

TOWARD MONOSEMANTIC CLINICAL EXPLANATIONS FOR ALZHEIMER’S DIAGNOSIS VIA ATTRIBUTION AND MECHANISTIC INTERPRETABILITY

Anonymous authors

Paper under double-blind review

A TECHNICAL APPENDICES

A.1 ATTRIBUTIONAL THEORY AND METHODS

Attribution explainability methods follow the framework of additive feature attribution, where the explanation model $g(f, \mathbf{x})$ is represented as a linear function of simplified input features:

$$g(f, \mathbf{x}) = \phi_0 + \sum_{i=1}^M \phi_i x_i \quad (1)$$

Here, f is the predictive model, $\phi_i \in \mathbb{R}$ is the attribution (importance) assigned to feature x_i , and M is the number of simplified input features.

For this study, we employed six well-established attributional interpretability methods applied to large language models (LLMs), denoted as $K = 6$: *Feature Ablation*, *Layer Activations* (which capture the embedding activation space of a specific layer of interest within the LLM), *Layer DeepLIFT SHAP*, *Layer Gradient SHAP* (Lundberg & Lee, 2017), *Layer Integrated Gradients* (Sundararajan et al., 2017), and *Layer Gradient \times Activation*.

To align these layer-wise interpretability methods with the additive feature attribution framework, we reinterpret the internal activations (i.e., latent units) of a network layer L as simplified input features. The objective is to estimate an attribution score ϕ_i for each unit, where $\phi_i \in \mathbb{R}$ quantifies the contribution of the corresponding neuron to the model’s prediction.

Layer SHAP implementations: This directly corresponds to the Shapley formulation:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)] \quad (2)$$

In practice, Deep SHAP approximates this using sampling and a chain-rule based linearization over network layers Lundberg & Lee (2017). *Gradient SHAP* assumes that input features are independent and that the explanation model is linear, allowing explanations to be expressed as an additive composition of feature contributions. Under these assumptions, SHAP values (Lundberg & Lee, 2017) can be approximated by computing the expected gradients over a distribution of perturbed inputs. Specifically, Gaussian noise is added to each input feature to generate multiple baseline samples, and the resulting gradients are averaged to approximate SHAP attributions.

Activation Attribution: This method treats the raw activation $a_i^L(\mathbf{x})$ as proportional to its importance in the output. In the additive form:

$$\phi_i = a_i^L(\mathbf{x}) \quad (3)$$

Assuming linearity between layer L and the output, activations themselves serve as proxy contributions.

Gradient \times Activation Attribution: This method computes the element-wise product between the activation values and the gradients of the model output with respect to those activations, thereby

capturing the first-order sensitivity of the output to the neurons in the layer. To this end, the method estimates the first-order sensitivity of the output with respect to the activation:

$$\phi_i = a_i^L(\mathbf{x}) \cdot \frac{\partial f}{\partial a_i^L}(\mathbf{x}) \quad (4)$$

This corresponds to a local linear approximation (first-order Taylor expansion) of the model at \mathbf{x} , akin to DeepLIFT and the SHAP linearization used in DeepLift SHAP (Lundberg & Lee, 2017).

Feature Ablation Attribution: This attributional interpretability technique is a perturbation-based approach to estimating attributions. It involves replacing the input or output values of a selected layer with a given baseline or reference value and computing the resulting change in the model’s output. By default, each neuron (i.e., scalar input or output value) within the layer is ablated independently. For neuron group $S \subseteq \{1, \dots, d_L\}$, the perturbed activation is:

$$\tilde{a}_i^L = \begin{cases} b_i^L & \text{if } i \in S, \\ a_i^L(\mathbf{x}) & \text{otherwise,} \end{cases} \quad (5)$$

and the attribution is the marginal effect:

$$\phi_S = f(\mathbf{x}; \tilde{\mathbf{a}}_S^L) - f(\mathbf{x}) \quad (6)$$

All attribution methods were applied to the final (22nd) layer of the MODERN-BERT LLM—the model variant that achieved the highest classification accuracy in our evaluations (see Supplementary material section 1.1). These formulations allow us to ground various neural attribution techniques within a unified additive explanation model, facilitating their comparison and hybridization under shared theoretical assumptions.

A.2 ATTRIBUTIONAL EXPLANATION OPTIMIZER FRAMEWORK

Let $\mathcal{A} = \{A_1, A_2, \dots, A_K\}$ denote the set of $K = 6$ attribution methods applied to the final layer L of the model f . Each method A_k generates an attribution vector $\boldsymbol{\phi}^{(k)} = [\phi_1^{(k)}, \phi_2^{(k)}, \dots, \phi_M^{(k)}]$, where M is the number of latent features (neurons) in layer L . The goal is to derive a unified attribution vector $\boldsymbol{\phi}$ that captures the consensus explanation across methods.

A.2.1 SCORING AND WEIGHTING ATTRIBUTION METHODS

Each attribution vector $\boldsymbol{\phi}^{(k)}$ is evaluated using the following quality metrics:

A.2.2 EVALUATION INTERPRETABILITY METRICS

We evaluate the robustness of each attribution method A_k using the following stability metrics:

Relative Input Stability (RIS):

$$M_{\text{RIS}}^{(k)} = \text{RIS}(f, \boldsymbol{\phi}^{(k)}; \mathbf{x}) = \frac{\|\mathbf{x}\|_p}{\|\boldsymbol{\phi}^{(k)}(\mathbf{x})\|_p} \max_{\mathbf{x}' \in \mathcal{N}_{\mathbf{x}}, \hat{y}_{\mathbf{x}'} = \hat{y}_{\mathbf{x}}} \frac{\|\boldsymbol{\phi}^{(k)}(\mathbf{x}) - \boldsymbol{\phi}^{(k)}(\mathbf{x}')\|_p}{\|\mathbf{x} - \mathbf{x}'\|_p} \quad (7)$$

Relative Output Stability (ROS):

$$M_{\text{ROS}}^{(k)} = \text{ROS}(f, \boldsymbol{\phi}^{(k)}; \mathbf{x}) = \frac{\|f(\mathbf{x})\|_p}{\|\boldsymbol{\phi}^{(k)}(\mathbf{x})\|_p} \max_{\mathbf{x}' \in \mathcal{N}_{\mathbf{x}}, \hat{y}_{\mathbf{x}'} = \hat{y}_{\mathbf{x}}} \frac{\|\boldsymbol{\phi}^{(k)}(\mathbf{x}) - \boldsymbol{\phi}^{(k)}(\mathbf{x}')\|_p}{\|f(\mathbf{x}) - f(\mathbf{x}')\|_p} \quad (8)$$

Here, $\mathcal{N}_{\mathbf{x}}$ denotes a neighborhood of perturbed inputs \mathbf{x}' around \mathbf{x} , and $\hat{y}_{\mathbf{x}}$ is the predicted class label. Both metrics measure the relative sensitivity of the attribution vector $\boldsymbol{\phi}^{(k)}$ to perturbations in the input or output space.

Sparseness Metric: We quantify the **sparseness** of the attribution vector $\boldsymbol{\phi}^{(k)} \in \mathbb{R}^d$ using the *Gini Index*, a measure of inequality that has been shown to satisfy several desirable properties for evaluating sparseness Chalasani et al. (2020). This formulation is adopted in the context of explaining neural network predictions Chalasani et al. (2020).

Let $v \in \mathbb{R}_{\geq 0}^d$ be a non-negative vector. Denote by $v_{(k)}$ the k -th smallest element in v after sorting it in non-decreasing order. Then, the **Gini Index** $G(v) \in [0, 1]$ is defined as:

$$G(v) = 1 - 2 \sum_{k=1}^d \frac{v_{(k)}}{\|v\|_1} \cdot \left(\frac{d - k + 0.5}{d} \right), \quad (9)$$

where $\|v\|_1 = \sum_{i=1}^d v_i$ is the ℓ_1 -norm of v . To evaluate the sparseness of an attribution vector $\boldsymbol{\phi}^{(k)}$, we apply the Gini Index to the vector of its absolute values:

$$\text{Sparseness}(\boldsymbol{\phi}^{(k)}) = G(|\boldsymbol{\phi}^{(k)}|),$$

where $|\boldsymbol{\phi}^{(k)}| = (|\phi_1^{(k)}|, |\phi_2^{(k)}|, \dots, |\phi_d^{(k)}|)$.

Higher values of $G(|\boldsymbol{\phi}^{(k)}|)$ indicate greater sparseness. In the extreme case, if only one component is non-zero, the Gini Index reaches its maximum value of 1, indicating perfect sparseness. If all components are equal, the Gini Index is 0.

A.2.3 AGGREGATION OF ATTRIBUTIONS

The weighted average attribution vector $\bar{\boldsymbol{\phi}}$ is calculated as:

$$\bar{\boldsymbol{\phi}} = \sum_{k=1}^K w_k \cdot \boldsymbol{\phi}^{(k)} \quad (10)$$

This vector serves as the target explanation for the optimization process.

A.2.4 EXPLANATION RECONSTRUCTION VIA ENCODER-DECODER MODELS

An encoder-decoder model is trained to generate a reconstructed explanation $\hat{\boldsymbol{\phi}}$ from the original input \mathbf{x} . Two architectures are considered the Diffusion UNet1D Ronneberger et al. (2015) and the x-transformer autoencoder Vaswani et al. (2017); Nguyen & Salazar (2019).

Diffusion model: The diffusion model follows the basic structure of a 1-dimensional U-Net and is trained using diffusion principles. In this framework, diffusion models Ho et al. (2020) are latent variable models in which the observed data $\phi_0^{(k)}$ is gradually corrupted through a forward noising process, producing a sequence of latent variables $\phi_{1:T}^{(k)}$. A corresponding reverse process is then learned to recover the original data from noise. The mathematical formulation is as follows:

FORWARD PROCESS: A fixed Markov chain progressively adds Gaussian noise to the data:

$$q(\phi_{1:T}^{(k)} | \phi_0^{(k)}) := \prod_{t=1}^T q(\phi_t^{(k)} | \phi_{t-1}^{(k)}), \quad q(\phi_t^{(k)} | \phi_{t-1}^{(k)}) := \mathcal{N}(\phi_t^{(k)}; \sqrt{1 - \beta_t} \phi_{t-1}^{(k)}, \beta_t \mathbf{I}) \quad (11)$$

Alternatively, sampling from the forward process at an arbitrary timestep t is possible in closed form:

$$q(\phi_t^{(k)} | \phi_0^{(k)}) = \mathcal{N}(\phi_t^{(k)}; \sqrt{\bar{\alpha}_t} \phi_0^{(k)}, (1 - \bar{\alpha}_t) \mathbf{I}), \quad (12)$$

where $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$.

REVERSE PROCESS: A learned time-reversal model with Gaussian transitions:

$$p_{\theta}(\phi_{0:T}^{(k)}) := p(\phi_T^{(k)}) \prod_{t=1}^T p_{\theta}(\phi_{t-1}^{(k)} | \phi_t^{(k)}), \quad p_{\theta}(\phi_{t-1}^{(k)} | \phi_t^{(k)}) := \mathcal{N}(\phi_{t-1}^{(k)}; \boldsymbol{\mu}_{\theta}(\phi_t^{(k)}, t), \Sigma_{\theta}(\phi_t^{(k)}, t)), \quad (13)$$

where $p(\phi_T^{(k)}) := \mathcal{N}(\phi_T^{(k)}; \mathbf{0}, \mathbf{I})$.

TRAINING OBJECTIVE: The training objective of diffusion models is based on a variational bound, which includes Kullback–Leibler (KL) divergence terms. The KL term comparing the true posterior from the forward process and the model’s learned reverse process is written as:

$$\text{KL}\left(q(\phi_{t-1}^{(k)} | \phi_t^{(k)}, \phi_0^{(k)}) \| p_\theta(\phi_{t-1}^{(k)} | \phi_t^{(k)})\right) \quad (14)$$

Both distributions are Gaussian:

$$q(\phi_{t-1}^{(k)} | \phi_t^{(k)}, \phi_0^{(k)}) = \mathcal{N}(\phi_{t-1}^{(k)}; \tilde{\mu}_t(\phi_t^{(k)}, \phi_0^{(k)}), \tilde{\beta}_t \mathbf{I}) \quad (15)$$

$$p_\theta(\phi_{t-1}^{(k)} | \phi_t^{(k)}) = \mathcal{N}(\phi_{t-1}^{(k)}; \mu_\theta(\phi_t^{(k)}, t), \sigma_t^2 \mathbf{I}) \quad (16)$$

The closed-form KL divergence between two Gaussians $\mathcal{N}(\mu_1, \sigma_1^2 \mathbf{I})$ and $\mathcal{N}(\mu_2, \sigma_2^2 \mathbf{I})$ in d -dimensions is:

$$\text{KL} = \frac{1}{2} \left[\log \left(\frac{\sigma_2^2}{\sigma_1^2} \right) + \frac{\sigma_1^2 + \|\mu_1 - \mu_2\|^2}{\sigma_2^2} - d \right] \quad (17)$$

In our setting, this term is computed for each timestep t and summed across all steps:

$$\mathcal{L}_{1:T-1} = \sum_{t=2}^T \mathbb{E}_{q(\phi_0^{(k)}, \phi_t^{(k)})} \left[\text{KL}\left(q(\phi_{t-1}^{(k)} | \phi_t^{(k)}, \phi_0^{(k)}) \| p_\theta(\phi_{t-1}^{(k)} | \phi_t^{(k)})\right) \right] \quad (18)$$

This forms a core part of the evidence lower bound (ELBO) optimized during training. Using variational inference, we minimize the negative ELBO:

$$\mathcal{L} = \mathbb{E}_q \left[-\log p(\phi_T^{(k)}) + \sum_{t=1}^T \text{KL}\left(q(\phi_{t-1}^{(k)} | \phi_t^{(k)}, \phi_0^{(k)}) \| p_\theta(\phi_{t-1}^{(k)} | \phi_t^{(k)})\right) - \log p_\theta(\phi_0^{(k)} | \phi_1^{(k)}) \right]. \quad (19)$$

Each KL term compares Gaussian distributions and can be computed in closed form. The posterior $q(\phi_{t-1}^{(k)} | \phi_t^{(k)}, \phi_0^{(k)})$ is also Gaussian:

$$q(\phi_{t-1}^{(k)} | \phi_t^{(k)}, \phi_0^{(k)}) = \mathcal{N}(\phi_{t-1}^{(k)}; \tilde{\mu}_t(\phi_t^{(k)}, \phi_0^{(k)}), \tilde{\beta}_t \mathbf{I}), \quad (20)$$

with:

$$\tilde{\mu}_t(\phi_t^{(k)}, \phi_0^{(k)}) = \frac{\sqrt{\tilde{\alpha}_t} \beta_t}{1 - \tilde{\alpha}_t} \phi_0^{(k)} + \frac{\sqrt{\tilde{\alpha}_t} (1 - \tilde{\alpha}_{t-1})}{1 - \tilde{\alpha}_t} \phi_t^{(k)}, \quad (21)$$

$$\tilde{\beta}_t = \frac{1 - \tilde{\alpha}_{t-1}}{1 - \tilde{\alpha}_t} \beta_t. \quad (22)$$

SIMPLIFIED TRAINING LOSS: The common parameterization rewrites the objective as denoising score matching:

$$\mathcal{L}_{\text{simple}}(\theta) := \mathbb{E}_{t, \phi_0^{(k)}, \epsilon} \left[\left\| \epsilon - \epsilon_\theta(\sqrt{\tilde{\alpha}_t} \phi_0^{(k)} + \sqrt{1 - \tilde{\alpha}_t} \epsilon, t) \right\|^2 \right], \quad (23)$$

where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and ϵ_θ is the neural network trained to predict noise.

