

Developing and Applying Machine Learning and Artificial Intelligence Approaches to Leverage
Multi-Omics to Advance Understanding of Alzheimer's Disease Pathology

By
Jerome J. Choi

A dissertation submitted in partial fulfillment of
the requirements for the degree of

Doctor of Philosophy
(Epidemiology)

at the
University of Wisconsin-Madison
2025

Date of final oral examination: 05/12/2025

This dissertation is approved by the following members of the Final Oral Committee:
Corinne D. Engleman, Professor, Population Health Sciences
Daifeng Wang, Associate Professor, Biostatistics and Medical Informatics
John Svaren, Professor, Comparative Biosciences
Shaneda Warren Anderson, Assistant Professor, Population Health Sciences
Megan Zuelsdorff, Assistant Professor, School of Nursing
Tianyuan Lu, Assistant Professor, Population Health Sciences

Table of Contents

Abstract.....	iv
Dedication	v
Acknowledgements	vi
Introduction.....	8
Understanding Complex Biological Systems	8
Alzheimer's Disease.....	9
Multi-Omics Approaches in Biological Research.....	10
Machine Learning for Multi-Omics Data Integration.....	12
Gene Regulation and Transcription Factors in Oligodendrocytes.....	14
Challenges in Integrating Multi-Omics Data for Disease Prediction	16
Application of Multi-Omics and Machine Learning in Alzheimer’s Disease	17
Bridging Biological Insights with Computational Approaches.....	18
Specific Aims	20
Significance.....	23
Innovation.....	25
References.....	26
CHAPTER 2: COSIME: Cooperative multi-view integration and Scalable and Interpretable Model Explainer	37
Abstract.....	37
Introduction	38
Results.....	40
Overview of COSIME.....	40
Simulation study.....	42
Classifying cognitive diagnosis for Alzheimer’s disease from transcriptomic (astrocytes) and metabolomic data	45

Classifying cognitive status (dementia) from transcriptomic (oligodendrocyte) and epigenomic (oligodendrocyte) data.....	49
Predicting Alzheimer’s disease progression scores from transcriptomic (microglia)and spatial transcriptomic MERFISH (astrocyte) data	53
Discussion.....	57
Methods	59
COSIME overview	59
Simulated data	59
Real data	65
Model design	68
Feature Importance and Interaction.....	79
Downstream analysesDown	83
Benchmarking methods	87
References	88
Supplementary information	101
CHAPTER 3: Multi-Omics Integration of Transcriptomics and Metabolomics with Machine Learning Uncovers Novel Risk Factors for Alzheimer's Disease	105
Abstract	105
BACKGROUND.....	107
METHODS.....	109
Study design	109
Study participants and cohort description	110
Alzheimer’s disease outcomes	111
Plasma metabolites collections and quality control.....	112
Genotyping, quality control, and imputation.....	113
Gene expression imputation	115
Statistical analyses.....	115
RESULTS.....	118
Descriptive statistics for the participants.....	118
Untargeted imputed transcriptomics and metabolomics	119
AD phenotype prediction	120
Biomarker importance and interactions.....	121
DISCUSSION.....	124
REFERENCES	128

Conclusion	134
Appendix (Publication) – CHAPTER 1: CoTF-reg reveals cooperative transcription factors in oligodendrocyte gene regulation using single-cell multi-omics	138
Abstract.....	138
Introduction	138
Results.....	139
Deep learning and single-cell multi-omics for identifying cooperative transcription factors in oligodendrocytes	
identification of the co-binding transcription factors in oligodendrocyte-specific regulatory regions	139
Identification of the co-binding transcription factors in oligodendrocyte-specific regulatory regions.....	139
Oligodendrocytes gene expression relationships between transcription factors and target genes	141
Deep learning and Shapley interaction scores to measure cooperativity of co-binding transcription factors .	142
Oligodendrocyte gene regulatory network analysis for cooperative TF pairs and transcription factor hierarchy	143
Independent validation for cooperative TFs.....	143
Independent validation for the prediction performance of the models.....	144
Methods	145
Validation of cooperative TF pairs.....	146
ChIP-seq enrichment analysis	146
Boolean cooperativity of TF pairs.....	146
Validation of the prediction performance evaluation study	146
Single-cell ATAC-seq data.....	146
Single-cell RNA-seq data	146
Differential expression testing.....	147
Position Frequency Matrices	147
Co-enrichment analysis	147
Key transcription factors	147
Deep learning models	147
Shapley interaction scores	147
Coefficient of variance	148
Hierarchy analysis	148
Statistics and reproducibility	148
References	148
Supplementary information	151

Abstract

Alzheimer's disease (AD) is a progressive neurodegenerative disorder driven by complex, multifactorial processes involving genetic, epigenetic, transcriptomic, and metabolic dysregulation. Understanding the molecular mechanisms underlying AD requires approaches that can capture the dynamic interactions across these biological layers. The disease is not the result of a single molecular event but emerges from disruptions in coordinated biological networks, making it difficult to identify clear causal pathways or biomarkers. Multi-omics technologies provide a means to interrogate these layers simultaneously, offering a comprehensive view of disease biology by linking genetic variation to downstream molecular and cellular phenotypes. However, the integration and interpretation of multi-omics data pose significant analytical challenges due to high dimensionality, heterogeneity, sparsity, and biological noise. Machine learning (ML) offers powerful tools to address these challenges by modeling non-linear relationships, uncovering hidden patterns, learning from high-dimensional data, and enabling accurate and interpretable phenotype prediction. This dissertation develops and applies novel ML-based frameworks for integrative multi-omics analysis in the context of AD. First, the roles of cooperative transcription factor (TF) pairs in regulating target genes within the gene regulatory networks (GRNs) of oligodendrocytes were characterized, utilizing multi-omics data and deep learning approaches. Second, COSIME: Cooperative Multi-view Integration and Scalable and Interpretable Model Explainer was developed with applications to Alzheimer's Disease (AD). Third, multi-omics integration of transcriptomics and metabolomics with machine learning was used to uncover novel risk factors for AD. By combining computational innovation with biological insight, this work advances systems-level understanding of AD and contributes new tools for precision medicine and biomarker discovery.

Dedication

As I reflect on my Ph.D. journey, I am deeply aware that none of this would have been possible without the love, support, and encouragement of those closest to me. The challenges and triumphs along the way have been shared with many who have been my unwavering source of strength. It is with immense gratitude that I dedicate this achievement to them.

I would like to express my deepest gratitude to my family. To my parents, thank you for your endless love, encouragement, and belief in me. Your support has been the foundation of everything I've achieved. To my wife, thank you for your unwavering love, patience, and strength. You have been my constant source of support, balance, and inspiration throughout this journey. To my brother and his family, thank you for your constant support and for being there throughout this journey. I would also like to express my heartfelt gratitude to my in-laws. Your warmth, kindness, and support have been a great source of comfort and strength, and I am so thankful for the family you've embraced me into. And to our beloved dogs, thank you for your comforting presence and loyal companionship during the many long days and late nights of this work.

Finally, I want to thank all those who have quietly supported me along the way. Your belief in me, your encouragement, and your constant presence have made this journey possible. I am deeply grateful for the role each of you has played in helping me reach this point.

Acknowledgements

I am deeply grateful to my co-advisors, Drs. Corinne Engelman and Daifeng Wang, for your invaluable guidance, mentorship, and support throughout my graduate studies. Your insight, high standards, and trust in my work have profoundly shaped my development as a researcher. I am also grateful to my additional supervisors, Drs. John Svaren and Tianyuan Lu, for their support, advice, and encouragement throughout this process. Their perspectives and expertise have been instrumental in broadening my scientific thinking. I also thank my committee members, Drs. Shaneda Warren Anderson, and Megan Zuelsdorff, for their thoughtful feedback and guidance, which have greatly enriched the direction and quality of this dissertation.

I am also incredibly thankful to my mentors, peers, and friends at the Engelman Lab, Wang Lab, and the Waisman Center, whose expertise, thoughtful discussions, and shared enthusiasm made this work not only possible but deeply rewarding. I have learned so much from working alongside such talented, curious, and generous individuals. Your willingness to share ideas, provide feedback, and collaborate on challenging problems has been invaluable to my growth as a scientist. I especially appreciate the supportive and collegial environment you created—one that encouraged exploration, resilience, and continuous learning. Your encouragement, insights, and friendship have meant a great deal to me, both professionally and personally, and I carry those experiences with deep gratitude.

Additionally, I would like to express my sincere gratitude to the funding agencies that financially supported this work, including R21 NS128761, RF1 MH128695, R01 AG067025, the National Science Foundation CAREER Award (2144475), and a core grant to the Waisman Center from NICHD (P50 HD105353). I am also deeply grateful to the participants, investigators, and

teams from the WRAP and Wisconsin ADRC, whose contributions were invaluable to this research. This work was further supported by grants from the National Institutes of Health, including R01 AG27161, P30 AG062715, and RF1 AG054047. Finally, I would like to acknowledge the Waisman Center for their fellowship support through the Wisconsin Distinguished Graduate Fellowships, which provided critical funding and resources that greatly contributed to the advancement of my research.

At last, I am deeply grateful for the collective contributions of all those who have supported and inspired me throughout this process. Your guidance, encouragement, and collaboration have been essential in shaping this dissertation and my academic journey. I will always be thankful for the role each of you has played in helping me reach this milestone, and I carry your impact with me.

Introduction

Understanding Complex Biological Systems

Biological systems are inherently complex, involving intricate networks of molecular interactions that govern cellular functions and organismal behavior¹. These systems operate across multiple scales, from individual molecules to entire organs, with a web of interconnected biological processes that ensure proper function. At the core of biological complexity lies the interplay between genes, proteins, metabolites, and other cellular components, each playing a vital role in regulating cellular states and responses to environmental signals².

To fully understand how biological systems function, we must consider not only individual components but also their dynamic interactions. This complexity is evident in numerous physiological processes, including development, immune responses, and disease mechanisms³⁻⁵. For example, the expression of genes in a specific cell type is not driven by isolated factors but is instead shaped by a coordinated network of regulatory mechanisms, such as transcription factors, epigenetic modifications, and chromatin remodeling⁶. Moreover, cellular responses are often governed by the integration of external stimuli and the internal state of cell, creating a continuously adaptive environment⁷.

As the understanding of these systems deepens, it becomes clear that studying biological processes in isolation does not provide a comprehensive understanding of their roles in health and disease. Traditional, reductionist approaches have focused on individual molecules or pathways, but the complexity of biological systems demands a more holistic perspective¹. This has led to the rise of multi-omics approaches, which enable the simultaneous study of several layers of biological information from genomics, transcriptomics, metabolomics, and epigenomics⁸⁻¹². By integrating

data from these different sources, multi-omics can provide a more complete picture of the biological processes at play and how they contribute to various states, including diseases.

In the context of complex diseases, such as Alzheimer's disease, understanding these systems becomes even more critical. Diseases often involve disruptions at multiple levels, from genetic mutations to changes in gene expression, protein function, and cellular metabolism^{13,14}. To understand these multifaceted processes and their contribution to disease progression, we must adopt a systems biology approach — one that integrates information across these diverse molecular layers.

Alzheimer's Disease

Alzheimer's disease (AD) is a prime example of a complex, multifactorial disorder, characterized by the progressive degeneration of cognitive functions, memory loss, and behavioral changes^{10,14-16}. Its pathophysiology is influenced by genetic, epigenetic, and metabolic factors, which collectively contribute to the onset and progression of the disease^{17,18}. At the molecular level, AD is marked by the accumulation of abnormal protein aggregates, such as amyloid-beta plaques and tau tangles, which disrupt normal cellular function¹⁹⁻²⁵. However, these hallmark features alone do not fully explain the disease's complexity, and additional molecular factors, including changes in gene expression and metabolism, play a critical role in disease progression.

The genetic component of AD is crucial, with several risk genes, such as *APOE*, contributing to susceptibility²⁶⁻²⁸. Yet, genetic risk alone does not account for the entire disease process, suggesting the involvement of intricate molecular networks. Beyond genetic mutations, disruptions in cellular processes, such as synaptic plasticity, neurotransmission, mitochondrial

function, and inflammation also play significant roles in disease progression²⁹⁻³⁵. These changes occur across multiple levels, from molecular signaling to tissue and organ dysfunction, making AD a disease that transcends simple genetic or biochemical explanation.

Furthermore, the role of the brain's immune system in AD, particularly the activation of microglia and astrocytes, has become increasingly evident^{34,36-41}. Chronic neuroinflammation, a key feature of AD, not only exacerbates neuronal damage but also alters the trajectory of disease progression^{35,42}. Additionally, alterations in metabolism, including impaired glucose metabolism in the brain, have been linked to cognitive decline, further highlighting the multifaceted nature of AD^{43,44}.

To better understand AD and develop effective treatments, it is essential to adopt a holistic approach that focuses on the complex interplay between genomics, transcriptomics, metabolomics, and epigenomics. By integrating these molecular layers, we can uncover new insights into the mechanisms underlying the disease and identify potential therapeutic targets. This complexity highlights the need for advanced methodologies capable of analyzing vast amounts of omics data to provide a more comprehensive understanding of AD.

Multi-Omics Approaches in Biological Research

The study of biological systems has historically relied on single-omics approaches, such as genomics or transcriptomics, to examine isolated aspects of cellular function^{45,46}. However, these approaches only offer a limited understanding of complex biological processes, as they capture only one layer of information in a multi-faceted network of biological interactions. To gain a more comprehensive understanding, multi-omics approaches have emerged as powerful tools to integrate and analyze data across several biological domains simultaneously. By incorporating

diverse types of omics data, such as genomics, transcriptomics, metabolomics, and epigenomics, we can explore the interconnectedness of genetic, molecular, and biochemical features within cells and tissues^{11,47-49}.

Each omic layer provides valuable insight into different biological processes: genomics reveals the genetic blueprint⁵⁰, transcriptomics identifies gene expression patterns⁵¹, epigenomics explores chromatin accessibility and regulatory element activity⁵², and metabolomics measures the small molecules involved in metabolic pathways⁵³. When integrated, these datasets offer a more complete and dynamic picture of cellular states, biological functions, and disease mechanisms⁸⁻¹².

The power of multi-omics lies in its ability to link molecular features across different levels, providing insight into how genetic variants lead to altered gene expression, how this affects protein function, and how changes in protein activity ultimately influence cellular processes and disease outcomes⁸⁻¹². By integrating multi-omics data, we can identify the connections between these molecular factors, gaining a deeper understanding of the mechanisms driving disease pathology and potential therapeutic targets.

However, integrating these diverse omics data presents its own set of challenges. The data from different omics layers can vary in their type, format, and dimensionality, making it difficult to combine them into a cohesive analysis. Additionally, the relationships between these different biological features are often complex and non-linear, requiring advanced computational tools to identify meaningful patterns and interactions. As a result, the need for novel computational frameworks and analytical tools that can handle multi-omics data has become increasingly evident.

Machine Learning for Multi-Omics Data Integration

Integrating multi-omics data presents significant challenges due to the varying types, formats, and dimensionalities of data across different omics layers^{8,9}. These differences make it difficult to combine them into a unified analysis. Additionally, the complex, non-linear relationships between biological features necessitate advanced computational tools to identify meaningful patterns and interactions. Machine learning (ML) techniques, including deep learning, offer powerful solutions for overcoming these obstacles^{54,55}. By leveraging these tools, we can integrate genomics, transcriptomics, metabolomics, and epigenomics data, providing a unified framework for analyzing complex biological systems. This integration enables the subsequent exploration of feature interactions, helping to uncover deeper insights into disease mechanisms and biological processes.

Through ML, we can develop predictive models that reveal hidden patterns and relationships in multi-omics data, identifying how various molecular features (such as genes, proteins, and metabolites) contribute to biological outcomes. By training models on large, multi-dimensional datasets, ML algorithms can identify complex associations that might be difficult to detect using traditional statistical methods⁵⁶. For instance, deep neural networks can learn non-linear relationships and intricate interactions between genomic, transcriptomic, metabolomic, and epigenomic features, providing a holistic understanding of cellular processes⁵⁷.

In the context of epidemiological and bioinformatics research, ML models are particularly valuable for predicting disease outcomes and identifying biomarkers^{58,59}. For example, in Alzheimer's disease, ML approaches can integrate genomic, transcriptomics metabolomic, and epigenomic data to predict disease progression or identify potential therapeutic targets^{60,61}. By

integrating multi-omics data, these models can identify complex biological signatures that are associated with disease phenotypes, offering more accurate predictions than single-omics approaches.

Additionally, ML can help identify important feature interactions within multi-omics data.^{62,63} These interactions, whether synergistic or antagonistic, are often crucial in understanding the underlying mechanisms of diseases. For example, certain gene mutations might have more significant effects when coupled with specific changes in protein expression or metabolite levels. Moreover, ML can assess the importance of individual features, highlighting those that have the greatest impact on disease outcomes^{64,65}. Through techniques like feature importance ranking and interaction modeling, ML can uncover these critical relationships, providing a deeper understanding of disease mechanisms.

Despite the potential of integrating multi-omics data with machine learning, several challenges remain. The first challenge is the integration of multi-omics data itself. Different omics layers vary significantly in their structure, scale, and dimensionality. These differences make it difficult to combine the data into a cohesive model that can be easily interpreted by machine learning algorithms. Existing methods, like Cooperative Learning⁶⁶, DIABLO: an integrative approach for identifying key molecular drivers from multi-omics assays⁶⁷, and MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data⁶⁸, have been used to address this issue by representing each omic layer in a unified form. However, these techniques often struggle to fully capture complex, non-linear relationships across different data types, limiting their ability to model the intricate interactions that drive disease mechanisms. The second challenge is model interpretability. While machine learning models can make accurate

predictions, understanding the biological relationships driving those predictions is essential, particularly when working with complex, high-dimensional multi-omics data. This understanding is crucial for ensuring that the model's outcomes are not only accurate but also biologically meaningful. Existing interpretability methods, such as SHAP (Shapley Additive Explanations)⁶⁹ and LIME (Local Interpretable Model-Agnostic Explanations)⁷⁰, provide ways to explain how individual features contribute to model predictions. These techniques quantify the contribution of each feature, helping to uncover how specific genetic, epigenetic, or metabolic factors influence disease outcomes. However, these methods have limitations— SHAP is computationally expensive when dealing with large datasets and cannot be applied to complex machine learning models. Meanwhile, LIME lacks the capability to compute feature interactions, which are essential for understanding the relationships between variables in multi-omics data. Despite these challenges, both methods remain crucial in enhancing the biological relevance of the findings by providing insights into the underlying mechanisms driving predictions.

The application of ML to multi-omics data integration offers substantial promise for advancing our understanding of complex biological systems and predicting disease outcomes. However, the existing methods for integrating and interpreting multi-omics data face significant limitations, particularly in handling the heterogeneity of data and ensuring model interpretability. To fully harness the power of ML in biological research and clinical applications, it is essential to develop new methods that address these challenges, allowing for more accurate and biologically meaningful insights.

Gene Regulation and Transcription Factors in Oligodendrocytes

Oligodendrocytes are specialized glial cells within the central nervous system (CNS) responsible for the formation and maintenance of the myelin sheath, which insulates neuronal axons and

facilitates rapid electrical signaling⁷¹⁻⁷³. This myelination process is essential for normal neural function, and defects in oligodendrocyte development or function can lead to a variety of neurological disorders, including AD⁷²⁻⁷⁴. Understanding the intricate gene regulatory networks that govern oligodendrocyte differentiation, maturation, and function is crucial for developing targeted therapies for these conditions.

Transcription factors (TFs) are proteins that bind to specific DNA sequences to regulate the expression of target genes (TGs)⁷⁵. In oligodendrocytes, a precise balance of TF activity is required for proper gene expression during the differentiation of oligodendrocyte precursor cells (OPCs) into mature oligodendrocytes^{76,77}. Several TFs have been identified as critical players in this process, including SOX10 and OLIG2, among others⁷⁸⁻⁸⁴. These transcription factors not only direct the expression of genes involved in oligodendrocyte differentiation but also can coordinate the expression of genes essential for myelination and maintaining oligodendrocyte function.

However, the role of transcription factors in gene regulation is not simply a matter of individual proteins acting in isolation. Increasing evidence suggests that TFs do not act alone but cooperate in complex TF-TF interactions within the regulatory regions of TGs^{85,86}. This cooperation can involve direct protein-protein interactions or indirect effects through the modulation of chromatin structure and accessibility. For instance, SOX10 and OLIG2 are known to work together to regulate a set of genes critical for oligodendrocyte differentiation and myelination^{78,81,87}. These cooperative interactions are critical for ensuring the temporal and spatial regulation of gene expression that supports the development and function of oligodendrocytes.

Understanding how TFs cooperate in oligodendrocytes to regulate gene expression is of great interest because disruptions in these regulatory networks can have profound effects on

oligodendrocyte function and CNS health^{88,89}. Dysregulation of TF interactions can lead to defective myelination, contributing to various neurodegenerative diseases, including AD⁹⁰. Moreover, recent studies suggest that the identification of cooperative TF pairs could provide valuable insights into the regulatory mechanisms underlying these diseases and might even offer new therapeutic targets^{32,91,92}.

To fully unravel these complex TF interactions, it is necessary to adopt advanced computational approaches that can integrate different types of genomic data, such as Assay for Transposase-Accessible Chromatin using sequencing (ATAC-seq) and RNA sequencing (RNA-seq), to identify TF pairs that cooperate to regulate oligodendrocyte-specific genes^{93,94}. These methods hold the potential to uncover new mechanisms of oligodendrocyte function and pathology, providing a deeper understanding of how transcriptional regulation contributes to CNS health and disease.

Challenges in Integrating Multi-Omics Data for Disease Prediction

Integrating multi-omics data is crucial for enhancing disease phenotype prediction, as it provides a comprehensive view of the molecular underpinnings of diseases. However, this integration is complex due to the variability in data types, formats, and scales. Genomic data provide static DNA sequences, transcriptomic data capture dynamic gene expression changes, and metabolomic data offers insights into protein abundance and metabolic activities. Standardizing these diverse data types is essential for effective integration.

Moreover, multi-omics datasets are often incomplete and noisy, due to technical limitations and variability in data collection methods⁹⁵. Addressing missing values and minimizing noise through advanced imputation techniques and robust preprocessing methods is critical for accurate

analysis. The high dimensionality of these datasets, where the number of features often exceeds the number of samples, also increases the risk of overfitting⁹⁶. Employing dimensionality reduction techniques and ensuring access to high-performance computing resources are essential for managing these challenges, as they help reduce computational complexity and enable the extraction of key features, such as using embeddings to represent high-dimensional data in a lower-dimensional space for more efficient and accurate analysis^{97,98}.

Understanding the complex relationships between the features in omics data is difficult due to the non-linear and multifaceted nature of biological systems⁹. To address this, sophisticated computational models are needed to capture these complexities and uncover meaningful patterns. Additionally, while ML models are effective in identifying patterns within complex datasets, their "black box" nature often limits interpretability⁹⁹. Improving model transparency is particularly important in clinical settings, where understanding the rationale behind predictions can guide decision-making.

To tackle these challenges, novel computational frameworks and analytical tools are needed to effectively integrate diverse omics data. Overcoming these obstacles will significantly enhance disease phenotype predictions and help identify potential therapeutic targets.

Application of Multi-Omics and Machine Learning in Alzheimer's Disease

Despite significant research efforts, the etiology and pathogenesis of AD remain poorly understood. A deeper understanding of the molecular factors involved in the disease is essential to identify potential biomarkers for early diagnosis. The integration of multi-omics data and machine learning approaches offers promising strategies to reveal these underlying mechanisms and provide new insights into AD progression¹⁰⁻¹².

In the context of AD, the integration of multi-omics data with clinical information—such as cognitive assessments and clinical outcomes—further enhances the accuracy of predictive models. This holistic approach enables more personalized disease predictions and better-targeted therapeutic strategies^{100,101}. By leveraging multi-omics data, machine learning algorithms can identify key molecular signatures, provide early diagnostic biomarkers, and predict disease phenotypes such as cognitive decline. This integrated strategy holds significant potential to advance our understanding of AD and improve patient care.

Bridging Biological Insights with Computational Approaches

While multi-omics technologies and ML have provided valuable insights into the molecular mechanisms of complex diseases, integrating these tools to deepen our understanding of disease biology remains a significant challenge. In this dissertation, I develop novel computational methods to enhance the integration of multi-omics data, specifically focusing on AD phenotype prediction. The second aim (Aim 2) introduces Cooperative Multi-view Integration and Scalable Interpretable Model Explainer (COSIME), a machine learning method that integrates multi-view omics data to improve disease phenotype prediction and assess feature importance and interactions, both within-view and across-view. The third aim (Aim 3) applies COSIME to identify key biomarkers and interactions between genomics and metabolomics, providing insights into AD risk factors and enhancing the understanding of disease progression. Collectively, these efforts advance bioinformatics and contribute to a systems-level understanding of the biological mechanisms underlying AD.

The integration of multi-omics data remains a significant challenge in understanding the molecular complexity of diseases like AD. This dissertation tackles this issue by developing ML

models that process diverse omics data types, including genomics, transcriptomics, and epigenomics. These models are designed to capture complex relationships between molecular layers, improving the accuracy of AD phenotype predictions and providing deeper insights into the disease's molecular mechanisms.

A key focus of this research is the development of algorithms that can effectively address challenges in multi-omics data, such as noise and high-dimensionality. The proposed computational frameworks integrate heterogeneous omics data while ensuring meaningful and interpretable biological insights. These ML methods not only enhance prediction accuracy but also help us better understand the regulatory dynamics of AD.

This dissertation also investigates the regulatory mechanisms in AD, specifically the role of transcription factors in gene regulation, through the integration of single-cell gene regulatory networks. By leveraging ML models to analyze chromatin accessibility and gene expression, the work delves into how these regulatory dynamics contribute to disease progression and offer new perspectives on genetic-environmental interactions in AD risk.

Ultimately, this research advances our understanding of AD by providing novel bioinformatics tools that integrate multi-omics data, enabling more accurate analysis of the disease's underlying biological mechanisms and identifying potential therapeutic targets.

Specific Aims

Studies of brain diseases, such as AD have implicated white matter alterations in their pathogenesis, and models of cell types, for instance, oligodendrocyte defects, have validated mechanistic connections with cognitive disorders¹⁰². Recent technological advances have enabled gene regulation detection through multi-omics (i.e., genomics, transcriptomics, proteomics). Particularly, emerging next-generation sequencing (NGS) technologies, such as single-cell RNA sequencing (scRNA-seq), allow us to study functional genomics and gene regulation at the cell-type level. Large collaborative projects, such as PsychENCODE¹⁰³ and Psych-AD¹⁰⁴, generate single-cell multi-omics data aiming to understand the molecular mechanisms of brain diseases, including AD. Metabolomics technology has emerged as a tool for studying small molecules influenced by factors, such as genetics and disease processes. Some plasma and cerebrospinal fluid (CSF) metabolites were identified as associated with AD^{105,106} and the relationships between genomics, metabolomics, and AD phenotypes need to be comprehensively understood.

Our long-term goal is to integrate and analyze large-scale multi-omics data at the population level to shed light on complex relationships involved in AD pathology. The overall objective of this research is to integrate multi-omics data, such as scRNA-seq and single-cell Assay for Transposase-Accessible Chromatin sequencing (scATAC-seq), to identify the activity of the co-binding TFs in gene regulatory networks (GRNs). We then aim to develop a new flexible machine learning model that optimally integrates multi-omics data to accurately select AD risk factors and predict AD phenotypes. Lastly, we aim to apply the new method to AD epidemiological multi-omics data to identify novel AD risk factors. Our central hypothesis is that functional genetic variants, metabolites, and their interactions will influence changes in pre-clinical AD biomarkers and cognitive function. The rationale for our proposed research is that integration of genomic and

metabolomic data will enable the identification and statistical modeling of the complex interplay of genes and metabolites involved in AD pathology, which is necessary to achieve the goal of precision medicine for AD. We will test our central hypothesis by executing the following aims:

Aim 1: Characterize the roles of cooperative transcription factor (TF) pairs in regulating target genes within the gene regulatory networks (GRNs) of oligodendrocytes, utilizing multi-omics data and deep learning approaches.

H1a: scATAC-seq data will be used to identify co-binding TF pairs in oligodendrocyte-specific regulatory regions, and deep learning models trained on scRNA-seq data will predict the expression of target genes regulated by these cooperative TF pairs.

H1b: eQTLs will independently validate the regulatory roles of cooperative TF pairs, confirming their involvement in the regulation of genes associated with oligodendrocyte function and neurodegenerative diseases, such as AD.

Aim 2: Develop COSIME: Cooperative Multi-view Integration and Scalable Interpretable Model Explainer with applications to AD.

H2a: Developing a predictive model using an unsupervised neural network combined with learnable optimal transport methods to integrate multi-omics data will enhance the accuracy of AD phenotype prediction.

H2b: COSIME will provide interpretable outputs by identifying feature importance and elucidating feature interactions within and across different omics modalities when applied to AD phenotypes.

Aim 3: Identify genomic and metabolomic AD risk factors and their interaction effects on AD phenotypes by applying COSIME to data from the Wisconsin Registry for Alzheimer's Prevention (WRAP) and Wisconsin Alzheimer's Disease Research Center (ADRC) cohorts.

H3: Genes and metabolites will be identified and replicated as potential AD risk factors by integrating multi-omics data using COSIME.

Through the proposed study, we expect to have an important positive impact because the identification of novel AD biomarkers, such as co-binding transcription factors, metabolic profiles, and interactions between those, using advanced bioinformatics tools will help understand AD pathology. This research will help prevent, diagnose, and treat dementia, and lays the groundwork for precision medicine.

Significance

scRNA-seq technology is a recently emerging tool to study functional genomics and gene regulation. Metabolomics technology is another tool that has emerged to study metabolites, small molecules influenced by factors, such as genetics and disease processes. Multi-omics offers a holistic view of human health and disease by providing an integrative perspective across multiple levels of biology (e.g., predicting GRNs and disease phenotype predictions). While several studies have identified transcription factors^{107,108} or metabolites^{109,110} associated with AD status, only a few focused on preclinical AD phenotypes^{111,112}, which may give insight into the pathophysiology of AD, point to therapeutic targets and discover genes and metabolites that can be utilized as early AD biomarkers.

The Wang lab has developed a computational pipeline of integrative multi-omics analyses for predicting cell-type specific disease genes and GRNs and a machine learning analysis found that cell-type specific disease genes improved clinical phenotype predictions, including those for AD.¹¹³ The Engelman lab has integrated genetic, lifestyle, metabolomic, and AD biomarker data to better understand these associations¹⁰⁶. What remains unclear is the comprehensive understanding of functional genomics and gene regulation, metabolomics, and their interactions in AD pathology. Developing and applying flexible and robust bioinformatics tools to optimally integrate multi-omics data and predict preclinical AD phenotypes is critical to prevent, diagnose, and treat AD.

This research will be significant because the identification of novel AD biomarkers, such as transcription factors, metabolic profiles, and interactions between those using advanced bioinformatics tools will help understand AD pathology. Machine learning approaches, including

deep learning, are expected to lead to identifying new pathways for targeting therapeutic agents and novel risk prediction and diagnostic tools for early AD pathology.

Innovation

The proposed research has several innovative aspects. First, it focuses on revealing the specific roles of co-binding TFs at regulatory elements and characteristics of regulatory hub enhancers. This is very novel because even though there are some studies about single-cell functional genomic research for oligodendrocytes—one of the major cell types for the central nervous system related to AD—as far as we know, no earlier studies predicting GRNs, TFs, gene-regulatory elements, and target genes using oligodendrocytes for AD exist.

The second unique aspect is multi-omics data integration and phenotype prediction. We are planning to develop a new machine learning model that will be flexible and suitable for handling the multi-omics data we use for phenotype prediction. Moreover, using a latent vector to find the best optimal transport plan to integrate multi-omics data and predict phenotypes is a new approach as far as we are aware.

Third, there are many AD phenotypes that are related to AD continuum categories of NIA-AA's A/T/N biomarker profiles¹¹⁴ in WRAP. Using cognitive functions as a proxy for preclinical AD diagnosis and ptau217 as a proxy for amyloid deposition will give us a relatively large sample size.

Lastly, our research will discover novel pathways in AD pathology by providing novel AD risk genes and metabolites. There could be overlaps between our findings and the results in other existing studies, however, since we will consider interactions between genes and metabolites in integrating multi-omics data for machine learning phenotype prediction, we may find novel AD risk factors. Furthermore, validating our findings in WRAP in an additional cohort will strongly support our hypotheses.

References

1. Ma'ayan, A. Complex systems biology. *J. R. Soc. Interface*. **14**, 20170391 (2017).
2. Wolf, Y. I., Katsnelson, M. I. & Koonin, E. V. Physical foundations of biological complexity. *Proc. Natl. Acad. Sci. U.S.A.* **115**, (2018).
3. Adami, C., Ofria, C. & Collier, T. C. Evolution of biological complexity. *Proc. Natl. Acad. Sci. U.S.A.* **97**, 4463–4468 (2000).
4. Chaplin, D. D. Overview of the immune response. *J Allergy Clin Immunol* **125**, S3-23 (2010).
5. Naylor, S. & Chen, J. Y. Unraveling human complexity and disease with systems biology and personalized medicine. *Per Med* **7**, 275–289 (2010).
6. Colvis, C. M. *et al.* Epigenetic mechanisms and gene networks in the nervous system. *J Neurosci* **25**, 10379–10389 (2005).
7. Poljšak, B. & Milisav, I. Clinical implications of cellular stress responses. *Bosn J Basic Med Sci* **12**, 122–126 (2012).
8. Hasin, Y., Seldin, M. & Lusk, A. Multi-omics approaches to disease. *Genome Biol* **18**, 83 (2017).
9. Subramanian, I., Verma, S., Kumar, S., Jere, A. & Anamika, K. Multi-omics Data Integration, Interpretation, and Its Application. *Bioinform Biol Insights* **14**, 1177932219899051 (2020).
10. Badhwar, A. *et al.* A multiomics approach to heterogeneity in Alzheimer's disease: focused review and roadmap. *Brain* **143**, 1315–1331 (2020).
11. Nativio, R. *et al.* An integrated multi-omics approach identifies epigenetic alterations associated with Alzheimer's disease. *Nat Genet* **52**, 1024–1035 (2020).
12. Kapoor, M. *et al.* Multi-omics integration analysis identifies novel genes for alcoholism with potential overlap with neurodegenerative diseases. *Nat Commun* **12**, 5071 (2021).

13. Costa, V., Aprile, M., Esposito, R. & Ciccodicola, A. RNA-Seq and human complex diseases: recent accomplishments and future perspectives. *Eur J Hum Genet* **21**, 134–142 (2013).
14. Andrews, S. J. *et al.* The complex genetic architecture of Alzheimer’s disease: novel insights and future directions. *eBioMedicine* **90**, 104511 (2023).
15. Kumar, A., Sidhu, J., Lui, F. & Tsao, J. W. Alzheimer Disease. in *StatPearls* (StatPearls Publishing, Treasure Island (FL), 2025).
16. DeTure, M. A. & Dickson, D. W. The neuropathological diagnosis of Alzheimer’s disease. *Mol Neurodegeneration* **14**, 32 (2019).
17. Wang, H., Yang, F., Zhang, S., Xin, R. & Sun, Y. Genetic and environmental factors in Alzheimer’s and Parkinson’s diseases and promising therapeutic intervention via fecal microbiota transplantation. *npj Parkinsons Dis.* **7**, 70 (2021).
18. Killin, L. O. J., Starr, J. M., Shiue, I. J. & Russ, T. C. Environmental risk factors for dementia: a systematic review. *BMC Geriatr* **16**, 175 (2016).
19. Gerrits, E. *et al.* Distinct amyloid- β and tau-associated microglia profiles in Alzheimer’s disease. *Acta Neuropathol* **141**, 681–696 (2021).
20. Griciuc, A. *et al.* Alzheimer’s disease risk gene CD33 inhibits microglial uptake of amyloid beta. *Neuron* **78**, 631–643 (2013).
21. Hampel, H. *et al.* The Amyloid- β Pathway in Alzheimer’s Disease. *Mol Psychiatry* **26**, 5481–5503 (2021).
22. O’Brien, R. J. & Wong, P. C. Amyloid precursor protein processing and Alzheimer’s disease. *Annu Rev Neurosci* **34**, 185–204 (2011).

23. Ramanan, V. K. *et al.* GWAS of longitudinal amyloid accumulation on ¹⁸F-florbetapir PET in Alzheimer's disease implicates microglial activation gene *IL1RAP*. *Brain* **138**, 3076–3088 (2015).
24. Sadigh-Eteghad, S. *et al.* Amyloid-beta: a crucial factor in Alzheimer's disease. *Med Princ Pract* **24**, 1–10 (2015).
25. Bloom, G. S. Amyloid- β and Tau: The Trigger and Bullet in Alzheimer Disease Pathogenesis. *JAMA Neurol* **71**, 505 (2014).
26. Ferrari-Souza, J. P. *et al.* *APOE* ϵ 4 associates with microglial activation independently of A β plaques and tau tangles. *Sci. Adv.* **9**, eade1474 (2023).
27. Raulin, A.-C. *et al.* ApoE in Alzheimer's disease: pathophysiology and therapeutic strategies. *Mol Neurodegeneration* **17**, 72 (2022).
28. Yin, Z. *et al.* APOE4 impairs the microglial response in Alzheimer's disease by inducing TGF β -mediated checkpoints. *Nat Immunol* **24**, 1839–1853 (2023).
29. D'Alessandro, M. C. B., Kanaan, S., Geller, M., Praticò, D. & Daher, J. P. L. Mitochondrial dysfunction in Alzheimer's disease. *Ageing Research Reviews* **107**, 102713 (2025).
30. Cai, Q. & Tammineni, P. Mitochondrial Aspects of Synaptic Dysfunction in Alzheimer's Disease. *J Alzheimers Dis* **57**, 1087–1103 (2017).
31. Skaper, S. D., Facci, L., Zusso, M. & Giusti, P. Synaptic Plasticity, Dementia and Alzheimer Disease. *CNSNDT* **16**, 220–233 (2017).
32. Kandimalla, R. & Reddy, P. H. Therapeutics of Neurotransmitters in Alzheimer's Disease. *J Alzheimers Dis* **57**, 1049–1069 (2017).
33. Kinney, J. W. *et al.* Inflammation as a central mechanism in Alzheimer's disease. *A&D Transl Res & Clin Interv* **4**, 575–590 (2018).

34. MRC CFAS *et al.* Microglial immunophenotype in dementia with Alzheimer's pathology. *J Neuroinflammation* **13**, 135 (2016).
35. Zhang, W., Xiao, D., Mao, Q. & Xia, H. Role of neuroinflammation in neurodegeneration development. *Sig Transduct Target Ther* **8**, 267 (2023).
36. Chen, Y. & Colonna, M. Microglia in Alzheimer's disease at single-cell level. Are there common patterns in humans and mice? *J Exp Med* **218**, e20202717 (2021).
37. Hansen, D. V., Hanson, J. E. & Sheng, M. Microglia in Alzheimer's disease. *Journal of Cell Biology* **217**, 459–472 (2018).
38. McQuade, A. & Blurton-Jones, M. Microglia in Alzheimer's Disease: Exploring How Genetics and Phenotype Influence Risk. *Journal of Molecular Biology* **431**, 1805–1817 (2019).
39. Sadick, J. S. *et al.* Astrocytes and oligodendrocytes undergo subtype-specific transcriptional changes in Alzheimer's disease. *Neuron* **110**, 1788–1805.e10 (2022).
40. Verkhratsky, A., Olanbarria, M., Noristani, H. N., Yeh, C.-Y. & Rodriguez, J. J. Astrocytes in Alzheimer's Disease. *Neurotherapeutics* **7**, 399–412 (2010).
41. González-Reyes, R. E., Nava-Mesa, M. O., Vargas-Sánchez, K., Ariza-Salamanca, D. & Mora-Muñoz, L. Involvement of Astrocytes in Alzheimer's Disease from a Neuroinflammatory and Oxidative Stress Perspective. *Front Mol Neurosci* **10**, 427 (2017).
42. Cherry, J. D., Olschowka, J. A. & O'Banion, M. K. Neuroinflammation and M2 microglia: the good, the bad, and the inflamed. *J Neuroinflammation* **11**, 98 (2014).
43. Butterfield, D. A. & Halliwell, B. Oxidative stress, dysfunctional glucose metabolism and Alzheimer disease. *Nat Rev Neurosci* **20**, 148–160 (2019).
44. Ryu, W.-I. *et al.* Brain cells derived from Alzheimer's disease patients have multiple specific innate abnormalities in energy metabolism. *Mol Psychiatry* **26**, 5702–5714 (2021).

45. Malone, A. F. *et al.* Harnessing Expressed Single Nucleotide Variation and Single Cell RNA Sequencing To Define Immune Cell Chimerism in the Rejecting Kidney Transplant. *JASN* **31**, 1977–1986 (2020).
46. Barbu, M. C. *et al.* Epigenetic prediction of major depressive disorder. *Mol Psychiatry* **26**, 5112–5123 (2021).
47. Desai, N. *et al.* Temporal and spatial heterogeneity of host response to SARS-CoV-2 pulmonary infection. *Nat Commun* **11**, 6319 (2020).
48. Miles, L. A. *et al.* Single-cell mutation analysis of clonal evolution in myeloid malignancies. *Nature* **587**, 477–482 (2020).
49. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).
50. Theissinger, K. *et al.* How genomics can help biodiversity conservation. *Trends in Genetics* **39**, 545–559 (2023).
51. Dinh, H. Q. *et al.* Single-cell transcriptomics identifies gene expression networks driving differentiation and tumorigenesis in the human fallopian tube. *Cell Reports* **35**, 108978 (2021).
52. Klemm, S. L., Shipony, Z. & Greenleaf, W. J. Chromatin accessibility and the regulatory epigenome. *Nat Rev Genet* **20**, 207–220 (2019).
53. Fessenden, M. Metabolomics: Small molecules, single cells. *Nature* **540**, 153–155 (2016).
54. Reel, P. S., Reel, S., Pearson, E., Trucco, E. & Jefferson, E. Using machine learning approaches for multi-omics data analysis: A review. *Biotechnology Advances* **49**, 107739 (2021).
55. Feldner-Busztin, D. *et al.* Dealing with dimensionality: the application of machine learning to multi-omics data. *Bioinformatics* **39**, btad021 (2023).

56. Rajula, H. S. R., Verlato, G., Manchia, M., Antonucci, N. & Fanos, V. Comparison of Conventional Statistical Methods with Machine Learning in Medicine: Diagnosis, Drug Development, and Treatment. *Medicina (Kaunas)* **56**, 455 (2020).
57. Leng, D. *et al.* A benchmark study of deep learning-based multi-omics data fusion methods for cancer. *Genome Biol* **23**, 171 (2022).
58. Ng, S., Masarone, S., Watson, D. & Barnes, M. R. The benefits and pitfalls of machine learning for biomarker discovery. *Cell Tissue Res* **394**, 17–31 (2023).
59. Park, D. J. *et al.* Development of machine learning model for diagnostic disease prediction based on laboratory tests. *Sci Rep* **11**, 7567 (2021).
60. Chu, S. H. *et al.* Integration of Metabolomic and Other Omics Data in Population-Based Study Designs: An Epidemiological Perspective. *Metabolites* **9**, 117 (2019).
61. Horgusluoglu, E. *et al.* Integrative metabolomics-genomics approach reveals key metabolic pathways and regulators of Alzheimer's disease. *Alzheimers Dement* **18**, 1260–1278 (2022).
62. Song, W. *et al.* AutoInt: Automatic Feature Interaction Learning via Self-Attentive Neural Networks. in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management* 1161–1170 (ACM, Beijing China, 2019). doi:10.1145/3357384.3357925.
63. Tsang, M., Rambhatla, S. & Liu, Y. How does this interaction affect me? Interpretable attribution for feature interactions. Preprint at <https://doi.org/10.48550/ARXIV.2006.10965> (2020).
64. Rengasamy, D. *et al.* Feature importance in machine learning models: A fuzzy information fusion approach. *Neurocomputing* **511**, 163–174 (2022).

65. Saarela, M. & Jauhiainen, S. Comparison of feature importance measures as explanations for classification models. *SN Appl. Sci.* **3**, 272 (2021).
66. Ding, D. Y., Li, S., Narasimhan, B. & Tibshirani, R. Cooperative learning for multiview analysis. *Proc. Natl. Acad. Sci. U.S.A.* **119**, e2202113119 (2022).
67. Singh, A. *et al.* DIABLO: an integrative approach for identifying key molecular drivers from multi-omics assays. *Bioinformatics* **35**, 3055–3062 (2019).
68. Argelaguet, R. *et al.* MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. *Genome Biol* **21**, 111 (2020).
69. Lundberg, S. & Lee, S.-I. A Unified Approach to Interpreting Model Predictions. Preprint at <https://doi.org/10.48550/ARXIV.1705.07874> (2017).
70. Ribeiro, M. T., Singh, S. & Guestrin, C. ‘Why Should I Trust You?’: Explaining the Predictions of Any Classifier. Preprint at <https://doi.org/10.48550/ARXIV.1602.04938> (2016).
71. Elbaz, B. & Popko, B. Molecular Control of Oligodendrocyte Development. *Trends in Neurosciences* **42**, 263–277 (2019).
72. Ertle, B., Schlachetzki, J. C. M. & Winkler, J. Oligodendroglia and Myelin in Neurodegenerative Diseases: More Than Just Bystanders? *Mol Neurobiol* **53**, 3046–3062 (2016).
73. Simons, M. & Nave, K.-A. Oligodendrocytes: Myelination and Axonal Support. *Cold Spring Harb Perspect Biol* **8**, a020479 (2016).
74. Quan, L., Uyeda, A. & Muramatsu, R. Central nervous system regeneration: the roles of glial cells in the potential molecular mechanism underlying remyelination. *Inflamm Regener* **42**, 7 (2022).
75. Lambert, S. A. *et al.* The Human Transcription Factors. *Cell* **172**, 650–665 (2018).

