



UNIVERSITY OF
CAMBRIDGE

Department of Engineering

Causal Graph Discovery from Genomic Data in Health and Alzheimer's Disease

Author Name: Yuan Zhewen

Supervisor: Prof Ramji Venkataramanan

Date: June 2, 2025

I hereby declare that, except where specifically indicated, the work submitted herein is my own original work.

Signed YUAN ZHEWEN date 2nd June 2025

1 Abstract

In Alzheimer’s disease (AD), neuroinflammation and amyloid- β accumulation fuel a damaging feedback loop, accelerating disease progression. Understanding this interplay is crucial because neuroinflammation is known as an early driver of AD, making it a key target for effective therapies aimed at halting this destructive cycle. Based on previous studies, we hypothesized that amyloid- β might act as a key mediator linking neuroinflammation, measured by soluble TREM2 (sTREM2) levels in cerebrospinal fluid, to downstream neurodegeneration. To test this, we performed a causal mediation analysis on longitudinal ADNI biomarker data, treating amyloid- β levels as a mediator of sTREM2’s effect on AD progression. The analysis yielded no evidence of a mediated effect: amyloid- β did not significantly carry sTREM2’s influence on disease development, and sTREM2 levels alone showed no clear distinction across AD stages. These negative findings indicate that neuroinflammation’s impact on AD progression cannot be explained solely by amyloid- β , suggesting that the relationship among inflammation, amyloid- β , and neurodegeneration is more complex than a simple linear cascade.

In light of this complexity, we shifted from a single-mediator causal mediation analysis to a more exploratory approach, applying a causal discovery pipeline to investigate the interplay between neuroinflammation and AD pathology. In particular, we turned to microglial gene expression to explore how neuroinflammatory states might causally intersect with AD’s molecular pathways. We analyzed single-cell RNA sequencing data from microglia isolated from the human prefrontal cortex, comparing cells from AD patients with those from healthy ones. To identify genes most strongly associated with AD status, we implemented a multi-method feature selection pipeline. Three distinct algorithms: random forest importance, sparse logistic regression, and gradient boosting were applied in parallel within a repeated subsampling framework to assess feature stability. This robust procedure yielded a consistent panel of differentially expressed genes predictive of AD diagnosis in microglia. Notably, the final gene panel included both well-established AD-related genes and lesser-known candidates, capturing pathways ranging from inflammatory activation to protein homeostasis and metabolic regulation. Throughout the feature selection process, we took care to avoid overfitting and evaluated each algorithm using standard binary classification performance metrics on held-out subsets, ensuring the generalizability of the selected features.

With the refined gene set, we applied multiple causal discovery algorithms, leveraging their complementary strengths to enhance confidence in the inferred relationships. We applied several causal discovery methods in parallel to the selected gene-expression data. The PC algorithm, a constraint-based approach, learns a directed acyclic graph (DAG) structure by testing conditional independencies. Greedy Equivalence Search (GES), a score-based method, refines that structure by optimizing a model-selection criterion. LiNGAM (Linear Non-Gaussian Acyclic Model) recovers causal ordering under linear, non-Gaussian assumptions. Finally, a variational autoencoder (VAE)-based technique enforces DAG constraints during training to uncover more complex, potentially nonlinear relationships. This model is capable of detecting complex, nonlinear interactions by learning latent representations while searching for a directed network. The causal graphs generated by

different discovery algorithms exhibit substantial structural differences. Therefore, the final integrative causal graph was constructed from the consensus edges - those that appeared robustly across algorithms. The causal discovery pipeline also offered insights into the suitability of various methods and the extent to which violations of some assumptions affect their performance in the context of high-dimensional gene expression data.

The resulting causal gene network sheds light on how microglial activation may drive or modulate AD-related pathways. Several genes emerged as hubs and bottleneck regulators, indicating central roles in network connectivity. Notably, well-established AD risk genes appeared as key regulators in multiple discovery algorithms, forming critical bridges between neuroinflammation, amyloid pathology, and mitochondrial or metabolic dysfunction. For example, the long-noncoding RNA **MALAT1** reduces neuronal injury and suppresses neuroinflammation; the transcription factor **MEF2C** serves as an off-switch to mitigate inflammatory cascades; and **LYN**, a potent immune regulator, critically modulates microglial function.

Our causal network notably identified several genes that have only recently garnered attention for their potential roles in AD. For instance, **B2M**, a component of MHC class I, emerged as a critical link between immune activity and amyloid pathology; **DHRSX**, a relatively novel protein, regulates autophagy; **CD53**, a cell-surface organizer associated with diverse immune functions, has only just been recognized as upregulated in disease-associated microglia; and **DOCK8**, known for controlling cell migration, has only recently been implicated in AD contexts. Identifying these emerging genes highlights our network's ability to reveal novel candidates and pinpoint previously underappreciated players.

Collectively, these results reveal a highly interconnected microglial gene network in AD. Rather than supporting a simple linear cascade, our results point to a more complex system involving multiple converging pathways. Neuroinflammatory signals, mediated by TREM2 and other factors, feed into transcriptional programs that in turn influence amyloid deposition dynamics through A β -related genes, as well as genes regulating mitochondrial and metabolic function. Together, these interconnected processes contribute to the overall progression of Alzheimer's disease.

Overall, this research advances our understanding of Alzheimer's disease by clarifying the complex role of microglia and uncovering potential therapeutic targets through a comprehensive analytical framework. We advance the field through three primary contributions: First, we challenge the prevailing notion of a simple, linear pathway in which neuroinflammation drives AD progression primarily through increased amyloid- β deposition, showing instead that this model fails to capture the disease's true complexity. Second, we develop a robust causal discovery pipeline tailored to the high dimensionality and sparsity of single-cell RNA sequencing data in microglia - an understudied genomic context in AD research. By integrating rigorous feature selection with multiple complementary causal inference algorithms, our approach addresses both statistical and computational challenges. Finally, our analysis reveals a set of emerging genes and strong causal edges, many of which have only recently been linked to AD, offering a valuable foundation for future experimental validation. By highlighting these overlooked pathways and molecular players, this work refines existing mechanistic models and points toward novel avenues for therapeutic development.

Contents

1 Abstract	2
2 Introduction	5
2.1 Motivation	5
2.2 Objectives and Outline	6
3 Causal Mediation Analysis	7
3.1 Definition	7
3.2 Neuroinflammation, Amyloid Beta and AD	7
3.3 Preliminary Results and Analysis	8
4 Feature Selection	10
4.1 Overview	10
4.1.1 Data	10
4.1.2 Evaluation Metrics	11
4.2 Methodology of Feature Selection	12
4.2.1 Random Forest	12
4.2.2 Sparse Logistic Regression	12
4.2.3 Gradient Boosting	13
4.3 Feature Selection Results and Discussion	13
4.3.1 Random Forest	13
4.3.2 Sparse Logistic Regression	15
4.3.3 Gradient Boosting	16
4.3.4 Performance of Feature Selection Algorithms	17
4.4 Biological Interpretation of Findings	19
5 Causal Discovery	20
5.1 Overview	20
5.2 Methodology of Casual Discovery	21
5.2.1 PC Algorithm	21
5.2.2 GES Method	21
5.2.3 LiNGAM Method	22
5.2.4 VAE Method	23
5.3 Causal Discovery Results and Discussion	24
5.3.1 Causal Discovery on Genes Selected via Gradient Boosting	26
5.3.2 Causal Discovery on Genes Selected via Sparse Logistic Regression .	31
5.3.3 Causal Discovery on Genes Selected via Random Forest	38
5.4 Performance of Causal Discovery Methods	41
6 Future Works	42
7 Conclusion	43

2 Introduction

2.1 Motivation

Alzheimer's disease (AD) stands as a progressive and devastating neurodegenerative disorder, representing the most prevalent form of dementia and posing a global public health crisis. Characterized by an insidious onset, AD leads to a decline in cognitive functions, memory loss, and profound behavioral alterations. The pathological landscape of AD has traditionally been defined by two key protein abnormalities: amyloid- β peptides, which are sticky fragments of a larger protein that clump together outside neurons, and tau protein, which normally helps stabilize the internal skeleton of neurons. In AD, amyloid- β accumulates into extraneuronal plaques, and tau becomes hyperphosphorylated and forms intraneuronal neurofibrillary tangles [1].

While these two features have long been central to AD research, it is increasingly recognized that they represent only part of a more intricate pathogenic cascade, rather than the exclusive drivers of the disease. Currently, there is no single, universally accepted genetic theory that comprehensively explains the pathogenesis of AD. This inherent complexity is mirrored in the limited success of therapeutic strategies narrowly targeting amyloid- β or tau pathology [2], compelling a broader understanding of underlying mechanisms. Indeed, the classical hallmarks are interwoven with other critical pathological processes, including the progressive demise of neurons and synapses, metabolic dysregulation, mitochondrial dysfunction, and aberrant autophagy [3]. Crucially, emerging from this complex interplay is neuroinflammation, now increasingly acknowledged as a pivotal third pathological pillar of Alzheimer's disease, alongside amyloid- β and tau. This neuroinflammatory response is primarily characterized by the chronic activation of glial cells - microglia and astrocytes - and their sustained release of a diverse array of inflammatory mediators.

While the exact role of neuroinflammation in AD - whether causative, accelerative, or secondary - remains debated, its potential to become detrimentally self-sustaining is a major concern. The "pro-inflammatory spiral" concept, for example, suggests aging brains can enter a vicious cycle where initial cellular injury incites neuroinflammation, leading to further damage and perpetuating the disease [4]. Significantly, neuroinflammatory processes can manifest early in AD pathogenesis, potentially hastening disease onset [5]. This inflammation actively interacts with proteins like amyloid- β , creating a detrimental feedback loop: inflammation can promote amyloid- β production and aggregation, while amyloid- β itself further fuels the inflammatory response, accelerating neuronal damage and disease progression. Consequently, understanding the specific inflammatory mediators and the overall role of neuroinflammation in AD is therapeutically crucial, offering new targets for intervention. Identifying ways to modulate or resolve this chronic neuroinflammation, such as by preventing excessive glial activation or neutralizing key inflammatory molecules, could potentially slow or halt the damaging cascade driving neuronal dysfunction and AD progression.

This understanding motivated our study, which aims to elucidate the detailed role of neuroinflammation in AD, specifically its causal relationship with the disease. Causality studies offer distinct advantages in biological research, primarily their ability to infer

cause-and-effect relationships beyond mere statistical correlations. This is crucial for analyzing complex, high-dimensional biological datasets, such as those from single-cell gene expression studies, where numerous potential interactions can obscure underlying mechanisms. By constructing causal networks, these methods can identify critical driver genes or molecules that initiate or propagate disease, highlighting novel therapeutic targets. Given that currently identified genetic factors do not fully explain AD heritability, suggesting many contributors remain unknown, we first investigated established causal relationships from existing research. We then explored causal relationships at the genetic level, leveraging large-scale single-cell gene expression datasets and employing causal discovery approaches.

2.2 Objectives and Outline

Given the current limited understanding of the intricate relationship between neuroinflammation and AD pathogenesis, this work aims to address the following three key research questions:

- How valid is the proposed linear model in which neuroinflammation drives Alzheimer’s disease progression through amyloid- β mediation, and what proportion of the total effect is attributable to amyloid- β in this pathway?
- If the simple linear model of neuroinflammation \rightarrow amyloid- β accumulation \rightarrow AD progression proves insufficient, then at the genetic level, what are the causal regulatory relationships among microglial genes that drive Alzheimer’s pathology?
- Which methodological approaches are most suitable and effective for causal discovery when analyzing high-dimensional single-cell RNA sequencing (scRNA-seq) data?

Our following analysis framework comprises two main components: a causal mediation analysis and a causal discovery pipeline. The causal mediation workflow is illustrated in Figure 1 and described in Section 3. Identifying limitations in the assumed causal relationships motivates the second component: the causal discovery pipeline. This pipeline is depicted in Figure 2 and detailed in Sections 4 (feature selection) and 5.1 (causal discovery). Definitions of terms appearing in the figures are provided in the corresponding sections.

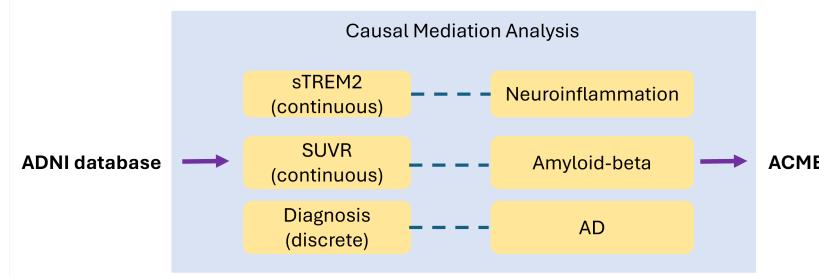


Figure 1: Causal mediation analysis framework

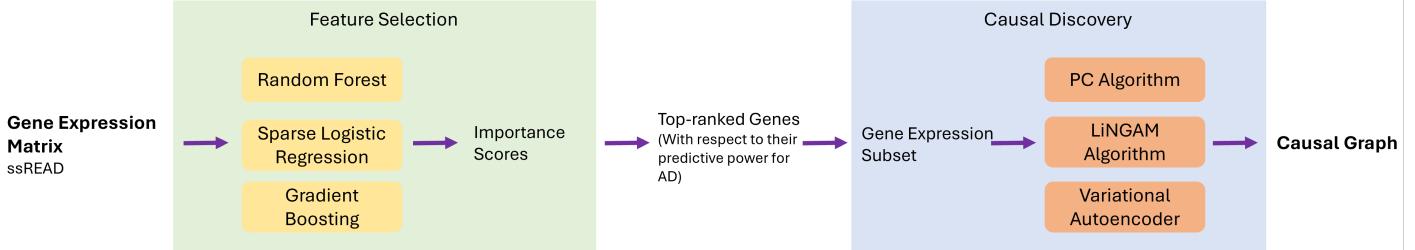


Figure 2: Causal graph discovery framework

3 Causal Mediation Analysis

3.1 Definition

Formulated in a potential outcome framework, causal mediation analysis provides a principled way to decompose the total effect of an exposure or treatment T on an outcome Y into two parts: the indirect effect, which is the portion of T 's effect that operates through a mediator M ; the direct effect, which is the portion of T 's effect that operates outside of M . This decomposition is illustrated in Figure 3. The key estimand for the indirect pathway is the Average Causal Mediation Effect (ACME), first formalized by [6]. For a binary treatment $T \in \{0, 1\}$, it is defined as

$$\bar{\delta}(t) = E [Y_i(t, M_i(1)) - Y_i(t, M_i(0))] ,$$

where $M_i(t)$ is the value that the mediator M would take for unit i under treatment $T = t$, $Y_i(t, m)$ is the outcome we would observe for unit i if we set $T = t$ and (counterfactually) force the mediator to m . The expectation is computed over all the units i , and is an average if all the units are weighted equally. Thus $\bar{\delta}(t)$ averages over the population the change in Y when we swap the mediator from its natural value under control ($M_i(0)$) to its value under treatment ($M_i(1)$) while holding T fixed at t .

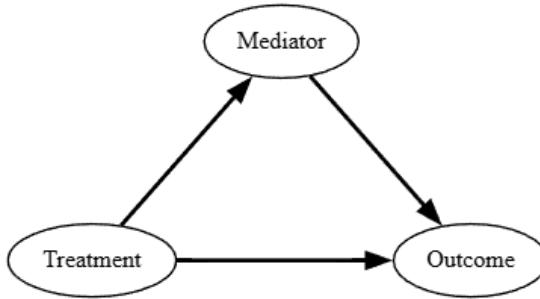


Figure 3: Causal mediation DAG

3.2 Neuroinflammation, Amyloid Beta and AD

Recent work [7] suggests that amyloid- β mediates the relationship between neuroinflammation and Alzheimer's disease (AD). Neuroinflammation arises from dysfunctional mi-

microglia - immune cells in the brain responsible for clearing cellular debris. When microglial function is impaired, amyloid- β accumulates, triggering inflammation and accelerating neuronal damage. This causal pathway is illustrated in Figure 4.

Building on these findings, we aimed to validate and quantify the mediating role of amyloid- β in neuroinflammation-driven AD using causal mediation analysis.

To quantify neuroinflammation, we used sTREM2 as a surrogate marker. sTREM2, the soluble form of the microglial TREM2 receptor, is released into cerebrospinal fluid (CSF) and plasma upon proteolytic cleavage and serves as a measurable indicator of microglial activity [8]. Amyloid- β deposition was visualized via PET imaging and quantified using the standardized uptake value ratio (SUVR). Neurodegeneration, representing AD progression, was categorized into three stages: Healthy, Mild Cognitive Impairment (MCI), and Alzheimer's disease (AD), with the latter indicating advanced disease.

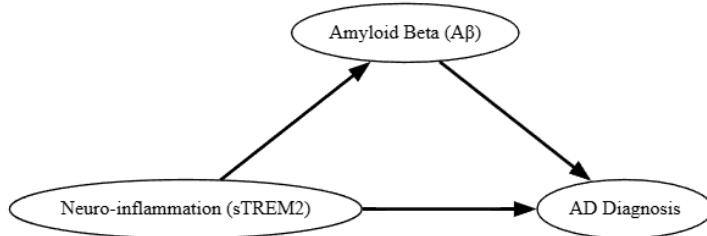


Figure 4: Causal relationship between neuroinflammation, amyloid- β and AD

3.3 Preliminary Results and Analysis

We used data from the Alzheimer's Disease Neuroimaging Initiative (ADNI), which provides a comprehensive, multimodal set of biomarker measurements and participant information across all stages of Alzheimer's disease.

We first conducted an initial causal mediation analysis using continuous amyloid- β levels as the mediator, continuous sTREM2 levels as the treatment, and AD diagnosis (a categorical outcome) as the outcome. The analysis was performed using the mediation package in R [9]. In every case - the ACME under control, the ACME under treatment, and the average ACME - the 95% confidence interval included zero and the p-value exceeded 0.05, indicating no statistically meaningful indirect (mediated) effect. Moreover, the estimated total effect of treatment on AD diagnosis was also small, had a confidence interval spanning zero, and was not statistically significant. These null results run counter to the pathways implied by our causal diagram and our original hypothesis.

To investigate whether sTREM2 and amyloid- β levels differ by AD stage, we performed separate one-way ANOVAs using diagnosis (cognitively normal, MCI, AD) as the grouping factor. The boxplots in Figure 5 visualize the distributions, and Table 2 reports the ANOVA summaries. For sTREM2, the ANOVA yielded no significant effect of diagnosis on mean levels, indicating similar sTREM2 concentrations across stages. By contrast, amyloid- β levels differed highly significantly between groups, consistent with the pronounced separation seen in Figure 5.

