

# ***An explainable graph neural network approach for integrating multi-omics data with prior knowledge to identify biomarkers from interacting biological domains.***

Rohit K. Tripathy<sup>a,1</sup>, Zachary Frohock<sup>a,1</sup>, Hong Wang<sup>a</sup>, Gregory A. Cary<sup>b</sup>, Stephen Keegan<sup>b</sup>, Gregory W. Carter<sup>b,2</sup>, Yi Li<sup>a,2</sup>

<sup>a</sup>The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA

<sup>b</sup>The Jackson Laboratory, Bar Harbor, ME, USA

<sup>1</sup>These authors contributed equally.

<sup>2</sup>Corresponding authors: Gregory W Carter, The Jackson Laboratory, 600 Main St., Bar Harbor, ME, 04609, USA. Tel: 207-288-6025, email: [Gregory.Carter@jax.org](mailto:Gregory.Carter@jax.org); Yi Li, The Jackson Laboratory, 10 Discovery Dr, Farmington, CT, 06032, USA. Tel: 860-280-8685, email: [Yi.Li@jax.org](mailto:Yi.Li@jax.org)

## **Abstract**

The rapid growth of multi-omics datasets, in addition to the wealth of existing biological prior knowledge, necessitates the development of effective methods for their integration. Such methods are essential for building predictive models and identifying disease-related molecular markers. We propose a framework for supervised integration of multi-omics data with biological priors represented as knowledge graphs. Our framework is based on the use of graph neural networks (GNNs) to model the relationships among features from high-dimensional 'omics data and set transformers to integrate low dimensional representations of 'omics features. Furthermore, our framework incorporates explainability methods to elucidate important biomarkers and extract interaction relationships between biological quantities of interest. We demonstrate the effectiveness of our approach by applying it to Alzheimer's disease (AD) multi-omics data from the ROSMAP cohort, showing that the integration of transcriptomics and proteomics data with AD biological domain network priors improves the prediction accuracy of AD status and highlights robust AD biomarkers.

## **Introduction**

Advances in high-throughput technologies have led to an explosion in the generation and availability of molecular data, encompassing the analysis of diverse biomolecules such as DNA, RNA, proteins and metabolites (Schneider, 2011). This has, consequently, enabled the study of fundamental processes such as gene expression (Hrdlickova, 2017) and DNA methylation (Chen Y. R., 2018), and opened new avenues for understanding complex biological systems and disease mechanisms. Profiling multiple 'omics modalities in a disease cohort can provide a more comprehensive understanding of how distinct molecular processes operate in tandem to contribute to disease development and progression. Deriving such insights necessitates the development of methods for multi-modal integration. Indeed, suitably designed integrative analysis can not only improve predictive outcomes but also help identify novel therapeutic targets, enabling the development of personalized medicine (Günther, 2012).

Integrating and analyzing multi-omics datasets poses significant computational challenges. These datasets are typically high-dimensional and heterogeneous, making data reduction and identification of shared patterns essential. Additionally, omics coverage may be incomplete, leading to missing data and potential biases. A variety of unsupervised methods have been proposed to address these challenges and derive insight from multi-omics datasets. The standard approach to dealing with the high dimensionality of multi-omics data is to employ matrix factorization techniques. Methods such as

multi-omics factor analysis (Argelaguet, 2018), iCluster (Shen, 2009) and iNMF (Gao, 2021) look for latent factors shared across data modalities. Another class of unsupervised methods attempts to produce a unified representation of heterogeneous 'omics modalities by clustering samples based on similarities shared between their omics profiles – see, for instance, similarity network fusion (SNF) (Wang B. M., 2014) and unsupervised graph kernel learning approaches (Speicher, 2015; Mariette, 2018).

In spite of their applications to a variety of bulk and single-cell multi-omics datasets for discovering molecular mechanisms and identifying biomarkers (Vahabi, 2022), unsupervised methods do not allow one to detect signals or patterns pertinent to a specific target phenotype, such as a particular disease of interest. Meanwhile, methods for integrating heterogeneous data in the supervised setting are relatively sparse, where the challenge posed by the high dimensionality of multi-omics data is further compounded by the small dataset size (i.e., the number of patient samples is significantly smaller than the total number of biological molecules profiled); particularly in the bulk 'omics setting. Existing methods for supervised integration seek to exploit structures in 'omics datasets: patients with similar 'omics profiles are likely to share similar disease diagnoses. Based on this principle, several methods have been proposed to leverage graph neural networks (GNNs) to pose the task of patient phenotype prediction as a graph node classification problem. MOGONET (Wang T. S., 2021) leverages GNN feature extractors by using empirically generated patient similarity networks. MoGCN (Li X. M., 2022), on the other hand, learns a unified GNN model using a patient similarity graph topology generated with SNF and low-dimensional node features learnt through an autoencoder. While methods based on patient-similarity structures can alleviate the computational challenges associated with high-dimensionality in data features and low sample size, they do not leave any room for exploiting structures in the feature space, i.e. prior information about the relationship between biomolecules being measured.

In this work, we propose a novel explainable GNN framework, or GNNRAI (GNN-derived representation alignment and integration), for supervised integration of multi-omics data. Unlike existing methods, such as MOGONET and MoGCN, which use networks to model relationships among samples, we use graphs to model relationships among modality features (for example, genes in transcriptomics and proteins in proteomics data). This enables us to encode prior biological knowledge as graph topology. Given  $k$  'omics modalities, each sample is represented as  $k$  graphs in our framework. We leverage supervised GNNs to learn modality-specific low-dimensional embeddings. These low-dimensional embeddings are first aligned to each other to enforce shared patterns, and then integrated using a set transformer (Lee & Teh, 2019). The integrated multi-omics representations are used to predict the target phenotype. Our model architecture allows us to incorporate samples with incomplete 'omics measurements and avoid a reduction in statistical power. To identify predictive modality features, we employ the method of integrated gradients (Sundararajan, 2017) which estimates the importance of each feature to model predictions. We demonstrate the effectiveness of our framework by applying it to the task of predicting Alzheimer's disease (AD) status by integrating transcriptomics and proteomics data from the Religious Order Study/Memory Aging project (ROSMAP) cohort. Given that proteomics data typically have a much smaller number of features relative to transcriptomics data, exacerbated by the smaller number of samples with proteomic data in the ROSMAP cohort, multi-omics integration methods might mask the role of the proteomics modality (Yang, 2023). Our results show that proteomics data are more predictive than transcriptome data in the ROSMAP cohort and the integration of the two data modalities using our GNNRAI improves upon the predictive performance of the two unimodal models. Graph topology for our 'omics-specific GNNs is derived from recent work on AD biological domains (biodomains, or BDs), which are expertly curated knowledge graphs for AD-associated endophenotypes (Cary, 2024). Our modeling framework is compared to the MOGONET approach on held-out validation data and shows improved prediction metrics. We derive important

'omics features within AD-associated biodomains via the integrated gradients method. Finally, we probe our trained single-biodomain models to derive interactions between these biodomains using a set transformer in a second modelling stage. Interpretation of biodomain interactions via the method of integrated Hessians (Janizek, 2021) allows us to gain further insight into AD biology.

## Results

### ***GNNRAI for supervised multi-omics integration and biomarker identification***

In this work we developed an AI framework, GNNRAI (GNN-derived representation alignment and integration), for performing supervised multi-omics data integration, accommodating potentially incomplete data and identifying informative biomarkers and biological interactions. The backbone of our proposed method consisted of GNN-based feature extractor modules. Omics data, coupled with prior knowledge graphs, were processed through these GNN-based feature extractors to produce low-dimensional embeddings. Modeling the relationships between markers reduced the training sample size burden since correlation structure reduces the effective dimensions in high dimensional omics data. Leveraging prior pathway knowledge and integrating multi-omics data maximized the likelihood that the identified informative features were functional.

A schematic of our end-to-end GNNRAI model is shown in Figure 1. The MLP classifiers were designed for samples with a single modality, while the set transformer module was used exclusively for samples with complete multi-omics measurements. This architecture facilitated efficient training on incomplete multi-omics datasets, as the feature extractor modules were updated by all samples regardless of the completeness of their omics data. To explain the predictions from our model, we leveraged the integrated gradients method (Sundararajan, 2017), a method for *post-hoc* interpretability of black-box models. Furthermore, we used the method of integrated Hessians (Janizek, 2021) to extract informative biological interactions between single-biodomain model representations. Though we only demonstrated the integration of two modalities, it is straightforward to extend to multiple modalities.

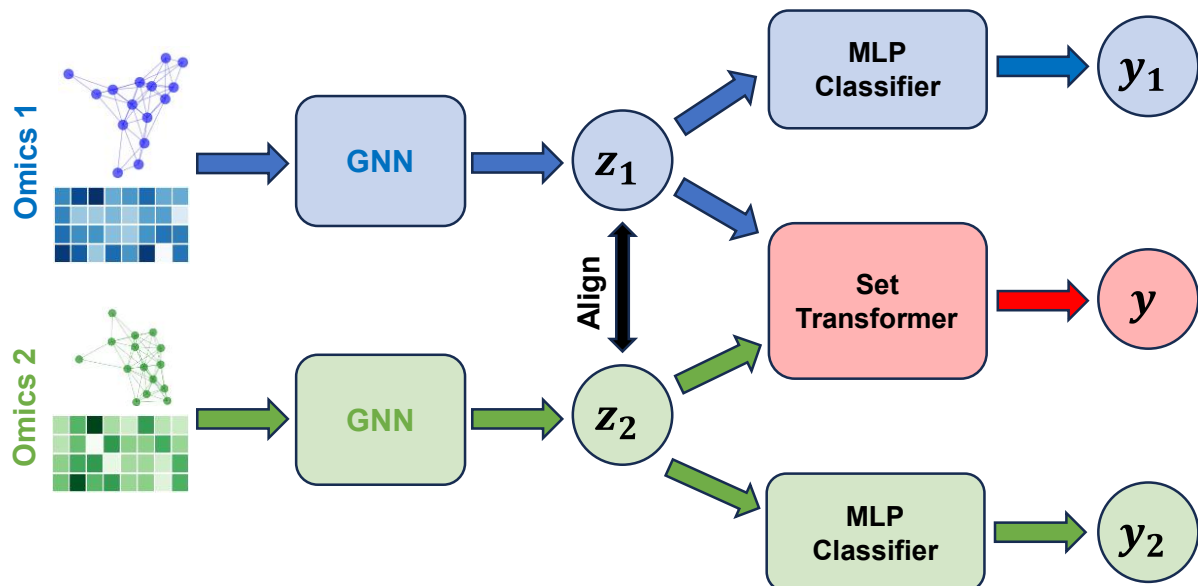


Figure 1: Schematic of our end-to-end integrative model GNNRAI. Data from individual 'omics modalities are processed in their respective GNN feature extractors to produce low-dimensional embeddings ( $z_1$  and  $z_2$ ).  $z_1$  and  $z_2$  are then aligned and integrated through a set transformer. They are also processed through separate MLP (multi-layer perceptron) classifiers to produce modality-specific predictions of the target when a sample has incomplete multi-omics measurements.