76. Tiane, A. *et al.* From OPC to Oligodendrocyte: An Epigenetic Journey. *Cells* **8**, 1236 (2019).
77. Marton, R. M. *et al.* Differentiation and maturation of oligodendrocytes in human three-dimensional neural cultures. *Nat Neurosci* **22**, 484–491 (2019).
78. Liu, Z. *et al.* Induction of oligodendrocyte differentiation by Olig2 and Sox10: Evidence for reciprocal interactions and dosage-dependent mechanisms. *Developmental Biology* **302**, 683–693 (2007).
79. Lopez-Anido, C. *et al.* Differential Sox10 genomic occupancy in myelinating glia. *Glia* **63**, 1897–1914 (2015).
80. Pozniak, C. D. *et al.* Sox10 directs neural stem cells toward the oligodendrocyte lineage by decreasing Suppressor of Fused expression. *Proc Natl Acad Sci U S A* **107**, 21795–21800 (2010).
81. Sock, E. & Wegner, M. Using the lineage determinants Olig2 and Sox10 to explore transcriptional regulation of oligodendrocyte development. *Developmental Neurobiology* **81**, 892–901 (2021).
82. Turnescu, T. *et al.* Sox8 and Sox10 jointly maintain myelin gene expression in oligodendrocytes. *Glia* **66**, 279–294 (2018).
83. Wang, J. *et al.* Olig2 Ablation in Immature Oligodendrocytes Does Not Enhance CNS Myelination and Remyelination. *J. Neurosci.* **42**, 8542–8555 (2022).
84. Yu, Y. *et al.* Olig2 targets chromatin remodelers to enhancers to initiate oligodendrocyte differentiation. *Cell* **152**, 248–261 (2013).
85. Oksuz, O. *et al.* Transcription factors interact with RNA to regulate genes. *Mol Cell* **83**, 2449–2463.e13 (2023).

86. Göös, H. *et al.* Human transcription factor protein interaction networks. *Nat Commun* **13**, 766 (2022).
87. Küspert, M., Hammer, A., Bösl, M. R. & Wegner, M. Olig2 regulates Sox10 expression in oligodendrocyte precursors through an evolutionary conserved distal enhancer. *Nucleic Acids Res* **39**, 1280–1293 (2011).
88. He, D. *et al.* Chd7 cooperates with Sox10 and regulates the onset of CNS myelination and remyelination. *Nat Neurosci* **19**, 678–689 (2016).
89. Santiago, C. & Bashaw, G. J. Transcription factors and effectors that regulate neuronal morphology. *Development* **141**, 4667–4680 (2014).
90. Meng, G. & Mei, H. Transcriptional Dysregulation Study Reveals a Core Network Involving the Progression of Alzheimer’s Disease. *Front Aging Neurosci* **11**, 101 (2019).
91. Radaeva, M., Ton, A.-T., Hsing, M., Ban, F. & Cherkasov, A. Drugging the ‘undruggable’. Therapeutic targeting of protein–DNA interactions with the use of computer-aided drug discovery methods. *Drug Discovery Today* **26**, 2660–2679 (2021).
92. Talukdar, P. D. & Chatterji, U. Transcriptional co-activators: emerging roles in signaling pathways and potential therapeutic targets for diseases. *Sig Transduct Target Ther* **8**, 427 (2023).
93. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).
94. Hao, Y. *et al.* Dictionary learning for integrative, multimodal and scalable single-cell analysis. *Nat Biotechnol* **42**, 293–304 (2024).
95. Flores, J. E. *et al.* Missing data in multi-omics integration: Recent advances through artificial intelligence. *Front Artif Intell* **6**, 1098308 (2023).

96. Charilaou, P. & Battat, R. Machine learning models and over-fitting considerations. *World J Gastroenterol* **28**, 605–607 (2022).
97. Kabir, M. F., Chen, T. & Ludwig, S. A. A performance analysis of dimensionality reduction algorithms in machine learning models for cancer prediction. *Healthcare Analytics* **3**, 100125 (2023).
98. Jia, W., Sun, M., Lian, J. & Hou, S. Feature dimensionality reduction: a review. *Complex Intell. Syst.* **8**, 2663–2693 (2022).
99. Hassija, V. *et al.* Interpreting Black-Box Models: A Review on Explainable Artificial Intelligence. *Cogn Comput* **16**, 45–74 (2024).
100. Clark, C., Dayon, L., Masoodi, M., Bowman, G. L. & Popp, J. An integrative multi-omics approach reveals new central nervous system pathway alterations in Alzheimer’s disease. *Alz Res Therapy* **13**, 71 (2021).
101. Meng, L. *et al.* Multi-omics analysis reveals the key factors involved in the severity of the Alzheimer’s disease. *Alz Res Therapy* **16**, 213 (2024).
102. Nasrabady, S. E., Rizvi, B., Goldman, J. E. & Brickman, A. M. White matter changes in Alzheimer’s disease: a focus on myelin and oligodendrocytes. *acta neuropathol commun* **6**, 22 (2018).
103. Emani, P. S. *et al.* Single-cell genomics and regulatory networks for 388 human brains. *Science* **384**, eadi5199 (2024).
104. Psych-AD: Neuropsychiatric Symptoms in Alzheimer’s Disease. <https://adknowledgeportal.synapse.org/Explore/Programs/DetailsPage?Program=Psych-AD>.
105. Vogt, N. M. *et al.* The gut microbiota-derived metabolite trimethylamine N-oxide is elevated in Alzheimer’s disease. *Alz Res Therapy* **10**, 124 (2018).

106. Darst, B. F., Lu, Q., Johnson, S. C. & Engelman, C. D. Integrated analysis of genomics, longitudinal metabolomics, and Alzheimer's risk factors among 1,111 cohort participants. *Genetic Epidemiology* **43**, 657–674 (2019).
107. Kodam, P., Sai Swaroop, R., Pradhan, S. S., Sivaramakrishnan, V. & Vadrevu, R. Integrated multi-omics analysis of Alzheimer's disease shows molecular signatures associated with disease progression and potential therapeutic targets. *Sci Rep* **13**, 3695 (2023).
108. Kapoor, M. *et al.* Multi-omics integration analysis identifies novel genes for alcoholism with potential overlap with neurodegenerative diseases. *Nat Commun* **12**, 5071 (2021).
109. Nho, K. *et al.* Altered bile acid profile in mild cognitive impairment and Alzheimer's disease: Relationship to neuroimaging and CSF biomarkers. *Alzheimer's & Dementia* **15**, 232–244 (2019).
110. Wang, J. *et al.* Peripheral serum metabolomic profiles inform central cognitive impairment. *Sci Rep* **10**, 14059 (2020).
111. Casanova, R. *et al.* Blood metabolite markers of preclinical Alzheimer's disease in two longitudinally followed cohorts of older individuals. *Alzheimer's & Dementia* **12**, 815–822 (2016).
112. Rai, S. N. *et al.* Therapeutic Potential of Vital Transcription Factors in Alzheimer's and Parkinson's Disease With Particular Emphasis on Transcription Factor EB Mediated Autophagy. *Front. Neurosci.* **15**, 777347 (2021).
113. Jin, T. *et al.* scGRNom: a computational pipeline of integrative multi-omics analyses for predicting cell-type disease genes and regulatory networks. *Genome Med* **13**, 95 (2021).
114. Jack, C. R. *et al.* A/T/N: An unbiased descriptive classification scheme for Alzheimer disease biomarkers. *Neurology* **87**, 539–547 (2016).

CHAPTER 2: COSIME: Cooperative multi-view integration and Scalable and Interpretable Model Explainer

1 COSIME: Cooperative multi-view integration and Scalable 2 and Interpretable Model Explainer

3 Jerome J. Choi^{1,2}, Noah Cohen Kalafut^{2,4}, Tim Gruenloh³, Corinne D. Engelman¹, Tianyuan Lu^{1,3,◊} &
4 Daifeng Wang^{2,3,4,◊,*}

5 ¹*Department of Population Health Sciences, University of Wisconsin-Madison School of Medicine and Public Health*

6 ²*Waisman Center, University of Wisconsin-Madison*

7 ³*Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison School of Medicine and
8 Public Health*

9 ⁴*Department of Computer Sciences, University of Wisconsin-Madison*

10
11 [◊]*Joint senior authors.*

12 ^{*}*Corresponding author. daifeng.wang@wisc.edu*

13 **Abstract**

14 Single-omics approaches often provide a limited view of complex biological systems, whereas multi-omics
15 integration offers a more comprehensive understanding by combining diverse data views. However, integrat-
16 ing heterogeneous data types and interpreting the intricate relationships between biological features—both
17 within and across different data views—remains a bottleneck. To address these challenges, we introduce
18 COSIME (Cooperative Multi-view Integration and Scalable Interpretable Model Explainer). COSIME uses
19 backpropagation of Learnable Optimal Transport (LOT) to deep neural networks, enabling the learning of
20 latent features from multiple views to predict disease phenotypes. In addition, COSIME incorporates Monte
21 Carlo sampling to efficiently estimate Shapley values and Shapley-Taylor indices, enabling the assessment
22 of both feature importance and their pairwise interactions—synergistically or antagonistically—in predict-
23 ing disease phenotypes. We applied COSIME to both simulated data and real-world datasets, including
24 single-cell transcriptomics, single-cell spatial transcriptomics, epigenomics, and metabolomics, specifically
25 for Alzheimer’s disease-related phenotypes. Our results demonstrate that COSIME significantly improves
26 prediction performance while offering enhanced interpretability of feature relationships. For example, we
27 identified that synergistic interactions between microglia and astrocyte genes associated with Alzheimer’s

28 disease are more likely to be active at the edges of the middle temporal gyrus as indicated by spatial
29 locations. Finally, COSIME is open-source and available for general use.

30 **Introduction**

31 Single-omics approaches, while valuable for providing insights into individual biological layers, often
32 provide a limited view of complex biological systems. Focusing on a single molecular layer—such as
33 genomics, transcriptomics, or metabolomics—provides a limited view of cellular processes, as these layers
34 do not function in isolation. Multi-view data integration overcomes this limitation by combining information
35 from multiple omic layers, offering a more holistic view of biology. This comprehensive approach allows
36 us to identify novel biomarkers, understand the underlying biology of diseases, and improve predictions
37 for disease phenotypes[1, 2]. Machine learning algorithms can efficiently process and analyze large-scale
38 multi-view data, identifying patterns and relationships that are often complex for traditional methods to
39 discern. Moreover, in biological systems, features—especially those from different omic layers—may
40 interact in complex ways that jointly influence phenotypes. Machine learning can help overcome these
41 challenges by modeling and quantifying complex feature interactions, including those across different omic
42 layers.

43 There are a few studies that provide multi-omics data integration and prediction. Cooperative learning[3]
44 is a supervised learning method with multiple sets of features and it combines the usual squared-error loss
45 of predictions with an agreement penalty that encourages input datasets to be similar. However, the strength
46 of the agreement penalty is controlled by a fixed parameter, which is typically set through a tuning, and
47 is not learned dynamically during training. This means that the model does not deeply integrate features
48 in a way that fully captures their complex interdependencies. Moreover, this approach is less flexible in
49 adapting to complex cross-view relationships since it does not learn interactions end-to-end during training.
50 Data Integration Analysis for Biomarker discovery using Latent components (DIABLO)[4] is a supervised
51 learning method for multi-omics integration that seeks for common information across different data types
52 through the selection of a subset of molecular features. Nonetheless, DIABLO assumes that the different
53 omic layers share some level of homogeneity in terms of their relationship to the outcome. In practice, this
54 may not always hold, especially when the omic views measure fundamentally different biological processes
55 that are not strongly linked to each other. Multi-Omics Factor Analysis v2 (MOFA+)[5] is a method for multi-
56 omics integration based on factor analysis and assumes that the latent structure of data can be captured with
57 linear models, designed to analyze and integrate multiple types of omics data in an unsupervised fashion.

58 However, MOFA+ does not have a built-in mechanism for analyzing feature interactions or for providing
59 direct insights into how individual features contribute to the final prediction. Importantly, none of these
60 methods are designed to handle non-linear relationships or feature interactions. Non-linear interactions are
61 especially relevant in biological data, where the effects of one variable (e.g., a gene) on an outcome (e.g.,
62 disease) might not be constant and might vary depending on other factors. Furthermore, deep learning
63 and optimal transport approaches have not been as widely adopted in multi-omics integration as traditional
64 statistical methods.

65 Several approaches have been developed for the interpretability of machine learning models. SHAP
66 (SHapley Additive exPlanations)[6] provides functions to compute feature importance for a particular
67 prediction in machine learning models. While SHAP has optimized algorithms for certain types of models
68 (like Tree explainer for decision trees), it can still be computationally expensive, especially for large datasets
69 and complex models. Moreover, the feature interaction matrix is only available as an output for tree-based
70 models. Additionally, dependence plots display the relationship between only two features, making it difficult
71 to capture the overall patterns of feature interactions in predictions. Local Interpretable Model-agnostic
72 Explanations (LIME)[7] offers localized interpretability, ideal for understanding individual predictions in
73 simple models. LIME works by approximating the model locally around a specific data point using a simple
74 interpretable model, such as linear regression. This means that LIME generates explanations for individual
75 predictions rather than providing a global understanding of the model's behavior. It may miss patterns that
76 are important across the entire dataset, leading to misleading or incomplete explanations if the behavior of
77 the model varies significantly across the feature space. Additionally, LIME does not account for feature
78 interactions, limiting its ability to capture more complex relationships between features.

79 To tackle these challenges, this study introduces Cooperative Multiview Integration and Scalable and
80 Interpretable Model Explainer (COSIME). COSIME features two key components. First, it integrates
81 multi-view data leveraging deep neural network encoders (deep encoders) and Learnable Optimal Transport
82 (LOT) techniques, combining both unsupervised and supervised learning; Second, COSIME implements
83 a mechanism for assessing feature importance within each view, as well as quantifying both within-view
84 and across-view interactions by estimating Shapley values and Shapley-Taylor indices. Through extensive
85 evaluations, we demonstrate the utility of COSIME using both simulated and real-world multi-view datasets.
86 On simulated data, we assessed binary classification and continuous prediction tasks under varying signal
87 levels, using both early and late fusion strategies[8]. For real-world applications, COSIME was applied
88 to Alzheimer's disease diagnosis using transcriptomics-metabolomics and transcriptomics-epigenomics
89 datasets, as well as to predicting Alzheimer's disease progression scores using transcriptomics-spatial

90 transcriptomics data. These analyses highlight the flexibility of COSIME, offering a comprehensive solution
91 for multi-view data analysis.

92 **Results**

93 **Overview of COSIME**

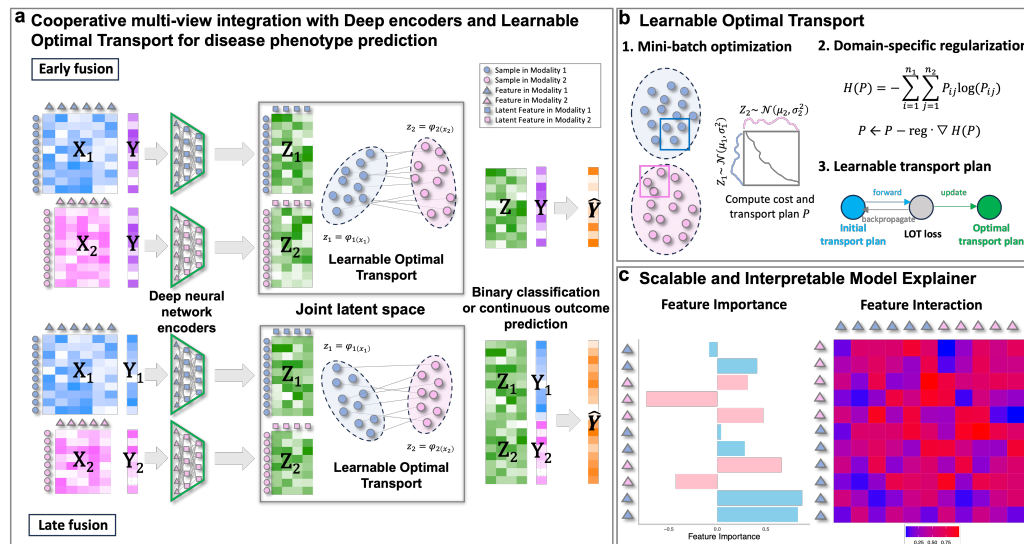
94 An overview of COSIME is illustrated in **Fig. 1**. The first component involves integrating multi-view data
95 for disease phenotype prediction by leveraging deep learning-based encoders, enabling the model to address
96 both linear and non-linear relationships, offering superior flexibility compared to traditional methods that
97 typically rely on linear assumptions. COSIME effectively captures the complex, multi-layered interactions
98 between different omic modalities—such as transcriptomics and epigenomics, transcriptomics and spatial
99 transcriptomics, and transcriptomics and metabolomics—while preserving the distinct features of each
100 data type. Learnable Optimal Transport (LOT) offers three main advantages over conventional alignment
101 methods. First, it uses a learnable transport plan, which allows for dynamic adaptation to complex,
102 heterogeneous data distributions, ensuring more accurate and flexible alignment. Second, LOT supports
103 mini-batch processing, making it scalable for large datasets while maintaining computational efficiency.
104 Finally, LOT can incorporate domain-specific regularization, enabling the model to better handle view-
105 specific differences and improve alignment in specialized contexts. These features make LOT more robust
106 than traditional methods, which often struggle with data misalignment and view-specific challenges.

107 The second component of COSIME focuses on computing feature importance values for each view,
108 as well as pairwise interaction values for both within-view and across-view interactions. To achieve this,
109 COSIME employs the Scalable and Interpretable Model Explainer, which leverages the Shapley values and
110 Shapley-Taylor indices[9] to compute both feature importance and pairwise feature interactions, respectively.
111 By applying Monte Carlo sampling and batch processing, COSIME efficiently estimates feature importance
112 and pairwise feature interactions, enabling faster and more scalable computation while enhancing model
113 interpretability. Additionally, it allows for the identification of the directionality of interactions—whether
114 they exhibit synergism (complementary effects) or antagonism (conflicting effects)—offering deeper insights
115 into how features interact in the model.

116 COSIME enhances the prediction accuracy by improving the modeling of across-view relationships
117 within a shared latent space. It effectively integrates heterogeneous data, enabling more accurate predictions
118 across different biological layers. Moreover, it also provides a powerful interpretability by identifying feature

119 importance and interactions, allowing for deeper insights into how individual features and their combinations
 120 contribute to the model's predictions. While these two main components of COSIME can be used together
 121 for both enhanced prediction and interpretability, they are also designed to be used independently depending
 122 on the specific needs of the analysis. This combination paves the way for identifying key biomarkers and
 123 understanding disease mechanisms. By integrating multi-omics data and assessing feature relationships,
 124 COSIME addresses the challenges of heterogeneous data and complex feature interactions, paving the way
 125 for identifying key biomarkers and understanding disease mechanisms.

Fig 1: Cooperative multi-view integration and Scalable and Interpretable Model Explainer (COSIME).



a COSIME integrates multi-omics data for disease phenotype prediction through a three-step process: (1) Each omic dataset is passed through separate deep neural network encoders (deep encoders). (2) Learnable Optimal Transport (LOT) aligns and merges these features into a joint latent space. (3) The integrated latent representation is then used to predict disease phenotypes. **b** LOT integrates heterogeneous datasets by offering three key advantages: (1) a learnable transport plan that adapts during training, (2) mini-batch processing for scalability, and (3) domain-specific regularization to enhance alignment between source and target distributions. **c** Scalable and Interpretable Model Explainer interprets model predictions by identifying feature importance and interactions. It quantifies both individual feature attributions and pairwise feature interactions, revealing synergistic or antagonistic effects on the model's output.

126 In the following sections, the performance metrics were compared across different methods using the
 127 mean and ± 1.96 times the standard deviation from 5-fold cross-validation. For the best-trained models
 128 across the different multi-view datasets, we computed feature importance for each data view and interaction

129 values both within and across views. All results for the multi-view predictions using COSIME models and
130 other methods are provided in **Supplementary Data 1 and 2**.

131 **Simulation study**

132 Multi-view data were generated by different signal levels (high and low) and types of outcome variables
133 (binary and continuous).

134 Compared to three benchmarking methods—Cooperative Learning (CL)[3], Data Integration Analysis
135 for Biomarker discovery using Latent components (DIABLO)[4], and Multi-Omics Factor Analysis v2
136 (MOFA+)[5] with logistic regression, COSIME early fusion performed best for binary outcome classification
137 with high-signal multi-view datasets (AUROC: 0.845 ± 0.026 , AUPRC: 0.854 ± 0.029 , accuracy: $0.754 \pm$
138 0.054). COSIME late fusion (AUROC: 0.828 ± 0.021 , AUPRC: 0.853 ± 0.021 , accuracy: 0.761 ± 0.019)
139 outperformed the three benchmarking methods (**Fig. 2a**). For binary outcome classification with low-signal
140 multi-view datasets, COSIME late fusion achieved the best performance (AUROC: 0.737 ± 0.017 , AUPRC:
141 0.736 ± 0.025 , accuracy: 0.640 ± 0.052), while COSIME early fusion (AUROC: 0.615 ± 0.038 , AUPRC:
142 0.662 ± 0.039 , accuracy: 0.573 ± 0.027) outperformed CL, DIABLO, and MOFA+ with logistic regression
143 (**Fig. 2a**).

144 We computed feature interaction values for each view using COSIME early fusion for binary outcome
145 prediction with high-signal multi-view datasets. Pairwise feature interactions were categorized into three
146 groups: (Group 1) the interaction terms artificially introduced during data generation, (Group 2) interactions
147 involving latent features, and (Group 3) other interactions. The scaled pairwise feature interaction values
148 for the first 50 features are shown in **Fig. 2b**. Pairwise Wilcoxon rank-sum two-sided tests were performed
149 between these three groups, and the results demonstrate that COSIME successfully captured significant
150 differences in interaction values across the groups. In view A, Group 1 exhibits significantly higher
151 interaction values than both Group 2 (Pairwise Wilcoxon rank-sum two-sided test p -value < 0.05) and
152 Group 3 (Pairwise Wilcoxon rank-sum two-sided test p -value < 0.0001). Additionally, Group 2 shows
153 significantly higher interaction values compared to Group 3 (Pairwise Wilcoxon rank-sum two-sided test
154 p -value < 0.0001). Similarly, in view B, Group 1 shows significantly higher interaction values than both
155 Group 2 (Pairwise Wilcoxon rank-sum two-sided test p -value < 0.001) and Group 3 (Pairwise Wilcoxon
156 rank-sum two-sided test p -value < 0.0001), while Group 2 also has significantly higher interaction values
157 than Group 3 (Pairwise Wilcoxon rank-sum two-sided test p -value < 0.0001).

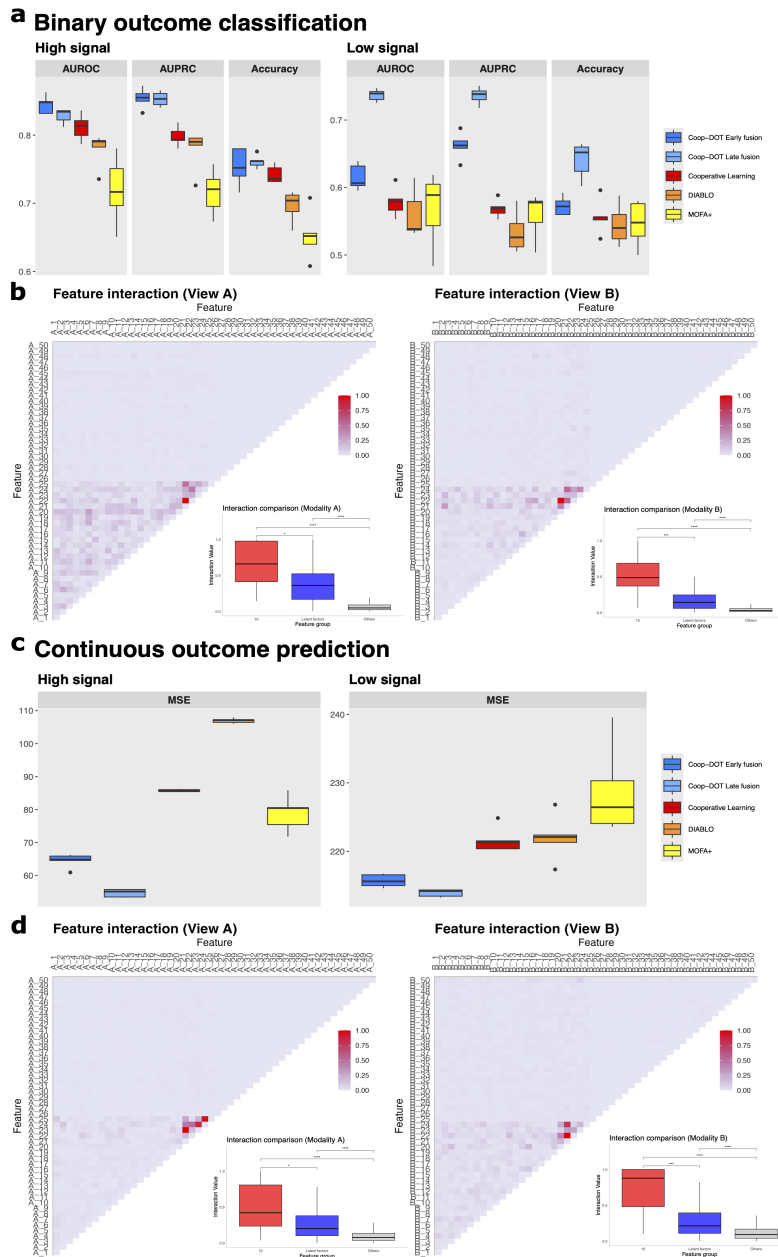
158 For continuous outcome prediction with high-signal multi-view datasets, COSIME delivered the most

159 accurate results (MSE: 54.701 ± 2.391). In comparison, COSIME early fusion outperformed CL, DIABLO,
160 and MOFA+ with regression (MSE: 64.490 ± 4.130) (**Fig. 2c**). When applied to low-signal multi-view
161 datasets, COSIME early fusion achieved the best performance (MSE: 215.707 ± 1.833), while COSIME
162 late fusion showed superior results over the benchmarking methods (MSE: 218.425 ± 1.099) (**Fig. 2c**).

163 Feature interaction values for each view were also calculated using COSIME early fusion for contin-
164 uous outcome prediction with high-signal data. The pairwise feature interactions were divided into three
165 categories: (Group 1) the interaction terms artificially added during data generation, (Group 2) interactions
166 involving latent features, and (Group 3) all other interactions. The scaled pairwise feature interaction values
167 for the first 50 features are shown in **Fig. 2d**. Pairwise Wilcoxon rank-sum two-sided tests were conducted
168 between the three groups, and the findings confirm that COSIME effectively identified significant differences
169 in interaction values across these categories. In view A, Group 1 showed significantly higher interaction val-
170 ues than both Group 2 (Pairwise Wilcoxon rank-sum two-sided test p -value < 0.05) and Group 3 (Pairwise
171 Wilcoxon rank-sum two-sided test p -value < 0.0001), with Group 2 also having significantly higher interac-
172 tion values than Group 3 (Pairwise Wilcoxon rank-sum two-sided test p -value < 0.0001). In view B, Group
173 1 displayed significantly higher interaction values than both Group 2 (Pairwise Wilcoxon rank-sum two-
174 sided test p -value < 0.001) and Group 3 (Pairwise Wilcoxon rank-sum two-sided test p -value < 0.0001),
175 while Group 2 also had significantly higher interaction values compared to Group 3 (Pairwise Wilcoxon
176 rank-sum two-sided test p -value < 0.0001).

177 Overall, both COSIME early fusion and late fusion models outperformed the benchmarking methods
178 in predicting the outcomes. Additionally, COSIME accurately captured the pairwise feature interactions
179 purposely introduced during data generation, validating its ability to identify interactions present in the data.

Fig 2: Prediction performance and feature interaction from simulated data.



a Prediction performance for binary outcome classification using COSIME early fusion, COSIME late fusion, CL, DIABLO, and MOFA+ with logistic regression on both high-signal and low-signal multi-view datasets. **b** Heatmaps depicting the pairwise interactions of the first 50 features from views A and B for binary outcome prediction, along with box plots and pairwise Wilcoxon rank-sum two-sided tests for three interaction groups: (Group 1) 10 interaction terms artificially introduced during data generation, (Group 2) 1,790 interactions involving latent features, and (Group 3) 3,150 other interactions. The box plots display the distribution of interaction values, including the first quartile, median, and third quartile. Asterisks indicate statistical significance: * for p -value < 0.05 and **** for p -value < 0.0001 . **c** Prediction performance for continuous outcome prediction using COSIME early fusion, COSIME late fusion, CL, DIABLO, and MOFA+ with regression on both high-signal and low-signal multi-view datasets. **d** Heatmaps depicting the pairwise interactions of the first 50 features from views A and B for continuous outcome prediction, along with box plots and pairwise Wilcoxon rank-sum two-sided tests for three interaction groups: (Group 1) 10 interaction terms artificially introduced during data generation, (Group 2) 1,790 interactions involving latent features, and (Group 3) 3,150 other interactions. The box plots display the distribution of interaction values, including the first quartile, median, and third quartile. Asterisks indicate statistical significance: *** p -value < 0.001 and **** for p -value < 0.0001 .

180 **Classifying cognitive diagnosis for Alzheimer’s disease from transcriptomic (astro-** 181 **cytes) and metabolomic data**

182 Multi-view data from scRNA-seq of astrocytes and metabolomics were used to test the prediction per-
183 formance of COSIME models. COSIME early fusion with matched samples ($n=2,286$) achieved the best
184 performance for binary classification of Alzheimer’s disease (AD) cognitive diagnosis using multi-view
185 datasets (astrocytes and metabolomics) (AUROC: 0.842 ± 0.011 , AUPRC: 0.864 ± 0.019 , accuracy: 0.773
186 ± 0.018). To leverage all available samples (4292 metacells in transcriptomics and 2286 samples in
187 metabolomics) across views, we also applied COSIME late fusion for unmatched samples, which yielded
188 comparable results (AUROC: 0.804 ± 0.007 , AUPRC: 0.829 ± 0.014 , accuracy: 0.732 ± 0.016) (**Fig. 3a**).

189 Feature importance values were used to prioritize features. The top 20 most important genes (based on
190 absolute values) from the transcriptomic data are displayed in **Fig. 3b**, and enrichment analysis reveals that
191 these genes are strongly associated with AD. Notably, processes such as long-term synaptic potentiation,
192 amyloid-beta ($A\beta$) formation, and amyloid precursor protein (APP) catabolism are strongly implicated in
193 AD. These processes that involve synaptic function, $A\beta$ production, and APP metabolism represent crucial
194 aspects of AD pathophysiology[10, 11].

195 **Fig 3c** shows the top 20 most important metabolites (ranked by absolute feature importance values)

196 from the metabolomic data. Enrichment analysis of these metabolites reveals metabolic pathways such as
197 pantothenate and CoA biosynthesis, arginine biosynthesis, nicotinate and nicotinamide metabolism, and
198 histidine metabolism. These pathways may contribute to AD by exacerbating the brain's energy deficit,
199 promoting oxidative damage, and driving neuroinflammation[12].

200 Pairwise feature interaction values were computed, and the top 50 across-view feature interaction values
201 (based on absolute feature interaction values) are shown in **Fig. 3d**. *PRCP* has synergistic effects with
202 metabolites such as urea and 5-methylthioadenosine (MTA). *PRCP* is involved in peptide degradation,
203 particularly those peptides related to inflammatory responses [13, 14]. Urea is a byproduct of protein
204 metabolism, typically excreted via the kidneys. It plays a role in nitrogen metabolism and helps in the detox-
205 ification of ammonia in the brain [15, 16]. MTA is a byproduct of the methionine salvage pathway and plays
206 a role in maintaining cellular methylation capacity. It is involved in the recycling of S-adenosylmethionine,
207 which is crucial for methylation reactions, including those related to neurotransmitters and DNA [17].
208 Erucate (22:1n9) is a monounsaturated fatty acid, which can influence lipid metabolism, and metabolic
209 homeostasis. The synergism between *PRCP* and MTA suggests a relationship between peptide degradation,
210 methylation capacity, and inflammatory pathways. Given that *PRCP* is involved in degrading inflammatory
211 peptides and MTA plays a role in maintaining methylation for neurotransmitter synthesis and DNA methy-
212 lation, this interaction may indicate a broader involvement of *PRCP* in neuroinflammatory processes and
213 cellular maintenance in AD. Similarly, the interaction between *PRCP* and Erucate (22:1n9) links peptide
214 degradation with lipid metabolism. Erucate (22:1n9), a key fatty acid in lipid metabolism, may interact with
215 *PRCP* to modulate cellular homeostasis, suggesting that *PRCP* might influence both neuroinflammatory
216 and lipid metabolic pathways critical for maintaining brain cell integrity in AD.

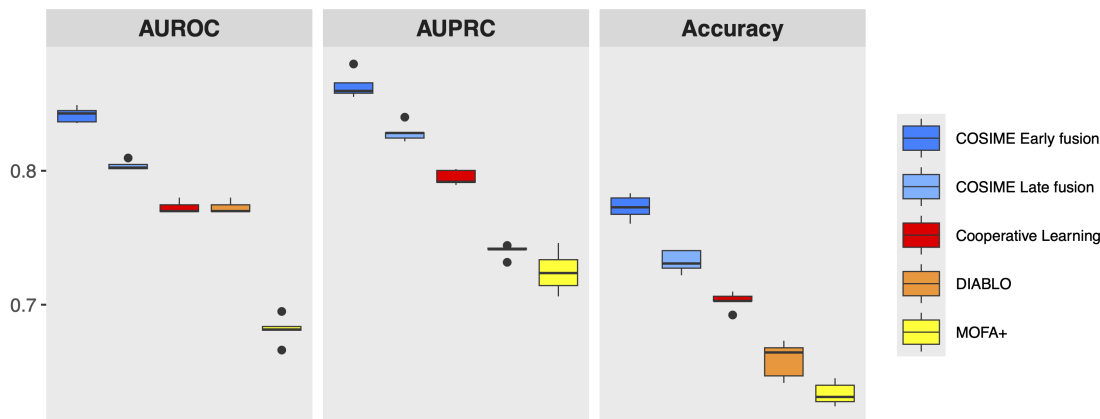
217 *AGBLA* shows antagonistic effects with several metabolites such as arachidonate (20:4n6), 2-O-
218 methylascorbic acid, and N-acetylaspartate (NAA) (**Fig. 3d**). *AGBLA* encodes a protein that is involved
219 in regulating actin dynamics and cell signaling, particularly in processes such as endocytosis, vesicle traf-
220 ficking, and cell motility[18]. Arachidonate (20:4n6) is a polyunsaturated fatty acid and is a precursor
221 for various eicosanoids, including prostaglandins, leukotrienes, and thromboxanes. These molecules play
222 crucial roles in inflammation, immune responses, and cell signaling[19]. 2-O-methylascorbic acid is a
223 derivative of ascorbic acid (vitamin C), and is known for its antioxidant properties[20]. NAA is a metabolite
224 primarily found in neurons and is considered a marker of neuronal integrity and function[21, 22]. *AGBLA*
225 and arachidonate (20:4n6) may exhibit antagonistic interactions, as *AGBLA* may disrupt pro-inflammatory
226 pathways triggered by arachidonate (20:4n6). By inhibiting actin remodeling and vesicular trafficking,
227 which are essential for eicosanoid function, *AGBLA* may affect immune responses and inflammation, pro-

228 cesses central to neurodegeneration. The antagonism between *AGBL4* and 2-O-methylascorbic acid may
229 arise from the involvement of *AGBL4* in promoting cellular processes such as inflammation and oxidative
230 stress, which counteract the antioxidant effects of 2-O-methylascorbic acid. *AGBL4* may influence actin
231 dynamics or cellular responses in a way that increases the need for antioxidants or reduces their efficacy.
232 Additionally, the antagonistic relationship between *AGBL4* and NAA may reflect the impact of *AGBL4*
233 on neuronal integrity. *AGBL4* may alter pathways involved in NAA synthesis or utilization, potentially
234 impairing neuronal health and neurotransmitter balance, as evidenced by reduced NAA levels. Collectively,
235 these interactions may contribute to neurodegenerative processes.

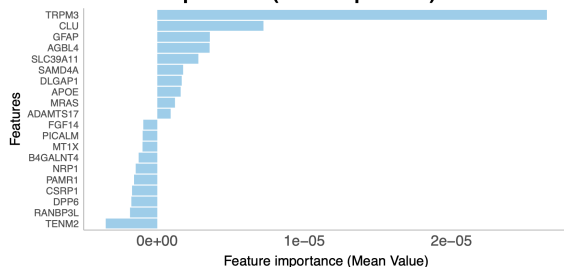
236 Furthermore, *PRCP* and *AGBL4* may play an important role as hub genes in the across-view network
237 (**Fig. 3e**). These two genes are likely involved in astrocyte function under both normal and pathological
238 conditions, particularly through their roles in inflammation, actin remodeling, and neurovascular interac-
239 tions[23, 24]. As hub genes, they may influence across-view networks, connecting pathways that involve
240 peptide signaling, cytoskeletal dynamics, neuroinflammation, and synaptic plasticity. These interactions
241 are crucial for understanding the role of astrocytes in diseases such as AD, where neuroinflammation, astro-
242 cyte reactivity, and synaptic dysfunction are key features[25, 26]. Astrocyte reactivity, often triggered by
243 neuroinflammation, is associated with the increased production of pro-inflammatory cytokines like $\text{TNF-}\alpha$,
244 $\text{IL-1}\beta$, and IL-6 , all of which can worsen neuronal loss and synaptic dysfunction[27, 28]. In this context,
245 *AGBL4* might shift the balance toward a more damaging inflammatory state, thereby increasing the risk of
246 AD.

Fig 3: Classifying cognitive diagnosis for Alzheimer’s disease from transcriptomic (astrocytes) and metabolomic data.

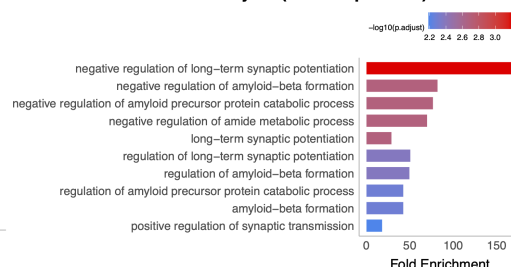
a Prediction performance



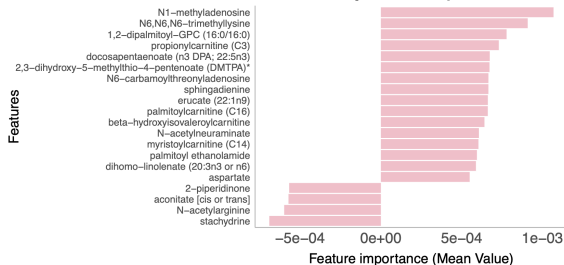
b Feature Importance (Transcriptomics)



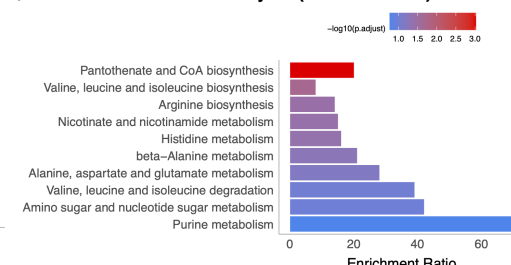
Enrichment analysis (Transcriptomics)



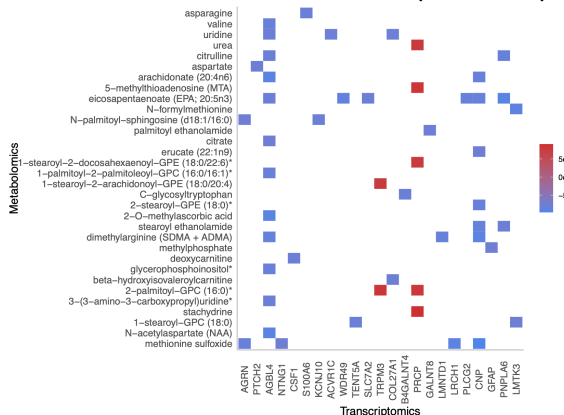
c Feature Importance (Metabolomics)



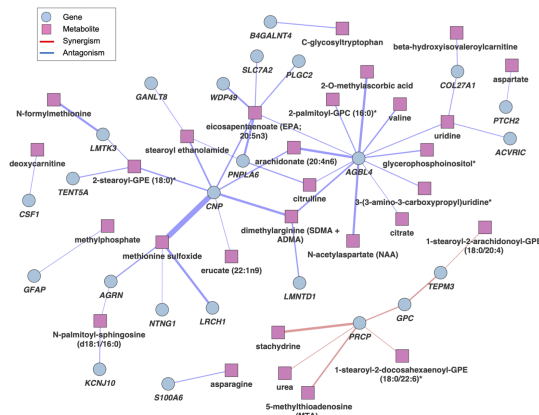
Enrichment analysis (Metabolomics)



d Feature interaction (Across-view)



e Across-view Interaction Network



a Prediction performance using COSIME early fusion, COSIME late fusion, CL, DIABLO, and MOFA+ with logistic regression. **b** Bar plots showing the top 20 prioritized genes (absolute feature importance values) and enrichment analysis for those genes. **c** Bar plots showing the top 20 prioritized metabolites (absolute feature importance values) and enrichment analysis for those metabolites. **d** Heatmap showing the top 50 pairwise across-view (gene-metabolite) interactions. **e** Across-view interaction network showing relationships between the top 50 pairwise across-view (gene-metabolite) interactions.

247 **Classifying cognitive status (dementia) from transcriptomic (oligodendrocyte) and** 248 **epigenomic (oligodendrocyte) data**

249 Multi-view data from both scRNA-seq and scATAC-seq of oligodendrocytes were utilized to evaluate
250 the predictive performance of COSIME models. COSIME early fusion performed the best for binary
251 classification of cognitive status (dementia) using multi-view datasets (oligodendrocytes) (AUROC: 0.874
252 ± 0.011 , AUPRC: 0.805 ± 0.013 , accuracy: 0.794 ± 0.041) (**Fig. 4a**).

253 We investigated pairwise interactions between the ten known key TFs involved in oligodendrocytes[29].
254 For instance, synergistic interactions were identified between SOX10 and MAZ as well as between MYRF
255 and MAZ. Specifically, SOX10 is critical for oligodendrocyte differentiation and can be involved in reactive
256 gliosis or attempts to repair myelin in AD and dementia[30]. MAZ is involved in regulating genes related to
257 cell proliferation, differentiation, and survival[31]. In dementia, even though MAZ is not oligodendrocyte-
258 specific, the presence of both SOX10 and MAZ may indicate synergism which is a disruption of normal
259 oligodendrocyte functions and myelin loss. MYRF is another key player in myelination that works closely
260 with driving factors of oligodendrocyte differentiation and myelin gene expression[32]. For similar reasons,
261 synergism between MYRF and MAZ may indicate that oligodendrocyte differentiation is promoted in
262 an attempt to repair myelin loss, a process that is common in dementia cases due to demyelination and
263 oligodendrocyte dysfunction (**Fig. 4b**).

264 Meanwhile, antagonistic interaction was identified between SOX8 and HIF1A. SOX8 is critical in
265 regulating cell fate decisions, development, and differentiation across various tissues[33] while HIF1A, a
266 stress and inflammatory-related factor[34], could suppress differentiation and antagonistically contributing
267 to the pathology of AD and dementia. SOX8 is also antagonistically interacting with PBX1, which is
268 important for the development and differentiation of neural cells, including oligodendrocytes[35, 36]. PBX1
269 could be involved in AD pathology by disrupting oligodendrocyte differentiation and myelin integrity, while
270 SOX8 might counteract these effects and support oligodendrocyte survival. Their antagonistic interaction

271 may suggest that an imbalance between these two TFs may contribute to the progression of the dementia
272 (**Fig. 4b**).

273 Based on the top fifty across-view absolute interaction values (TF expression levels and ATAC peak
274 counts), we identified several TFs that are either synergistically or antagonistically interacting with multiple
275 oligodendrocyte-specific regulatory regions. We found synergistic interactions between MAZ and ATAC
276 peak regions that map to multiple genes, including *PTPRF*[37], *ERBB3*[38], *GRID1*[39], *LIMCH1*[40],
277 *FLNC*[41], and *DUSP7*[42] (based on genomic annotations). Since increased chromatin openness likely
278 promotes the expression of these genes associated with neurodegeneration, myelin repair, and cell stress
279 responses, the identified synergistic interactions further support the potential role of MAZ in dementia
280 progression. Furthermore, ZNF135 is a zinc finger TF known for its involvement in regulating chromatin
281 structure and transcriptional activation, particularly in response to cellular stress[43]. Synergistic inter-
282 actions were identified between ZNF135 and ATAC peak regions that map to genes involved in neural
283 development, synaptic signaling, axon guidance, and neurotransmitter regulation, such as *NKX2-2*[44],
284 *PTPRF*[37], *CNTN2*[45], *PSD2*[46], and *SLC6A9*[47], which may reflect the activation of disease-related
285 pathways during the progression of neurodegeneration (**Fig. 4c**).

286 In contrast, PBX1 acts antagonistically with several open chromatin regions, as its upregulation may
287 reduce chromatin accessibility, which is critical for oligodendrocyte function and myelin maintenance. This
288 antagonistic action may prevent the activation of protective genes such as *CPS1*[48], and *CYTH1*[49] in
289 those open chromatic regions, leading to impaired neuroinflammatory responses and exacerbating neurode-
290 generative processes. The ability of PBX1 to close chromatin at disease-associated regions may contribute
291 to a negative predictive signal for dementia, as it correlates with the loss of protective mechanisms in
292 oligodendrocytes (**Fig. 4c**).

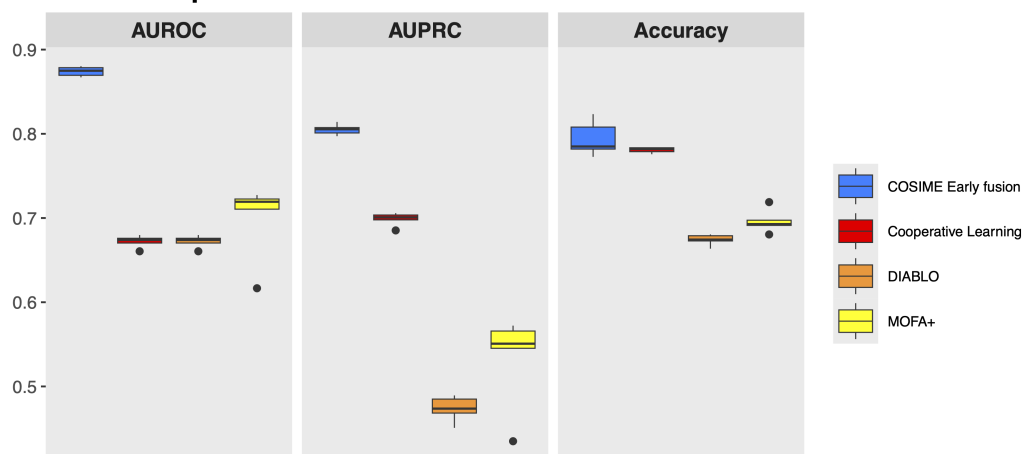
293 It is crucial to examine the regulatory mechanisms underlying the interactions between TFs, ATAC
294 peak regions, target genes, and dementia. Both STAT3 and JUN are synergistic TFs, and they also
295 exhibit synergistic activity in the genomic region chr17:78,351,242-78,361,401, which contains the *SOCS3*
296 gene. STAT3 is a critical TF involved in the JAK-STAT signaling pathway, regulating immune responses,
297 inflammation, cell survival, and differentiation[50, 51]. JUN, a component of the AP-1 transcription factor
298 complex, regulates genes involved in cell proliferation, apoptosis, and stress responses, particularly in
299 response to inflammatory signals[52, 53]. STAT3 and JUN are known to regulate *SOCS3*, a gene that is
300 upregulated in dementia. *SOCS3* plays a crucial role in regulating the JAK-STAT pathway by inhibiting
301 STAT3 activation[54, 55]. These connections suggest that the coordinated activity of STAT3 and JUN
302 in this region may contribute to the dysregulation of *SOCS3*, linking it to the neuroinflammatory and

303 neurodegenerative pathways associated with dementia (**Fig. 4d**).