Table 1: Causal mediation analysis result of sTREM2 (mediator), A β (treatment) and AD (outcome)

	Estimate	95% CI lower	95% CI upper	p-value
ACME (control)	0.00309	-0.03186	0.04	0.87
ACME (treated)	0.00422	-0.04023	0.05	0.85
ACME (average)	0.00366	-0.03420	0.04	0.85
Total Effect	0.03992	-0.02532	0.11	0.23

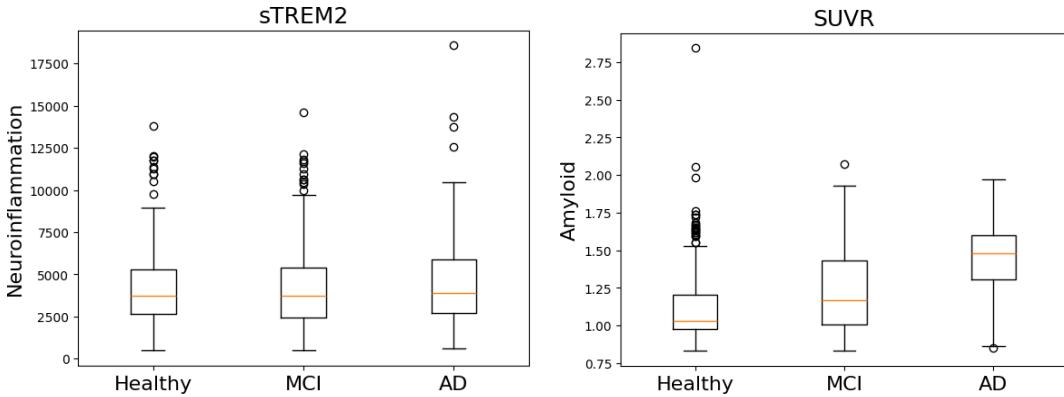


Figure 5: Boxplots of sTREM2 levels among participants (left) and amyloid- β levels among participants (right)

Table 2: ANOVA results

	Sum of squares	df	F	p-value
Diagnosis and STREM2	1.99×10^7	2.0	1.89	0.15
Diagnosis and Amyloid- β	11.65	2.0	98.67	7.57×10^{-40}

Our preliminary results present a more nuanced picture than initially hypothesized, suggesting that the proposed causal relationship involving sTREM2 may not operate as straightforwardly as anticipated. Several factors likely contribute to this discrepancy, underscoring the intricate role of sTREM2 in AD pathology.

Firstly, sTREM2 exhibits a complex, non-linear relationship with AD progression. As reported by [8], sTREM2 levels are not static; they often increase during the early stages of AD but tend to decline in more advanced phases. Further illustrating this complexity, another study by [10] observed a decrease in sTREM2 levels during the earliest asymptomatic stages if amyloid pathology was present without concomitant tau pathology or neurodegeneration. These stage-dependent dynamics, along with some contradictory findings in the literature regarding sTREM2 level changes in AD, highlight the challenges in defining a simple linear association.

More importantly, a direct causal link from baseline sTREM2 levels to the initiation of

$A\beta$ accumulation is questionable. Evidence from [8] suggests that alterations in sTREM2 levels are more likely a response to existing neuropathology rather than baseline sTREM2 concentrations directly instigating $A\beta$ deposition. Supporting this perspective, [11] proposed that the sTREM2-associated microglial response occurs subsequent to the initial fibrillization of $A\beta$ and may even facilitate later pathological developments, such as p-tau accumulation, in the early stages of AD. This positions sTREM2 as a component within a developing pathological cascade, rather than a primary trigger for $A\beta$ pathology from baseline.

Furthermore, the actual biological interactions are likely more intricate than initially presumed, involving a broader network of genes beyond TREM2 alone. For instance, TREM2 signaling is critically dependent on adaptor proteins like DAP12 and DAP10, whose encoding genes are themselves integral components of this pathway. Additionally, the transformation of microglia into disease-associated microglia (DAM) states involves significant shifts in the expression of numerous other genes, including Apoe, Lpl, P2ry12, and Tmem119, as detailed by [12].

Given these considerations—the non-linear and stage-specific behavior of sTREM2, its apparent role as a reactive marker to ongoing pathology, and the involvement of a wider array of genes—there is currently no clear consensus on a simple, direct causal relationship between baseline sTREM2 levels and $A\beta$ pathology. Consequently, to better delineate the complex causal pathways implicated in AD, we have decided to employ causal discovery methodologies using genome-wide single-cell gene expression data from human brain samples.

4 Feature Selection

4.1 Overview

In our raw dataset, the gene-expression matrix contains 12,762 genes. To identify a compact set of genes with the strongest predictive power for AD and thus suitable for downstream causal-discovery analyses, we applied three complementary feature-selection methods in tandem: random forest, sparse logistic regression, and gradient boosting. We will now describe the three feature-selection methods applied to the gene-expression data. All three algorithms were implemented using Scikit-Learn (v1.5; [13]).

4.1.1 Data

To investigate causal mechanisms underlying AD development using gene expression data, we obtained single-cell RNA and spatial transcriptomics data from postmortem brain samples from healthy donors and patients diagnosed with AD from the database "A single-cell and spatial RNA-seq database for Alzheimer's disease" (ssREAD) introduced by [14], which compiles datasets from 18 human and mouse brain studies. For our downstream analysis, we focus on single-cell RNA data from microglia in the human prefrontal cortex (PFC), in order to capture the causal networks of the neuroinflammatory compound contributing to disease development. The PFC is also known to play a critical role in

higher cognitive functions and becomes significantly perturbed early in AD. To enhance the robustness of our results, we combined data across multiple studies and applied batch effect correction, centering the data and normalizing by standard deviation prior to further analysis.

The subsequent data processing pipeline is as follows: First, to eliminate class imbalance bias, we downsampled the majority class (healthy controls) so that AD cases and healthy controls appear in a 1:1 ratio. Next, to balance statistical rigour against computational cost, we adopted a repeated subsampling strategy: on each of 20 iterations, we drew a stratified random subset of 1,000 samples (500 AD, 500 controls) and ran all three feature selectors independently. Finally, we set a consistency threshold - retaining only those genes chosen in more than 5 out of 20 runs to ensure that our final gene panel reflects stable, reproducible signals rather than one-off artifacts.

4.1.2 Evaluation Metrics

We will assess the performance of different feature selection algorithms by checking the following metrics: area under the receiver operating characteristic curve (AUROC), area under the precision-recall curve (AUPRC), accuracy, recall, precision and F1-score. The formal definitions of all evaluation metrics can be found in Appendix 7.

Precision quantifies the fraction of predicted positives that are actually correct, while recall measures the fraction of true positives that the model successfully identifies. The F1-score, defined as the harmonic mean of precision and recall, balances these two metrics into a single summary value that is especially useful when class frequencies differ.

To assess performance across all possible classification thresholds, we also generated two complementary curves. The Receiver Operating Characteristic (ROC) curve plots the true positive rate against the false positive rate as the threshold varies; its area under curve (AUROC) reflects how well the model separates positives from negatives overall. The Precision-Recall (PR) curve, by contrast, plots precision versus recall at each threshold and yields the area under the PR curve (AUPRC), which emphasizes performance on the positive (often minority) class. Because AUROC can remain high simply by correctly classifying the abundant negative class, AUPRC is generally more informative for severely imbalanced datasets. In practice, we compute model scores for every example, sweep the threshold from the highest score to the lowest, calculate true/false positives and negatives at each step, derive the corresponding rates (true positive rate, false positive rate, precision, recall), and then interpolate those points to form the continuous ROC and PR curves.

As mentioned, we perform 20 bootstrap iterations to derive a stable set of selected features. In each iteration, we reserve 20% of the data (using a fixed random seed shared across all algorithms) as the test set, fit the model on the remaining 80%, and then compute both the ROC and Precision-Recall curves on the held-out data. After all 20 runs, we report the mean and standard deviation of each metric across thresholds. This procedure ensures that our performance estimates and the resulting AUROC and AUPRC curves reflect both the average behavior of the model and its variability under resampling.

4.2 Methodology of Feature Selection

4.2.1 Random Forest

Random Forest is a supervised learning algorithm that constructs an ensemble of decision trees to improve predictive accuracy and guard against overfitting. Each tree is grown on a different bootstrap sample of the training data and, at each split, considers only a random subset of features when selecting the optimal axis-aligned threshold that minimizes the weighted Gini impurity of the resulting child nodes [15]. Trees continue splitting until a stopping criterion - such as maximum depth, minimum samples per leaf, or perfectly pure nodes - is reached. In classification tasks, each tree casts a vote for the class of a new observation, and the forest's output is determined by majority vote. By combining bagging with feature subsampling, Random Forest decorrelates the individual trees, markedly reducing variance while maintaining strong performance and allowing for straightforward measures of feature importance.

4.2.2 Sparse Logistic Regression

Logistic regression is a discriminative classifier. It maps the input feature vector \mathbf{x} to the predicted class labels y through the probability $P(y|x)$. We posit a linear score

$$z = \mathbf{w}^\top \mathbf{x} + b,$$

where w is the weight vector and b is the bias (intercept). The predicted probability that $y = 1$ is obtained using the sigmoid function:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

To make the model more selective and reduce the risk of overfitting, we added a ℓ_1 penalty to the objective function - a technique commonly referred to as sparse logistic regression. This regularization encourages sparsity by driving more weight coefficients to zero, which is especially beneficial in high-dimensional settings.

Assuming independent samples $\left\{(\mathbf{x}^{(i)}, y^{(i)})\right\}_{i=1}^n$, the final loss function used is the binary cross entropy with ℓ_1 regularization:

$$\mathcal{L}\left(\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n, \mathbf{w}\right) = -\frac{1}{n} \sum_{i=1}^n \left[y^{(i)} \ln \sigma(z^{(i)}) + (1 - y^{(i)}) \ln(1 - \sigma(z^{(i)})) \right] + \lambda \sum_{j=1}^d |w_j|$$

We apply gradient descent to update the weights w_j and bias b . After training, we predict class labels by applying a threshold of 0.5:

$$\hat{y} = \begin{cases} 1, & \sigma(\mathbf{w}^\top \mathbf{x} + b) \geq 0.5, \\ 0, & \text{otherwise.} \end{cases}$$

Choosing the regularization strength λ is critical to striking the right balance between bias and variance in an ℓ_1 -penalized logistic regression model. To avoid arbitrary selection

and to ensure that our model generalizes well, we employ a three-fold cross-validation procedure combined with a grid search over candidate λ values. Specifically, we partition the training data into three equal folds and, for each λ in our grid, we train the model on two folds and evaluate its performance on the remaining fold. We repeat this process so that each fold serves once as the validation set, then average the chosen performance metric across all three validations. The value of λ that achieves the highest average score is deemed optimal, and we then retrain the final model on the entire training set using that selected λ . This approach ensures that our regularization strength is tuned in a data-driven way, guarding against both overfitting and underfitting.

4.2.3 Gradient Boosting

Gradient Boosting is a supervised learning algorithm that builds an ensemble of decision trees sequentially, with each tree correcting the errors of the previous ones. Unlike Random Forest, where trees are built independently, Gradient Boosting fits each new tree to the negative gradient of a loss function, effectively focusing on instances where the current model performs poorly. This iterative approach improves performance by minimizing the loss at each step [15].

At each step m , we compute the pseudo-residuals r_{im} , the negative gradients of the loss function $L(y_i, F_{m-1}(x_i))$ with respect to the current predictions $F_{m-1}(x_i)$. A new decision tree $h_m(x)$ is then trained to predict these residuals, targeting the instances where the ensemble performs poorly. Its contribution is scaled by a learning rate η and added to the model:

$$F_m(x) = F_{m-1}(x) + \eta h_m(x)$$

To prevent overfitting, we reserved a validation set and applied early stopping based on validation loss. The learning rate was also tuned to maximize AUROC.

4.3 Feature Selection Results and Discussion

4.3.1 Random Forest

Figure 6 shows a bar plot of the genes most frequently selected by the Random Forest algorithm. A histogram of all genes selected at least once is presented in Figure 7. The corresponding ROC and PR curves are presented in Figure 8.

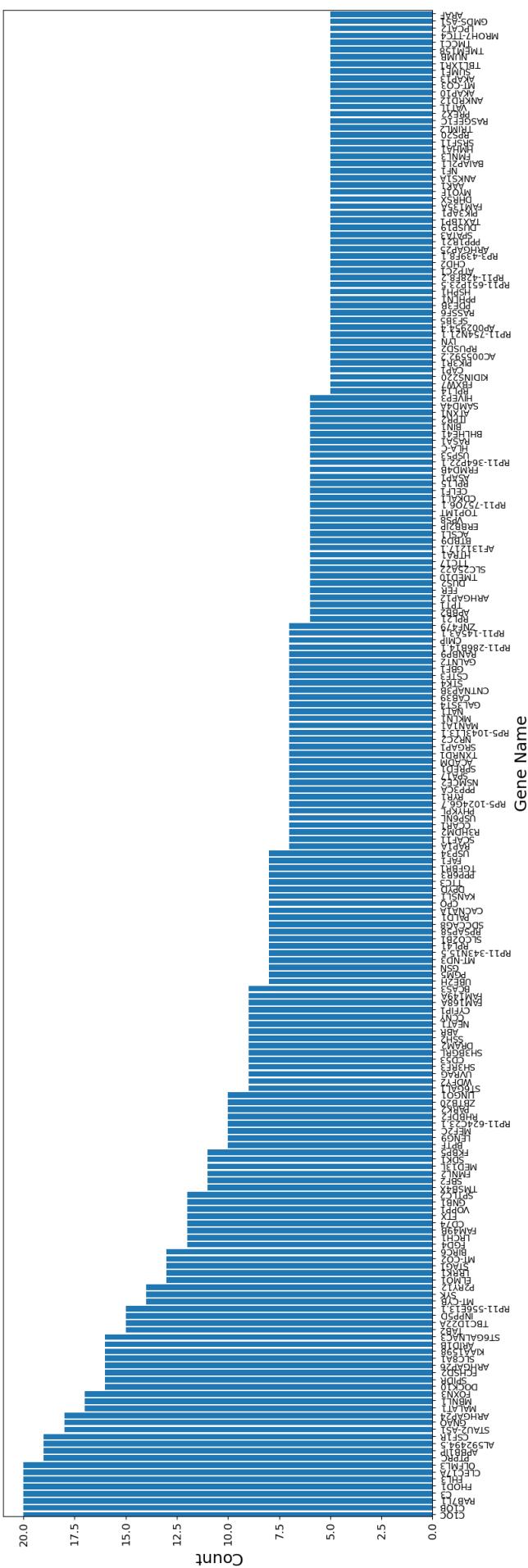


Figure 6: Genes Repeatedly Selected (≥ 5 Times) by Random Forest.

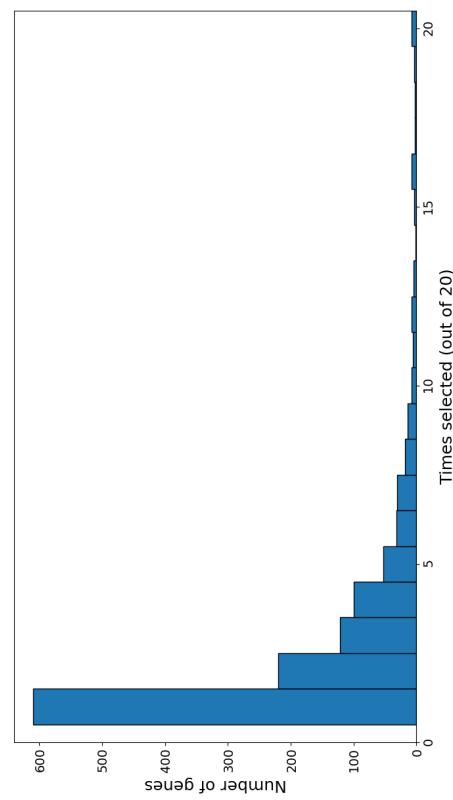


Figure 7: Histogram of gene selection counts by Random Forest (genes selected at least once)

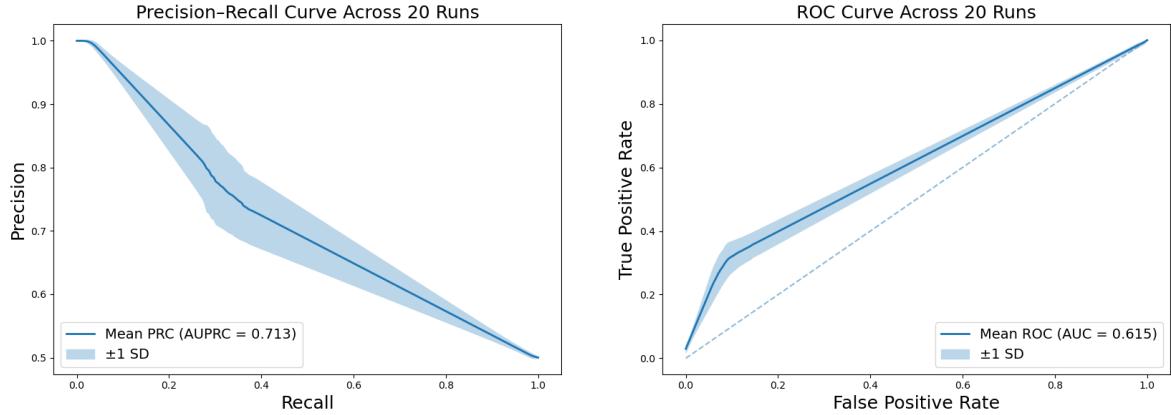


Figure 8: Precision-recall curve (left) and ROC curve (right) of Random Forest

4.3.2 Sparse Logistic Regression

Figure 9 shows a bar plot of the genes most frequently selected by the Random Forest algorithm. A histogram of all genes selected at least once is presented in Figure 10. The corresponding ROC and PR curves are presented in Figure 11.