### ***Alzheimer's disease patient classification datasets.***

In our study, we implemented the GNNRAI framework to integrate transcriptomic and proteomic data for the binary classification of Alzheimer's disease within the ROSMAP cohort. We analyzed gene and protein data from the dorsolateral prefrontal cortex (DLPFC) brain region. The processing of the data and the criteria for AD diagnosis were elaborated in the Methods section. Building on Cary's 2024 research on AD biodomains (Cary, 2024), we created 16 datasets for AD classification, each representing a different biodomain. These datasets were complemented with knowledge graphs derived from querying the Pathway Commons database. Refer to the Methods section for details on biodomains and graph sizes for each biodomain. After data processing, we had 228 samples with both transcriptomic and proteomic data, 59 with only proteomic data, and 336 with only transcriptomic data. Our GNNRAI models were trained on each of these 16 biodomain-specific datasets. The datasets consisted of graphs with nodes representing genes or proteins from a biodomain, with their expression or abundance values as node features, structured by the biodomain's knowledge graph from querying the Pathway Commons. Each sample was labeled with a binary indicator to denote whether it was from an AD patient or a healthy control.

### ***Proposed GNNRAI AI framework outperformed benchmark MOGONET method on AD/control classification.***

The *multi-omics graph convolutional network*, or MOGONET (Wang T. S., 2021), framework was a supervised learning framework for integrating multi-omics data using GNNs. MOGONET processed individual modalities separately by constructing patient similarity networks using the cosine distance metric to assign edges. Graph neural networks operated on these patient similarity networks to make modality-specific predictions, which were then integrated through a view correlation discovery network (VCDN). In contrast to MOGONET, our approach imposed a network topology over the space of input features within each modality. The MOGONET architecture made it implausible to incorporate priors on the space of features (such as AD biodomains). Furthermore, MOGONET required samples to have complete measurements (i.e., no missing modalities). We trained our unimodal and integrative models on the set of samples with both transcriptomics and proteomics measurements for each of the 16 BDs and compared their validation predictive performance to that of MOGONET trained on the same datasets. A comparison of the validation performance between these models is shown in Figure 2. We observed that when trained on an equal number of samples, unimodal proteomics models

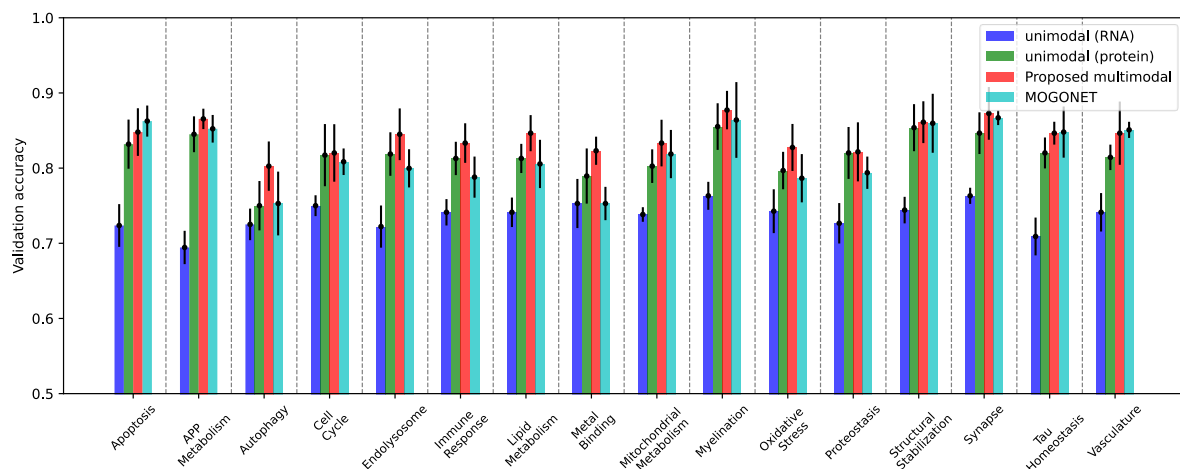


Figure 2: Validation performance of our proposed integrative model (red) compared to validation performance of the benchmark MOGONET model (cyan) on the set of common samples having both proteomics and transcriptomics measurements. The performance of the integrative models is also compared to unimodal GNN transcriptomics (blue) and proteomics (green) models.

consistently outperformed unimodal transcriptomics models. Despite having fewer features (see Table 3), proteomics data were more predictive in the ROSMAP cohort, aligning with (Johnson, 2022). Our integrative models outperformed the integrative MOGONET models in 13 out of 16 BD datasets (except for *apoptosis*, *Tau homeostasis* and *vasculature*). Additionally, for seven BD datasets (*cell cycle*, *endolysosome*, *immune response*, *lipid metabolism*, *metal binding*, *oxidative stress* and *proteostasis*), our unimodal proteomics models surpassed the multimodal MOGONET models. This was likely because proteomics and transcriptomics data were not always consistent, and MOGONET integrated modality-specific predictions rather than modality representations. In contrast, our framework's integration of transcriptomics and proteomics modalities improved the unimodal predictive performance across all 16 BD datasets, demonstrating the effective integration.

***Multi-omics GNNRAI models outperformed unimodal GNNRAI models trained on transcriptomics and proteomics alone.***

For samples with complete measurements, our model's performance was evaluated using the held-out validation set predictions from the integrative component (i.e., the set transformer in Figure 1). For samples with incomplete measurements, predictions were based on their respective unimodal classifiers. We found that integrating the two modalities resulted in better performing classifiers compared to the unimodal counterparts. This finding was significant given that we had 564 samples with transcriptomic data but only 287 with proteomic data. A larger number of less predictive transcriptomic samples could obscure the superior performance of proteomic samples if the useful information from both modalities was not effectively aligned and integrated.

To ensure a fair comparison, we used two sets of validation samples for evaluating the performance. The first validation dataset comprised samples with transcriptomics measurements (used to test the unimodal transcriptomics GNN models). The second validation dataset included samples with

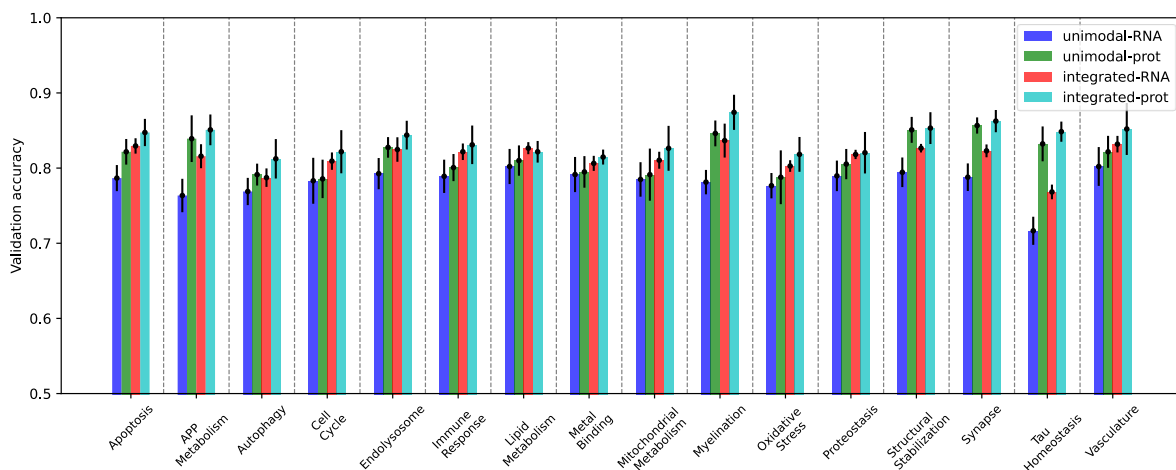


Figure 3: Validation performance of the integrative multimodal model (red and cyan) compared to performance of unimodal GNN models (blue and green) trained on the incomplete multi-omics datasets for 16 AD biodomains. The performance of the multimodal model is calculated on two sets of validation samples – the set of all validation samples with transcriptomics measurements and that with proteomics measurements.

proteomics measurements (used to test the unimodal proteomics models). Figure 3 compares the validation accuracy of GNN models trained solely on transcriptomics and proteomics samples with the end-to-end multi-modal models trained on all samples. Validation performance for unimodal models were denoted as ‘unimodal-RNA’ and ‘unimodal-prot’, while the two validation scores from the integrative model were denoted as ‘integrated-RNA’ and ‘integrated-prot’ respectively. In spite of the smaller training dataset, ‘unimodal-prot’ consistently outperformed ‘unimodal-RNA’ across all 16 BDs,

reiterating that proteomics data provided more AD-predictive information than transcriptome data in the ROSMAP cohort. When we compared the multimodal to unimodal performance, ‘integrated-RNA’ consistently surpassed ‘unimodal-RNA’ across all 16 BDs, demonstrating that integrating proteomics with transcriptomics data enhanced classification performance. Similarly, there was a consistent performance improvement from ‘unimodal-prot’ to ‘integrated-prot’, albeit less pronounced than in the RNA modality. Furthermore, ‘integrated-prot’ was generally better than ‘integrated-RNA’ except for *lipid metabolism* BD, despite the fact that the proteome-specific classifier was trained on only 59 samples compared to 336 samples for the transcriptome-specific classifier. This suggests that the target-predictive signals from transcriptomic and proteomic embeddings were aligned and integrated effectively, resulting in smaller performance differences between samples with transcriptomic data and proteomic data than in unimodal transcriptomic models.