In our implementation we compute the total loss for the diffusion model as:

$$\mathcal{L}_{\text{similarity}}(\hat{\phi}, \bar{\phi}) = \mathcal{L}_{\text{similarity}}(\theta) = \frac{1}{K+1} \sum_{l=0}^K \mathcal{L}_{\text{simple}}^{(l)}(\theta) \quad (24)$$

x-Transformer: Let the input sequence be:

$$\phi^{(k)} = [\phi_1^{(k)}, \phi_2^{(k)}, \dots, \phi_T^{(k)}] \in \mathbb{R}^{T \times d_{\text{in}}}$$

where $d_{\text{in}} = 7$ is the input dimensionality and $T = 512$ is the sequence length. We consider a Transformer-based encoder-decoder architecture operating on input sequences $\Phi^{(k)} \in \mathbb{R}^{B \times T \times d_{\text{in}}}$

at diffusion step k , where: B is the batch size, T is the sequence length, d_{in} is the input feature dimension, and $\Phi^{(k)}$ is the input sequence at step k .

The processing pipeline is mathematically formulated as follows:

INPUT PROJECTION AND POSITIONAL ENCODING: We first project the input to the model dimension d and add positional encodings:

$$\mathbf{X}_0 = \mathbf{W}_{\text{in}} \Phi^{(k)} + \mathbf{P}, \quad \mathbf{X}_0 \in \mathbb{R}^{B \times T \times d} \quad (25)$$

where: $\mathbf{W}_{\text{in}} \in \mathbb{R}^{d_{\text{in}} \times d}$ is a learnable linear projection matrix, and $\mathbf{P} \in \mathbb{R}^{1 \times T \times d}$ is a learnable positional embedding matrix.

ENCODER: MULTI-HEAD SELF-ATTENTION LAYERS: The encoder consists of L_e stacked multi-head self-attention (MHSA) layers:

$$\mathbf{H}_{\text{enc}} = \text{MHSA}_{L_e} \circ \dots \circ \text{MHSA}_1(\mathbf{X}_0) \quad (26)$$

where each MHSA layer performs:

$$\text{MHSA}(\mathbf{X}) = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}} \right) \mathbf{V} \quad (27)$$

with: $\mathbf{Q}, \mathbf{K}, \mathbf{V}$: Query, Key, and Value matrices obtained via learned linear projections, and d_h : the dimensionality of each attention head.

DECODER INPUT PROJECTION: During training, the decoder may receive the ground-truth output $\Phi_{\text{target}}^{(k)} \in \mathbb{R}^{B \times T \times 1}$:

$$\mathbf{Y}_0 = \mathbf{W}_{\text{dec}} \Phi_{\text{target}}^{(k)} + \mathbf{P} \quad (28)$$

where $\mathbf{W}_{\text{dec}} \in \mathbb{R}^{1 \times d}$ is a projection matrix.

If no decoder input is available (e.g., during inference), $\Phi_{\text{target}}^{(k)}$ is initialized to a zero tensor.

DECODER MHSA + CROSS-ATTENTION LAYERS: The decoder consists of L_d layers of MHSA followed by cross-attention (CA) using the encoder context:

$$\mathbf{H}_{\text{dec}} = \text{CA}_{L_d} \circ \dots \circ \text{CA}_1(\text{MHSA}_{L_d} \circ \dots \circ \text{MHSA}_1(\mathbf{Y}_0) \mid \mathbf{H}_{\text{enc}}) \quad (29)$$

Each cross-attention (CA) layer uses the decoder hidden state as the query and encoder output as the key and value:

$$\text{CA}(\mathbf{Y}, \mathbf{H}_{\text{enc}}) = \text{Softmax} \left(\frac{\mathbf{Q}_{\text{dec}} \mathbf{K}_{\text{enc}}^\top}{\sqrt{d_h}} \right) \mathbf{V}_{\text{enc}} \quad (30)$$

OUTPUT PROJECTION: Finally, the decoder output is projected back to the target dimension:

$$\hat{\Phi}^{(k)} = \mathbf{W}_{\text{out}} \mathbf{H}_{\text{dec}}, \quad \hat{\Phi}^{(k)} \in \mathbb{R}^{B \times T \times 1} \quad (31)$$

where $\mathbf{W}_{\text{out}} \in \mathbb{R}^{d \times 1}$ is a linear projection matrix.

The similarity cost function is given by the Mean Squared Error (MSE) loss between the predicted output of the x-Transformer and the target weighted attribution vector as follow:

$$\mathcal{L}_{\text{similarity}}(\hat{\Phi}, \bar{\Phi}) = \mathcal{L}_{\text{MSE}} = \frac{1}{T} \sum_{t=1}^T \|\hat{\Phi}_t - \bar{\Phi}_t\|^2, \quad (32)$$

A.3 THE TOTAL COST FUNCTION OF THE OPTIMIZER

As previously highlighted, the reconstruction of the optimal explanation and the associated cost function adhere to the same principles and architectural design outlined in Mamalakis et al. (2025). The cost function consists of three key components: sparseness, as defined in ?; ROS and RIS scores Agarwal et al. (2022); and similarity. The integration of these components ensures a robust and interpretable evaluation. The total cost function for training the reconstruction model is:

$$\begin{aligned} \mathcal{L}_{\text{total}}(\Phi^{(k)}, \hat{\Phi}) = & \lambda_1 \cdot \frac{1}{M_{\text{RIS}}(f, \hat{\Phi})} + \lambda_2 \cdot \frac{1}{M_{\text{ROS}}(f, \hat{\Phi})} \\ & + \lambda_3 \cdot M_{\text{sparse}}(f, \hat{\Phi}) + \lambda_4 \cdot \mathcal{L}_{\text{similarity}}(\hat{\Phi}, \bar{\Phi}) \end{aligned} \quad (33)$$

where: $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ are hyperparameters controlling the influence of each loss term. This formulation enables a principled and quantitative integration of multiple attribution methods, optimizing toward a robust and interpretable explanation.

A.4 THE UMAP EXTRACTION AND THE LINEAR CONSTRAIN

Given a dataset $\hat{\Phi} = \{\hat{\phi}_1, \hat{\phi}_2, \dots, \hat{\phi}_n\} \subset \mathbb{R}^D$, UMAP aims to find a low-dimensional embedding $U = \{u_1, u_2, \dots, u_n\} \subset \mathbb{R}^d$ where typically $d = 2$ or $d = 3$, such that the local topological structure of the data in $\hat{\Phi}$ is preserved in U .

HIGH-DIMENSIONAL GRAPH CONSTRUCTION: First, the algorithm constructs a k-nearest neighbors graph in the high-dimensional space $\hat{\Phi}$. The distance metric used to calculate the pairwise distances is typically Euclidean:

$$d(\hat{\phi}_i, \hat{\phi}_j) = \|\hat{\phi}_i - \hat{\phi}_j\|_2$$

Next, a conditional probability is defined between points $\hat{\phi}_i$ and $\hat{\phi}_j$ using a Gaussian distribution:

$$p_{ij} = \exp\left(-\frac{\|\hat{\phi}_i - \hat{\phi}_j\|^2}{\sigma_i^2}\right)$$

where σ_i is the bandwidth for the Gaussian distribution, determined through a binary search to match a fixed perplexity.

The graph is symmetrized:

$$P_{ij} = \frac{p_{ij} + p_{ji}}{2}$$

LOW-DIMENSIONAL EMBEDDING GRAPH: In the low-dimensional space, a similar probability is defined between points u_i and u_j :

$$q_{ij} = \frac{1}{1 + a\|u_i - u_j\|^{2b}}$$

where a and b are hyperparameters that control the shape of the distribution, and $\|u_i - u_j\|_2$ is the Euclidean distance between points in the low-dimensional embedding.

OBJECTIVE FUNCTION: The optimization process involves minimizing the cross-entropy between the high-dimensional and low-dimensional probability distributions:

$$\mathcal{L} = \sum_{i < j} [P_{ij} \log(Q_{ij}) + (1 - P_{ij}) \log(1 - Q_{ij})]$$

This loss function encourages points that are close in the high-dimensional space to be close in the low-dimensional space, and points that are distant to remain distant.

OPTIMIZATION PROCESS: The optimization is carried out using stochastic gradient descent (SGD), updating the embedding points $\{u_i\}$ iteratively based on the gradient of the loss function \mathcal{L} . The gradient updates for the low-dimensional embedding u_i are computed as follows:

$$\frac{\partial \mathcal{L}}{\partial u_i} = - \sum_{j \neq i} (P_{ij} - Q_{ij}) \frac{u_i - u_j}{\|u_i - u_j\|_2^2}$$

REGULARIZATION CONSTRAINT: To prevent the embedding from collapsing to a single point, we introduce a variance constraint to ensure that the variance of the embedding does not approach zero:

$$\text{Var}(U) = \frac{1}{n} \sum_{i=1}^n \|u_i - \bar{u}\|_2^2 \geq \epsilon$$

where $\bar{U} = \frac{1}{n} \sum_{i=1}^n u_i$ is the mean of the embeddings, and $\epsilon > 0$ is a small constant that enforces a lower bound on the variance.

APPLICATION OF UMAP IN OUR PROBLEM: To obtain a comparable low-dimensional representation of the attribution scores across all tokenizer features, we applied a feature-wise UMAP projection procedure to the normalized attribution matrix. For each attribution method, the attribution tensor has shape $\mathbb{R}^{M \times T}$, where M denotes the number of test samples in the evaluation cohort and T corresponds to the dimensionality of the tokenizer embedding space of the input text. For each feature $j \in \{1, \dots, T\}$, we first applied min-max normalization to the feature-specific attribution vector

$$\mathbf{x}^{(j)} \in \mathbb{R}^M,$$

and subsequently performed a one-dimensional UMAP projection to obtain a two-dimensional embedding

$$\mathbf{y}^{(j)} \in \mathbb{R}^{M \times 2}.$$

The resulting coordinates were then normalized to the interval $[0, 1]$ to ensure that all feature-wise embeddings share a common bounded range. This procedure preserves the relative neighborhood structure of the M -sample attribution distribution for each feature while mapping all T features into a comparable two-dimensional representation space.

The motivation for applying UMAP independently to each of the T tokenizer features is to ensure that all attribution methods are projected into an aligned and comparable representation space. Since each attribution method produces values defined over the same token embedding dimensions, a feature-wise nonlinear projection enables consistent cross-method comparison of attribution patterns within the shared tokenizer feature space.

LINEAR CONSTRAINT FOR EQUAL COMPONENTS IN UMAP: Let $u_i = (u_{i1}, u_{i2}, \dots, u_{id})$ denote the embedding of the i -th data point in a d -dimensional space. The requirement that the first and second embedding components are equal can be written as:

$$u_{i1} = u_{i2} \quad \forall i \in \{1, 2, \dots, n\}.$$

Equivalently, this can be expressed as the linear equality constraint:

$$u_{i1} - u_{i2} = 0 \quad \forall i \in \{1, 2, \dots, n\}.$$

This constraint enforces that, for each data point i , the first and second coordinates of the embedding vector u_i are identical.

Within the total objective $\mathcal{L}_{\text{total}}(\boldsymbol{\phi}^{(k)}, \hat{\boldsymbol{\phi}})$ of Eq. 35, an additional penalty term may be introduced to enforce this constraint. The penalty can be written as:

$$\lambda_5 \sum_{i=1}^n (u_{i1} - u_{i2})^2,$$

where λ_5 is a regularization parameter controlling the strength of the constraint. This term encourages the first and second components of each reconstructed embedding point from the optimizer ($\hat{\boldsymbol{\phi}}$) to be equal, while still allowing flexibility depending on the value of λ_5 .

A.5 THE SUPERPOSITION AND THE MONOSEMANTIC REPRESENTATIONS

We model an embedding space as a real vector space \mathbb{R}^d , where a hidden activation vector $\mathbf{h} \in \mathbb{R}^d$ represents a combination of underlying semantic features. By the linear representation hypothesis, each interpretable feature corresponds to a fixed direction in \mathbb{R}^d (Olah et al. (2020b); Elhage et al. (2022b)).

Let $\mathbf{a} \in \mathbb{R}^F$ be a sparse feature activation vector and $W \in \mathbb{R}^{d \times F}$ be a linear transformation such that:

$$\mathbf{h} = W\mathbf{a} = \sum_{i=1}^F a_i \mathbf{w}_i,$$

where \mathbf{w}_i denotes the i -th column of W , corresponding to the direction of the feature i .