Figure 9: Genes Repeatedly Selected (≥ 5 Times) by the Sparse Logistic Regression

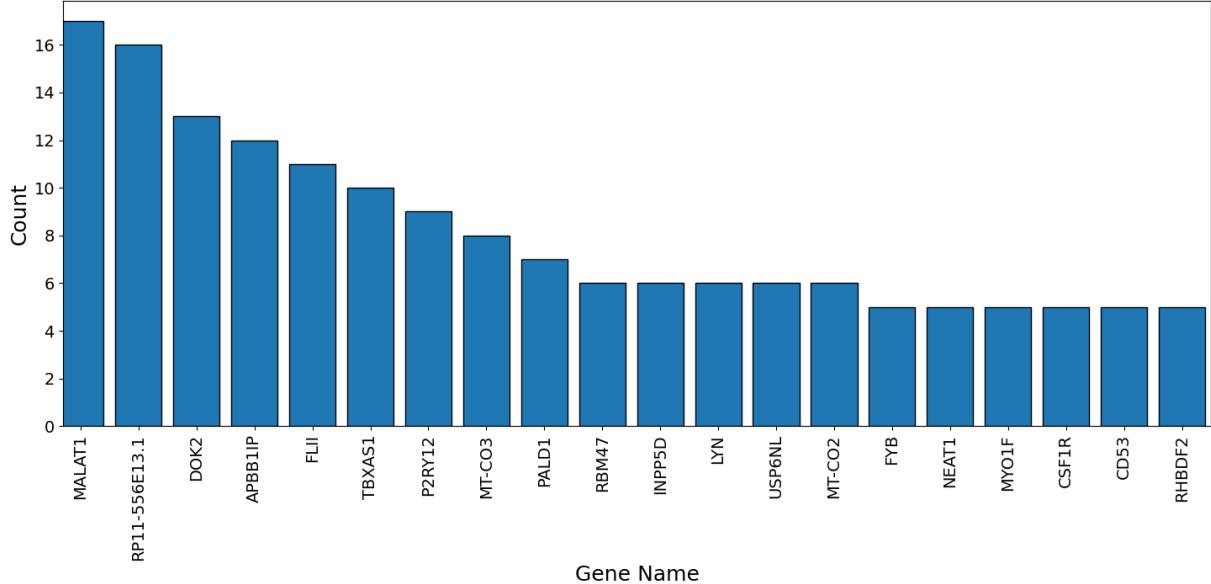


Figure 10: Histogram of gene selection counts by Sparse Logistic Regression Algorithm (genes selected at least once)

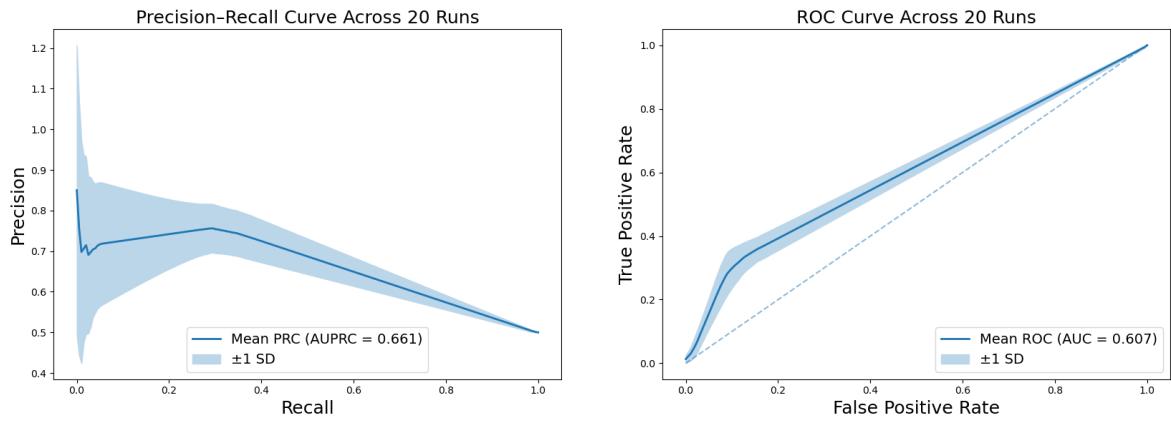
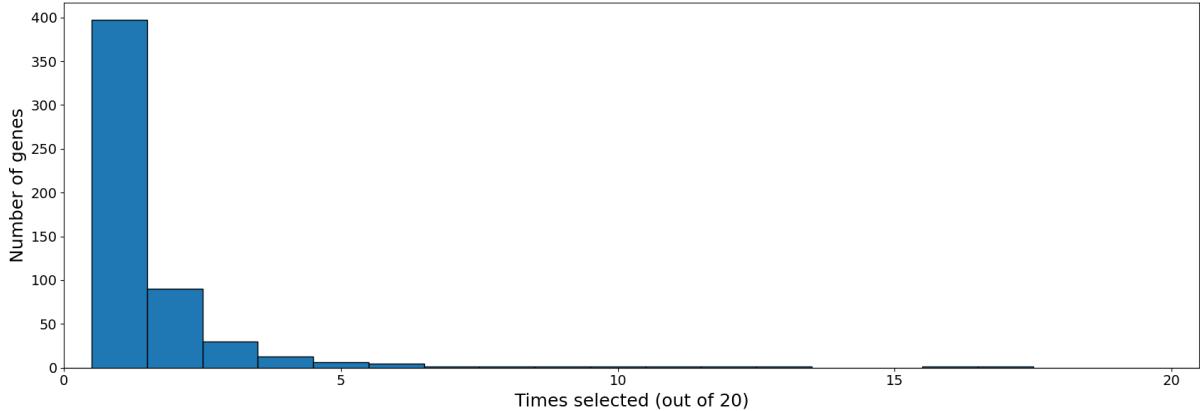


Figure 11: Precision-recall curve (left) and ROC curve (right) of Sparse Logistic Regression

4.3.3 Gradient Boosting

Figure 12 shows a bar plot of the genes most frequently selected by the Random Forest algorithm. A histogram of all genes selected at least once is presented in Figure 13.

Figure 12: Genes Repeatedly Selected (≥ 5 Times) by the Gradient Boosting

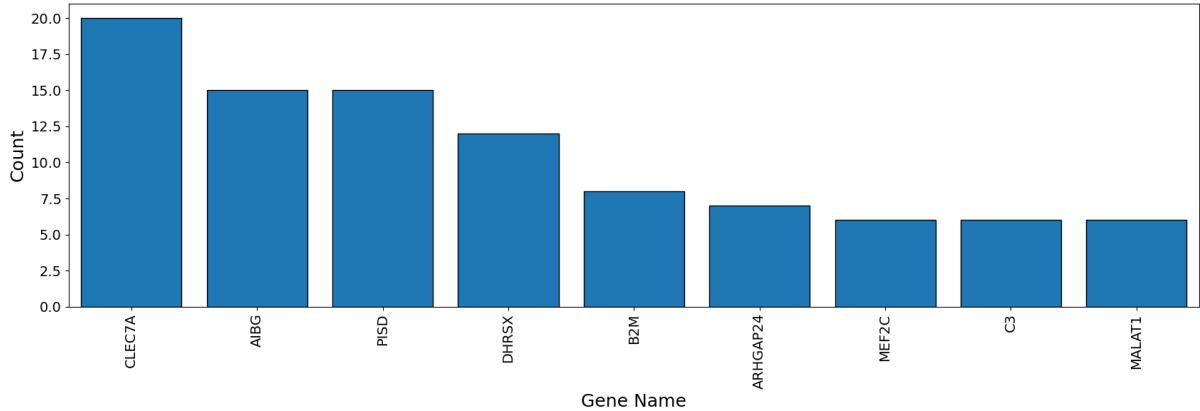
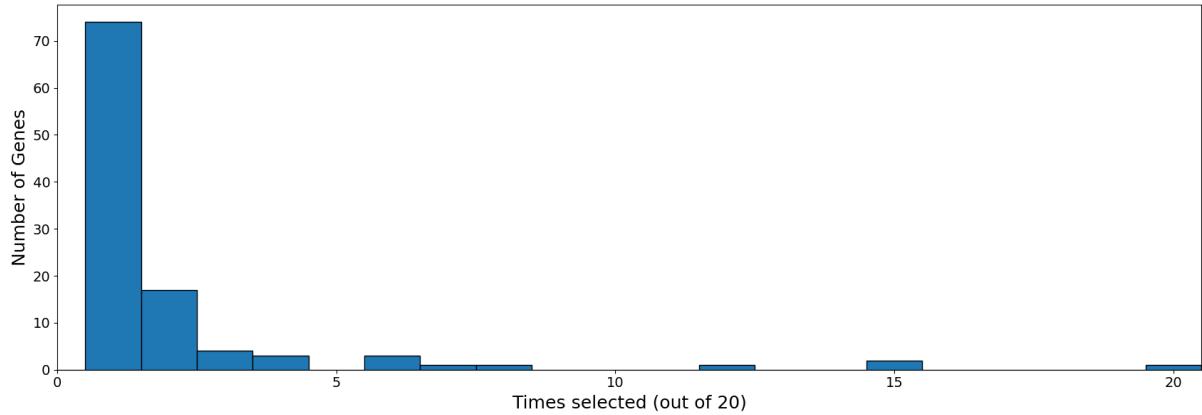


Figure 13: Histogram of gene selection counts by Gradient Boosting (genes selected at least once)



The corresponding ROC and PR curves are presented in Figure 14.

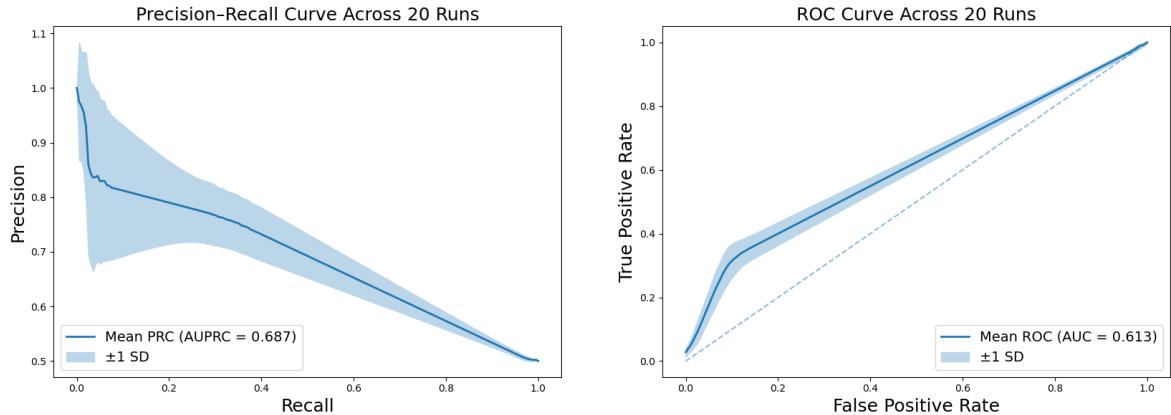


Figure 14: Precision-recall curve (left) and ROC curve (right) of gradient boosting

4.3.4 Performance of Feature Selection Algorithms

Table 3: Performance comparison of Random Forest, Logistic Regression, and Gradient Boosting. Results are in the format of mean (std).

Model	Model Specific Parameters	AUROC	AUPRC	Accuracy	Precision	Recall	F1
Random Forest	Gini impurity: 0.001	0.615 _(0.022)	0.602 _(0.020)	0.612 _(0.020)	0.775 _(0.047)	0.315 _(0.028)	0.447 _(0.034)
Logistic Regression	Cross-validation: 3 Inverse-regularization: 10^{-3} to $10^{-0.3}$ Convergence tolerance: 5×10^{-4}	0.607 _(0.024)	0.584 _(0.024)	0.607 _(0.023)	0.752 _(0.054)	0.321 _(0.029)	0.449 _(0.035)
Gradient Boosting	Early-stopping: True Learning rate: 0.1 Convergence tolerance: 5×10^{-4}	0.612 _(0.022)	0.596 _(0.020)	0.612 _(0.021)	0.759 _(0.047)	0.328 _(0.031)	0.457 _(0.036)
Baseline genes		0.380 _(0.0419)	0.509 _(0.034)	0.575 _(0.028)	0.752 _(0.082)	0.225 _(0.042)	0.345 _(0.056)

To assess our models’ added value, we compare them against a simple baseline: a logistic regression classifier trained only on the two most established AD biomarkers, APOE and TREM2. The baseline’s ROC and precision-recall curves are presented in Appendix 7.

The results are summarized in Table 3. It can be seen that all three machine learning models for feature selection demonstrate a clear improvement over the baseline model using only APOE and TREM2. This suggests that the additional features and the learning algorithms are capturing meaningful predictive signals from the data. The baseline model, relying solely on APOE and TREM2, offers minimal discriminative power in our dataset. In comparison, machine learning-based approaches achieve AUROC values around 0.61 and AUPRC values ranging from 0.59 to 0.60, indicating a significant improvement over the baseline of approximately 10-15%.

Table 4: Precision, recall, and F1 of Random Forest, Logistic Regression, and Gradient Boosting with optimized decision threshold. Results are in the format of mean (std).

Model	F1	Recall	Precision
Random Forest	0.644 _(0.069)	0.932 _(0.023)	0.518 _(0.055)
Logistic Regression	0.666 _(0.001)	0.997 _(0.002)	0.500 _(0.001)
Gradient Boosting	0.639 _(0.074)	0.916 _(0.025)	0.519 _(0.056)

Among the three models, performance differences are minimal. Gradient Boosting slightly outperforms the others across most metrics, while Random Forest and Sparse Logistic Regression show nearly identical results, marginally trailing Gradient Boosting. Notably, Gradient Boosting also required less computational time than Sparse Logistic Regression, making it the most promising model in terms of both performance and efficiency.

A significant limitation across all models is the low recall, which ranges from 0.315 to 0.328. This means the models correctly identify only about 32% of actual positive cases, a concerning outcome in the context of AD diagnosis, where high false negative rates are unacceptable. On the other hand, the models exhibit relatively high precision, suggesting that when a positive case is predicted, it is likely correct. This highlights a trade-off: while precision is strong, the models fail to capture the majority of true positive cases, underscoring the need for further improvement in sensitivity.

We optimized the decision threshold for each model by sweeping from 0.0 to 1.0 and selecting the cutoff that maximized F1 while maintaining average precision ≥ 0.50 . Although the default threshold is 0.5, the grid search identified 0.42 as the optimal cutoff for all three algorithms. At this new operating point, recall and F1 both improve markedly, with precision still held above our 0.5 floor. The recalibrated performance metrics are shown in Table 4. Note that changing the decision threshold will not alter the ROC or PR curves; it simply moves to a different point along each curve.

Another important feature to highlight is the characteristic shape of the ROC curves observed across all three algorithms (Figure 8, 11, 14), which consistently display a pronounced kink followed by a long, straight segment. This is because each test fold contains only 200 positives and 200 negatives: every false positive increases the false positive rate

by 0.005, and every true positive increases the true positive rate by the same increment. In theory this yields up to 201×201 possible points, but most probability thresholds collapse into just a handful of distinct steps. As a result, the curve typically exhibits one large initial jump, corresponding to the handful of easy positives that the model ranks above all negatives, followed by a long segment with slope almost equal to 1 in which lowering the threshold admits true and false positives at roughly the same rate. This combination of a pronounced kink and an extended straight line tail is therefore a direct consequence of the limited test set size and the discrete nature of the model's output scores.

Similarly, the averaged PR curve exhibits a sharp kink followed by a long, nearly linear descent (Figure 8, 11, 14). At a very high threshold, only a few of the easiest positives are identified, so recall jumps from zero to a small value in a single step, and precision is initially very high. This forms the initial kink. As the threshold is lowered beyond that first group, each additional true positive is typically accompanied by a false positive, since the remaining scores do not clearly separate positives from negatives. As a result, precision decreases almost linearly as recall increases. When all test points are eventually included (recall = 1), precision converges to the overall prevalence, which explains why the curve ends at (1, 0.5).

The wide shaded band at low recall reflects the high variability across runs in capturing those few easy positives - small differences in predictions can cause large swings in precision when only a handful of examples are involved. Adding or missing a single false positive can shift precision by tens of percentage points. As recall increases and the average is taken over hundreds of predictions, these fluctuations diminish, and the band narrows, resulting in a smooth and predictable decline toward the baseline.

4.4 Biological Interpretation of Findings

There are 6 genes selected by all three algorithms: **ARHGAP24**, **C3**, **CLEC17A**, **DHRSX**, **MALAT1**, and **MEF2C**. Their specific roles and interactions in AD will be explored in the causal discovery analysis that follows.

To identify the biological relevance of features selected by each of the three algorithms, we performed Gene Ontology analysis. This analysis reveals the biological processes, cellular components, and molecular functions significantly impacted in the studied condition (see Appendix 7 for detailed plots). Genes selected by gradient boosting are primarily involved in modulating and executing immune system processes, with B2M and MEF2L identified as key players. Features chosen by sparse logistic regression are critically implicated in monocyte and leukocyte differentiation, chemotaxis, and migration; key genes include DOCK8, INF2, P2RY12, LYN, MALAT1, NEAT1, INPP5D, and CSMD1. Finally, genes selected by random forest are significantly involved in organizing and regulating the cytoskeleton and Rho protein signaling, with FHOD1, FCI, C1QTNF1, TRPC4, ARHGAP24, and CSMD1 as likely key contributors. These identified genes play an active role in cell motility, morphogenesis, and potential cell-cell or cell-environment interactions.

Overall, the gene ontology identified biological processes highlight key pathological pillars of AD: neuroinflammation, immune dysregulation involving both brain-resident and

infiltrating immune cells, and disruptions in neuronal and glial cell biology, particularly cytoskeletal integrity and dynamic cellular responses.

It is worth noting that the baseline performance suggests TREM2 and APOE exhibit little to no predictive power for AD, an unexpected result given their well-established biological relevance. This discrepancy may stem from several factors. For APOE, AD risk is primarily linked to genetic variants - namely, the $\varepsilon 2$, $\varepsilon 3$, and $\varepsilon 4$ alleles [16] - rather than differences in gene expression. The $\varepsilon 4$ variant, in particular, increases AD risk in a dose-dependent manner. Consequently, using APOE expression in microglia fails to capture the true genetic risk conferred by $\varepsilon 4$. As for TREM2, our analysis of its levels in cerebrospinal fluid (CSF) showed no statistically significant differences across AD diagnostic stages. This lack of separation suggests limited prognostic value, which aligns with its minimal predictive power observed in the baseline model. In addition, our dataset aggregates single-cell RNA sequencing data across multiple studies, introducing variability due to inconsistent protocols and normalization procedures. Finally, by focusing exclusively on microglia, we may have overlooked disease-relevant expression signatures in other cell types, further limiting our feature selection and overall model performance.

5 Causal Discovery

5.1 Overview

After identifying the most predictive genes for AD, we retained only their corresponding expression data and used it as input for the downstream causal discovery models. Since each feature selection algorithm yielded a different set of selected genes, we applied the causal discovery algorithms separately to each set. Additionally, to illustrate the impact of data dimensionality on causal discovery performance, we present results in the order of increasing number of selected genes: Gradient Boosting, Sparse Logistic Regression, and Random Forest.

A causal relationship between two variables, A and B, exists when intervening on A (the cause) induces a change in B (the effect) under the same conditions. Causal discovery is the process of inferring the structure of these relationships from data by constructing a causal graph - typically a directed acyclic graph (DAG). In such a graph, each node denotes a variable of interest, and each directed edge $X \rightarrow Y$ signifies that X is a direct cause of Y . Together, the DAG's pattern of connections encodes both direct and indirect causal influences across the system.

Various methodological approaches have been proposed, including constraint-based methods, score-based methods, and functional causal models. In this study, we employed the constraint-based PC algorithm, the score-based Greedy Equivalence Search (GES) algorithm, and the Linear Non-Gaussian Acyclic Model (LiNGAM), a representative functional causal model. The functional model was further integrated with a variational autoencoder (VAE) framework, enabling the recovery of causal relationships under less restrictive assumptions.

We now describe in detail the causal discovery methods applied to the gene expression data of the selected genes. The PC, GES and LiNGAM algorithms are implemented with

the Causal-Learn package [17].

5.2 Methodology of Casual Discovery

5.2.1 PC Algorithm

The PC algorithm [18] is a constraint-based procedure for recovering the Markov equivalence class of a causal DAG from purely observational data. It is founded on four key assumptions. First, the Causal Markov Condition requires that each variable be independent of its non-descendants given its direct parents. Second, Faithfulness demands that every and only the conditional independencies in the data correspond to d-separation relations in the true graph. Third, Causal Sufficiency assumes there are no unmeasured common causes among the observed variables, so that all confounding arises from measured nodes. Finally, acyclicity stipulates that the causal structure contains no directed cycles. The PC algorithm proceeds in three phases:

- Initialize with a complete undirected graph over all variables. Iteratively, for each adjacent pair $X - Y$, search for a conditioning set S among their current neighbors that renders $X \perp\!\!\!\perp Y | S$. Whenever such a S of size n is found, remove the edge $X - Y$ and record S as the sepset for (X, Y) . Increase n and repeat until no further edges can be deleted.
- For each triplet $X - Y - Z$, where X and Z are not directly connected, if Y is not in the sepset set of (X, Z) , orient Y as a collider: $X \rightarrow Y \leftarrow Z$. Otherwise, leave the edges undirected.
- We repeatedly apply the following orientation rules until convergence: If $A \rightarrow B$ and $B - C$ is undirected (with A and C nonadjacent), orient $B - C$ as $B \rightarrow C$. If there exists a directed path $A \rightarrow \dots \rightarrow B$ and an undirected edge $A - B$, orient it as $A \rightarrow B$.