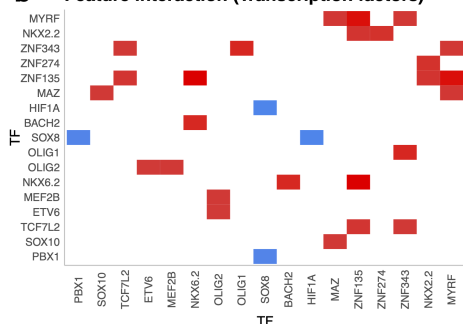
304 The four key TFs involved in oligodendrocyte function — SOX10[30], MYRF[30], OLIG1[56], and
305 OLIG2[56]— exhibit synergistic relationships, except for SOX10 and OLIG1, which interact antagonisti-
306 cally. SOX10, MYRF, OLIG1, and OLIG2 are essential TFs that regulate oligodendrocyte differentiation,
307 myelination, and myelin maintenance. In the context of dementia, their synergistic interactions suggest that
308 the activities of these TFs contribute to increased myelin damage and oligodendrocyte dysfunction, which
309 accelerates the progression of neurodegeneration. Notably, all four of these TFs antagonistically interact
310 with a ATAC peak in a specific genomic region, chr18:77,049,711-77,059,094, where the gene *MBP* is
311 located. These interactions suggest that these TFs may suppress the expression of *MBP*, which is essential
312 for the formation and maintenance of the myelin sheath in oligodendrocytes[57] (**Fig. 4d**). Taken together,
313 the characterization of these interactions provides valuable insights into the underlying mechanisms that
314 may contribute to the pathogenesis and progression of dementia.

Fig 4: Classifying dementia from transcriptomic (oligodendrocyte) and epigenomic (oligodendrocyte) data

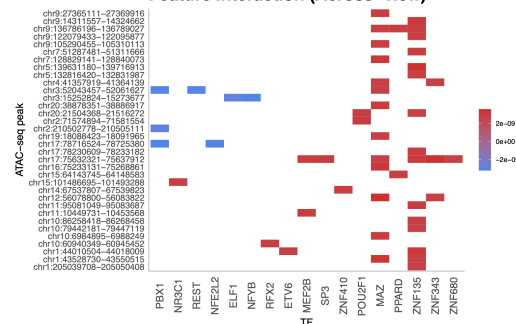
a Prediction performance



b Feature interaction (Transcription factors)

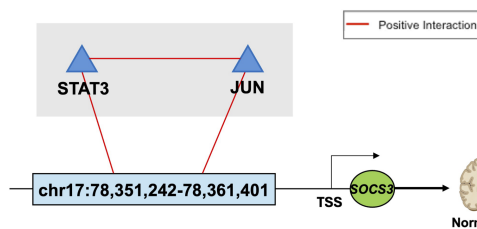


c Feature interaction (Across-view)

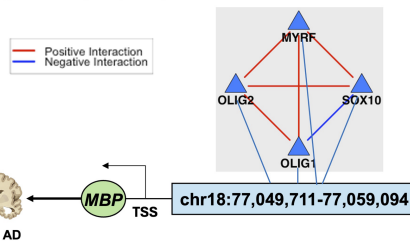


d

Synergistic (TF-peak) interaction



Antagonistic (TF-peak) interaction



a Prediction performance using COSIME early fusion, CL, DIABLO, and MOFA+ with logistic regression. **b** Heatmap showing the top 50 pairwise TF interactions. **c** Heatmap showing the top 50 pairwise across-view (TF-ATAC peak) interactions. **d** Diagrams illustrating examples of synergistic and antagonistic relationships between TF-ATAC peak-AD phenotype.

315 **Predicting Alzheimer’s disease progression scores from transcriptomic (microglia)**
 316 **and spatial transcriptomic MERFISH (astrocyte) data**

317 To assess the predictive capability of COSIME models, multi-view data combining scRNA-seq of microglia
 318 and single-cell spatial transcriptomics (MERFISH) of astrocytes were applied. COSIME early fusion
 319 performed better than the other benchmarking methods for continuous outcome prediction of Alzheimer’s
 320 disease progression scores using multi-view datasets (microglia and astrocytes) (MSE:0.033 ± 0.005) (**Fig.**
 321 **5a**).

322 **Fig. 5b** shows the top 20 important genes (ranked by mean absolute feature importance values) in
 323 each view. In microglia, there are genes closely related to AD through their roles in neuroinflammation,
 324 $A\beta$ clearance, and microglial activation, among those with a positive feature importance. *SPP1* is a
 325 well-established microglial marker and plays a critical role in the activation and inflammatory response
 326 of microglia. It has been shown to regulate phagocytosis of amyloid plaques and can contribute to
 327 both neuroprotection and neurodegeneration[58]. *APOE* genotype is a major genetic risk factor for AD,
 328 particularly the E4 allele, where microglial *APOE* is involved in $A\beta$ clearance. It has been shown that *APOE*
 329 e4 (the risk allele) impairs this function and enhances neuroinflammation[59, 60]. *CIQB* encodes one of the
 330 subunits of C1q, which is part of the complement system and a major player in immune responses[61, 62].
 331 C1q marks apoptotic or damaged neurons for phagocytosis by microglia. These genes are likely upregulated
 332 in AD and often reflect microglial dysfunction, contributing to both increased neuroinflammation and
 333 impaired $A\beta$ clearance. On the other hand, the genes that have negative feature importance values may
 334 be downregulated in microglia and they are crucial to the immune response, microglial surveillance, and
 335 synaptic pruning in AD. For example, *BINI* has been implicated in AD and is involved in regulating synaptic
 336 function and endocytosis[63].

337 In the astrocyte MERFISH data, the genes that have negative feature importance values may suggest
 338 dysfunctional astrocytes with impaired synaptic regulation, reduced neuroprotection, and failure to clear
 339 neurotoxic debris in AD. For instance, *RYR3* is involved in calcium signaling in astrocytes, which plays a
 340 role in synaptic regulation and glutamate uptake[64, 65]. Downregulation of *RYR3* could impair calcium
 341 signaling and neurotransmitter homeostasis, contributing to synaptic dysfunction in AD. In contrast, the
 342 genes that have positive feature importance values are involved in neuroinflammation, astrocyte differentia-
 343 tion, synaptic function, and neurodegeneration, and their upregulation in AD suggests activated astrocytes
 344 contributing to neuroinflammatory responses and synaptic damage. *DGKG* is involved in signal trans-
 345 duction through lipid metabolism and inflammatory response. Its upregulation in astrocytes may promote

346 increased neuroinflammation in AD. Also, *DGKG* has been linked to glutamate receptor signaling, which
347 could disrupt synaptic communication, a hallmark of AD (**Fig. 5b**)[66]. Enrichment analyses for the top 20
348 prioritized genes (absolute feature importance values) for both microglia and astrocytes were implemented
349 (**Supplementary Figure 6**). The top 10 categories in each analysis show that most of them are closely
350 related to AD.

351 We investigated the spatial distribution of cells with the expression levels and feature importance values
352 of *RYR3* (**Fig. 5c**). As a result, cells with the highest 25% of feature importance were more likely to
353 be located in the central regions of the middle temporal gyrus, where the astrocyte MERFISH data were
354 collected, compared to the remaining 75% (T-statistic: 18.783, two-tailed t-test p -value < 0.0001), based
355 on the Euclidean distance from each cell to the midpoint of the bounding box. Similarly, when comparing
356 the lowest 25% of feature importance values to the remaining 75% for *DGKG*, we found that cells with
357 highly negative feature importance values were more likely to be located at the center of the middle temporal
358 gyrus. The two-tailed t-test (T-statistic: -2.861 and two-tailed t-test p -value < 0.01) revealed a significant
359 difference, suggesting that genes with high negative effects are preferentially active in central regions. These
360 findings further support the observation that upregulated genes associated with AD tend to be more active at
361 the edges of the middle temporal gyrus, while downregulated genes are concentrated in the central regions
362 (**Fig. 5c**).

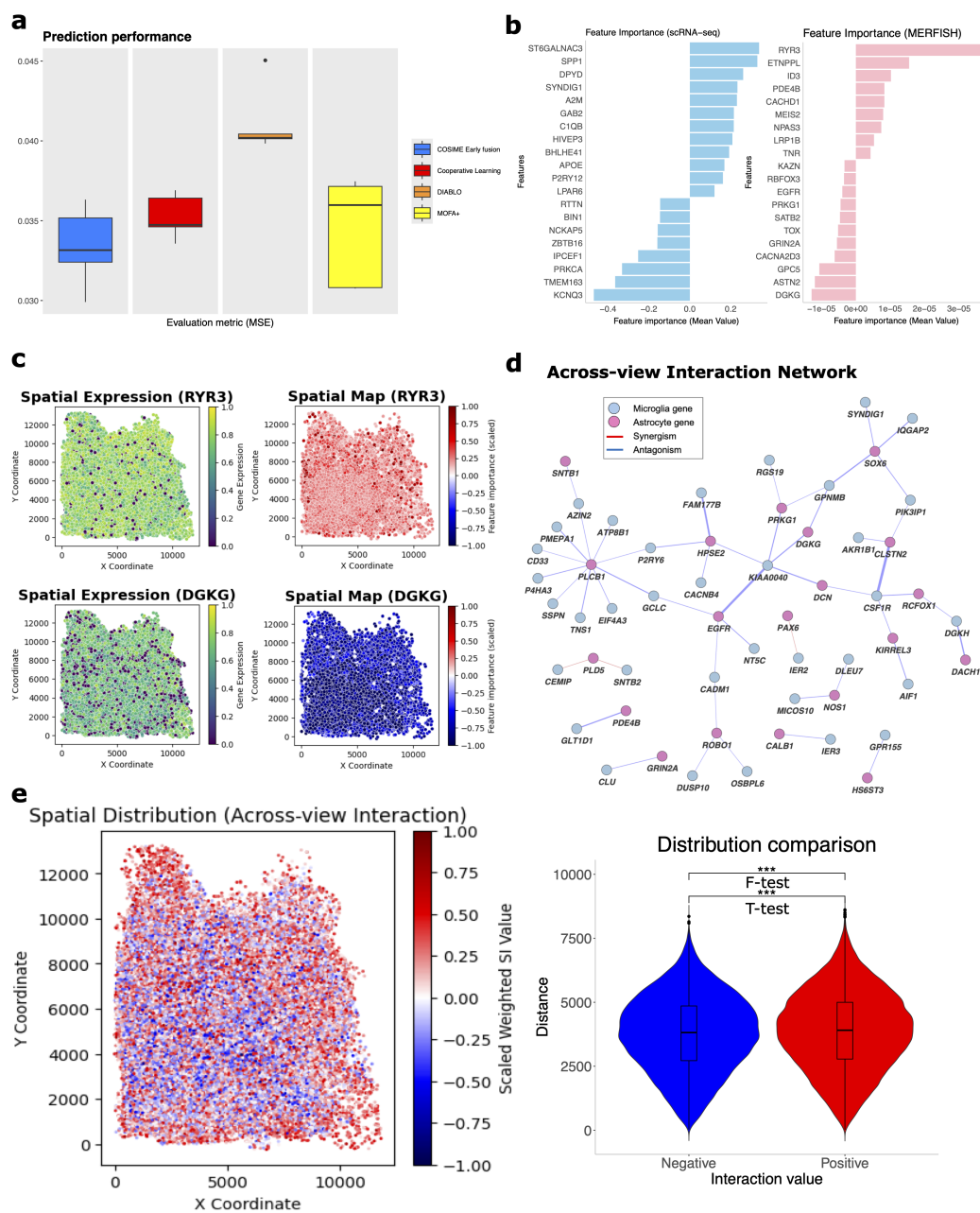
363 We constructed a gene-gene interaction network based on the top 50 absolute across-view interaction
364 values. Several astrocyte-specific MERFISH genes, such as *PLCB1* and *HPSE2*, were identified as key hub
365 genes, linking multiple microglia genes. *PLCB1*, in particular, is involved in critical signaling pathways
366 that regulate cellular functions such as inflammation, synaptic transmission, and cell differentiation. As an
367 astrocyte-specific gene, *PLCB1* plays a central role in intracellular signaling, including phosphatidylinositol
368 signaling, which could modulate microglial activity and other immune-related processes in the context
369 of AD[67]. Similarly, *HPSE2*, another astrocyte-associated gene, regulates heparan sulfate metabolism, a
370 process involved in cellular signaling, extracellular matrix remodeling, and inflammation[68]. Its interaction
371 with microglial genes suggests that *HPSE2* may influence neuroinflammatory responses and contribute to
372 the pathogenesis of AD (**Fig. 5d**).

373 A spatial map was generated using all across-view gene-gene interaction values (**Fig. 5e**). Similar
374 to the approach used in the analyses for **Fig. 5c**, we computed the Euclidean distance from each cell to
375 the midpoint of the bounding box and performed a two-tailed t-test to assess the null hypothesis that the
376 distances between cells in positive (synergistic) interaction values and the midpoint of the bounding box are
377 equal to the distances between negative (antagonistic) interaction values) and the midpoint of the bounding

378 box were the same. We also conducted a two-tailed F-test to assess the null hypothesis that the variances of
379 the two groups were equal. The null hypothesis was rejected in both tests (T-statistic: 18.783 (two-tailed
380 t-test p -value < 0.001) and F-statistic: 1.242 (two-tailed t-test p -value < 0.001)) (**Fig. 5e**). These results
381 suggest that synergistic interactions between microglia and astrocyte genes associated with AD are more
382 likely to be active at the edges of the middle temporal gyrus, while antagonistic interactions tend to be
383 localized more toward the center.

384 In AD, it has been shown that certain brain regions are more vulnerable to neurodegeneration, including
385 parts of the temporal lobe, which includes the middle temporal gyrus[69]. The middle temporal gyrus
386 is an important region for memory processing, semantic processing, and cognitive functions[70]. It is
387 particularly vulnerable to early pathological changes in AD, such as amyloid plaque deposition and tau
388 tangles[71]. The edges of middle temporal gyrus, like the cortical periphery, can be more exposed to
389 various factors such as inflammation, glial activation, and synaptic dysfunction, which are key drivers of
390 neurodegeneration in AD[72]. As a result, it is plausible that genes at the edges of the middle temporal gyrus
391 are involved in triggering AD progression, especially through their potential role in glial activation (e.g.,
392 neuroinflammation), synaptic dysfunction, and neurodegeneration. These regions are often early hotspots
393 for AD pathology, and the genes expressed there may contribute to the progression of AD.

Fig 5: Predicting Alzheimer's disease progression scores from transcriptomic (microglia) and spatial transcriptomic MERFISH (astrocyte) data.



a Prediction performance using COSIME early fusion, CL, DIABLO, and MOFA+ with regression. **b** Bar plots showing the top 20 prioritized genes (absolute feature importance values) in each view. **c** Spatial plots showing expression level and importance values across cells for *RYR3* and *DGKG*. **d** Across-view interaction network showing relationships between the top 50 pairwise across-view (gene-gene) interactions. **e** Spatial distribution of all across-view interaction values colored by the directionality and the distribution of 80,072 positive and 68,619 negative interaction values, with statistical tests (two-tailed t-test and two-tailed F-test) applied to assess differences between two groups. The violin plots including box plots display the distribution of interaction values in the two groups. Asterisk indicate statistical significance: *** for p -value < 0.001 .

394 Discussion

395 We introduce COSIME (Cooperative Multiview Integration and Scalable Interpretable Model Explainer)
 396 and apply it to various multi-view data sets, including simulated data, transcriptomics-metabolomics,
 397 transcriptomics-epigenomics, and transcriptomics-spatial transcriptomics. COSIME has two key compo-
 398 nents: first, it integrates multi-view data for disease phenotype prediction by leveraging deep encoders and
 399 LOT techniques. COSIME combines both unsupervised and supervised learning. Second, it computes
 400 feature importance values for each view and pairwise feature interaction values for both within-view and
 401 across-view.

402 COSIME addresses several challenges in the existing literature. First, COSIME can handle both linear
 403 and nonlinear relationships as well as feature interactions within the data. The deep encoders transform the
 404 data into meaningful embeddings without relying on linear assumptions. Once encoded, these embeddings
 405 can be used with either linear or nonlinear methods for disease phenotype prediction. Second, COSIME
 406 accommodates both matched and unmatched samples across two data views. Each view is processed through
 407 a separate deep encoder. For matched samples, COSIME supports either early fusion or late fusion, while
 408 for unmatched samples, it uses late fusion. LOT optimally aligns embeddings from the two views into a
 409 joint latent space. Third, LOT is a novel technique that aligns the latent space distributions of two distinct
 410 data views by treating the transport plan as a learnable parameter, which is dynamically optimized during
 411 training. This allows LOT to adapt more flexibly to complex data structures, align latent spaces in a shared
 412 representation space, and adjust its alignment strategy to the specific characteristics of the data, making it
 413 particularly effective for downstream tasks such as prediction and classification. Fourth, COSIME efficiently
 414 computes feature importance scores for individual features and their pairwise interactions through Monte
 415 Carlo sampling and batch processing in any complex machine learning models. Lastly, COSIME can identify

416 the directionality of feature interactions, distinguishing between synergism and antagonism. This enhances
417 model interpretability and can be applied to a wide range of machine learning models. COSIME outputs
418 feature importance and interaction matrices that are both straightforward and intuitive.

419 To the best of our knowledge, COSIME is the first machine learning method that incorporates the
420 LOT concept to integrate two heterogeneous multi-view datasets for predictive modeling. Given that the
421 relationships in omics data—such as transcriptomics, metabolomics, epigenomics, and spatial transcrip-
422 tomics—are often nonlinear, its capacity to manage such complexity is crucial. Existing methods such as
423 Cooperative Learning[3], DIABLO[4], and MOFA+[5] are not designed to capture nonlinear relationships.
424 We have demonstrated that COSIME models (both early fusion and late fusion) outperform these methods
425 (Cooperative learning, DAILBO, and MOFA+ with regression models) in predicting disease phenotypes.
426 Moreover, unlike other existing methods, COSIME is able to output both feature importance and interaction
427 values from any complex machine learning models. SHAP[6] has a function to output a feature interaction
428 matrix only available for tree-based models and LIME[7] does not account for feature interactions and is
429 more suitable for interpreting individual predictions in simpler models, such as linear regression. COSIME
430 is user-friendly, allowing users to either train a model using their data, compute feature importance and
431 interaction values with a pre-trained model and their data, or do both. Moreover, users can choose a fusion
432 method (early or late) based on their research objectives and data, and can also choose to compute feature
433 importance, feature interactions, or both.

434 We demonstrate that COSIME outperforms existing well-known multi-omics models in integrating
435 and predicting outcomes, while also providing valuable feature importance and interaction information.
436 However, there are several areas where COSIME can be further improved and updated. First, while
437 COSIME currently computes pairwise feature interaction values, there may be meaningful interactions
438 involving more than two features. Expanding the model to account for higher-order feature interactions
439 could enhance predictive power but would require more sophisticated algorithms and optimizations. Second,
440 determining the optimal number of hidden layers in the deep encoders, along with the size of the latent
441 space, is a complex task. Although the dimensionality of the joint latent space has been fine-tuned, exploring
442 the number of hidden layers will be addressed in future work to further optimize model performance. To
443 keep pace with rapidly advancing technologies in single-cell sequencing, omics data generation, and the
444 development of state-of-the-art methods, we will continue to update COSIME in future work. Lastly,
445 COSIME can be applied to additional brain-related diseases, such as neuropsychiatric disorders[73, 74],
446 which were not explored in this study, once large-scale multi-omics data become available in the future.
447 Moreover, COSIME can be integrated with other methods[75–77] when multi-view imputation for missing

448 data is needed, expanding its versatility in handling diverse datasets.

449 **Methods**

450 **COSIME overview**

451 Cooperative Multiview Integration and Scalable and Interpretable Model Explainer (COSIME) is a machine
452 learning model that integrates multi-view data for disease phenotype prediction and computes feature
453 importance and interaction scores. By leveraging deep learning-based encoders, COSIME effectively
454 captures the complex, multi-layered interactions between different omic modalities while preserving the
455 unique characteristics of each data type. The integration of LOT techniques aligns and merges heterogeneous
456 datasets, improving the accuracy of modeling across-view relationships in the joint latent space. In addition,
457 COSIME leverages the Shapley-Taylor Interaction Index[9] to compute feature importance and interaction
458 values, allowing for a deeper understanding of how individual features and their interactions contribute to
459 the predictions (**Fig. 1**).

460 **Simulated Data**

461 We generated synthetic data by creating two distinct datasets, \mathbf{X}_1 and \mathbf{X}_2 , each with a specified number of
462 features, to simulate a scenario for multi-view analysis. The function `generate_data` was used to simulate
463 these datasets. Latent factors, which are unobserved variables that affect the data, were introduced with
464 specified standard deviations (`std`, `std1`, and `std2`) and were either correlated or uncorrelated depending on
465 the experimental design. The strength of the influence of these latent factors on the multi-view datasets
466 was controlled by the scaling factors `scale1` and `scale2`. Additionally, interaction effects between the latent
467 factors of \mathbf{X}_1 and \mathbf{X}_2 were optionally included to introduce non-linear relationships, with the weight of these
468 interactions controlled by the `interaction_weight` parameter.

469 **Latent Factors and Dataset Generation**

470 We generate the latent factors for the datasets \mathbf{X}_1 and \mathbf{X}_2 based on the correlation setting. If the latent
471 factors are correlated, we use a shared latent factor matrix \mathbf{U} while if the latent factors are uncorrelated, we
472 use separate latent factor matrices \mathbf{U}_1 and \mathbf{U}_2 .

473 **Latent Factor Generation:** The latent factors for both \mathbf{X}_1 and \mathbf{X}_2 are drawn from normal distributions:

$$\mathbf{U}_1 \sim \mathcal{N}(0, \sigma_1^2)_{n \times f_{\text{latent}}} \quad (1)$$

474

$$\mathbf{U}_2 \sim \mathcal{N}(0, \sigma_2^2)_{n \times f_{\text{latent}}} \quad (2)$$

475 Where:

476 • n is the number of samples.477 • f_{latent} is the number of latent factors.478 • σ_1^2, σ_2^2 are the variances of the latent factors affecting \mathbf{X}_1 and \mathbf{X}_2 , respectively.479 **Correlated Latent Factors:** If the latent factors in the two data views are correlated, we use a shared480 latent factor matrix \mathbf{U} , and both datasets \mathbf{X}_1 and \mathbf{X}_2 are generated as:

$$\mathbf{X}_1 = \mathbf{U} \cdot \beta_1 + \epsilon_1 \quad (3)$$

481

$$\mathbf{X}_2 = \mathbf{U} \cdot \beta_2 + \epsilon_2 \quad (4)$$

482 Where:

483 • \mathbf{U} is the shared latent factor matrix.484 • β_1, β_2 are the weights for \mathbf{X}_1 and \mathbf{X}_2 , respectively.485 • ϵ_1, ϵ_2 represent the noise for \mathbf{X}_1 and \mathbf{X}_2 , respectively, and follow:

$$\epsilon_1 \sim \mathcal{N}(0, \sigma_{\epsilon_1}^2) \quad (5)$$

486

$$\epsilon_2 \sim \mathcal{N}(0, \sigma_{\epsilon_2}^2) \quad (6)$$

487 **Uncorrelated Latent Factors:** If the latent factors in the two data views are uncorrelated, we generate488 separate latent factor matrices \mathbf{U}_1 and \mathbf{U}_2 for \mathbf{X}_1 and \mathbf{X}_2 , and the datasets are generated as:

$$\mathbf{X}_1 = \mathbf{U}_1 \cdot \beta_1 + \epsilon_1 \quad (7)$$

489

$$\mathbf{X}_2 = \mathbf{U}_2 \cdot \beta_2 + \epsilon_2 \quad (8)$$

490 Where:

- 491 • \mathbf{U}_1 and \mathbf{U}_2 are the uncorrelated latent factor matrices for \mathbf{X}_1 and \mathbf{X}_2 .
- 492 • β_1, β_2 are the weights for \mathbf{X}_1 and \mathbf{X}_2 , respectively.
- 493 • ϵ_1, ϵ_2 represent the noise for \mathbf{X}_1 and \mathbf{X}_2 , respectively, and follow:

$$\epsilon_1 \sim \mathcal{N}(0, \sigma_{\epsilon_1}^2) \quad (9)$$

494

$$\epsilon_2 \sim \mathcal{N}(0, \sigma_{\epsilon_2}^2) \quad (10)$$

495 Interaction Effects

496 **Correlated Latent Factors:** When the latent factors in the two data views are correlated, the interaction
 497 term is generated by considering pairwise combinations of the latent factors. Specifically, the interaction
 498 term $\mathbf{z}_{\text{interaction}}$ is calculated as:

$$\mathbf{z}_{\text{interaction}} = \sum_{i=f_{\text{latent}}+1}^{f_{\text{latent}}+5} \sum_{j=f_{\text{latent}}+1}^{f_{\text{latent}}+5} w \cdot U_i \cdot U_j \quad (11)$$

499 Where:

- 500 • U_i and U_j are the elements of the latent factor matrix \mathbf{U} .
- 501 • The interaction term involves the product of latent factors U_i and U_j for each pair.
- 502 • The sum is performed over the indices i and j , which range from $f_{\text{latent}} + 1$ to $f_{\text{latent}} + 5$, where the
 503 interaction is modeled between latent factors.
- 504 • w is the interaction weight, a constant that determines the strength of the interaction between latent
 505 factors.

506 The interaction term $\mathbf{z}_{\text{interaction}}$, which is given by equation 11, is then added to the feature matrices \mathbf{X}_1
 507 and \mathbf{X}_2 as follows:

$$x_{1,i} = x_{1,i} + \mathbf{z}_{\text{interaction}} \quad (12)$$

$$x_{2,j} = x_{2,j} + \mathbf{z}_{\text{interaction}} \quad (13)$$

508 Where equations 12 and 13 show how the interaction term $\mathbf{z}_{\text{interaction}}$ is added to the features x_1 and x_2
 509 in the correlated case.

510 **Outcome Variable**

511 The outcome variables y_1 and y_2 are generated based on the feature matrices \mathbf{X}_1 and \mathbf{X}_2 , taking into account
512 whether interaction terms are included or not.

513 **Correlated Latent Factors** For the case where the latent factors are correlated (i.e., a shared latent factor
514 matrix \mathbf{U} is used for both \mathbf{X}_1 and \mathbf{X}_2), the outcome generation is as follows:

515 • **With Interaction:** When interaction terms are included, the expected outcome μ_{all} is computed as
516 the linear combination of the latent factors \mathbf{U} and the latent factor weights β , along with the pairwise
517 interaction terms between the additional latent factors. The interaction term $\mathbf{z}_{\text{interaction}}$, given by
518 equation 11 is used, which is defined as:

$$\mu_{\text{all}} = \mathbf{U}\beta + \mathbf{z}_{\text{interaction}} \quad (14)$$

519 • **Without Interaction:** When no interaction terms are included, the expected outcome is simply the
520 linear combination of the latent factors \mathbf{U} weighted by β :

$$\mu_{\text{all}} = \mathbf{U}\beta \quad (15)$$

521 **Uncorrelated Latent Factors** For the case where the latent factors are uncorrelated (i.e., separate latent
522 factor matrices \mathbf{U}_1 and \mathbf{U}_2 are used for \mathbf{X}_1 and \mathbf{X}_2 , respectively), the outcome generation is as follows:

523 • The expected outcome is the sum of the linear combinations of \mathbf{U}_1 and \mathbf{U}_2 with their respective
524 weights $\beta_{\mathbf{U}_1}$ and $\beta_{\mathbf{U}_2}$:

$$\mu_{\text{all}} = \mathbf{U}_1\beta_{\mathbf{U}_1} + \mathbf{U}_2\beta_{\mathbf{U}_2} \quad (16)$$

525 **Outcome Variable Generation:** The outcome variables y_1 and y_2 are generated using the expected
526 outcomes μ_{all} computed above, along with Gaussian noise for continuous outcomes or using a logistic
527 function for binary outcomes.

528 • **Continuous Outcomes:** For continuous outcomes, the outcome variables y_1 and y_2 are generated by
529 adding Gaussian noise to the expected values:

$$\mathbf{y}_1 = \mu_{\text{all}} + \epsilon_1, \quad \epsilon_1 \sim \mathcal{N}(0, \sigma^2) \quad (17)$$

530

$$\mathbf{y}_2 = \mu_{\text{all}} + \epsilon_2, \quad \epsilon_2 \sim \mathcal{N}(0, \sigma^2) \quad (18)$$

531 Where ϵ_1 and ϵ_2 are Gaussian noise terms with mean 0 and variance σ^2 .

532 • **Binary Outcomes:** For binary outcomes, the predicted probabilities \hat{y}_1 and \hat{y}_2 are computed using
533 the logistic function:

$$\hat{y}_1 = \frac{1}{1 + \exp(-\mu_{\text{all}})} \quad (19)$$

534

$$\hat{y}_2 = \frac{1}{1 + \exp(-\mu_{\text{all}})} \quad (20)$$

535 The outcome variables are then generated based on Bernoulli trials with added Gaussian noise to the
536 predicted probabilities:

$$\mathbf{y}_1 \sim \text{Bernoulli}(p = \hat{y}_1 + \epsilon_3), \quad \epsilon_3 \sim \mathcal{N}(0, \sigma^2) \quad (21)$$

537

$$\mathbf{y}_2 \sim \text{Bernoulli}(p = \hat{y}_2 + \epsilon_4), \quad \epsilon_4 \sim \mathcal{N}(0, \sigma^2) \quad (22)$$

538 Where ϵ_3 and ϵ_4 are Gaussian noise terms with mean 0 and variance σ^2 .

539 **Experimental Conditions**

540 We generated a total of eight experimental conditions, divided between binary and continuous outcome
541 variables. Each condition included two distinct datasets to facilitate multi-view analysis. The conditions
542 were based on two factors: strength of signal in the data (high vs. low) and the fusion method (early vs.
543 late). For each dataset, we used different random seeds to introduce variability while keeping all other
544 parameters consistent, particularly the latent factor structure \mathbf{U} . In all datasets, 10 interaction terms were
545 introduced within each view to capture the dependencies between selected features. These interactions
546 model the relationships between certain features within the same view.

547 The experimental conditions were as follows:

548 • **Signal Strength and Fusion Type**

- 549 – **High Signal (Early Fusion)**: Data were generated with a strong latent factor influence, including
 550 interaction effects between features. The latent factor structure \mathbf{U} and interaction terms remained
 551 consistent across both datasets.
- 552 – **High Signal (Late Fusion)**: The same strong latent factor influence was used as in the early
 553 fusion condition. However, the datasets were generated with different samples, while maintaining
 554 consistency in latent factor structure \mathbf{U} and interaction effects.
- 555 – **Low Signal (Early Fusion)**: Data were generated with a weaker latent factor influence. Inter-
 556 action effects were included, and the latent factor structure \mathbf{U} and interaction terms remained
 557 consistent across both datasets.
- 558 – **Low Signal (Late Fusion)**: Data were generated with weak latent factor influence. The two
 559 datasets were generated with different samples, while maintaining consistency in latent factor
 560 structure \mathbf{U} and interaction effects.
- 561 • **Outcome Type**
- 562 – **Binary Outcomes**: Interaction effects were included, and the latent factor structure \mathbf{U} and
 563 interaction terms remained consistent across datasets, based on the signal strength and fusion
 564 type.
- 565 – **Continuous Outcomes**: Similar to binary outcomes, but the outcome variables were generated
 566 using Gaussian noise as detailed in the main text.

567 **Parameters used to generate the datasets:**

568 The following are the parameters used to generate the datasets, based on the notations introduced in the
 569 previous sections.

- 570 • Sample size: $n = 1000$
- 571 • Number of features: $f_{\mathbf{X}_1} = f_{\mathbf{X}_2} = 100$
- 572 • Number of latent factors: $f_{\text{latent}} = 25$ (Including five for interacting features)
- 573 • Scaling factors: $\text{scale}_1 = \text{scale}_2 = 5$ for binary high signal and continuous high signal; $\text{scale}_1 =$
 574 $\text{scale}_2 = 10$ for binary low signal and continuous low signal

- 575 • Latent factor standard deviations: $\text{std} = 1$ (shared standard deviation for latent factors affecting both
576 views if correlated), $\text{std}_1 = 1$ (standard deviation for latent factors influencing \mathbf{X}_1 if uncorrelated),
577 $\text{std}_2 = 1$ (standard deviation for latent factors influencing \mathbf{X}_2 if uncorrelated)
- 578 • Factor strength: $\text{latent_strength} = 2$ for high signal, $\text{latent_strength} = 1$ for low signal
- 579 • Noise standard deviation: $\sigma = 5$ for strong signal conditions, $\sigma = 10$ for weak signal conditions
580 (binary outcomes), $\sigma = 15$ for weak signal conditions (continuous outcomes)
- 581 • Correlation between two datasets: TRUE if \mathbf{X}_1 and \mathbf{X}_2 are correlated, FALSE if they are uncorrelated
- 582 • Data type: Binary or continuous outcomes for \mathbf{y}_1 and \mathbf{y}_2
- 583 • Interaction effects: TRUE if interactions between features of \mathbf{X}_1 and \mathbf{X}_2 are included, with interaction
584 weight $w = 10$
 - 585 – Interaction effects between the 21st and 25th features within each view (\mathbf{X}_1 and \mathbf{X}_2) were
586 artificially introduced.

587 Where:

- 588 – \mathbf{X}_1 and \mathbf{X}_2 represent the feature matrices for the two datasets.
- 589 – \mathbf{U} is the shared latent factor matrix when the latent factors are correlated, and \mathbf{U}_1 and \mathbf{U}_2 are the
590 separate latent factor matrices for the uncorrelated case.
- 591 – β , $\beta_{\mathbf{U}_1}$, and $\beta_{\mathbf{U}_2}$ are the latent factor weights used to generate the expected outcome variables.
- 592 – latent_strength indicates the strength of the latent factors' influence on the outcome variables, with
593 higher values corresponding to stronger signal conditions.
- 594 – σ denotes the standard deviation of the noise added to the outcome variables.
- 595 – w is the interaction weight, a constant that controls the influence of interaction terms between features
596 from \mathbf{X}_1 and \mathbf{X}_2 , particularly for the 21st to 25th features in each view.

597 **Real Data**

598 **Transcriptomics-Metabolomics (ROSMAP)**

599 The *Religious Orders Study* and *Rush Memory and Aging Project* (ROSMAP) are ongoing longitudinal
600 clinical-pathologic cohort studies of aging and AD [78]. The ROSMAP data include a wealth of clinical,

601 cognitive, neuroimaging, genetic, and neuropathological information from older adults, providing invaluable
602 insights into the early stages of AD, cognitive decline, and related neurodegenerative processes.

603 **Transcriptomics** Single-cell RNA-seq (scRNA-seq) for astrocytes was downloaded from a published
604 study [79]. From 149,558 cells and 17,817 genes, we projected 4,292 metacells and 16,718 genes for
605 AD pathologic diagnosis using *Metacell-2* [80]. Then, we performed differential expression testing for
606 Alzheimer’s disease (AD) pathological diagnosis using *Seurat* and selected the top 280 differentially ex-
607 pressed genes (p -value adjusted < 0.1 and $\log_2(\text{fold change}) > 0.2$). Additionally, 25 genes that are
608 either AD-risk or astrocyte-specific (*APOE*[59, 60], *CLU*[81], *SORLI*[82], *TREM2*[82], *ABCA7*[82],
609 *PICALM*[82], *BINI*[63], *PLD3*[82], *CD33*[83], *CD2AP*[82], *EPHA1*[82], *INPP5D*[82], *FERMT2*[82],
610 *CRI*[82], *PTK2B*[82], *SLC24A4*[82], *CASS4*[82], *ACE*[83], *SORCSI*[82], *GAB2*[84], *CDK5RI*[85],
611 *PRNP*[86], *IL1RAP*[87], *CNP*[68], *TOMM40*[88]) were also selected.

612 **Metabolomics** ROSMAP metabolomic data processed by *Metabolon HD4* was downloaded from *AD*
613 *knowledge Portal* - backend (syn10235592). The dataset comprises a total of 1,055 biochemicals, 971
614 compounds of known identity (named biochemicals) and 84 compounds of unknown structural identity
615 (unnamed biochemicals) for 514 brain samples. We selected 305 biochemicals that do not have any
616 missingness in their samples.

617 **Multi-view Data** For early fusion, we merged transcriptomic and metabolomic data and achieved 2,286
618 metacells (samples) and 305 genes and 305 metabolites in each view. Metabolites were matched to the
619 donors for metacells in transcriptomic data. For late fusion, we used full samples for each view whose AD
620 pathologic diagnosis is known. 4,292 metacells and 2,286 samples for transcriptomics and metabolomics,
621 respectively were used.

622 **Transcriptomics-Epigenomics (SEA-AD)**

623 The *Seattle Alzheimer’s Disease Brain Cell Atlas* (SEA-AD) consortium studies deep molecular and cellular
624 understanding of the early pathogenesis of AD [89]. By integrating neuropathological data, single-cell and
625 spatial genomics, and longitudinal clinical metadata, SEA-AD serves as a unique resource for exploring the
626 mechanisms underlying Alzheimer’s disease and related dementias.

627 **Transcriptomics** scRNA-seq data for oligodendrocytes from middle temporal gyrus (111,194 cells) was
628 downloaded from the *Registry of Open Data on AWS* as AnnData objects (h5ad format). 3,927 metacells

629 and 17,368 genes were projected for cognitive status (dementia). We selected expressions for the 206
630 transcription factors that are known to co-bind or interact each other in oligodendrocytes[29].

631 **Epigenomics** Single-cell ATAC-seq (scATAC-seq) data for all nuclei from middle temporal gyrus was
632 downloaded from the *Registry of Open Data on AWS* as AnnData objects (h5ad format). There were 35,925
633 cells and 218,882 peaks for oligodendrocytes. Among those peaks, we chose 702 peaks that are overlapped
634 with the known oligodendrocyte-specific peaks[29].

635 **Multi-view Data** We merged the scRNA-seq and scATAC-seq datasets by cell. The 35,925 cells from the
636 scATAC-seq data were then grouped according to the metacells they belonged to in the scRNA-seq data,
637 and the average counts for each metacell were computed.

638 **scRNA-seq-Single-cell Spatial Transcriptomic (MERFISH) (SEA-AD)**

639 **scRNA-seq** scRNA-seq data for microglia was downloaded from *CellxGene* [90] for a published study
640 [89]. From 40,000 cells and 18,279 genes, we projected 1,009 metacells and 17,348 genes for AD
641 progression score using *Metacell-2*. Then, we implemented differential expression testing for AD pro-
642 gression score using the Differential Gene Expression analysis based on the negative binomial distribution
643 (*DESeq2*) [91] and selected 400 differentially expressed genes (p -value adjusted < 0.05 and $\log_2(\text{fold}$
644 $\text{change}) > 0.75$). Additionally, 28 genes that are either AD-risk or microglia-specific were identi-
645 fied. (*CSF1R*[92], *TREM2*[82], *SORLI*[82], *CD68*[92], *TNF*[93], *CLEC7A*[92], *IL6*[108], *APOE*[59, 60],
646 *CIQA*[92], *CIQB*[92], *CIQC*[92], *SPP1*[92], *MAPK8*[94], *PTGS2*[95], *VEGFA*[96], *FOS*[92], *CLU*[81],
647 *CR1*[82], *PICALM*[82], *CD33*[83], *SORCSI*[82], *GAB2*[84], *CASS4*[82], *CDK5R1*[92], *PRNPL*[86],
648 *IL1RAP*[92], *CNP*[68], *TOMM40*[88]) were also selected.

649 **Spatial transcriptomic (MERFISH)** Astrocyte single-cell MERFISH data for the whole taxonomy col-
650 lected from the middle temporal gyrus was obtained through the *Open Data Registry on AWS* as AnnData
651 objects (h5ad format) for a published study [89]. From 1,321,191 cells and 135 genes, 8,693 metacells
652 and 134 genes were projected for AD progression score using *Metacell-2*. We 1,016 metacells were for
653 astrocytes.

654 **Multi-view Data** Among the 1,016 metacells in the MERFISH data, we randomly selected 1,009 meta-
655 cells. The AD progression score was highly correlated between the scRNA-seq and MERFISH datasets (r

656 = 0.986). For each dataset, we computed the mean AD progression score, which was highly correlated with
 657 the corresponding scores in both datasets ($r = 0.997$ for scRNA-seq and $r = 0.996$ for MERFISH).

658 **Model Design**

659 Let $\mathbf{X}_1 \in \mathbb{R}^{n \times p_{x_1}}$ and $\mathbf{X}_2 \in \mathbb{R}^{n \times p_{x_2}}$ represent two data views, and $\mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^n$ be real-valued labels that
 660 represent the same disease phenotype corresponding to \mathbf{X}_1 and \mathbf{X}_2 , respectively. COSIME inputs each
 661 data view into a deep neural network encoder (deep encoder), performs dimensionality reduction to learn
 662 K -dimensional features, and creates embeddings in an unsupervised fashion. Then, the Learnable Optimal
 663 Transport (LOT) aligns the distributions of two latent spaces in a joint latent space. We use the joint latent
 664 vector to classify or predict disease phenotypes.

665 **Deep Neural Network Encoders (Deep encoders)**

666 COSIME is designed to process dual-view data using separate probabilistic deep encoders for each view.
 667 Each view $x_1 \in X_1$ and $x_2 \in X_2$ is passed through its own deep encoder to learn a latent representation.
 668 These encoders approximate the posterior distributions $p(z_1|x_1)$ and $p(z_2|x_2)$, where z_1 and z_2 are the latent
 669 variables corresponding to views 1 and 2, respectively. COSIME utilizes backpropagation of Learnable
 670 Optimal Transport (LOT) to these deep neural networks, enabling the learning of latent features from
 671 multiple views, which is then used for predicting disease phenotypes.

672 The deep encoders output the mean (μ_i) and log-standard deviation ($\log(\sigma_i)$) for the posterior distribution
 673 of each view. The encoder for each view consists of multiple fully connected layers with LeakyReLU
 674 activations. Batch Normalization is applied to stabilize learning and accelerate convergence, while Dropout
 675 is used as a regularization technique to mitigate overfitting. The final output of each deep encoder is a latent
 676 vector z_i that represents the compressed information from the corresponding input data of a view.

677 To regularize the learned posterior distributions and prevent overfitting, the model incorporates the
 678 Kullback-Leibler divergence (KLD) loss during the training of the deep encoders. This loss encourages the
 679 posterior distributions of each view to remain close to a standard Gaussian prior distribution $p(z_i) = \mathcal{N}(0, I)$,
 680 which helps improve the ability of the model to generalize and prevents overfitting. The KLD loss for each
 681 view is computed as:

$$682 \quad \mathcal{L}_{\text{KLD}_1} = D_{\text{KLD}_1}(\psi_1(z_1|x_1) \parallel p(z_1)) \quad (23)$$

$$\mathcal{L}_{\text{KLD}_2} = D_{\text{KLD}_2}(\psi_2(z_2|x_2) \parallel p(z_2)) \quad (24)$$

683 where ψ_1 and ψ_2 are the encoder functions for views A and B , respectively, D_{KLD_A} and D_{KLD_B} are the
 684 KLD measures for the latent representations of views A and B , respectively, and $p(z_1)$ and $p(z_2)$ are the
 685 prior distributions for the latent variables z_1 and z_2 corresponding to views A and B . The KLD loss for each
 686 view i is defined as:

$$D_{\text{KLD}_i}(\psi_i(z_i|x_i) \parallel p(z_i)) = \frac{1}{2} \left(\sigma_i^2 + \mu_i^2 - \log \sigma_i^2 - 1 \right) \quad (25)$$

687 where μ_i and σ_i are the mean and standard deviation of the posterior distribution for view i .

688 **Learnable Optimal Transport (LOT)**

689 COSIME utilizes Learnable Optimal Transport (LOT), a novel method designed to align the latent space
 690 distributions of two distinct data views. There are several key advantages of LOT.

691 First, unlike traditional optimal transport methods, Learnable Optimal Transport (LOT) treats the
 692 transport plan as a learnable parameter, optimized dynamically during training rather than relying on a
 693 fixed or pre-defined plan. This allows LOT to adapt more flexibly to complex data structures and align
 694 latent spaces in a joint latent space. By learning the transport plan iteratively, LOT can adjust its alignment
 695 strategy to the specific characteristics of the data, making it particularly effective for downstream tasks such
 696 as prediction and classification.

697 Second, LOT can scale effectively with large datasets. By incorporating mini-batch Sinkhorn[97] up-
 698 dates, LOT accelerates the optimization process, reduces memory usage, and ensures efficient computation.
 699 This makes LOT particularly well-suited for high-dimensional data from multiple sources, enabling it to
 700 handle large-scale tasks without sacrificing performance.

701 Third, LOT integrates domain-specific regularization through the coupling matrix, which allows the
 702 model to incorporate prior knowledge about the relationship between the views. This is particularly beneficial
 703 when working with heterogeneous data in the multi-view analysis, where domain-specific relationships can
 704 enhance the quality of the learned latent space. We apply entropy regularization to stabilize the optimization
 705 process, ensuring smooth convergence and better-converging embeddings, which are crucial for downstream
 706 tasks.

707 **Learnable Optimal Transport (LOT) Process** LOT treats the transport plan as a learnable parameter,
 708 which is iteratively updated through backpropagation. This enables LOT to adaptively align the latent
 709 spaces of different views in a shared representation space.

710 Both the source and target distributions in LOT are modeled as Gaussian distributions, each param-
 711 eterized by their means (μ) and standard deviations (σ). The transport plan is learned by minimizing
 712 the pairwise transport cost between these distributions. The cost reflects the distance between their latent
 713 vectors, and the transport plan is updated during training to reduce this cost, aligning the distributions in a
 714 shared latent space.