If $F > d$, the map W cannot be invertible, and thus different combination of characteristics can map to the same embedding. This gives rise to superposition, where multiple semantic features are embedded into shared subspaces or overlapping neuron activations Elhage et al. (2022b).

MONOSEMANTIC REPRESENTATIONS: A representation is called monosemantic when each neuron corresponds to a single interpretable feature Olah et al. (2020b). Mathematically, this corresponds to the case where W is full-rank and aligned with the identity matrix (or a rotation of it):

$$W = I \Rightarrow \mathbf{h} = \mathbf{a}.$$

This implies that each feature a_i is represented by a unique dimension h_i , with no overlap. Each neuron responds to a single, isolated concept, akin to “grandmother cells” in neuroscience Quiroga et al. (2005).

POLYSEMANTIC REPRESENTATIONS: In contrast, polysemantic neurons represent multiple, distinct concepts. Formally, if neuron h_j computes:

$$h_j = \sum_{i=1}^F W_{j,i} a_i,$$

and two or more $W_{j,i} \neq 0$, then neuron j encodes multiple features simultaneously, exhibiting polysemanticity Elhage et al. (2022b); Bills et al. (2023a).

More generally, a polysemantic embedding may be viewed as a mixture:

$$\mathbf{h} = \sum_{k=1}^K \alpha_k \mathbf{c}_k, \quad K > 1,$$

where \mathbf{c}_k are concept vectors and α_k are scalar weights.

This behavior is prevalent in both neural network activations and in biological neurons that exhibit mixed selectivity Rigotti et al. (2013).

Monosemantic representations arise from disentangled bases, where neurons correspond to isolated features. Superposition emerges from dimensionality compression and necessarily leads to polysemantic neurons, each encoding a combination of features. Sparse auto-encoder is a way to try to solve the polysemantic neurons—each encoding problem.

A.6 THE SAE APPROACH AND ARCHITECTURES

Sparse Autoencoder (SAE) architectures have advanced our understanding of how language and vision models represent features Gorton (2024). Neural network behavior is often explained via *computational circuits*—collections of neurons that together compute meaningful functions. Classical circuit analysis has identified key components such as edge detectors Olah et al. (2020a) or word-copying units Olsson et al. (2022). By using features derived from SAEs rather than raw neurons, researchers have improved the interpretability of circuits related to complex behaviors Marks et al. (2024).

Feature discovery can involve visual analysis McDougall (2024), manual inspection Bricken et al. (2023), and even assistance from large language models Bills et al. (2023b). Their causal role is often validated via activation interventions: modifying a feature activation vector \mathbf{a} and observing predictable changes in model output Templeton et al. (2024).

The mathematical formulation situates SAE architectures within the theoretical framework of superposition and semantic disentanglement. By expressing hidden states as sparse linear combinations of interpretable features, SAEs bridge the gap between low-level activations and human-understandable concepts.

LINEAR FORMULATION OF SAEs: Let $\mathbf{x} \in \mathbb{R}^d$ denote a layer’s neuron activation vector in a pretrained model. A Sparse Autoencoder learns a sparse feature representation $\mathbf{a} \in \mathbb{R}^F$ such that:

$$\hat{\mathbf{x}} = W\mathbf{a} + \mathbf{b}, \tag{34}$$

where $W \in \mathbb{R}^{d \times F}$ is the decoder (dictionary) matrix and $\mathbf{b} \in \mathbb{R}^d$ is a learned bias term. Each column $W_{\cdot,i}$ represents the direction of feature i in neuron space, and a_i is its activation. This linear mapping enables complex activations to be expressed as combinations of more interpretable features.

If $F > d$, then the feature space is overcomplete, and W cannot be full-rank. This leads to superposition, where multiple features overlap in the same subspace, and individual neurons encode multiple unrelated concepts Elhage et al. (2022b). If W is invertible and aligned to a basis, each neuron corresponds to a single feature. The representation is monosemantic and disentangled Olah et al. (2020b). When W has overlapping columns, neurons can respond to multiple features, yielding polysemantic behavior. That is, for some j , $x_j = \sum_i W_{j,i} a_i$ involves multiple nonzero terms Bills et al. (2023a).

VARIANTS OF SAEs: Variants of SAEs like TopK, JumpReLU, and Gated-SAEs offer increasingly precise control over the mapping between low-level activations and human-understandable concepts, enabling fine-grained analysis and intervention.

TopK-SAEs: Instead of using a soft sparsity constraint (e.g., L1 regularization), TopK-SAEs enforce hard sparsity using a top- K activation function:

$$\mathbf{a} = \text{TopK}(W_{\text{enc}}(\mathbf{x} - \mathbf{b}_{\text{dec}})), \quad (35)$$

which retains only the K largest entries of the preactivation and zeros out the rest. This promotes discrete sparsity and avoids complex hyperparameter tuning.

JumpReLU-SAEs: JumpReLU replaces ReLU with a thresholded step function:

$$\text{JumpReLU}_{\theta}(x) = x \cdot H(x - \theta), \quad (36)$$

where $H(\cdot)$ is the Heaviside step function and θ is a learnable threshold. This allows neurons to activate only above a semantic threshold, aligning with binary behavior observed in some interpretable features. However, the discontinuity makes training difficult due to non-differentiability.

Gated-SAEs: Gated-SAEs introduce a gating mechanism that decouples activation magnitude and presence. Let W_{mag} and W_{gate} be two encoders. Then the feature activation is computed as:

$$\mathbf{a} = (W_{\text{mag}}(\mathbf{x})) \odot H(W_{\text{gate}}(\mathbf{x}) - \theta), \quad (37)$$

where \odot denotes elementwise multiplication. This enables better control over when and how strongly a feature activates, making them easier to train than JumpReLU-SAEs Rajamanoharan et al. (2024).

In this study we utilize two different architectures of SAEs the standard SAE and TopK-SAE.

A.7 ATTRIBUTION FROM SPARSE FEATURE SPACE TO INPUT TOKENS

Let $\mathbf{x}_{\text{input}} \in \mathbb{R}^{d_{\text{input}}}$ denote the input embedding vector (e.g., LLM token embeddings), $\mathbf{x} = f(\mathbf{x}_{\text{input}}) \in \mathbb{R}^d$ the hidden layer activation of the LLM, $\mathbf{a} = \text{Encoder}(\mathbf{x}) \in \mathbb{R}^F$ the SAE sparse feature vector, and $\hat{\mathbf{x}} = W\mathbf{a} + \mathbf{b}$ the reconstructed activation from the SAE decoder. Now suppose we have a sparse attribution vector ψ_i over features \mathbf{a} , i.e., $\psi \in \mathbb{R}^F$, where each ψ_i reflects the importance of SAE feature a_i . We aim to assign importance Φ_k to each input token dimension $x_{\text{input},k}$.

ATTRIBUTION FLOW THROUGH THE ENCODER: We propagate the feature attributions backward through the encoder to the input. Using the chain rule:

$$\Phi_k = \sum_{i=1}^F \psi_i \cdot \frac{\partial a_i}{\partial x_{\text{input},k}} = \sum_{i=1}^F \psi_i \cdot \frac{\partial a_i}{\partial \mathbf{x}} \cdot \frac{\partial \mathbf{x}}{\partial x_{\text{input},k}} \quad (38)$$

where $\frac{\partial a_i}{\partial \mathbf{x}}$ is the encoder Jacobian (SAE layer), and $\frac{\partial \mathbf{x}}{\partial x_{\text{input},k}}$ is the LLM gradient from input token to hidden layer.

This gives us a scalar attribution $\Phi_k \in \mathbb{R}$ for each token/input embedding dimension k .

This represents how much each input token contributes to the sparse SAE features that have been identified as important. In this way, we evaluate the contribution of input features based on the monosemantic behavior of the trained network’s mechanism. Based on our study thus far, we will apply the six attribution methods previously discussed at two levels: from the SAE feature space to the encoder layer, and from the encoder layer to the input embedding space. This dual-level attribution analysis enables us to investigate how interpretable sparse features relate to model internals and ultimately influence the input-level representations.

To this end, we define a two-step attribution mechanism:

STEP 1: ATTRIBUTION FROM SPARSE FEATURES TO ENCODER LAYER

Let $\boldsymbol{\psi} \in \mathbb{R}^F$ represent the importance scores of sparse features (obtained via attribution methods). We propagate these to the encoder layer as:

$$\boldsymbol{\phi}^{\text{enc}} = W\boldsymbol{\psi} \in \mathbb{R}^d, \quad (39)$$

where $\boldsymbol{\phi}^{\text{enc}}$ quantifies the contribution of each encoder neuron to the important SAE features.

STEP 2: ATTRIBUTION FROM ENCODER LAYER TO INPUT

To assign attribution scores to input dimensions, we propagate $\boldsymbol{\phi}^{\text{enc}}$ to the input embedding via the gradient of the encoder:

$$\boldsymbol{\phi}^{\text{input}} = \left(\frac{\partial \mathbf{x}}{\partial \mathbf{x}_{\text{input}}} \right)^\top \boldsymbol{\phi}^{\text{enc}} \in \mathbb{R}^{d_{\text{input}}}. \quad (40)$$

Alternatively, attribution methods (e.g., Integrated Gradients, SHAP) can directly estimate:

$$\boldsymbol{\phi}^{\text{input}} = \text{AttributionMethod}(f, \mathbf{x}_{\text{input}}, \boldsymbol{\phi}^{\text{enc}})$$

This dual-level attribution analysis allows us to connect semantically meaningful sparse features to the raw input representation space.

B SUPPLEMENTARY MATERIAL

B.0 RELATED WORK

B.0.1 ATTRIBUTIONAL INTERPRETABILITY

Attributional interpretability (AtI), a branch of explainable AI (XAI), focuses on explaining model outputs by tracing predictions back to individual input contributions, often using gradient-based methods Bereska & Gavves (2024). While gradients provide insights into the relationship between inputs and outputs, they can be sensitive to perturbations or discontinuities, posing challenges for reliable interpretation.

AtI encompasses various methods for interpreting complex, nonlinear models, including techniques like Local Interpretable Model-agnostic Explanations (LIME; Ribeiro et al. (2016)) and SHapley Additive exPlanations (SHAP; Lundberg & Lee (2017)). In medical imaging, popular attribution techniques include SHAP, Layer-wise Relevance Propagation (LRP; Bach et al. (2015)), and gradient-based methods like GRAD-CAM (Singh et al. (2020)). These methods aim to enhance trust in models and provide valuable insights into decision-making processes. However, they face limitations. For instance, LRP emphasizes positive preactivations, often yielding less precise explanations, while SHAP is computationally intensive due to the complexity of calculating Shapley values Lundberg & Lee (2017). Adaptations like Monte Carlo methods and stratified sampling (e.g., SVARM) have improved the efficiency and precision of certain techniques Kolpaczki et al. (2024).

B.0.2 MECHANISTIC INTERPRETABILITY AND SPARSE AUTOENCODER

Mechanistic interpretability (MI), a key area of explainable AI (XAI), focuses on understanding the internal activation patterns of AI models by analyzing their fundamental components, such as features, neurons, layers, and connections. Unlike AtI, MI takes a bottom-up approach, aiming to uncover the causal relationships and precise computations that transform inputs into outputs. This method identifies specific neural circuits driving behavior and provides a reverse-engineering perspective. Insights from fields like physics, neuroscience, and systems biology further guide the development of transparent and value-aligned AI systems.

A core principle of MI is the concept of polysemanticity, where individual neurons encode multiple concepts, contrasted with monosemanticity, where neurons correspond to a single semantic concept. Polysemanticity reduces interpretability, as neurons represent overlapping features. Structures like sparse autoencoders (SAEs) address this by leveraging the superposition hypothesis, which posits that neural networks use high-dimensional spaces to represent more features than the number of neurons, encoding them in nearly orthogonal directions. SAEs decompose embeddings from deep layers, such as MLPs or transformer attention layers, into higher-dimensional monosemantic representations, aligning activation patterns with specific concepts of interest Cunningham et al. (2023); Elhage et al. (2022a).

Sparse Autoencoder architectures have significantly advanced our understanding of feature representations in language and vision models Gorton (2024). Neural network behavior is often interpreted through *computational circuits*—groups of neurons that compute meaningful functions, such as edge detectors Olah et al. (2020a) or word-copying units Olsson et al. (2022). Leveraging SAE-derived features instead of raw neurons has improved the interpretability of circuits associated with complex behaviors Marks et al. (2024). This shift enables clearer mappings between neuron activations and high-level functions, facilitating validation of model behavior Bereska & Gavves (2024). By aligning internal representations with privileged basis directions—distinct semantic vectors within network layers—researchers further enhance monosemanticity and advance the interpretability of deep models.