We assessed conditional independence using Fisher's test. Applying PC to gene-expression data faces several key challenges. First, the dimensionality of genomic datasets makes the sheer number of conditional-independence tests both computationally prohibitive and statistically underpowered. Second, the assumption of causal sufficiency is quite strong and is often violated in biological contexts. Third, the algorithm recovers only a Markov-equivalence class of DAGs, so many distinct causal structures remain indistinguishable. Finally, PC's accuracy hinges on reliable independence tests - often assuming approximate multivariate Gaussianity, which rarely holds across the skewed, heavy-tailed distributions typical of gene-expression measurements. Moreover, unlike other approaches, the PC algorithm does not estimate effect sizes and therefore cannot quantify the strength of causal relationships.

5.2.2 GES Method

The Greedy Equivalence Search (GES) algorithm, proposed by [19], is a two-stage score-based causal discovery method that operates on Markov equivalence classes. It alternates

between adding and removing edges to maximize a model selection score, typically the Bayesian Information Criterion (BIC).

In the forward phase, GES starts from an empty graph and iteratively considers all eligible edges $X - Y$ not present in the DAG. For each, it evaluates all acyclic orientations, selects the one that yields the greatest score improvement, and adds it. This process continues until no further single-edge insertion improves the score. In the backward phase, GES considers each existing edge, evaluates the score after its removal, and deletes the edge that provides the largest score gain. This continues until no further deletion improves the score.

There are several possible reasons why GES fails to extract sufficient edges from the gene expression matrix. Gene expression data are typically high-dimensional and represent densely connected graphs, which pose challenges for score-based methods. The greedy search heuristic used by GES can easily become trapped in suboptimal local maxima. Moreover, as noted by [20], limited sample sizes can make likelihood-based scores noisy and unreliable. Since GES modifies one edge at a time to improve the score, an early suboptimal change in a dense graph can commit the algorithm to a poor equivalence class, from which it cannot recover.

5.2.3 LiNGAM Method

The LiNGAM method, introduced by [21], assumes the absence of unobserved confounders and requires the disturbance terms to follow non-Gaussian distributions. In the context of gene expression data, such disturbances may arise from intrinsic transcriptional noise, such as bursts in mRNA production, fluctuations in cellular states, or technical measurement errors. The disturbance term is e_i and c_i is an optional constant term.

Assume the observed variables x_i , where $i \in \{1, \dots, m\}$ can be arranged in a causal order denoted by $k(i)$. The value assigned to each variable x_i is modelled as

$$x_i = \sum_{k(j) < k(i)} b_{ij} x_j + z_i + c_i.$$

In vector form, the system of equations can be written as

$$\mathbf{x} = \mathbf{B}\mathbf{x} + \mathbf{z},$$

The solution of x can be written in form of

$$\mathbf{x} = \mathbf{A}\mathbf{z},$$

where

$$\mathbf{A} = (\mathbf{I} - \mathbf{B})^{-1} \tag{1}$$

This formulation essentially reduces to an unmixing problem, which can be addressed using Independent Component Analysis (ICA). The canonical ICA model is expressed as $\mathbf{x} = \mathbf{As}$, where the objective is to recover both the mixing matrix A and the latent, statistically independent source components s . The algorithm proceeds as follows:

1. We first mean-center the gene-by-sample data matrix \mathbf{X} , then whiten it to obtain uncorrelated components with unit variance. This involves computing the sample covariance

$$\text{Cov}(\mathbf{X}) = \mathbf{E}\Lambda\mathbf{E}^T,$$

where \mathbf{E} contains the eigenvectors and Λ is the diagonal matrix of eigenvalues. The whitened data matrix is

$$\mathbf{X}_{\text{white}} = \Lambda^{-1/2}\mathbf{E}^T\mathbf{X},$$

satisfying $\text{Cov}(\mathbf{X}_{\text{white}}) = \mathbf{I}$.

2. After whitening, the remaining mixing corresponds to an unknown orthogonal rotation. We estimate an orthogonal matrix Q that maximizes a non-Gaussianity measure across the rows of $Q\mathbf{X}_{\text{white}}$. The resulting unmixing matrix is

$$\mathbf{W} = Q\Lambda^{-1/2}\mathbf{E}^T.$$

This two-step process defines the standard ICA workflow.

3. ICA returns the rows of \mathbf{W} in arbitrary order. We identify the unique row permutation matrix P that avoids zeros on the main diagonal, typically by minimizing $\sum_i 1/|\widetilde{W}_{ii}|$. Each row of $\widetilde{\mathbf{W}}$ is then normalized by its diagonal element to obtain $\widetilde{\mathbf{W}}'$ with unit diagonal. The estimated matrix of \mathbf{B} is $\widehat{\mathbf{B}} = \mathbf{I} - \widetilde{\mathbf{W}}'$.
4. The causal ordering is determined by finding a permutation matrix Π such that $\Pi\mathbf{B}\Pi^T$ is as close as possible to strictly lower-triangular. This Π defines the inferred gene ordering.

5.2.4 VAE Method

We also implemented the deep, graph-based generative framework of [22] to handle richer data distributions and capture nonlinear parent–child mappings. Building on our structural equation (1), we parameterize both the noise transformation and the final reconstruction with neural networks. Our VAE decoder takes a latent variable Z through a two-stage mapping:

$$X = f_2((I - B)^{-1}f_1(Z)) \tag{2}$$

where $f_1 : R^m \rightarrow R^m$ models potentially nonlinear effects in the disturbance space, and $f_2 : R^m \rightarrow R^d$ (with d the data dimension) reconstructs the observed X .

The encoder mirror images the decoder’s two-stage mapping, but in reverse:

$$Z = f_4((I - B)f_3(X)), \tag{3}$$

where $f_3 : R^d \rightarrow R^m$ first transforms the data into the “disturbance” space, the linear operator $(I - B)$ then unmixes according to our DAG weights, and $f_4 : R^m \rightarrow R^m$ completes the mapping into the latent code Z .

Under this architecture, the decoder defines the likelihood $p_\theta(X | Z)$ and the encoder defines the approximate posterior distribution $q_\theta(Z|X)$, enabling us to learn both the DAG weights B and the neural-network parameters θ jointly via the usual VAE evidence lower bound (ELBO), given by

$$\mathcal{L}_{\text{ELBO}} \equiv -D_{\text{KL}}(q(Z | X) \| p(Z)) + E_{q(Z|X)}[\log p(X | Z)]. \quad (4)$$

A schematic representation of the architecture is given in Figure 25. MLP in the figure refers to a multilayer perceptron.

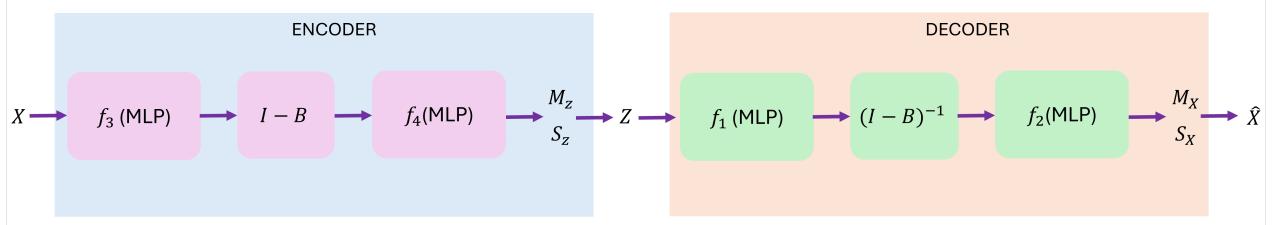


Figure 15: Graph-based neural architecture of the Variational Autoencoder

An essential component of our GNN-based causal discovery framework is the explicit enforcement of acyclicity. From algebraic graph theory [23], it is known that for a non-negative adjacency matrix B , the (i, j) -th entry of B^k is strictly positive if there exists at least one directed path of length k from node i to node j . In particular, the diagonal entries can be used to detect cycles: $(B^k)_{ii} > 0$ if and only if there exists a closed walk of length k that starts and ends at node i .

Leveraging this property, we impose a more practical and computationally convenient constraint. As introduced by [24], for any $\alpha > 0$, a graph is acyclic if and only if

$$\text{tr}[(I + \alpha B)^m] - m = 0 \quad (4)$$

where $B \in R^{m \times m}$. The proof is provided in the Appendix 7. The learning problem thus becomes:

$$\begin{aligned} \min_{B, \theta} \quad & f(B, \theta) \equiv -\mathcal{L}_{\text{ELBO}} \\ \text{s.t.} \quad & h(B) \equiv \text{tr}[(I + B)^m] - m = 0, \end{aligned} \quad (5)$$

This leads to the following Lagrangian formulation and the new loss function:

$$\mathcal{L}_c(B, \theta, \lambda) = f(B, \theta) + \lambda h(B) + \frac{c}{2} |h(B)|^2, \quad (6)$$

where λ is the Lagrange multiplier and c is the penalty parameter. This optimization problem is solved using the Adam optimizer via gradient descent.

5.3 Causal Discovery Results and Discussion

Using the methodology described above, we generated nine causal graphs by combining three feature selection algorithms with three causal discovery methods. To analyze the resulting graphs, we perform formal graph analysis and identify key genes triggering

and/or mediating microglial action in AD through identification of hubs and bottlenecks. Nodes with a high number of outgoing connections, quantified by degree centrality-are termed hubs [25]. In biological systems, hubs often correspond to key regulatory elements due to their broad influence. Bottlenecks are nodes that lie on a large proportion of shortest paths between node pairs, as measured by betweenness centrality. These nodes serve as critical chokepoints for information flow within the network. As an illustrative example, a sample is shown in Figure 16.

As the random forest algorithm produces a very large causal graph with too many edges to display clearly, we provide only a description of the graph. The corresponding adjacency matrix is available in the GitHub repository.

Each hub or bottleneck gene will be further analyzed if it has not been previously identified as a hub and/or bottleneck by other algorithms.

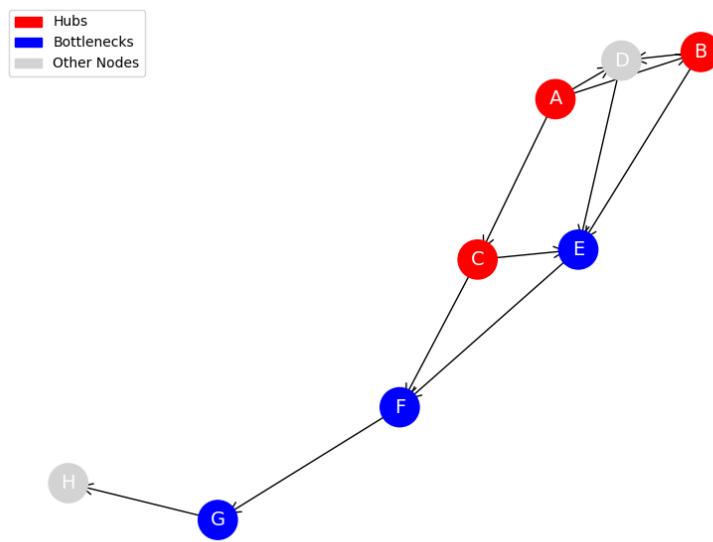


Figure 16: Example network highlighting hub and bottleneck nodes

Graph Analysis

To compare and analyze the DAGs generated by different causal discovery methods, we mainly deployed two graph analysis techniques: by checking the consensus graph and the Structural Hamming Distance (SHD).

The consensus causal graph is constructed by integrating multiple individual causal graphs represented as adjacency matrices, to identify commonly supported relationships. Each input adjacency matrix undergoes a binarization process, where existing connections are standardized to represent the presence of a causal link, while their original strengths are disregarded. The consensus graph is formed by identifying the intersection of these binarized graphs. Because an edge in the consensus graph must be unanimously supported by all contributing causal discovery algorithms or datasets, it signifies a relationship that is stable and consistently detected despite variations in analytical approaches or data samples. These edges are considered more robust.

Structural Hamming Distance (SHD) is a standard metric for assessing the similarity of two causal graphs and is widely adopted in the literature [26]. It measures the minimum number of elementary edits required to transform one graph into another by counting edge deletions (edges present in G_1 but not in G_2), insertions (edges present in G_2 but not in G_1), and orientation flips (when an edge appears in both graphs with reversed direction, e.g. $i \rightarrow j$ versus $j \rightarrow i$, which are treated as two edits).

5.3.1 Causal Discovery on Genes Selected via Gradient Boosting

PC Algorithm on Genes Selected via Gradient Boosting

Figure 17 shows the causal diagram inferred by the PC algorithm from the genes selected via gradient boosting.

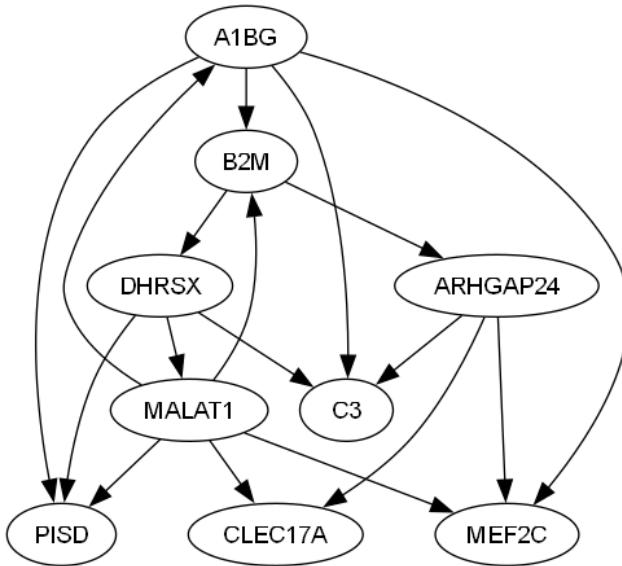


Figure 17: Causal graph constructed using PC algorithm on genes selected via Gradient Boosting

The identified hubs are **MALAT1**, **A1BG**, and **ARHGAP24**, while the bottlenecks are **B2M**, **MALAT1**, and **DHRSX**.

Metastasis-Associated Lung Adenocarcinoma Transcript 1 (MALAT1) is a long non-coding RNA (lncRNA) that regulates gene expression and exhibits neuroprotective functions [27]. It reduces neuronal injury, inhibits apoptosis, and suppresses neuroinflammation in certain neurodegenerative and injury models. MALAT1 often acts as a molecular "sponge" for microRNAs (miRNAs), sequestering them to prevent gene silencing. This regulation influences key survival pathways, such as PI3K/AKT. Its central role in modulating multiple downstream targets and pathways explains its identification as both a hub (high connectivity) and a bottleneck (critical for signal flow) in the causal graph.

Beta-2-Microglobulin (B2M), a component of MHC class I molecules, plays a key role in immune response and is increasingly implicated in AD. B2M is gaining significant attention in AD research. It can co-aggregate with amyloid-beta peptides, enhancing their neurotoxicity, and is found at elevated levels in the brains of AD patients [28].

As B2M can cross the blood-brain barrier, it may link peripheral immune activity with neuroinflammation. Its identification as a bottleneck suggests a central role in amyloid pathology and cognitive decline. Notably, clearing peripheral B2M has shown therapeutic potential in mouse models [29].

Dehydrogenase/Reductase SDR Family Member on Chromosome X (DHRSX) has been increasingly linked to AD. It is a novel protein involved in regulating autophagy - the cellular process for clearing waste [30]. Autophagy dysfunction is a hallmark of AD, contributing to the accumulation of toxic aggregates such as amyloid-beta and tau. The identification of DHRSX as a bottleneck suggests it may play a pivotal role in this pathway; its impairment could hinder waste clearance, exacerbating AD pathology.

The other two hub genes, Rho GTPase Activating Protein 24 (ARHGAP24) and Alpha-1-B Glycoprotein (A1BG), are less well studied and have more indirect links to Alzheimer's disease (AD). ARHGAP24 regulates Rho GTPases, which control the actin cytoskeleton-critical for neuronal structure and synaptic function [31]. Its identification as a hub suggests a potential role in the cytoskeletal remodeling and synaptic dysfunction associated with AD. A1BG is a plasma protein with an unclear primary function, but it interacts with cysteine proteinase inhibitors involved in amyloid precursor protein (APP) processing. By modulating protease activity, A1BG may indirectly affect amyloid-beta production [31]. It also interacts with alpha-2-macroglobulin, a key player in the immune response-another major component of AD pathology [32].

LiNGAM Algorithm on Genes Selected via Gradient Boosting

Figure 18 shows the causal diagram inferred by the LiNGAM algorithm from the genes selected via gradient boosting.

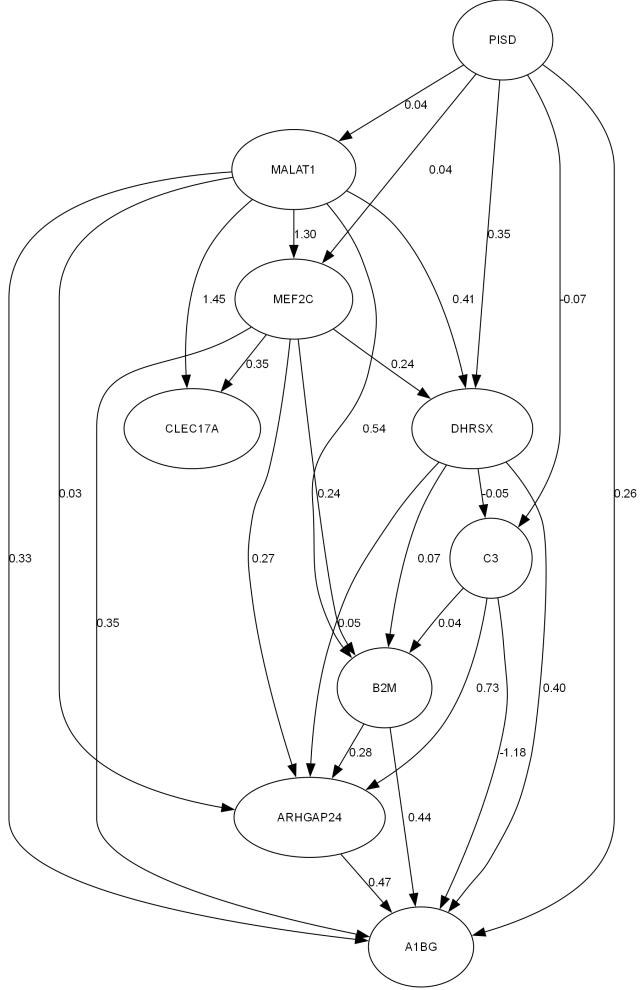


Figure 18: Causal graph constructed using LiNGAM on genes selected via Gradient Boosting

The identified hubs are **MALAT1**, **MEF2C** and **PISD**, while the bottlenecks are **DHRSX**, **MALAT1** and **MEF2C**.

Myocyte Enhancer Factor 2C (MEF2C) is a master-regulator transcription factor essential for neuronal survival, synaptic plasticity, and memory. Its direct relevance to AD is highlighted by the consistent finding that its expression is significantly reduced in patient brains, a decline that correlates with the severity of both amyloid and tau pathology [33]. Reflecting its control over a vast array of downstream genes, MEF2C's status as a hub gene at the center of AD-related biological pathways is identified in other studies as well [34]. This central position means its failure creates a functional bottleneck, where the loss of MEF2C activity triggers a cascade of detrimental effects, including exacerbated amyloid pathology, synaptic failure, and increased oxidative stress.

This critical role, however, also illuminates MEF2C's therapeutic potential. A study by [35] revealed that elevated MEF2C levels are a key feature of cognitive resilience, found in individuals who remain mentally sharp despite harboring significant AD brain pathology. Crucially, they demonstrated that boosting MEF2C expression in a mouse model of neurodegeneration was sufficient to rescue cognitive function and protect neurons.