#### **Validation on ROSMAP, Mount Sinai Brain Bank (MSBB) and Mayo Clinic transcriptomics and proteomics data**

To validate the predictive ability of our models trained on ROSMAP DLPFC samples, we curated samples with transcriptomics and/or proteomics measurements from the following studies and brain regions -

- 1) ROSMAP samples from the *anterior cingulate cortex* (ACC) and *posterior cingulate cortex* (PCC) brain regions with transcriptomic measurements.
- 2) MSBB samples from the *parahippocampal gyrus* (PHG), *frontal pole* (FP), *inferior frontal gyrus* (IFG), and *superior temporal gyrus* (STG) brain regions. Only PHG tissue had both transcriptomic and proteomic data, while the remaining tissues had transcriptomics only.
- 3) Mayo Clinic samples from the *temporal cortex* (TCX) brain region. Although TCX tissue had both transcriptomic and proteomic data, the proteomics measurements were acquired by label-free quantification, different from the tandem mass tag (TMT) quantification platform used in ROSMAP and MSBB. Hence, we did not validate ROSMAP-derived models on Mayo proteomics data.

Table 1 shows the sample counts of the curated validation datasets. The procedure for annotating MSBB and Mayo samples with ground truth labels were described in the Methods section. We noted

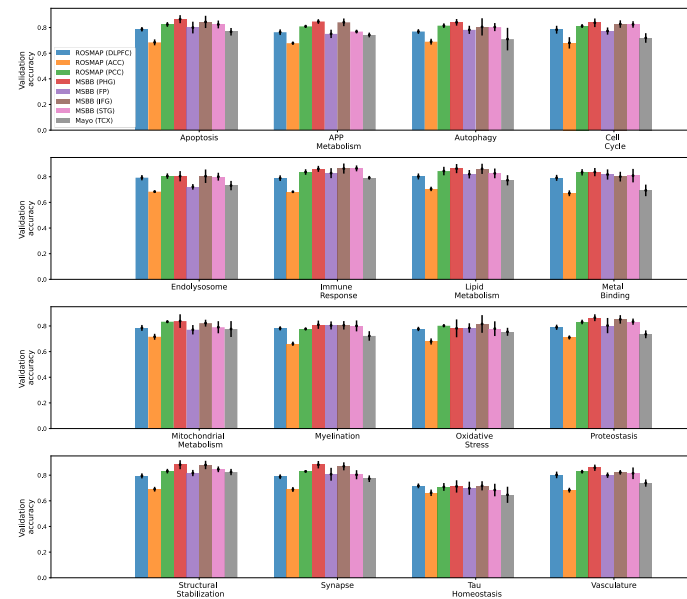
**Table 1: Transcriptomics and proteomics sample count in validation data**

Tissue	ROSMAP (ACC)	ROSMAP (PCC)	MSBB (PHG)	MSBB (FP)	MSBB (IFG)	MSBB (STG)	Mayo (TCX)
RNA	354	332	122	106	106	108	62
Protein	-	-	122	-	-	-	



that MSBB and ROSMAP adopted similar diagnostic criteria, while Mayo clinic cohort followed Mayo neurologist guidelines (McKhann, 1984).

We first applied the ROSMAP DLPFC-trained transcriptomics model to transcriptomics data from ROSMAP ACC and PCC brain regions, MSBB PHG, FP, IFG and STG brain regions and Mayo TCX brain region. For comparison, we also included the predictive accuracy on the ROSMAP DLPFC validation dataset (Figure 4). Figure 4 shows that the same GNN model has different predictive performance on different brain regions. Generally, ROSMAP PCC had slightly higher predictive accuracy than DLPFC, which was in turn higher than the predictive accuracy of ACC. For MSBB, PHG had the highest predictive accuracy, followed by IFG and STG, which predicted better than FP. MSBB PHG had the

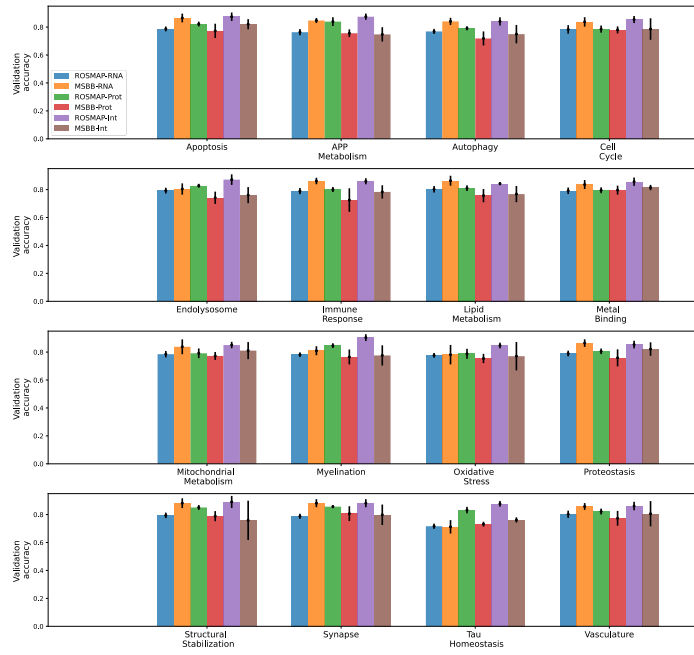


*Figure 4: Predictive performance of applying unimodal transcriptomics model trained on ROSMAP DLPFC training dataset to validation transcriptomics samples from various cohorts and brain tissues.*

highest predictive accuracy across the three cohorts. MSBB IFG and STG had comparable performance with ROSMAP PCC, while MSBB FP had comparable performance with ROSMAP DLPFC, which predicted better than Mayo TCX. ROSMAP ACC had the lowest predictive accuracy on average. Nevertheless, ROSMAP ACC predictive accuracy was above 0.6 across the 16 AD biodomains, better than a random guess. The different predictive performance on transcriptomic data from different brain regions might be explained in terms of neuropathological burdens – the FP and DLPFC regions are impacted at a similar disease stage, whereas the PHG tends to be affected much earlier in disease progression. This could be further exacerbated by the differences in cohort

sample selection between ROSMAP and MSBB. MSBB samples were selected for multi-omics profiling based on the presence of remarkable AD neuropathology, whereas the ROSMAP study was a longitudinal cohort and did not pre-select samples for extreme neuropathology. Thus, disease signatures learned in the ROSMAP DLPFC samples could be amplified in the MSBB samples, leading to an increase in the predictive accuracy of our models. This also implied that transcriptomic signatures might be translational across relevant brain tissues.

Next, we applied the ROSMAP DLPFC-trained proteomics and integrative models to proteomics and multi-omics data from MSBB PHG brain region respectively (Figure 5). We also included ROSMAP DLPFC and MSBB PHG transcriptomics performance in Figure 5 for comparison. MSBB PHG had lower predictive accuracy than ROSMAP DLPFC for proteomics and integrative models. Unlike in ROSMAP, where proteomics data had higher predictive accuracy than transcriptomics data and the integrative model improved upon the two unimodal models, proteomics data were less predictive than transcriptomics data, and the integrative model had lower predictive accuracy than the proteomics model for 9 AD biodomains in MSBB. This might be due to the fact that the GNN models were trained



*Figure 5: Predictive performance of applying unimodal and multimodal models trained on ROSMAP DLPFC training dataset to samples from ROSMAP DLPFC validation and MSBB PHG validation datasets. Prot: protein; Int: integrative.*

permuting the training labels and trained our integrative models on the permuted datasets. The resulting models were called null models. The importance scores for all genes/proteins in a graph from each sample of the 300 null models were used as background test statistics. Since we calculated FDR

on a much smaller set of ROSMAP proteomics data ( $287 \times 2/3 = 191$  samples) than ROSMAP transcriptomics data ( $564 \times 2/3 = 376$  samples), causing poor generalization performance on unseen data.

### ***Identification of biomarkers relevant to AD***

We applied the integrated gradients method to our trained multi-omics GNNRAI models to derive importance scores on input graph nodes (genes and proteins). We used a permutation-based approach to determine the importance score threshold by controlling the false discovery rate (FDR) to be below 0.05. Like standard permutation procedure for multiple hypothesis testing, we treated the original importance scores as the observed test statistics, generated 300 permuted datasets by randomly



rather than corrected p-values, the estimated empirical FDR was confidently accurate for B=100-200 (Millstein, 2013). See Methods for the procedure to calculate permutation-based FDR.