715 The process involves the following steps:

716 1. Step 1: Mini-Batching Processing

717 Before performing the computation, the source and target distributions are split into mini-batches to
 718 ensure efficient processing, especially for large datasets. Let the total number of source samples be
 719 N_{src} and the total number of target samples be N_{tgt} . In the mini-batching procedure, we split the
 720 source distribution $\{\mu_{\text{src}}, \sigma_{\text{src}}\}$ and the target distribution $\{\mu_{\text{tgt}}, \sigma_{\text{tgt}}\}$ into smaller batches of size B_{src}
 721 and B_{tgt} , respectively. Denote the i -th mini-batch of the source and target distributions as:

$$\mathcal{B}_{\text{src}}^i = \{\mu_{\text{src}}^i, \sigma_{\text{src}}^i\}, \quad \mathcal{B}_{\text{tgt}}^i = \{\mu_{\text{tgt}}^i, \sigma_{\text{tgt}}^i\} \quad (26)$$

722 for $i = 1, 2, \dots, N_{\text{src}}/B_{\text{src}}$ and $j = 1, 2, \dots, N_{\text{tgt}}/B_{\text{tgt}}$, respectively. This mini-batch processing
 723 reduces memory consumption and allows the transport plan to be optimized iteratively.

724 Each mini-batch is processed separately, and the transport plan is updated incrementally based on
 725 these smaller subsets of the source and target distributions.

726 2. Step 2: Pairwise Distance Computation

727 The pairwise distances between the means μ_{src}^i and μ_{tgt}^j and the standard deviations σ_{src}^i and σ_{tgt}^j of
 728 the mini-batches are computed using the Euclidean distance:

$$\text{dist}_{\mu}(i, j) = \|\mu_{\text{src}}^i - \mu_{\text{tgt}}^j\|_2, \quad \text{dist}_{\sigma}(i, j) = \|\sigma_{\text{src}}^i - \sigma_{\text{tgt}}^j\|_2 \quad (27)$$

729 where i and j refer to the samples from the source and target mini-batches, respectively.

730 These pairwise distances are used to construct the cost matrix for each mini-batch, which will then
 731 be used in the subsequent steps.

732 3. Step 3: Cost Matrix Construction

733 After computing the pairwise distances, the total cost matrix for each mini-batch is constructed as
 734 follows:

$$\text{cost}_{total}^i(i, j) = \text{dist}_\mu(i, j) + \text{dist}_\sigma(i, j) + \epsilon \quad (28)$$

735 where ϵ is a small constant for numerical stability.

736 4. Step 4: Sinkhorn Algorithm (Learnable Transport Plan)

737 The transport matrix \mathbf{P}^i for each mini-batch is computed iteratively using the Sinkhorn algorithm.
 738 The transport plan is initialized as a learnable parameter, \mathbf{P}^i , where each element of \mathbf{P}^i represents the
 739 transport probability between a sample from the source and a sample from the target. The transport
 740 plan is updated during each iteration to minimize the total transportation cost.

741 The update rule for the transport plan $\mathbf{P}^i(i, j)$ for each mini-batch is:

$$\mathbf{P}^i(i, j) = \frac{\exp\left(-\frac{\text{cost}_{total}^i(i, j)}{\text{reg}}\right)}{Z_{\text{src}}^i \cdot Z_{\text{tgt}}^i} \quad (29)$$

742 where Z_{src}^i and Z_{tgt}^i are the normalizing factors for the source and target distributions, respectively:

$$Z_{\text{src}}^i = \sum_j \exp\left(-\frac{\text{cost}_{total}^i(i, j)}{\text{reg}}\right), \quad Z_{\text{tgt}}^i = \sum_i \exp\left(-\frac{\text{cost}_{total}^i(i, j)}{\text{reg}}\right) \quad (30)$$

743 Here, the learnable transport plan allows gradients to flow through the transport matrix, enabling the
 744 optimization process using backpropagation.

745 **Handling Different Sample Sizes** The Sinkhorn algorithm can effectively handle source and target
 746 distributions with differing sample sizes, making it particularly useful for COSIME late fusion. To
 747 ensure that transport plans remain valid, we apply normalization to the weights of the source and
 748 target distributions during each mini-batch update as follows:

$$\mathbf{A}(i) = \frac{\mathbf{A}(i)}{\sum_i \mathbf{A}(i)}, \quad \mathbf{B}(j) = \frac{\mathbf{B}(j)}{\sum_j \mathbf{B}(j)} \quad (31)$$

749 This adjustment ensures that the transport plan remains meaningful even when the source and target
 750 distributions contain different numbers of samples.

751 The Sinkhorn algorithm iteratively refines the transport plan \mathbf{P}^i until convergence or until a predefined
752 number of iterations.

753 5. Step 5: Dual Variables and Marginal Constraints

754 Dual variables \mathbf{u}^i and \mathbf{v}^i are introduced to enforce the marginal constraints of the transport problem.
755 These variables ensure that the transport plan satisfies the required marginal distributions for the
756 source and target distributions. The dual variables are updated iteratively for each mini-batch:

$$\mathbf{u}^i(i) = \frac{\mathbf{A}(i)}{Z_{\text{src}}^i} \cdot \exp\left(-\frac{\text{cost}_{\text{total}}^i(i, j)}{\text{reg}}\right) \quad (32)$$

$$\mathbf{v}^i(j) = \frac{\mathbf{B}(j)}{Z_{\text{tgt}}^i} \cdot \exp\left(-\frac{\text{cost}_{\text{total}}^i(i, j)}{\text{reg}}\right) \quad (33)$$

757 where the dual variables ensure that the transport plan matches the marginal distributions of the source
758 and target distributions.

759 6. Step 6: Optimizing the Transport Plan

760 During each mini-batch, the transport plan is optimized by performing backpropagation to minimize
761 the total loss. This involves computing the gradients of the transport plan with respect to the loss
762 function and updating the transport plan parameters via gradient descent.

763 7. Step 7: Final LOT Loss Calculation

764 After optimizing the transport plan, the final LOT loss is computed as the sum of the product of the
765 cost matrix and the transport plan:

$$\mathcal{L}_{\text{LOT}} = \sum_{i,j} \text{cost}_{\text{total}}(i, j) \cdot \mathbf{P}(i, j) \quad (34)$$

766 This loss function quantifies the minimum cost required to transform the source distribution into the
767 target distribution, thus aligning the two views in a joint latent space.

768 Phenotype Prediction/Classification

769 COSIME employs the aligned latent representations from each view, obtained via separate probabilistic deep
770 encoders. After the deep encoders generate the latent vectors for each view, LOT aligns these embeddings
771 into a joint latent space. The model then uses these aligned embeddings for phenotype prediction.

772 Two fusion strategies are employed to combine the embeddings from the two views: early fusion and
 773 late fusion. Both fusion strategies are applied for binary outcome classification and continuous outcome
 774 prediction tasks.

775 **Early Fusion** The fused representation is obtained by taking the mean of the representations $z_1 \in \mathbb{R}^d$ and
 776 $z_2 \in \mathbb{R}^d$ from each view:

$$z_{\text{fusion}} = \frac{z_1 + z_2}{2} \quad (35)$$

777 **Late Fusion** The latent representations $z_1 \in \mathbb{R}^d$ and $z_2 \in \mathbb{R}^d$ from each view are vertically concatenated
 778 to form the fused representation z_{fusion} :

$$z_{\text{fusion}} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad (36)$$

779 where $z_{\text{fusion}} \in \mathbb{R}^{2d}$ is the fused vector obtained by stacking z_1 and z_2 vertically.

780 **Binary Outcome Prediction** The fused representation z_{fusion} is then passed through a linear layer followed
 781 by a sigmoid activation function to compute the predicted probability \hat{p} :

$$\hat{p} = \sigma(W_{\text{binary}} \cdot z_{\text{fusion}} + b) \quad (37)$$

782 where $\sigma(\cdot)$ is the sigmoid function, and W_{binary} and b are the learned weights and bias parameters of
 783 the prediction layer.

784

785 **Loss function**

786 The **BCEWithLogitsLoss** is computed for the early fusion setup, which combines the sigmoid activation
 787 and binary cross-entropy in a single, more numerically stable operation. This loss function is defined as:

$$\mathcal{L}_{\text{BCE}}^{\text{early}} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\sigma(\hat{p}_i)) + (1 - y_i) \log(1 - \sigma(\hat{p}_i))] \quad (38)$$

788 where \hat{p}_i is the predicted probability of the i -th sample belonging to the positive class (i.e., the probability
 789 output by the model), and y_i is the true binary label for the i -th sample, where $y_i \in \{0, 1\}$.

790 **Continuous Phenotype Prediction** For continuous phenotype prediction, the fused representation z_{fusion}
 791 is passed through a linear transformation:

$$\hat{y} = W_{\text{reg}} \cdot z_{\text{fusion}} + b \quad (39)$$

792 where W_{reg} and b are the learned weights and bias of the regression layer.

793
 794 **Loss for Regression** The Mean Squared Error (MSE) loss for continuous outcome prediction using a
 795 linear model is computed as:

$$\mathcal{L}_{\text{MSE}}^{\text{early}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (40)$$

796 where \hat{y}_i is the predicted continuous value for the i -th sample, representing the model's prediction,
 797 and y_i is the true continuous value for the i -th sample, which is the actual target value the model aims to
 798 predict. The loss measures the squared difference between the predicted and true values for each sample
 799 and averages this over all N samples.

800
 801 **Loss for Neural Network** A neural network (NN) model can be used for continuous outcome prediction,
 802 the output \hat{y}_i is predicted by the neural network, and the loss function remains MSE:

$$\mathcal{L}_{\text{MSE}}^{\text{NN, early}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i^{\text{NN}})^2 \quad (41)$$

803 where \hat{y}_i^{NN} is the continuous value predicted by the neural network for the i -th sample, which is the
 804 output of the neural network model for that particular input, and y_i is the true continuous value for the
 805 i -th sample, which is the actual target value that the model is attempting to predict. The loss measures the
 806 squared difference between the predicted value and the true value for each sample, and the average squared
 807 error is computed over all N samples.

808 **Prediction Loss:** The prediction loss term $\mathcal{L}_{\text{pred}}$ is defined based on the task at hand:

$$\mathcal{L}_{\text{pred}} = \begin{cases} \mathcal{L}_{\text{BCE}}^{\text{early}}, & \text{for binary outcome classification (early fusion)} \\ \mathcal{L}_{\text{MSE}}^{\text{early}}, & \text{for continuous outcome prediction (early fusion, regression)} \\ \mathcal{L}_{\text{MSE}}^{\text{NN, early}}, & \text{for continuous outcome prediction (early fusion, NN)} \\ \mathcal{L}_{\text{BCE}}^{\text{late}}, & \text{for binary outcome classification (late fusion)} \\ \mathcal{L}_{\text{MSE}}^{\text{late}}, & \text{for continuous outcome prediction (late fusion, regression)} \\ \mathcal{L}_{\text{MSE}}^{\text{NN, late}}, & \text{for continuous outcome prediction (late fusion, NN)} \end{cases}$$

809 Total Weighted Loss

810 During the training phase, the contributions of each loss term can vary in magnitude, depending on the
 811 relative scale of the individual losses. To balance the influence of each term on the total loss, we normalize
 812 the loss terms dynamically. This dynamic normalization ensures that no individual loss term has an outsized
 813 influence on the total loss, allowing for a more balanced optimization process.

814 We first compute the inverse of each loss term:

$$\text{Inv}_{\mathcal{L}_{\text{KLD}_A}} = \frac{1}{\mathcal{L}_{\text{KLD}_A}(\theta) + 1e^{-8}} \quad (42)$$

$$\text{Inv}_{\mathcal{L}_{\text{KLD}_B}} = \frac{1}{\mathcal{L}_{\text{KLD}_B}(\theta) + 1e^{-8}} \quad (43)$$

$$\text{Inv}_{\mathcal{L}_{\text{LOT}}} = \frac{1}{\mathcal{L}_{\text{LOT}}(\theta) + 1e^{-8}} \quad (44)$$

$$\text{Inv}_{\mathcal{L}_{\text{pred}}} = \frac{1}{\mathcal{L}_{\text{pred}}(\theta) + 1e^{-8}} \quad (45)$$

818 Next, we compute the total inverse magnitude:

$$\text{Total_Inverse} = \text{Inv}_{\mathcal{L}_{\text{KLD}_A}} + \text{Inv}_{\mathcal{L}_{\text{KLD}_B}} + \text{Inv}_{\mathcal{L}_{\text{LOT}}} + \text{Inv}_{\mathcal{L}_{\text{pred}}} \quad (46)$$

819 Now, we normalize the inverse magnitudes to get their relative adjusted weights:

$$\text{Norm}_{\mathcal{L}_{\text{KLD}_A}} = \frac{\text{Inv}_{\mathcal{L}_{\text{KLD}_A}}}{\text{Total_Inverse}} \quad (47)$$

$$\text{Norm}_{\mathcal{L}_{\text{KLD}_B}} = \frac{\text{Inv}_{\mathcal{L}_{\text{KLD}_B}}}{\text{Total_Inverse}} \quad (48)$$

$$\text{Norm}_{\mathcal{L}_{\text{LOT}}} = \frac{\text{Inv}_{\mathcal{L}_{\text{LOT}}}}{\text{Total_Inverse}} \quad (49)$$

$$\text{Norm}_{\mathcal{L}_{\text{pred}}} = \frac{\text{Inv}_{\mathcal{L}_{\text{pred}}}}{\text{Total_Inverse}} \quad (50)$$

822 Finally, the total weighted loss $\mathcal{L}_{\text{Total_weighted}}(\theta)$ is computed as:

$$\begin{aligned} \mathcal{L}_{\text{Total_weighted}}(\theta) = & KLD_{A_{\text{weight}}} \cdot \text{Norm}_{\mathcal{L}_{\text{KLD}_A}} \cdot \mathcal{L}_{\text{KLD}_A}(\theta) + \\ & KLD_{B_{\text{weight}}} \cdot \text{Norm}_{\mathcal{L}_{\text{KLD}_B}} \cdot \mathcal{L}_{\text{KLD}_B}(\theta) + \\ & OT_{\text{weight}} \cdot \text{Norm}_{\mathcal{L}_{\text{LOT}}} \cdot \mathcal{L}_{\text{LOT}}(\theta) + \\ & CL_{\text{weight}} \cdot \text{Norm}_{\mathcal{L}_{\text{pred}}} \cdot \mathcal{L}_{\text{pred}}(\theta) \end{aligned} \quad (51)$$

823 where λ_{KLD_A} , λ_{KLD_B} , λ_{LOT} , and λ_{pred} are the weights that control the importance of each normalized
824 loss term, specifically $\text{Norm}_{\mathcal{L}_{\text{KLD}_A}}$, $\text{Norm}_{\mathcal{L}_{\text{KLD}_B}}$, $\text{Norm}_{\mathcal{L}_{\text{LOT}}}$, and $\text{Norm}_{\mathcal{L}_{\text{pred}}}$, respectively. These normalized
825 loss terms are the contributions of each individual loss term $\mathcal{L}_{\text{KLD}_A}(\theta)$, $\mathcal{L}_{\text{KLD}_B}(\theta)$, $\mathcal{L}_{\text{LOT}}(\theta)$, and $\mathcal{L}_{\text{pred}}(\theta)$
826 to the total loss, which are normalized to ensure that no single loss dominates the optimization process.

827

828 To optimize the model parameters, we minimize the total weighted loss:

$$\theta^* = \underset{\theta}{\text{argmin}} \mathcal{L}_{\text{Total_weighted}}(\theta) \quad (52)$$

829 This method prevents overemphasis on any one loss component, allowing the model to effectively
830 balance view alignments and prediction accuracy. It enables flexible adjustment of task importance, which
831 is crucial when working with multi-view data, where different objectives (e.g., alignment, prediction) may
832 compete for optimization.

833 **Model training process**

834 The input data for each view was split into 75% for training and 25% for testing (holdout). 5-fold cross-
835 validation was applied to the training set. The best-performing model was selected based on the minimum
836 classification (or prediction) loss during the training phase and was then evaluated on the testing set. The
837 **Algorithm 1: COSIME Model Training** is detailed in the **Supplementary Information**.

838 The prediction performance of the models was evaluated using holdout data with metrics such as Area
839 Under the Receiver Operating Characteristic Curve (AUROC), Area Under the Precision-Recall Curve
840 (AUPRC), and accuracy for binary outcome classification, and Mean Squared Error (MSE) for continuous

841 outcome prediction. The performance metrics were compared across different methods using the mean and
 842 ± 1.96 times the standard deviation. For the best-trained models across the different multi-view datasets, we
 843 computed feature importance for each data view and interaction values both within and across views.

844 **Hyperparameter tuning**

845 We tuned several hyperparameters such as the learning rate, learning rate decay (gamma), loss weights,
 846 dropout rate, joint latent space dimensionality, early stopping patience, minimum change for early stopping,
 847 and weight decay in the optimizer. Hyperparameter optimization was performed using *Ray Tune* [98]
 848 with Distributed Asynchronous Hyperparameter Optimization (*Hyperopt*) [99]. The Adam optimizer was
 849 employed to minimize the total loss, with a maximum of 300 epochs. Early stopping was used to prevent
 850 overfitting, based on performance improvements on the validation set in the training phase.

851 **Optimization and Loss Minimization**

852 During the training phase, the model updates the hyper-parameters by minimizing the total loss. The
 853 optimization process involves computing gradients with respect to the total loss function and updating the
 854 parameters using gradient descent or its variants.

855 By minimizing the total loss, COSIME simultaneously learns to:

- 856 • Regularize the posterior distributions of the latent representations for each view through the KLD
 857 losses,
- 858 • Align the embeddings from both views into a shared latent space using LOT,
- 859 • Predict or classify the disease phenotype by minimizing classification (or prediction) loss.

860 **Evaluating the Performance of the Model Predictions**

861 To evaluate the performance of the best model selected in the training phase for both binary outcome
 862 classification and continuous outcome prediction, several standard metrics are employed:

863 **Binary Outcome Classification**

- 864 1. **Area Under the Receiver Operating Characteristic Curve (AUROC):** AUROC measures the
 865 ability of the model to discriminate between the positive and negative classes. It is calculated by
 866 assessing the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold

867 settings. The AUROC value ranges from 0 to 1, where 1 indicates perfect classification, and 0.5 is
 868 the expected performance of a random classifier. The formula for AUROC is:

$$\text{AUROC} = \int_0^1 \text{TPR}(f) d\text{FPR}(f) \quad (53)$$

869 where $\text{TPR}(f)$ and $\text{FPR}(f)$ are the True Positive Rate and False Positive Rate at each decision
 870 threshold f .

871 **2. Area Under the Precision-Recall Curve (AUPRC):** AUPRC is another useful metric for evaluating
 872 binary classification performance, particularly when dealing with imbalanced datasets. It measures
 873 the area under the curve created by plotting precision against recall at various threshold settings. The
 874 formula for AUPRC is:

$$\text{AUPRC} = \int_0^1 \text{Precision}(r) d\text{Recall}(r) \quad (54)$$

875 where $\text{Precision}(r)$ and $\text{Recall}(r)$ are the precision and recall at each decision threshold r .

876 **3. Accuracy:** Accuracy is the proportion of correctly classified instances (both true positives and true
 877 negatives) out of the total number of instances. It is calculated as:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (55)$$

878 where TP is True Positives, TN is True Negatives, FP is False Positives, and FN is False Negatives.

879 **Continuous Outcome Prediction**

880

881 **Mean Squared Error (MSE):** MSE measures the average squared difference between the true values
 882 and the predicted values. It is a common loss function for regression tasks and is used to quantify the error
 883 in continuous prediction tasks. The formula for MSE is:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (56)$$

884 where y_i is the true continuous label for the i -th sample, \hat{y}_i is the predicted value for the i -th sample,
 885 and n is the total number of samples.

886

887 These evaluation metrics assess the effectiveness of the COSIME model in both binary and continuous
888 phenotype prediction tasks. Higher AUROC and AUPRC values, along with higher accuracy and lower
889 MSE, indicate better model performance.

890 **Feature Importance and Interaction**

891 COSIME introduces an advanced approach for computing both feature importance and feature interactions,
892 building upon the Shapley-Taylor framework[9]. By using Monte Carlo sampling and a batching process,
893 COSIME efficiently approximates Shapley values and Shapley-Taylor indices, making it highly scalable
894 and interpretable for large, high-dimensional datasets while maintaining memory efficiency. This approach
895 ensures COSIME is well-suited for complex real-world applications.

896 COSIME employs Monte Carlo sampling to efficiently approximate both Shapley values and feature
897 interactions. This approach significantly enhances its scalability, interpretability, and flexibility, especially
898 when handling large datasets and high-dimensional feature spaces. Instead of relying on exact calculations,
899 which become computationally expensive as the number of features or data points increases, Monte Carlo
900 sampling allows COSIME to estimate Shapley values and Shapley-Taylor indices efficiently. This method
901 also offers precise and interpretable calculations of pairwise feature interactions by directly simulating fea-
902 ture combinations. By estimating these interactions in this way, COSIME avoids complex approximations or
903 indirect assumptions, providing clearer insights into how combinations of features jointly affect predictions.
904 Furthermore, COSIME is better equipped to handle complex feature interactions and nonlinear relationships
905 more effectively.

906 Additionally, COSIME is memory-efficient due to its batching process, which breaks the data into
907 smaller subsets to minimize memory usage. This enables COSIME to handle large datasets on systems
908 with limited memory, making it suitable for real-world applications with resource constraints. The batching
909 process ensures that COSIME can scale efficiently while maintaining performance across high-dimensional
910 feature spaces.

911 COSIME also provides the ability to calculate signed feature interactions, which indicate whether pairs
912 of features interact synergistically (positive interaction) or antagonistically (negative interaction). This is
913 achieved by directly simulating feature combinations during the Monte Carlo sampling process, allowing
914 COSIME to estimate both the magnitude and direction of interactions. The sign of an interaction is
915 determined by comparing the joint effect of the features to their individual effects, revealing whether their

916 combined influence on the prediction. This enhancement adds a significant layer of interpretability, enabling
 917 users to understand not only the strength but also the directionality of feature interactions.

918 Finally, COSIME allows users to customize the Monte Carlo sample size and maximum memory
 919 usage, providing flexibility in the computational process. When the memory usage is specified, COSIME
 920 automatically determines the optimal batch size, making the method both user-friendly and adaptable
 921 to different system configurations. The details of **Algorithm 2: COSIME Feature Importance and**
 922 **Interaction Computation** are provided in the **Supplementary Information**.

923 Shapley Value (Feature Importance Computation)

924 The Shapley value ϕ_i for a given feature i is computed by considering the marginal contribution of feature i
 925 across all possible combinations of other features. This is approximated using Monte Carlo sampling over
 926 multiple iterations. The Shapley value for each feature i and each sample s is computed as:

$$\phi_p(s) = \frac{1}{M} \sum_{m=1}^M \left(f(\mathbf{x}_s^{(p-1)}) - f(\mathbf{x}_s^{(p)}) \right) \quad (57)$$

927 where:

- 928 • $f(\mathbf{x}_s)$ represents the prediction using all features,
- 929 • $f(\mathbf{x}_s^{(p)})$ represents the prediction when feature p is masked,
- 930 • $f(\mathbf{x}_s^{(p-1)})$ represents the prediction when all features except for feature p are masked.

931 The Shapley values for all features are stored in a matrix $\mathcal{S} \in \mathbb{R}^{N \times F}$, where N is the number of samples
 932 and F is the number of features.

$$\mathcal{S}_{ij} = \phi_j(s_i) \quad (58)$$

933 where i indexes the i -th sample ($i \in \{1, 2, \dots, N\}$) and j indexes the j -th feature ($j \in \{1, 2, \dots, F\}$),
 934 and $\phi_j(s_i)$ represents the Shapley value for feature j with respect to sample i .

935 Feature Interaction Effects

936 The interaction effect between two features p and q quantifies the joint contribution of features p and q
 937 compared to when they are individually absent. The interaction effect \mathcal{I}_{pq} is computed as:

$$\mathcal{I}_{pq}(s) = f(\mathbf{x}_s) - f(\mathbf{x}_s^{(p)}) - f(\mathbf{x}_s^{(q)}) + f(\mathbf{x}_s^{(pq)}) \quad (59)$$

938 where:

- 939 • $f(\mathbf{x}_s)$: the prediction of the model using all features,
- 940 • $f(\mathbf{x}_s^{(p)})$ and $f(\mathbf{x}_s^{(q)})$: the predictions of the model with features p and q individually masked,
- 941 • $f(\mathbf{x}_s^{(pq)})$: the prediction of the model when both features p and q are masked.

942 For pairwise interactions, the interaction effect is averaged over M Monte Carlo iterations:

$$\mathcal{I}_{ij}(s) = \frac{1}{M} \sum_{m=1}^M \left[f(\mathbf{x}_s) - f(\mathbf{x}_s^{(i)}) - f(\mathbf{x}_s^{(j)}) + f(\mathbf{x}_s^{(ij)}) \right] \quad (60)$$

943 The interaction effects for all pairs of features are stored in the interaction matrix $\mathcal{I} \in \mathbb{R}^{F \times F}$, where each
944 entry \mathcal{I}_{ij} represents the interaction between features i and j .

$$\mathcal{I}_{ij} = \mathcal{I}_{ji} \quad (61)$$

945 Monte Carlo Approach and Interaction Algorithm

946 Given the model f , input data $X \in \mathbb{R}^{N \times F}$ (where N is the number of samples and F is the number of
947 features), and target labels $y \in \mathbb{R}^N$, we compute the Shapley values and interaction effects as follows:

- 948 • **Shapley Values:** For each feature i , compute the marginal contribution to each sample over M Monte
949 Carlo iterations. Store the results in the matrix \mathcal{S} .
- 950 • **Interaction Effects:** For each pair of features (i, j) , compute the interaction effect using the formula
951 above, averaging over M Monte Carlo iterations. Store the results in the interaction matrix \mathcal{I} .

952 The algorithm processes the data in batches to manage memory efficiently. If the batch size is not
953 provided, it is dynamically computed based on the available memory. The number of batches B is given by:

$$B = \left\lceil \frac{N}{\text{batch size}} \right\rceil \quad (62)$$

954 The batch processing ensures that we do not exceed the maximum memory usage during computation.
955 Considering the available resources and estimated computing time, we used 100 Monte Carlo iterations for
956 the simulated multi-view data and 10 iterations for the real multi-view data. For each multi-view dataset,
957 the test (holdout) data was used.

958 **Memory Management**

959 To ensure that the memory usage remains within a specified limit, we compute the memory required for
960 each batch as:

$$\text{Batch memory (GB)} = \frac{\text{batch size} \times \text{element size of } X \times F}{1 \times 10^9} \quad (63)$$

961 If the calculated batch memory exceeds the available memory, the batch size is adjusted accordingly.
962 The dynamic batch size calculation ensures that the computation does not exceed the memory limit.

963 **Output**

964 The Shapley values are stored in a matrix \mathcal{S} with dimensions $N \times F$, and the interaction effects are stored
965 in a matrix \mathcal{I} with dimensions $F \times F$. These matrices represent the individual contributions of each feature
966 and their pairwise interactions in the prediction process.

967 In summary, our new approach provides a comprehensive methodology for assessing both individual
968 feature importance (Shapley values) and pairwise feature interactions in machine learning models. By
969 leveraging these tools, we enable a deeper understanding of how features individually and collectively
970 influence model predictions. The computational framework we present, along with the accompanying code,
971 facilitates the application of this methodology to various prediction tasks, enhancing model interpretability
972 and providing valuable insights into feature relationships.

973 **Downstream analyses**

974 **Gene Enrichment Analysis**

975 The function `enrichGO` from the `clusterProfiler` [100] R package, a universal enrichment tool for
976 interpreting omics data, was used to perform enrichment analyses for the 20 prioritized genes by absolute
977 feature importance (FI) values in each view of real multi-view datasets. The top 10 categories were displayed
978 by $\log_{10}(\text{FDR})$ (**Fig. 3b** and **Supplementary Figure 1**).

979 **Metabolite Enrichment Analysis**

980 `MetaboAnalyst 6.0` [101] was used to implement a metabolite set enrichment analysis for 20 prioritized
981 metabolites by absolute feature FI values in the metabolomic data for classifying cognitive diagnosis for
982 Alzheimer's disease (**Fig. 3c**).

983 **Interaction Network**

984 For the top 50 across-view interaction values, the edges and nodes for the interacting genes (and genes and
985 metabolites) were formatted in R, and the network was generated using Cytoscape 3.10.3 [102] (Fig. 3e
986 and Fig. 5d).

987 **Spatial mapping and statistical tests**

988 Spatial expression plots for *RYS3* and *DGKG* were generated in Python, employing matplotlib and
989 seaborn (Fig. 5c). To assess whether there are significant spatial differences between cells with high FI
990 values and low feature importance values for each gene, the following steps were taken:

- 991 • Compute the midpoint of the bounding box: It is calculated by finding the midpoint of the bounding
992 box defined by the x and y coordinates. Specifically, the center is determined by averaging the
993 maximum and minimum values of both the x and y coordinates:

$$\text{center}_x = \frac{\max(x) + \min(x)}{2}, \quad \text{center}_y = \frac{\max(y) + \min(y)}{2} \quad (64)$$

994 This central point serves as the reference from which the distances of each point are measured.

- 995 • Compute the Euclidean distance from each point to the midpoint of the bounding box: Once the
996 midpoint of the bounding box is defined, the Euclidean distance from each point (x_i, y_i) to the
997 computed midpoint of the bounding box $(\text{center}_x, \text{center}_y)$ is calculated using the formula:

$$\text{distance}_i = \sqrt{(x_i - \text{center}_x)^2 + (y_i - \text{center}_y)^2} \quad (65)$$

998 This step computes a scalar value for each point that quantifies how far each point is from the center.

- 999 • Perform a two-sample two-tailed t-test: The dataset is divided into two groups based on the feature
1000 importance:

- 1001 – Group 1: Points with FI values greater than or equal to the 75th percentile (i.e., the top 25%).
- 1002 – Group 2: Points with FI values below the 75th percentile.

1003 We performed a two-sample two-tailed t-test to evaluate whether the distances from the midpoint of
1004 the bounding box differed between the top 25% of cells with the highest feature importance values

1005 and the remaining 75% of cells. The null hypothesis (H_0) for the two-tailed t-test is:

$$H_0 : \mu_{\text{Group1}_{\text{FI}}} = \mu_{\text{Group2}_{\text{FI}}}$$

1006 where $\mu_{\text{Group1}_{\text{FI}}}$ and $\mu_{\text{Group2}_{\text{FI}}}$ represent the means of the distances for the high and low FI groups,
1007 respectively. The alternative hypothesis is:

$$H_1 : \mu_{\text{Group1}_{\text{FI}}} \neq \mu_{\text{Group2}_{\text{FI}}}$$

1008 The test statistic t and the associated p -value are calculated to determine whether the difference in means
1009 is statistically significant.

1010

1011 Similarly, the spatial distribution of across-view interaction values was visualized using the same tools
1012 (**Fig. 5e**). The following steps were taken to generate the spatial distribution of across-view interaction values
1013 and assess whether there are significant spatial differences between cells with positive feature interaction
1014 values and negative feature interaction values, the subsequent steps were carried out:

- 1015 • Distance calculation: The distance of each cell from the center of the tissue sample was calculated
1016 using its spatial coordinates (X, Y) . This distance quantifies how far each cell is from the central point
1017 of the tissue section. The Euclidean distance from the center is computed as:

$$\text{Distance to center} = \sqrt{(X - X_{\text{center}})^2 + (Y - Y_{\text{center}})^2} \quad (66)$$

1018 where X_{center} and Y_{center} are the coordinates of the center of the tissue.

- 1019 • Cell-level interaction score: For each cell, the interaction score was computed using the MERFISH
1020 gene expression levels and pairwise across-view interaction values. Let $\mathbf{E}_{\text{MERFISH}}^{(i)}$ denote the gene
1021 expression vector for cell i , which is of length G_{MERFISH} . The interaction matrix \mathbf{M} has dimen-
1022 sions $G_{\text{MERFISH}} \times G_{\text{scRNA-seq}}$, where each element M_{g_1, g_2} represents the pairwise interaction between
1023 MERFISH gene g_1 and scRNA-seq gene g_2 . The cell-level interaction score $S^{(i)}$ for each cell i was
1024 calculated as the dot product of the gene expression vector for that cell and the interaction matrix:

$$S^{(i)} = \mathbf{E}_{\text{MERFISH}}^{(i)} \cdot \mathbf{M} \quad (67)$$

1025 In expanded form, the score for cell i with respect to scRNA-seq gene j is:

$$S_j^{(i)} = \sum_{g_1=1}^{G_{\text{MERFISH}}} \mathbf{E}_{\text{MERFISH}}^{(i)}[g_1] \cdot M_{g_1,j} \quad (68)$$

1026 This score quantifies the interaction between the MERFISH gene expression profile of cell i and
 1027 the scRNA-seq gene expression profile, enabling grouping of cells based on their interaction scores.
 1028 Cells were then grouped based on the sign of their interaction scores: those with positive scores
 1029 were assigned to the Positive group, while those with negative scores were assigned to the Negative
 1030 group. This grouping allowed for further comparison of spatial distributions and other characteristics
 1031 between the two groups.

1032 • Perform a two-sample two-tailed t -test: To determine whether there were significant differences in
 1033 the distances of cells from the center between the two groups, an independent two-sample t -test was
 1034 performed. The hypotheses tested were:

1035 – Null hypothesis (H_0): The mean distance from the center of the Positive and Negative groups
 1036 are equal, i.e.,

$$H_0 : \mu_{\text{Positive}} = \mu_{\text{Negative}}$$

1037 – Alternative hypothesis (H_1): The mean distance from the center of the Positive and Negative
 1038 groups are not equal, i.e.,

$$H_1 : \mu_{\text{Positive}} \neq \mu_{\text{Negative}}$$

1039 • Perform a two-tailed F -test: In addition to comparing the means, an F -test was performed to compare
 1040 the variances of the two groups. The hypotheses tested were:

1041 – Null hypothesis (H_0): The variances of the Positive and Negative groups are equal, i.e.,

$$H_0 : \sigma_{\text{Positive}}^2 = \sigma_{\text{Negative}}^2$$

1042 – Alternative hypothesis (H_1): The variances of the Positive and Negative groups are not equal,
 1043 i.e.,

$$H_1 : \sigma_{\text{Positive}}^2 \neq \sigma_{\text{Negative}}^2$$

1044 **Benchmarking methods**

1045 **Cooperative Learning**

1046 Cooperative learning is a framework for fitting machine learning models to multi-view data[3]. It presents
1047 a loss function that minimizes the prediction error across all views, while encouraging the predictions from
1048 the views to agree. This framework with linear regression, which is what we use for benchmarking.

1049 We implemented the Cooperative Learning framework using the `multiview` R package. Following the
1050 examples provided on the Cooperative Learning GitHub repository, we preprocessed the input data using
1051 the `preProcess` function to center and scale it. For model evaluation, we conducted 5-fold cross-validation
1052 using the `cv.multiview` function, with a maximum of 300 iterations. For binary outcome classification,
1053 we set `family = binomial()` and `type.measure = "class"`. For continuous outcome prediction, we
1054 used the default `family = gaussian()` and `type.measure = "mse"`. The hyperparameter values for
1055 the regularization parameter were chosen from the set $\{0, 0.2, 0.4, 0.6, 0.8, 1\}$, while the lambda values were
1056 left at their default settings.

1057 **Data Integration Analysis for Biomarker discovery using Latent components (DIABLO)**

1058 Data Integration Analysis for Biomarker discovery using Latent components (DIABLO) is a supervised
1059 extension of canonical correlation analysis[4]. It aims to identify linear combinations of variables from
1060 each dataset that maximize the correlation between datasets.

1061 We implemented DIABLO using the `mixOmics` R package. Following the tutorial on the mixOmics
1062 website, we created a data block using the `matrix` function. For model training, we performed 5-fold
1063 cross-validation using the `block.splsda` function. The number of variables was tested within the set
1064 $\{10, 20, 30, 40, 50\}$.

1065 **Multi-Omics Factor Analysis v2 (MOFA+)**

1066 Multi-Omics Factor Analysis v2 (MOFA+) is an unsupervised matrix decomposition method that can be
1067 viewed as a generalization of sparse principal components analysis[6].

1068 We used the `MOFA2` R package to obtain latent representations. Following the tutorial on the MOFA+
1069 downstream analysis in R, we trained a MOFA model using the `run_mofa` function. To project unseen data
1070 from a new view into the latent space Z_i , we computed the pseudo-inverse of the weight matrix W_i . The
1071 latent representations Z_i for all views were then concatenated. For prediction, we applied logistic regression
1072 for binary outcomes and linear regression for continuous outcomes. To classify and predict new, unseen

1073 samples, we used the concatenated latent representation matrix Z in PyTorch (Python 3.10). We evaluated
 1074 different latent factor dimensions from the set $\{1, 5, 10, 15\}$.

1075 **Data Availability**

1076 Simulated data and all data supporting the results are included in Supplementary Data 1–9 and are publicly
 1077 available on GitHub (<https://github.com/daifengwanglab/COSIME>). ROSMAP scRNA-seq was sourced
 1078 from a website (https://compbio.mit.edu/ad_aging_brain/#loading-the-raw-data) and metabolomic
 1079 data can be downloaded from the AD Knowledge Portal - backend (syn10235592). SEA-AD scATAC-seq,
 1080 scRNA-seq, and MERFISH data were obtained from SEA-AD: Seattle Alzheimer’s Disease Brain Cell
 1081 Atlas (<https://cellxgene.cziscience.com/collections/1ca90a2d-2943-483d-b678-b809bf464c30>).

1082 **Code Availability**

1083 The entire COSIME code is implemented in Python. The COSIME Python package and examples can be
 1084 accessed publicly at <https://github.com/daifengwanglab/COSIME>.

1085 **References**

- 1086 1. Hasin, Y., Seldin, M. & Lusi, A. Multi-omics approaches to disease. en. *Genome Biology* **18**, 83. issn:
 1087 1474-760X. [https://genomebiology.biomedcentral.com/articles/10.1186/s13059-](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1215-1)
 1088 [017-1215-1](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1215-1) (2024) (Dec. 2017).
- 1089 2. Subramanian, I., Verma, S., Kumar, S., Jere, A. & Anamika, K. Multi-omics Data Integration,
 1090 Interpretation, and Its Application. eng. *Bioinformatics and Biology Insights* **14**, 1177932219899051.
 1091 issn: 1177-9322 (2020).
- 1092 3. Ding, D. Y., Li, S., Narasimhan, B. & Tibshirani, R. Cooperative learning for multiview analysis. en.
 1093 *Proceedings of the National Academy of Sciences* **119**, e2202113119. issn: 0027-8424, 1091-6490.
 1094 <https://pnas.org/doi/full/10.1073/pnas.2202113119> (2024) (Sept. 2022).
- 1095 4. Singh, A. *et al.* DIABLO: an integrative approach for identifying key molecular drivers from multi-
 1096 omics assays. en. *Bioinformatics* **35** (ed Birol, I.) 3055–3062. issn: 1367-4803, 1367-4811. [https:](https://academic.oup.com/bioinformatics/article/35/17/3055/5292387)
 1097 [//academic.oup.com/bioinformatics/article/35/17/3055/5292387](https://academic.oup.com/bioinformatics/article/35/17/3055/5292387) (2024) (Sept. 2019).

- 1098 5. Argelaguet, R. *et al.* MOFA+: a statistical framework for comprehensive integration of multi-modal
1099 single-cell data. en. *Genome Biology* **21**, 111. issn: 1474-760X. <https://genomebiology.biomedcentral.com/articles/10.1186/s13059-020-02015-1> (2024) (Dec. 2020).
- 1101 6. Lundberg, S. & Lee, S.-I. *A Unified Approach to Interpreting Model Predictions* Version Number: 2.
1102 2017. <https://arxiv.org/abs/1705.07874> (2024).
- 1103 7. Ribeiro, M. T., Singh, S. & Guestrin, C. "Why Should I Trust You?": *Explaining the Predictions of*
1104 *Any Classifier* Version Number: 3. 2016. <https://arxiv.org/abs/1602.04938> (2024).
- 1105 8. Gadzicki, K., Khamsehashari, R. & Zetsche, C. *Early vs Late Fusion in Multimodal Convolutional*
1106 *Neural Networks in 2020 IEEE 23rd International Conference on Information Fusion (FUSION)* (IEEE, Rustenburg, South Africa, July 2020), 1–6. ISBN: 978-0-578-64709-8. <https://ieeexplore.ieee.org/document/9190246/> (2024).
- 1109 9. Dhamdhere, K., Agarwal, A. & Sundararajan, M. *The Shapley Taylor Interaction Index* Version
1110 Number: 2. 2019. <https://arxiv.org/abs/1902.05622> (2024).
- 1111 10. O'Brien, R. J. & Wong, P. C. Amyloid precursor protein processing and Alzheimer's disease. eng.
1112 *Annual Review of Neuroscience* **34**, 185–204. issn: 1545-4126 (2011).
- 1113 11. Hampel, H. *et al.* The Amyloid- Pathway in Alzheimer's Disease. en. *Molecular Psychiatry* **26**,
1114 5481–5503. issn: 1359-4184, 1476-5578. [https://www.nature.com/articles/s41380-021-](https://www.nature.com/articles/s41380-021-01249-0)
1115 [01249-0](https://www.nature.com/articles/s41380-021-01249-0) (2024) (Oct. 2021).
- 1116 12. Kinney, J. W. *et al.* Inflammation as a central mechanism in Alzheimer's disease. en. *Alzheimer's &*
1117 *Dementia: Translational Research & Clinical Interventions* **4**, 575–590. issn: 2352-8737, 2352-8737.
1118 <https://alz-journals.onlinelibrary.wiley.com/doi/10.1016/j.trci.2018.06.014>
1119 (2024) (Jan. 2018).
- 1120 13. Mallela, J., Yang, J. & Shariat-Madar, Z. Prolylcarboxypeptidase: A cardioprotective enzyme. en.
1121 *The International Journal of Biochemistry & Cell Biology* **41**, 477–481. issn: 13572725. <https://linkinghub.elsevier.com/retrieve/pii/S1357272508001003> (2024) (Mar. 2009).
- 1123 14. Ngo, M.-L., Mahdi, F., Kolte, D. & Shariat-Madar, Z. Upregulation of prolylcarboxypeptidase (PRCP)
1124 in lipopolysaccharide (LPS) treated endothelium promotes inflammation. eng. *Journal of Inflammation*
1125 *(London, England)* **6**, 3. issn: 1476-9255 (Jan. 2009).

- 1126 15. Weiner, I. D., Mitch, W. E. & Sands, J. M. Urea and Ammonia Metabolism and the Control of
1127 Renal Nitrogen Excretion. eng. *Clinical journal of the American Society of Nephrology: CJASN* **10**,
1128 1444–1458. issn: 1555-905X (Aug. 2015).
- 1129 16. Gropman, A. L., Summar, M. & Leonard, J. V. Neurological implications of urea cycle disorders.
1130 eng. *Journal of Inherited Metabolic Disease* **30**, 865–879. issn: 1573-2665 (Nov. 2007).
- 1131 17. Moreno, B. *et al.* Methylthioadenosine reverses brain autoimmune disease. eng. *Annals of Neurology*
1132 **60**, 323–334. issn: 0364-5134 (Sept. 2006).
- 1133 18. Ye, B. *et al.* Klf4 glutamylation is required for cell reprogramming and early embryonic development
1134 in mice. en. *Nature Communications* **9**, 1261. issn: 2041-1723. [https://www.nature.com/
1135 articles/s41467-018-03008-2](https://www.nature.com/articles/s41467-018-03008-2) (2024) (Mar. 2018).
- 1136 19. Wang, B. *et al.* Metabolism pathways of arachidonic acids: mechanisms and potential therapeutic
1137 targets. en. *Signal Transduction and Targeted Therapy* **6**, 94. issn: 2059-3635. [https://www.
1138 nature.com/articles/s41392-020-00443-w](https://www.nature.com/articles/s41392-020-00443-w) (2024) (Feb. 2021).
- 1139 20. Gegotek, A. & Skrzydlewska, E. en. in *Vitamins and Hormones* 247–270 (Elsevier, 2023). isbn: 978-
1140 0-443-15768-4. [https://linkinghub.elsevier.com/retrieve/pii/S0083672922000838
1141](https://linkinghub.elsevier.com/retrieve/pii/S0083672922000838) (2024).
- 1142 21. Moffett, J. R., Ross, B., Arun, P., Madhavarao, C. N. & Namboodiri, A. M. A. N-Acetylaspartate in
1143 the CNS: from neurodiagnostics to neurobiology. eng. *Progress in Neurobiology* **81**, 89–131. issn:
1144 0301-0082 (Feb. 2007).
- 1145 22. Rebelos, E. *et al.* Circulating N-Acetylaspartate does not track brain NAA concentrations, cognitive
1146 function or features of small vessel disease in humans. en. *Scientific Reports* **12**, 11530. issn: 2045-
1147 2322. <https://www.nature.com/articles/s41598-022-15670-0> (2024) (July 2022).
- 1148 23. Lima, F. R. S. *et al.* Cellular prion protein expression in astrocytes modulates neuronal survival
1149 and differentiation. en. *Journal of Neurochemistry* **103**, 2164–2176. issn: 0022-3042, 1471-4159.
1150 <https://onlinelibrary.wiley.com/doi/10.1111/j.1471-4159.2007.04904.x> (2024)
1151 (Dec. 2007).
- 1152 24. Salter, M. *et al.* Initial Identification of a Blood-Based Chromosome Conformation Signature for
1153 Aiding in the Diagnosis of Amyotrophic Lateral Sclerosis. en. *EBioMedicine* **33**, 169–184. issn:
1154 23523964. <https://linkinghub.elsevier.com/retrieve/pii/S2352396418302226> (2024)
1155 (July 2018).