B.1 ALZHEIMER DATASET AND PREPROCESSING

B.1.1 PREPROCESSING

The ADNI data Mueller et al. (2005) was downloaded from the Image & Data Archive (IDA) Neu et al. (2023), run by the Laboratory of Neuro Imaging (LONI) at the USC Mark and Mary Stevens Neuroimaging and Informatics Institute. The download comprised folders including information about participants' enrollment, biospecimen, assessments, medical history, imaging and study information. In this work, only baseline ('bl') visit data was extracted, that is - the first visit the patient underwent when joining each study. The number of unique participant's RIDs (subject's roster ID) was then recorded, and the intersection of such identifiers across the baseline datasets was calculated through an overlap matrix assessing participant coverage by considering datasets symmetrically. The obtained result, underwent precise analysis and filtering. Non-informative and administrative columns (i.e.: SOURCE, update_stamp, SITEID, etc.) were removed across all datasets, to then perform a column-wise completeness check to retain only variables with at least 80% of values present and to balance data availability with feature retention. By prioritizing datasets with the highest number of unique RIDs at baseline, pairwise merging based on shared RIDs was performed (i.e.: inner joins), considering the following files: ADAS, NEUROBAT, FAQ, VITALS, DXSUM. Diagnosis data was sorted chronologically according to EXAMDATE and de-duplicated so as to obtain the first - baseline - diagnosis per subject. Moreover, to ensure robust classification, this was complemented by matching data from adni_diagnosisDXSUM files. For data augmentation purposes, demographics data was obtained from adni_demographic_PTDEMOG and merged according to matching RIDs. Biospecimen and medication data were filtered, cleaned and aggregated by participant - however, due to high sparsity and no adherence of column data to the completeness threshold, such information was not included in the final merge. Similarly, no genetic data was included, due to the lack of relevant biological variables with enough completeness, as remaining columns were primarily collection metadata. The final merged dataset - after excluding administrative columns - comprised 2791 unique participant RIDs with comprehensive neuropsychological, clinical, biospecimen, vital sign,

594 and demographic data at baseline, with the following diagnosis count: 1207 patients diagnosed
 595 with Early Mild Cognitive Impairment (EMCI), 441 with Late Mild Cognitive Impairment (LMCI),
 596 and 1143 control subjects. For the binary classification task, EMCI and LMCI subjects were unified
 597 into a unique MCI cohort - mimicking AD vs CN classification, while for the three-class task,
 598 all three subsets were retained, considering only 440 subjects per class, for balancing purposes.
 599 Variables from the obtained merged dataset, were mapped to their descriptions and categorical
 600 values, according to the DATADIC_adni123GO dictionary from ADNI Mueller et al. (2005). Text
 601 was then generated by iterating through each subject row, replacing column names with their
 602 description and appending the corresponding column value for the specific patient. Whereby
 603 categorical values were present, they were replaced with their corresponding textual value (i.e.:
 604 " 'sex': 0 " - was transformed into "The patient's sex is: male"). Two distinct datasets - one for
 605 training and one for testing - were generated from the obtained final datasets, and they were split
 606 into training, testing and validation sets.

607 Another dataset was utilized for further model refinement and finetuning. Specifically, the ad-
 608 ditional data was extrapolated from MRI files from the Latin American Brain Health Institute
 609 (BrainLat) dataset, a multi-site initiative that provides neuroimaging, cognitive, and clinical data
 610 across several countries in the Latin American region Prado et al. (2023). The data included cogni-
 611 tion, demographic and records information of 780 subjects. A pre-processing pipeline similar to
 612 that employed for ADNI, was followed. Namely, after filtering throughout all MRI files, 760 unique
 613 and common MRI IDs - representing each subject - were identified. After dropping subjects with a
 614 higher proportion of data missing, and columns not fulfilling the completeness threshold, median
 615 imputation based on diagnosis group mean was applied for variables with less than 30% of data
 616 missing (such as 'Age' and 'years of education' for example) with the goal of obtaining a more
 617 complete dataset. After dropping administrative and non-informative columns, the final merged
 618 dataset comprised variables deriving from cognitive tests (MOCA - Montreal Cognitive Assess-
 619 ment test and the IFS - INECO Frontal Screening) and participants' demographics. The diagnosis
 620 distribution of the obtained dataset was the following: 101 control subjects (CN), 109 diagnosed
 621 with Fronto-Temporal Dementia (FTD), and 118 subjects with AD. For the binary classification
 622 task, here AD and Fronto-Temporal dementia were unified into a unique cognitively impaired
 623 cohort, similar as to what obtained for ADNI, while for the three-class task, the original labels were
 624 retained. The same process as for ADNI was followed to obtain textual descriptions of BrainLat
 625 patients' data, considering the related dictionary from Prado et al. (2023). Finally, training and
 626 testing files were obtained, whereby each class had 50 representative samples each, both for the
 627 binary and for the three-class classification. The handling of the final split into training, testing,
 628 and validation sets was handled as for ADNI. Throughout the manuscript, the label 'AD' is used
 629 for convenience to denote the MCI cohort in ADNI (in both the binary and three-class settings),
 630 and the AD+FTD cohort in BrainLat (in the binary setting). This choice is purely notational, as the
 631 term functions as a class label rather than a clinical diagnosis, and the emphasis is on the model's
 632 ability to discriminate between the defined classes.

632 B.1.2 DEMOGRAPHIC COMPARISON OF ALZHEIMER'S COHORTS AND MATCHED CONTROLS

633 To ensure demographic comparability and reduce confounding in downstream analyses, we
 634 examined age and sex distributions across each Alzheimer's disease (AD) cohort and control
 635 groups.

636 Considering the cohorts for the binary classification from ADNI Mueller et al. (2005), AD subjects
 637 ($n = 1207$) and the control group ($n = 1143$), it is worth noting that both groups consider subjects
 638 who were born between a range that goes from the 1930s to the 1960s with comparable distribu-
 639 tions. The AD group exhibits sharper age peaks, (Figure 36(a)), while the control group shows
 640 a more uniform spread. A similar pattern is evident from the three-class classification cohorts
 641 (Figure 36(c)), whereby patients diagnosed with LMCI and MCI tend to be demonstrate higher
 642 density at certain points, whereas healthy subjects' birth year distribution tends to be flatter.

643 The gender distribution is uniform across groups, both in binary and three-class classification
 644 (Figures 36(b) and 1(d)), with a slight predominance of female participants in AD groups, but
 645 overall disparity suggests minimal risk of demographic bias.

646 Regarding the BrainLat dataset Prado et al. (2023), similar patterns are evident. Control subjects
 647 are, on average, younger than subjects diagnosed with AD by 4 years, although the distribution

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

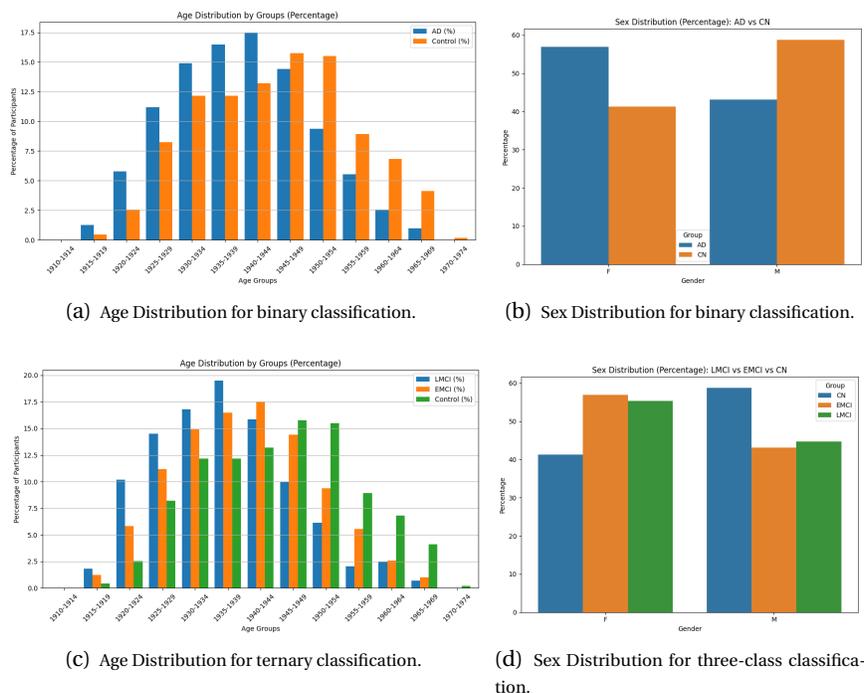


Figure 1: Demographic distributions (age and gender) for Alzheimer’s cohorts and control groups for both binary and ternary classification tasks. The top row refers to the binary task, while the bottom row analyzes cohorts for the ternary classification task.

for AD tends to be more coherently spread than the one for CN (AD cohort mean age: 71, with a standard deviation of 8.7, CN cohort mean age: 67, with standard deviation of 8.5). In the cohorts obtained for the three-class classification task, the age difference remains the same - as AD subjects tend to be the oldest, followed by those belonging to the FTD cohort and CN cohort respectively. Age variability in this case, becomes more comparable between the different diagnoses. Similarly to what was found for ADNI, gender-wise, the data distribution tends to be more skewed toward female participants, both in the AD and in the CN cohorts. The same is found for the subsets obtained for the three-class classification task, whereby female patients diagnosed with AD and FTD represent a higher number than male ones.

B.1.3 PHENOTYPIC AND LIFESTYLE PROFILING

To characterize the ADNI cohorts beyond age and sex, we analyzed phenotypic and lifestyle variables spanning physical health (e.g., systolic and diastolic blood pressure, respiratory and pulse rate, height, weight, body temperature, dominant hand) and behavioral and lifestyle factors (e.g., living situation, marital status, primary language). These features were compared across all four groups to identify significant inter-group differences Mueller et al. (2005).

In the comparison between AD and CN cohorts for the binary classification, a significant difference was found in subjects’ pulse rate ($p < 0.05$) based on independent samples t-test - consistent with the nervous system dysfunction that Alzheimer’s involves. Instead, no significance was found for systolic and diastolic blood pressure, respiratory rate, body temperature and weight. In terms of behavioral and lifestyle factors, a significant difference in marital status - based on Fisher’s exact test - was observed between the two groups. Although most of subjects in the AD and CN groups were married, widowed individuals made up a larger proportion than divorced individuals in the AD group, while the opposite was true for CN subjects. Moreover, the CN group had a higher percentage of individuals who had never been married. Subjects also differed for living situation (Fisher’s exact test). Most subjects diagnosed with AD, lived in a house and smaller proportions lived in - respectively - a condo, an apartment, and a mobile home, with the lowest percentages

702 residing in a retirement community and in an assisted living facility. Although CN subjects also
 703 predominantly lived in a house, they were more likely than AD subjects to live in an apartment or
 704 a condo, followed by a mobile home, an assisted living facility and lastly, a retirement community.
 705

706 B.1.4 ASSESSING COMPATIBILITY BETWEEN IID AND OOD COHORTS

707
 708 The selection of ADNI (IID) and BrainLat (OOD) cohorts was motivated by their demographic
 709 comparability and complementary clinical profiles. As described in Section B.1.2, both datasets
 710 show overlapping age and sex distributions, with balanced ratios and only minor female predomi-
 711 nance. These similarities minimize confounding, ensuring that performance differences reflect
 712 domain shifts rather than demographic bias.

713 Phenotypic and lifestyle profiling (Section B.1.3) revealed moderate inter-group differences, such
 714 as in pulse rate and marital status, consistent with disease-specific traits. ADNI primarily rep-
 715 resents the Alzheimer’s continuum (CN, MCI, LMCI), whereas BrainLat includes FTD, AD, and
 716 controls. Despite differing diagnostic labels, these groups share clinical overlap: FTD often ex-
 717 hibits MCI-like cognitive decline, and LMCI represents a prodromal AD stage Petersen et al. (1999);
 718 Jack et al. (2018); Gorno-Tempini et al. (2011).

719 This overlap establishes a natural testbed for generalization, challenging models trained on
 720 IID data to transfer to OOD settings with related but non-identical diagnoses. The IID/OOD
 721 pairing thus provides a rigorous, clinically meaningful framework to evaluate the adaptability and
 722 robustness of LLM-based diagnostic systems.

723 B.1.5 MODALITIES SUBGROUP EXTRACTIONS

724
 725 Table 1: Variables with character counts, generation order, and estimated tokens (tokens
 726 $\approx \lceil \text{chars}/4 \rceil$).
 727

728 Variable	Description	Chars	Order	729 Tokens (est.)
730 PTGENDER	The participant’s sex is	30	1	8
731 PTDOB	Their Date of Birth is	31	2	8
732 PTDOBY	Their Year of Birth is	28	3	7
733 PTHAND	Their Handedness is	26	4	7
734 PTMARRY	Their Marital status at baseline is	44	5	11
735 PTEDUCAT	Their education in years is	31	6	8
736 PTNOTRT	Participant Retired?	25	7	7
737 PTHOME	Type of Participant residence	54	8	14
738 PTTLANG	Language to be used for testing the Participant	56	9	14
739 PTPLANG	Participant’s Primary Language	39	10	10
740 PTETHCAT	The participant’s Ethnicity is	54	11	14
741 PTRACCAT	Trail Making Test: Race	28	12	7
742 PTSOURCE	Information Source	37	13	10
743 VSWEIGHT	The participant’s weight is	32	14	8
744 VSWTUNIT	The weight was measured in	34	15	9
745 VSBPSYS	The participant’s Systolic - mmHg	40	16	10
746 VSBPDIA	The participant’s Diastolic - mmHg	40	17	10
747 VSPULSE	The participant’s Seated Pulse Rate (per minute) is	56	18	14
748 VSRESP	The participant’s Respirations (per minute) are	51	19	13
749 VSTEMP	The participant’s Temperature is	37	20	10
750 VSTMPSRC	The Temperature Source was	32	21	8
751 VSTMPUNT	The Temperature Units were	38	22	10
752 DXDEP	Depressive symptoms present?	32	23	8
753 CLOCKCIRC	On the Clock Drawing Test the participant answered the follow- ing questions in this way: Approximately circular face	126	24	32
754 CLOCKSYP	Symmetry of number placement	39	25	10
755 CLOCKNUM	Correctness of numbers	31	26	8