We determined whether a gene/protein was informative only in correctly predicted validation samples. To rank the informative genes/proteins identified across these analyses, we added the total number of non-overlapping sample IDs for which a given gene/protein was identified as informative across modalities and AD biodomains for each study, divided by the total number of correctly predicted validation samples in each study, then calculated the average fraction of informative samples across studies/brain tissues, based on which genes/proteins were ranked. The top 20 AD-predictive

Rank	Gene Symbol	Mean fraction samples	Total samples (ROSMAP DLPFC)	Total samples (MSBB PHG)
1	<i>MDK</i>	0.979	129	81
2	<i>IL1B</i>	0.700	92	58
3	<i>VGF</i>	0.616	65	61
4	<i>FLNA</i>	0.563	41	67
5	<i>ICAM1</i>	0.496	46	53
6	<i>DCN</i>	0.473	35	56
7	<i>COL1A1</i>	0.473	43	51
8	<i>NTN1</i>	0.451	39	50
9	<i>APOB</i>	0.439	39	48
10	<i>CD44</i>	0.433	34	50
11	<i>IQGAP3</i>	0.426	50	39
12	<i>LTF</i>	0.423	38	46
13	<i>PRKCD</i>	0.399	38	42
14	<i>GFAP</i>	0.395	37	42
15	<i>LGMN</i>	0.393	51	33
16	<i>OLFM4</i>	0.393	38	41
17	<i>TAC1</i>	0.387	25	48
18	<i>APP</i>	0.382	61	25
19	<i>NOX4</i>	0.380	46	34
20	<i>CYP51A1</i>	0.375	38	38

Table 2: Top 20 informative genes across modalities, tissues, and AD biodomains. Total samples are the count of unique sample identifiers for which the gene is identified as informative to model prediction.

genes/proteins are shown in Table 2. The top ranked gene in these analyses was *MDK*, which was informative in the binary classification task for 210 unique validation samples (average of 97.9% of correctly predicted validation samples). *MDK* is a secreted growth factor that has consistently been identified among a suite of matrisome proteins that associate with A $\beta$  plaques (i.e., Module M42 in (Johnson, 2022; Drummond, 2022)), and has been shown to influence the aggregation of amyloid-beta, both *in vitro* and *in vivo* (Levites, 2023). *VGF* was another top-ranked gene in these analyses and

has consistently been identified as a robust biomarker for AD (Tandon, 2023; Watson, 2023), as well as being a top predicted regulator of multiscale AD networks (Beckmann, 2020).

The three genes in the top 20 that had the highest integrated AD Target Risk Score (TRS (Cary, 2024)), were *APP*, *LGMN* and *LTF*. Each of these genes were informative for over 80 total samples and had TRS in the top 2% of all scored genes. *APP* is a well-known disease gene that is the proteolytic precursor of the A $\beta$  peptide, which is a major component of one of the hallmark neuropathologies of the disease, and variants within *APP* are causal for rare autosomal dominantly inherited forms of AD. *LTF*, or lactotransferrin, has recently been identified as a predictor of A $\beta$  burden (Tsatsanis, 2021). *LGMN*, also known as  $\delta$ -secretase, is an asparagine endopeptidase that is involved in the cleavage of both tau (Zhang, 2014) and *APP* (Yao, 2021) proteolysis, which is linked to increased pathogenicity in each case. This corresponded with the findings from the individual gene/protein analyses where *LGMN* was the most impactful in Tau Homeostasis and APP Metabolism biodomains.

There were also two genes among the top 20 that were novel candidate biomarkers and did not have prior publications directly linking their functions to AD pathogenesis. For example, *IQGAP3* was ranked #11 in this analysis, had a TRS in the top 10%, and was differentially expressed in both transcriptomics and proteomic samples. Despite having no publications where *IQGAP3* is implicated in AD, it was linked

with cytoskeletal maintenance and neurite outgrowth (Wang S. W., 2007) which is consistent with its role in the Structural Stabilization domain in these analyses. *OLFM4* was another example of a highly ranked gene in these analyses (#16), with limited evidence in the literature linking it with AD.

At least seven genes among the top 20 were strongly related to AD biology, showing that our integrative method identified functional features due to the integrated prior biological pathway knowledge.

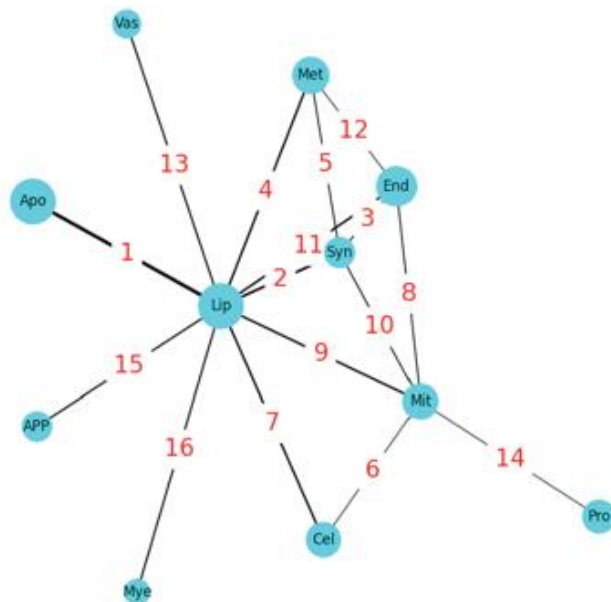


Figure 6: Interactions between biodomains based on Integrated Hessians analysis on multi-modal model trained on samples with complete multi-omics. The nodes represent biodomains and edges represent interactions with the edge annotations being the rank of the interaction. The biodomain names based on the 3-letter abbreviations are: apoptosis (apo), lipid metabolism (lip), synapse (syn), endolysosome (end), metal binding (met), cell cycle (cel), mitochondrial metabolism (mit), myelination (mye), vasculature (vas), proteostasis (pro) and APP metabolism (app).

### Detecting interactions among biodomains

For a multi-modal integrative model trained on a given biodomain, the class token representation from the final set transformer was a single low-dimensional embedding that unified information across modalities within the biodomain. We, therefore, collected class token representations in integrative models from all 16 BDs, and trained an auxiliary set transformer to

integrate information across biodomains. Integrated Hessians (Janizek, 2021) was applied to this second set transformer model to derive interaction scores between its input tokens. The biodomains partition gene functions into distinct molecular endophenotypes. However, these endophenotypes can and do interact during the etiology of the disease. Therefore, a primary goal of utilizing the biodomain framework is to identify interactions between domains that could broaden our understanding of the disease development.

The top interactions detected are shown as a graph where each node is a BD in Figure 6. Interactions were ranked by repeating the model training and informative interaction identification process ten times with different random weight initializations. From each iteration, the top ten percent of interactions, determined by ranking the number of samples for which each interaction was informative, were noted. Interactions present in the top ten percent three or more times out of ten were recorded. Rank was determined first by the number of appearances in the top ten percent, then, in the case of ties, by the total number of samples in which the interaction was informative. The domain nodes with the largest degree were Lipid Metabolism (degree = 9), followed by Mitochondrial Metabolism (degree = 5), Synapse and Endolysosome (degree = 4, each). The centrality of these domains was supported by the observation that Synapse, Lipid Metabolism, and Mitochondrial Metabolism were among the top risk-enriched biodomains (Cary, 2024). The observation that Lipid Metabolism was a hub in this graph suggested that aspects of Lipid Metabolism influenced many other disease processes. The centrality of Lipid Metabolism to AD pathogenesis was supported by myriad observations from recent decades, including genetic studies that implicate Lipid Metabolism associated genes (e.g. *APOE*, *CLU*, *ABCA7*, *SORL1*) in driving AD risk (Bellenguez, 2022), the observation that amyloid-beta production occurs in lipid raft membrane microdomains (Ehehalt, 2003), recent lipidomic studies that identify changes in lipid species that are specific to the disease (Baloni, 2022; Batra, 2023), and many more.

Other interaction relationships represented in this graph were informative and supported by evidence from the literature. As mentioned above, the edge between APP Metabolism and Lipid Metabolism was supported by the influence of lipid rafts on amyloid-beta production. The link between Lipid Metabolism and Mitochondrial Metabolism was also very well supported given that  $\beta$ -oxidation of fatty acids, which is the primary catabolic pathway, occurs in the mitochondrial matrix. Given that mitochondria provide the requisite energy and precursor metabolites for cell cycle progression, and that mitochondrial biogenesis and dynamics are influenced by cell cycle regulators, the interdependence of Cell Cycle and Mitochondrial Metabolism was also well supported. The top ranked edge between Lipid Metabolism and Apoptosis evoked ferroptosis, which is an iron-dependent cell death mechanism that is distinct from apoptosis but involves the accumulation of peroxidated lipid species (Yan H. F., 2021) and is the focus of newly emerging hypotheses of disease pathogenesis (Wang F. W., 2022). Further, supporting this association is the observation that 7 of the top 20 identified informative genes – i.e. *APOB* (Wu, 2024), *LGMN* (Chen C. A., 2021; Yan L. H., 2023), *LTF* (Xiao, 2022; Wang Y. L., 2020), *NOX4* (Park, 2021), *CD44* (Liu, 2019; Ye, 2024), *PRKCD* (Lv, 2024), and *CYP51A1* (Li Y. R., 2024) – are associated with regulating aspects of ferroptosis in diverse contexts. It was noted that the immune response biodomain, which is strongly involved in AD, is absent in Figure 6. The interplay between the immune response and lipid metabolism BDs ranked seventeen in these analyses – the highest ranking interaction not included in Figure 6. The reason for this omission in our list of top pairwise interactions remains unclear and warrants further investigation which is beyond the scope of our current work.

## Discussion

In this work, we proposed an end-to-end AI framework, or GNNRAI, for supervised alignment and integration of multi-omics data with prior information expressed as knowledge graphs. Our method was based on the use of GNNs for learning low-dimensional embeddings from high-dimensional data and could accommodate samples with missing modalities. Using the ROSMAP data, we curated 16 binary classification datasets – each dataset comprising a view of the ROSMAP gene expression and protein abundance data within an AD-associated biodomain. We noted that the size of biological domains varied significantly from smallest to largest domains and therefore allowed us to test the

robustness of our approach to varying input dimensionality. Our approach has validated its efficacy in the task of integrating transcriptomics and proteomics data from the ROSMAP cohort. It has outshined the benchmark MOGNET method in 13 of the 16 BDs and shown improvements over the two unimodal models for all 16 BDs. This outcome is noteworthy considering the disparity in sample sizes, with 564 transcriptomic samples and only 287 proteomic samples. The abundance of less predictive transcriptomic samples could potentially conceal the enhanced performance of the proteomic samples unless the valuable insights from both data types are properly synchronized.