- 1156 25. Price, B. R., Johnson, L. A. & Norris, C. M. Reactive astrocytes: The nexus of pathological and clinical
1157 hallmarks of Alzheimer's disease. eng. *Ageing Research Reviews* **68**, 101335. ISSN: 1872-9649 (July
1158 2021).
- 1159 26. González-Reyes, R. E., Nava-Mesa, M. O., Vargas-Sánchez, K., Ariza-Salamanca, D. & Mora-
1160 Muñoz, L. Involvement of Astrocytes in Alzheimer's Disease from a Neuroinflammatory and Oxida-
1161 tive Stress Perspective. eng. *Frontiers in Molecular Neuroscience* **10**, 427. ISSN: 1662-5099 (2017).
- 1162 27. Lawrence, J. M., Schardien, K., Wigdahl, B. & Nonnemacher, M. R. Roles of neuropathology-
1163 associated reactive astrocytes: a systematic review. eng. *Acta Neuropathologica Communications* **11**,
1164 42. ISSN: 2051-5960 (Mar. 2023).
- 1165 28. Zhang, W., Xiao, D., Mao, Q. & Xia, H. Role of neuroinflammation in neurodegeneration devel-
1166 opment. en. *Signal Transduction and Targeted Therapy* **8**, 267. ISSN: 2059-3635. <https://www.nature.com/articles/s41392-023-01486-5> (2024) (July 2023).
- 1168 29. Choi, J. J., Svaren, J. & Wang, D. *Single-cell multi-omics analysis reveals cooperative transcription*
1169 *factors for gene regulation in oligodendrocytes* en. June 2024. <http://biorxiv.org/lookup/doi/10.1101/2024.06.19.599799> (2024).
- 1171 30. Hornig, J. *et al.* The Transcription Factors Sox10 and Myrf Define an Essential Regulatory Network
1172 Module in Differentiating Oligodendrocytes. en. *PLoS Genetics* **9** (ed Barres, B. A.) e1003907. ISSN:
1173 1553-7404. <https://dx.plos.org/10.1371/journal.pgen.1003907> (2024) (Oct. 2013).
- 1174 31. Triner, D. *et al.* Myc-Associated Zinc Finger Protein Regulates the Proinflammatory Response in
1175 Colitis and Colon Cancer via STAT3 Signaling. eng. *Molecular and Cellular Biology* **38**, e00386–18.
1176 ISSN: 1098-5549 (Nov. 2018).
- 1177 32. Bujalka, H. *et al.* MYRF Is a Membrane-Associated Transcription Factor That Autoproteolytically
1178 Cleaves to Directly Activate Myelin Genes. en. *PLoS Biology* **11** (ed french-Constant, C.) e1001625.
1179 ISSN: 1545-7885. <https://dx.plos.org/10.1371/journal.pbio.1001625> (2024) (Aug.
1180 2013).
- 1181 33. Schmidt, K., Glaser, G., Wernig, A., Wegner, M. & Rosorius, O. Sox8 Is a Specific Marker for Muscle
1182 Satellite Cells and Inhibits Myogenesis. en. *Journal of Biological Chemistry* **278**, 29769–29775. ISSN:
1183 00219258. <https://linkinghub.elsevier.com/retrieve/pii/S0021925820842434> (2024)
1184 (Aug. 2003).

- 1185 34. Sarkar, A. & Hochedlinger, K. The sox family of transcription factors: versatile regulators of stem
1186 and progenitor cell fate. eng. *Cell Stem Cell* **12**, 15–30. issn: 1875-9777 (Jan. 2013).
- 1187 35. McGettrick, A. F. & O'Neill, L. A. The Role of HIF in Immunity and Inflammation. en. *Cell*
1188 *Metabolism* **32**, 524–536. issn: 15504131. [https://linkinghub.elsevier.com/retrieve/
1189 pii/S1550413120304150](https://linkinghub.elsevier.com/retrieve/pii/S1550413120304150) (2024) (Oct. 2020).
- 1190 36. Hau, A.-C. *et al.* Transcriptional cooperation of PBX1 and PAX6 in adult neural progenitor cells.
1191 eng. *Scientific Reports* **11**, 21013. issn: 2045-2322 (Oct. 2021).
- 1192 37. Lee, J.-R. Protein tyrosine phosphatase PTPRT as a regulator of synaptic formation and neuronal
1193 development. eng. *BMB reports* **48**, 249–255. issn: 1976-670X (May 2015).
- 1194 38. Sithanandam, G. & Anderson, L. M. The ERBB3 receptor in cancer and cancer gene therapy. eng.
1195 *Cancer Gene Therapy* **15**, 413–448. issn: 1476-5500 (July 2008).
- 1196 39. Fossati, M. *et al.* Trans-Synaptic Signaling through the Glutamate Receptor Delta-1 Mediates In-
1197 hibitory Synapse Formation in Cortical Pyramidal Neurons. eng. *Neuron* **104**, 1081–1094.e7. issn:
1198 1097-4199 (Dec. 2019).
- 1199 40. Lin, Y.-H. *et al.* LIMCH1 regulates nonmuscle myosin-II activity and suppresses cell migration. eng.
1200 *Molecular Biology of the Cell* **28**, 1054–1065. issn: 1939-4586 (Apr. 2017).
- 1201 41. Zhou, A.-X., Hartwig, J. H. & Akyürek, L. M. Filamins in cell signaling, transcription and organ
1202 development. en. *Trends in Cell Biology* **20**, 113–123. issn: 09628924. [https://linkinghub.
1203 elsevier.com/retrieve/pii/S0962892409002876](https://linkinghub.elsevier.com/retrieve/pii/S0962892409002876) (2024) (Feb. 2010).
- 1204 42. Kondoh, K. & Nishida, E. Regulation of MAP kinases by MAP kinase phosphatases. en. *Biochimica*
1205 *et Biophysica Acta (BBA) - Molecular Cell Research* **1773**, 1227–1237. issn: 01674889. <https://linkinghub.elsevier.com/retrieve/pii/S0167488906004538> (2024) (Aug. 2007).
- 1206
- 1207 43. Lomniczi, A. *et al.* Epigenetic regulation of puberty via Zinc finger protein-mediated transcriptional
1208 repression. en. *Nature Communications* **6**, 10195. issn: 2041-1723. [https://www.nature.com/
1209 articles/ncomms10195](https://www.nature.com/articles/ncomms10195) (2024) (Dec. 2015).
- 1210 44. Briscoe, J. *et al.* Homeobox gene Nkx2.2 and specification of neuronal identity by graded Sonic hedge-
1211 hog signalling. en. *Nature* **398**, 622–627. issn: 0028-0836, 1476-4687. [https://www.nature.
1212 com/articles/19315](https://www.nature.com/articles/19315) (2024) (Apr. 1999).

- 1213 45. Chatterjee, M. *et al.* Contactin-2, a synaptic and axonal protein, is reduced in cerebrospinal fluid and
1214 brain tissue in Alzheimer's disease. *eng. Alzheimer's Research & Therapy* **10**, 52. ISSN: 1758-9193
1215 (June 2018).
- 1216 46. Citri, A. & Malenka, R. C. Synaptic Plasticity: Multiple Forms, Functions, and Mechanisms. *en.*
1217 *Neuropsychopharmacology* **33**, 18–41. ISSN: 0893-133X, 1740-634X. [https://www.nature.com/
1218 articles/1301559](https://www.nature.com/articles/1301559) (2024) (Jan. 2008).
- 1219 47. Marques, B. L. *et al.* Neurobiology of glycine transporters: From molecules to behavior. *en. Neuro-*
1220 *science & Biobehavioral Reviews* **118**, 97–110. ISSN: 01497634. [https://linkinghub.elsevier.
1221 com/retrieve/pii/S0149763420304826](https://linkinghub.elsevier.com/retrieve/pii/S0149763420304826) (2024) (Nov. 2020).
- 1222 48. McGowan, M. *et al.* The Application of Neurodiagnostic Studies to Inform the Acute Management
1223 of a Newborn Presenting With Sarbamoyl Shosphate Synthetase 1 Deficiency. *eng. Child Neurology*
1224 *Open* **8**, 2329048X20985179. ISSN: 2329-048X (2021).
- 1225 49. Rak, J. *et al.* Cytohesin 1 regulates homing and engraftment of human hematopoietic stem and
1226 progenitor cells. *eng. Blood* **129**, 950–958. ISSN: 1528-0020 (Feb. 2017).
- 1227 50. Steelman, A. J. *et al.* Activation of oligodendroglial Stat3 is required for efficient remyelination. *eng.*
1228 *Neurobiology of Disease* **91**, 336–346. ISSN: 1095-953X (July 2016).
- 1229 51. Jain, M. *et al.* Role of JAK/STAT in the Neuroinflammation and its Association with Neurological
1230 Disorders. *eng. Annals of Neurosciences* **28**, 191–200. ISSN: 0972-7531 (July 2021).
- 1231 52. Zhou, C. *et al.* JUN is a key transcriptional regulator of the unfolded protein response in acute myeloid
1232 leukemia. *eng. Leukemia* **31**, 1196–1205. ISSN: 1476-5551 (May 2017).
- 1233 53. Caprariello, A. V., Mangla, S., Miller, R. H. & Selkirk, S. M. Apoptosis of oligodendrocytes in the
1234 central nervous system results in rapid focal demyelination. *eng. Annals of Neurology* **72**, 395–405.
1235 ISSN: 1531-8249 (Sept. 2012).
- 1236 54. Yan, M. *et al.* SOCS modulates JAK-STAT pathway as a novel target to mediate the occurrence of neu-
1237 roinflammation: Molecular details and treatment options. *en. Brain Research Bulletin* **213**, 110988.
1238 ISSN: 03619230. [https://linkinghub.elsevier.com/retrieve/pii/S0361923024001217
1239 \(2024\) \(July 2024\).](https://linkinghub.elsevier.com/retrieve/pii/S0361923024001217)

- 1240 55. Dominguez, E., Mauborgne, A., Mallet, J., Desclaux, M. & Pohl, M. SOCS3-mediated blockade of
1241 JAK/STAT3 signaling pathway reveals its major contribution to spinal cord neuroinflammation and
1242 mechanical allodynia after peripheral nerve injury. eng. *The Journal of Neuroscience: The Official*
1243 *Journal of the Society for Neuroscience* **30**, 5754–5766. ISSN: 1529-2401 (Apr. 2010).
- 1244 56. Zhou, Q. & Anderson, D. J. The bHLH Transcription Factors OLIG2 and OLIG1 Couple Neuronal
1245 and Glial Subtype Specification. en. *Cell* **109**, 61–73. ISSN: 00928674. <https://linkinghub.elsevier.com/retrieve/pii/S0092867402006773> (2024) (Apr. 2002).
- 1247 57. Galiano, M. *et al.* Myelin basic protein functions as a microtubule stabilizing protein in differentiated
1248 oligodendrocytes. en. *Journal of Neuroscience Research* **84**, 534–541. ISSN: 0360-4012, 1097-4547.
1249 <https://onlinelibrary.wiley.com/doi/10.1002/jnr.20960> (2024) (Aug. 2006).
- 1250 58. Mendiola, A. S. *et al.* Defining blood-induced microglia functions in neurodegeneration through
1251 multiomic profiling. en. *Nature Immunology* **24**, 1173–1187. ISSN: 1529-2908, 1529-2916. <https://www.nature.com/articles/s41590-023-01522-0> (2024) (July 2023).
- 1253 59. Ferrari-Souza, J. P. *et al.* APOE 4 associates with microglial activation independently of A plaques
1254 and tau tangles. en. *Science Advances* **9**, eade1474. ISSN: 2375-2548. <https://www.science.org/doi/10.1126/sciadv.ade1474> (2024) (Apr. 2023).
- 1256 60. Yin, Z. *et al.* APOE4 impairs the microglial response in Alzheimer’s disease by inducing TGF-
1257 mediated checkpoints. en. *Nature Immunology* **24**, 1839–1853. ISSN: 1529-2908, 1529-2916. <https://www.nature.com/articles/s41590-023-01627-6> (2024) (Nov. 2023).
- 1259 61. Schäfer, M. K.-H. *et al.* Complement C1q Is Dramatically Up-Regulated in Brain Microglia in
1260 Response to Transient Global Cerebral Ischemia. en. *The Journal of Immunology* **164**, 5446–5452.
1261 ISSN: 0022-1767, 1550-6606. [https://journals.aai.org/jimmunol/article/164/10/](https://journals.aai.org/jimmunol/article/164/10/5446/32714/Complement-C1q-Is-Dramatically-Up-Regulated-in)
1262 [5446/32714/Complement-C1q-Is-Dramatically-Up-Regulated-in](https://journals.aai.org/jimmunol/article/164/10/5446/32714/Complement-C1q-Is-Dramatically-Up-Regulated-in) (2024) (May 2000).
- 1263 62. Fraser, D. A., Pisalyaput, K. & Tenner, A. J. C1q enhances microglial clearance of apoptotic neurons
1264 and neuronal blebs, and modulates subsequent inflammatory cytokine production. eng. *Journal of*
1265 *Neurochemistry* **112**, 733–743. ISSN: 1471-4159 (Feb. 2010).
- 1266 63. Voskobiyuk, Y. *et al.* Alzheimer’s disease risk gene BIN1 induces Tau-dependent network hyperex-
1267 citability. eng. *eLife* **9**, e57354. ISSN: 2050-084X (July 2020).

- 1268 64. Matyash, M., Matyash, V., Nolte, C., Sorrentino, V. & Kettenmann, H. Requirement of functional
1269 ryanodine receptor type 3 for astrocyte migration. en. *The FASEB Journal* **16**, 1–25. ISSN: 0892-6638,
1270 1530-6860. <https://onlinelibrary.wiley.com/doi/10.1096/fj.01-0380fje> (2024) (Jan.
1271 2002).
- 1272 65. Zhang, Y. *et al.* Purification and Characterization of Progenitor and Mature Human Astrocytes Reveals
1273 Transcriptional and Functional Differences with Mouse. eng. *Neuron* **89**, 37–53. ISSN: 1097-4199
1274 (Jan. 2016).
- 1275 66. Massart, J. & Zierath, J. R. Role of Diacylglycerol Kinases in Glucose and Energy Homeostasis.
1276 en. *Trends in Endocrinology & Metabolism* **30**, 603–617. ISSN: 10432760. [https://linkinghub.
1277 elsevier.com/retrieve/pii/S1043276019301201](https://linkinghub.elsevier.com/retrieve/pii/S1043276019301201) (2024) (Sept. 2019).
- 1278 67. Lorente-Gea, L., García, B., Martín, C., Quirós, L. M. & Fernández-Vega, I. Heparan sulfate proteo-
1279 glycans and heparanases in Alzheimer’s disease: current outlook and potential therapeutic targets.
1280 eng. *Neural Regeneration Research* **12**, 914–915. ISSN: 1673-5374 (June 2017).
- 1281 68. Sadick, J. S. *et al.* Astrocytes and oligodendrocytes undergo subtype-specific transcriptional changes
1282 in Alzheimer’s disease. en. *Neuron* **110**, 1788–1805.e10. ISSN: 08966273. [https://linkinghub.
1283 elsevier.com/retrieve/pii/S0896627322002446](https://linkinghub.elsevier.com/retrieve/pii/S0896627322002446) (2024) (June 2022).
- 1284 69. Pandya, V. A. & Patani, R. Region-specific vulnerability in neurodegeneration: lessons from normal
1285 ageing. eng. *Ageing Research Reviews* **67**, 101311. ISSN: 1872-9649 (May 2021).
- 1286 70. Davey, J. *et al.* Exploring the role of the posterior middle temporal gyrus in semantic cognition:
1287 Integration of anterior temporal lobe with executive processes. eng. *NeuroImage* **137**, 165–177. ISSN:
1288 1095-9572 (Aug. 2016).
- 1289 71. Chen, S. *et al.* Spatially resolved transcriptomics reveals genes associated with the vulnerability
1290 of middle temporal gyrus in Alzheimer’s disease. en. *Acta Neuropathologica Communications* **10**,
1291 188. ISSN: 2051-5960. [https://actaneurocomms.biomedcentral.com/articles/10.1186/
1292 s40478-022-01494-6](https://actaneurocomms.biomedcentral.com/articles/10.1186/s40478-022-01494-6) (2024) (Dec. 2022).
- 1293 72. Harrison, N. A., Doeller, C. F., Voon, V., Burgess, N. & Critchley, H. D. Peripheral inflammation
1294 acutely impairs human spatial memory via actions on medial temporal lobe glucose metabolism. eng.
1295 *Biological Psychiatry* **76**, 585–593. ISSN: 1873-2402 (Oct. 2014).

- 1296 73. Emani, P. S. *et al.* Single-cell genomics and regulatory networks for 388 human brains. en. *Science*
 1297 **384**, eadi5199. issn: 0036-8075, 1095-9203. [https://www.science.org/doi/10.1126/](https://www.science.org/doi/10.1126/science.adi5199)
 1298 [science.adi5199](https://www.science.org/doi/10.1126/science.adi5199) (2025) (May 2024).
- 1299 74. Gupta, C. *et al.* *Network-based drug repurposing for psychiatric disorders using single-cell genomics*
 1300 en. Dec. 2024. <http://medrxiv.org/lookup/doi/10.1101/2024.12.01.24318008> (2025).
- 1301 75. Chandrashekar, P. B. *et al.* DeepGAMI: deep biologically guided auxiliary learning for multimodal
 1302 integration and imputation to improve genotype–phenotype prediction. en. *Genome Medicine* **15**,
 1303 88. issn: 1756-994X. [https://genomemedicine.biomedcentral.com/articles/10.1186/](https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-023-01248-6)
 1304 [s13073-023-01248-6](https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-023-01248-6) (2025) (Oct. 2023).
- 1305 76. Cohen Kalafut, N., Huang, X. & Wang, D. Joint variational autoencoders for multimodal imputation
 1306 and embedding. en. *Nature Machine Intelligence* **5**, 631–642. issn: 2522-5839. [https://www.](https://www.nature.com/articles/s42256-023-00663-z)
 1307 [nature.com/articles/s42256-023-00663-z](https://www.nature.com/articles/s42256-023-00663-z) (2025) (May 2023).
- 1308 77. Alatkari, S. A. & Wang, D. CMOT: Cross-Modality Optimal Transport for multimodal inference. en.
 1309 *Genome Biology* **24**, 163. issn: 1474-760X. [https://genomebiology.biomedcentral.com/](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-023-02989-8)
 1310 [articles/10.1186/s13059-023-02989-8](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-023-02989-8) (2025) (July 2023).
- 1311 78. Bennett, D. A. *et al.* Religious Orders Study and Rush Memory and Aging Project. eng. *Journal of*
 1312 *Alzheimer's disease: JAD* **64**, S161–S189. issn: 1875-8908 (2018).
- 1313 79. Mathys, H. *et al.* Single-cell atlas reveals correlates of high cognitive function, dementia, and
 1314 resilience to Alzheimer's disease pathology. en. *Cell* **186**, 4365–4385.e27. issn: 00928674. [https:](https://linkinghub.elsevier.com/retrieve/pii/S009286742300973X)
 1315 [//linkinghub.elsevier.com/retrieve/pii/S009286742300973X](https://linkinghub.elsevier.com/retrieve/pii/S009286742300973X) (2024) (Sept. 2023).
- 1316 80. Ben-Kiki, O., Bercovich, A., Lifshitz, A. & Tanay, A. Metacell-2: a divide-and-conquer metacell
 1317 algorithm for scalable scRNA-seq analysis. en. *Genome Biology* **23**, 100. issn: 1474-760X. [https:](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-022-02667-1)
 1318 [//genomebiology.biomedcentral.com/articles/10.1186/s13059-022-02667-1](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-022-02667-1) (2024)
 1319 (Dec. 2022).
- 1320 81. Endo, F. *et al.* Molecular basis of astrocyte diversity and morphology across the CNS in health and
 1321 disease. en. *Science* **378**, eadc9020. issn: 0036-8075, 1095-9203. [https://www.science.org/](https://www.science.org/doi/10.1126/science.adc9020)
 1322 [doi/10.1126/science.adc9020](https://www.science.org/doi/10.1126/science.adc9020) (2024) (Nov. 2022).
- 1323 82. Karch, C. M., Cruchaga, C. & Goate, A. M. Alzheimer's Disease Genetics: From the Bench to the
 1324 Clinic. en. *Neuron* **83**, 11–26. issn: 08966273. [https://linkinghub.elsevier.com/retrieve/](https://linkinghub.elsevier.com/retrieve/pii/S0896627314004851)
 1325 [pii/S0896627314004851](https://linkinghub.elsevier.com/retrieve/pii/S0896627314004851) (2024) (July 2014).

- 1326 83. Malik, M. *et al.* CD33 Alzheimer's Risk-Altering Polymorphism, CD33 Expression, and Exon 2
1327 Splicing. en. *The Journal of Neuroscience* **33**, 13320–13325. ISSN: 0270-6474, 1529-2401. <https://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.1224-13.2013> (Aug.
1328 2013).
- 1330 84. Reiman, E. M. *et al.* GAB2 Alleles Modify Alzheimer's Risk in APOE 4 Carriers. en. *Neu-*
1331 *ron* **54**, 713–720. ISSN: 08966273. [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S0896627307003790)
1332 [S0896627307003790](https://linkinghub.elsevier.com/retrieve/pii/S0896627307003790) (2024) (June 2007).
- 1333 85. Wilkaniec, A., Gassowska-Dobrowolska, M., Strawski, M., Adamczyk, A. & Czapski, G. A. In-
1334 hibition of cyclin-dependent kinase 5 affects early neuroinflammatory signalling in murine model
1335 of amyloid beta toxicity. en. *Journal of Neuroinflammation* **15**, 1. ISSN: 1742-2094. [https://](https://jneuroinflammation.biomedcentral.com/articles/10.1186/s12974-017-1027-y)
1336 jneuroinflammation.biomedcentral.com/articles/10.1186/s12974-017-1027-y
1337 (2024) (Dec. 2018).
- 1338 86. Victoria, G. S., Arkhipenko, A., Zhu, S., Syan, S. & Zurzolo, C. Astrocyte-to-neuron intercellular
1339 prion transfer is mediated by cell-cell contact. en. *Scientific Reports* **6**, 20762. ISSN: 2045-2322.
1340 <https://www.nature.com/articles/srep20762> (2024) (Feb. 2016).
- 1341 87. Lau, S.-F., Cao, H., Fu, A. K. Y. & Ip, N. Y. Single-nucleus transcriptome analysis reveals dysregula-
1342 tion of angiogenic endothelial cells and neuroprotective glia in Alzheimer's disease. en. *Proceedings*
1343 *of the National Academy of Sciences* **117**, 25800–25809. ISSN: 0027-8424, 1091-6490. [https://](https://pnas.org/doi/full/10.1073/pnas.2008762117)
1344 pnas.org/doi/full/10.1073/pnas.2008762117 (2024) (Oct. 2020).
- 1345 88. Roses, A. *et al.* Understanding the genetics of APOE and TOMM40 and role of mitochondrial
1346 structure and function in clinical pharmacology of Alzheimer's disease. en. *Alzheimer's & Dementia*
1347 **12**, 687–694. ISSN: 1552-5260, 1552-5279. [https://alz-journals.onlinelibrary.wiley.](https://alz-journals.onlinelibrary.wiley.com/doi/10.1016/j.jalz.2016.03.015)
1348 [com/doi/10.1016/j.jalz.2016.03.015](https://alz-journals.onlinelibrary.wiley.com/doi/10.1016/j.jalz.2016.03.015) (2024) (June 2016).
- 1349 89. Gabitto, M. I. *et al.* Integrated multimodal cell atlas of Alzheimer's disease. en. *Nature Neuroscience*.
1350 ISSN: 1097-6256, 1546-1726. <https://www.nature.com/articles/s41593-024-01774-5>
1351 (2024) (Oct. 2024).
- 1352 90. Megill, C. *et al.* *cellxgene: a performant, scalable exploration platform for high dimensional sparse*
1353 *matrices* en. Apr. 2021. <http://biorxiv.org/lookup/doi/10.1101/2021.04.05.438318>
1354 (2024).

- 1355 91. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq
1356 data with DESeq2. en. *Genome Biology* **15**, 550. issn: 1474-760X. <https://genomebiology.biomedcentral.com/articles/10.1186/s13059-014-0550-8> (2014) (Dec. 2014).
1357
- 1358 92. Chen, Y. & Colonna, M. Microglia in Alzheimer's disease at single-cell level. Are there common
1359 patterns in humans and mice? eng. *The Journal of Experimental Medicine* **218**, e20202717. issn:
1360 1540-9538 (Sept. 2021).
- 1361 93. Bhaskar, K. *et al.* Microglial derived tumor necrosis factor- drives Alzheimer's disease-related
1362 neuronal cell cycle events. en. *Neurobiology of Disease* **62**, 273–285. issn: 09699961. <https://linkinghub.elsevier.com/retrieve/pii/S0969996113002829> (2024) (Feb. 2014).
1363
- 1364 94. Ajoolabady, A., Lindholm, D., Ren, J. & Pratico, D. ER stress and UPR in Alzheimer's disease:
1365 mechanisms, pathogenesis, treatments. en. *Cell Death & Disease* **13**, 706. issn: 2041-4889. <https://www.nature.com/articles/s41419-022-05153-5> (2024) (Aug. 2022).
1366
- 1367 95. Rangaraju, S. *et al.* Identification and therapeutic modulation of a pro-inflammatory subset of disease-
1368 associated-microglia in Alzheimer's disease. en. *Molecular Neurodegeneration* **13**, 24. issn: 1750-
1369 1326. <https://molecularneurodegeneration.biomedcentral.com/articles/10.1186/s13024-018-0254-8> (2024) (Dec. 2018).
1370
- 1371 96. Sun, N. *et al.* Human microglial state dynamics in Alzheimer's disease progression. en. *Cell* **186**,
1372 4386–4403.e29. issn: 00928674. <https://linkinghub.elsevier.com/retrieve/pii/S0092867423009716> (2024) (Sept. 2023).
1373
- 1374 97. Damodaran, B. B., Flamary, R., Seguy, V. & Courty, N. An Entropic Optimal Transport loss for
1375 learning deep neural networks under label noise in remote sensing images. en. *Computer Vision*
1376 *and Image Understanding* **191**, 102863. issn: 10773142. <https://linkinghub.elsevier.com/retrieve/pii/S1077314219301559> (2024) (Feb. 2020).
1377
- 1378 98. Liaw, R. *et al.* *Tune: A Research Platform for Distributed Model Selection and Training* Version
1379 Number: 1. 2018. <https://arxiv.org/abs/1807.05118> (2025).
- 1380 99. Bergstra, J., Yamins, D. & Cox, D. D. Hyperopt: Distributed asynchronous hyper-parameter opti-
1381 mization. *Astrophysics Source Code Library*. ADS Bibcode: 2022ascl.soft05008B, ascl:2205.008.
1382 <https://ui.adsabs.harvard.edu/abs/2022ascl.soft05008B> (2025) (May 2022).
- 1383 100. Yu, G. Thirteen years of clusterProfiler. en. *The Innovation* **5**, 100722. issn: 26666758. <https://linkinghub.elsevier.com/retrieve/pii/S2666675824001607> (2025) (Nov. 2024).
1384

- 1385 101. Pang, Z. *et al.* MetaboAnalyst 6.0: towards a unified platform for metabolomics data processing,
 1386 analysis and interpretation. en. *Nucleic Acids Research* **52**, W398–W406. ISSN: 0305-1048, 1362-
 1387 4962. <https://academic.oup.com/nar/article/52/W1/W398/7642060> (2025) (July
 1388 2024).
- 1389 102. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular interac-
 1390 tion networks. eng. *Genome Research* **13**, 2498–2504. ISSN: 1088-9051 (Nov. 2003).

1391 **Acknowledgements**

1392 This work was supported by National Institutes of Health grants, R21 NS128761, RF1MH128695,
 1393 R01AG067025, R21NS128761, RF1AG054047, and National Science Foundation Career Award 2144475,
 1394 and a core grant to the Waisman Center from NICHD (P50 HD105353).

1395 **Author Information**

1396 **Authors and Affiliations**

1397 **Department of Population Health Sciences, University of Wisconsin-Madison, Madison, WI 53726,**
 1398 **USA**

1399 Jerome J. Choi, Corinne D. Engelman, and Tianyuan Lu

1400

1401 **Waisman Center, University of Wisconsin-Madison, Madison, WI 53705, USA**

1402 Jerome J. Choi, Noah Cohen Kalafut, and Daifeng Wang

1403

1404 **Department of Computer Sciences, University of Wisconsin-Madison, Madison, WI 5706,**
 1405 **USA**

1406 Noah Cohen Kalafut and Daifeng Wang

1407

1408 **Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, Madison,**
 1409 **WI 53076, USA**

1410 Tim Gruenloh and Daifeng Wang

1411 **Contributions**

1412 Study conception, D.W.; Conceptualization, J.J.C. and D.W.; Methodology, J.J.C. and N.C.K.; Formal Anal-
1413 ysis, J.J.C. and T.G.; Investigation, J.J.C., T.L., and D.W.; Writing–Original Draft, J.J.C.; Writing–Review
1414 and Edit, J.J.C., N.C.K., T.G., C.D.E., T.L., and D.W.; Supervision, T.L. and D.W.; Funding Acquisition,
1415 C.D.E. and D.W..

1416 **Corresponding Authors**

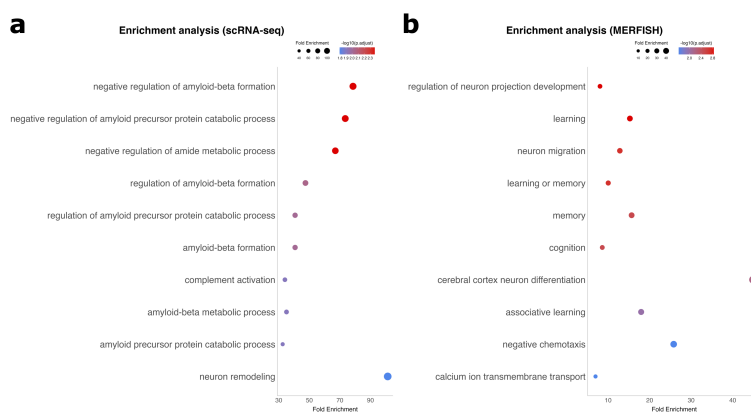
1417 Correspondence to Daifeng Wang.

1418 **Ethics Declarations**

1419 **Competing interests**

1420 The authors declare no competing interests.

Supplementary information



Supplementary Figure 1: Enrichment analyses prioritized for each view.

a The top 20 prioritized genes (absolute feature importance values) in microglia **b** The top 20 prioritized genes (absolute feature importance values) in astrocytes.

Supplementary Data 1

Table showing the prediction performance of COSIME models and benchmark models for binary outcome classification.

Supplementary Data 2

Table showing the prediction performance of COSIME models and benchmark models for continuous outcome classification.

Supplementary Data 3

Feature importance and interactions for the simulated binary high signal early fusion.

Supplementary Data 4

Feature importance and interactions for simulated binary low signal early fusion.

Supplementary Data 5

Feature importance and interactions for the simulated continuous high signal early fusion.

Supplementary Data 6

Feature importance and interactions for the simulated continuous low signal early fusion.

Supplementary Data 7

Feature importance and interactions for the ROSMAP binary classification early fusion.

Supplementary Data 8

Feature importance and interactions for the SEA-AD scRNA-seq and scATAC-seq (oligodendrocytes) early fusion.

Supplementary Data 9

Feature importance and interactions for the SEA-AD scRNA-seq (microglia) and MERFISH (astrocytes) early fusion.

Algorithm 1: COSIME Model Training

Training Loopfor $m = 1$ to M do**Input:**

- Training: $X_1 \in \mathbb{R}^{n \times p_{x1}}, X_2 \in \mathbb{R}^{n \times p_{x2}}$
- Validation: $(X_1^{\text{val}}, X_2^{\text{val}}, y_1^{\text{val}}, y_2^{\text{val}})$
- Test: $(X_1^{\text{test}}, X_2^{\text{test}}, y_1^{\text{test}}, y_2^{\text{test}})$

Initialization:

- Set training iterations M , optimizer, and early stopping criteria
- Set tuning hyperparameters:
batch_size, learning_rate, learning_gamma, $\lambda_{\text{KLD.A.weight}}, \lambda_{\text{KLD.B.weight}}, \lambda_{\text{OT.weight}},$
 $\lambda_{\text{CL.weight}},$ dropout, dim, $p_{\text{earlystop.patience}}, \delta,$ decay

Forward Pass:

- Compute KLD Losses, LOT Loss, and Prediction Loss: $\mathcal{L}_{\text{KLD}_A}(\theta), \mathcal{L}_{\text{KLD}_B}(\theta), \mathcal{L}_{\text{LOT}}(\theta),$
 $\mathcal{L}_{\text{pred}}(\theta)$
- Compute Total Weighted Loss: $\mathcal{L}_{\text{Total.weighted}}(\theta)$

Backpropagation:

- Backpropagate gradients for all parameters, including the transport plan for \mathcal{L}_{LOT}
- Perform optimization step: $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{Total.weighted}}(\theta)$

Early Stopping:Compute validation prediction loss: $\mathcal{L}_{\text{pred}}(\theta)$ for validation**If** $\mathcal{L}_{\text{pred}}(\theta) < \mathcal{L}_{\text{pred}}(\theta_{\text{best}})$ (**improvement in validation loss**):Update $\theta_{\text{best}} = \theta$ Reset patience counter: $p_{\text{counter}} = 0$ **Else If** $p_{\text{counter}} > p_{\text{earlystop.patience}}$ (**early stopping**):

Break training loop

Else:Increment patience counter: $p_{\text{counter}} \leftarrow p_{\text{counter}} + 1$ **End If**

end for

Holdout Evaluation

Evaluate best model on the test set: compute losses and metrics;

OutputTrained best COSIME model and the evaluation metrics;

Algorithm 2: COSIME Feature Importance and Interaction Computation

Input

$X \in \mathbb{R}^{N \times F}$: Input feature matrix,
 $y \in \mathbb{R}^N$: Target labels

Initialization

- Define the model f and input matrix X with N samples and F features.
- Set number of Monte Carlo iterations M .
- Define the batch size, or compute it based on available memory.
- Compute the number of batches: $B = \lceil \frac{N}{\text{batch size}} \rceil$.
- Compute the memory required for each batch:

$$\text{Batch memory (GB)} = \frac{\text{batch size} \times \text{element size of } X \times F}{1 \times 10^9}.$$

- If batch memory exceeds available memory, adjust the batch size dynamically.

Feature Importance Computation

1. For each feature i in X , compute the marginal contribution over M Monte Carlo iterations.
2. Store feature importance values in matrix $\mathcal{S} \in \mathbb{R}^{N \times F}$:
 $\mathcal{S}_{ij} = \phi_j(s_i)$, where $\phi_j(s_i)$ is the feature importance value for feature j for sample s_i .

Feature Interaction Computation

1. For each pair of features (i, j) , compute the interaction effect over M Monte Carlo iterations.
2. Store interaction effects in matrix $\mathcal{I} \in \mathbb{R}^{F \times F}$:
 $\mathcal{I}_{ij} = \mathcal{I}_{ji}$, where \mathcal{I}_{ij} is the interaction effect between features i and j .

Output

$\mathcal{S} \in \mathbb{R}^{N \times F}$: Feature importance matrix,
 $\mathcal{I} \in \mathbb{R}^{F \times F}$: Feature interaction matrix.

CHAPTER 3: Multi-Omics Integration of Transcriptomics and Metabolomics with Machine Learning Uncovers Novel Risk Factors for Alzheimer's Disease

1 **Multi-Omics Integration of Transcriptomics and Metabolomics with** 2 **Machine Learning Uncovers Novel Risk Factors for Alzheimer's Disease**

3

4

Jerome J. Choi¹, Corinne D. Engelman¹, Tianyuan Lu^{1,2}

5

6

¹Department of Population Health Sciences, University of Wisconsin-Madison

7

²Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison

8

9

10 **Abstract**

11 INTRODUCTION

12 Alzheimer's disease (AD) is a neurodegenerative disorder marked by cognitive decline, memory
13 impairment, and functional deterioration. Its complex pathogenesis involves factors such as
14 amyloid plaques, tau tangles, neuroinflammation, and synaptic dysfunction, but the precise
15 mechanisms remain unclear, hindering effective treatment. Genetic, environmental, and lifestyle
16 factors contribute to AD risk, yet their interactions are poorly understood. Recent advances in
17 transcriptomics and metabolomics have shed light on the molecular underpinnings of AD, with
18 gene expression alterations and metabolic disruptions implicated in disease progression. However,
19 the interplay between these omics layers suggests a need for integrative approaches.

20 METHODS

21 This study utilizes the Cooperative multi-view integration and Scalable and Interpretable Model
22 Explainer (COSIME) machine learning method to combine imputed whole blood transcriptomics

23 and measured plasma metabolomics data. A predictive model was developed through 5-fold cross-
24 validation to identify genes and metabolites associated with AD-related phenotypes, including
25 Preclinical Alzheimer's Cognitive Composite 3 (PACC3) and plasma phosphorylated tau 217
26 (ptau217), from the Wisconsin Registry for Alzheimer's Prevention (WRAP) cohort (N = 1,046
27 for both PACC3 and ptau217). The model performance was validated using an independent dataset
28 from the Wisconsin Alzheimer's Disease Research Center (ADRC) cohort (N = 85 for PACC3 and
29 N = 89 for ptau217). Feature importance and interactions were assessed to identify potential risk
30 factors.

31 RESULTS

32 The normalized root mean square error (NRMSE) was 0.654 ± 0.004 for predicting PACC3 and
33 1.361 ± 0.046 for predicting ptau217 in the WRAP cohort. The NRMSE was 1.095 for predicting
34 PACC3 and 1.795 for predicting ptau217 in the Wisconsin ADRC cohort. Biomarker importance
35 and pairwise interaction values were computed. Several genes known to down-regulate cognitive
36 function (e.g., *HOXA4* and *LRFN4*) and others associated with the up-regulation of ptau217 (e.g.,
37 *ACAD9* and *GPAT3*) were highly ranked based on our biomarker importance scores. Similarly,
38 metabolites that down-regulate cognitive function (e.g., sphingomyelin (d17:2/16:0, d18:2/15:0)
39 and N-carbamoylvaline) and those that up-regulate ptau217 (e.g., S-adenosylhomocysteine and
40 propionylcarnitine (C3)) were also identified as highly ranked. Some metabolites, such as
41 sphingomyelin (d17:2/16:0, d18:2/15:0) and tetradecadienoate (14:2), may interact synergistically,
42 while sphingomyelin (d18:2/16:0, d18:1/16:1) and choline may interact antagonistically.

43 DISCUSSION

44 There is significant potential for using integrated transcriptomics and metabolomics data to
45 identify novel biomarkers and predict AD phenotypes. Genes and metabolites associated with
46 cognitive decline and tau pathology were identified, offering insights into additional molecular
47 mechanisms underlying AD. The identified interactions between certain metabolites provide
48 further clues to the complex metabolic dysregulation that may contribute to disease pathogenesis,
49 underscoring the importance of integrative approaches for understanding AD.

50 **1 BACKGROUND**

51 Alzheimer's disease (AD) is a neurodegenerative disease characterized by progressive cognitive
52 decline, memory impairment, and functional deterioration, representing a major public health
53 challenge¹. The pathogenesis of AD is multifactorial, involving the accumulation of amyloid
54 plaques, neurofibrillary tangles, neuroinflammation, and synaptic dysfunction.²⁻⁵ Despite
55 significant advances in understanding the molecular underpinnings of the disease, the precise
56 mechanisms remain unclear, and the complexity of AD presents a challenge for early diagnosis
57 and effective therapeutic development. Genetic, environmental, and lifestyle factors all contribute
58 to AD risk,^{6,7} but the interactions between these factors are poorly understood. As a result, there is
59 a pressing need for more comprehensive approaches to elucidate the molecular landscape of AD
60 and identify novel biomarkers that could aid in early detection and intervention.

61

62 In recent years, transcriptomics and metabolomics have emerged as key omics approaches for
63 investigating the molecular basis of AD.⁸ Transcriptomics, which analyzes gene expression
64 patterns, provides insights into the dysregulated cellular processes in AD, such as neuronal loss,
65 synaptic dysfunction, and inflammatory responses.⁹ Alterations in the expression of genes related

66 to amyloid processing, tau metabolism, and neuroinflammation have been linked to disease
67 progression, highlighting the role of transcriptional dysregulation in AD pathophysiology.¹⁰ On
68 the other hand, metabolomics focuses on the small molecules involved in cellular metabolism,
69 offering a unique window into the biochemical changes that occur in AD.¹¹ Studies have shown
70 that metabolites associated with energy metabolism, oxidative stress, and neurotransmitter
71 pathways are altered in AD, suggesting that disrupted metabolic networks contribute to disease
72 onset and progression.¹² While both transcriptomics and metabolomics provide valuable insights
73 into AD, the interplay between genes and metabolites suggests that a more integrative approach is
74 required.¹³ The interactions between genetic and metabolic alterations play a critical role in AD
75 pathophysiology,¹⁴ and understanding how these molecular networks influence each other is
76 essential for unraveling the complexity of AD mechanisms.

77 Multi-view integration of omics data provides a comprehensive understanding of systems biology
78 and the complex interactions between biomarkers, particularly in complex diseases such as
79 AD.^{15,16} By combining data from multiple omics layers, it is possible to capture a more holistic
80 view of the molecular mechanisms underlying AD. Moreover, this integrative approach may allow
81 for the identification of key biomarkers and pathways that may be missed when examining a single
82 omics layer in isolation.¹⁷ The vast and complex nature of multi-omics data requires advanced
83 computational tools to analyze and interpret the intricate relationships between genes and
84 metabolites. Machine learning techniques offer powerful methods for uncovering patterns within
85 large, high-dimensional datasets.¹⁸ By leveraging machine learning algorithms, we can integrate
86 diverse omics data to predict disease phenotypes, identify novel biomarkers, and better understand
87 the underlying pathophysiology of AD.¹⁹ This study applies a state-of-art machine learning method,
88 Cooperative multi-view integration and Scalable and Interpretable Model Explainer (COSIME)²⁰

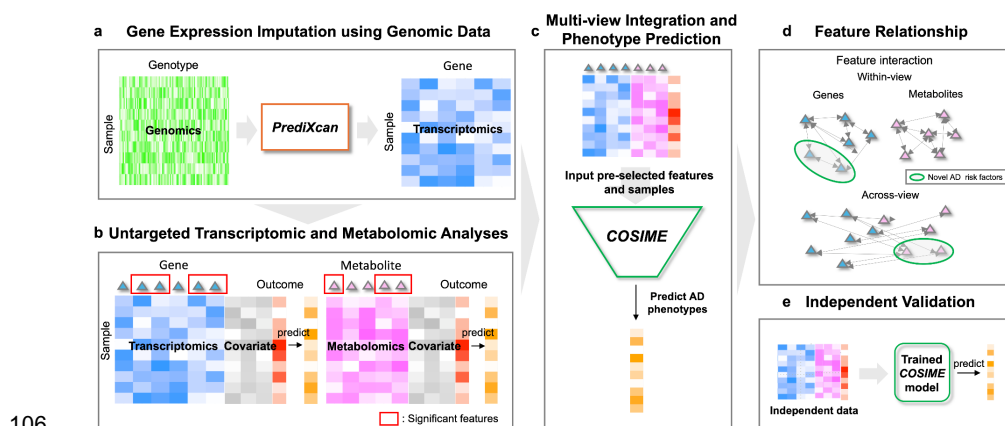
89 to integrate imputed transcriptomics and metabolomics, predict AD-related phenotypes, and
90 identify potential AD-risk factors from the Wisconsin Registry for Alzheimer's Prevention
91 (WRAP)²¹ cohort. The cohort from the Wisconsin Alzheimer's Disease Research Center (ADRC)
92 is used for independent validation.

93 **2 METHODS**

94 **2.1 Study design**

95 An overview of the study design is presented in **FIGURE 1**. Gene expression data were imputed
96 using the genotypes of WRAP and Wisconsin ADRC participants via PrediXcan²². To reduce the
97 dimensionality, untargeted transcriptomic and metabolomic analyses were performed, including
98 age, sex, *APOE* score, and years of education (for Preclinical Alzheimer's Cognitive Composite 3
99 (PACC3) only) as covariates, to predict AD-related phenotypes in WRAP using linear regression.
100 Subsequently, for each AD-related phenotype, nominally significant genes and metabolites (p-
101 value < 0.05) were input into a COSIME model, which was trained to predict the outcome. After
102 training the COSIME model, biomarker importance for individual genes and metabolites, as well
103 as their interaction values, were computed. Lastly, the Wisconsin ADRC cohort was used to
104 independently validate the prediction performance of the COSIME²⁰ models.

105



106

107

108 **FIGURE 1**

109 Workflow diagram. **a** Gene expression data for the WRAP participants were imputed using PrediXcan. **b** Untargeted
 110 transcriptomics and metabolomics adjusting for covariates were implemented to pre-select significant biomarkers. **c**
 111 Multi-omics data were integrated and AD phenotypes were predicted using COSIME. **d** Biomarker importance and
 112 interaction values for both within-view and across-view were computed by COSIME. **e** Independent validation for
 113 predicting AD phenotypes using the Wisconsin ADRC cohort as independent multi-view data was performed.

114 **2.2 Study participants and cohort description**

115 The cohort for the main analysis, WRAP, is an ongoing longitudinal study designed to identify
 116 midlife factors associated with the development of AD. Participant enrollment began in 2001, with
 117 the first follow-up visit occurring 2 to 4 years after the baseline visit, and all subsequent visits
 118 taking place at 2-year intervals thereafter. WRAP participants were dementia-free at enrollment
 119 (mean age 54 years). All study procedures were approved by the University of Wisconsin School
 120 of Medicine and Public Health Institutional Review Board and are in accordance with the
 121 Declaration of Helsinki.

122

123 Wisconsin ADRC participants were recruited from multiple sources such as memory clinic
124 providers, events, media, partner organizations, and other studies. Adults aged 45 years or older
125 with decisional capacity and English fluency were eligible for enrollment. Exclusion criteria
126 included active major medical or psychiatric illnesses and lack of a study partner. This study
127 included participants who have normal cognitive status based on their diagnoses.

128 **2.3 Alzheimer’s disease outcomes**

129 PACC3 is a global cognitive composite score used to assess the cognitive performance of
130 participants in the WRAP study. During each study visit, participants complete a detailed cognitive
131 battery, which is described in more detail elsewhere²¹. Longitudinal global cognitive performance
132 is evaluated using a three-test version of the modified PACC²³, which includes the Rey Auditory
133 Verbal Learning Test (AVLT; Trials 1–5)²⁴, the Wechsler Memory Scale Logical Memory II²⁵,
134 and the Wechsler Adult Intelligence Scale-Revised Digit Symbol Substitution.²⁶ The Logical
135 Memory II test was introduced at Visit 2 for most participants (94%), so the PACC3 baseline is
136 generally at Visit 2, with Visit 1 being the baseline for recent enrollees where it was administered
137 at Visit 1. Each test score is converted into a Z-score, using the mean and standard deviation from
138 baseline scores of cognitively unimpaired individuals.