Continued on next page

	Variable	Description	Chars	Order	Tokens (est.)
756					
757					
758	CLOCKHAND	Presence of the two hands	34	27	9
759	CLOCKTIME	Presence of the two hands, set to ten after eleven	59	28	15
760	CLOCKSCOR	Clock Drawing Test: Total Score	36	29	9
761	COPYCIRC	On the Clock copying task the participant scored as follows: Approximately circular face	95	30	24
762	COPYSYM	Symmetry of number placement	39	31	10
763	COPYNUM	Correctness of numbers	31	32	8
764	COPYHAND	Presence of the two hands	34	33	9
765	COPYTIME	Presence of the two hands, set to ten after eleven	59	34	15
766	COPYSCOR	Clock copying task: Total Score	36	35	9
767	AVTOT1	On the Auditory Verbal Learning Test the participant scored as follows in each trial: Trial 1 Total	104	36	26
768					
769	AVERR1	Total Intrusions	19	37	5
770	AVTOT2	Trial 2 Total	16	38	4
771	AVERR2	Total Intrusions	19	39	5
772	AVTOT3	Trial 3 Total	16	40	4
773	AVERR3	Total Intrusions	19	41	5
774	AVTOT4	Trial 4 Total	16	42	4
775	AVERR4	Total Intrusions	19	43	5
776	AVTOT5	Trial 5 Total	16	44	4
777	AVERR5	Total Intrusions	19	45	5
778	AVTOT6	Trial 6 Total	16	46	4
779	AVERR6	Total Intrusions	19	47	5
780	AVTOTB	List B Total	15	48	4
781	AVERRB	Total Intrusions	19	49	5
782	CATANIMSC	On the Category Fluency Test Animals the scores were: - Total Correct	73	50	19
783	CATANPERS	Perseverations	17	51	5
784	CATANINTR	Intrusions	13	52	4
785	TRAASCOR	Part A - Time to Complete	29	53	8
786	TRAAERRCOM	Errors of Commission	23	54	6
787	TRAAERROM	Errors of Omission	21	55	6
788	TRABSCOR	Part B - Time to complete	30	56	8
789	TRABERRCOM	Errors of Commission	23	57	6
790	TRABERROM	Errors of Omission	21	58	6
791	AVDEL30MIN	On the Auditory Verbal Learning Test the participant scored as follows: 30 Minute Delay Total	96	59	24
792	AVDELERR1	Total Intrusions	19	60	5
793	AVDELTOT	Recognition Score	20	61	5
794	AVDELERR2	Total Intrusions	19	62	5
795	ANARTERR	American National Adult Reading Test: ANART Total Score (Total # of errors)	81	63	21
796	FAQFINAN	For the Functional Activities Questionnaire the participant scored as follows for each question: Writing checks, paying bills, or balancing checkbook.	151	64	38
797					
798	FAQFORM	Assembling tax records, business affairs, or other papers.	58	65	15
799	FAQSHOP	Shopping alone for clothes, household necessities, or groceries.	64	66	16
800	FAQGAME	Playing a game of skill such as bridge or chess, working on a hobby.	68	67	17
801					
802	FAQBEVG	Heating water, making a cup of coffee, turing off the stove.	60	68	15
803	FAQMEAL	Preparing a balanced meal.	26	69	7
804	FAQEVENT	Keeping track of current events.	32	70	8
805	FAQTV	Paying attention to and understanding a TV program, book, or magazine.	70	71	18
806					
807	FAQREM	Remembering appointments, family occasions, holidays, medi- cations.	66	72	17
808					
809					

Continued on next page

Variable	Description	Chars	Order	Tokens (est.)
FAQTRAVL	Trail Making Test for FAQ score: Traveling out of the neighborhood, driving, or arranging to take public transportation.	121	73	31
FAQTOTAL	Total Score for FAQ is	26	74	7

Based on Table 1, we extracted nine subgroups as follows: Demographics, Vital Signs, Clock Drawing Test, Clock Copying Test, Auditory Verbal Learning Test (version 1), Category Fluency – Animal Test, Auditory Verbal Learning Test (version 2), American National Adult Reading Test, and Functional Activities Questionnaire.

B.1.6 DATASETS CLAIMS

Data used in the preparation of this article was obtained from the Alzheimer’s Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu) on August 8th 2025 (version: "08Aug2025") and it included all ADNI phases. ADNI was launched in 2003 as a public-private partnership, led by Principal Investigator Michael W. Weiner, MD. It aimed at testing whether cognitive, imaging, genetic, clinical, neuropsychological assessment and other biological markers, can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimer’s disease (AD). The goals also include the validation of biomarkers for clinical trials, and the provision of data concerning the diagnosis and progression of Alzheimer’s disease to the scientific community. For up-to-date information, see adni.loni.usc.edu.

B.2 SUMMARY OF TRAINING OUTCOMES FOR LLM ENCODERS ON IID AND OOD DATASETS

We systematically compared the performance of different *fine-tuned encoder models* (BERT, RoBERTa, DistilBERT, ALBERT, BioBERT, ModernBERT) on ADNI (in-domain, IID) and evaluated cross-dataset generalization to BRAINLAT (out-of-domain, OOD). On ADNI, ModernBERT is the strongest encoder across all metrics: *Binary*—Acc: 0.7237, F1: 0.7589, ROC-AUC: 0.8395, AUC-PR: 0.8641. *Three-class*—Acc: 0.6505, F1: 0.6880, ROC-AUC: 0.7867, AUC-PR: 0.7848. BioBERT and RoBERTa are the most competitive baselines but remain below ModernBERT. In the zero-shot transfer from ADNI to BRAINLAT and the binary classification task, ModernBERT achieved modest performance with an average accuracy of approximately 0.55. In a representative run, the model reached 0.53 accuracy, 0.52 precision, 0.70 recall, an F1 score of 0.60, and both *ROC-AUC* and *AUC-PR* near 0.58. These results highlight a conservative decision threshold and the difficulty of domain transfer without adaptation. Introducing few-shot supervision improved performance moderately. In the K -shot regime, accuracy increased by up to 0.10 compared to zero-shot, reaching approximately 0.62 at $K = 10$, with parallel gains in F1. The ROC-AUC and AUC-PR metrics remained high and stable, suggesting that limited supervision can partially mitigate domain shift but does not fully bridge the gap. LoRA-based parameter-efficient adaptation produced results comparable to few-shot training, offering efficiency in training without substantial additional gains in predictive performance. By contrast, full fine-tuning of all pretrained weights on BRAINLAT yielded the strongest improvements, with accuracy rising to 0.84 and consistent gains across F1, ROC-AUC, and AUC-PR. These results demonstrate that full supervised adaptation remains the most effective approach to address domain shift when sufficient labeled data are available. In the three-class BRAINLAT setting, zero-shot transfer from ADNI yielded limited generalization (*Accuracy*=0.40, *F1*=0.41, *ROC-AUC*=0.44), reflecting the challenge of domain and class shifts. Few-shot adaptation ($K = 10$) moderately improved performance (*Accuracy*=0.49, *F1*=0.48), while LoRA-based fine-tuning achieved comparable results (*Accuracy*=0.50, *F1*=0.48). Full fine-tuning produced the strongest gains, reaching *Accuracy*=0.69, *F1*=0.73, and *ROC-AUC*=0.81. These findings confirm that, although limited supervision aids adaptation, full parameter optimization is essential for robust multi-class generalization across cohorts. However, this setting is outside the scope of this work: we focus on explanation performance under OOD conditions without training on the OOD cohort (i.e., without full fine-tuning).

Therefore, for all downstream analyses we *stick with ModernBERT*: in the IID setting we use *ModernBERT* fine-tuned on ADNI (best overall on in-domain tasks), and in the OOD setting we use *ModernBERT* in a zero-shot configuration on BRAINLAT (best overall under out-of-domain

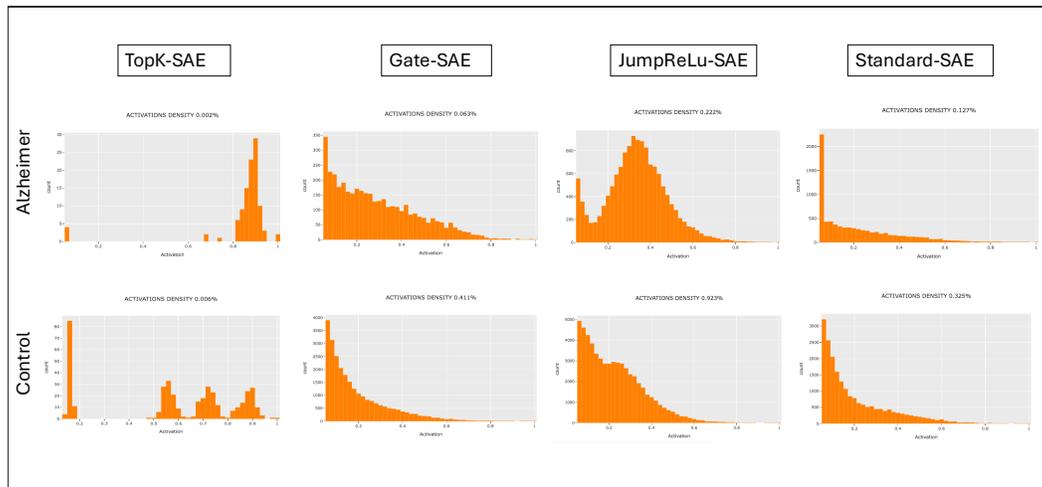


Figure 2: Latent space projections from four SAE variants (TopK-SAE, Gate-SAE, JumpReLU-SAE, Standard-SAE) applied to Alzheimer’s and Control groups. TopK-SAE shows the clearest group separation, highlighting its superior ability to extract interpretable, clinically relevant features.

conditions). All subsequent explainability analyses were conducted using the final (22nd) layer of *ModernBERT*.

B.3 HYPERPARAMETER TUNING FOR THE RECONSTRUCTION OPTIMIZER AND SAE MODELS.

A thorough hyperparameter tuning process was conducted for each simulation (Figures 3, 4, 5). The explanation optimizer was trained with learning rates of $2e^{-2}$, $2e^{-3}$, $2e^{-4}$, and $2e^{-5}$, with the best performance observed at $2e^{-4}$. Various combinations of the weighting parameters $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ were tested—for example, (0.3, 0.2, 0.25, 0.25)—with the optimal configuration found to be (0.1, 0.3, 0.1, 0.5). For the UMAP constraints, subgroup levels were evaluated across several scales: no UMAP, every $4\times$ batch size, $10\times$ batch size, and full cohort level. The best performance was achieved at the $4\times$ batch size level. Regarding the SAE (Sparse Autoencoder), different model variants were evaluated, including *Standard*, *TopK*, *JumpReLU*, and *GATE*, as described in the Methods section. Among these, the *TopK* variant achieved the best results. Feature space depths of $16\times$, $32\times$, and $64\times$ were tested, with $32\times$ providing the best trade-off between sparseness and reconstruction performance. The final simulation and training settings included the Adam optimizer (Kingma & Ba, 2014) with a learning rate of $2e^{-4}$, a batch size of 64, and 200 total training steps, using a 50/50 train-validation split. The learning rate schedule followed a fixed-step approach with a step size of 150 and a decay factor (gamma) of 0.95. For the SAE training, we used 6,000 training steps, 200,000 training tokens, a learning rate of $5e^{-5}$, and a model dimension of 768, consistent with the 22-layer *Modern-BERT* architecture. The context size was 512, with warm-up steps of 1,000, learning rate decay steps of 1,200, and L1 warm-up steps of 300. Finally, explanation metrics such as ROS, RIS, and sparseness were computed using default configurations from the *quantus* Python package (Hedström et al., 2023). Figure 3 presents a comparative visualization of activation patterns projections generated by different Sparse Autoencoder (SAE) variants—TopK-SAE, Gate-SAE, JumpReLU-SAE, and Standard-SAE—applied to two subject groups: Alzheimer’s and Control. While the specific axes and metrics are not labeled, the separation between the two groups provides insight into the effectiveness of each SAE in producing disentangled, semantically meaningful representations. Among the models, the TopK-SAE exhibits the clearest separation between the Alzheimer’s and Control cohorts, suggesting superior performance in capturing clinically relevant patterns. This visual evidence supports the paper’s central claim that monosemantic representations enhance interpretability and robustness in clinical applications of LLMs.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

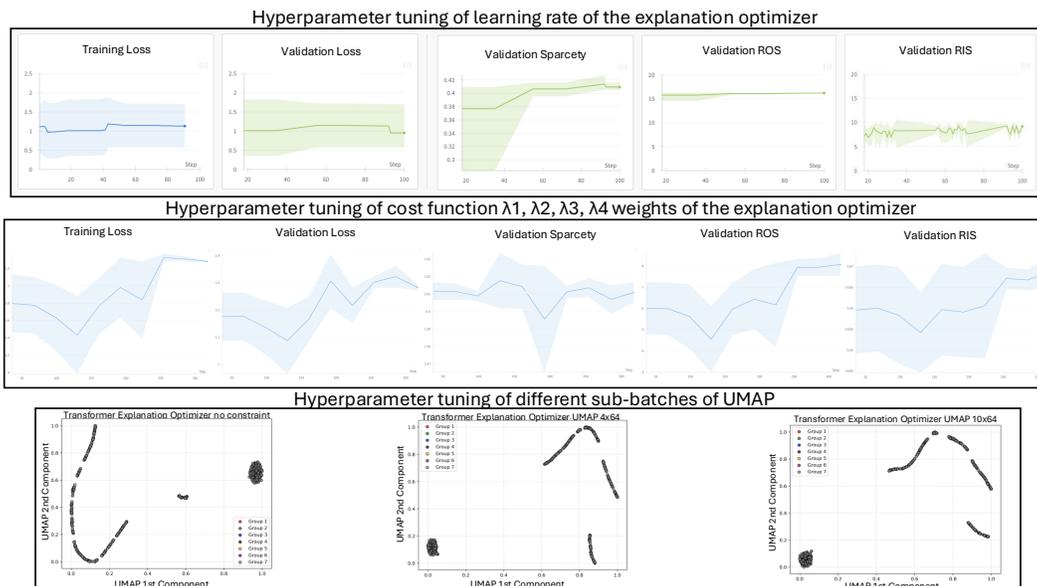


Figure 3: **Hyperparameter tuning of the explanation optimizer and UMAP settings.** Top row: impact of learning rate on training loss, validation loss, sparseness, ROS, and RIS metrics. Middle row: sensitivity analysis of the explanation cost weights $\lambda_1, \lambda_2, \lambda_3,$ and λ_4 , showing trade-offs between attribution sparseness and robustness. Bottom row: UMAP projections of token-level attribution spaces under different sub-batch configurations, revealing how UMAP resolution influences the geometric structure of explanations.