The task of integrating multi-omics data is computationally challenging for several reasons. We demonstrated that the curse of dimensionality arising from the large number of 'omics features relative to the sample size could be overcome by leveraging the correlation structure in graphs and message passing in GNNs. Additionally, we showed that the separation of feature extraction modules and set transformer-based integration allowed us to utilize samples with missing modalities – a characteristic feature of multi-omics datasets.

Our framework allows one to integrate modalities where prior information about the relationship between input features can be expressed in the form of knowledge graphs. We leveraged existing work on AD biodomains to extract network topologies for transcriptomics and proteomics modalities. Our approach did not make a distinction in the structure of the knowledge graphs used for these modalities, thereby implicitly imposing a simplified assumption that network relationships between transcripts is exactly reproduced in the proteins they code for. Furthermore, we did not incorporate data from other modalities within the ROSMAP study, such as methylomics and metabolomics, due to the current unavailability of direct prior knowledge graphs for these modalities. However, existing transcriptomic and proteomic networks can be leveraged for the construction of gene-centered methylomics and metabolomic knowledge graphs. For instance, metabolites catalyzed by the same genes/proteins may be determined to share a relationship. For genes that have an edge in a transcriptomic network, CpG sites within their regulatory regions (promoters, enhancers etc.) may be determined to share the same edge within the corresponding methylation network. Finally, we identified informative features through the model-agnostic method of integrated gradients which derived importance scores on individual graph nodes independently. Interpretation of GNN predictions can, in theory, be enhanced by using an explanation method to identify informative subgraphs, or *motifs*. While methods in this direction did exist (Ying, 2019), our experience was that we were unable to extract meaningful subgraph/motifs through the application of such methods. How to identify correlated informative features efficiently is one of the important future research directions for GNNs.

## Methods

### *Data preprocessing*

To investigate AD mechanisms, we adopted a combination of clinical and neuropathological criteria used in (Johnson, 2022) to assign ground truth labels (AD case or control) to patients within the ROSMAP cohort. In particular, we used clinical cognitive tests, such as MMSE (the Mini-Mental State Examination (Folstein, 1975)), or CDR (Clinical Dementia Rating) to assess dementia: MMSE score  $\leq 24$  or CDR  $\geq 1$ . Neuropathological assessment of patients was conducted post-mortem using Braak staging (Braak, 2006) and CERAD (The Consortium to Establish a Registry for Alzheimer's Disease) scoring (Wolfsgruber, 2014) to reflect AD hallmarks. CERAD scores 0-3 correspond to no AD/none, possible/sparse, probable AD/moderate, and definite/frequent, respectively. The Braak score, or seven Braak stagings, classifies the severity and distribution of tau pathology in the brain. Cases with CERAD 0-1 and Braak 0-3 without dementia at last evaluation were annotated as controls (if Braak score

equals 3, then CERAD must equal 0); cases with CERAD 2–3 and Braak 3–6 with dementia at last evaluation were annotated as AD.

Downloaded RNA-Seq count data were log2 transformed and corrected for age, sex and postmortem interval (PMI) covariates. Downloaded protein abundance data were log2 transformed and median zero centered per feature. Finally, age, sex and postmortem interval (PMI) were regressed out.

For validation MSBB samples, we used similar diagnostic criteria, except MMSE was replaced with CDR (Wang M. B., 2018) since MSBB did not provide MMSE information. MSBB gene expression/protein abundance data were processed similarly to ROSMAP data. In addition, patient race was regressed out.

For patients in the Mayo study, AD and controls were taken to be the reported diagnosis according to Mayo neurologist guidelines, as described in (McKhann, 1984). In contrast to ROSMAP and MSBB, which made use of tandem mass tag (TMT) quantification, Mayo proteomics data were acquired with label-free quantification, hence we did not validate our models on Mayo proteomics data. Mayo gene expression data were processed similarly to ROSMAP data.

### ***Network priors from Alzheimer’s disease biological domains***

Table 3: Summary of AD biodomain knowledge graphs				
Biological domain	Nodes (mRNA)	Edges (mRNA)	Nodes (Proteins)	Edges (Proteins)
APP Metabolism	147	316	87	154
Apoptosis	958	10963	451	3261
Autophagy	505	2665	304	1096
Cell Cycle	728	6561	337	1779
Endolysosome	1096	8909	665	3874
Immune Response	1507	16354	688	4683
Lipid Metabolism	1671	13515	891	4366
Metal Binding and Homeostasis	801	4822	441	1398
Mitochondrial Metabolism	1273	7092	858	3493
Myelination	189	348	120	132
Oxidative Stress	346	2080	185	637
Proteostasis	2675	29812	1477	13323
Structural Stabilization	2245	23850	1297	12163
Synapse	2067	21056	1234	9686
Tau Homeostasis	45	90	41	81
Vasculature	756	7098	356	2008

The prior biological knowledge ascribed to nodes and edges in the knowledge graphs used for the analysis was derived from publicly accessible biological databases. These graphs provided a topological organization to the biological domains, which were 19 AD-associated endophenotypic descriptors, such as immune response and mitochondrial metabolism (Cary, 2024). The biological domains were lists of functional biological definitions describing aspects of AD, and were defined with suites of relevant Gene Ontology (GO) terms (Ashburner, 2000). Each GO term was annotated with a set of genes, and biological processes within a domain that were enriched for composite metrics of disease risk (Cary, 2024) could be identified using standard enrichment procedures. We used significantly-enriched GO terms (gene set enrichment analysis adjusted p-value  $< 0.01$  and normalized enrichment score  $> 1.7$ ) – 16 of the 19 biological domains had GO terms that met these criteria – and extracted the leading edge genes from each term to seed knowledge graph generation through a pathway reconstruction pipeline. We performed a shortest path reconstruction among all risk-enriched genes for each domain using protein-protein interaction (PPI) edge annotations from the Pathway Commons database (Cerami, 2010), version 13. Given the nonlinear relationships implicated in most biological interactions, the shortest path to connect two genes was selected for creating an edge between two nodes. For a given protein, expressed by a gene in the biological domain, an edge was derived from the larger PPI network. The final network object consisted of edges, which were the PPI, and nodes, which were the GO term-derived gene list. A summary of the sizes of the transcriptomics and proteomics prior networks within each biodomain is shown in Table 3.

### Modeling framework for multi-omics integration

In this work, we proposed an end-to-end framework for supervised integration of incomplete multi-omics data. Our modeling framework comprised two key components: 1) graph neural network-based feature extractors and 2) feature alignment as well as set transformer-based feature integration among modalities.

#### Graph Neural Network-based feature extractors

Let  $x_i \in R^{d_i}$  be the set of features for the  $i^{th}$  modality and  $\mathcal{G}$  be an undirected graph with  $d_i$  nodes,

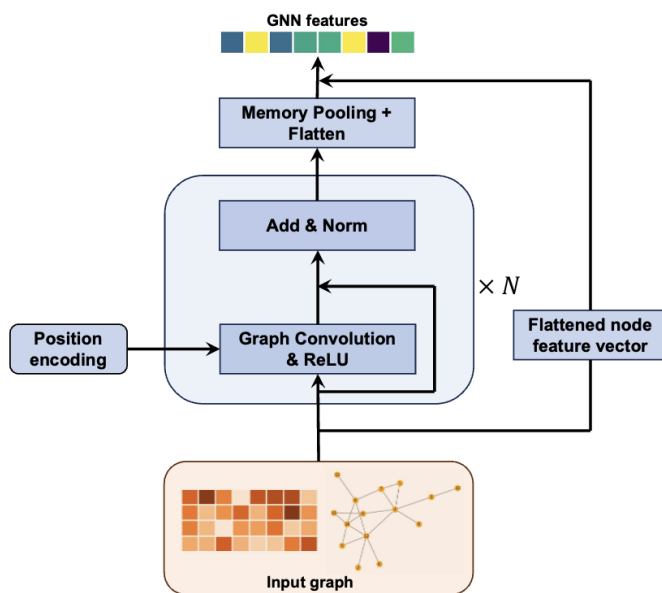


Figure 7: Schematic of GNN-based feature extractor. The feature extractor comprised of a stack of graph convolution blocks, a memory pooling layer and a residual connection between the original input node features and the output of the memory pooling layer.

with each node representing a feature in the  $i^{th}$  modality.  $\mathcal{G}$  may be constructed empirically via binarizing the matrix of correlation coefficients among features or may represent prior knowledge on the space of features defining known pairwise relationships. Let  $\mathcal{E}$  be the list of edges in  $\mathcal{G}$  and  $\mathcal{W}$  be an optional list of corresponding vector-valued edge weights such that  $|\mathcal{W}| = |\mathcal{E}|$ . Given a sample of 'omics measurements  $x_i$  and an associated graph topology  $\mathcal{G}$ , we set up a graph neural network (GNN) and learned  $z_i = g(x_i, \mathcal{E}, \mathcal{W})$ , where  $z_i \in R^m$  is a vector of low-dimensional embeddings. A schematic of the GNN-based feature extractor is shown in Figure 7.

