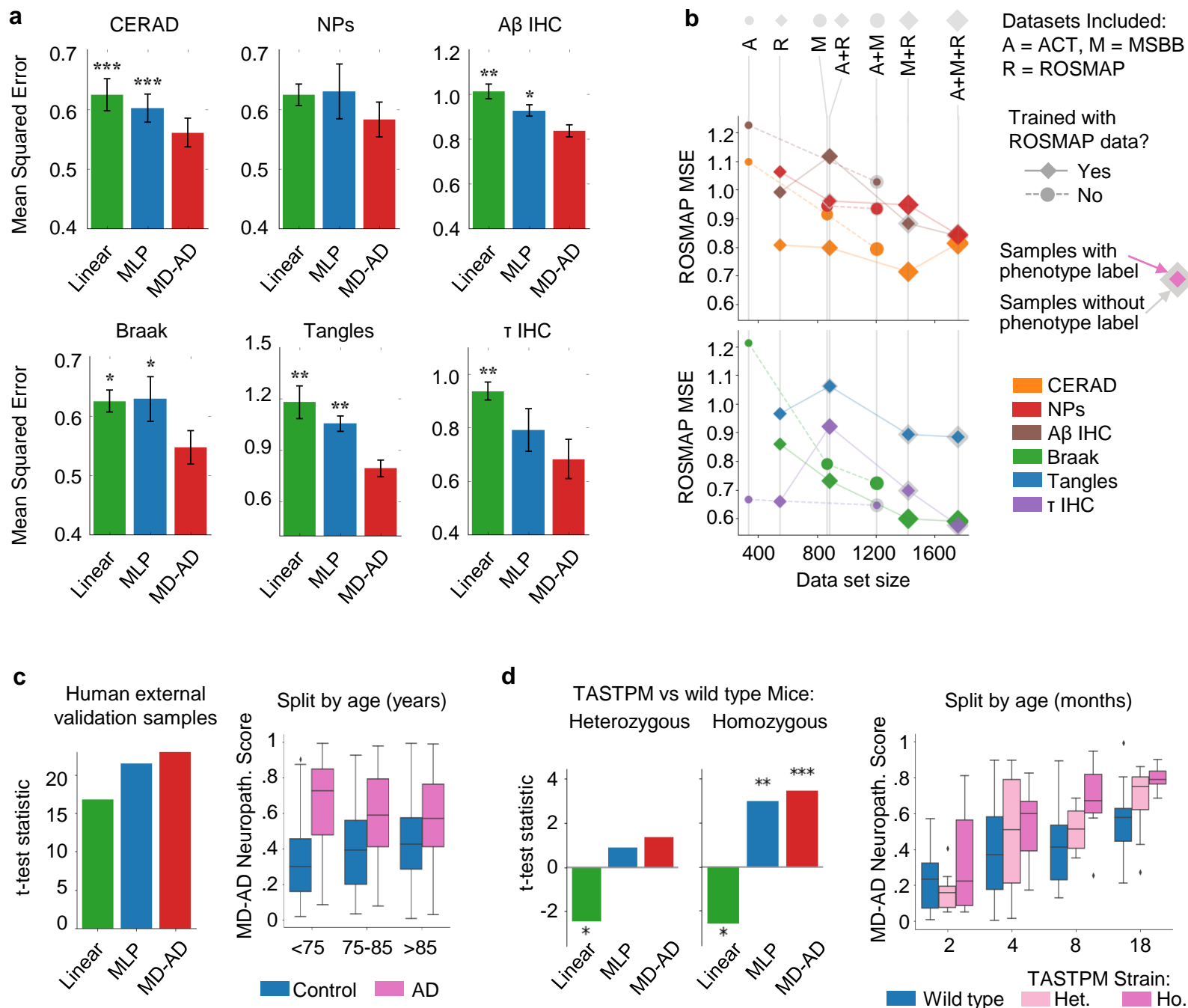
805 17.   Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression
806         data using empirical Bayes methods. *Biostatistics* **8**, 118–127 (2007).

807 18.   Mirra, S. S. *et al.* The consortium to establish a registry for Alzheimer's disease
808         (CERAD). Part II. Standardization of the neuropathologic assessment of Alzheimer's
809         disease. *Neurology* **41**, 479–486 (1991).

810 19.   Braak, H. & Braak, E. Neuropathological stageing of Alzheimer-related changes. *Acta
811         Neuropathol.* **82**, 239–259 (1991).