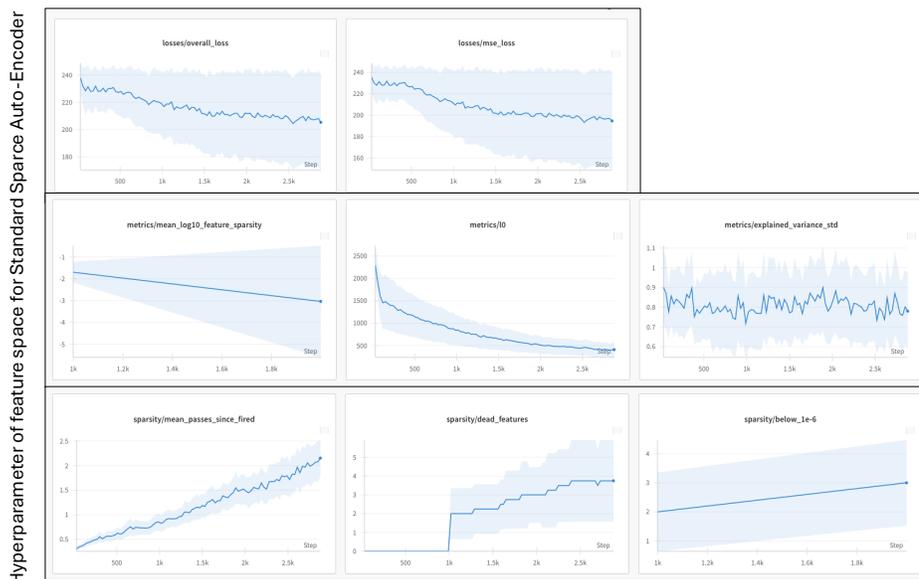


Figure 4: **Hyperparameter tuning of the feature space for the Standard Sparse Autoencoder (SAE).** The plots track training dynamics and sparseness characteristics across training steps. Top row: loss trends for overall and reconstruction loss. Middle row: log-sparsity metric, Kullback–Leibler divergence (KL), and explained variance standard deviation. Bottom row: progression of sparsity across mean-poisson stem-freed features, fixed features, and a threshold-based view ($1e-6$). These results guide optimal SAE configurations for producing monosemantic feature representations.

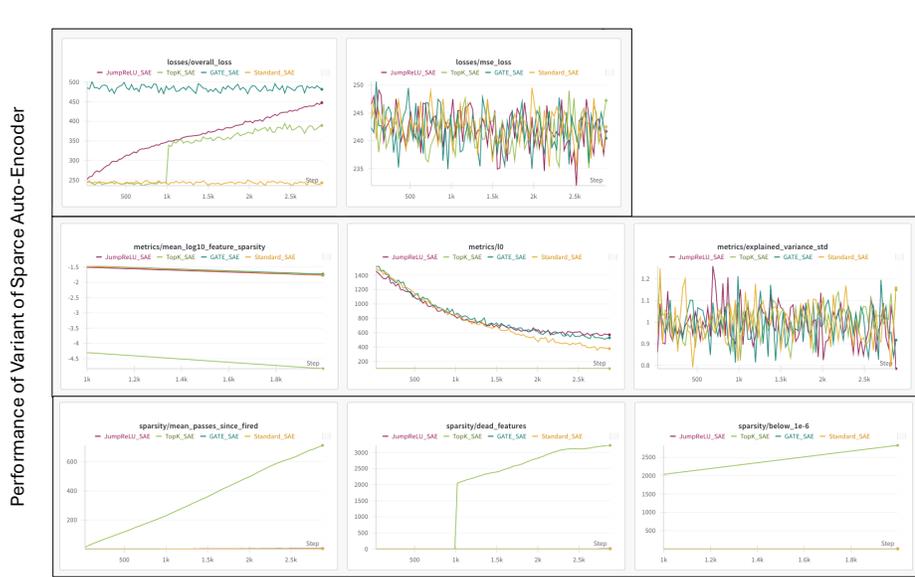


Figure 5: Performance comparison of different variants of the Sparse Autoencoder (SAE). Top row: overall and reconstruction loss across training steps for JumpInit-SAE, Top- k SAE, Gated-SAE, and Standard-SAE. Middle row: log-sparsity metric, KL divergence, and explained variance standard deviation, showing divergence in regularization behavior. Bottom row: sparsity progression for mean-poisson stem-freed features, fixed feature count, and a thresholded view ($1e-6$). JumpInit-SAE shows early convergence in sparsity, while Gated-SAE maintains tighter control over variance. These results highlight trade-offs between sparsity enforcement mechanisms and attributional stability.

B.4 EXTRA RESULTS

Figure 6 compares training dynamics and interpretability metrics on ADNI for binary (top row) and three-class (bottom row) classification. Each subfigure shows three variants of the Explanation Optimizer: *without* SAE (black), *with* SAE (brown), and *with* SAE + linear UMAP constraint (grey). Each row in the figure presents the model’s behavior over training steps across six key metrics: train loss, validation loss, UMAP reconstruction error (MSE), Relative Output Stability (ROS), Relative Input Stability (RIS), and sparsity. Across both tasks, all training ROS and RIS values for the SAE-based variants (brown/grey) are consistently lower than the no-SAE baseline (black), indicating improved attributional robustness. While sparseness does decrease when introducing SAE, the reduction is modest; adding the linear UMAP constraint (grey) achieves a better balance, maintaining relatively high sparsity while keeping low RIS/ROS. Finally, the training and validation curves track closely and remain smooth for the SAE variants, providing no evidence of overfitting: validation loss follows training loss without widening gaps in either the binary or three-class setting. Overall, the SAE-enhanced Explanation Optimizer demonstrates significantly improved performance across all interpretability metrics, supporting the hypothesis that enforcing monosemantic representations improves explanation clarity and reliability—especially in high-stakes clinical contexts like Alzheimer’s disease classification.

Across IID (ADNI) and OOD (BrainLat) settings, and for both binary (Alzheimer vs. Control) and three-class (Control/LMCI/MCI) tasks, the tables reveal a consistent stability–sparsity frontier driven by the proposed explanation optimizers and the presence of a monosemantic bottleneck (SAE). In the binary IID case (Table 2), SAE substantially improves stability for explainers that learn features—most notably Layer Conductance and especially TEO—with large drops in RIS/ROS for both Alzheimer and Control, while Activation with SAE increases RIS/ROS and is therefore less robust. In the binary OOD case (Table 3), these patterns persist and even strengthen: TEO with SAE bottleneck attains the lowest RIS/ROS overall, demonstrating strong cross-dataset stability, whereas TEO–UMAP recovers higher sparseness (>0.40) at the cost of higher RIS/ROS than TEO

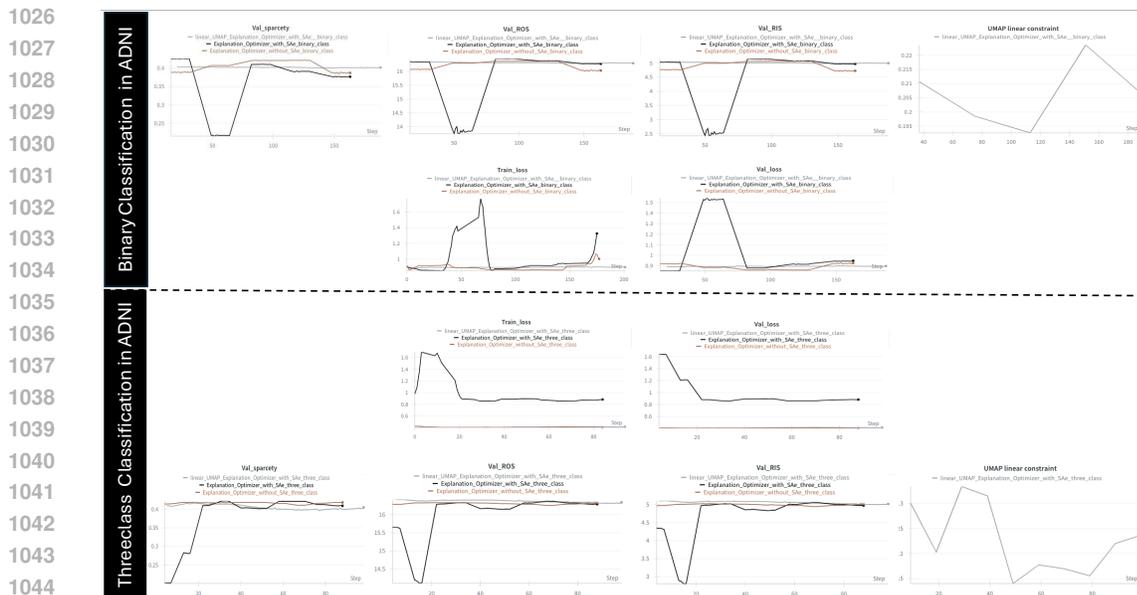


Figure 6: Training and interpretability dynamics on ADNI for binary (top) and three-class (bottom) classification. *Each subfigure includes three variants: Explanation Optimizer without SAE (black), with SAE (brown), and with SAE + linear UMAP constraint (grey.* For each variant we plot train loss, validation loss, UMAP reconstruction MSE (linear UMAP constrain), Relative Output Stability (ROS), Relative Input Stability (RIS), and sparcety, cohorts shown separately. SAE reduces volatility and lowers UMAP MSE and RIS/ROS versus the no-SAE baseline; adding a linear UMAP constraint on top of SAE further improves manifold structure and attribution stability, at a minor cost in sparsity.

with SAE, offering a tunable sparsity–stability trade-off. In the three-class IID setting (Table 4), Feature Ablation is the sparsity leader across Control/LMCI/MCI (0.52–0.53) with moderate, steady RIS/ROS; Layer Conductance with SAE markedly reduces RIS/ROS for LMCI/MCI; and TEO with SAE again delivers the most stable attributions across all classes (lowest RIS/ROS), albeit with reduced sparseness. The same rank ordering holds OOD (Table 5): TEO with SAE remains the stability winner for Control/LMCI/MCI, TEO–UMAP trades some stability for additional sparsity, and Feature Ablation remains the simplest high-sparsity baseline. Throughout all tables, gradient-formulaic methods (Grad-SHAP, Guided Backprop, Integrated Gradients) show near-invariant RIS/ROS (5.6/16.93) regardless of SAE, class, or domain, indicating that SAE chiefly benefits learned-attribution methods. Collectively, Tables 2, 3, 4, and 5 support three conclusions: (i) adding an SAE bottleneck reliably lowers RIS/ROS where explanations are learned (Layer Conductance, TEO), (ii) TEO with SAE is the default when stability is paramount, while TEO–UMAP is preferred when higher sparsity is required, and (iii) the class-wise and IID to OOD behaviors are consistent, underscoring the robustness of monosemantic representations for clinical explanation.

B.5 STATISTICAL ANALYSIS

We conducted both parametric and non-parametric statistical tests on the binary and three-class classification performance of all classes and tasks in the ADNI cohort to assess the significance of differences introduced by the monosemantic bottleneck (SAE) in traditional attribution techniques, focusing on the metrics of Sparseness, RIS, and ROS.

For the binary classification task, in both the Control and Alzheimer’s groups, paired testing demonstrated that SAE produced robust and statistically significant reductions in attribution-based measures and Complexity, while effects on RIS were smaller but still reliable, and changes in ROS were modest and often non-significant after correction. In the Control group, Complexity decreased from 0.3377 ± 0.0017 (no-SAE) to 0.3140 ± 0.0010 (SAE), yielding $t(29) = 64.0$, $p = 1.5 \times 10^{-47}$ (FDR $q < 10^{-46}$), and RIS declined from 9.313 ± 0.143 to 9.175 ± 0.109 , $t(29) = 4.22$,

Table 2: Evaluation scores with and without SAE. Values are mean \pm std. Classes: Alzheimer and Control. Metrics: Sparseness (higher is better), RIS (lower is better), ROS (lower is better). Column order: Alzheimer (No SAE), Alzheimer (SAE), Control (No SAE), Control (SAE). All evaluation metrics were calculated on 200 randomly selected patients from each class (binary-class classification task) in the ADNI testing cohort (IID). Abbreviations, DEO: Diffusion Explanation Optimizer, TEO: Transformer Explanation Optimizer, TEO-UMAP: Transformer Explanation Optimizer with UMAP constraint.