First, we constructed a positional encoding of  $d_i$  features into  $k$  ( $k < d_i$ ) learnable communities using a standard lookup table embedding function in *Pytorch* which were softmax-normalized. The input node feature vector  $x_i$  to  $g$  was linearly transformed by using the  $k$  communities as bases, followed by passing through a stack of graph convolution blocks. Each block (round-corner rectangular box in Figure 7) comprised two sub-blocks. The first sub-block was a message passing graph convolution layer, followed by a ReLU activation. The second sub-block was a residual connection adding input node features to the output of ReLU, followed by a batch normalization. Input node features which passed through  $N$  graph convolution blocks were transformed into latent node features  $\tilde{x}_i \in \mathbb{R}^{d_i}$ . The associated graph topology remained unchanged since we did not employ edge-updating in our model. The transformed node feature vector  $\tilde{x}_i$  passed through a memory pooling layer (Hosein Khasahmadi, 2020) which learned a coarse graph representation through soft cluster assignments. Subsequently,  $\tilde{x}_i$  was reduced to  $\tilde{z}_i \in \mathbb{R}^{d_i \times d_c}$ , where  $d_c$  is the number of clusters and  $d_i'$  is the number of features in a cluster, which was then flattened. Lastly, a residual connection was employed adding the linearly transformed original input node feature vector to the flattened  $\tilde{z}_i$ , which was subsequently linearly mapped to  $z_i \in \mathbb{R}^m$ .

### Graph convolution layer

Let  $h_{in} \in \mathbb{R}^{R \times C_{in}}$  denote the input matrix to our graph convolution layer, where  $R$  was the number of nodes, and  $C_{in}$  represented the number of features in an input node. Let  $C_{out}$  be the number of features in an output node of a graph convolution layer. Suppose  $R$  nodes were mapped to  $k$  communities, and the softmax-normalized positional encoding were  $p \in \mathbb{R}^{R \times k}$ . We first mapped the input node features linearly:  $h' = h_{in}W$ , where  $W \in \mathbb{R}^{C_{in} \times \hat{C}_{out}}$  were learnable weights,  $\hat{C}_{out} = k \cdot C_{out}$ ,  $h' \in \mathbb{R}^{R \times \hat{C}_{out}}$ . Linearly transformed node features  $h'$  were then reshaped to  $h' \in \mathbb{R}^{R \times k \times C_{out}}$ . The positional encoding was then used to linearly weight node features along each output channel -  $\hat{h}_{:,i} = \sum h'_{:,i} \odot \text{unsqueeze}(p)$ , where  $\odot$  is the element-wise product operation and the summation was among  $k$  communities. The resulting node features were  $\hat{h} \in \mathbb{R}^{R \times C_{out}}$ . Finally, a message-passing update on the node features was performed:

$$h_{out,j} = b + \sum_{r \in \{j\} \cup \mathcal{N}_j} \frac{\text{unsqueeze}(\hat{h}_j) e_{r,j}}{\hat{d}_r \hat{d}_j},$$

where  $\hat{h}_j \in \mathbb{R}^{C_{out}}$  was the feature vector of the  $j^{th}$  node,  $\mathcal{N}_j$  was the set of indices of the neighbor nodes to node  $j$ ,  $\hat{d}_j$  and  $\hat{d}_r$  were the degrees of nodes  $j$  and  $r$  respectively, and  $e_{r,j} \in \mathbb{R}^{n_e}$  was a vector of  $n_e$  attributes for the edge connecting nodes  $j$  and  $r$ . Note that our graph convolution layer was an extension of the standard graph convolution of (Kipf, 2016) where we allowed for the inclusion of positional embeddings and vector-valued edge attributes.

### Aligning embeddings among modalities

Before integration, we calculated the absolute pairwise correlation (for samples with complete multi-omics data) or the absolute pairwise cross-correlation (for samples with incomplete multi-omics data) among modality representations and took its negative value as a loss component.

### Integration using set transformers

Let  $z_i \in \mathbb{R}^m$  be the GNN embeddings corresponding to the  $i^{th}$  'omics modality. We collected embeddings from GNNs corresponding to each individual modality into a set  $Z = (z_0, z_1, \dots, z_K)^T \in \mathbb{R}^{(K+1) \times m}$ , where  $K(\geq 1)$  was the number of modalities and  $z_0 \in \mathbb{R}^m$  was a set of learnable parameters called the *class token*. We then used the standard transformer encoder architecture



described in (Vaswani, 2017) to set up an integrative classifier. The prediction from the integrative classifier,  $\hat{y} \in \mathbb{R}^C$  (where  $C$  was the number of classes) was given by -

$$\hat{y} = \text{MLP}(\text{Encoder}(Z)_0),$$

where  $\text{Encoder}(Z)_0 \in \mathbb{R}^m$  was the latent representation of the input learnable class token  $z_0$ , and  $\text{MLP}(\cdot): \mathbb{R}^m \rightarrow \mathbb{R}^C$  was a fully connected neural network which mapped the final class token representation to the target label space. The transformer encoder was a composition of  $n$  encoder blocks i.e.,  $\text{Encoder}(Z) = E_n \circ E_{n-1} \dots \circ E_1(Z)$ . The sequence of operations within each encoder block was as follows:

1. Multi-head self-attention (MSA), residual connection and layer normalization (LN) (Ba, 2016) on the set of 'omics tokens -  $Z' = \text{LN}(Z + \text{MSA}(Z))$ ,
2. Position-specific feedforward network followed by a residual connection and layer normalization (LN) -  $Z_{out} = \text{LN}(Z' + \text{FFN}(Z'))$ , where  $\text{FFN}(\cdot): \mathbb{R}^m \rightarrow \mathbb{R}^m$  was a one hidden layer MLP with ReLU activation and operated on each token individually, i.e.,  $\text{FFN}(Z) = \text{FFN}(z_0, z_1, \dots, z_K) = (\text{FFN}(z_0), \text{FFN}(z_1), \dots, \text{FFN}(z_K))^T$ .

### Full model architecture and training

Our full integrative multi-omics model was expressed as:

$$y = h(g_1(x_1, \mathcal{E}_1, \mathcal{W}_1), g_2(x_2, \mathcal{E}_2, \mathcal{W}_2), \dots, g_K(x_K, \mathcal{E}_K, \mathcal{W}_K)),$$

where,  $h(\cdot)$  was a set transformer which integrated feature representations generated by the 'omics GNNs  $g_i$ s,  $x_i$ s were the node features for the  $i^{th}$  'omics modality and  $\mathcal{E}_i, \mathcal{W}_i$  were the edge index list and edge attribute list associated with the  $i^{th}$  modality. A schematic of the architecture is shown in **Error! Reference source not found.**. Given multi-omics data for a given sample, measurements for each modality were processed through their respective GNN modules. The parameters of the GNN were trained by classifying the embeddings to the corresponding target labels using multi-layer perceptrons (MLP) or a set transformer which collected embeddings for all available modalities, integrated them, then made a prediction on the target label.

### Training with complete samples

We first describe the training of our model when the multi-omics data set had complete samples, i.e., we had measurements in all available 'omics modalities for all patients. We split the total dataset into three folds, maintaining the ratio of sample labels in each fold as in the whole dataset. Samples from two of these splits were chosen as training samples while samples from the remaining split were used as validation samples. Since samples in our dataset had binarized labels (AD/control) we therefore used the binary cross entropy loss function. We set up our end-to-end integrative model as detailed in the previous sections and trained it using stochastic gradient descent optimization, specifically the Adam optimization method (Kingma, 2014), using mini-batches of training data. The objective function we optimized was as follows –

$$\mathcal{L} = \sum_{i=1}^K \mathcal{L}_i + \lambda_1 \mathcal{L}_{\text{int}} + \lambda_2 \sum_{i>j} \mathcal{L}_{i,j}^{\text{align}} + \lambda_3 \mathcal{L}_{\text{reg}},$$

where  $\mathcal{L}_i$  was the loss incurred on the predictions made by the MLP classifier on the embeddings of the  $i^{th}$  modality,  $\mathcal{L}_{\text{int}}$  was the loss incurred on the predictions made by the integrative module, i.e., the set transformer,  $\mathcal{L}_{i,j}^{\text{align}}$  was the alignment loss between the  $i^{th}$  and  $j^{th}$  modalities, and  $\mathcal{L}_{\text{reg}}$  were

norm penalties on the learnable parameters in the model.  $\lambda_1, \lambda_2, \lambda_3$  were trade-off hyper-parameters. Furthermore, we applied sample weights on our prediction losses ( $\mathcal{L}_i$  and  $\mathcal{L}_{\text{int}}$ ) to account for any potential class imbalance.

### ***Training with incomplete samples***

To account for incomplete samples (i.e., samples with missing measurements in one or more modalities), we split our total training dataset into disjoint subsets based on modality representation. For our transcriptomics and proteomics datasets, we split the total dataset into three subsets – the set of common samples and two additional sets comprising samples having measurements in transcriptomics or proteomics alone. Each epoch of training for our integrative model now comprised a single epoch through each disjoint subset. The common sample subset was trained with the full loss function as described in the previous section. The loss function was modified for training data subsets with missing modalities. For a subset of the training data with measurements in a single modality only, the components  $\mathcal{L}_{\text{int}}$  and  $\mathcal{L}_{i,j}^{\text{align}}$  were set to 0. The order in which the disjoint data subsets were processed to update the model weights during training was randomized from epoch to epoch.

### ***Hyperparameter selection***

As described in the previous section, we created 3 stratified splits of our dataset, from which we picked one fold for validation and used the remaining two for training. For any given combination of hyperparameters, we trained our models 3 times per split, each with a different random initialization, and cross-validated on the validation dataset. In total, we trained the model 9 times with the same hyperparameters. We averaged the validation performance of our model on each of these 9 trials and reported it as the predictive performance corresponding to a given setting of hyperparameters. We performed a grid search over our hyperparameters and picked the hyperparameters with the highest average validation accuracy over 9 trials.

### ***Biomarker identification and interaction analysis***

Given the complex, nonlinear nature of deep compositional models, explaining model predictions and thereby elucidating important features and feature interactions is nontrivial. Here we described the method of integrated gradients for extracting important biomarkers and the method of integrated Hessians for deriving informative interactions between AD biological domains.