Phosphatidylserine Decarboxylase (PISD) is an enzyme localized to the inner mitochond-

drial membrane, where it performs an indispensable role in lipid metabolism by producing phosphatidylethanolamine-a key component of mitochondrial membranes. Its function is crucial for maintaining mitochondrial structure and respiratory efficiency. The link between PISD and AD stems directly from this role, as mitochondrial dysfunction is a well-established early feature of AD pathogenesis, disrupting neuronal energy production and calcium balance [36]. While transcriptome-wide network analyses of AD brain tissue do not typically identify PISD as a top hub gene, its classification as such could arise from its centrality within the phospholipid biosynthesis pathway.

VAE Algorithm on Genes Selected via Gradient Boosting

Figure 19 shows the causal diagram inferred by the VAE algorithm from the genes selected via gradient boosting.

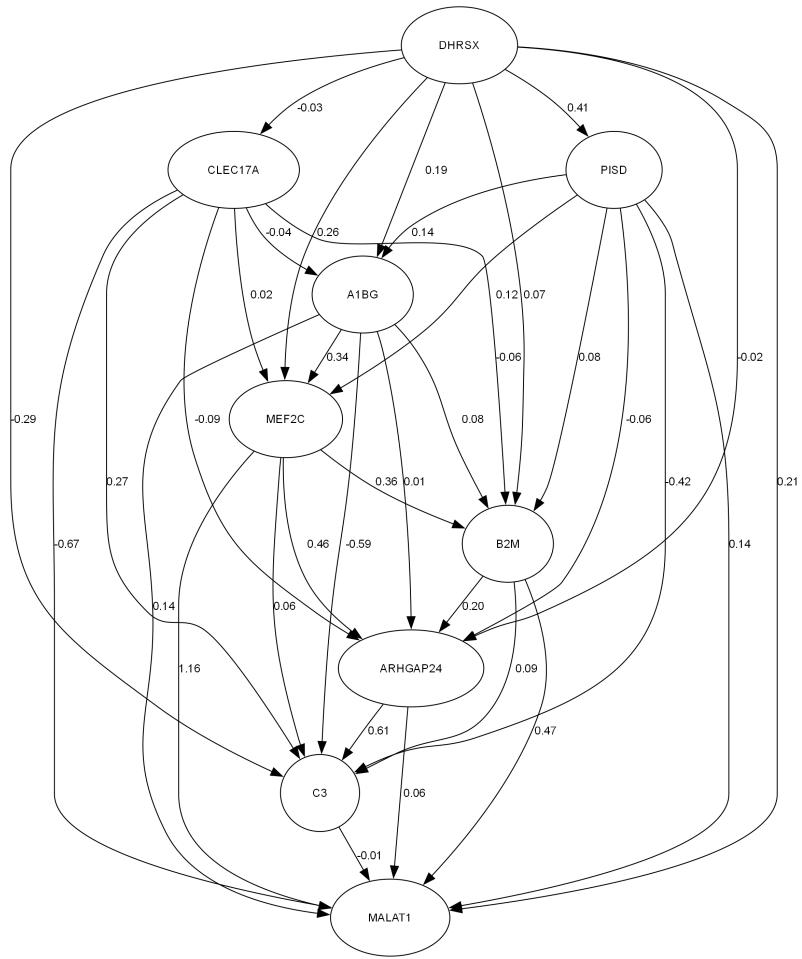


Figure 19: Causal graph constructed using VAE on genes selected via Gradient Boosting

The identified hubs are **DHRSX**, **MEF2C** and **PISD**, while the bottlenecks are **MEF2C**, **CLEC17A** and **B2M**.

C-type Lectin Domain Containing 17A (CLEC17A) is a cell surface receptor protein whose structure places it in the C-type lectin family, which is deeply involved in the immune response. While its specific role in the brain is not yet well-defined, its link to

AD can be plausibly inferred through the mechanism of neuroinflammation. The brain’s resident immune cells, microglia, utilize a wide array of pattern recognition receptors to detect signs of brain injury and disease, including damage-associated molecular patterns released from distressed cells and protein aggregates like amyloid-beta. Receptors in the C-type lectin family are key players in this process [37]. As a bottleneck in a causal network, it means while CLEC17A may not interact with a large number of proteins, it functions as a crucial information gateway. It could be a specific receptor that, upon activation by a damage signal in the AD brain, initiates a distinct downstream inflammatory cascade.

Comparative Analysis of the Three Algorithms

The consensus causal graph is shown in Figure 20, while the SHD comparison is presented in Table 5. The high normalized SHD values suggest that the three methods identified substantially different causal structures.

Table 5: Structural Hamming Distance (SHD) of genes selected via gradient boosting: Results include both absolute and normalized SHD.

Algorithm	Edges	SHD with PC	SHD with LiNGAM	SHD with VAE
PC	17	–	32(72.7%)	26(50.0%)
LiNGAM	27	–	–	40(64.5%)
VAE	35	–	–	–

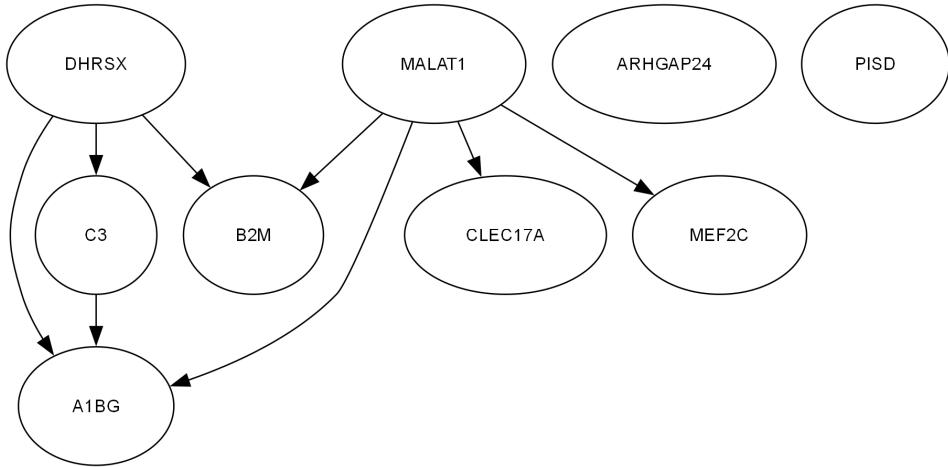


Figure 20: Consensus graph generated from three causal discovery methods applied to genes selected via Gradient Boosting

Notably, some edges exhibit particularly large weights in both the LiNGAM and VAE methods and are also identified by conditional independence testing. These include the edges between MALAT1–MEF2C, MALAT1–CLEC17A, and C3–A1BG.

A study by [35] established a direct regulatory link between MEF2C and MALAT1, identifying MALAT1 as a downstream target gene of MEF2C in excitatory neurons. This finding is highly relevant to AD, where both MEF2C function and MALAT1 levels are

known to be compromised. The downregulation of MEF2C by neuroinflammation in AD would therefore lead to reduced expression of MALAT1, diminishing its neuroprotective effects, including miRNA sponging and blood-brain barrier maintenance.

This regulatory axis may also converge on mitochondrial health, a central pillar of AD pathology. Both MEF2C and MALAT1 are independently linked to promoting neuronal survival and maintaining mitochondrial integrity [38]. Consequently, the loss of MEF2C in AD could exacerbate mitochondrial dysfunction by suppressing MALAT1 expression.

While the literature does not directly link MALAT1 to CLEC17A, a plausible regulatory relationship can be hypothesized. The strongest evidence for this is MALAT1's established ability to regulate other C-type lectins; for instance, it boosts DC-SIGN expression in dendritic cells by sponging miR-155 [39]. This provides a mechanistic precedent for how MALAT1 could similarly influence CLEC17A. In the context of AD, if MALAT1 regulates CLEC17A, the known dysregulation of MALAT1 in AD could lead to aberrant CLEC17A activity. Given CLEC17A's role in immune cell adhesion and communication, this could alter the neuroinflammatory milieu by affecting how these cells interact within the brain.

Similarly although a direct protein interaction between C3 and A1BG in AD has not been established, evidence suggests they may be co-regulated during inflammatory states. For example, a study on cervical neoplasia-a condition marked by inflammation-found that both C3 and A1BG were significantly overexpressed in patient serum [40]. Given that neuroinflammation is a key feature of AD, this finding supports a potential indirect association where C3 and A1BG levels are linked through shared inflammatory pathways.

Summary

Biologically, our causal graph recovered several genes already known to play key roles in AD, including **MALAT1** and **MEF2C**. We also identified genes that have recently attracted growing attention in AD research, such as **B2M** and **DHRSX**. Although **CLEC17A** and **PISD** are not yet well characterized in the context of AD, both participate in biological processes closely linked to AD pathogenesis-innate immune signaling for CLEC17A and mitochondrial phospholipid metabolism for PISD. In the consensus graph, we observed the **MALAT1–MEF2C** connection, a regulatory relationship supported by previous studies. The other two edges similarly suggest plausible regulatory interactions that merit experimental validation. Taken together, these findings highlight both well-established and emerging gene targets and point toward several promising avenues for future research.

5.3.2 Causal Discovery on Genes Selected via Sparse Logistic Regression

PC Algorithm on Genes Selected via Sparse Logistic Regression

Figure 21 shows the causal diagram inferred by the PC algorithm from the genes selected via Sparse Logistic Regression.

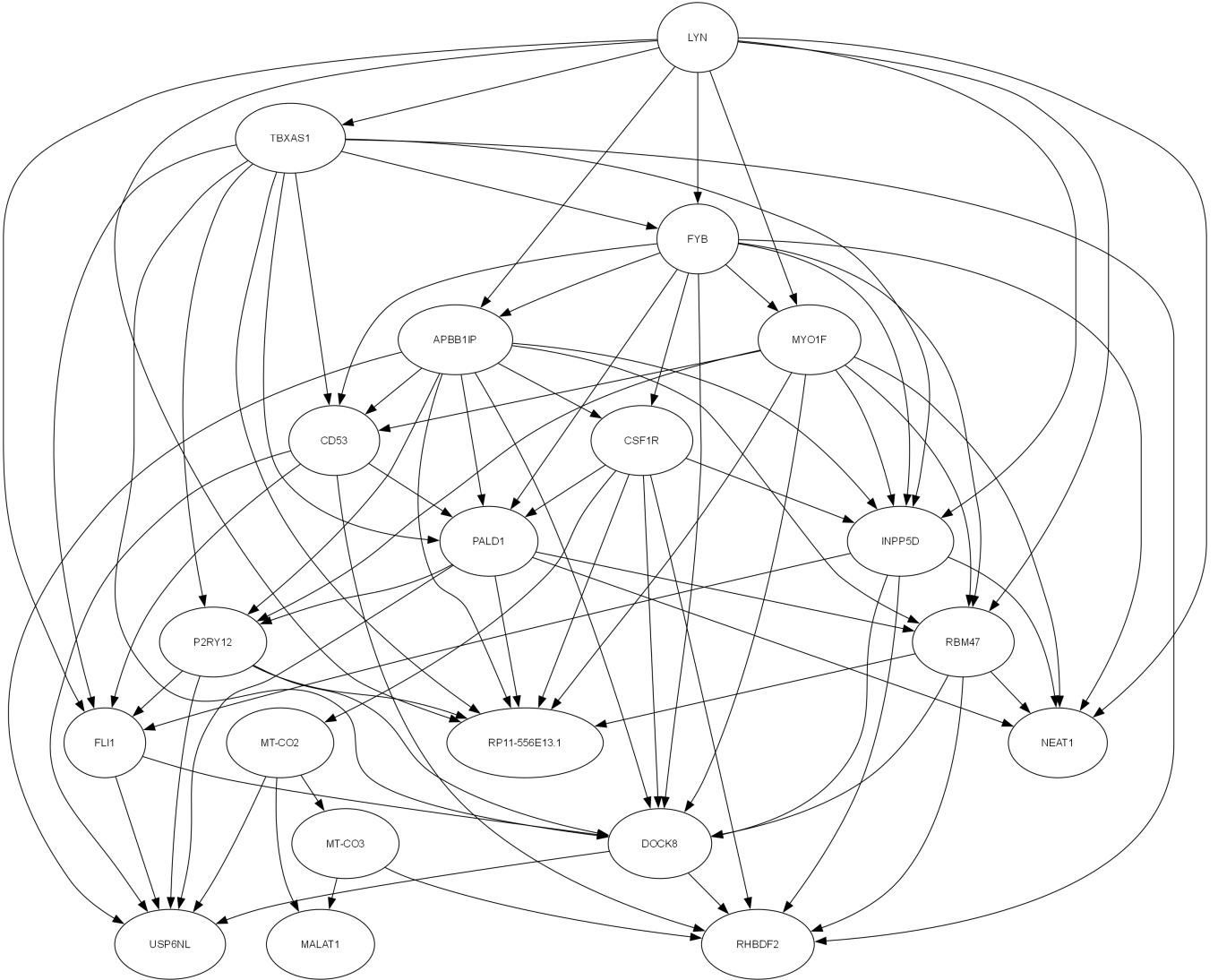


Figure 21: Causal graph constructed using PC algorithm on genes selected via Sparse Logistic Regression

The identified hubs are **APBB1IP**, **FYB**, **LYN**, **TBXAS1** and **MYO1F**, while the bottlenecks are **CSF1R**, **MT-CO2**, **FYB**, **PALD1** and **DOCK8**.

APBB1IP is an adaptor protein that interacts directly with the Amyloid Precursor Protein (APP), the source of the toxic amyloid-beta peptide [41]. It regulates both APP processing and its downstream signaling. Its identification as a hub suggests it is a key organizer that influences the generation of toxic amyloid beta and links the core amyloid pathology to numerous other cellular functions.

The primary role of FYN Binding Protein 1 (FYB) is to bind and regulate FYN kinase, an enzyme with a crucial role in AD. In the AD brain, FYN kinase is activated by amyloid-beta, and it, in turn, phosphorylates the Tau protein. This action links the two main pathologies of the disease and is a key driver of amyloid-induced synaptic damage [42]. FYB's position as both a hub and bottleneck stems directly from this function, placing it at a central intersection that controls the toxic signaling from amyloid-beta to downstream synaptic injury.

LYN, a Src family tyrosine kinase, is a powerful immune regulator that acts as a critical modulator of microglial function in Alzheimer's disease (AD). It interacts with Toll-like receptor 4 (TLR4), a key microglial receptor that recognizes amyloid-beta. A study by [43] revealed a detrimental role for this kinase; its absence in AD mouse models enhanced the ability of microglia to clear amyloid-beta via phagocytosis, leading to reduced neuronal damage. These findings indicate that LYN negatively regulates the protective functions of microglia in response to AD pathology.

Thromboxane A Synthase 1 (TBXAS1) synthesizes thromboxane A2, a potent inflammatory molecule that also causes blood vessel constriction [44]. Its parent pathway (arachidonic acid metabolism) is heavily involved in neuroinflammation [45]. Furthermore, genetic studies have identified a variant in TBXAS1 that is associated with a decreased risk for late-onset AD [46]. As a hub, TBXAS1 likely connects a network of genes involved in both inflammation and blood flow regulation, two processes known to be impaired in the AD brain.

As the master regulator of microglial survival and proliferation, CSF1R (Colony Stimulating Factor 1 Receptor) is essential for maintaining the brain's immune cell population. This critical role has made CSF1R a major therapeutic target in AD, where its inhibition is used as a strategy to deplete or modulate disease-associated microglia [47]. This non-redundant function explains its identification as a bottleneck: the viability of the entire microglial population is dependent on the signaling that flows through this single receptor.

The gene MT-CO2 is a critical component of cellular energy production, encoding a core subunit of Complex IV in the mitochondrial respiratory chain. This context explains why it acts as a network bottleneck in Alzheimer's disease, where mitochondrial failure is a known early event [48]. A defect in MT-CO2 creates a chokepoint in ATP synthesis, crippling the energy supply and leaving neurons, with their high metabolic demands, vulnerable to catastrophic failure.

PALD1 (Paladin 1) is expressed in the endothelial cells that form the blood-brain barrier (BBB). A recent epigenetics study found that the methylation of PALD1 is altered in the entorhinal cortex of AD patients—one of the first brain regions affected by the disease [49]. The BBB is a critical gatekeeper for the brain.

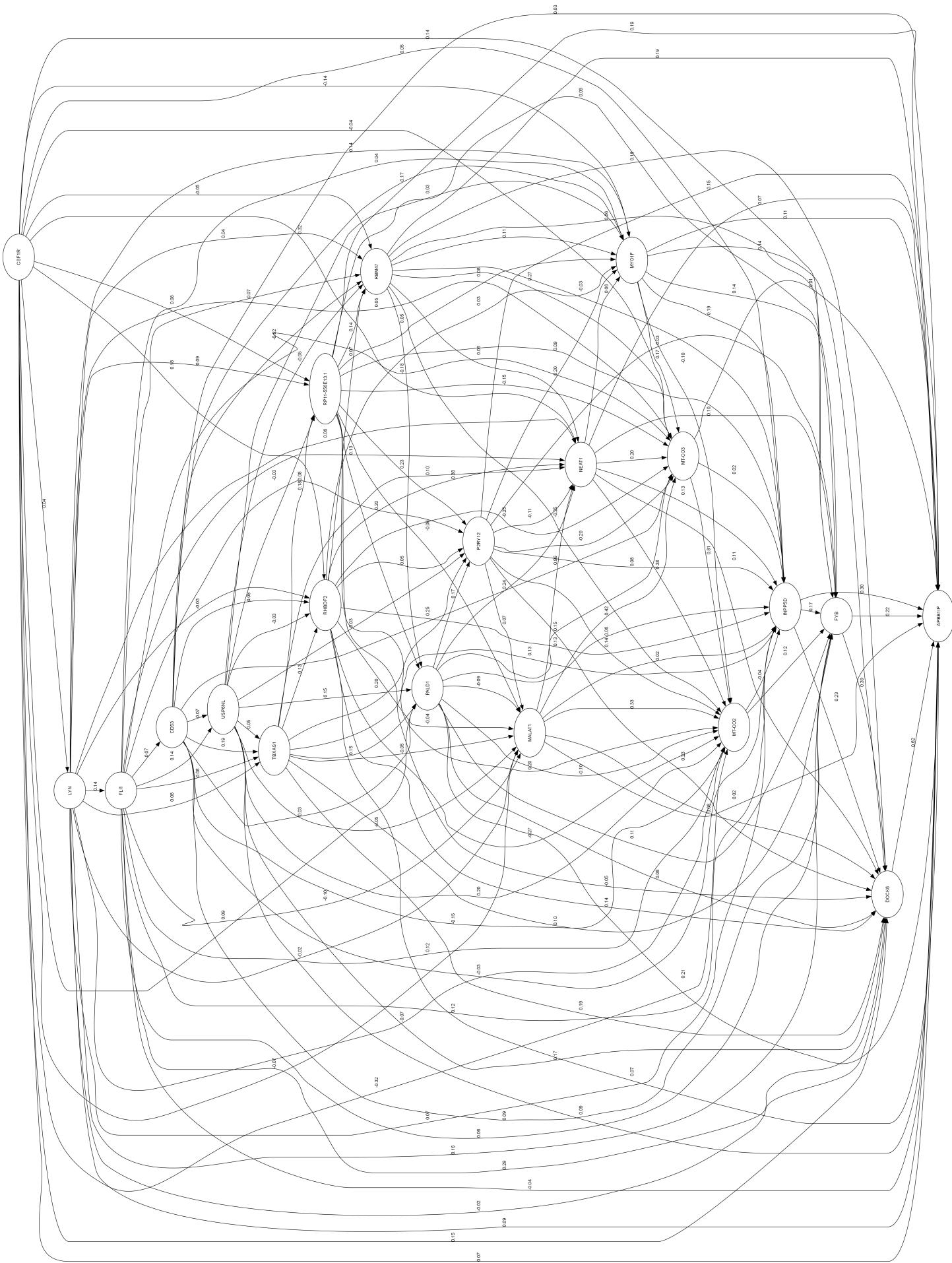
DOCK8 (Dedicator of Cytokinesis 8) is involved in controlling cell migration. Recent studies have revealed that it is expressed in microglia, and its levels are elevated in AD models [50]. Critically, these studies showed that DOCK8 is required for microglia to migrate towards amyloid beta plaques and that knocking down DOCK8 was protective.