Method	Metric	Alzheimer		Control	
		No SAE	SAE	No SAE	SAE
Activation	Sparseness	0.316364045 \pm 0.007573187	0.296615553 \pm 0.007063087	0.256150148 \pm 0.01759176	0.251987915 \pm 0.004701258
	RIS	14.30227 \pm 0.368612837	21.3084024 \pm 0.311021506	14.23653893 \pm 0.338127875	19.32752421 \pm 0.932339244
	ROS	25.54851914 \pm 0.482826721	32.61738364 \pm 0.307928076	25.54874248 \pm 0.328628877	30.6394217 \pm 0.933402318
Layer Conduct	Sparseness	0.396588773 \pm 0.026146226	0.391508745 \pm 0.007549659	0.374476778 \pm 0.007092037	0.247974319 \pm 0.007908586
	RIS	12.39850671 \pm 2.640648847	5.628509596 \pm 0.023609264	5.650216593 \pm 0.0391015	5.614140617 \pm 0.018407327
	ROS	23.146556 \pm 1.686524073	16.94708243 \pm 0.010343602	16.9614573 \pm 0.031757755	16.93014311 \pm 0.005191254
Feature Ablation	Sparseness	0.523581491 \pm 0.009806381	0.523492234 \pm 0.010441342	0.525551296 \pm 0.011012696	0.526520447 \pm 0.008837301
	RIS	23.15233791 \pm 0.819793812	23.56094785 \pm 0.103321598	22.56110559 \pm 0.288403561	23.62208862 \pm 0.093292383
	ROS	33.90884482 \pm 0.161311922	34.92976979 \pm 0.074540131	33.90759033 \pm 0.374703897	34.9726397 \pm 0.10327352
Gradinet-SHAP	Sparseness	0.319169309 \pm 0.004303288	0.082047681 \pm 0.015464469	0.433346255 \pm 0.003004736	0.133912362 \pm 0.009943283
	RIS	5.623104979 \pm 0.023650419	5.621792106 \pm 0.022658996	5.632513362 \pm 0.023548461	5.619618915 \pm 0.019004514
	ROS	16.93566464 \pm 0.001800535	16.93449562 \pm 1.54604E-05	16.94606355 \pm 0.002246403	16.93473306 \pm 5.99931E-06
Gradient Activation	Sparseness	0.327713636 \pm 0.03837053	0.203454702 \pm 0.011686877	0.249969957 \pm 0.023039465	0.16678783 \pm 0.00722766
	RIS	5.614858128 \pm 0.019338754	5.625220574 \pm 0.021340754	5.616992957 \pm 0.021780592	5.617270064 \pm 0.022114697
	ROS	16.93028085 \pm 0.003427604	16.93433453 \pm 6.78428E-05	16.934673 \pm 1.43645E-14	16.93473306 \pm 4.02532E-05
Integrated-Gradient	Sparseness	0.298289818 \pm 0.008006549	0.121161362 \pm 0.005775061	0.430359021 \pm 0.006572262	0.064427234 \pm 0.005909053
	RIS	5.620585979 \pm 0.018022403	5.622360787 \pm 0.017750318	5.62793468 \pm 0.018989203	5.621391149 \pm 0.016872007
	ROS	16.93257772 \pm 0.001464829	16.93453232 \pm 1.20129E-05	16.94336632 \pm 0.002408932	16.93456653 \pm 8.33237E-06
DEO	Sparseness	0.338261111 \pm 0.003260844	0.337375 \pm 0.00290587	0.337742857 \pm 0.001742551	0.314044444 \pm 0.001036523
	RIS	9.283888889 \pm 0.080010212	9.279 \pm 0.064555158	9.313125 \pm 0.142722049	9.175 \pm 0.108803655
	ROS	20.63421053 \pm 0.086558637	20.615 \pm 0.088049029	20.61588235 \pm 0.202578961	20.515 \pm 0.129878486
TEO	Sparseness	0.421975723 \pm 0.000305212	0.267210213 \pm 0.001025675	0.419939638 \pm 0.000480888	0.268167213 \pm 0.000728522
	RIS	5.051961362 \pm 0.019221728	1.622662574 \pm 0.17080061	5.06881834 \pm 0.01838977	0.996401319 \pm 0.263922792
	ROS	16.35285123 \pm 0.00563874	12.92504253 \pm 0.170261034	16.37765691 \pm 0.001096906	12.29830928 \pm 0.261259725
TEO-UMAP	Sparseness	N/A	0.39891406 \pm 0.000414208	N/A	0.40566988 \pm 0.00031341
	RIS	N/A	5.439373370 \pm 0.033211570	N/A	5.47087230 \pm 0.17460810
	ROS	N/A	16.3036705 \pm 0.00333634	N/A	16.21021807 \pm 0.0078926

$p = 9.5 \times 10^{-5}$ (FDR $q = 1.9 \times 10^{-4}$), both clearly rejecting the null hypothesis, whereas ROS decreased slightly from 20.616 ± 0.203 to 20.515 ± 0.131 , $t(29) = 2.30$, $p = 0.026$ (FDR $q = 0.026$), a marginal result that did not withstand correction. Attribution metrics showed the largest SAE effects: Grad-SHAP dropped from 0.4333 ± 0.0030 to 0.1339 ± 0.0099 ($p < 10^{-50}$), Guided Backprop from 0.2500 ± 0.0230 to 0.1668 ± 0.0072 ($p < 10^{-19}$), Integrated Gradients from 0.4304 ± 0.0066 to 0.0644 ± 0.0059 ($p < 10^{-80}$), and Optimizer from 0.4199 ± 0.0005 to 0.2682 ± 0.0007 ($p < 10^{-100}$), all leading to decisive rejection of the null. For the Alzheimer’s group, the same direction of effects was observed: Complexity decreased by -0.024 ($p < 10^{-10}$), RIS by -0.12 ($p = 4.6 \times 10^{-4}$), both rejecting the null, while ROS declined by -0.09 but did not reach significance ($p = 0.073$, FDR $q = 0.11$). Attribution metrics again showed dramatic reductions under SAE, with Grad-SHAP ($p < 10^{-45}$), Guided Backprop ($p = 3.2 \times 10^{-7}$), Integrated Gradients ($p < 10^{-55}$), and Optimizer ($p < 10^{-95}$) all supporting strong rejection of the null. Together these results demonstrate that SAE reliably improves attribution stability and reduces Complexity and RIS in both groups, with ROS showing only weak or inconsistent improvement.

For the three-class classification task, we evaluated whether SAE changed the three target metrics (Complexity, RIS, ROS) relative to no-SAE using paired t -tests and Wilcoxon signed-rank tests for each clinical group (Control, MCI, LMCI), applying Benjamini-Hochberg FDR to control multiplicity. For the MCI group, ROS showed the clearest and most consistent improvement with SAE: the paired t -test yielded $t(17) = -10.12$, $p = 1.30 \times 10^{-8}$ (FDR $q = 3.90 \times 10^{-8}$), and the Wilcoxon test yielded $W = 0$, $p = 8.0 \times 10^{-6}$ (FDR $q = 2.3 \times 10^{-5}$), with a very large paired Cohen’s $d = -2.39$ and rank-biserial correlation $r_{rb} = -1.00$, indicating markedly lower ROS under SAE (mean difference -0.904 ; SAE 20.672 vs. no-SAE 21.576). RIS in MCI also decreased with SAE by non-parametric testing: the paired t -test did not reach significance ($t(18) = -0.785$, $p = 0.443$, FDR $q = 0.443$), whereas the Wilcoxon test detected a reduction ($W = 19$, $p = 0.00117$, FDR $q = 0.00117$), with small effect sizes ($d = -0.18$, $r_{rb} = -0.80$; mean difference -0.481 ; SAE 10.528 vs. no-SAE 11.010). In contrast, Complexity in MCI increased with SAE according to

Table 3: Evaluation scores with and without SAE. Values are mean \pm std. Classes: Alzheimer and Control. Metrics: Sparseness (higher is better), RIS (lower is better), ROS (lower is better). Column order: Alzheimer (No SAE), Alzheimer (SAE), Control (No SAE), Control (SAE). All evaluation metrics were calculated on 50 randomly selected patients from each class (binary-class classification task) in the **BrainLat** testing cohort (OOD). Abbreviations, TEO: Transformer Explanation Optimizer, TEO-UMAP: Transformer Explanation Optimizer with UMAP constraint.

Method	Metric	Alzheimer (No SAE)	Alzheimer (SAE)	Control (No SAE)	Control (SAE)
Activation	Sparseness	N/A	0.1533105 \pm 0.010287697	N/A	0.3965415 \pm 0.030322127
	RIS	N/A	19.162526 \pm 0.364196762	N/A	18.2411505 \pm 0.539197156
	ROS	N/A	31.28270375 \pm 1.541354976	N/A	29.04058325 \pm 0.473071038
Layer Conduct	Sparseness	N/A	0.239227 \pm 0.029777681	N/A	0.25431875 \pm 0.020993978
	RIS	N/A	6.1620695 \pm 0.149481666	N/A	6.21494775 \pm 0.20764754
	ROS	N/A	16.94383425 \pm 0.007089038	N/A	16.9402505 \pm 0.004966063
Feature Ablation	Sparseness	N/A	0.5287845 \pm 0.00700955	N/A	0.52849725 \pm 0.004358269
	RIS	N/A	23.5834065 \pm 0.064538961	N/A	24.14743175 \pm 0.115957516
	ROS	N/A	34.65309725 \pm 0.252584222	N/A	34.961296 \pm 0.220544503
Gradinet-SHAP	Sparseness	N/A	0.120076 \pm 0.014392456	N/A	0.057057 \pm 0.027064338
	RIS	N/A	6.04400225 \pm 0.039573077	N/A	6.030265 \pm 0.047136969
	ROS	N/A	16.93474475 \pm 5.76852E-05	N/A	16.93475925 \pm 5.7373E-06
Gradient Activation	Sparseness	N/A	0.1139535 \pm 0.01766843	N/A	0.062973 \pm 0.006903384
	RIS	N/A	6.032837 \pm 0.027736802	N/A	6.0338695 \pm 0.039792395
	ROS	N/A	16.93468825 \pm 3.59398E-06	N/A	16.934848 \pm 3.74789E-05
Integrated-Gradient	Sparseness	N/A	0.0642685 \pm 0.005166108	N/A	0.0143455 \pm 0.000312693
	RIS	N/A	6.05793275 \pm 0.045559192	N/A	6.0275535 \pm 0.033936686
	ROS	N/A	16.9347635 \pm 7.76745E-06	N/A	16.934873 \pm 1.06145E-05
TEO	Sparseness	N/A	0.26914625 \pm 0.001645095	N/A	0.272516 \pm 0.000382866
	RIS	N/A	0.683544 \pm 0.667616072	N/A	0.47335295 \pm 0.280125046
	ROS	N/A	11.52356 \pm 0.659063208	N/A	11.213036 \pm 0.51496551
TEO-UMAP	Sparseness	N/A	0.398914 \pm 0.00047836	N/A	0.40425175 \pm 0.002851775
	RIS	N/A	5.43937375 \pm 0.038349421	N/A	5.42815425 \pm 0.194389712
	ROS	N/A	16.303675 \pm 0.003852471	N/A	16.157664 \pm 0.105405246

Table 4: Evaluation scores with and without SAE. Values are mean \pm std. Classes: Control, LMCI, and MCI. Metrics: Sparseness (higher is better), RIS (lower is better), ROS (lower is better). All the evaluation metrics were computed on 100 randomly selected patients from each class (three-class classification task) in the testing cohort in **ADNI** dataset (IID). Column order: Control (No SAE), Control (SAE), LMCI (No SAE), LMCI (SAE), MCI (No SAE), MCI (SAE). Abbreviations, TEO: Transformer Explanation Optimizer, TEO-UMAP: Transformer Explanation Optimizer with UMAP constraint.

Method	Metric	Control		LMCI		MCI	
		No SAE	SAE	No SAE	SAE	No SAE	SAE
Activation	Sparseness	0.302953824 \pm 0.037699004	0.345031647 \pm 0.009533629	0.271540783 \pm 0.038404292	0.264362053 \pm 0.062996342	0.262558 \pm 0.03794219	0.309052111 \pm 0.060296425
	RIS	14.40424929 \pm 0.165969844	18.99683335 \pm 4.483970104	15.07860783 \pm 1.975358824	18.42309221 \pm 2.351829228	16.65684576 \pm 2.82076491	19.42272406 \pm 3.474528306
	ROS	25.721685 \pm 0.170545275	30.30262794 \pm 0.263575977	26.39512961 \pm 1.972727766	29.73331658 \pm 4.847156437	27.96063871 \pm 2.823294834	30.74189461 \pm 6.104484774
Layer Conduct	Sparseness	0.231524647 \pm 0.009587223	0.331457882 \pm 0.006159197	0.362282261 \pm 0.00636308	0.246395316 \pm 0.062800978	0.305312706 \pm 0.007627518	0.292950278 \pm 0.057776928
	RIS	5.626004529 \pm 0.020886163	5.622249412 \pm 0.014860465	13.14289435 \pm 0.325483065	5.623582526 \pm 0.909878534	6.608303235 \pm 2.236340343	5.629072111 \pm 1.130688807
	ROS	16.93496553 \pm 0.004532066	16.93904894 \pm 0.010478577	24.5003352 \pm 0.411948845	16.93375379 \pm 2.745745658	17.91907076 \pm 2.258068808	16.93837556 \pm 3.406698804
Feature Ablation	Sparseness	0.523915176 \pm 0.006693817	0.526105 \pm 0.01204565	0.522595565 \pm 0.009666847	0.526753263 \pm 0.083976662	0.522188941 \pm 0.009398367	0.525710278 \pm 0.104819481
	RIS	23.32498194 \pm 0.410939712	23.07664533 \pm 0.140278365	22.24471861 \pm 0.162941689	21.97939553 \pm 3.680221107	23.49843053 \pm 0.458719189	23.00058106 \pm 4.540229242
	ROS	34.66532071 \pm 0.457995824	34.41792582 \pm 0.146270917	33.60637257 \pm 0.161507737	33.30711989 \pm 5.499688014	34.87369265 \pm 0.440653725	34.31625272 \pm 6.805012455
Gradinet-SHAP	Sparseness	0.231029588 \pm 0.020644371	0.184435333 \pm 0.014766351	0.129241348 \pm 0.032633243	0.301055444 \pm 0.072119068	0.089142118 \pm 0.013138265	0.288114875 \pm 0.061769427
	RIS	5.61894294 \pm 0.013910382	5.621940533 \pm 0.025299231	5.615247522 \pm 0.014396146	5.621698722 \pm 0.946218825	5.629231 \pm 0.018682697	5.618608438 \pm 1.166983293
	ROS	16.93382388 \pm 0.002076107	16.934845 \pm 0.0000146059	16.92551835 \pm 0.001442338	16.93477544 \pm 2.862502663	16.93917776 \pm 0.002118976	16.93476019 \pm 3.524101451
Guided Backprop	Sparseness	0.269674235 \pm 0.006145842	0.229587733 \pm 0.00357389	0.383890696 \pm 0.017652114	0.431001667 \pm 0.115600632	0.291671824 \pm 0.020033511	0.257909625 \pm 0.109537561
	RIS	5.629017941 \pm 0.022520137	5.621027267 \pm 0.019431586	5.627154783 \pm 0.021213565	5.629664833 \pm 0.947801076	5.626876882 \pm 0.019349509	5.617237063 \pm 1.166361854
	ROS	16.934673 \pm 0	16.9348278 \pm 0.0000122544	16.93392839 \pm 0.000750952	16.93466433 \pm 2.862491635	16.93401118 \pm 0.00064598	16.93471594 \pm 3.524085042
Integrated Gradient	Sparseness	0.045146294 \pm 0.007116898	0.263864 \pm 0.004189771	0.10839887 \pm 0.026164396	0.3889465 \pm 0.084086627	0.110163824 \pm 0.015744148	0.266043 \pm 0.090525251
	RIS	5.62077706 \pm 0.021492346	5.6209158 \pm 0.020949259	5.609370435 \pm 0.017753051	5.628201333 \pm 0.947585532	5.628253647 \pm 0.020904475	5.620282063 \pm 1.166468694
	ROS	16.93310612 \pm 0.001097405	16.93475353 \pm 0.0000121647	16.92756339 \pm 0.000606244	16.93434683 \pm 2.862457357	16.93384506 \pm 0.001969998	16.93460281 \pm 3.524039
TEO	Sparseness	0.391835667 \pm 0.000814648	0.268163938 \pm 0.064942517	0.413087063 \pm 0.000325772	0.285971625 \pm 0.037421241	0.390886118 \pm 0.004742559	0.283782105 \pm 0.052291259
	RIS	4.807986067 \pm 0.018432249	1.546787813 \pm 0.11712595	5.093836 \pm 0.18806243	2.264221 \pm 0.487706388	4.828254294 \pm 0.037276868	2.161698368 \pm 0.454717751
	ROS	16.1172194 \pm 0.008959191	12.858954 \pm 0.117943866	16.40431106 \pm 0.002382358	15.56455925 \pm 2.274547046	16.13535541 \pm 0.032411959	13.46760021 \pm 2.764087874
TEO-UMAP	Sparseness	N/A	0.39734881 \pm 0.07492051	N/A	0.41611163 \pm 0.08696112	N/A	0.41175421 \pm 0.237175073
	RIS	N/A	5.10662522 \pm 0.20827341	N/A	5.10165242 \pm 0.16974677	N/A	5.11160575 \pm 0.1072146
	ROS	N/A	16.412282 \pm 0.84387466	N/A	16.4031485 \pm 3.86138978	N/A	16.4088329 \pm 0.482439623