812 20.   Matarin, M. *et al.* A Genome-wide gene-expression analysis and database in transgenic
813         mice during development of amyloid or tau pathology. *Cell Rep.* **10**, 633–644 (2015).

814 21.   van der Maaten, L. & Hinton, G. Visualizing Data using t-SNE. *J. Mach. Learn. Res.* **9**,
815         2579–2605 (2008).

816 22.   Daly, M. J. *et al.* PGC-1α-responsive genes involved in oxidative phosphorylation are
817         coordinately downregulated in human diabetes. *Nat. Genet.* **34**, 267–273 (2003).

818 23.   Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based approach for
819         interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15545–
820         15550 (2005).

821 24.   Qiu, Y.-Q. KEGG Pathway Database. in *Encyclopedia of Systems Biology* (eds. Dubitzky,
822         W., Wolkenhauer, O., Cho, K.-H. & Yokota, H.) 1068–1069 (Springer New York, 2013).
823         doi:10.1007/978-1-4419-9863-7_472

824 25.   Trabzuni, D. *et al.* Widespread sex differences in gene expression and splicing in the adult
825         human brain. *Nat. Commun.* **4**, (2013).

826 26.   Olah, M. *et al.* A transcriptomic atlas of aged human microglia. *Nat. Commun.* **9**, 1–8
827         (2018).

828 27.   Chan, G. *et al.* CD33 modulates TREM2: Convergence of Alzheimer loci. *Nat. Neurosci.*
829         **18**, 1556–1558 (2015).

830 28.   Wang, Y. *et al.* TREM2 lipid sensing sustains the microglial response in an Alzheimer's
831         disease model. *Cell* **160**, 1061–1071 (2015).

832 29.   Mathys, H. *et al.* Temporal Tracking of Microglia Activation in Neurodegeneration at
833         Single-Cell Resolution. *Cell Rep.* **21**, 366–380 (2017).

834 30.   Olah, M. *et al.* Single cell RNA sequencing of human microglia uncovers a subset that is
835         associated with Alzheimer's disease. (2020).

836 31.   Lovestone, S. *et al.* AddNeuroMed - The european collaboration for the discovery of
837         novel biomarkers for alzheimer's disease. *Ann. N. Y. Acad. Sci.* **1180**, 36–46 (2009).

838 32.   Johnson, E. C. B. *et al.* Large-scale proteomic analysis of Alzheimer's disease brain and
839         cerebrospinal fluid reveals early changes in energy metabolism associated with microglia
840         and astrocyte activation. *Nat. Med.* **26**, 769–780 (2020).

841 33.   Lambert, J. C. *et al.* Genome-wide association study identifies variants at CLU and CR1
842         associated with Alzheimer's disease. *Nat. Genet.* **41**, 1094–1099 (2009).

843 34.   Farfel, J. M. *et al.* Relation of genomic variants for Alzheimer disease dementia to
844         common neuropathologies. *Neurology* **87**, 489–496 (2016).

845   35.   Chibnik, L. B. *et al.* CR1 is associated with amyloid plaque burden and age-related
846          cognitive decline. **69**, 560–569 (2011).

847   36.   Thambisetty, M. *et al.* Effect of complement CR1 on brain amyloid burden during aging
848          and its modification by APOE genotype. *Biol. Psychiatry* **73**, 422–428 (2013).

849   37.   Patrick, E. *et al.* A cortical immune network map identifies distinct microglial
850          transcriptional programs associated with beta-amyloid and Tau pathologies. *Press*

851   38.   Mirra, S. S. *et al.* The Consortium to Establish a Registry for Alzheimer's Disease
852          (CERAD): Part II. Standardization of the neuropathologic assessment of Alzheimer's
853          disease. *Neurology* **41**, (1991).

854   39.   Latimer, C. S. *et al.* Resistance and resilience to Alzheimer's disease pathology are
855          associated with reduced cortical pTau and absence of limbic-predominant age-related
856          TDP-43 encephalopathy in a community-based cohort. *Acta Neuropathol. Commun.* **7**, 9
857          (2019).

858   40.   Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–
859          D503 (2020).

860

**Figure 1**

| | # Expression Samples | CERAD | NPs | Aβ IHC | Braak | Tangles | τ IHC |
|---|---|---|---|---|---|---|---|
| **ACT** | 337 | ✔ | | ✔ | ✔ | | ✔ |
| **MSBB** | 879 | ✔ | ✔ | | ✔ | | |
| **ROSMAP** | 524 | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |

**Figure 2**

**Figure 4**

**Figure 6**