MYO1F (Myosin IF) are motor proteins that work with the actin cytoskeleton to generate movement and force. While MYO1F is not extensively studied in AD directly, myosins are essential for microglial function, particularly their ability to move towards, change shape, and engulf debris—a process called phagocytosis [51].

LiNGAM Algorithm on Genes Selected via Sparse Logistic Regression

Figure 22 shows the causal diagram inferred by the LiNGAM algorithm from the genes selected via Sparse Logistic Regression.

Figure 22: Causal graph constructed using LiNGAM on genes selected via Sparse Logistic Regression



The identified hubs are **CSF1R**, **FLI1**, **LYN**, **CD53** and **RP11-556E13.1**, while the bottlenecks are **PALD1**, **FLI1**, **LYN**, **TBXAS1** and **INPP5D**.

As a transcription factor, Friend Leukemia Integration 1 (FLI1) controls the expression of numerous other genes. It is highly expressed in two cell types critical to AD: microglia and the endothelial cells that form the blood-brain barrier (BBB). This is significant because FLI1 regulates inflammatory activation in microglia while also being positioned to control the integrity of the BBB, which is compromised in AD [48]. This dual role, linking both neuroinflammatory and vascular pathologies, explains its identification as both a hub and a bottleneck.

CD53, a tetraspanin protein highly expressed on microglia, functions as a cell-surface organizer for a wide range of immune activities, including cell adhesion, migration, and the regulation of key receptors. Its importance in Alzheimer's disease (AD) is highlighted by single-cell RNA sequencing studies that consistently identify CD53 as a key upregulated gene in disease-associated microglia [52]. This role as a central coordinator of signaling platforms perfectly explains its identification as a hub, as it connects numerous pathways to influence diverse microglial functions, from environmental sensing to phagocytosis.

As one of the most well-established genetic risk factors for late-onset Alzheimer's disease, the enzyme Inositol Polyphosphate-5-Phosphatase D (INPP5D) functions as a key negative regulator of microglia. It acts as a crucial "brake" on microglial activation and phagocytosis, directly opposing signals from activating receptors like TREM2 [53]. Consequently, loss of INPP5D function can lead to a hyper-responsive inflammatory state. This explains its identification as a bottleneck, as it serves as a single, critical control point that dampens the entire microglial activation cascade.

RP11-556E13.1 is the genomic identifier for a long non-coding RNA and is not yet well-characterized. However, studies have shown that lncRNAs as a class are profoundly dysregulated in the AD brain and are involved in controlling fundamental processes like neuronal survival, inflammation, and APP processing [54].

VAE Algorithm on Genes Selected via Sparse Logistic Regression

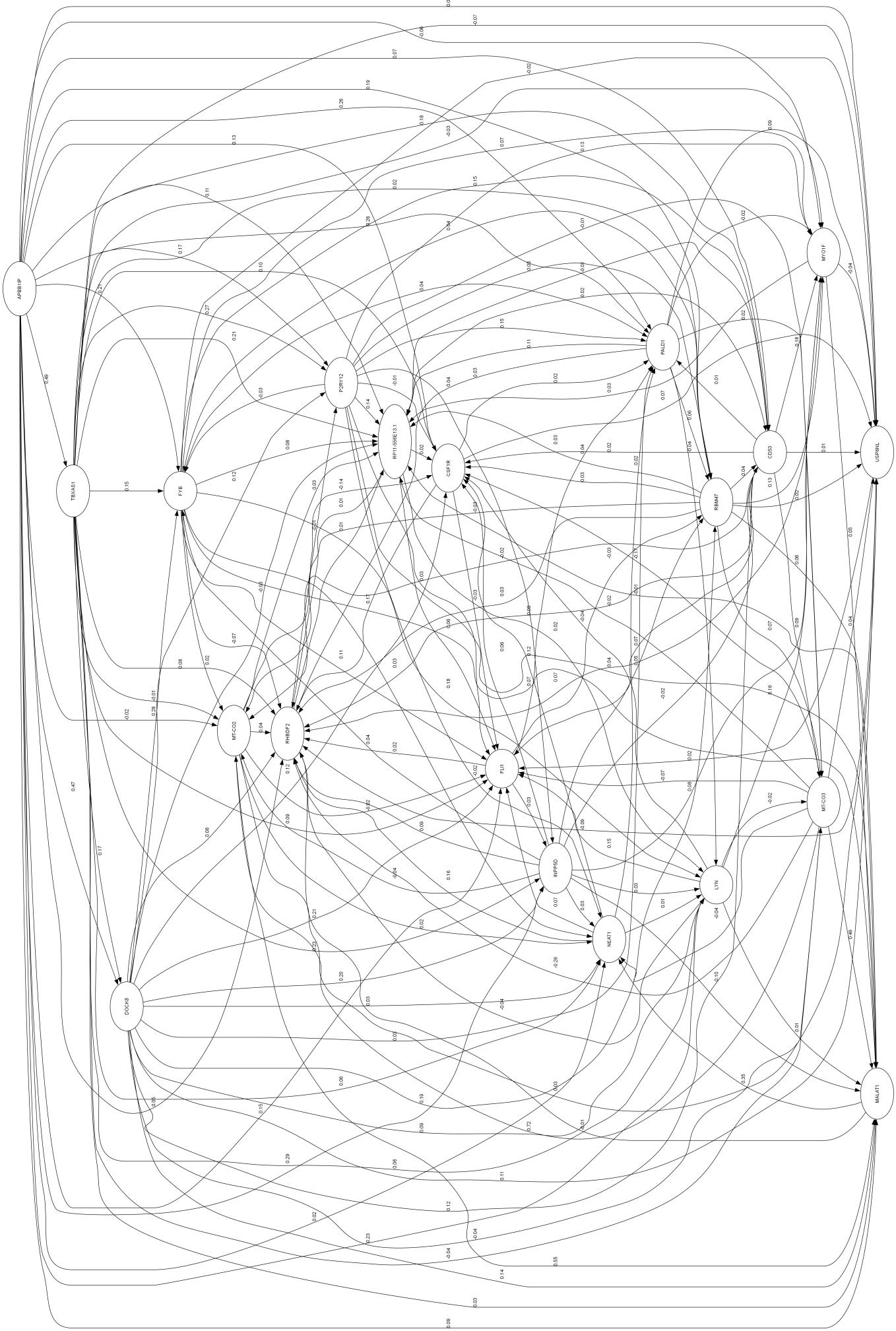
Figure 23 shows the causal diagram inferred by the VAE algorithm from the genes selected via Sparse Logistic Regression.

The identified hubs are **TBXAS1**, **APBB1IP**, **DOCK8**, **FYB** and **P2RY12**, while the bottlenecks are **FLI1**, **RHBDF2**, **RP11-556E13.1**, **P2RY12** and **PALD1**.

The purinergic receptor P2RY12 is the most definitive molecular signature of homeostatic microglia, the healthy, surveilling state essential for brain maintenance. A key pathological event in AD is the transformation of these cells into a "disease-associated microglia" state, a switch universally characterized by the profound loss of P2RY12 expression [55]. This loss is functionally critical, as P2RY12 is the primary sensor for nucleotides that guide microglia to sites of injury. It is a hub for the entire genetic program defining microglial health, and a bottleneck whose downregulation represents the critical gateway into the chronic neuroinflammatory state of AD.

The protein RHBDF2 (also known as iRhom2), a genetic risk factor for late-onset Alzheimer's disease (AD), functions as a master regulator of protein availability on the

Figure 23: Causal graph constructed using VAE on genes selected via Sparse Logistic Regression



surface of microglia. Very recent study by [56] shows that it achieves this by chaperoning the sheddase enzyme ADAM17, which acts as molecular scissors. This positions RHBDL2 at a critical nexus in AD pathology, as it simultaneously controls the shedding of two opposing molecules: it promotes the cleavage and inactivation of the protective receptor TREM2, while also enabling the release of the pro-inflammatory cytokine TNF-alpha. This explains that loss of RHBDL2 function is protective, as it boosts TREM2 levels and enhances amyloid-beta clearance.

Comparative Analysis of the Three Algorithms

The consensus causal graph is shown in Figure 24, while the SHD comparison is presented in Table 6. The high normalized SHD values suggest that the three methods identified substantially different causal structures.

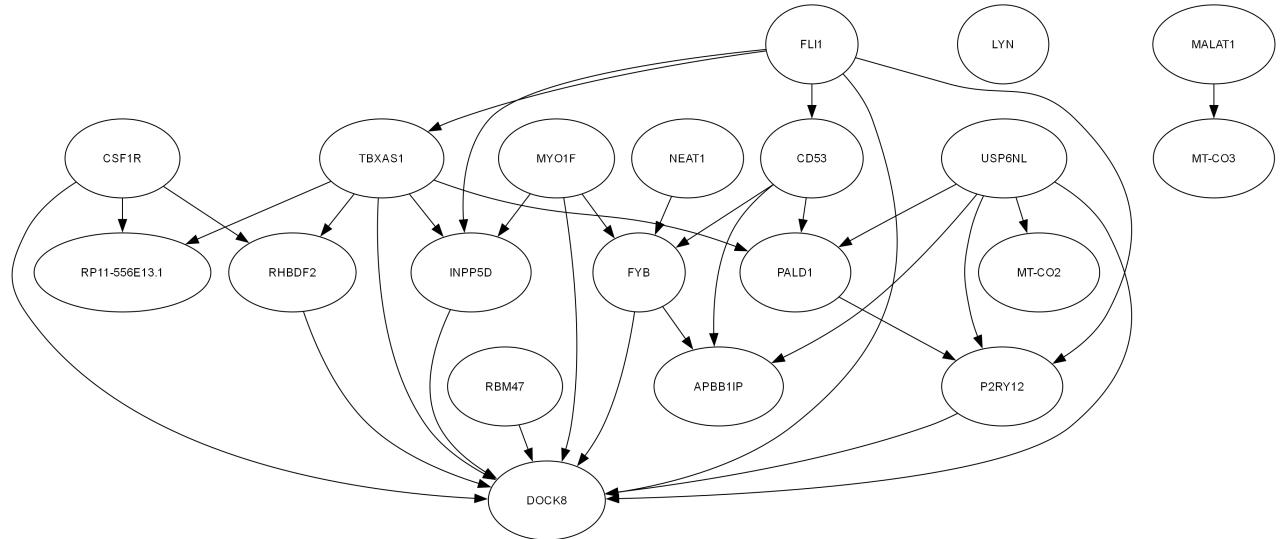


Figure 24: Consensus graph generated from three causal discovery methods applied to genes selected via Sparse Logistic Regression

Edges with particularly large weights in both the LiNGAM and VAE methods-and also identified by conditional independence testing-include MALAT1-MT-CO3 and DOCK8-FYB. Although a direct molecular interaction or co-regulatory pathway for these two gene pairs is not well-studied, these speculative relationships warrant further investigation.

Table 6: Structural Hamming Distance (SHD) of genes selected via logistic regression: Results include both absolute and normalized SHD.

Algorithm	Edges	SHD with PC	SHD with LiNGAM	SHD with VAE
PC	79	—	157(68.0%)	168(70.6%)
LiNGAM	152	—	—	167(53.7%)
VAE	159	—	—	—

Summary

Biologically, our causal graph recovered multiple genes already known to play key roles in AD—namely **LYN**, **CSF1R**, **INPP5D**, and **MT-CO2**. We also identified genes not yet extensively characterized in AD but recently proposed as potential disease markers, such as **PALD1**, **CD53** and **DOCK8**. Finally, several other genes have been highlighted for their therapeutic potential despite sparse functional follow-up, including **TBXAS1** and **RHBDL2**.

5.3.3 Causal Discovery on Genes Selected via Random Forest

PC Algorithm on Genes Selected via Random Forest

Due to the large number of genes involved in this causal discovery process, the PC algorithm became computationally intractable and failed to run. Given the time constraints and the need to maintain consistency in the PC algorithm implementation across different gene sets, we currently do not have results for causal discovery on the genes selected by the random forest method.

LiNGAM Algorithm on Genes Selected via Random Forest

In our network analysis, we identified ten key hubs: **USP6NL**, **LYN**, **RHBD2**, **TGFBR1**, **SYK**, **C1QC**, **CD74**, **PIK3AP1**, **PALD1**, and **INPP5D**; and ten bottlenecks: **MROH7-TTC4**, **SPATA3**, **RPUSD2**, **CELF1**, **RP11-428F8.2**, **MAN1A1**, **DUSP19**, **TMEM158**, **SAMD4A**, and **NR2C2**.

There are several genes that are less characterized in the context of AD, including **TMEM158**, **DUSP19**, **RPUSD2**, **SPATA3**, **MROH7-TTC4** and **RP11-428F8.2**.

The intracellular kinase **SYK** (Spleen Tyrosine Kinase) is critically important in AD because it serves as the direct downstream signaling partner for **TREM2**, a major AD risk gene. When **TREM2** is activated on microglia, it is **SYK** that transmits the signal onward, initiating a cascade that controls essential functions like phagocytosis, inflammation, and cell survival [57]. This position as the central intracellular transducer for signals from **TREM2** and other receptors explains its identification as a hub, as it connects upstream activation to a vast network of downstream cellular actions.

PIK3AP1 (Phosphoinositide-3-Kinase Adaptor Protein 1) functions as a crucial adaptor protein that activates the PI3K signaling pathway, which is essential for cell survival and metabolism. In microglia, the PI3K pathway is a critical downstream component of the **TREM2-SYK** signaling axis. By acting as a bridge, **PIK3AP1** relays the activation signal from **SYK** to the PI3K complex, thereby promoting essential microglial functions like survival and phagocytic activity [58].

The protein **C1QC** forms a crucial part of the classical complement cascade. A key pathological event in the early AD brain is the aberrant deposition of C1q onto vulnerable synapses. This deposition effectively "tags" these synapses for elimination by microglia in a destructive process of inappropriate synaptic pruning. This loss of synapses, driven by a misguided immune response, is now thought to be a primary cause of the cognitive

decline that precedes widespread neuron death [59]. Its identification as a hub is therefore a reflection of its biological role: C1QC acts as the physical bridge connecting the entire complement system to the microglial machinery that deconstructs the brain's connections.

CD74 is a protein that acts as the invariant chain for MHC class II molecules, which are essential for presenting antigens to other immune cells. CD74 expression is dramatically upregulated on microglia in the AD brain [55].

TGFBR1 is the primary receptor for the TGF-beta signaling pathway, which plays a crucial anti-inflammatory role in the healthy brain and is essential for maintaining homeostasis. In Alzheimer's disease, this protective pathway is known to be severely dysregulated. Consequently, alterations affecting TGFBR1 can cripple the brain's innate ability to suppress neuroinflammation, allowing damaging inflammatory processes to proceed unchecked and contribute to disease progression [60].

USP6NL is a key protein that regulates endocytosis and receptor trafficking, the fundamental process by which cells internalize molecules [61]. This function is highly relevant to AD because core pathological events-including the processing of Amyloid Precursor Protein (APP) and the clearance of amyloid-beta by microglia are heavily dependent on these endocytic pathways [62]. Its role as a central regulator of this crucial cellular machinery explains its identification as a hub.

CELF1 is an RNA-binding protein that controls the splicing and stability of messenger RNAs. It is directly implicated in AD through its specific regulation of the Tau protein (MAPT) mRNA. By binding to this transcript, CELF1 influences its splicing-a process that, when aberrant, can generate the toxic Tau isoforms that aggregate into neurofibrillary tangles [63]. This role as a key post-transcriptional regulator for a core AD pathology explains its identification as a bottleneck.

NR2C2 (Nuclear Receptor Subfamily 2 Group C Member 2) is a nuclear receptor. It is known to regulate genes involved in lipid metabolism, glucose homeostasis, and inflammation [64]. Dysregulation of lipid and glucose metabolism are well-established features of AD.

MAN1A1 (Mannosidase Alpha Class 1A Member 1) is involved in N-linked glycosylation, a process that modifies newly made proteins. Proper glycosylation is essential for the correct folding and function of countless cell surface receptors and secreted proteins, including APP and components of the inflammatory system [65].

The RNA-binding protein SAMD4A acts as a translational repressor, preventing specific messenger RNAs from being converted into their corresponding proteins [66]. This function is highly relevant to neuroinflammation, as its known targets include key proteins involved in inflammatory responses.

VAE Algorithm on Genes Selected via Random Forest

We identified 10 hub genes: **RP11-286B14.1**, **TOP1MT**, **ATXN1**, **AF131217.1**, **HSPH1**, **TXNRD1**, **VPS8**, **VAT1L**, **STAG1**, and **SPA17**; and ten bottlenecks: **CCAR1**, **ACADM**, **MT-ND3**, **FAM149A**, **ARHGAP26**, **AAK1**, **RPSAP58**, **VOPP1**, **DUS2**, and **RASGEF1C**.

There are several genes that are less characterized in the context of AD, including RP11-286B14.1, AF131217.1, STAG1, VAT1L, SPA17 and VPS8.

The TOP1MT gene encodes a mitochondrial topoisomerase essential for maintaining mitochondrial DNA integrity, a process known to fail early in A). As a hub gene, its high connectivity means that its malfunction can trigger widespread mitochondrial instability. A deficiency in TOP1MT directly impairs the expression of mitochondrial genes, crippling the cell's primary energy-generating pathway, the oxidative phosphorylation (OXPHOS) system [66]. Bioenergetic failure is a hallmark of AD.

The ATXN1 gene encodes Ataxin-1, a protein that regulates gene expression. While known for its role in spinocerebellar ataxia (SCA1), its link to AD risk involves a different mechanism centered on amyloid-beta production. Studies show that a loss of normal ATXN1 function initiates a specific pathological cascade. Reduced ATXN1 levels decrease the transcriptional repressor CIC, which then unleashes the ETV4/5 transcription factors. These factors directly increase the expression of Bace1, the gene encoding the critical beta-secretase enzyme. The resulting overabundance of BACE1 accelerates the amyloidogenic processing of amyloid precursor protein (APP), leading to increased production of toxic amyloid-beta peptides [67].

The HSPH1 gene encodes Hsp105/110, a molecular chaperone that directly counters the toxic protein aggregation central to AD. Its protective function in AD is primarily linked to the clearance of pathological tau, a hallmark of the disease. HSPH1 works in concert with other chaperones (like Hsp70) to form a complex that specifically recognizes and binds to toxic, modified tau. This machinery then targets the aberrant tau for degradation through the ubiquitin-proteasome system, effectively removing it from the neuron [68].