139

140 The Wisconsin ADRC uses a slightly different set of neuropsychological tests compared to WRAP.
141 To maximize the use of available data, we collaborated with neuropsychologists at the Wisconsin
142 ADRC to create a PACC-3-Trail Making Test (TMT) score for the replication analysis.
143 Specifically, we converted the Craft Story score into an estimated Logical Memory score using a

144 published crosswalk table and followed previous methodology by substituting the Digit Symbol
145 score with the total completion time from the TMT-B test.^{23,27} Since the crosswalk table only
146 provides a conversion for the Craft Story score to the Logical Memory score for the first five visits,
147 we limited the replication analysis to data from these five visits. The final PACC-3-TMT score
148 was calculated by standardizing and averaging the results from the AVLT, the estimated Logical
149 Memory score, and the reversed TMT-B test scores.

150

151 Plasma ptau217 assay results are from University of Wisconsin ADRC Biofluid Lab and were
152 analyzed using the ALZpath pTau217 v2 assay on the Quanterix HD-X. Samples were collected
153 between Aug 2011 and June 2023. 5,082 samples were analyzed using kit lot# 999008 and 309
154 samples were analyzed using kit lot# 999024. A bridging study was performed using 35 samples
155 to determine the relationship between lots. The Pearson correlation coefficient was 0.97 ($p < 2.2e-$
156 16) and linear regression was used to standardize the data from lot# 999024 to lot# 999008 ($y =$
157 $0.65514x + 0.02667$).

158 **2.4 Plasma metabolite collections and quality control**

159 2.4.1 Collection

160 Participants underwent venipuncture, most after fasting ≥ 8 hours, and provided thirty mL of blood
161 into 3×10 mL lavender top EDTA tubes (BD 366643; Franklin Lakes, New Jersey, USA). Samples
162 were mixed gently by inverting 10–12 times and were centrifuged 15 min at 2000 g at room
163 temperature within 1 h of collection. Plasma samples were aliquoted into 2 mL cryovials (Wheaton
164 Cryolite W985863; Millville, New Jersey, USA). Aliquoted plasma was frozen at -80°C within
165 90 min and stored until overnight shipment to Metabolon, Inc. (Durham, NC), which similarly

166 kept samples frozen at -80° C until analysis. Participants who fasted <8 hours prior to venipuncture
167 were excluded (n = 298).

168 2.4.2 Quality control

169 Samples with missing values for more than 40% of metabolites were removed (0 plasma samples)
170 prior to analysis. Non-xenobiotic metabolites with missing values for >75% (n=43), and xenobiotic
171 metabolites with missing values for >99% (n=75) of all samples were removed.

172 2.4.3 Fasting status

173 Samples from individuals who fasted less than 8 hours prior to venipuncture, or who did not
174 provide fasting status were removed.

175 **2.5 Genotyping, quality control, and imputation**

176 2.5.1 Genotyping array data

177 DNA was extracted from whole blood samples and genotyped at the University of Wisconsin
178 Biotech Center in three batches (2017, 2021, 2024) using the Illumina Multi-Ethnic Genotyping
179 Array (MEGA).

180 2.5.2 Processing and Quality Control

181 Initial quality control was performed separately for each batch using PLINK v1.9²⁹. SNPs mapped
182 to chromosome 0 were removed, followed by variants missing in >5% of samples and samples
183 with >5% missing variants. Additional filters were applied to remove variants missing in >2% of
184 remaining samples and samples with >2% missing variants. SNPs within the ACTG nucleotide

185 bases were retained, and sex discrepancies were checked using X chromosome homozygosity,
186 removing 10 samples with inconsistent reported sex. The data were remapped from hg37 to hg38
187 using the UCSC liftover tool²⁹. The HRC checker tool was used to process the data before
188 imputation, removing duplicates, mismatches, and monomorphic SNPs. Data were sorted by
189 chromosome and saved as VCF files using BCFtools³¹. Imputation was done via the TOPMed
190 Imputation Server using Eagle v2.4 and the TOPMed reference panel r3³¹⁻³³. After imputation,
191 variants with a low INFO score ($R^2 < 0.8$) were removed. Data were merged and filtered for
192 genotyping rate $>98\%$ and MAF >0.001 . The final dataset consisted of 14,980,298 SNPs shared
193 across all datasets.

194 2.5.3 Relatedness and genetic ancestry

195 A pruned set of the genotyped SNPs, using PLINK v1.9 --indep-pairwise 50 10 0.1, was used to
196 estimate relatedness and genetic similarities to 1000 Genomes superpopulations³⁵. Relationship
197 inference was performed using KING v2.3.1³⁶ to identify related samples. The pruned SNPs, in a
198 subset of unrelated samples, were used in a supervised learning model in ADMIXTURE v1.3.0,
199 using the default maximum likelihood estimation, to determine proportions of genetic similarities
200 between the Wisconsin samples and the 1000 Genomes defined superpopulations (AFR, AMR,
201 EAS, EUR, SAS). Then the subset of related samples was projected onto the resulting population
202 structure output from ADMIXTURE³⁷.

203 2.5.4 *APOE* score

204 The *APOE* score is a weighted risk score for the effect of the six *APOE* genotypes on AD
205 neuropathology. It is calculated using the natural log (ln) of the odds ratios (OR) for the six

206 different *APOE* genotypes, based on a study of AD neuropathology³⁹. The *APOE* genotypes
207 considered are $\epsilon 2\epsilon 2$, $\epsilon 2\epsilon 3$, $\epsilon 3\epsilon 3$, $\epsilon 2\epsilon 4$, $\epsilon 3\epsilon 4$, and $\epsilon 4\epsilon 4$, with $\epsilon 3\epsilon 3$ serving as the reference group.⁴⁰

208 **2.6 Gene expression imputation**

209 The recently published pre-trained database, which used GTEx whole blood data, was obtained
210 from the PredictDB Data Repository⁴¹. Genotypes from WRAP and Wisconsin ADRC participants
211 were used to impute whole blood gene expression levels for the participants via PrediXcan.
212 Specifically, genetic risk scores for predicting tissue-specific gene expression levels were pre-
213 trained based on GTEx, which contains measured individual-level genetics and transcriptomics
214 data. PrediXcan then calculates genetically predicted whole blood gene expression levels using
215 genotypes of WRAP and Wisconsin ADRC participants, although transcriptomics data for these
216 participants were not measured.

217 **2.7 Statistical analyses**

218 **2.6.1 Sample selection**

219 Participants were selected based on the availability of plasma metabolite measurements, genotypes,
220 outcome data, and covariates such as age, sex, years of education, and *APOE* score. Among
221 multiple visits, we selected the last observation for each participant, where the outcome was
222 measured at or within six months after the plasma metabolites. Participants were then filtered based
223 on their European genetic ancestry admixture (>0.5). For the Wisconsin ADRC cohort, participants
224 diagnosed with dementia were excluded.

225 **2.7.1. Pre-selecting potential AD-risk factors**

226 A multivariate linear regression model was fitted for each biomarker (6,678 genes and 1,125
227 plasma metabolites) and each outcome in WRAP. In each model, a biomarker and covariates such
228 as age, gender, *APOE* score, and years of education (for PACC3 only). were used as predictors to
229 predict an outcome as follows:

230

$$231 \quad PACC3 = age + gender + APOE \text{ score} + years \text{ of education} + gene$$

$$232 \quad PACC3 = age + gender + APOE \text{ score} + years \text{ of education} + metabolite$$

$$233 \quad ptau217 = age + gender + APOE \text{ score} + gene$$

$$234 \quad ptau217 = age + gender + APOE \text{ score} + metabolite$$

235

236 Years of education were included only for PACC3, as PACC3 and years of education are correlated.
237 For each outcome, genes and metabolites with a p-value < 0.05 were considered nominally
238 significant biomarkers. These untargeted imputed transcriptomic and metabolomic analyses were
239 performed to improve the computational efficiency of downstream machine learning modelling by
240 focusing on potentially relevant associations.

241 **2.7.2. Multi-omics integration and phenotype prediction**

242 Significant biomarkers, pre-selected from untargeted imputed transcriptomics and metabolomics
243 in WRAP, were input into the multi-omics models using COSIME²⁰, Cooperative Multi-view
244 Integration and Scalable Interpretable Model Explainer. COSIME integrates multi-view data
245 leveraging deep neural network encoders and Learnable Optimal Transport techniques and

246 implements a mechanism for assessing feature importance within each view, as well as quantifying
 247 feature interactions by estimating Shapley values and Shapley-Taylor indices.

248

249 Two models were trained for predicting PACC3 and ptau217, respectively. Hyperparameter
 250 optimization was performed using *Ray Tune* with Distributed Asynchronous Hyperparameter
 251 Optimization (*Hyperopt*). Early stopping was employed to prevent overfitting, based on
 252 performance improvements on the test data during the training phase. 5-fold cross-validation was
 253 used to evaluate the prediction performance of the best model for each fold. The best model was
 254 chosen based on the minimum prediction loss.

255

256 The prediction performance of the models was assessed using NRMSE to compare different scales
 257 of prediction. The NRMSE standardizes the error by normalizing it with respect to the interquartile
 258 range (IQR) of the observed data. The equation for NRMSE using IQR is given by:

259

$$260 \quad NRMSE = \frac{\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}}{IQR(y)}$$

261

262 Where y_i represents the actual values, \hat{y}_i denotes the predicted values, and n is the number of
 263 participants. The term $IQR(y)$ refers to the interquartile range of actual values, which is calculated
 264 as $Q_3 - Q_1$, where Q_3 and Q_1 are the third and first quartiles of the actual values, respectively.

265 **2.7.3 Computing biomarker importance and interaction**

266 Individual biomarker importance and interactions within each omics data type and across-omics
 267 data types in WRAP were computed using a trained model based on 75% of the samples and
 268 holdout data including the remaining 25% of the samples for each outcome. These calculations
 269 were conducted using COSIME.

270 **3 RESULTS**

271 **3.1. Descriptive statistics for the participants**

272 This study included 1,046 participants from the WRAP cohort for both outcomes: the PACC3 and
 273 ptau217, for the main analysis. Additionally, 85 participants for PACC3 and 89 participants for
 274 ptau217 were included from the Wisconsin ADRC cohort for independent validation. The
 275 demographic characteristics of individuals across cohorts and outcomes are described in **TABLE**
 276 **1.**

277

278 **TABLE 1.** Demographic of participants in two cohorts.

Cohort	WRAP		Wisconsin ADRC	
	PACC3	ptau217	PACC3	ptau217
Outcome	PACC3	ptau217	PACC3	ptau217
Total participants	1,046	1,046	85	89
Age (mean [SD])	67.04 (6.89)	67.20 (6.97)	69.14 (7.65)	68.99 (7.70)
Females	716 (68.5%)	716 (68.5%)	50 (58.8%)	54 (60.7%)

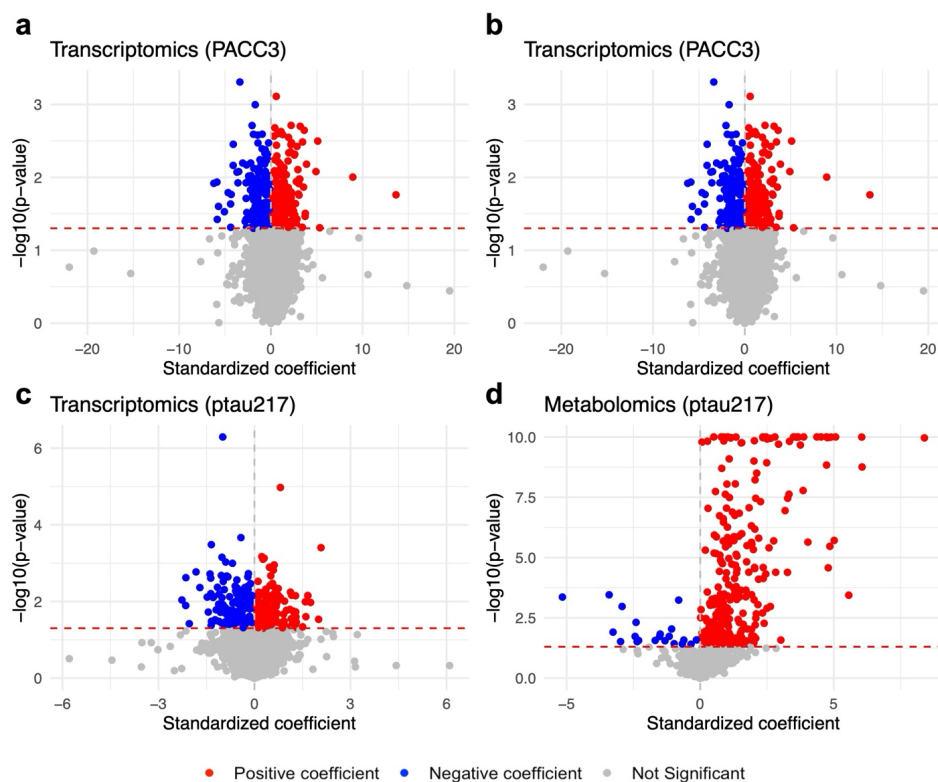
	330	330	35	35
Males	(31.5%)	(31.5%)	(41.2%)	(39.3%)
<i>APOE</i> score (mean [SD])	0.63 (1.04)	0.63 (1.04)	0.61 (1.03)	0.64 (1.03)
Years of education (mean [SD])	15.80 (2.23)	-	16.88 (2.16)	-

279

280 3.2 Untargeted imputed transcriptomics and metabolomics

281 Out of 6,678 genes with imputed whole blood expression levels, 359 were significant for PACCC3
 282 and 348 for ptau217. Among 1,125 measured plasma metabolites, 150 were significant for PACCC3
 283 and 301 for ptau217 (**FIGURE 2**). These genes and metabolites were included in downstream
 284 machine learning modelling using COSIME.

285



286

287 **FIGURE 2**

288 Volcano plots summarizing results of untargeted imputed transcriptomic and metabolomic analyses. P-values (y-axis)
 289 and magnitudes of associations (x-axis) for **a** untargeted imputed transcriptomics predicting PACC3, **b** untargeted
 290 metabolomics predicting PACC3, **c** untargeted imputed transcriptomics predicting ptau217, and **d** untargeted
 291 metabolomics predicting ptau217.

292 **3.3 AD phenotype prediction**

293 For the main analysis, the prediction performance of each model was evaluated five times (using
 294 5-fold cross-validation) on holdout test data from the WRAP cohort. As an independent validation,

295 the trained COSIME models for the outcomes were used to assess their prediction performance on
 296 the Wisconsin ADRC cohort. As a result, the NRMSE was 0.654 ± 0.004 for predicting PACC3
 297 and 1.361 ± 0.046 for predicting ptau217 in the WRAP cohort (**TABLE 2**). The NRMSE was
 298 1.095 for predicting PACC3 and 1.795 for predicting ptau217 in the Wisconsin ADRC cohort
 299 (**TABLE 2**).

300

301 **TABLE 2.** Normalized Root Mean Squared Error (NRMSE) for each cohort and their outcomes.

Cohort	WRAP		Wisconsin ADRC	
Outcome	PACC3	ptau217	PACC3	ptau217
NRMSE	0.654 ± 0.004	1.361 ± 0.046	1.095	1.795

302 WRAP: mean \pm 1.96 * standard deviation.

303 **3.4 Biomarker importance and interactions**

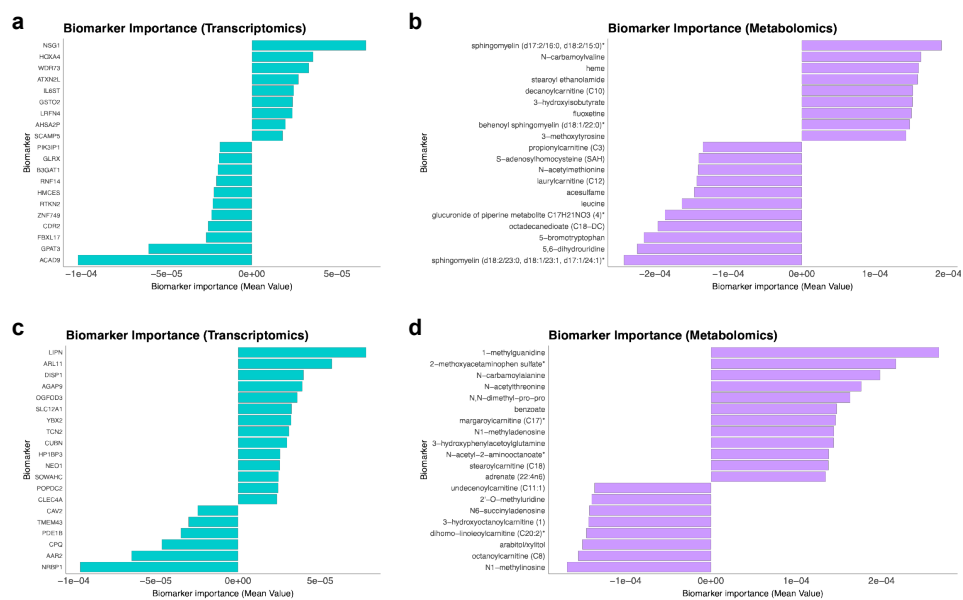
304 **FIGURE 3a** displays the top 20 genes ranked by their absolute biomarker importance values for
 305 predicting PACC3, while **FIGURE 3b** shows the top 20 metabolites with the highest biomarker
 306 importance values for predicting PACC3. Biomarkers with positive values for PACC3 may be
 307 down-regulated in AD or dementia-related pathways, whereas those with negative values may be
 308 up-regulated. **FIGURE 3c** presents the top 20 genes based on absolute biomarker importance
 309 values for predicting ptau217, and **FIGURE 3d** shows the top 20 metabolites for ptau217
 310 prediction. For ptau217, biomarkers with positive values may be up-regulated in AD or dementia-
 311 related pathways, and those with negative values may be down-regulated. Several genes known to
 312 down-regulate cognitive function (e.g., *HOXA4* and *LRFN4*) and others associated with the up-
 313 regulation of ptau217 (e.g., *ACAD9* and *GPAT3*) were highly ranked based on our biomarker

314 importance scores. Several metabolites that down-regulate cognitive function (e.g., sphingomyelin
 315 (d17:2/16:0, d18:2/15:0) and N-carbamoylvaline) and those that up-regulate ptau217 (e.g., S-
 316 adenosylhomocysteine and propionylcarnitine (C3)) were also identified as highly ranked.

317

318 For PACC3 (**FIGURE 3a**), *HOXA4* and *LRFN4* are downregulated, while *ACAD9* and *GPTA3*
 319 were upregulated. In terms of metabolites (**FIGURE 3b**), sphingomyelin (d17:2/16:0, d18:2/15:0)
 320 and N-carbamoylvaline were downregulated, whereas S-adenosylhomocysteine and
 321 propionylcarnitine (C3) were upregulated. For ptau217 (**FIGURE 3c**), *LIPN* and *ARL11* were
 322 upregulated, while *CLEC4A* and *PDE1B* were downregulated. Lastly, for metabolites (**FIGURE**
 323 **3d**), margoylcarnitine and N-acetyl-2-aminooctanoate were upregulated, while
 324 Undecenoylcarnitine (C11:1) and 2'-O-methyluridine were downregulated.

325



326

327 **FIGURE 3**

328 Biomarker importance values for PACC3 and ptau217. **a** Gene importance values for PACC3. **b** Metabolite
329 importance values for PACC3. **c** Gene importance values for ptau217. **d** Metabolite importance values for ptau217.

330

331 The top 50 metabolite-metabolite pairs ranked by their absolute biomarker interaction values for
332 predicting PACC3 are displayed in **FIGURE 4a** and the top 50 metabolite-metabolite pairs ranked

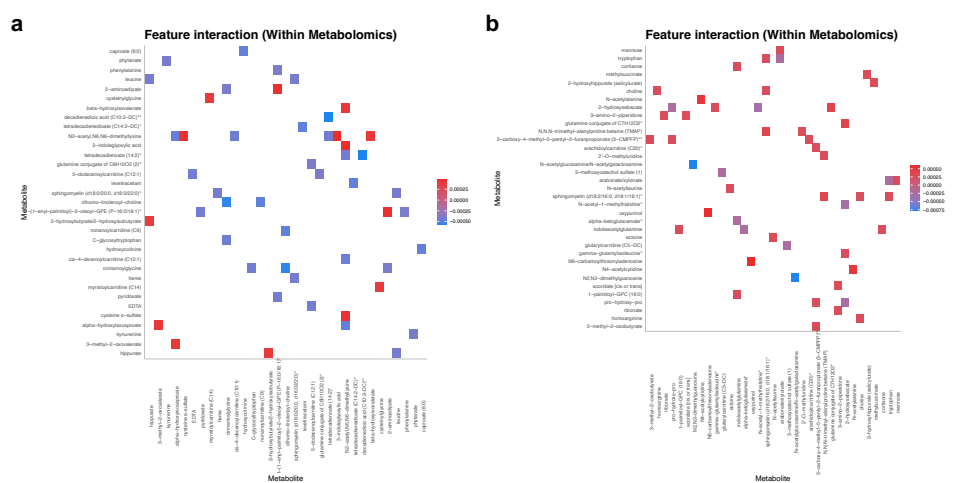
333 by their absolute biomarker interaction values for predicting ptau217 are displayed in **FIGURE**

334 **4b**. Several pairs were notable, such as the potential antagonistic interaction between

335 decadienedioic acid (C10:2-DC)** and tetradecadienoate (14:2) for PACC3 and the potential

336 synergistic interaction between sphingomyelin (d18:2/16:0, d18:1/16:1) and choline for ptau217.

337



338

339 **FIGURE 4**

340 Metabolite pairwise interaction values for PACC3 and ptau217. **a** Metabolite pairwise interaction values for PACC3.

341 **b** Metabolite pairwise interaction values for ptau217.

342 **4 DISCUSSION**

343 In this study, we aimed to predict two critical AD-related phenotypes—PACC3 and plasma
344 ptau217—using multi-omics data, including imputed transcriptomics and metabolomics,
345 combined with COSIME, a novel machine learning approach, for the WRAP cohort. We then
346 assessed the individual biomarker importance and explored their pairwise interactions within the
347 WRAP cohort. To validate the prediction performance of the trained models, we conducted
348 independent validation by applying the same models to predict outcomes in the Wisconsin ADRC
349 cohort. The analyses of biomarker importance and pairwise interactions provide valuable insights
350 into potential AD risk factors and highlights key molecular features associated with disease
351 progression.

352 Genes supporting neuronal growth and synaptic signaling are disrupted in AD. *HOXA4*, a key
353 player in early brain development and neuronal differentiation, may be downregulated in AD,
354 potentially impairing neuronal growth and connectivity⁴⁰. Similarly, *LRFN4*, which supports
355 synaptic communication, may also be reduced, leading to synaptic dysfunction—a core feature of
356 AD⁴¹. *PDE1B*, important for cyclic nucleotide signaling and memory formation, could also be
357 downregulated, disrupting synaptic plasticity and contributing further to cognitive deficits⁴². These
358 downregulations are associated with lower PACC3 scores, while the absence of corresponding
359 upregulation limits any compensatory effect linked to ptau217.

360 Disruptions in lipid metabolism impact membrane structure and function in AD. *GPAT3* supports
361 lipid metabolism essential for maintaining cellular membranes and neuronal function⁴³. In AD,
362 dysregulation in lipid homeostasis may compromise these processes. *LIPN*, another lipid
363 metabolism gene, might be upregulated in response to tau pathology, possibly acting as a

364 compensatory mechanism to preserve membrane integrity⁴⁴. Sphingomyelins, such as
365 sphingomyelin (d17:2/16:0, d18:2/15:0), are critical for membrane structure and signaling, and
366 their downregulation in AD reflects impaired lipid processing⁴⁵. Tetradecadienoate (14:2), which
367 helps maintain membrane fluidity, may also be downregulated, further destabilizing membranes
368 and exacerbating neurodegeneration⁴⁶. The observed downregulations correspond with reduced
369 PACC3 scores, while *LIPN* upregulation may indicate a ptau217-related compensatory stress
370 response.

371 Mitochondrial dysfunction and impaired energy metabolism are consistent hallmarks of AD.
372 *ACAD9* plays a central role in mitochondrial fatty acid metabolism, and its function is critical for
373 sustaining neuronal energy needs⁴⁷. Carnitines like propionylcarnitine (C3) and
374 margaroylcarnitine (C17) facilitate fatty acid transport into mitochondria⁴⁸. While C3 may be
375 downregulated in AD, C17 could be upregulated in an attempt to support mitochondrial energy
376 production. Similarly, undecenoylcarnitine (C11:1) may be reduced, indicating disrupted fatty acid
377 metabolism and diminished neuronal energy capacity⁴⁹. These mitochondrial expression patterns
378 show that downregulations are associated with poorer PACC3 performance, whereas upregulated
379 components such as C17 may be linked to ptau217 elevation.

380 Altered amino acid metabolism may reflect compensatory responses or metabolic breakdown. N-
381 carbamoylvaline, involved in amino acid metabolism, could be downregulated in AD, reflecting
382 general metabolic dysfunction⁵⁰. In contrast, N-acetyl-2-aminooctanoate may be upregulated,
383 potentially as an adaptive response to the altered metabolic environment in AD-affected neurons⁵⁰.
384 These changes suggest that amino acid metabolite downregulation is tied to reduced PACC3, while
385 upregulation is related to increased ptau217 expression.

386 Epigenetic and RNA-processing disruptions contribute to AD pathology. S-adenosylhomocysteine
387 (SAH), a key component in the methylation cycle, may be reduced in AD, suggesting impaired
388 gene regulation and DNA repair capacity⁵¹. Similarly, 2'-O-methyluridine, a modified nucleoside
389 involved in RNA metabolism, could be downregulated, pointing to stress in RNA processing and
390 contributing to cognitive impairment⁵². These RNA and epigenetic disruptions are associated with
391 decreased PACC3 scores, while the lack of compensatory upregulation suggests a minimal
392 ptau217-related counter-response.

393 Immune system alterations may limit the brain's ability to respond to AD pathology. *CLEC4A*,
394 involved in immune signaling and microglial function, may be downregulated in AD, which could
395 impair the clearance of toxic proteins like tau⁵³. On the other hand, *ARL11*, associated with cellular
396 trafficking and signaling, might be upregulated in response to cellular stress from tau accumulation,
397 reflecting a potential coping mechanism within the neuron⁵⁴. These immune-related changes link
398 *CLEC4A* downregulation to diminished PACC3 and *ARL11* upregulation to heightened ptau217
399 levels.

400 Metabolite interactions can amplify membrane instability and cognitive dysfunction.
401 decadienedioic acid (C10:2-DC)**⁵⁵ and tetradecadienoate (14:2)⁴⁶ may interact antagonistically;
402 their joint downregulation could destabilize both the structure and fluidity of membranes,
403 accelerating neurodegeneration. In contrast, sphingomyelin (d18:2/16:0, d18:1/16:1)⁵⁶ and
404 choline⁵⁷ might interact synergistically in the context of elevated ptau217, with disruptions in both
405 potentially impairing acetylcholine synthesis and amplifying cognitive decline. Together, these
406 metabolite changes show that downregulation contributes to lower PACC3 scores, while
407 upregulated or synergistically altered molecules correlate with elevated ptau217.

408

409 While this study provides valuable insights into potential biomarkers and their interactions in AD,
410 several limitations must be considered. First, the gene expression levels in this study were imputed
411 based on genotypes rather than directly measured. The accuracy of the imputation depends on how
412 precise the pre-trained models could capture true gene expression patterns, which vary strongly
413 across genes and tissues. Therefore, although this approach enabled a more comprehensive
414 analysis when measured transcriptomics data are unavailable, we acknowledge that it may
415 introduce errors. Future studies could further improve accuracy by integrating more diverse
416 reference datasets and refining the imputation models. Second, the use of the Wisconsin ADRC
417 cohort as an independent validation set for the trained models may not provide a robust assessment
418 of predictive performance due to its relatively small sample size. Utilizing larger independent
419 cohorts with characteristics similar to WRAP and Wisconsin ADRC for validation would help
420 improve the robustness of the models' predictive performance. To further deepen our
421 understanding of AD, integrating a wider array of biomarkers—such as those from proteomics—
422 would provide a more holistic view of its pathology. Furthermore, while this study focuses on
423 specific biomarkers and their pairwise interactions, AD is a complex and multifactorial disease.
424 Therefore, other factors, such as environmental influences and epigenetic changes, which were not
425 captured in this analysis, could also play critical roles in disease progression. In the future study,
426 investigating the complex interplay between genetic, environmental, and epigenetic factors will
427 offer valuable insights into the multifactorial nature of AD progression.

428 **ACKNOWLEDGMENTS**

429 The authors especially thank the WRAP and Wisconsin ADRC participants and staff for their
430 contributions to the studies. Without their efforts, this research would not be possible.

431 This study was supported by the National Institutes of Health (NIH) grants R01AG027161
432 (Wisconsin Registry for Alzheimer Prevention: Biomarkers of Preclinical AD), P30 AG062715
433 (the Wisconsin Alzheimer's Disease Research Center), and RF1AG054047 (Genomic and
434 Metabolomic Data Integration in a Longitudinal Cohort at Risk for Alzheimer's Disease.
435 Computational resources were supported by core grants from the Center for Demography and
436 Ecology (P2CHD047873) and the Center for Demography of Health and Aging (P30AG017266).

437 **CONFLICT OF INTEREST STATEMENT**

438 The authors declare no competing interests.

439 **REFERENCES**

- 440 1. DeTure, M. A. & Dickson, D. W. The neuropathological diagnosis of Alzheimer's disease. *Mol*
441 *Neurodegeneration* **14**, 32 (2019).
- 442 2. O'Brien, R. J. & Wong, P. C. Amyloid precursor protein processing and Alzheimer's disease. *Annu*
443 *Rev Neurosci* **34**, 185–204 (2011).
- 444 3. Sadigh-Eteghad, S. *et al.* Amyloid-beta: a crucial factor in Alzheimer's disease. *Med Princ Pract* **24**,
445 1–10 (2015).
- 446 4. Shankar, G. M. & Walsh, D. M. Alzheimer's disease: synaptic dysfunction and Aβeta. *Mol*
447 *Neurodegener* **4**, 48 (2009).

- 448 5. Cuddy, L. K. *et al.* A β -accelerated neurodegeneration caused by Alzheimer's-associated *ACE* variant
449 R1279Q is rescued by angiotensin system inhibition in mice. *Sci. Transl. Med.* **12**, eaaz2541 (2020).
- 450 6. Killin, L. O. J., Starr, J. M., Shiue, I. J. & Russ, T. C. Environmental risk factors for dementia: a
451 systematic review. *BMC Geriatr* **16**, 175 (2016).
- 452 7. Gatz, M. *et al.* Role of Genes and Environments for Explaining Alzheimer Disease. *Arch Gen*
453 *Psychiatry* **63**, 168 (2006).
- 454 8. Horgusluoglu, E. *et al.* Integrative metabolomics-genomics approach reveals key metabolic pathways
455 and regulators of Alzheimer's disease. *Alzheimers Dement* **18**, 1260–1278 (2022).
- 456 9. Lau, S.-F., Cao, H., Fu, A. K. Y. & Ip, N. Y. Single-nucleus transcriptome analysis reveals
457 dysregulation of angiogenic endothelial cells and neuroprotective glia in Alzheimer's disease. *Proc.*
458 *Natl. Acad. Sci. U.S.A.* **117**, 25800–25809 (2020).
- 459 10. Rajmohan, R. & Reddy, P. H. Amyloid-Beta and Phosphorylated Tau Accumulations Cause
460 Abnormalities at Synapses of Alzheimer's disease Neurons. *J Alzheimers Dis* **57**, 975–999 (2017).
- 461 11. Batra, R. *et al.* The landscape of metabolic brain alterations in Alzheimer's disease. *Alzheimer's &*
462 *Dementia* **19**, 980–998 (2023).
- 463 12. Tönnes, E. & Trushina, E. Oxidative Stress, Synaptic Dysfunction, and Alzheimer's Disease. *JAD* **57**,
464 1105–1121 (2017).
- 465 13. Chu, S. H. *et al.* Integration of Metabolomic and Other Omics Data in Population-Based Study Designs:
466 An Epidemiological Perspective. *Metabolites* **9**, 117 (2019).
- 467 14. Karch, C. M. & Goate, A. M. Alzheimer's disease risk genes and mechanisms of disease pathogenesis.
468 *Biol Psychiatry* **77**, 43–51 (2015).
- 469 15. Hasin, Y., Seldin, M. & Lusis, A. Multi-omics approaches to disease. *Genome Biol* **18**, 83 (2017).
- 470 16. Subramanian, I., Verma, S., Kumar, S., Jere, A. & Anamika, K. Multi-omics Data Integration,
471 Interpretation, and Its Application. *Bioinform Biol Insights* **14**, 1177932219899051 (2020).
- 472 17. Garg, M. *et al.* Disease prediction with multi-omics and biomarkers empowers case-control genetic
473 discoveries in the UK Biobank. *Nat Genet* **56**, 1821–1831 (2024).

- 474 18. Hoffmann, J. *et al.* Machine learning in a data-limited regime: Augmenting experiments with synthetic
475 data uncovers order in crumpled sheets. *Sci Adv* **5**, eaau6792 (2019).
- 476 19. Qiu, S. *et al.* Multimodal deep learning for Alzheimer's disease dementia assessment. *Nat Commun* **13**,
477 3404 (2022).
- 478 20. Choi, J. J. *et al.* COSIME: Cooperative multi-view integration and Scalable and Interpretable Model
479 Explainer. Preprint at <https://doi.org/10.1101/2025.01.11.632570> (2025).
- 480 21. Johnson, S. C. *et al.* The Wisconsin Registry for Alzheimer's Prevention: A review of findings and
481 current directions. *Alzheimers Dement (Amst)* **10**, 130–142 (2018).
- 482 22. Gamazon, E. R. *et al.* A gene-based association method for mapping traits using reference
483 transcriptome data. *Nat Genet* **47**, 1091–1098 (2015).
- 484 23. Jonaitis, E. M. *et al.* Measuring longitudinal cognition: Individual tests versus composites. *Alz &*
485 *Dem Diag Ass & Dis Mo* **11**, 74–84 (2019).
- 486 24. Vakil, E. & Blachstein, H. Rey auditory-verbal learning test: Structure analysis. *J. Clin. Psychol.* **49**,
487 883–890 (1993).
- 488 25. Atkinson, L. The Wechsler Memory Scale-Revised: Abnormality of selected index differences.
489 *Canadian Journal of Behavioural Science / Revue canadienne des sciences du comportement* **24**, 537–
490 539 (1992).
- 491 26. Wechsler, D. Wechsler Adult Intelligence Scale--Fourth Edition. <https://doi.org/10.1037/t15169-000>
492 (2012).
- 493 27. Monsell, S. E. *et al.* Results From the NACC Uniform Data Set Neuropsychological Battery Crosswalk
494 Study. *Alzheimer Disease & Associated Disorders* **30**, 134–139 (2016).
- 495 28. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets.
496 *Gigascience* **4**, s13742-015-0047–8 (2015).
- 497 29. Nassar, L. R. *et al.* The UCSC Genome Browser database: 2023 update. *Nucleic Acids Research* **51**,
498 D1188–D1195 (2023).
- 499 30. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience* **10**, giab008 (2021).

- 500 31. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat Genet* **48**, 1284–1287
501 (2016).
- 502 32. Fuchsberger, C., Abecasis, G. R. & Hinds, D. A. minimac2: faster genotype imputation. *Bioinformatics*
503 **31**, 782–784 (2015).
- 504 33. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature*
505 **590**, 290–299 (2021).
- 506 34. Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the expanded 1000 Genomes
507 Project cohort including 602 trios. *Cell* **185**, 3426–3440.e19 (2022).
- 508 35. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies.
509 *Bioinformatics* **26**, 2867–2873 (2010).
- 510 36. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated
511 individuals. *Genome Res.* **19**, 1655–1664 (2009).
- 512 37. Reiman, E. M. *et al.* Exceptionally low likelihood of Alzheimer’s dementia in APOE2 homozygotes
513 from a 5,000-person neuropathological study. *Nat Commun* **11**, 667 (2020).
- 514 38. Deming, Y. *et al.* Neuropathology-based APOE genetic risk score better quantifies Alzheimer’s risk.
515 *Alzheimer’s & Dementia* **19**, 3406–3416 (2023).
- 516 39. Liang, Y., Nyasimi, F. & Im, H. K. Pervasive polygenicity of complex traits inflates false positive rates
517 in transcriptome-wide association studies. Preprint at <https://doi.org/10.1101/2023.10.17.562831>
518 (2023).
- 519 40. Li, Q. S. *et al.* Association of peripheral blood DNA methylation level with Alzheimer’s disease
520 progression. *Clin Epigenet* **13**, 191 (2021).
- 521 41. de Wit, J. & Ghosh, A. Control of neural circuit formation by leucine-rich repeat proteins. *Trends*
522 *Neurosci* **37**, 539–550 (2014).
- 523 42. Ahmad, N., Lesa, K. N., Sudarmanto, A., Fakhruddin, N. & Ikawati, Z. The role of Phosphodiesterase-
524 1 and its natural product inhibitors in Alzheimer’s disease: A review. *Front Pharmacol* **13**, 1070677
525 (2022).

- 526 43. Wang, Y. *et al.* Midlife Chronological and Endocrinological Transitions in Brain Metabolism: System
527 Biology Basis for Increased Alzheimer's Risk in Female Brain. *Sci Rep* **10**, 8528 (2020).
- 528 44. Helbecque, N., Cottel, D. & Amouyel, P. Low-density lipoprotein receptor-related protein 8 gene
529 polymorphisms and dementia. *Neurobiology of Aging* **30**, 266–271 (2009).
- 530 45. Baloni, P. *et al.* Multi-Omic analyses characterize the ceramide/sphingomyelin pathway as a
531 therapeutic target in Alzheimer's disease. *Commun Biol* **5**, 1074 (2022).
- 532 46. Hammond, T. C. *et al.* Human Gray and White Matter Metabolomics to Differentiate APOE and Stage
533 Dependent Changes in Alzheimer's Disease. *J Cell Immunol* **3**, 397–412 (2021).
- 534 47. He, M. *et al.* A New Genetic Disorder in Mitochondrial Fatty Acid β -Oxidation: ACAD9 Deficiency.
535 *The American Journal of Human Genetics* **81**, 87–103 (2007).
- 536 48. Nho, K. *et al.* Serum metabolites associated with brain amyloid beta deposition, cognition and dementia
537 progression. *Brain Communications* **3**, fcab139 (2021).
- 538 49. Hudson, S. A. & Tabet, N. Acetyl-L-carnitine for dementia. *Cochrane Database of Systematic Reviews*
539 (2003) doi:10.1002/14651858.CD003158.
- 540 50. Griffin, J. W. D. & Bradshaw, P. C. Amino Acid Catabolism in Alzheimer's Disease Brain: Friend or
541 Foe? *Oxid Med Cell Longev* **2017**, 5472792 (2017).
- 542 51. Panza, F. *et al.* Possible Role of S-Adenosylmethionine, S-Adenosylhomocysteine, and
543 Polyunsaturated Fatty Acids in Predementia Syndromes and Alzheimer's Disease. *JAD* **16**, 467–470
544 (2009).
- 545 52. Zhang, X. *et al.* Small RNA modifications in Alzheimer's disease. *Neurobiol Dis* **145**, 105058 (2020).
- 546 53. Kaur, G. *et al.* DNA Methylation: A Promising Approach in Management of Alzheimer's Disease and
547 Other Neurodegenerative Disorders. *Biology* **11**, 90 (2022).
- 548 54. Yin, R.-H., Yu, J.-T. & Tan, L. The Role of SORL1 in Alzheimer's Disease. *Mol Neurobiol* **51**, 909–
549 918 (2015).
- 550 55. Schuck, P. F. *et al.* cis-4-decenoic acid provokes mitochondrial bioenergetic dysfunction in rat brain.
551 *Life Sciences* **87**, 139–146 (2010).

- 552 56. Javaid, S. *et al.* Dynamics of Choline-Containing Phospholipids in Traumatic Brain Injury and
553 Associated Comorbidities. *Int J Mol Sci* **22**, 11313 (2021).
- 554 57. Yuan, J. *et al.* Is dietary choline intake related to dementia and Alzheimer’s disease risks? Results from
555 the Framingham Heart Study. *Am J Clin Nutr* **116**, 1201–1207 (2022).
- 556

Conclusion

The integration of multi-omics data holds immense potential for unraveling the complex molecular mechanisms underlying AD. In this dissertation, I have successfully executed three ambitious aims that leverage cutting-edge techniques to address fundamental questions in AD research. The work presented here not only provides valuable insights into the molecular and genetic architecture of AD but also contributes to the advancement of precision medicine by identifying novel biomarkers and risk factors for AD.

In Aim 1, we identified cooperative TF pairs that regulate target genes in oligodendrocytes. By combining scRNA-seq and scATAC-seq data, we uncovered key TF interactions that play a pivotal role in the regulation of oligodendrocyte-specific genes. We identified several highly cooperative TF pairs, such as SOX10 and OLIG2, which are already known, and additionally discovered previously unreported cooperative pairs, such as SOX10 and NKX2.2. The use of deep learning models allowed for accurate predictions of target gene expression, revealing critical insights into how gene regulation is disrupted in AD. Additionally, the independent validation of these TF pairs through eQTLs further strengthened our findings and demonstrated the relevance of these regulatory interactions to neurodegenerative diseases like AD. This framework provides a deeper understanding of the GRNs involved in oligodendrocyte function and lays the foundation for future therapeutic strategies targeting transcriptional dysregulation in AD.

Aim 2 focused on developing COSIME (Cooperative Multi-view Integration and Scalable Interpretable Model Explainer), a novel machine learning approach designed to integrate multi-omics data and predict AD phenotypes. The integration of unsupervised neural networks with optimal transport methods proved to be an effective strategy for combining diverse data types,

such as genomics and transcriptomics, improving the accuracy of AD phenotype prediction. COSIME also provides a level of interpretability, allowing us to gain insights into the relative importance of features and their interactions within and across different omics modalities. Using COSIME, we identified that synergistic interactions between microglia and astrocyte genes associated with AD are more likely to be active at the edges of the middle temporal gyrus, as revealed through spatial transcriptomics data. We also constructed a gene-gene interaction network, which uncovered several astrocyte-specific MERFISH genes—such as *PLCB1* and *HPSE2*—as key hub genes linking multiple microglia genes. *PLCB1*, in particular, is involved in phosphatidylinositol signaling and plays a central role in inflammation and synaptic function, while *HPSE2* modulates heparan sulfate metabolism and may influence neuroinflammatory pathways in AD. This novel machine learning method has the potential to transform how we analyze complex, high-dimensional biological data, offering a robust framework for predicting disease outcomes and identifying novel AD risk factors.

In Aim 3, the application of COSIME to clinical multi-omics data from the Wisconsin Registry for Alzheimer’s Prevention (WRAP) and the Wisconsin Alzheimer’s Disease Research Center (ADRC) cohorts led to the identification and replication of key genomic and metabolomic risk factors for AD. By integrating transcriptomic and metabolomic data, we aimed to uncover how interactions between genes and metabolites influence AD phenotypes. For example, metabolite interactions such as the antagonistic relationship between sphingomyelin (d17:2/16:0, d18:2/15:0) and tetradecadienoate (14:2), which destabilize membrane structure and fluidity, could accelerate neurodegeneration. In contrast, sphingomyelin (d18:2/16:0, d18:1/16:1) and choline may interact synergistically in the presence of elevated ptau217, potentially impairing acetylcholine synthesis and amplifying cognitive decline. These metabolite changes are associated

with lower PACC3 scores when downregulated and higher ptau217 levels when upregulated. These findings have important implications for understanding the molecular interactions underlying AD, including gene-gene, metabolite-metabolite, and gene-metabolite relationships and could serve as the basis for the development of personalized treatments by enabling targeted therapies that address individual molecular profiles and specific interaction patterns unique to each patient. The ability to predict AD risk based on multi-omics data is a significant step toward precision medicine, with further work needed to validate these findings in larger, diverse cohorts and translate them into targeted therapeutic interventions tailored to the specific molecular signatures identified in high-risk individuals, as it enables the identification of individuals at high risk of developing AD before the onset of clinical symptoms.

The results of this dissertation contribute to a deeper understanding of the genetic and metabolic factors that drive AD and provide a new computational framework for integrating and analyzing multi-omics data. By identifying cooperative TFs, discovering novel biomarkers, and uncovering complex gene-metabolite interactions, this research paves the way for more targeted interventions in AD. For instance, the identification of interacting biomarkers—such as gene and metabolite biomarkers that co-regulate critical pathways in neurodegeneration—can allow for the development of combination therapies that target these pathways simultaneously, improving treatment efficacy by addressing multiple disease mechanisms at once. As we continue to improve the integration of multi-omics data, we move closer to the ultimate goal of precision medicine: providing personalized, effective treatments for individuals with AD.

Ultimately, this work demonstrates that a systems biology approach, integrating genomics, transcriptomics, and metabolomics, is crucial for advancing our understanding of AD pathology.

The application of novel computational tools like COSIME will allow for more accurate predictions, better identification of biomarkers, and the development of targeted therapies that address the complex nature of AD. This dissertation sets the stage for future research in neurodegenerative diseases, providing new directions for the diagnosis, prevention, and treatment of AD and other related cognitive disorders.

Building upon the findings and methodologies presented in this dissertation, several avenues for future research can be explored to further enhance our understanding of AD and improve the predictive capabilities of multi-omics integration models. One important direction is the exploration of high-order feature interactions, which could reveal more intricate relationships between genes, metabolites, and other molecular factors. Advanced clustering techniques could also be employed to identify subgroups within the AD population, potentially leading to the discovery of novel disease subtypes or risk profiles. Another promising development would be the incorporation of causal inference techniques into the COSIME model explainer. This would allow for a deeper understanding of the causal relationships between molecule factors in AD, thereby enabling more accurate predictions of disease progression. Finally, extending COSIME to handle longitudinal data represents a crucial next step, as it would allow for the analysis of temporal changes in multi-omics data, offering new insights into how AD evolves over time. These advancements would significantly improve our ability to predict disease onset and progression, leading to more personalized and effective therapeutic strategies.