#### ***Integrated gradients***

Let  $x = (x_1, x_2, \dots, x_d)^T \in \mathbb{R}^d$  be the input features to a deep learning model,  $f: \mathbb{R}^d \rightarrow \mathbb{R}^C$ , where  $C$  was the number of classes of the target label. Let  $f_c(x)$  be the model output score for the  $c^{\text{th}}$  class. The components of the gradient vector,  $\nabla_x f_c$ , represented the sensitivity of the class score to small perturbations of the input features, and the magnitude of its components may be interpreted as a proxy for the importance of input features. The integrated gradient attribution of the  $i^{\text{th}}$  input feature on the  $c^{\text{th}}$  class was defined as follows:

$$\phi_i = (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial f_c(x' + \alpha(x - x'))}{\partial x_i} d\alpha,$$

where  $x'$  was a user-defined baseline input. The integrand represented the model evaluation at an input constructed by the linear interpolation between the baseline  $x'$  and the true input  $x$ . The choice of the baseline was task dependent. For instance, in image classification tasks, where the input  $x$  is a tensor representing pixel intensities, it is customary to pick the zero tensor as the baseline, i.e.,  $x' =$

0. In our multi-omics model, inputs were node features representing gene expression/protein abundance levels. We set our baseline as the average expression/abundance level in our training control samples. Features with large magnitude attribution scores on the disease class label (i.e., class label 1) were implicated as informative disease biomarkers.

### **Integrated Hessians**

Integrated Hessians for pairwise feature interactions (Janizek, 2021) is a natural extension of the method of integrated gradients. Let  $\phi_i: \mathbb{R}^d \rightarrow \mathbb{R}$  be the integrated gradient attributions on the  $i^{th}$  input feature. Applying integrated gradient explanations on the function  $\phi_i(x)$  resulted in a new  $d$ -dimensional vector whose  $j^{th}$  component  $\Gamma_{i,j}(x) = \phi_j(\phi_i(x))$  explained the contribution of feature  $x_j$  to the model attributions on feature  $x_i$ . We interpreted the quantity  $\Gamma_{i,j}(x)$  as the interaction between features  $i$  and  $j$ . For  $i \neq j$ , the feature interaction scores were given by:

$$\Gamma_{i,j}(x) = (x_i - x'_i)(x_j - x'_j) \times \int_{\alpha=0}^1 \int_{\beta=0}^1 \alpha\beta \frac{\partial^2 f_c(x' + \alpha\beta(x - x'))}{\partial x_i \partial x_j} d\alpha d\beta,$$

where  $x'$  was a baseline input. The self-interaction term  $\Gamma_{i,i}(x)$  was given by:

$$\Gamma_{i,i}(x) = \phi_i(x) - \sum_{i \neq j} \Gamma_{i,j}(x),$$

where the  $i^{th}$  feature integrated score  $\phi_i(x)$  represented the marginal contribution of  $x_i$  to model prediction. The self-interaction score  $\Gamma_{i,i}(x)$  was, thus, defined as the difference between the marginal contribution of  $x_i$  and every pairwise interaction involving  $x_i$ .

### **Identify informative markers or marker interactions with specified false discovery rate (FDR)**

To determine the importance score threshold above which a marker or a marker interaction was viewed informative, we adopted the permutation approach to compute the empirical FDR. Specifically, we randomly permuted the order of the ground truth labels to generate  $B$  permuted datasets of ground truth labels. We then trained the GNN model on each of the  $B$  permuted datasets and computed importance scores for markers using integrated gradients or marker interactions using integrated Hessians. Following the non-parametric procedure outlined in (Xie, 2005), we calculated the false discovery rate for a given threshold  $d$  as:

$$FDR(d) = \pi_0 \frac{\left( \sum_{b=1}^B \frac{\#\{i: z_i^{(b)} > d\}}{B} \right)}{\#\{i: Z_i > d\}},$$

where,  $Z_i$  was the importance score of marker or interaction  $i$  in the unpermuted data,  $z_i^{(b)}$  was the importance score of marker or interaction  $i$  in the  $b^{th}$  permutation,  $\pi_0$  was the prior probability that a marker or an interaction was uninformative. For AD study,  $\pi_0$  is close to 1, and we took the value of 0.97. The scores obtained under the null hypothesis (permuted data) could be used by different analyses from the same dataset – for instance, analyses from different random initializations and different training dataset folds. For an FDR threshold of  $\leq 0.05$ , the estimated empirical FDR was confidently accurate for 100 to 200 permutations (Millstein, 2013).

## Author Contributions

YL and GWC designed the project and analysis methods. HW, ZF, RT and YL developed the AI software, ZF processed data, ZF, RT and HW performed analyses using the AI framework, GAC, SK and GWC interpreted the analysis results in the context of AD. RT and YL drafted the manuscript, all authors revised the manuscript.

## Acknowledgments and Funding Sources

We would like to thank L Health for her advice on ROSMAP proteomics data quality control, and JC Wiley for his help in BD interaction interpretations. This study was supported by NIA grants R21 AG083299 and U54 AG065187. Data used in this project were obtained from the Accelerating Medicines Partnership Program for Alzheimer's Disease (AMP-AD) Consortium members below:

**Religious Orders Study/Memory and Aging Project (ROSMAP):** We are grateful to the participants in the Religious Order Study and the Memory and Aging Project. This work was supported by the US National Institutes of Health (U01 AG046152, R01 AG043617, R01 AG042210, R01 AG036042, R01 AG036836, R01 AG032990, R01 AG18023, RC2 AG036547, P50 AG016574, U01 ES017155, KL2 RR024151, K25 AG041906-01, R01 AG30146, P30 AG10161, R01 AG17917, R01 AG15819, K08 AG034290, P30 AG10161, and R01 AG11101).

**Mount Sinai Brain Bank (MSBB):** This work was supported by grants R01AG046170, RF1AG054014, RF1AG057440, and R01AG057907 from the NIH/NIA. R01AG046170 is a component of the AMP-AD Target Discovery and Preclinical Validation Project. Brain tissue collection and characterization was supported by NIH HHSN271201300031C.

**Mayo RNAseq Study:** Study data were provided by the following sources: The Mayo Clinic Alzheimer's Disease Genetic Studies, led by Dr. Nilufer Ertekin-Taner and Dr. Steven G. Younkin, Mayo Clinic, Jacksonville, FL, using samples from the Mayo Clinic Study of Aging, the Mayo Clinic Alzheimer's Disease Research Center, and the Mayo Clinic Brain Bank. Data collection was supported through funding by NIA grants P50 AG016574, R01 AG032990, U01 AG046139, R01 AG018023, U01 AG006576, U01 AG006786, R01 AG025711, R01 AG017216, R01 AG003949, NINDS grant R01 NS080820, CurePSP Foundation, and support from Mayo Foundation. Study data include samples collected through the Sun Health Research Institute Brain and Body Donation Program of Sun City, Arizona. The Brain and Body Donation Program is supported by the National Institute of Neurological Disorders and Stroke (U24 NS072026 National Brain and Tissue Resource for Parkinson's Disease and Related Disorders), the NIA (P30 AG19610 Arizona Alzheimer's Disease Core Center), the Arizona Department of Health Services (contract 211002, Arizona Alzheimer's Research Center), the Arizona Biomedical Research Commission (contracts 4001, 0011, 05- 901, and 1001 to the Arizona Parkinson's Disease Consortium), and the Michael J. Fox Foundation for Parkinson's Research.

## Data Availability

The ROSMAP transcriptomics and proteomics data were downloaded from the AD Knowledge Portal (<https://adknowledgeportal.synapse.org/>) (Greenwood, 2020). The ROSMAP study data can be found on the portal ([syn3219045](#)). Bulk brain RNASeq data can be found at [syn3388564](#) and TMT proteomics data can be found at [syn28723049](#). Alzheimer's disease biological domain networks can be found at [syn51739831](#). MSBB gene expression data can be found at [syn7391833](#). The full MSBB study can be found at [syn3159438](#).

## Code Availability

Code will be made available upon publication.

## References

- Argelaguet, R. V. (2018). Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets. . *Molecular systems biology*, 14(6), e8124.
- Ashburner, M. B. (2000). Gene ontology: tool for the unification of biology. *Nature genetics*, 25(1), 25-29.
- Ba, J. L. (2016). Layer Normalization. *arXiv preprint* .
- Baloni, P. A.-D. (2022). Multi-Omic analyses characterize the ceramide/sphingomyelin pathway as a therapeutic target in Alzheimer's disease. . *Communications Biology*, 5(1), 1074.
- Batra, R. K. (2023). Comparative brain metabolomics reveals shared and distinct metabolic alterations in Alzheimer's disease and progressive supranuclear palsy. . *medRxiv*.
- Beckmann, N. D. (2020). Multiscale causal networks identify VGF as a key regulator of Alzheimer's disease. *Nature communications*, 11(1), 3942.
- Bellenguez, C. K.-G. (2022). New insights into the genetic etiology of Alzheimer's disease and related dementias. *Nature genetics*, 54(4), 412-436.
- Braak, H. A. (2006). Staging of Alzheimer disease-associated neurofibrillary pathology using paraffin sections and immunocytochemistry. *Acta neuropathologica*, 112(4), 389-404.
- Bunne, C. S. (2023). Learning single-cell perturbation responses using neural optimal transport. *Nature Methods*, 20(11), 1759-1768.
- Cary, G. A. (2024). Genetic and multi-omic risk assessment of Alzheimer's disease implicates core associated biological domains. . *Alzheimer's & Dementia: Translational Research & Clinical Interventions*, 10(2), e12461.
- Cerami, E. G. (2010). Pathway Commons, a web resource for biological pathway data. *Nucleic acids research*.
- Chen, C. A. (2021). Legumain promotes tubular ferroptosis by facilitating chaperone-mediated autophagy of GPX4 in AKI. *Cell death & disease*, 12(1), 65.
- Chen, Y. R. (2018). Profiling DNA methylation using bisulfite sequencing (BS-Seq). *Plant chromatin dynamics: methods and protocols*, 31-43.
- Dammer, E. B. (2023). Batch correction and harmonization of—Omics datasets with a tunable median polish of ratio. *Frontiers in systems biology*, 3, 1092341.
- Demetci, P. S. (2022). SCOT: single-cell multi-omics alignment with optimal transport. *Journal of computational biology*, 29(1), 3-18.
- Drummond, E. K.-A. (2022). The amyloid plaque proteome in early onset Alzheimer's disease and Down syndrome. *Acta neuropathologica communications*, 10(1), 53.
- Ehehalt, R. K. (2003). Amyloidogenic processing of the Alzheimer  $\beta$ -amyloid precursor protein depends on lipid rafts. *The Journal of cell biology*, 160(1), 113-123.