The TXNRD1 gene encodes Thioredoxin Reductase 1, a vital cytoplasmic selenoenzyme that maintains cellular redox balance. As a master regulator of the antioxidant defense system, it protects cells from oxidative stress, a key pathological feature in many diseases. Its connection to AD is twofold. First, as a primary defender against oxidative damage, any dysfunction in TXNRD1 can lower the threshold for neuronal injury. More directly, emerging research reveals that TXNRD1 also plays a critical role in promoting "inflammaging" - the chronic, low-grade inflammation that fuels AD progression. It achieves this by driving the pro-inflammatory Senescence-Associated Secretory Phenotype (SASP). This dual role positions TXNRD1 at the crucial intersection of oxidative stress and neuroinflammation, two core pathogenic pathways in Alzheimer's disease.

Comparative Analysis of the Three Algorithms

The SHD comparison is presented in Table 7. The high normalized SHD values suggest that the two methods identified substantially different causal structures.

There are 2,075 common edges shared between the two working algorithms. Among these, a few edges exhibit strong causal relationships, notably: LPCAT2–RP11-364P22.1 and APBB1IP–RP11-364P22.1.

RP11-364P22.1 is classified as a long intergenic non-protein coding RNA (lincRNA), meaning it is transcribed from a genomic region that does not overlap with protein cod-

ing genes. Its association with AD is currently not well understood. However, long non-coding RNAs (lncRNAs) have been increasingly recognized for their roles in AD, influencing disease progression through various mechanisms-such as interactions with chromatin-modifying complexes and regulation of mRNA splicing, stability, and translation [69]. Due to the limited understanding of RP11-364P22.1’s role in AD, these two causal relationships have not yet been reported in the literature.

Among the other two genes identified, LPCAT2 plays a central role in lipid metabolism and is known to modulate inflammatory responses by downregulating inflammatory cytokine production [70]. Given that neuroinflammation is a key contributor to AD pathogenesis, LPCAT2 may be significantly involved in the disease process. Furthermore, elevated levels of lysophosphatidylcholine-lipid species regulated directly or indirectly by LPCAT enzymes-have been associated with neurodegenerative processes.

The involvement of APBB1IP in AD is also still emerging. Previous studies have shown that changes in APBB1IP expression occur in animal models exposed to chronic stress, which is a condition relevant to AD risk [71].

Table 7: Structural Hamming Distance (SHD) of genes selected via random forest: Results include both absolute and normalized SHD.

Algorithm	Edges	SHD with PC	SHD with LiNGAM	SHD with VAE
PC	–	–	–	–
LiNGAM	16271	–	–	17511(80.8%)
VAE	5390	–	–	–

Summary

Biologically, the hub and bottleneck genes identified by causal discovery algorithms applied to the random forest-selected genes are less well characterized in Alzheimer’s disease pathology and appear to play relatively peripheral roles.. Our causal graph recovered several genes already implicated in AD, albeit as secondary actors-namely **SYK**, **C1QC**, **CD74**, **TGFBR1**, **CELF1**, **MT-CO2**, **ATXN1**, and **HSPH1**. We also identified a set of genes that, while not extensively studied in AD, have recently been proposed as potential disease markers: **PIK3AP1**, **USP6NL**, **NR2C2**, **MAN1A1**, **SAMD4A**, **TOP1MT**, and **TXNRD1**.

5.4 Performance of Causal Discovery Methods

Although we lack a definitive ground truth for our causal structures and therefore cannot directly assess accuracy, we can still compare the three methods in terms of computational efficiency, stability of their estimates, and interpretability of their outputs.

For very high-dimensional datasets, the classic PC algorithm often proves unsuitable because its computational cost and statistical requirements grow prohibitively large. In the worst case, where the true graph is dense, the number of conditional independence tests increases exponentially with the number of variables, quickly becoming intractable.

Moreover, each test requires inverting the covariance matrix of the conditioning set; whenever the size of that set exceeds the sample size, the covariance submatrix is singular and no partial correlation can be computed [72]. Finally, PC’s theoretical guarantees depend on a strong-faithfulness assumption: every nonzero partial correlation stays bounded away from zero regardless of the conditioning set. In high dimensions, distributions violating this assumption proliferate and the probability of satisfying strong-faithfulness decays exponentially with p , undermining both the algorithm’s consistency and its practical reliability.

In contrast, LiNGAM tends to yield more stable estimates than VAE-based causal discovery methods under high dimensionality. LiNGAM exploits non-Gaussianity in the data to identify a unique linear ordering of variables and estimate causal strengths via simple regression, which remains robust even when the number of variables is large relative to the sample size [21]. By contrast, VAE-based approaches must learn nonlinear latent representations through stochastic optimization, and their edge-strength estimates shrink towards zero as the model’s capacity is strained by many input dimensions. In our experiments on gene expression data selected by random forests, the strongest edge weight from the VAE model is 0.16, and the total number of statistically significant edges was markedly lower than with LiNGAM—indicating higher variance in edge-strength estimate.

Additionally, across all three methods, high dimensionality exacerbates challenges in interpretability and consistency of the inferred graphs. As the number of variables grows, the proportion of robustly identified edges declines sharply [73]. The biological plausibility of the top causal edges also diminishes: many inferred relations become difficult to reconcile with known pathways, and different methods (PC, LiNGAM, VAE) often produce divergent network structures. This divergence not only hampers scientific interpretation but also makes it hard to decide which edges, if any, to trust when formulating downstream hypotheses or intervention strategies.

Finally, it is important to note that all the causal discovery methods we employed assume a strictly DAG, which precludes any feedback loops. In contrast, actual gene regulatory networks often involve cycles—transcription factors activate downstream genes that, in turn, regulate the original transcription factors. By enforcing acyclicity, these algorithms can distort or entirely miss bidirectional and autoregulatory interactions. One way to overcome this limitation is to adopt a longitudinal design, measuring each gene’s expression at multiple time points. The temporal ordering offers a built-in mechanism for distinguishing cause from effect: if gene A’s transcription repeatedly precedes an increase in gene B’s levels over successive time points, we can infer a directed edge from A to B with greater confidence. If gene B later influences gene A in subsequent intervals, that feedback loop becomes visible.

6 Future Works

Our computationally derived causal graphs provide a strong foundation for generating testable hypotheses regarding the complex molecular interactions in Alzheimer’s Disease (AD), particularly within microglia. Building upon our causal discovery results, three

primary directions for future research seem promising.

Firstly, future efforts could prioritize the experimental validation of these inferred causal relationships and the functional roles of identified key hub and bottleneck genes, such as MALAT1. Gene perturbation studies, notably employing the Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) system for precise gene editing in human iPSC-derived microglia (iMGLs), will be useful in our setting [74]. These experiments will allow us to meticulously assess the impact of altering these key genes on downstream gene expression and AD-relevant cellular phenotypes, thereby confirming their positions and influence within the computationally derived network.

Beyond validating the intra-microglial network, it is crucial to expand our focus to the broader brain ecosystem. AD pathogenesis involves intricate, often dysregulated, inter-cellular communication between microglia and other cell types like neurons, astrocytes, and oligodendrocytes - all active contributors to the disease, involved in neuroinflammation, amyloid- β and tau processing, and myelin degeneration [75]. We propose leveraging specialized computational algorithms, such as CellChat [76] or NicheNet [77], to predict and analyze signaling pathways and ligand-receptor interactions. This will help elucidate how the identified microglial hub genes might modulate communication with other neural cells, thereby influencing the overall brain environment in AD.

Finally, to achieve a more comprehensive and mechanistically robust understanding, future research should integrate diverse data types. Our current causal discovery, primarily based on transcriptomics, can be significantly enhanced by incorporating other omics layers, including proteomics, metabolomics, and epigenomics. This multi-omics approach will facilitate the construction of more holistic, multi-layered causal models of AD pathogenesis. Furthermore, incorporating important covariates such as age, gender, race, APOE genotype, and lifestyle factors (e.g., smoking habits) into these integrative causal network analyses will allow for the development of more robust, context-specific models that can better account for the heterogeneity of AD and guide more personalized therapeutic strategies.

7 Conclusion

In conclusion, we addressed the three key research questions posed at the outset. Our causal mediation analysis using ADNI biomarker data found no significant mediated effect of amyloid- β on the relationship between neuroinflammation and AD progression. Similarly, ANOVA showed no significant differences in sTREM2 levels across diagnostic stages, challenging the widely assumed linear cascade linking neuroinflammation, amyloid- β , and AD.

Given the limitations of this simplified model, we developed a comprehensive analytical pipeline combining robust feature selection - via Random Forest, Sparse Logistic Regression, and Gradient Boosting - with multiple causal discovery algorithms. Each feature selection method identified distinct sets of genes predictive of AD, with gene ontology analysis confirming enrichment in AD-relevant functions. Six genes (ARHGAP24, C3, CLEC17A, DHRSX, MALAT1, and MEF2C) were consistently selected across all meth-

ods. All models modestly outperformed a baseline classifier, with Gradient Boosting yielding the best performance.

We then applied four causal discovery algorithms (PC, GES, LiNGAM, and VAE) to each gene set. While PC was computationally intensive and GES struggled with local optima in high-dimensional settings, the VAE method captured non-linear effects but showed attenuation at scale. LiNGAM demonstrated promise, and using all methods in parallel allowed us to identify robust and consistent causal edges.

This work addresses key challenges in analyzing high-dimensional, observational scRNA-seq data, enhances the stability of feature selection, and explores the biological implications of inferred causal structures. In the absence of a known ground truth, these findings offer a valuable foundation for future experimental validation. By uncovering novel genetic targets and regulatory mechanisms, our results ultimately support the development of more precise diagnostic tools and targeted therapies for Alzheimer’s disease.

Appendix

Baseline Model Performance Curve and Evaluation Metrics

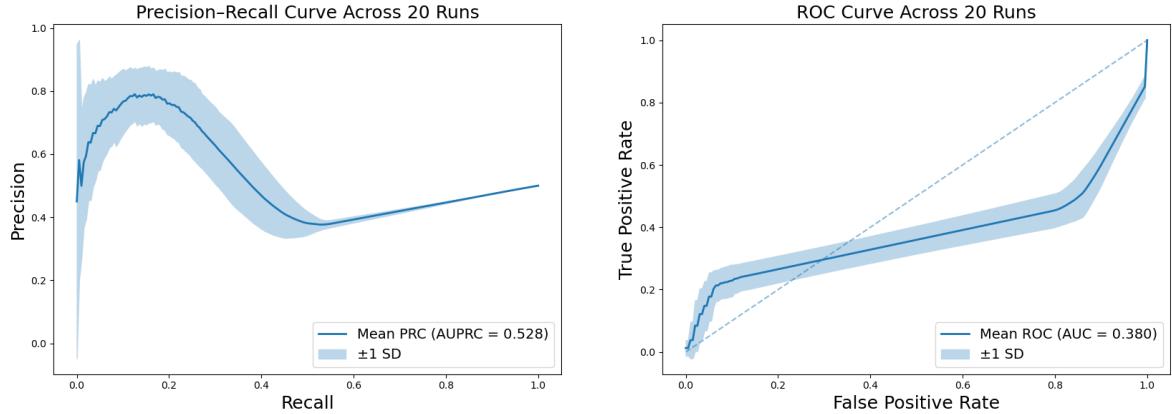


Figure 25: Precision-recall curve (left) and ROC curve (right) of baseline

We denote true positives, false positives, false negatives, and true negatives by TP, FP, FN, and TN, respectively.

Precision is defined as

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Recall is defined as

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

F1-Score is defined as

$$\text{F1} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

AUROC is defined as

$$\text{AUROC} = \int_0^1 \text{TP}\left(\text{FP}^{-1}(x)\right) dx,$$

AUPRC is defined as

$$\text{AUPRC} = \int_0^1 \text{Precision}\left(\text{Recall}^{-1}(x)\right) dx,$$

Gene Onotology Analysis Result

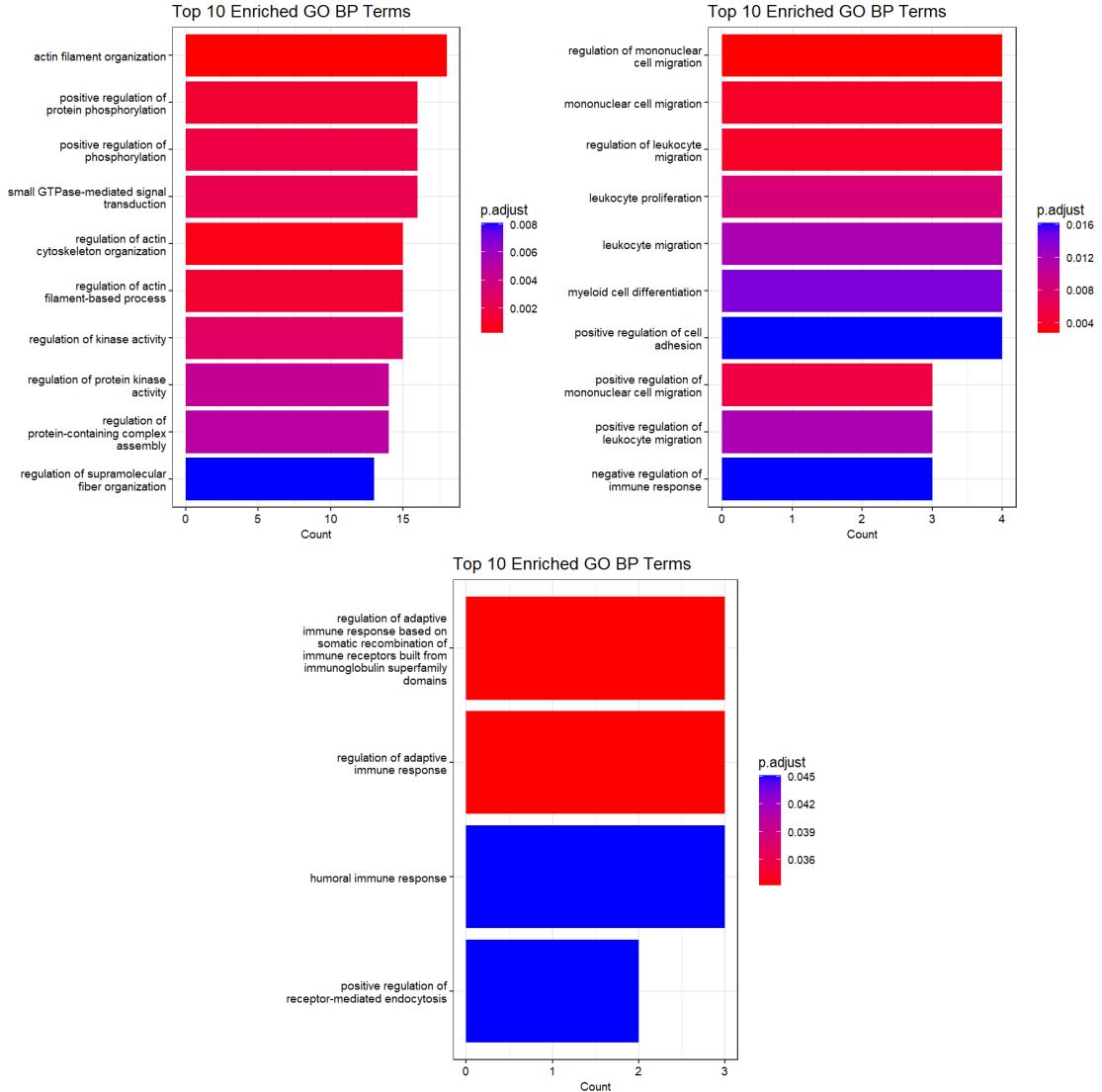


Figure 26: GO analysis results of genes selected via random forest(upper left), sparse logistic regression(upper right) and gradient boosting(lower middle)

Proof of Equation (4)

The matrix exponential is defined by its Taylor series:

$$\exp(B) = I + B + \frac{1}{2!}B^2 + \frac{1}{3!}B^3 + \dots$$

Since B is nonnegative, every power B^k has nonnegative entries, and so does $\exp(B)$. The trace of $\exp(B)$ is given by:

$$\text{tr}(\exp(B)) = \text{tr}(I) + \text{tr}(B) + \frac{1}{2!}\text{tr}(B^2) + \frac{1}{3!}\text{tr}(B^3) + \dots$$

The term $\text{tr}(I) = m$ comes from the zeroth power I . For $k \geq 1$, $\text{tr}(B^k)$ is exactly the sum of all length- k cycles (counted with weight). If the graph has no cycles at all, then for every $k \geq 1$, there are no length- k closed walks, so

$$\text{tr}(B^k) = 0, \quad k \geq 1,$$

and consequently

$$\text{tr}(\exp(B)) = m$$

Data and Code Availability

The full implementation for our approach used in this work can be found at this GitHub repository: <https://github.com/JennyYuanZW/masterdissertation.git>.

Acknowledgements

I would like to express my sincere gratitude to Prof. Ramji Venkataramanan and Prof. Valeriya Malysheva for their valuable guidance and support throughout the development of this report. Their insights and feedback were instrumental in shaping the direction and quality of this work. I would also like to express my sincere gratitude to Margo Bellanger and Lunki Sucipto for their support and valuable advice during the experimental phase of this project.

Risk Assessment Retrospective

There is no risk associated with this project.

References

- [1] A. Serrano-Pozo, M. P. Frosch, E. Masliah, and B. T. Hyman, “Neuropathological alterations in alzheimer disease,” *Cold Spring Harbor perspectives in medicine*, vol. 1, no. 1, p. a006189, 2011.
- [2] J. L. Cummings, T. Morstorf, and K. Zhong, “Alzheimer’s disease drug-development pipeline: few candidates, frequent failures,” *Alzheimer’s research & therapy*, vol. 6, pp. 1–7, 2014.
- [3] Q. Cai and P. Tammineni, “Mitochondrial aspects of synaptic dysfunction in alzheimer’s disease,” *Journal of Alzheimer’s disease*, vol. 57, no. 4, pp. 1087–1103, 2017.
- [4] W. J. Streit, R. E. Mrak, and W. S. T. Griffin, “Microglia and neuroinflammation: a pathological perspective,” *Journal of neuroinflammation*, vol. 1, pp. 1–4, 2004.