Appendix (Publication) – CHAPTER 1: CoTF-reg reveals cooperative transcription factors in oligodendrocyte gene regulation using single-cell multi-omics

communications biology

Article



<https://doi.org/10.1038/s42003-025-07570-6>

CoTF-reg reveals cooperative transcription factors in oligodendrocyte gene regulation using single-cell multi-omics

Check for updates

Jerome J. Choi^{1,2}, John Svaren^{1,3,6} & Daifeng Wang^{1,4,5,6} ✉

Oligodendrocytes are the myelinating cells within the central nervous system, but the mechanisms by which transcription factors (TFs) cooperate for gene regulation in oligodendrocytes remain unclear. We introduce coTF-reg, an analytical framework that integrates scRNA-seq and scATAC-seq data to identify cooperative TFs co-regulating the target gene (TG). First, we identify co-binding TF pairs in the same oligodendrocyte-specific regulatory regions. Next, we train a deep learning model to predict each TG expression using the co-binding TFs' expressions. Shapley interaction scores reveal high interactions between co-binding TF pairs, such as SOX10-TCF12. Validation using oligodendrocyte eQTLs and their eGenes that are regulated by these cooperative TFs show potential regulatory roles for genetic variants. Experimental validation using ChIP-seq data confirms some cooperative TF pairs, such as SOX10-OLIG2. Prediction performance of our models is evaluated through holdout data and additional datasets, and an ablation study is also conducted. The results demonstrate stable and consistent performance.

Oligodendrocytes play key functional roles in the central nervous system (CNS) function, including that they are responsible for myelination^{1,2}. Myelination is a complex neurodevelopmental process that begins during brain development in the third trimester of pregnancy and increases steadily during childhood, but it can also be dynamically regulated in the context of learning and diseases affecting the mature CNS^{3,4}. Also, Oligodendrocyte dysfunction and myelin abnormalities have been reported in CNS disorders^{2,5,6}. Multidirectional interactions between neuronal and glial cells are required for CNS function⁷, including interactions between oligodendrocytes and neurons through myelination⁸. Therefore, it is critical to better understand the functions and roles of oligodendrocytes and myelin.

Gene expression of oligodendrocyte development from oligodendrocyte progenitor cells (OPC) is governed by complex gene regulatory mechanisms involving transcription factors (TFs)^{3,4}. TFs often work in a combinatorial fashion to regulate gene expression from regulatory elements^{9,10}. For example, some TFs such as SOX10 and OLIG2 cooperate during the induction of genes for differentiation and myelin formation^{11–14}. Enhancers can increase transcription levels from promoters and

transcription start sites (TSS), and much of the regulatory code that drives cell type-specific gene expression resides in these distal regulatory elements. Especially, some active enhancers are associated with the gene expression that characterizes cell identity and functions¹⁵. Thus, it is important to identify active oligodendrocyte-specific enhancers as well as promoters and the co-binding TFs that are responsible for their activity.

Next-generation sequencing technologies, including single-cell RNA sequencing (scRNA-seq) and the assay for transposase-accessible chromatin sequencing (scATAC-seq), have provided important insights into cell-type-specific gene regulation. Recent functional genomic resources such as PsychENCODE2¹⁶ and GTEx¹⁷, and emerging tools for integrating multi-omics data enable creating cell-type-level gene regulatory networks (GRNs) linking TFs and their binding sites (TFBS), regulatory elements to target genes (TGs). Those networks can reveal the cell-type-specific regulatory roles of TFs via regulatory elements. Moreover, additional bioinformatic tools such as SCENIC¹⁸, Signac¹⁹, and scGRNom²⁰ predict cell-type-specific gene regulatory networks to explain potential TF-TG relationships. However, most of these studies and tools focus on relationships between

¹Waisman Center, University of Wisconsin-Madison, Madison, WI, USA. ²Department of Population Health Sciences, University of Wisconsin-Madison, Madison, WI, USA. ³Department of Comparative Biosciences, School of Veterinary Medicine, University of Wisconsin-Madison, Madison, WI, USA. ⁴Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, Madison, WI, USA. ⁵Department of Computer Sciences, University of Wisconsin-Madison, Madison, WI, USA. ⁶These authors jointly supervised this work: John Svaren, Daifeng Wang. ✉e-mail: daifeng.wang@wisc.edu

individual TFs and TGs instead of TF-TF interactions and their effects on TG expression. Consequently, due to the lack of tools, the mechanistic roles of cooperative TFs in establishing cell type-specific gene regulation remain uncharacterized.

To tackle these challenges, we introduce coTF-reg, an analytical framework that integrates scRNA-seq and scATAC-seq data to identify cooperative TFs co-regulating the TG. coTF-reg identifies cooperative co-binding TFs along with active regulatory elements for gene regulation as hallmarks of active oligodendrocyte-specific regulatory elements. First, it identifies co-binding TF pairs in these regulatory regions. Second, a deep learning model is trained to predict TG expression based on the expression profiles of co-binding TFs. Third, Shapley interaction scores are computed to evaluate the interactions between TF pairs. Our findings reveal high interactions between co-binding TF pairs, such as SOX10-TCF12. Validation using oligodendrocyte eQTLs and their eGenes that are regulated by these cooperative TFs showed potential regulatory roles for genetic variants. Experimental validation using ChIP-seq data confirmed some cooperative TF pairs, such as SOX10-OLIG2 and SOX10-NKX2.2. Prediction performance of our models was evaluated through holdout data and additional datasets, and an ablation study was also conducted. The results demonstrated stable and consistent performance. Overall, our results create an analytic framework in which co-binding TF pairs cooperatively activate the TG expression through oligodendrocyte-specific regulatory elements.

Results

Deep learning and single-cell multi-omics for identifying cooperative transcription factors in oligodendrocytes

In order to predict cooperative TFs involved in oligodendrocyte gene regulation, we designed coTF-reg, which integrates scRNA-seq and scATAC-seq data to identify the cooperative TFs that co-regulate the target gene (TG) expression in oligodendrocytes (Fig. 1, Methods and Materials). Briefly, we first used scATAC-seq data with peak-to-gene links²¹. Second, among the regulatory regions for various cell types, we focused on those specific to oligodendrocytes. We then identified transcription factor binding sites

(TFBSs) and co-binding TF pairs through motif co-occurrence and co-enrichment analyses. Third, we trained deep neural networks (DNNs) to predict the expression levels of the TGs and measure interaction effects between co-binding TFs on the expression levels of TGs using gene expression from scRNA-seq data²² and computed Shapley interaction (SI)^{23,24} scores for co-binding TF pairs and found cooperative TF pairs. Fourth, we built a gene regulatory network based on SI scores for co-binding TF pairs. Lastly, as an independent validation, to validate the cooperative TF pairs we found, we mapped oligodendrocyte eQTLs onto the regulatory regions where cooperative TF pairs exist, performed Liftover analysis and co-enrichment analysis using ChIP-seq data, and applied Boolean rules to characterize the cooperativity of regulatory factors. To evaluate the prediction performance of our models, we used other publicly available datasets and conducted an ablation study by generating random TF sets to predict TG expressions.

Identification of the co-binding transcription factors in oligodendrocyte-specific regulatory regions

First, we identified a set of 787 oligodendrocyte-differentially accessible and oligodendrocyte-specific regulatory regions by comparison of oligodendrocyte scATAC-seq data to other brain cell types. In this set, we identified 958 motifs for inferred TFBSs using the JASPAR database. Second, we used co-occurrence analysis and co-enrichment analysis to identify 8101 co-binding TF pairs out of 458,403 possible TF pairs ('Methods and Materials: Co-enrichment analysis' for more details). We removed TF pairs from the same families and applied a cutoff (<0.1) for false discovery rate (FDR) yielding 8101 co-binding TF pairs. There were 206 TFs that have co-binding TFs linked to 445 TGs (Supplementary Data) that are oligodendrocyte specific in 643 regulatory regions (Fig. 2a). We annotated the regulatory regions to categorize them into promoters (32.5%) and enhancers (67.5%) (Fig. 2b).

The density plots show the distributions of the number of co-binding TF pairs, the number of TGs, and the number of peaks for individual TFs that are co-bound to other TFs. Most of the TFs have 50 to 103 co-binding TFs (median = 78). The distribution of the number of TGs for TFs is right-

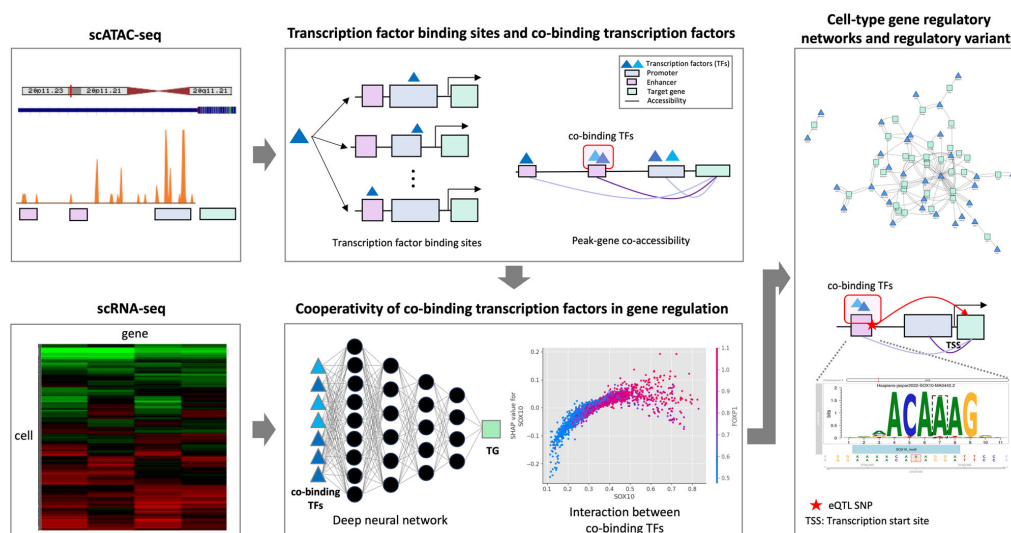


Fig. 1 | Deep learning and single-cell multi-omics for identifying cooperative transcription factors in oligodendrocytes. Inputs for the coTF-reg pipeline are scATAC-seq peak-gene links and scRNA-seq. It infers transcription factor binding sites (TFBSs) in regulatory regions and identifies co-binding TF pairs. Then, it

measures cooperativity of co-binding TFs by predicting TF-TG relationships for the levels of expression using deep learning models and Shapley interaction scores. It outputs a gene regulatory network linking co-binding TF pairs with their TGs and regulatory variants on the regulatory regions where co-TFs have their binding sites.

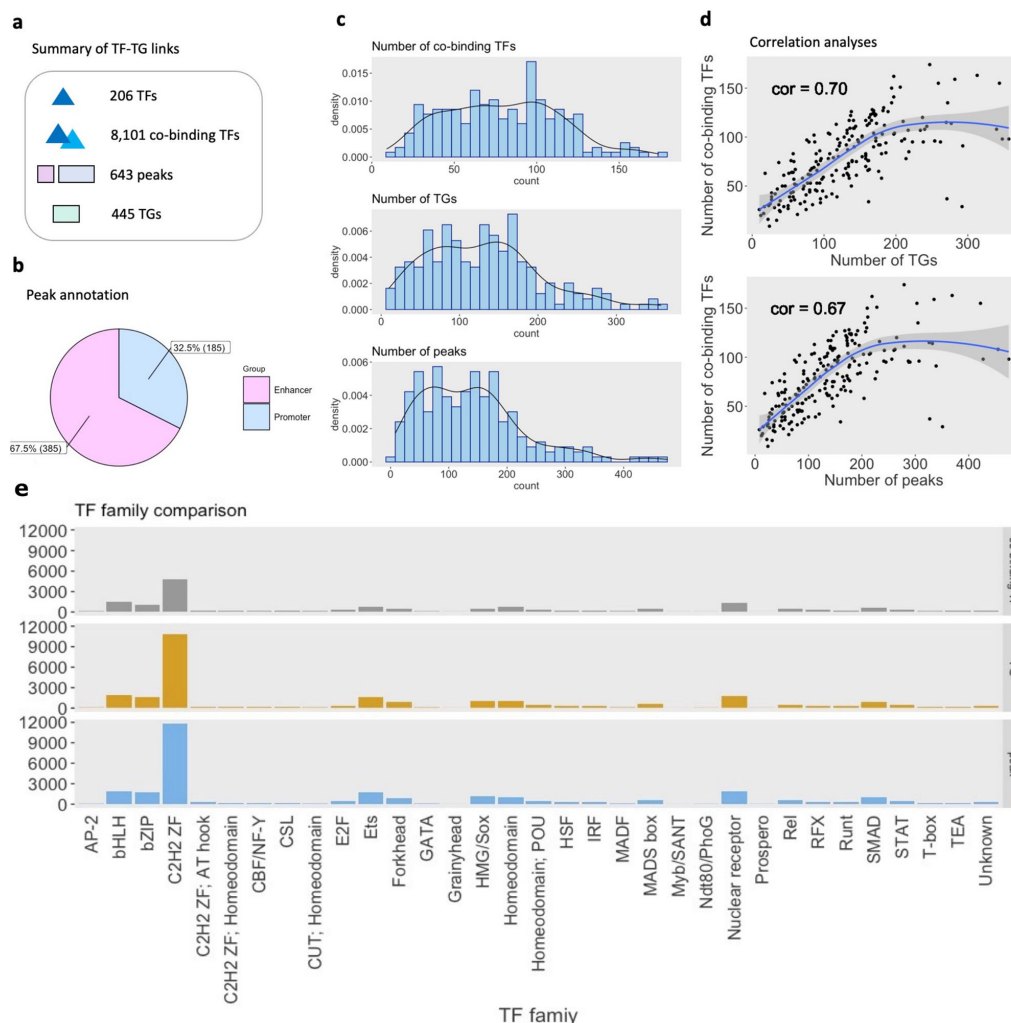


Fig. 2 | Distribution and correlation of the numbers of co-binding transcription factors, target genes, and peaks for individual transcription factors, peak annotation, and summary statistics for transcription factor-target gene links. a Summary statistics of transcription factor (TF)-target gene (TG) links. **b** Peak annotation. **c** Distributions of the numbers of co-binding TFs, TGs, and peaks for

individual transcription factors. **d** Correlations between the numbers of co-binding transcription factors and target genes, as well as the numbers of co-binding TFs and peaks. **e** Distribution of the numbers of co-binding TF pairs, TGs, and peaks for transcription factors by families.

skewed, and many TFs have 76 to 172 TGs linked. The distribution of the number of peaks for TFs is also right skewed, and the most frequent intervals were between 75 and 180 peaks (Fig. 2c). The distributions of the numbers of TGs and peaks for co-binding TF pairs are approximately normal. On average, co-binding TF pairs have 60 TGs linked and 59 peaks (Supplementary Fig. 1). Additionally, other density plots show the distributions of the number of TGs and the number of peaks for co-binding TF pairs and bar plots display the numbers of co-binding TFs, TGs, and peaks for individual TFs by their family categories (Fig. 2d). Co-binding TFs have 4 to 115 TGs (median = 59) and 4 to 123 peaks (median = 56) and the most frequent motifs are associated with TF families with C2H2 zinc finger (ZF), bZip, and bHLH DNA-binding domains (Fig. 2e). C2H2 ZF proteins are a large family and C2H2 ZF TFs (e.g., ZNF24²⁵ and KLF9/13²⁶) are known to play

significant roles in the development and function of oligodendrocytes, which are the myelinating cells of the CNS. These TFs can regulate the expression of genes essential for oligodendrocyte differentiation, survival, and myelination processes^{25,27,28}.

We computed Pearson correlation coefficient (r) to measure correlations between the number of co-binding TFs and the number of TGs and the number of co-binding TFs and the number of peaks for individual TFs (Fig. 2d). The number of co-binding TFs and the number of TGs for individual TFs are strongly positively correlated ($r = 0.70$). It suggests that TFs that are co-bound to other TFs tend to have more TGs linked to them. The number of co-binding TFs and the number of peaks for individual TFs are also strongly positively correlated ($r = 0.67$). It shows that co-binding TFs may exist in many different peaks.

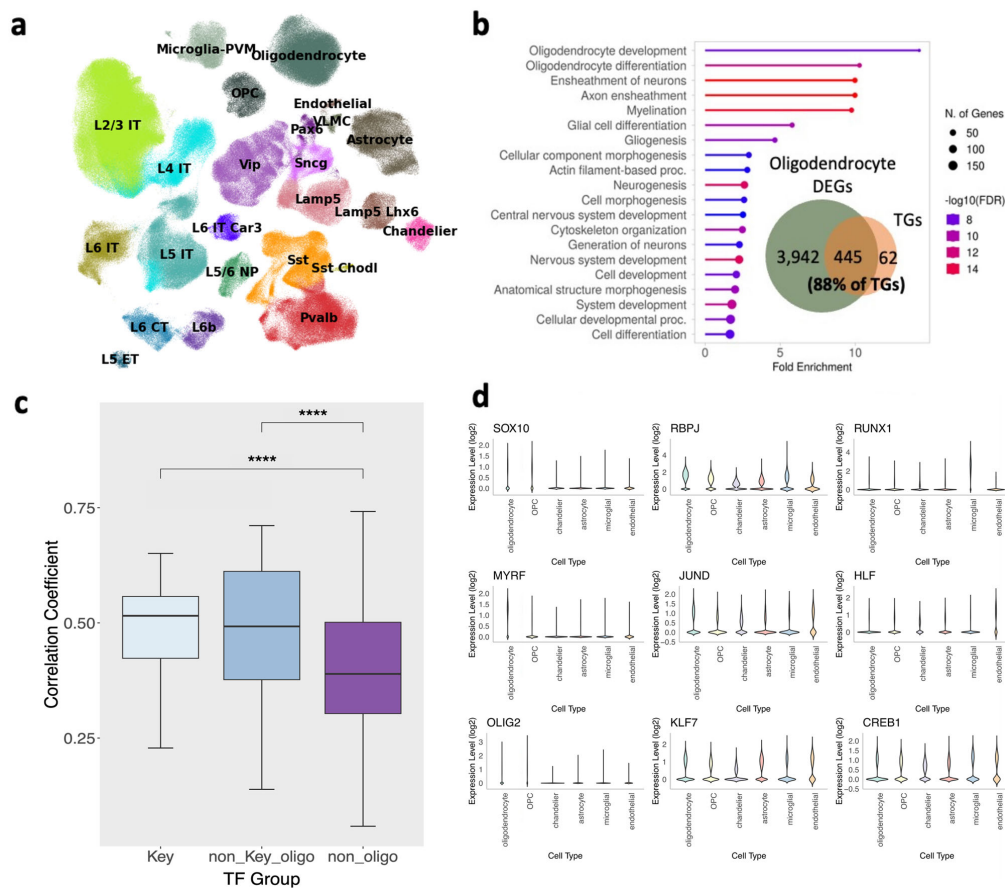


Fig. 3 | Oligodendrocyte gene expression relationships between transcription factors and target genes. **a** UMAP for eighteen cell type in middle temporal gyrus region, **b** Enrichments analysis for target genes that are oligodendrocyte-specific, **c** Pairwise two-sided t-tests for the correlation (between the expression of the TFs and

their TGs) comparison for three TF groups: 10 oligodendrocyte key transcription factors, 83 oligodendrocyte-specific non-key transcription factors, and 103 non-oligodendrocyte-specific transcription factors, and **d** Boxplots for the expression levels of the three categories in (c) (each column is an example for each category).

In the following sections, we incorporate RNA-seq data to explore gene expression relationships between TFs and TGs, train deep learning models to predict TG expression using co-binding TFs, and compute TF interaction scores using the trained models.

Oligodendrocyte gene expression relationships between transcription factors and target genes

A single cell study identified the unique gene expression profile of oligodendrocytes compared to other brain cell types²², as shown by the two dimensional Uniform Manifold Approximation and Projection (UMAP) space after computing latent representations of the neighborhood graph (Fig. 3a). The UMAP embeddings reveal that oligodendrocytes exhibit a distinct expression profile compared to other cell types. This separation suggests that oligodendrocytes have unique transcriptional programs that differentiate them from neighboring cell types. The distinct clustering of oligodendrocytes in the UMAP space indicates specialized functional roles and may reflect their involvement in myelination and maintenance of neural integrity. In order to focus on oligodendrocyte-specific mechanisms of gene regulation, we conducted differential expression testing using 17,946

genes and 20,191 metacells and identified 4387 differentially expressed genes (DEGs) for oligodendrocytes. We found 445 TGs out of 507 TGs of oligodendrocyte-specific regulatory elements (88%) were DEGs for oligodendrocytes. Subsequently, we conducted enrichment analysis for these 445 TGs revealing their involvement in crucial biological processes for oligodendrocytes, such as oligodendrocyte development, oligodendrocyte differentiation, and myelination (Fig. 3b).

We categorized TFs into oligodendrocyte key TFs, oligodendrocyte-specific non-key TFs, and non-oligodendrocyte-specific TFs using oligodendrocyte expression level and the list of key TFs (see 'Methods and Materials: Key TFs' for more details). 'Oligodendrocyte-specific key TFs' are oligodendrocyte differentially expressed TFs and key TFs, 'oligodendrocyte-specific non-key TFs' are oligodendrocyte differentially expressed TFs but not key TFs, and 'non-oligodendrocyte-specific TFs' are neither oligodendrocyte differentially expressed TFs nor key TFs. The key oligodendrocyte TFs were defined based on mouse loss-of-function studies that have shown that specific TFs are critical for oligodendrocyte differentiation. The key TFs include SOX10²⁹, SOX2^{30,31}, SOX8³², MYRF³³, OLIG1³⁴, OLIG2³⁵, TCF7L2^{36,37}, ZNF24²⁵, NKX2.2³⁸, and NKX6.2³⁹.

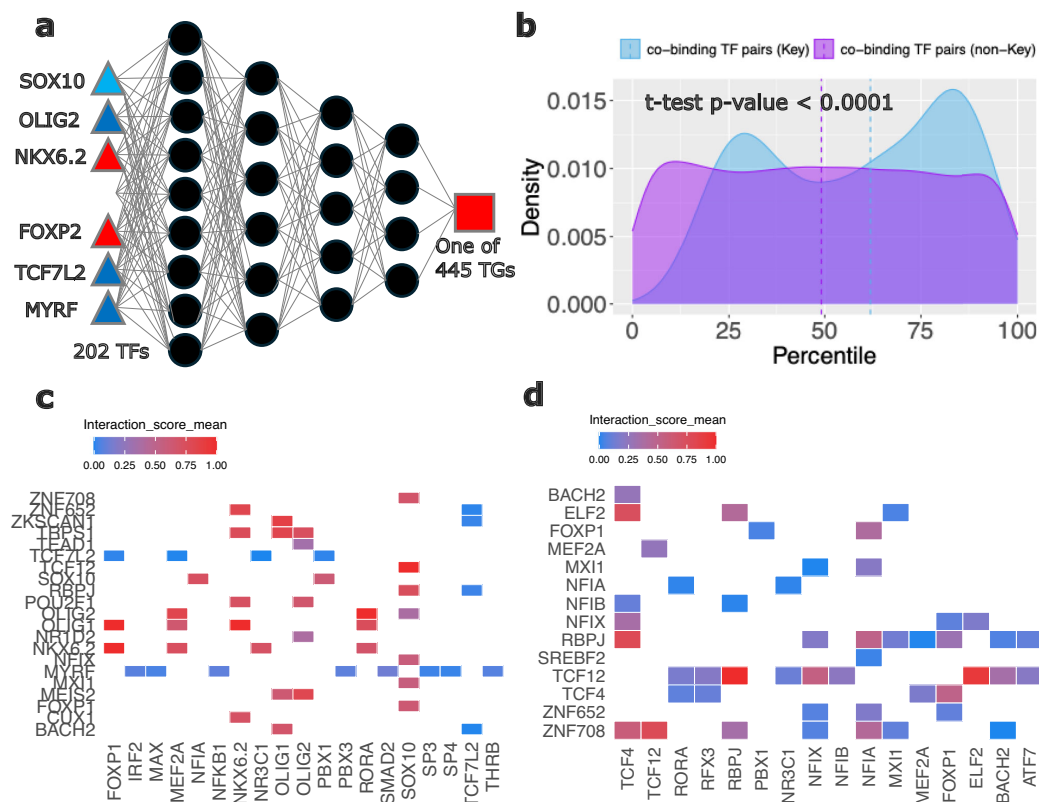


Fig. 4 | Cooperative transcription factor pairs by Shapley interaction scores. **a** Deep learning architecture, **b** Percentile distribution of interaction scores for 577 key-transcription factor pairs and 7029 non-key transcription factor pairs, **c** Top

forty-eight SI interaction scores for key transcription factor pairs, and **d** Top forty-eight SI interaction scores for non-key transcription factor pairs.

Each of 206 TF, who have co-binding TFs, regulates a different set of TGs, and we computed correlations between the expression of the TFs and their TGs in the three categories. Pairwise two-sided t-tests show that correlations between TFs and TGs in oligodendrocyte key TF pairs and those in non-oligodendrocyte-specific TF pairs are significantly different ($p < 0.001$). It also indicates that correlations between TFs and TGs in oligodendrocyte-specific non-key TF pairs and those in non-oligodendrocyte-specific TF pairs are significantly different ($p < 0.001$) (Fig. 3c). The results for differential expression testing show that the six TFs in the two categories, oligodendrocyte key TF pairs and oligodendrocyte-specific non-key TF pairs are all significant and up-regulated (Fig. S8).

We color-coded the UMAP embeddings based on the expression level of the TFs (Fig. 3a) and selected three TFs as examples for each category. Oligodendrocyte-specific key TFs such as SOX10, MYRF, and OLIG2 are specifically highly expressed in oligodendrocytes. Oligodendrocyte-enriched non-key TFs, including RBPJ, JUND, and KLF7, are expressed in multiple cell types but are more highly expressed in oligodendrocytes. Non-oligodendrocyte-specific TFs, such as RUNX1, HLF, and CREB1, are not specifically expressed in oligodendrocytes (Fig. 3d).

Deep learning and Shapley interaction scores to measure cooperativity of co-binding transcription factors

To understand the complex relationships between TFs for predicting TGs, we built deep learning models. We trained a deep learning model for each of the 445 TG. Each model used the expression levels for the 206 TFs that have

co-binding TFs to predict a TG expression level. We used seven hidden layers in each DNN (Fig. 4a). We excluded co-binding TF-TG pairs that exhibited high variability in their SI scores (coefficient of variance > 0.5). Using a trained model and a hold-out test dataset, we computed SI scores for TFs in each DNN. Additionally, we determined the percentile SI score for all co-binding TF pairs. Then, a two-sided t-test to compare the mean values for the percentile SI scores of key co-binding TF pairs and non-key co-binding TF pairs revealed a significant difference between the two groups ($p < 0.0001$) (Fig. 4b).

To emphasize the several important key co-binding TF pairs, we selected the top forty-eight interacting pairs for each key co-binding TF pair, such as SOX10, MYRF, OLIG1, OLIG2, NKX6.2, and TCF7L2, and generated a heatmap for their SI scores scaled from 0 to 1 (Fig. 4c). Similarly, we chose the top forty-eight interacting co-binding TF pairs for non-key TFs and created another heatmap for their SI scores scaled from 0 to 1 (Fig. 4d). We noticed that the SI scores for key-TF co-binding pairs have higher values than those for non-key co-binding TF pairs.

We also validated our model prediction performance for one TG, myelin basic protein (*MBP*), using additional data (Supplementary Fig. 2)⁴⁰. We regressed the scaled actual values on the scaled predicted values. For our primary dataset, we obtained an R-squared of 0.81 and a r of 0.90 (Supplementary Fig. 2a). Furthermore, when analyzing another dataset, we observed an R-squared of 0.69 and a r of 0.83, affirming the predictive capability of our model architecture (Supplementary Fig. 2b).

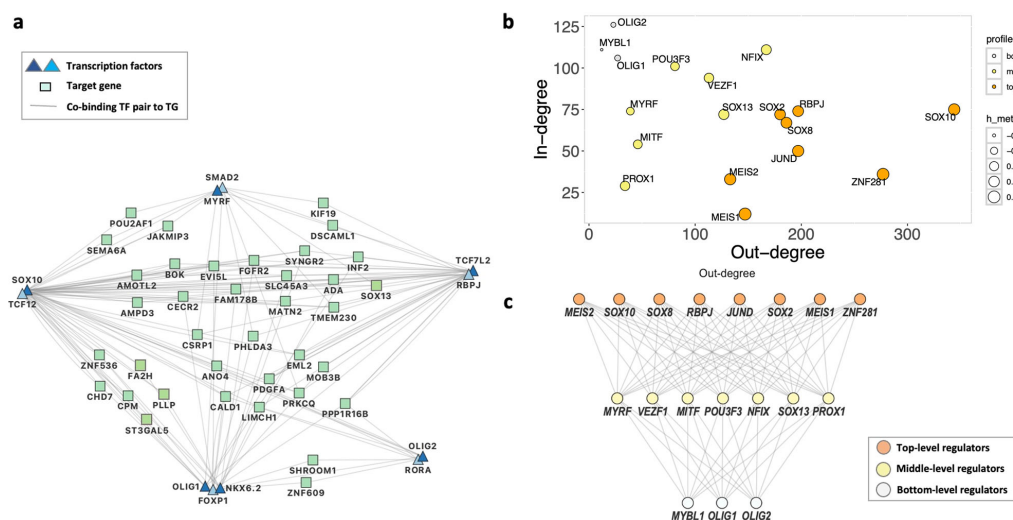


Fig. 5 | Gene regulatory network and transcription factor hierarchy. **a** Gene regulatory network for six key cooperative transcription factor pairs with the highest interaction scores, **b** A plot of in-degree (I) vs. out-degree (O) for the transcription factors that have I and O in the gene regulatory network, and **c** Transcription factor

hierarchy. Each node depicts a transcription factor. In (b, c), the edges colored in orange are the top-level (master) regulators, the edges colored in yellow are the middle-level regulators, and the edges colored in white are the bottom-level regulators.

Oligodendrocyte gene regulatory network analysis for cooperative TF pairs and transcription factor hierarchy

We chose one pair of co-binding TF with the highest interaction scores from six key co-binding TF pairs, including SOX10, MYRF, OLIG1, OLIG2, NKX6.2, and TCF7L2. We built a gene regulatory network (GRN) for these cooperative TF pairs and their TGs that are co-regulated by them (Fig. 5a). We found that a TG, *CALD1*, is co-regulated by three key cooperative TF pairs, SOX10-TCF12, RORA-OLIG2, and FOXP1-NKX6.2 and another TG, *PPP1R16B*, is co-regulated by three key cooperative TF pairs, RORA-OLIG2, FOXP1-NKX6.2, and FOXP1-OLIG1. There are other TGs, such as *AMOTL2*, *BOK*, *CALD1*, *FA2H*, and *CPM*, in the GRN that are co-regulated by two pairs of cooperative TFs.

We computed in-degree and out-degree for eighteen TFs that can also be TGs at the same time since TF feed forward and feedback loops are common (Fig. 5b). Then, we conducted a TF hierarchy analysis and found eight top-level regulators (Fig. 5c), called ‘Master regulators’, including SOX10, SOX2, and SOX8, which are key TFs that are known to play critical roles in oligodendrocyte differentiation^{30–32}. The other five master regulators, MEIS1, MEIS2, RBPJ, JUND, and ZNF281, are categorized as oligodendrocyte-specific non-key TFs in Fig. 3c. All TFs that are middle-level regulators and bottom-level regulators, except for MYRF, are categorized as oligodendrocyte-specific non-key TFs. MYRF is one of the key TFs which is specifically activated in myelinating oligodendrocytes. PROX1 has been identified as being important for oligodendrocyte differentiation^{41,42}. Most of these eighteen TFs are expressed in both oligodendrocytes and OPCs (Supplementary Fig. 4). It provides evidence that oligodendrocyte differentiation is pre-set in OPCs⁴³.

Independent validation for cooperative TFs

eQTL mapping. As an independent assessment of the regulatory regions we mapped oligodendrocyte eQTLs⁴⁴ onto oligodendrocyte-specific regulatory regions to explain the causal relationships between the expression levels of the co-binding TF pairs we identified and their target genes (TGs). Using chromosome and position of eQTL SNPs (eSNPs) from oligodendrocyte eQTLs, eSNPs integrated with a total of 643 oligodendrocyte-specific regulatory regions (Fig. 2a). This integration

facilitates the identification of potential regulatory connections between the eSNPs and the co-binding TFs in these regions, enhancing our understanding of how genetic variations influence the expression levels of the identified co-binding TF pairs and their corresponding TGs. Notably, it provides evidence of causation if the eQTL genes and TGs are identical where co-binding TF pairs occur, indicating that these co-binding TF pairs are co-regulating TG expressions.

First, among 4.8 million oligodendrocyte eQTLs, we filtered 2 million significant ($FDR < 0.05$) eQTLs. Second, we mapped these significant eSNPs onto oligodendrocyte-specific regulatory regions (Fig. 6a). In total, 383 eSNPs and 159 eGenes were mapped onto 188 regulatory regions. Among these, 373 eSNPs and 153 eGenes (and TGs) were found in 179 regulatory regions associated with key TF pairs. Enrichment analysis for TGs indicates their strong involvement in biological processes such as oligodendrocyte development, myelination, and oligodendrocyte differentiation. (Fig. 6b)

Validation of cooperative TF pairs. The model generated from human epigenome and expression data predicted a number of enriched TF pairs within oligodendrocyte-specific TF regulatory elements. In order to test if the coordination occurs as predicted, we utilized rat oligodendrocyte ChIP-seq data that were available for selected transcription factors. One predicted pair was OLIG2/SOX10, which had previously been shown to be extensively colocalized in analyses of rat oligodendrocytes⁴⁵. To visualize the preferential binding of SOX10 on a global scale, a read density plot for SOX10 ChIP-seq reads⁴¹ was generated centered on the previously defined OLIG2 peaks⁴⁵ in oligodendrocytes (Fig. 6c). In line with previous analysis, the average read density of SOX10 is highly enriched over OLIG2 bound sites. A newly found pair predicted by the model was that of NKX2.2 and SOX10, and we generated a similar plot of SOX10 ChIP-seq reads over a defined set of NKX2.2 ChIP-seq peaks in oligodendrocytes⁴⁶, and we found a similarly high enrichment of SOX10 binding on ~40% of NKX2.2 binding sites (Fig. 6d). An example of the colocalization is shown for the *MBP* gene, which *MBP* is a crucial TG in oligodendrocytes as a key component of the myelin sheath^{47,48}. Expression of *MBP* is essential for the differentiation and maturation of oligodendrocytes^{49,50}, and *MBP* maintains the structure and integrity of

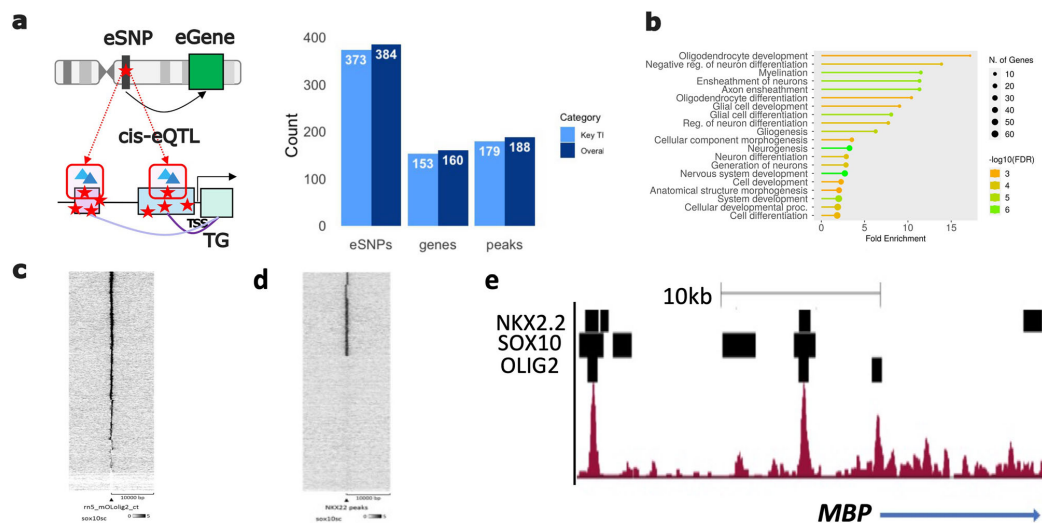


Fig. 6 | Independent validation using eQTLs and ChIP-seq data. **a** eQTL mapping onto oligodendrocytes regulatory regions, **b** ChIP-seq peak for cooperative TF pairs that co-regulate *MBP*. The ChIP-seq profile is for SOX10, and solid blocks indicate peaks for the specified transcription factors. **c** Gene ontology enrichment

analysis of target genes associated with oligodendrocyte eQTL's, and **d,e** Heatmaps show distribution of SOX10 ChIP-seq reads centered on the previously defined OLIG2 and NKX2.2 sites.

the myelin sheath⁵¹. As shown in Fig. 6e, there are at least 2 sites upstream of *MBP* where there is colocalization of SOX10 with NKX2.2 and OLIG2.

Boolean cooperativity of TF pairs. We applied a logic circuit to characterize Boolean cooperativity of TFs using Loregic⁵². A total of 206 TFs that form 8101 co-binding TF pairs were input. 6660 (82.2%) out of the 8101 co-binding TF pairs have consistent triplets—matching the same logic gate across all targets, demonstrating strong cooperation between the activities of the two TFs on the TGs. More than half of the TF1-TF2-TG pairs are categorized as “AND” indicating a positive correlation between TG expression and the expression of both TF1 and TF2 (Fig. S6a). We also achieved permutation scores to remove logic gates chosen by random. Still, 6092 TF pairs have consistent triplets and 64% of triplets are categorized as “AND” (Fig. S6b).

Independent validation for the prediction performance of the models

Model prediction validation and ablation study. Using Multi-omics scRNA-seq data²¹ from the same cells as the scATAC-seq data in the main analysis, we trained deep learning models and computed SI scores. Forty-eight SI scores for key-TF pairs in Fig. 4c were selected. The correlation between the SI scores computed from the main data and the Multi-omics data are shown in Fig. S5a. We also ran a two-sided t-test to compare the mean values for the percentile SI scores of key co-binding TF pairs and non-key co-binding TF pairs as we did for the main data (Fig. 4b). There was a significant difference between the two groups ($p < 0.0001$) (Fig. S5b).

Model performance was evaluated using the holdout data. Additionally, we included three more publicly available scRNA-seq datasets: Multi-omics, ROSMAP⁴⁰, and Cross-disorder⁵³, and validated the prediction performance of our model for each TG. Here, TGs were predicted using the trained models and the entire datasets. The holdout data, Multi-omics, and ROSMAP show consistently low normalized root mean squared error (NRMSE), while more than 75% of predictions in Cross-disorder also have low NRMSE. NRMSE can be compared across genes (Fig. S7a).

We also conducted an ablation study to compare the prediction performance of our models. Another dataset for 206 random TFs that are neither co-binding nor cooperative was created and their prediction performance was compared to that of 206 co-binding TFs (Fig. S7b). The model prediction performance is much better overall when 206 TFs used for predicting TGs are either co-binding, cooperative, or both.

Discussion

With resources provided by advances in single-cell sequencing, some studies^{54–57} have elucidated the roles of several TFs, enabling the construction of cell type-specific gene regulatory networks to explain potential TF-TG relationships using bioinformatic tools. However, most of these studies and tools primarily focus on relationships between independent TFs and TGs.

This study introduces an analytical framework, coTF-reg, which identifies co-binding TFs and their TGs in oligodendrocyte-specific regulatory regions. Deep learning models predicted TG expression levels using the expression levels of co-binding TF pairs, and we computed TF SI scores to define highly interacting co-binding TF pairs as ‘cooperative’ TFs that co-regulate TG expression levels. We found that the key co-binding TF pairs tend to highly interact with each other compared to non-key co-binding TF pairs for predicting TG expression levels. Independent validation, such as mapping eQTLs onto the regulatory regions, provides evidence for causal relationships between co-binding TF pairs and TGs. Additionally, converting these regions to the rat genome assembly coordinates and measuring the density of ChIP-seq signals for key cooperative TFs show that many of these TF pairs are enriched in the regulatory regions, indicating their collaborative role in co-regulating TG expression levels. We defined specific key TFs and examined co-binding TF pairs containing them, along with their interactions in predicting TG expression levels. We then compared these results with those of non-key TF pairs. Overall, co-binding TF pairs with known regulators of oligodendrocyte development exhibit higher SI scores, suggesting that they not only regulate TG expression individually but also cooperatively. We identified several highly cooperative TF pairs, such as SOX10 and OLIG2^{12,58}, which are already known. Additionally, we

discovered previously unreported cooperative pairs, such as SOX10 and NKX2.2.

Our study demonstrates several strengths. First, we concentrate on interactions between co-binding TF pairs and their impact on TG expression using deep learning approaches. Deep learning can elucidate complex TF relationships and their effects on TG expression levels. Second, the coTF-reg pipeline can be used by general users with any scATAC-seq and scRNA-seq data. The code for coTF-reg is openly available on GitHub, allowing users to input their scATAC-seq and scRNA-seq data for specific purposes. Third, we provide a comprehensive analytical framework that incorporates analyses utilizing co-bindings by motif and expression levels. We define 'cooperative' TF pairs as TF pairs significantly co-enriched across regulatory regions, exhibiting high SI scores in terms of expression when predicting TG expression. The term cooperativity has often been applied to co-bindings of TFs to nearby sites that facilitates stabilized binding due to protein-protein interactions, but in our model, we use TF pairs that can bind to sites in the same regulatory regions, since TFs can coordinately activate enhancers without direct interactions.

Nevertheless, there are some limitations to our study. To begin with, it's important to note that more than two TFs can co-regulate TG expression^{59,60}. However, our current tool is limited to analyzing interactions between two co-binding TFs. In future research, developing or applying more sophisticated methods capable of handling clusters of TFs that co-regulate the same TG expression will be informative. Moreover, our method for identifying binding sites relies on the position frequency matrices in the motif database. While both SOX10 and MYRF are key TFs for oligodendrocytes, we encountered difficulty in obtaining sufficient binding sites for MYRF. Consequently, we had to supplement with a different motif for MYRF based on our prior knowledge. More generally, the definition of TF motifs relies on disparate methods, and limitations of motif generation and analysis have been noted previously. Nonetheless, our analysis provided TF-TF coordination that we could validate using data from previous studies. We predict that future analysis can be used to determine if the predicted TF pairing plays a role in oligodendrocyte differentiation, since reliance on single factor studies is not able to recapitulate the important combinatorial functions of TFs in generating cell type-specific gene expression patterns. Lastly, there can be alternative methods for establishing cooperative relationships between TFs, such as Boolean rules⁶¹⁻⁶⁵. Logic-based models are also powerful tools for understanding the complex interactions among regulatory TFs in gene regulation. Developing new tools that incorporate Boolean rules and machine learning approaches will help us effectively infer more intricate TF relationships, paving the way for future research aimed at unraveling the complexities of gene regulation.

Methods

coTF-reg pipeline workflow

First, published scATAC-seq data with peak-to-gene links²¹ is inputted into the coTF-reg pipeline. Second, transcription factor binding sites (TFBSs) and co-binding TF pairs in the oligodendrocyte regulatory regions are identified through motif co-occurrence and co-enrichment analyses. Third, deep neural networks (DNNs) to predict the expression levels of the TGs are trained and the interaction effects between co-binding TFs on the expression levels of TGs using gene expression from scRNA-seq data²² are measured by computing Shapley interaction (SI)^{23,24} scores. Fourth, a gene regulatory network is built based on SI scores for co-binding TF pairs. Fifth, a TF hierarchy analysis is used to define TFs as regulators in three categories. Lastly, as an independent validation, to validate the cooperative TF pairs: 1. The oligodendrocyte eQTLs are mapped onto the regulatory regions where cooperative TF pairs exist, 2. Liftover analysis and co-enrichment analysis using ChIP-seq data are conducted, 3. Boolean rules are applied to characterize the cooperativity of regulatory factors. To evaluate the prediction performance of our models: 1. Other publicly available datasets are used as validation data to

predict TG expressions, 2. Ablation study is implemented by generating random TF sets to predict TG expressions.

Step 1: Infer transcription factor binding sites. We inferred transcription factor binding sites (TFBSs) in 787 scATAC-seq peak regions that have linkages with TGs.

- The R package *GenomicRanges* was used to format the ATAC-seq peaks into genomic ranges.
- Position frequency matrices (PFMs) for the 949 motifs in *JASPAR2022* database⁶⁶ were set in R, along with nine additional PFMs for the important modified motifs based on our prior knowledge.
- TFBSs in the scATAC-seq peak regions were inferred using a R package, *motifmatchr*⁶⁷.

Step 2: Identify co-binding transcription factor pairs. We identified co-binding TF pairs using the inferred TFBSs in Step 1.

- All possible TF-TF pairs with binding sites in the scATAC-seq peak regions were considered.
- TF pairs from the same families were excluded.
- Co-enrichment analysis: Co-occurrence analysis was conducted to find TF pairs that have overlapping regions. We then conducted hypergeometric tests to find significantly enriched TF pairs in the same regions. We used multiple testing corrections via FDR and applied FDR < 0.1 cutoff. We define the TF pairs that are co-enriched (FDR < 0.1) as 'co-binding' TF pairs.
- Gene regulatory networks (GRNs) were constructed for TG-co-TF pair-peak links and matched TGs and co-TF pairs to the scRNA-seq data.
- Lowly expressed TGs and TFs were removed from the GRNs by applying a cutoff, median expression level > 1; more than half of the cells are expressed, from the GRNs.
- Differential expression testing was implemented using *Seurat*⁶⁸ and selected TGs that are oligodendrocyte specific in the GRNs.
- Peaks were annotated as promoters or enhancers using *annotatr*⁶⁹.

Step 3: Measure cooperativity of co-binding transcription factors. Gene expression levels of the co-binding TF pairs from scRNA-seq data were incorporated into deep learning models to predict the expression levels of the TGs and measure interaction effects between co-binding TFs on the expression levels of TGs using Shapley interaction (SI) scores.

- Metacells for the cells in scRNA-seq data were projected using a Python package, *metacells*⁷⁰.
- Expression levels of TFs that have co-binding TFs and TGs were used to construct deep learning models for each TG using *PyTorch*⁷¹ in Python.
- SI scores for TF pairs were computed in each deep learning model.
- Interaction matrices for the SI scores were generated in deep learning models and the mean interaction scores for co-binding TF pairs were calculated.
- Coefficients of variation (CV)⁷² of the interaction scores for each co-binding TF pair were computed and the pairs with CV values higher than 0.5 were removed.

Step 4: Gene regulatory network and TF hierarchy analysis. A gene regulatory network was built for six key cooperative TFs.

- One cooperative TF pair for each of the six key TFs was selected based on the top interaction scores.
- A gene regulatory network was built linking cooperative TF pairs to TGs.
- TGs co-regulated by cooperative TF pairs were selected.
- A network plot was generated using *Cytoscape*⁷³.

Step 5: TF hierarchy analysis. TFs that can be TGs were chosen, and we implemented hierarchy analysis²⁴ for those TFs.