- Folstein, M. F. (1975). "Mini-mental state": a practical method for grading the cognitive state of patients for the clinician. *Journal of psychiatric research*, 12(3), 189-198.
- Gao, C. L. (2021). Iterative single-cell multi-omic integration using online learning. *Nature biotechnology*, 39(8), 1000-1007.
- Greenwood, A. K. (2020). The AD knowledge portal: a repository for multi-omic data on Alzheimer's disease and aging. . *Current protocols in human genetics*, 108(1), e105.
- Günther, O. P. (2012). A computational pipeline for the development of multi-marker bio-signature panels and ensemble classifiers. *BMC bioinformatic*, 1-18.
- Hosein Khasahmadi, A. H. (2020). Memory-Based Graph Networks. . *arXiv e-prints*.
- Hrdlickova, R. T. (2017). RNA-Seq methods for transcriptome analysis. *Wiley Interdisciplinary Reviews: RNA*.
- Janizek, J. D. (2021). Explaining explanations: Axiomatic feature interactions for deep networks. *Journal of Machine Learning Research*, 22(104), 1-54.
- Johnson, E. C. (2022). Large-scale deep multi-layer analysis of Alzheimer's disease brain reveals strong proteomic disease-related changes not observed at the RNA level. . *Nature neuroscience*, 25(2), 213-225.
- Khosla, P. T. (2020). Supervised contrastive learning. *Advances in neural information processing systems*, (pp. 33, 18661-18673).
- Kingma, D. P. (2014). Adam: A method for stochastic optimization. . *arXiv preprint*.
- Kipf, T. N. (2016). Semi-supervised classification with graph convolutional networks. *arXiv preprint*.
- Kontras, K. C. (2024). Improving Multimodal Learning with Multi-Loss Gradient Modulation. *arXiv preprint* .
- Korotin, A. S. (2022). Neural optimal transport. . *arXiv preprint* .
- Lee, J. L., & Teh, Y. (2019). Set transformer: A framework for attention-based permutation-invariant neural networks. *International conference on machine learning*, (pp. 3744-3753).
- Levites, Y. D. (2023). Aβ Amyloid Scaffolds the Accumulation of Matrisome and Additional Proteins in Alzheimer's Disease. *bioRxiv*.
- Li, X. M. (2022). MoGCN: a multi-omics integration method based on graph convolutional network for cancer subtype analysis. *Frontiers in Genetics*, 13, 806842.
- Li, Y. R. (2024). 7-Dehydrocholesterol dictates ferroptosis sensitivity. *Nature*, 626(7998), 411-418.
- Liu, T. J. (2019). The deubiquitylase OTUB1 mediates ferroptosis via stabilization of SLC7A11. . *Cancer research*, 79(8), 1913-1924.
- Lv, T. J. (2024). Sevoflurane causes neurotoxicity and cognitive impairment by regulating Hippo signaling pathway-mediated ferroptosis via upregulating PRKCD. . *Experimental Neurology*, 377, 114804.
- Mariette, J. &.-V. (2018). Unsupervised multiple kernel learning for heterogeneous data integration. *Bioinformatics*, 34(6), 1009-1015.

- McKhann, G. D. (1984). Clinical diagnosis of Alzheimer's disease: Report of the NINCDS-ADRDA Work Group\* under the auspices of Department of Health and Human Services Task Force on Alzheimer's Disease. . *Neurology*, 34(7), 939-939.
- Millstein, J. &. (2013). Computationally efficient permutation-based confidence interval estimation for tail-area FDR. *Frontiers in genetics*, 4, 179.
- Park, M. W. (2021). NOX4 promotes ferroptosis of astrocytes by oxidative stress-induced lipid peroxidation via the impairment of mitochondrial metabolism in Alzheimer's diseases. *Redox biology*, 41, 101947.
- Schiebinger, G. S. (2019). Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*, 176(4), 928-943.
- Schneider, M. V. (2011). Omics technologies, data and bioinformatics principles. *Bioinformatics for Omics Data: Methods and Protocols*, 3-30.
- Shen, R. O. (2009). Integrative clustering of multiple genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioinformatics*, 25(22), 2906-2912.
- Simonyan, K. V. (2013). Deep inside convolutional networks: Visualising image classification models and saliency maps. . *arXiv preprint*.
- Speicher, N. K. (2015). Integrating different data types by regularized unsupervised multiple kernel learning with application to cancer subtype discovery. *Bioinformatics*, 31(12), i268-i275.
- Storey, J. D. (2003). Statistical significance for genomewide studies. . *Proceedings of the National Academy of Sciences*, 100(16), 9440-9445.
- Sundararajan, M. T. (2017). Axiomatic attribution for deep networks. . *International conference on machine learning*, (pp. 3319-3328).
- Tandon, R. L. (2023). Machine learning selection of most predictive brain proteins suggests role of sugar metabolism in Alzheimer's disease. *Journal of Alzheimer's Disease*, 92(2), 411-424.
- Tsatsanis, A. M. (2021). The acute phase protein lactoferrin is a key feature of Alzheimer's disease and predictor of A $\beta$  burden through induction of APP amyloidogenic processing. *Molecular psychiatry*, 26(10), 5516-5531.
- Vahabi, N. &. (2022). Unsupervised multi-omics data integration methods: a comprehensive review. . *Frontiers in genetics*, 13, 854752.
- Vaswani, A. S. (2017). Attention is all you need. *Advances in neural information processing systems*.
- Villani, C. (2009). *Optimal transport: old and new*. Berlin: springer.
- Wang, B. M. (2014). Similarity network fusion for aggregating data types on a genomic scale. *Nature methods*, 11(3), 333-337.
- Wang, F. W. (2022). Iron dyshomeostasis and ferroptosis: a new Alzheimer's disease hypothesis? *Frontiers in aging neuroscience*, 14, 830569.
- Wang, M. B. (2018). The Mount Sinai cohort of large-scale genomic, transcriptomic and proteomic data in Alzheimer's disease. . *Scientific data*, 5(1), 1-16.



- Wang, S. W. (2007). IQGAP3, a novel effector of Rac1 and Cdc42, regulates neurite outgrowth. . *Journal of cell science*, 120(4), 567-577.
- Wang, T. S. (2021). MOGONET integrates multi-omics data using graph convolutional networks allowing patient classification and biomarker identification. . *Nature communications*, 12(1), 3445.
- Wang, W. T. (2020). What makes training multi-modal classification networks hard?. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, (pp. 12695-12705).
- Wang, Y. L. (2020). NEDD4L-mediated LTF protein degradation limits ferroptosis. *Biochemical and Biophysical Research Communications*, 531(4), 581-587.
- Watson, C. M. (2023). Quantitative mass spectrometry analysis of cerebrospinal fluid protein biomarkers in Alzheimer's Disease. *Scientific Data*, 10(1), 261.
- Wolfsgruber, S. J. (2014). The CERAD neuropsychological assessment battery total score detects and predicts Alzheimer disease dementia with high diagnostic accuracy. *The American Journal of Geriatric Psychiatry*, 22(10), 1017-1028.
- Wu, C. J. (2024). Contribution of ApoB-100/SORT1-Mediated Immune Microenvironment in Regulating Oxidative Stress, Inflammation, and Ferroptosis After Spinal Cord Injury. *Molecular Neurobiology*, 1-13.
- Xiao, Z. S. (2022). Reduction of lactoferrin aggravates neuronal ferroptosis after intracerebral hemorrhagic stroke in hyperglycemic mice. *Redox biology*, 50, 102256.
- Xie, Y. P. (2005). A note on using permutation-based false discovery rate estimates to compare different analysis methods for microarray data. *Bioinformatics*, 21(23), 4280-4288.
- Yan, H. F. (2021). Ferroptosis: mechanisms and links with diseases. *Signal transduction and targeted therapy*, 6(1), 49.
- Yan, L. H. (2023). Integrative analysis of TBI data reveals Lgm1 as a key player in immune cell-mediated ferroptosis. . *BMC genomics*, 24(1), 747.
- Yang, M. M.-L. (2023). Multi-omic integration via similarity network fusion to detect molecular subtypes of ageing. *Brain Communications*, 5(2), fcad110.
- Yao, Y. K. (2021). A delta-secretase-truncated APP fragment activates CEBPB, mediating Alzheimer's disease pathologies. *Brain*, 144(6), 1833-1852.
- Ye, H. H. (2024). Involvement of CD44 and MAPK14-mediated ferroptosis in hemorrhagic shock. *Apoptosis*, 29(1), 154-168.
- Ying, Z. B. (2019). Gnnexplainer: Generating explanations for graph neural networks. *Advances in neural information processing systems*.
- Zhang, Z. S. (2014). Cleavage of tau by asparagine endopeptidase mediates the neurofibrillary pathology in Alzheimer's disease. *Nature medicine*, 20(11), 1254-1262.