- [5] C. M. Henstridge, B. T. Hyman, and T. L. Spires-Jones, “Beyond the neuron–cellular interactions early in alzheimer disease pathogenesis,” *Nature Reviews Neuroscience*, vol. 20, no. 2, pp. 94–108, 2019.
- [6] K. Imai, L. Keele, and T. Yamamoto, “Identification, inference and sensitivity analysis for causal mediation effects,” 2010.
- [7] S. Kaji, S. A. Berghoff, L. Spieth, L. Schlaphoff, A. O. Sasmita, S. Vitale, L. Büschgens, S. Kedia, M. Zirngibl, T. Nazarenko *et al.*, “Apolipoprotein e aggregation in microglia initiates alzheimer’s disease pathology by seeding β -amyloidosis,” *Immunity*, vol. 57, no. 11, pp. 2651–2668, 2024.
- [8] M. Suárez-Calvet, G. Kleinberger, M. Á. Araque Caballero, M. Brendel, A. Rominger, D. Alcolea, J. Fortea, A. Lleó, R. Blesa, J. D. Gispert *et al.*, “strem 2 cerebrospinal fluid levels are a potential biomarker for microglia activity in early-stage alzheimer’s disease and associate with neuronal injury markers,” *EMBO molecular medicine*, vol. 8, no. 5, pp. 466–476, 2016.
- [9] D. Tingley, T. Yamamoto, K. Hirose, L. Keele, and K. Imai, “Mediation: R package for causal mediation analysis,” *Journal of statistical software*, vol. 59, pp. 1–38, 2014.
- [10] M. Suárez-Calvet, E. Morenas-Rodríguez, G. Kleinberger, K. Schlepckow, M. Á. Araque Caballero, N. Franzmeier, A. Capell, K. Fellerer, B. Nuscher, E. Eren *et al.*, “Early increase of csf strem2 in alzheimer’s disease is associated with tau related-neurodegeneration but not with amyloid- β pathology,” *Molecular neurodegeneration*, vol. 14, pp. 1–14, 2019.
- [11] D. Biel, M. Suárez-Calvet, P. Hager, A. Rubinski, A. Dewenter, A. Steward, S. Roemer, M. Ewers, C. Haass, M. Brendel *et al.*, “strem2 is associated with amyloid-related p-tau increases and glucose hypermetabolism in alzheimer’s disease,” *EMBO Molecular Medicine*, vol. 15, no. 2, p. e16987, 2023.
- [12] R. Fujikawa and M. Tsuda, “The functions and phenotypes of microglia in alzheimer’s disease,” *Cells*, vol. 12, no. 8, p. 1207, 2023.
- [13] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, “Scikit-learn: Machine learning in python,” *the Journal of machine Learning research*, vol. 12, pp. 2825–2830, 2011.
- [14] C. Wang, D. Acosta, M. McNutt, J. Bian, A. Ma, H. Fu, and Q. Ma, “A single-cell and spatial rna-seq database for alzheimer’s disease (ssread),” *Nature Communications*, vol. 15, no. 1, p. 4710, 2024.
- [15] T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*. Springer, 2009, vol. 2.
- [16] S. Cervantes, L. Samaranch, J. M. Vidal-Taboada, I. Lamet, M. J. Bullido, A. Frank-García, F. Coria, A. Lleó, J. Clarimón, E. Lorenzo *et al.*, “Genetic variation in apoe cluster region and alzheimer’s disease risk,” *Neurobiology of aging*, vol. 32, no. 11, pp. 2107–e7, 2011.
- [17] Y. Zheng, B. Huang, W. Chen, J. Ramsey, M. Gong, R. Cai, S. Shimizu, P. Spirtes, and K. Zhang, “Causal-learn: Causal discovery in python,” *Journal of Machine Learning Research*, vol. 25, no. 60, pp. 1–8, 2024.

- [18] P. Spirtes, C. N. Glymour, and R. Scheines, *Causation, prediction, and search*. MIT press, 2000.
- [19] D. M. Chickering, “Optimal structure identification with greedy search,” *Journal of machine learning research*, vol. 3, no. Nov, pp. 507–554, 2002.
- [20] A. Chattopadhyay and T.-P. Lu, “Gene-gene interaction: the curse of dimensionality,” *Annals of translational medicine*, vol. 7, no. 24, p. 813, 2019.
- [21] S. Shimizu, P. O. Hoyer, A. Hyvärinen, A. Kerminen, and M. Jordan, “A linear non-gaussian acyclic model for causal discovery.” *Journal of Machine Learning Research*, vol. 7, no. 10, 2006.
- [22] Y. Yu, J. Chen, T. Gao, and M. Yu, “Dag-gnn: Dag structure learning with graph neural networks,” in *International conference on machine learning*. PMLR, 2019, pp. 7154–7163.
- [23] J. A. Bondy and U. S. R. Murty, *Graph theory*. Springer Publishing Company, Incorporated, 2008.
- [24] X. Zheng, B. Aragam, P. K. Ravikumar, and E. P. Xing, “Dags with no tears: Continuous optimization for structure learning,” *Advances in neural information processing systems*, vol. 31, 2018.
- [25] C. Nithya, M. Kiran, and H. A. Nagarajaram, “Hubs and bottlenecks in protein-protein interaction networks,” *Reverse engineering of regulatory networks*, pp. 227–248, 2023.
- [26] S. Acid and L. M. de Campos, “Searching for bayesian network structures in the space of restricted acyclic partially directed graphs,” *Journal of artificial intelligence research*, vol. 18, pp. 445–490, 2003.
- [27] L. Li, Y. Xu, M. Zhao, and Z. Gao, “Neuro-protective roles of long non-coding rna malat1 in alzheimer’s disease with the involvement of the microrna-30b/cnr1 network and the following pi3k/akt activation,” *Experimental and molecular pathology*, vol. 117, p. 104545, 2020.
- [28] Y. Ke, P. Chen, C. Wu, Q. Wang, K. Zeng, and M. Liang, “ β 2-microglobulin and cognitive decline: unraveling the mediating role of the dunedin pace of aging methylation,” *Frontiers in Aging Neuroscience*, vol. 17, p. 1505185, 2025.
- [29] L. K. Smith, Y. He, J.-S. Park, G. Bieri, C. E. Snethlage, K. Lin, G. Gontier, R. Wabl, K. E. Plambeck, J. Udeochu *et al.*, “ β 2-microglobulin is a systemic pro-aging factor that impairs cognitive function and neurogenesis,” *Nature medicine*, vol. 21, no. 8, pp. 932–937, 2015.
- [30] G. Zhang, Y. Luo, G. Li, L. Wang, D. Na, X. Wu, Y. Zhang, X. Mo, and L. Wang, “Dhrsx, a novel non-classical secretory protein associated with starvation induced autophagy,” *International journal of medical sciences*, vol. 11, no. 9, p. 962, 2014.
- [31] S. Bolognin, E. Lorenzetto, G. Diana, and M. Buffelli, “The potential role of rho gtpases in alzheimer’s disease pathogenesis,” *Molecular neurobiology*, vol. 50, no. 2, pp. 406–422, 2014.

- [32] D. Blacker, M. A. Wilcox, N. M. Laird, L. Rodes, S. M. Horvath, R. C. Go, R. Perry, B. Watson, S. S. Bassett, M. G. McInnis *et al.*, “Alpha-2 macroglobulin is genetically associated with alzheimer disease,” *Nature genetics*, vol. 19, no. 4, pp. 357–360, 1998.
- [33] J. Ren, S. Zhang, X. Wang, Y. Deng, Y. Zhao, Y. Xiao, J. Liu, L. Chu, and X. Qi, “Mef2c ameliorates learning, memory, and molecular pathological changes in alzheimer’s disease in vivo and in vitro: Neuroprotective effects of mef2c,” *Acta Biochimica et Biophysica Sinica*, vol. 54, no. 1, p. 77, 2021.
- [34] Z. Hao, X. Guo, J. Wu, and G. Yang, “Revisiting the benefits of exercise for alzheimer’s disease through the lens of ferroptosis: A new perspective,” *Aging and disease*, 2025.
- [35] S. J. Barker, R. M. Raju, N. E. Milman, J. Wang, J. Davila-Velderrain, F. Gunter-Rahman, C. C. Parro, P. L. Bozzelli, F. Abdurrob, K. Abdelaal *et al.*, “Mef2 is a key regulator of cognitive potential and confers resilience to neurodegeneration,” *Science translational medicine*, vol. 13, no. 618, p. eabd7695, 2021.
- [36] N. Zhao, Y. Ren, Y. Yamazaki, W. Qiao, F. Li, L. M. Felton, S. Mahmoudiandehkordi, A. Kueider-Paisley, B. Sonoustoun, M. Arnold *et al.*, “Alzheimer’s risk factors age, apoe genotype, and sex drive distinct molecular pathways,” *Neuron*, vol. 106, no. 5, pp. 727–742, 2020.
- [37] M. Bermejo-Jambrina, J. Eder, L. C. Helgers, N. Hertoghs, B. M. Nijmeijer, M. Stunnenberg, and T. B. Geijtenbeek, “C-type lectin receptors in antiviral immunity and viral escape,” *Frontiers in immunology*, vol. 9, p. 590, 2018.
- [38] M. Lisek, O. Przybyszewski, L. Zylinska, F. Guo, and T. Boczek, “The role of mef2 transcription factor family in neuronal survival and degeneration,” *International journal of molecular sciences*, vol. 24, no. 4, p. 3120, 2023.
- [39] J. Wu, H. Zhang, Y. Zheng, X. Jin, M. Liu, S. Li, Q. Zhao, X. Liu, Y. Wang, M. Shi *et al.*, “The long noncoding rna malat1 induces tolerogenic dendritic cells and regulatory t cells via mir155/dendritic cell-specific intercellular adhesion molecule-3 grabbing nonintegrin/il10 axis,” *Frontiers in immunology*, vol. 9, p. 1847, 2018.
- [40] N. A. G. Canales, V. M. Marina, J. S. Castro, A. A. Jiménez, G. Mendoza-Hernández, E. L. McCarron, M. B. Roman, and J. I. Castro-Romero, “A1bg and c3 are overexpressed in patients with cervical intraepithelial neoplasia iii,” *Oncology Letters*, vol. 8, no. 2, pp. 939–947, 2014.
- [41] D. Santiard-Baron, D. Langui, M. Delehedde, B. Delatour, B. Schombert, N. Touchet, G. Tremp, M.-F. Paul, V. Blanchard, N. Sergeant *et al.*, “Expression of human fe65 in amyloid precursor protein transgenic mice is associated with a reduction in β -amyloid load,” *Journal of neurochemistry*, vol. 93, no. 2, pp. 330–338, 2005.
- [42] H. B. Nygaard, C. H. van Dyck, and S. M. Strittmatter, “Fyn kinase inhibition as a novel therapy for alzheimer’s disease,” *Alzheimer’s research & therapy*, vol. 6, pp. 1–8, 2014.
- [43] R. Islam, H. H. Choudhary, F. Zhang, H. Mehta, J. Yoshida, A. J. Thomas, and K. Hanafy, “Microglial tlr4-lyn kinase is a critical regulator of neuroinflammation, $\alpha\beta$ phagocytosis, neuronal damage, and cell survival in alzheimer’s disease,” *Scientific Reports*, vol. 15, no. 1, p. 11368, 2025.

- [44] D. Rucker and A. S. Dhamoon, “Physiology, thromboxane a2,” 2019.
- [45] M. H. Thomas and J. L. Olivier, “Arachidonic acid in alzheimer’s disease,” *Journal of Neurology & Neuromedicine*, vol. 1, no. 9, 2016.
- [46] C. Huang, R. Zhou, X. Huang, F. Dai, and B. Zhang, “Integrative analysis of single-nucleus rna sequencing and mendelian randomization to explore novel risk genes for alzheimer’s disease,” *Medicine*, vol. 103, no. 46, p. e40551, 2024.
- [47] A. Olmos-Alonso, S. T. Schetters, S. Sri, K. Askew, R. Mancuso, M. Vargas-Caballero, C. Holscher, V. H. Perry, and D. Gomez-Nicola, “Pharmacological targeting of csf1r inhibits microglial proliferation and prevents the progression of alzheimer’s-like pathology,” *Brain*, vol. 139, no. 3, pp. 891–907, 2016.
- [48] J. M. Perez Ortiz and R. H. Swerdlow, “Mitochondrial dysfunction in alzheimer’s disease: Role in pathogenesis and novel therapeutic opportunities,” *British journal of pharmacology*, vol. 176, no. 18, pp. 3489–3507, 2019.
- [49] Y. Sommerer, V. Dobricic, M. Schilling, O. Ohlei, S. S. Sabet, T. Wesse, J. Fuß, S. Franzenburg, A. Franke, L. Parkkinen *et al.*, “Entorhinal cortex epigenome-wide association study highlights four novel loci showing differential methylation in alzheimer’s disease,” *Alzheimer’s Research & Therapy*, vol. 15, no. 1, p. 92, 2023.
- [50] W. Zhang, F. Teng, X. Lan, P. Liu, A. Wang, F. Zhang, Z. Cui, J. Guan, and X. Sun, “A novel finding relates to the involvement of atf3/dock8 in alzheimer’s disease pathogenesis,” *Journal of Alzheimer’s Disease*, p. 13872877251336266, 2024.
- [51] A. E. Dart, S. Tollis, and R. G. Endres, “Investigating forces for uptake and cup closure—the role of myosins in phagocytosis.”
- [52] C. S. Frigerio, L. Wolfs, N. Fattorelli, N. Thrupp, I. Voytyuk, I. Schmidt, R. Mancuso, W.-T. Chen, M. E. Woodbury, G. Srivastava *et al.*, “The major risk factors for alzheimer’s disease: age, sex, and genes modulate the microglia response to a β plaques,” *Cell reports*, vol. 27, no. 4, pp. 1293–1306, 2019.
- [53] J. D. Samuels, K. A. Moore, H. E. Ennerfelt, A. M. Johnson, A. E. Walsh, R. J. Price, and J. R. Lukens, “The alzheimer’s disease risk factor inpp5d restricts neuroprotective microglial responses in amyloid beta-mediated pathology,” *Alzheimer’s & Dementia*, vol. 19, no. 11, pp. 4908–4921, 2023.
- [54] Z. Lan, Y. Chen, J. Jin, Y. Xu, and X. Zhu, “Long non-coding rna: insight into mechanisms of alzheimer’s disease,” *Frontiers in molecular neuroscience*, vol. 14, p. 821002, 2022.
- [55] H. Keren-Shaul, A. Spinrad, A. Weiner, O. Matcovitch-Natan, R. Dvir-Szternfeld, T. K. Ulland, E. David, K. Baruch, D. Lara-Astaiso, B. Toth *et al.*, “A unique microglia type associated with restricting development of alzheimer’s disease,” *Cell*, vol. 169, no. 7, pp. 1276–1290, 2017.
- [56] G. Jocher, G. Ozcelik, S. A. Müller, H.-E. Hsia, M. L. Osua, L. I. Hofmann, M. Aßfalg, L. Dinkel, X. Feng, K. Schlepckow *et al.*, “The late-onset alzheimer’s disease risk factor rhbdf2 is a modifier of microglial trem2 proteolysis,” *Life science alliance*, vol. 8, no. 5, 2025.

- [57] L. Zhang, X. Xiang, Y. Li, G. Bu, and X.-F. Chen, “Trem2 and strem2 in alzheimer’s disease: from mechanisms to therapies,” *Molecular Neurodegeneration*, vol. 20, no. 1, p. 43, 2025.
- [58] F. L. Yeh, D. V. Hansen, and M. Sheng, “Trem2, microglia, and neurodegenerative diseases,” *Trends in molecular medicine*, vol. 23, no. 6, pp. 512–533, 2017.
- [59] S. Hong, V. F. Beja-Glasser, B. M. Nfonoyim, A. Frouin, S. Li, S. Ramakrishnan, K. M. Merry, Q. Shi, A. Rosenthal, B. A. Barres *et al.*, “Complement and microglia mediate early synapse loss in alzheimer mouse models,” *Science*, vol. 352, no. 6286, pp. 712–716, 2016.
- [60] B. Spittau, N. Dokalis, and M. Prinz, “The role of tgf β signaling in microglia maturation and activation,” *Trends in immunology*, vol. 41, no. 9, pp. 836–848, 2020.
- [61] J. Zhang, Z. Jiang, and A. Shi, “Rab gtpases: The principal players in crafting the regulatory landscape of endosomal trafficking,” *Computational and structural biotechnology journal*, vol. 20, pp. 4464–4472, 2022.
- [62] Y. H. Qureshi, P. Baez, and C. Reitz, “Endosomal trafficking in alzheimer’s disease, parkinson’s disease, and neuronal ceroid lipofuscinosis,” *Molecular and Cellular Biology*, vol. 40, no. 19, pp. e00262–20, 2020.
- [63] J.-M. Gallo and C. Spickett, “The role of celf proteins in neurological disorders,” *RNA biology*, vol. 7, no. 4, pp. 474–479, 2010.
- [64] R. Hiwa, J. F. Brooks, J. L. Mueller, H. V. Nielsen, and J. Zikherman, “Nr4a nuclear receptors in t and b lymphocytes: Gatekeepers of immune tolerance,” *Immunological Reviews*, vol. 307, no. 1, pp. 116–133, 2022.
- [65] S. Schedin-Weiss, B. Winblad, and L. O. Tjernberg, “The role of protein glycosylation in alzheimer disease,” *The FEBS journal*, vol. 281, no. 1, pp. 46–62, 2014.
- [66] X.-Y. Wang and L.-N. Zhang, “Rna binding protein samd4: current knowledge and future perspectives,” *Cell & Bioscience*, vol. 13, no. 1, p. 21, 2023.
- [67] J. Suh, D. M. Romano, L. Nitschke, S. P. Herrick, B. A. DiMarzio, V. Dzhala, J.-S. Bae, M. K. Oram, Y. Zheng, B. Hooli *et al.*, “Loss of ataxin-1 potentiates alzheimer’s pathogenesis by elevating cerebral bacel transcription,” *Cell*, vol. 178, no. 5, pp. 1159–1175, 2019.
- [68] Y. Dong, T. Li, Z. Ma, C. Zhou, X. Wang, and J. Li, “Hspa1a, hspa2, and hspa8 are potential molecular biomarkers for prognosis among hsp70 family in alzheimer’s disease,” *Disease Markers*, vol. 2022, no. 1, p. 9480398, 2022.
- [69] M. Garofalo, C. Pandini, D. Sproviero, O. Pansarasa, C. Cereda, and S. Gagliardi, “Advances with long non-coding rnas in alzheimer’s disease as peripheral biomarker,” *Genes*, vol. 12, no. 8, p. 1124, 2021.
- [70] S. K. Jackson, W. Abate, J. Parton, S. Jones, and J. L. Harwood, “Lysophospholipid metabolism facilitates toll-like receptor 4 membrane translocation to regulate the inflammatory response,” *Journal of Leucocyte Biology*, vol. 84, no. 1, pp. 86–92, 2008.

- [71] P. Jungke, G. Ostrow, J.-L. Li, S. Norton, K. Nieber, O. Kelber, and V. Butterweck, “Profiling of hypothalamic and hippocampal gene expression in chronically stressed rats treated with st. johnâs wort extract (stw 3-vi) and fluoxetine,” *Psychopharmacology*, vol. 213, pp. 757–772, 2011.
- [72] S. Mohseni-Sehdeh and W. Saad, “Induced covariance for causal discovery in linear sparse structures,” *arXiv preprint arXiv:2410.01221*, 2024.
- [73] M. Kalisch and P. Bühlman, “Estimating high-dimensional directed acyclic graphs with the pc-algorithm.” *Journal of Machine Learning Research*, vol. 8, no. 3, 2007.
- [74] S. Bhardwaj, K. K. Kesari, M. Rachamalla, S. Mani, G. M. Ashraf, S. K. Jha, P. Kumar, R. K. Ambasta, H. Dureja, H. P. Devkota *et al.*, “Crispr/cas9 gene editing: New hope for alzheimer’s disease therapeutics,” *Journal of Advanced Research*, vol. 40, pp. 207–221, 2022.
- [75] R. Ziar, P. J. Tesar, and B. L. Clayton, “Astrocyte and oligodendrocyte pathology in alzheimer’s disease,” *Neurotherapeutics*, p. e00540, 2025.
- [76] S. Jin, C. F. Guerrero-Juarez, L. Zhang, I. Chang, R. Ramos, C.-H. Kuan, P. Myung, M. V. Plikus, and Q. Nie, “Inference and analysis of cell-cell communication using cellchat,” *Nature communications*, vol. 12, no. 1, p. 1088, 2021.
- [77] R. Browaeys, W. Saelens, and Y. Saeys, “Nichenet: modeling intercellular communication by linking ligands to target genes,” *Nature methods*, vol. 17, no. 2, pp. 159–162, 2020.