- In-degree (I) and out-degree (O) for the TFs were calculated.
- Hierarchy height metrics for the TFs were computed.
- TFs were classified as top-regulator, middle-regulator, or bottom-regulator.

Step 6: Independent validation. We implemented eQTL mapping, ChIP-seq enrichment analysis, and Boolean cooperativity analysis for validating cooperative TF pairs and model prediction validation and an ablation study for validating the prediction performance of the models.

Validation of cooperative TF pairs

eQTL mapping. We mapped the significant (FDR < 0.05) oligodendrocyte eQTLs onto the scATAC-seq peak regions.

- Publicly available oligodendrocyte eQTL data⁴⁴ were downloaded and the significant (FDR < 0.05) eQTLs were extracted.
- The significant eQTLs were mapped to the scATAC-seq peak regions in the GRNs.
- The results were verified by comparing the number of eQTLs mapped onto the peak regions for key-TF pairs and non-key-TF pairs.

ChIP-seq enrichment analysis

We performed the LiftOver analysis to convert genome coordinates for rat to human hg38 assembly using UCSC Genome Browser⁵.

- Genome coordinates for human hg38 assembly were converted to the rn5 rat genome coordinates for human (hg19) assembly.
- Overlapping genome coordinates between conserved (from hg38 to rn5) assembly and the regulatory regions in the GRN were identified.
- Cooperative TF pairs in the overlapping regions identified, along with the TGs they co-regulate.

Using the results from the LiftOver analysis, we tried to find signals in co-enriched binding sites for cooperative key TF pairs in rat oligodendrocyte ChIP-seq data. Heatmaps were created via EAs⁶. ChIP-seq tracks were visualized using UCSC genome browser. Previous ChIP-seq datasets for SOX10, OLIG2, and NKX2.2 are available at GEO accession numbers: GSE64703, GSE42447 and GSM1906296.

Boolean cooperativity of TF pairs

We applied a logic circuit to characterize Boolean cooperativity of TFs using Loregic³². Loregic is a computational tool, integrating gene expression and regulatory network, to characterize the cooperativity of regulatory factors. It uses 16 possible two-input-one-output logic gates (e.g. AND) to describe triplets of two factors regulating a common target. The GRN was inputted including co-binding TFs-TG links. Then, we binarized the gene expression levels to Boolean values 1 and 0 to represent high and low gene expression, respectively, using BoolNet²⁷. BoolNet assigned Boolean values to expression data on the basis of modular co-expression patterns by K-means clustering across inputted samples and therefore accounts for differences in the dynamic ranges of expression among genes in the input data. The triplet gene expression data was extracted and matched to all possible logic gates. We selected consistent logic gates. We also ran 100 permutation tests to find significant logic gates.

Validation of the prediction performance

Model prediction validation. To verify the performance of deep learning model architectures, we trained a deep learning model for predicting a TG, *MBP* using another data⁴⁰. The trained model was used to predict the expression level of *MBP* and compared the results with the model for *MBP* using the main data.

Using Multi-omics scRNA-seq data²¹ from the same cells as the scATAC-seq data in the main analysis, we trained deep learning models and computed SI scores, following the same processes we did in coTF-reg pipeline for identifying cooperative TFs in oligodendrocyte gene regulation ('Step 2 Measure cooperativity of co-binding TFs') for the main scRNA-seq data.

Model performance was evaluated using the SEA-AD²³ holdout data. We also include three more publicly available scRNA-seq datasets: Multi-omics²¹, ROSMAP⁴⁰, and Cross-disorder⁵³, and validate the prediction performance of our model for each TG. Here, TGs were predicted using the trained models and the entire datasets. Normalized root mean squared error (NRMSE) is used to compare the performance across different datasets.

Ablation study

It is important to assess whether the 206 co-binding TFs effectively predict their TGs. Another dataset with 206 random TFs that are neither co-binding nor cooperative was generated to evaluate the prediction performance of our models. We used our trained models to predict holdout data for random TFs and compared their prediction performance to that of 206 co-binding TFs.

Single-cell ATAC-seq data

Chromatin accessibility data²¹ was used for the main analyses. Brain samples were selected and eight thousand nuclei from each sample were subjected to the Chromium Next GEM Single-Cell Multiome ATAC-seq. We filtered oligodendrocyte-specific peak-gene links for our analyses. 930 peaks and 606 genes were initially chosen.

Single-cell RNA-seq data

SEA-AD (Main analysis). The data for the whole taxonomy collected from dorsolateral prefrontal cortex (1,395,601 cells) were downloaded through the Open Data Registry on AWS as AnnData objects (h5ad format)³². The cells for disease were excluded and only the controls were retained. Then, we projected metacells for the whole taxonomy and found 2004 metacells and 17,946 genes for oligodendrocytes.

Multi-omics. The normalized and quality controlled data was gained from the CELLxGENE (RRID:SCR_021059) portal. Brain samples were selected and eight thousand nuclei from each sample were subjected to the Gene Expression protocol (10x Genomics). We filtered 5459 cells for oligodendrocyte.

ROSMAP. The processed count matrix for oligodendrocyte was downloaded from a supplementary website for 'Single-cell atlas reveals correlates of high cognitive function, dementia, and resilience to Alzheimer's disease pathology'⁴⁰. We projected metacells for the controls only and found 7072 metacells and 16,707 genes.

Cross-disorder. Post quality control filtered data was obtained from the CELLxGENE portal. We projected metacells for oligodendrocyte controls and found 1004 metacells and 21,248 genes for oligodendrocytes.

Uniform manifold approximation and projection for dimension reduction

We gained scRNA-seq data for the whole taxonomy collected from dorsolateral prefrontal cortex through the Open Data Registry on AWS as AnnData objects (h5ad format)³². There were 1,395,601 cells across 18 sub-cell types. A total of 18,431 hg38 protein-coding genes, obtained via BioMart³, were selected from 36,517 genes. We normalized the data to a depth of 10,000 and log1 transformed it using Scanpy²⁹ in Python. Then, the highly variable genes (HVGs) were identified using dispersion-based methods³⁰ to normalize dispersion, obtained by scaling with the mean and standard deviation of the dispersions for genes falling into a given bin for mean expression of genes. The cutoffs for the mean dispersions for genes were a minimum of 0.0125 and a maximum of 3, and for the minimum dispersion was 0.5. We identified 3032 HVGs and scaled each gene to unit

variance to clip values exceeding standard deviation of 10. To reduce the dimensionality of the data, we ran principal component analysis and used top 30 PCs to compute the neighborhood graph of the cells. Finally, we embedded the neighborhood graph with 20 neighbors in two dimensions using Uniform Manifold Approximation and Projection for Dimension Reduction (UMAP)⁸¹.

Differential expression testing

We inputted metacells for all cell types to identify oligodendrocyte-specific genes using Seurat v4 in R. We used the Poisson likelihood ratio test in FindMarkers function assuming that gene expression follows the negative binomial distribution. Oligodendrocytes and oligodendrocyte precursor cells (OPCs), and astrocytes that are known as major cell types among glia in the CNS were grouped and the other fifteen cell types were compared. We used a cutoff, FDR (<0.05) to select differentially expressed genes in the oligodendrocyte group.

Position Frequency Matrices

Position frequency matrices (PFMs) for the 949 motifs in JASPAR2022 were used to infer TF binding sites. We added PFMs for MYRF, SP7, and OLIG2 that are one of the key TFs from another study⁸², *Mus musculus* in JASPAR2022⁶⁶, and HOCOMOCO v12⁸³, respectively. We also included shorter motifs for other key TFs, such as SOX10, MYRF, ZNF24, NKX2.2, and SP7, considering their importance in oligodendrocytes (Supplementary Fig. 3).

Co-enrichment analysis

We used a hypergeometric test to assess whether a number of overlaps in the binding sites for two TFs follows a hypergeometric distribution. Specifically, given that a random variable X represents the possible outcomes of a hypergeometric process, the probability of getting k or more overlapping binding sites between two TFs inside a particular chosen set, as a hypergeometric random process, is

$$\Pr(X \geq k | n; N; m) = \sum_{x=k}^{\min(n,m)} \frac{\binom{m}{x} \binom{N-m}{n-x}}{\binom{N}{n}} \quad (1)$$

where N is the total number of transcription binding sites for all TFs, m is the number of binding sites for TF1, n is the number of binding sites for TF2, and x is the number of overlapping binding (co-occurrence) sites between TF1 and TF2. We applied an FDR adjusted p -value as a cutoff (<0.1) for all possible TF pairs and chose co-binding TF pairs.

Key transcription factors

We defined ten key TFs that are oligodendrocyte marker genes based on mouse loss-of-function studies that have shown that specific TFs are critical for oligodendrocyte differentiation. This includes SOX10²⁹, SOX2^{30,31}, SOX8³², MYRF³³, OLIG1³⁴, OLIG2³⁵, TCF7L2^{36,37}, ZNF24²⁵, NKX2.2³⁸, and NKX6.2³⁹ were chosen as key TFs. Ten 'Oligodendrocyte-specific key TFs' are oligodendrocyte differentially expressed TFs and key TFs, eighty-three 'oligodendrocyte-specific non-key TFs' are oligodendrocyte differentially expressed TFs but not key TFs, and a hundred-thirteen 'non-oligodendrocyte-specific TFs' are neither oligodendrocyte differentially expressed TFs nor key TFs. Especially, ten 'Oligodendrocyte-specific key TFs' play crucial roles in the development and differentiation of oligodendrocytes. They regulate various stages of oligodendrocyte maturation and promote the expression of myelin genes; essentially, they are key players in the process of myelination within the CNS.

Deep learning models

We inputted expression levels of TFs that have co-binding pairs into the deep neural network (DNN) models to predict TG expression levels. 2004 metacells (samples), 206 TFs (features), and a TG expression level (label)

were used in the DNN models. A DNN for each TG was built to predict a TG expression level. The mean squared error (MSE) between predicted TG expression and actual TG expression was used as the loss function in DNN models. We cross-validated the training dataset (80% of the input samples) with 5-fold cross-validation and validated the best trained model on the 20% of hold-out validation dataset for the best use of data and to achieve reliable model performance. We used an early stopping function with patience 10 and determined the number of epochs and we set the batch size to 32. Adam with a learning rate 0.001 was used for training the models. The structure of our neural network model can be written as

$$Z_i = f(W_i \cdot X + b_i) \quad (2)$$

where X denotes the input data and f represents the activation function, specifically the LeakyReLU function. TF expression levels serve as the input data, while Z_i represents the output of the i^{th} hidden layer. The final output of the model is the predicted TG expression level, and W_i and b_i are the weight matrix and bias vector for the i^{th} layer, respectively.

To evaluate the performance of our neural network model, we utilize the Mean Squared Error (MSE) loss function. The MSE quantifies the average squared difference between the predicted outputs of the model, Z and the true labels in our regression task. Mathematically, we can express the MSE as follows:

$$MSE = \frac{1}{N} \sum_{i=1}^N (Z_i - Y_{true})^2, \quad (3)$$

where Z_i represents the predicted output for the i^{th} sample, and Y_{true} denotes the true label corresponding to the i^{th} sample.

Shapley interaction scores

We denote the set of all TFs by F , a feature $i \in F$, and a feature set $S \subseteq F$. We define the interaction effect between TF i and j , with feature set S , of a neural network f at a data point X_k to be

$$\delta_{ij}^f = f(X_k; S \cup \{i, j\}) - f(X_k; S \cup \{i\}) - f(X_k; S \cup \{j\}) + f(X_k; S), \quad (4)$$

where $f(X_k; S)$ is the prediction at X_k when only TFs in S are used, which often requires retraining the NN multiple times. A common approximation is to replace the absent features (i.e., $F \setminus S$) by the corresponding values in a baseline $C_{F \setminus S}$, such that

$$f(X_k; S) \approx f(X_k; S; C_{F \setminus S}) \quad (5)$$

The baseline is set as the empirical mean of each feature. The Shapley interaction score $S_{ij}^f(X_k)$ is the expectation of $\delta_{ij}^f(X_k; S)$,

$$S_{ij}^f(X_k) = E_{p(S)} \left[\delta_{ij}^f(X_k; S) \right], \quad (6)$$

over a uniformly random chosen feature set S from F . We use Monte-Carlo procedure⁸⁴ to approximate $S_{ij}^f(X_k)$ by a small number of samples of S . To aggregate the local interaction effect at different data points into a global interaction effect, we take the expectation $|S_{ij}^f(X_k)|$ of w.r.t. the empirical data distribution $p(X)$, such that

$$S_{ij}^f = E_{p(X)} \left[|S_{ij}^f(X)| \right] \quad (7)$$

For our deep ensemble of deep learning models, we utilize a posterior distribution of functions $q(f)$ induced by the ensemble distribution of the weights $q(w)$, as outlined in Eq. (2). This ensemble approach involves

training multiple instances of the model, each initialized with different random weights to promote diverse learning paths.

The weights w are drawn from a Gaussian prior, reflecting our initial uncertainty about their values. After training, we apply Bayesian inference techniques to update our beliefs about these weights and compute the posterior distribution $q(w)$. This posterior captures the uncertainty in the model parameters, providing a more comprehensive understanding of the model's behavior.

The function $q(f)$ represents the expected output of the model across this ensemble of weights. To compute the interaction score, we take the expectation of the interaction score SI_{ij} with respect to $q(f)$. This is estimated by averaging N_f samples drawn from the ensemble:

$$SI_{ij} = E_{q(f)} \left[SI_{ij}^f \right] \approx \frac{1}{N_f} \sum_{k=1}^{N_f} SI_{ij}^{f_k}. \quad (8)$$

We compute Shapley interaction scores^{23,24} for the co-binding TF pairs, TF i and TF j using the trained DNN models and validation datasets. We calculate mean values for co-binding TF pairs using interaction matrices. We rank them by percentile and scaled them to 0 and 1 for easier interpretation.

Coefficient of variance

The coefficient of variation (CV) is a statistical measure of the dispersion of data points in a data series around the mean. The CV represents the ratio of the standard deviation to the mean, and it is a useful statistic for comparing the degree of variation from one data series to another, even if the means are drastically different from one another. The CV is defined as the ratio of standard deviation to the mean as follows:

$$CV = \frac{\sigma}{\mu} \quad (9)$$

Hierarchy analysis

We computed connectivity statistics, out-degree (O) and in-degree (I), for individual TFs to get a 'hierarchy height' metric (h), a normalized value of the difference between O and I for each TF. The h is calculated as

$$h = \frac{O - I}{O + I} \quad (10)$$

We defined TFs as top-regulator ($h > 0.33$), middle-regulator ($-0.33 < h < 0.33$), and bottom-regulator ($h < -0.33$) by their h values.

Statistics and reproducibility

Data manipulation and analyses were performed using Python 3.10.14 and R 4.3.1. All relevant information including the sample sizes in the groups for statistical tests are included in the figure legends. The plots in this study are generated by Scanpy⁷⁹ (v1.10.3), and seaborn (v0.13.2) in Python and ggplot2 (v3.5.1) in R.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

All data supporting the results are included in Supplementary Data 1–5 and are publicly available on GitHub (<https://github.com/daifengwanglab/coTF-reg>). For the main analyses, scATAC-seq data were obtained from Supplementary Materials of the Multi-omics study²¹ and scRNA-seq data was sourced from SEA-AD: Seattle Alzheimer's Disease Brain Cell Atlas (<https://cellxgene.cziscience.com/collections/1ca90a2d-2943-483d-b678-b809bf464c30>). For the independent validation, scRNA-seq data were

acquired from the following websites: (1) Multiome (<https://cellxgene.cziscience.com/collections/ceb895f4-ff9f-403a-b7c3-187a9657ac2c>); (2) ROSMAP (https://compbio.mit.edu/ad_aging_brain/#loading-the-raw-data); (3) Cross-disorder (<https://cellxgene.cziscience.com/collections/c53573b2-eff4-4c5e-9ad0-b24d422df9b>).

Code availability

The code for the analyses and figures is available at <https://github.com/daifengwanglab/coTF-reg>.

Received: 19 June 2024; Accepted: 17 January 2025;

Published online: 05 February 2025

References

- Bercury, K. K. & Macklin, W. B. Dynamics and mechanisms of CNS myelination. *Dev. Cell* **32**, 447–458 (2015).
- Nasrabady, S. E., Rizvi, B., Goldman, J. E. & Brickman, A. M. White matter changes in Alzheimer's disease: a focus on myelin and oligodendrocytes. *Acta Neuropathol. Commun.* **6**, 22 (2018).
- Singh, D. K., Ling, E. & Kaur, C. Hypoxia and myelination deficits in the developing brain. *Int. J. Dev. Neurosci.* **70**, 3–11 (2018).
- Emery, B. & Lu, Q. R. Transcriptional and epigenetic regulation of oligodendrocyte development and myelination in the central nervous system. *Cold Spring Harb. Perspect. Biol.* **7**, a020461 (2015).
- Valdés-Tovar, M. et al. Insights into myelin dysfunction in schizophrenia and bipolar disorder. *WJP* **12**, 264–285 (2022).
- Maitre, M. et al. Myelin in Alzheimer's disease: culprit or bystander? *Acta Neuropathol. Commun.* **11**, 56 (2023).
- Quan, L., Uyeda, A. & Muramatsu, R. Central nervous system regeneration: the roles of glial cells in the potential molecular mechanism underlying remyelination. *Inflamm. Regen.* **42**, 7 (2022).
- Simons, M. & Nave, K.-A. Oligodendrocytes: myelination and axonal support. *Cold Spring Harb. Perspect. Biol.* **8**, a020479 (2016).
- Eibaz, B. & Popko, B. Molecular control of oligodendrocyte development. *Trends Neurosci.* **42**, 263–277 (2019).
- Ibarra, I. L. et al. Mechanistic insights into transcription factor cooperativity and its impact on protein-phenotype interactions. *Nat. Commun.* **11**, 124 (2020).
- Lopez-Anido, C. et al. Differential Sox10 genomic occupancy in myelinating glia. *Glia* **63**, 1897–1914 (2015).
- Sock, E. & Wegner, M. Using the lineage determinants Olig2 and Sox10 to explore transcriptional regulation of oligodendrocyte development. *Dev. Neurobiol.* **81**, 892–901 (2021).
- Bujalka, H. et al. MYRF is a membrane-associated transcription factor that autoproteolytically cleaves to directly activate myelin genes. *PLoS Biol.* **11**, e1001625 (2013).
- Weider, M. et al. Nfat/calcineurin signaling promotes oligodendrocyte differentiation and myelination by transcription factor network tuning. *Nat. Commun.* **9**, 899 (2018).
- Heinz, S., Romanoski, C. E., Benner, C. & Glass, C. K. The selection and function of cell type-specific enhancers. *Nat. Rev. Mol. Cell Biol.* **16**, 144–154 (2015).
- Emani, P. S. et al. Single-cell genomics and regulatory networks for 388 human brains. *Science* **384**, eadi5199 (2024).
- Lonsdale, J. et al. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
- Bravo González-Bias, C. et al. SCENIC+: single-cell multiomic inference of enhancers and gene regulatory networks. *Nat. Methods* **20**, 1355–1367 (2023).
- Stuart, T., Srivastava, A., Madad, S., Lareau, C. A. & Satija, R. Single-cell chromatin state analysis with Signac. *Nat. Methods* **18**, 1333–1341 (2021).
- Jin, T. et al. scGRNom: a computational pipeline of integrative multi-omics analyses for predicting cell-type disease genes and regulatory networks. *Genome Med.* **13**, 95 (2021).

21. Zhu, K. et al. Multi-omic profiling of the developing human cerebral cortex at the single-cell level. *Sci. Adv.* **9**, eadg3754 (2023).
22. Gabitto, M. L., Travaglini, K. J., Rachleff, V. M. et al. Integrated multimodal cell atlas of Alzheimer's disease. *Nat. Neurosci.* **27**, 2366–2383 (2024).
23. Cui, T. et al. Gene-gene interaction detection with deep learning. *Commun. Biol.* **5**, 1238 (2022).
24. Dhamdhere, K., Agarwal, A. & Sundararajan, M. The Shapley Taylor Interaction Index. <https://doi.org/10.48550/ARXIV.1902.05622>. (2019).
25. Elbaz, B. et al. Phosphorylation sState of ZFP24 controls oligodendrocyte differentiation. *Cell Rep.* **23**, 2254–2263 (2018).
26. Bernhardt, C. et al. KLF9 and KLF13 transcription factors boost myelin gene expression in oligodendrocytes as partners of SOX10 and MYRF. *Nucleic Acids Res.* **50**, 11509–11528 (2022).
27. Howng, S. Y. B. et al. ZFP191 is required by oligodendrocytes for CNS myelination. *Genes Dev.* **24**, 301–311 (2010).
28. Al-Naama, N., Mackeh, R. & Kino, T. C2H2-Type Zinc finger proteins in brain development, neurodevelopmental, and other neuropsychiatric disorders: systematic literature-based analysis. *Front. Neurol.* **11**, 32 (2020).
29. Aprato, J. et al. Myrf guides target gene selection of transcription factor Sox10 during oligodendroglial development. *Nucleic Acids Res.* **48**, 1254–1270 (2020).
30. Zhao, C. et al. Sox2 sustains recruitment of oligodendrocyte progenitor cells following CNS demyelination and primes them for differentiation during remyelination. *J. Neurosci.* **35**, 11482–11499 (2015).
31. Zhang, S. et al. Sox2 is essential for oligodendroglial proliferation and differentiation during postnatal brain myelination and CNS remyelination. *J. Neurosci.* **38**, 1802–1820 (2018).
32. Turnescu, T. et al. Sox8 and Sox10 jointly maintain myelin gene expression in oligodendrocytes. *Glia* **66**, 279–294 (2018).
33. Garnai, S. J. et al. Variants in myelin regulatory factor (MYRF) cause autosomal dominant and syndromic nanophthalmos in humans and retinal degeneration in mice. *PLoS Genet.* **15**, e1008130 (2019).
34. Chen, Y. et al. The oligodendrocyte-specific G protein-coupled receptor GPR17 is a cell-intrinsic timer of myelination. *Nat. Neurosci.* **12**, 1398–1406 (2009).
35. Wang, J. et al. Olig2 ablation in immature oligodendrocytes does not enhance CNS myelination and remyelination. *J. Neurosci.* **42**, 8542–8555 (2022).
36. Zhao, C. et al. Dual regulatory switch through interactions of Tcf7l2/Tcf4 with stage-specific partners propels oligodendroglial maturation. *Nat. Commun.* **7**, 10883 (2016).
37. Ye, F. et al. HDAC1 and HDAC2 regulate oligodendrocyte differentiation by disrupting the beta-catenin-TCF interaction. *Nat. Neurosci.* **12**, 829–838 (2009).
38. Qi, Y. et al. Control of oligodendrocyte differentiation by the Nkx2.2 homeodomain transcription factor. *Development* **128**, 2723–2733 (2001).
39. Southwood, C. et al. CNS myelin paranodes require Nkx6-2 homeoprotein transcriptional activity for normal structure. *J. Neurosci.* **24**, 11215–11225 (2004).
40. Mathys, H. et al. Single-cell atlas reveals correlates of high cognitive function, dementia, and resilience to Alzheimer's disease pathology. *Cell* **186**, 4365–4385.e27 (2023).
41. Bunk, E. C. et al. Prox1 is required for oligodendrocyte cell identity in adult neural stem cells of the subventricular zone. *Stem Cells* **34**, 2115–2129 (2016).
42. Kato, K. et al. Prox1 inhibits proliferation and is required for differentiation of the oligodendrocyte cell lineage in the mouse. *PLoS ONE* **10**, e0145334 (2015).
43. Suzuki, N. et al. Differentiation of oligodendrocyte precursor cells from Sox10-venus mice to oligodendrocytes and astrocytes. *Sci. Rep.* **7**, 14133 (2017).
44. Bryois, J. et al. Cell-type-specific cis-eQTLs in eight human brain cell types identify novel risk genes for psychiatric and neurological disorders. *Nat. Neurosci.* **25**, 1104–1112 (2022).
45. Yu, Y. et al. Olig2 targets chromatin remodelers to enhancers to initiate oligodendrocyte differentiation. *Cell* **152**, 248–261 (2013).
46. Aguado, L. C. et al. microRNA function is limited to cytokine control in the acute response to virus infection. *Cell Host Microbe* **18**, 714–722 (2015).
47. Galiano, M. R. et al. Myelin basic protein functions as a microtubule stabilizing protein in differentiated oligodendrocytes. *J. Neurosci. Res.* **84**, 534–541 (2006).
48. Aber, E. R. et al. Oligodendroglial macroautophagy is essential for myelin sheath turnover to prevent neurodegeneration and death. *Cell Rep.* **41**, 111480 (2022).
49. Ehrlich, M. et al. Rapid and efficient generation of oligodendrocytes from human induced pluripotent stem cells using transcription factors. *Proc. Natl Acad. Sci. USA* **114**, E2243–E2252 (2017).
50. Smirnova, E. V. et al. Comprehensive Atlas of the myelin basic protein interaction landscape. *Biomolecules* **11**, 1628 (2021).
51. Snaidero, N. et al. Antagonistic Functions of MBP and CNP establish cytosolic channels in CNS Myelin. *Cell Rep.* **18**, 314–323 (2017).
52. Wang, D. et al. Logic: a method to characterize the cooperative logic of regulatory factors. *PLoS Comput. Biol.* **11**, e1004132 (2015).
53. Rexach, J. E. et al. Cross-disorder and disease-specific pathways in dementia revealed by single-cell genomics. *Cell* **187**, 5753–5774.e28 (2024).
54. Yashar, W. M. et al. Predicting transcription factor activity using prior biological information. *iScience* **27**, 109124 (2024).
55. Duren, Z. et al. Sc-compReg enables the comparison of gene regulatory networks between conditions using single-cell data. *Nat. Commun.* **12**, 4763 (2021).
56. Ferrari, C., Manosalva Pérez, N. & Vandepoel, K. MINI-EX: integrative inference of single-cell gene regulatory networks in plants. *Mol. Plant* **15**, 1807–1824 (2022).
57. Duren, Z., Chen, X., Jiang, R., Wang, Y. & Wong, W. H. Modeling gene regulation from paired expression and chromatin accessibility data. *Proc. Natl Acad. Sci. USA* **114**, E4914–E4923 (2017).
58. Liu, Z. et al. Induction of oligodendrocyte differentiation by Olig2 and Sox10: Evidence for reciprocal interactions and dosage-dependent mechanisms. *Dev. Biol.* **302**, 683–693 (2007).
59. Kim, J. et al. The co-regulation mechanism of transcription factors in the human gene regulatory network. *Nucleic Acids Res.* **40**, 8849–8861 (2012).
60. Fuda, N. J., Ardehali, M. B. & Lis, J. T. Defining mechanisms that regulate RNA polymerase II transcription in vivo. *Nature* **461**, 186–192 (2009).
61. Buchler, N. E., Gerland, U. & Hwa, T. On schemes of combinatorial transcription logic. *Proc. Natl Acad. Sci. USA* **100**, 5136–5141 (2003).
62. Silva-Rocha, R. & De Lorenzo, V. Mining logic gates in prokaryotic transcriptional regulation networks. *FEBS Lett.* **582**, 1237–1244 (2008).
63. Krumsiek, J., Marr, C., Schroeder, T. & Theis, F. J. Hierarchical differentiation of myeloid progenitors is encoded in the transcription factor network. *PLoS ONE* **6**, e22649 (2011).
64. Malekpour, S. A., Shahdoust, M., Aghdam, R. & Sadeghi, M. wpLogicNet: logic gate and structure inference in gene regulatory networks. *Bioinformatics* **39**, btad072 (2023).
65. Malekpour, S. A., Haghverdi, L. & Sadeghi, M. Single-cell multi-omics analysis identifies context-specific gene regulatory gates and mechanisms. *Brief. Bioinform.* **25**, bbae180 (2024).
66. Castro-Mondragon, J. A. et al. JASPAR 2022: the 9th release of the open-access database of transcription factor binding profiles. *Nucleic Acids Res.* **50**, D165–D173 (2022).
67. Alicia, S. motifmatcher: Fast Motif Matching in R. R package version 1.28.0. <https://doi.org/10.18129/B9.bioc.motifmatcher> (2024).

68. Hao, Y. et al. Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573–3587.e29 (2021).
69. Cavalcante, R. G. & Sartor, M. A. annotatr: genomic regions in context. *Bioinformatics* **33**, 2381–2383 (2017).
70. Ben-Kiki, O., Bercovich, A., Lifshitz, A. & Tanay, A. Metacell-2: a divide-and-conquer metacell algorithm for scalable scRNA-seq analysis. *Genome Biol.* **23**, 100 (2022).
71. Paszke, A. et al. PyTorch: An imperative style, high-performance deep learning library. <https://doi.org/10.48550/ARXIV.1912.01703>. (2019).
72. Koopmans, L. H., Owen, D. B. & Rosenblatt, J. I. Confidence intervals for the coefficient of variation for the normal and log normal distributions. *Biometrika* **51**, 25–32 (1964).
73. Shannon, P. et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
74. Gerstein, M. B. et al. Architecture of the human regulatory network derived from ENCODE data. *Nature* **489**, 91–100 (2012).
75. Kent, W. J. et al. The human genome browser at UCSC. *Genome Res.* **12**, 996–1006 (2002).
76. Lerdrup, M., Johansen, J. V., Agrawal-Singh, S. & Hansen, K. An interactive environment for agile analysis and visualization of ChIP-seq data. *Nat. Struct. Mol. Biol.* **23**, 349–357 (2016).
77. Müssel, C., Hopfensitz, M. & Kestler, H. A. BoolNet—an R package for generation, reconstruction and analysis of Boolean networks. *Bioinformatics* **26**, 1378–1380 (2010).
78. Smedley, D. et al. BioMart – biological queries made easy. *BMC Genom.* **10**, 22 (2009).
79. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* **19**, 15 (2018).
80. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**, 495–502 (2015).
81. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. <https://doi.org/10.48550/ARXIV.1802.03426>. (2018).
82. Kim, D. et al. Homo-trimerization is essential for the transcription factor function of Myrf for oligodendrocyte differentiation. *Nucleic Acids Res.* **45**, 5112–5125 (2017).
83. Vorontsov, I. E. et al. HOCOMOCO in 2024: a rebuild of the curated collection of binding models for human and mouse transcription factors. *Nucleic Acids Res.* **52**, D154–D163 (2024).
84. Štrumbelj, E. & Kononenko, I. Explaining prediction models and individual predictions with feature contributions. *Knowl. Inf. Syst.* **41**, 647–665 (2014).

Acknowledgements

This work was supported by National Institutes of Health grants, R21 NS128761, RF1MH128695, R01AG067025, R21NS128761, and National

Science Foundation Career Award 2144475, and a core grant to the Waisman Center from NICHD (P50 HD105353).

Author contributions

Conceptualization, J.S. and D.W.; Methodology, J.C., J.S., and D.W.; Formal Analysis, J.C.; Investigation, J.C., J.S., and D.W.; Writing – Original Draft, J.C.; Writing – Review & Editing, J.C., J.S., and D.W.; Supervision, J.S., and D.W.; Funding Acquisition, J.S. and D.W.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s42003-025-07570-6>.

Correspondence and requests for materials should be addressed to Daifeng Wang.

Peer review information *Communications Biology* thanks Seyed Malekpour and Ping-Han Hsieh for their contribution to the peer review of this work. Primary Handling Editors: Chien-Yu Chen and Mengtan Xing.

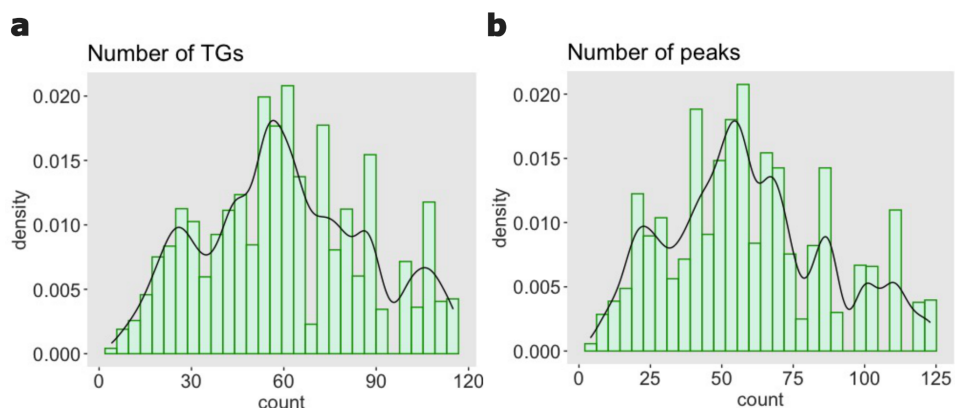
Reprints and permissions information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

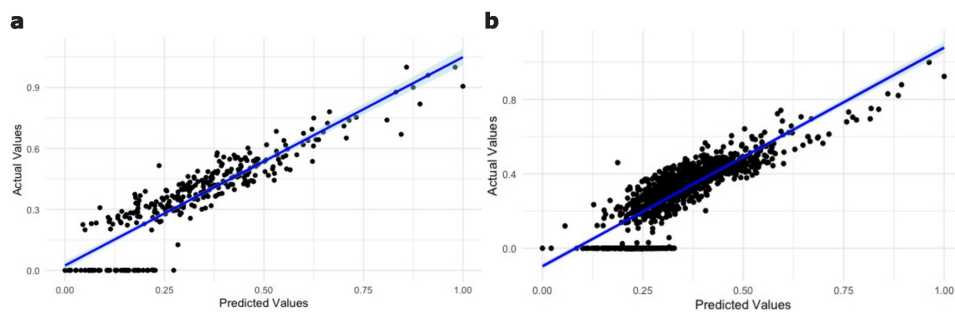
Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025

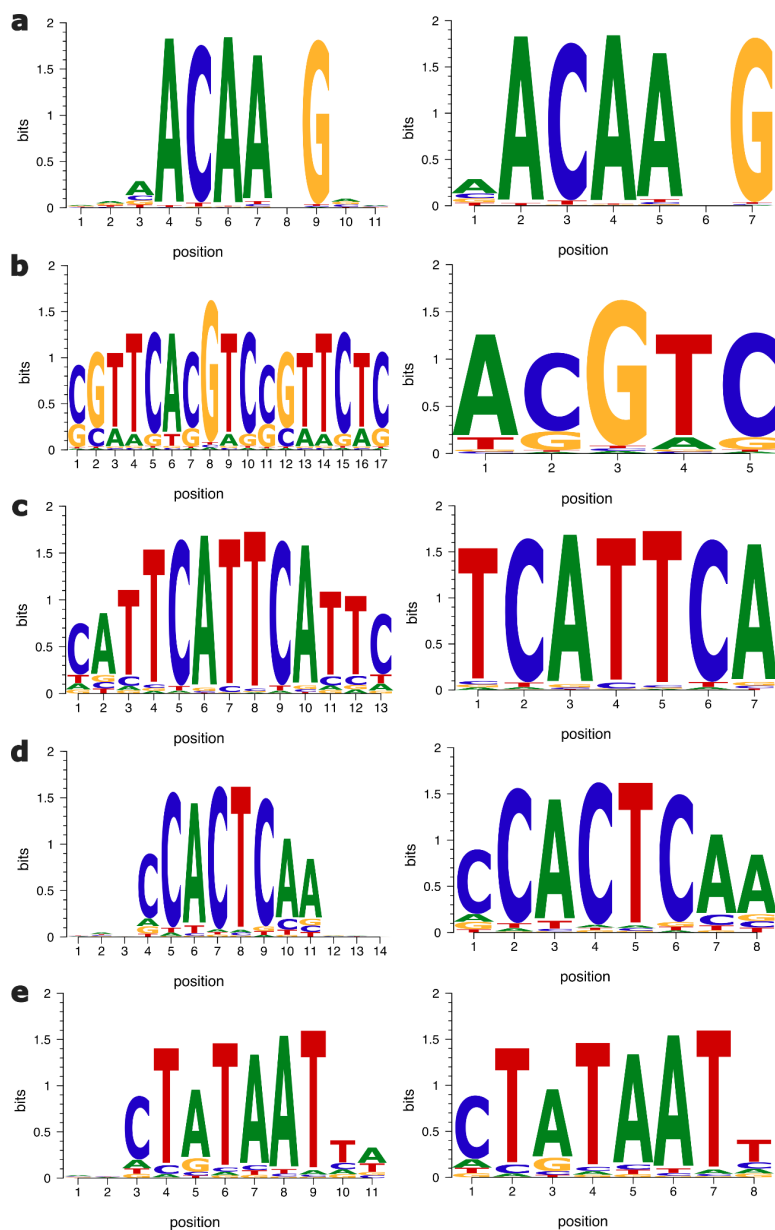
Supplementary information



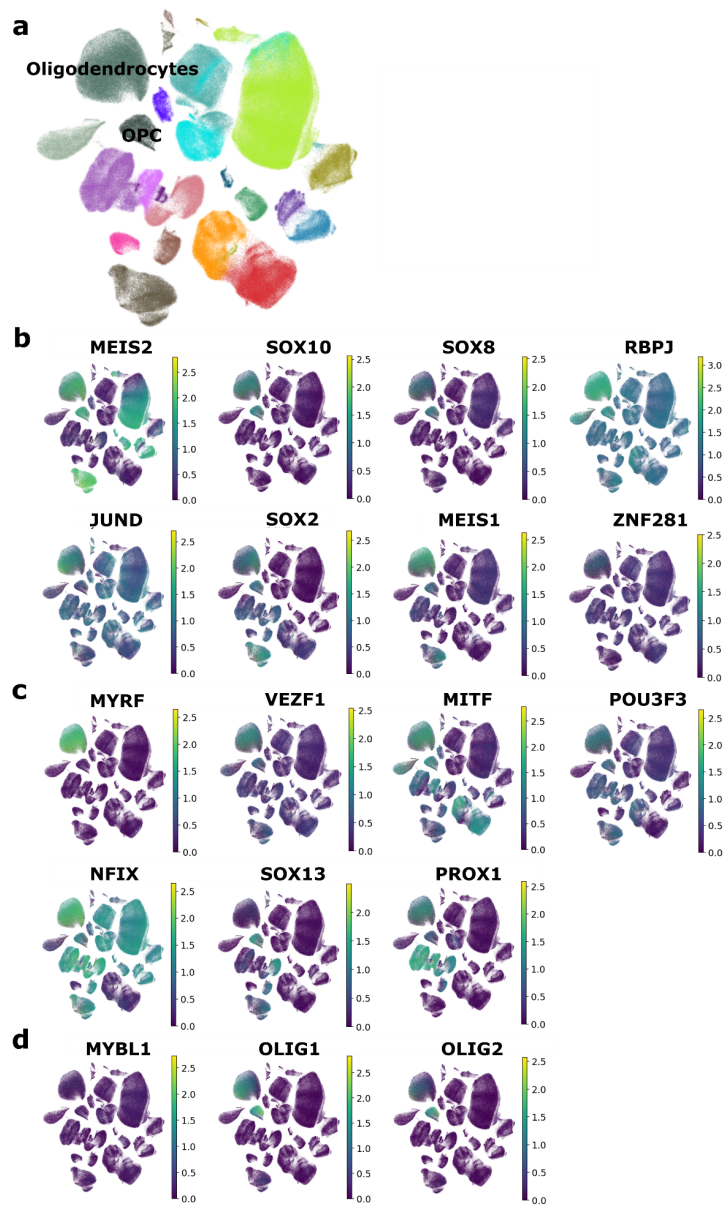
Supplementary Figure 1: Distributions of the numbers of target genes and peaks for co-binding TF pairs. **a** Distribution of the number of target genes (TGs) for co-binding transcription factor pairs. **b** Distribution of the number of peaks for co-binding transcription factor pairs.



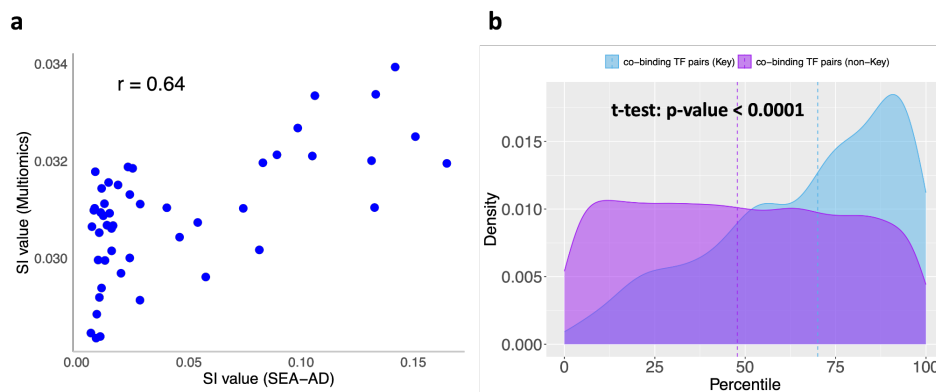
Supplementary Figure 2: Validation for the deep learning prediction. **a** Comparison between predicted values and actual labels in a deep learning model for predicting *MBP* using co-binding transcription factors in *Integrated Multimodal Cell Atlas of Alzheimer's Disease* data. **b** Comparison between predicted values and actual labels in a deep learning model for predicting *MBP* using co-binding transcription factors in *Single-cell atlas reveals correlates of high cognitive function, dementia, and resilience to Alzheimer's disease pathology* data.



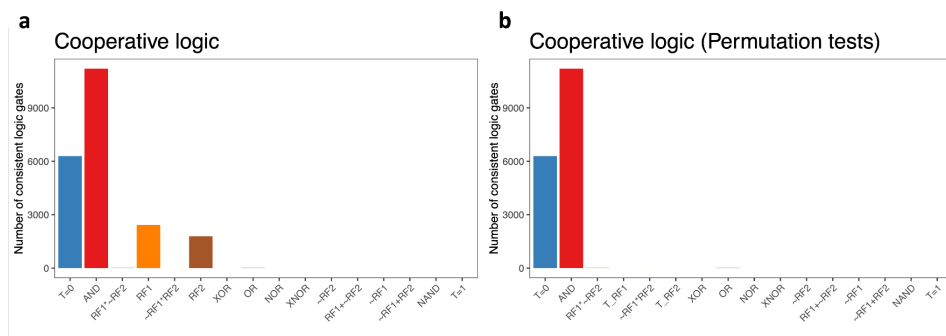
Supplementary Figure 3: Logo plots for the original and shorter motifs (The first column shows the long motifs and the second column shows the short motifs) a SOX10. b MYRF. c ZNF24. d NKX2.2. e SP7.



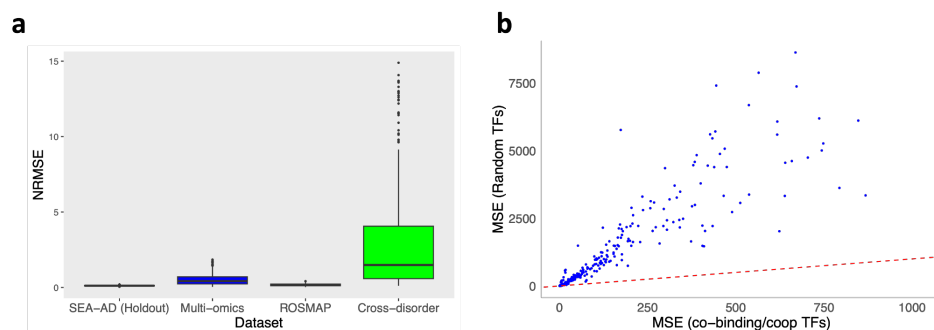
Supplementary Figure 4: UMAPs for all cell types and regulators. a UMAP for all cell types **b** top-level regulators ('Master regulators'). **c** middle-level regulators. **d** bottom-level regulators.



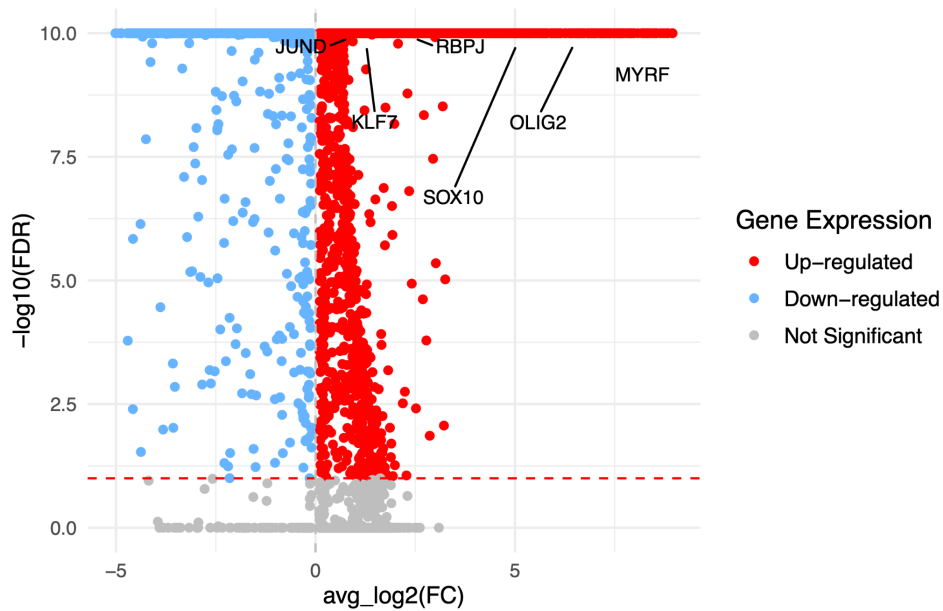
Supplementary Figure 5: TF interaction based on Multi-omics RNA-seq data. a Comparison for Top SI values for key-TF computed by two different scRNA-seq data, **b** T-test for the interaction score distributions of key-transcription factor pairs and non-key transcription factor pairs.



Supplementary Figure 6: Application of Boolean logic to find cooperative TFs. a Distribution of logic gates of gate-consistent triples. **b** Distribution of logic gates of gate-consistent triples after removing spurious logic gates by permutation tests.



Supplementary Figure 7: Holdout evaluation performance and model's prediction performance. a The x-axis represents four groups and the y-axis show normalized root mean squared error (NRMSE) **b** Mean squared error comparison for 206 co-binding/cooperative TFs and 206 random TFs that are neither co-binding nor cooperative.



Supplementary Figure 8: Differential gene expression testing for the 17,946 genes in SEA-AD data. The X-axis represents average log2 fold change values and the y-axis shows $-\log_{10}$ FDR values. Above the red dashed line, which is a cutoff for the FDR, the most up-regulated and significant genes are towards the

right (red), the most down-regulated and significant genes are towards the left (blue), and the statistically non-significant genes are under the red dashed line (grey).

Supplementary Data 1

Supplementary Data 2

Supplementary Data 3

Supplementary Data 4

Supplementary Data 5

© Copyright by Jerome J. Choi 2025
All Rights Reserved