the Wilcoxon test ($W = 17$, $p = 7.90 \times 10^{-4}$, FDR $q = 0.00117$), while the paired t -test was non-significant ($t(18) = 1.112$, $p = 0.281$, FDR $q = 0.421$); effect sizes were small-to-moderate ($d = 0.26$, $r_{\text{fb}} = 0.821$; mean difference $+0.510$; SAE 1.175 vs. no-SAE 0.665). For the **Control** and **LMCI** groups, the pasted records contained incomplete pairs that prevented reliable paired testing

Table 5: Evaluation scores with and without SAE. Values are mean \pm std. Classes: Control, LMCI, and MCI. Metrics: Sparseness (higher is better), RIS (lower is better), ROS (lower is better). Column order: Control (No SAE), Control (SAE), LMCI (No SAE), LMCI (SAE), MCI (No SAE), MCI (SAE). All evaluation metrics were calculated on 50 randomly selected patients from each class (three-class classification task) in the **BrainLat** testing cohort (OOD). Abbreviations, TEO: Transformer Explanation Optimizer, TEO-UMAP: Transformer Explanation Optimizer with UMAP constraint.

Method	Metric	Control		LMCI		MCI	
		No SAE	SAE	No SAE	SAE	No SAE	SAE
Activation	Sparseness	N/A	0.4504596 \pm 0.037517308	N/A	0.1907182 \pm 0.001584043	N/A	0.140032667 \pm 0.012393238
	RIS	N/A	19.0939584 \pm 0.179649958	N/A	18.6406378 \pm 0.911739674	N/A	18.07866133 \pm 0.031343749
	ROS	N/A	29.9240256 \pm 0.253797244	N/A	29.5628004 \pm 0.89780058	N/A	28.858315 \pm 0.045164984
Layer Conduct	Sparseness	N/A	0.3252352 \pm 0.014541625	N/A	0.185706 \pm 0.007343355	N/A	0.200561 \pm 0.011873024
	RIS	N/A	6.2120432 \pm 0.245036193	N/A	6.054585 \pm 0.047710839	N/A	6.268400333 \pm 0.040128533
	ROS	N/A	16.9581856 \pm 0.017257696	N/A	16.9636788 \pm 0.007069011	N/A	17.01458633 \pm 0.020606806
Feature Ablation	Sparseness	N/A	0.5280516 \pm 0.00583783	N/A	0.526218 \pm 0.005664214	N/A	0.529347333 \pm 0.011164491
	RIS	N/A	22.640559 \pm 0.033054016	N/A	23.5692968 \pm 0.057715633	N/A	23.59159233 \pm 0.111776179
	ROS	N/A	33.485311 \pm 0.233803422	N/A	34.5168758 \pm 0.073725495	N/A	34.37202467 \pm 0.080383672
Gradinet-SHAP	Sparseness	N/A	0.195138 \pm 0.026410904	N/A	0.0637004 \pm 0.026472282	N/A	0.113682167 \pm 0.042008189
	RIS	N/A	6.122703167 \pm 0.161879331	N/A	6.0314946 \pm 0.029788567	N/A	6.127377667 \pm 0.120700168
	ROS	N/A	16.93462983 \pm 4.27103E-05	N/A	16.9348698 \pm 7.57377E-05	N/A	16.93466933 \pm 8.86649E-05
Gradient Activation	Sparseness	N/A	0.177242 \pm 0.011243388	N/A	0.1835882 \pm 0.001632398	N/A	0.430289167 \pm 0.00215046
	RIS	N/A	6.123393833 \pm 0.191171312	N/A	6.0269246 \pm 0.030177701	N/A	6.144976167 \pm 0.120931658
	ROS	N/A	16.93457917 \pm 1.47434E-05	N/A	16.9347678 \pm 2.58844E-06	N/A	16.934534 \pm 2.79285E-05
Integrated-Gradient	Sparseness	N/A	0.067058 \pm 0.012083245	N/A	0.0071952 \pm 0.000900684	N/A	0.036059 \pm 0.004866926
	RIS	N/A	6.1224575 \pm 0.150190804	N/A	6.035594 \pm 0.01894045	N/A	6.147797667 \pm 0.09243145
	ROS	N/A	16.93462633 \pm 1.30486E-05	N/A	16.9347694 \pm 1.14018E-06	N/A	16.934618 \pm 8.89944E-06
TEO	Sparseness	N/A	0.416191667 \pm 0.002863111	N/A	0.3715978 \pm 0.000948703	N/A	0.42242125 \pm 0.000173513
	RIS	N/A	5.752004667 \pm 0.364536772	N/A	4.9396128 \pm 0.014829469	N/A	5.5420845 \pm 0.061116734
	ROS	N/A	16.379228 \pm 0.003422144	N/A	15.8121004 \pm 0.00994492	N/A	16.277279 \pm 0.001043319
TEO-UMAP	Sparseness	N/A	0.423819167 \pm 0.00056124	N/A	0.423865 \pm 5.01946E-05	N/A	0.424567429 \pm 0.000241638
	RIS	N/A	5.552539167 \pm 0.186236877	N/A	5.458319 \pm 0.029725596	N/A	5.557568429 \pm 0.095350875
	ROS	N/A	16.3661035 \pm 0.00731861	N/A	16.372555 \pm 0.001689458	N/A	16.357192 \pm 0.003475647

and FDR-adjusted inference in the same aggregate framework; consequently, we do not report hypothesis tests for these groups here to avoid bias from unmatched rows. Overall, across the three groups, the most robust and reproducible effect we could quantify was the *reduction in ROS under SAE* (clearly demonstrated in MCI with converged paired comparisons), while RIS showed a smaller SAE-related decrease by non-parametric testing and Complexity tended to increase under SAE for MCI.

B.6 INDIVIDUAL-LEVEL EXPLANATIONS AND PATTERNS

Figures 7–16 present qualitative local attribution examples for the binary (Control and Alzheimer) and three-class classification task (Control, LMCI, MCI) of ADNI cohorts across six explanation methods, each evaluated without (Figures 7, 9, 11, 13, 15) and with (Figures 8, 10, 12, 14, 16) the Sparse Autoencoder (SAE) layer. Each cell shows token-level attributions using colour-coded highlights (green = positive relevance; red = negative relevance). In general, higher Sparseness is associated with a more balanced distribution of positive and negative highlights (i.e., less diffuse maps), particularly for Layer Conduction, Feature Ablation, Gradient SHAP, and Integrated Gradient. For the Control class (see Figures 7 and 8), the qualitative highlighting patterns are broadly consistent across the six attribution techniques, Activation, Layer Conduction, Feature Ablation, Gradient SHAP, Gradient Activation, and Integrated Gradient—with no marked visual discrepancies. Notably, Feature Ablation, despite exhibiting the strongest Sparseness in the box plots, shows poorer stability (higher variability in inputs/outputs; elevated RIS/ROS), and the addition of the SAE layer tends to worsen this by exposing a larger set of features due to the decoder “decompression” effect; a similar trend is observed for Activation. For the Alzheimer’s class (Figures 9 and 10), Layer Conduction demonstrates a reduction in Sparseness with the SAE but a gain in stability (decreased RIS/ROS). Comparable improvements in stability with SAE are also observed for Gradient Activation, Integrated Gradient, and Gradient SHAP. In contrast, Activation and Feature Ablation perform worst under SAE, again exposing many more features and yielding less stable explanations. Across the remaining examples (Figures 11–17), similar patterns hold: instances with low Sparseness and high RIS/ROS tend to produce saturated red/green

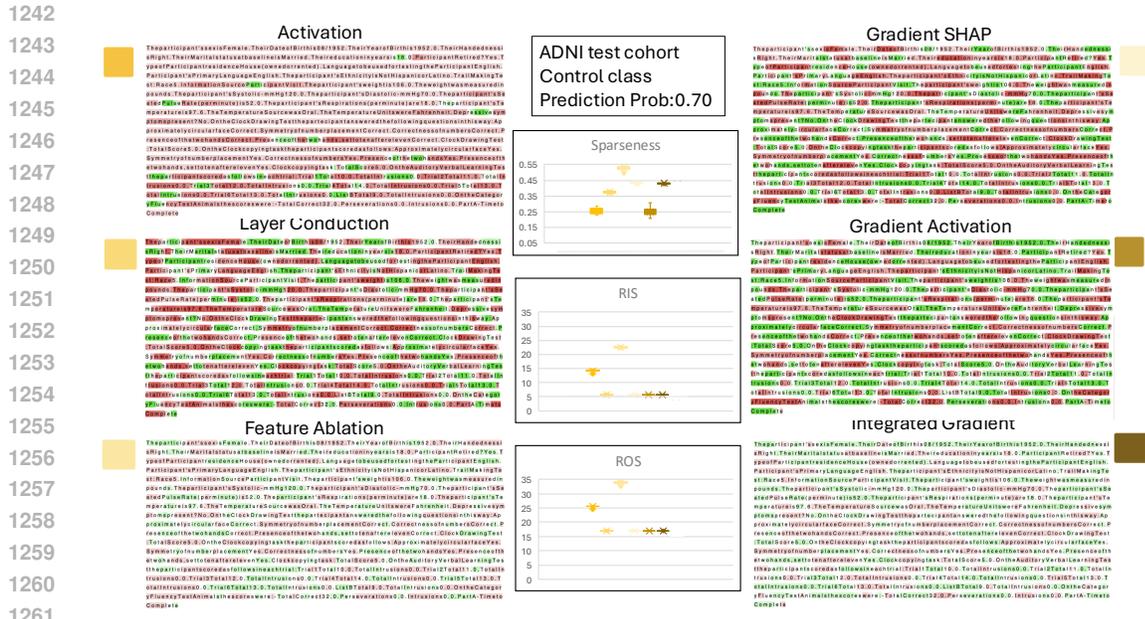


Figure 7: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a binary classification (Alzheimer’s disease vs Control) on the ADNI cohort; the examples shown here are from the Control class.

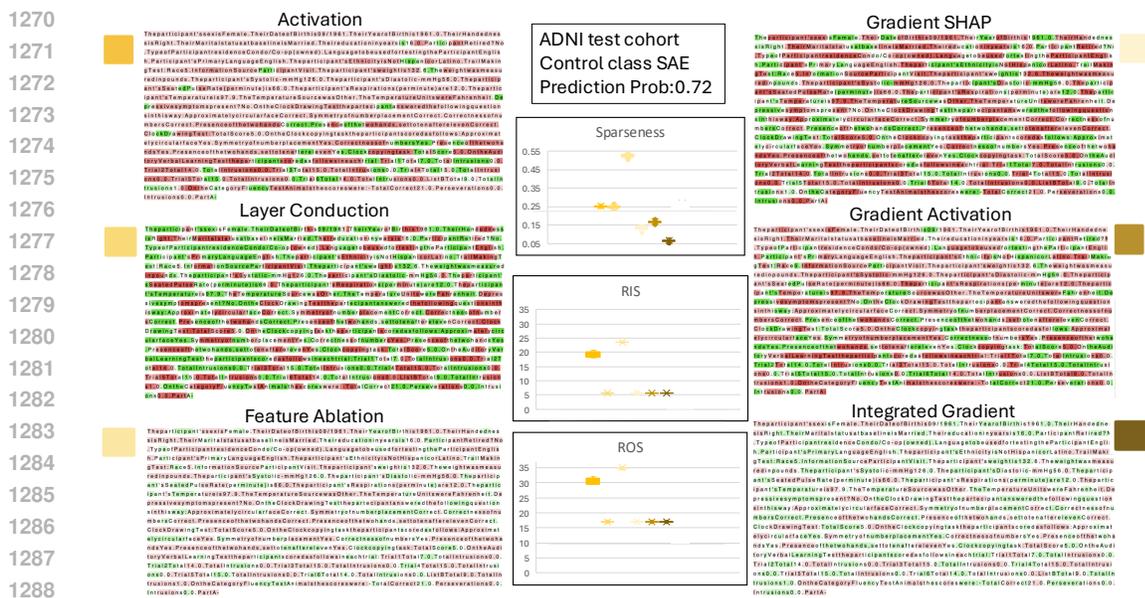


Figure 8: Local attribution examples across different explanation methods with the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a binary classification (Alzheimer’s disease vs Control) on the ADNI cohort; the examples shown here are from the Control class.

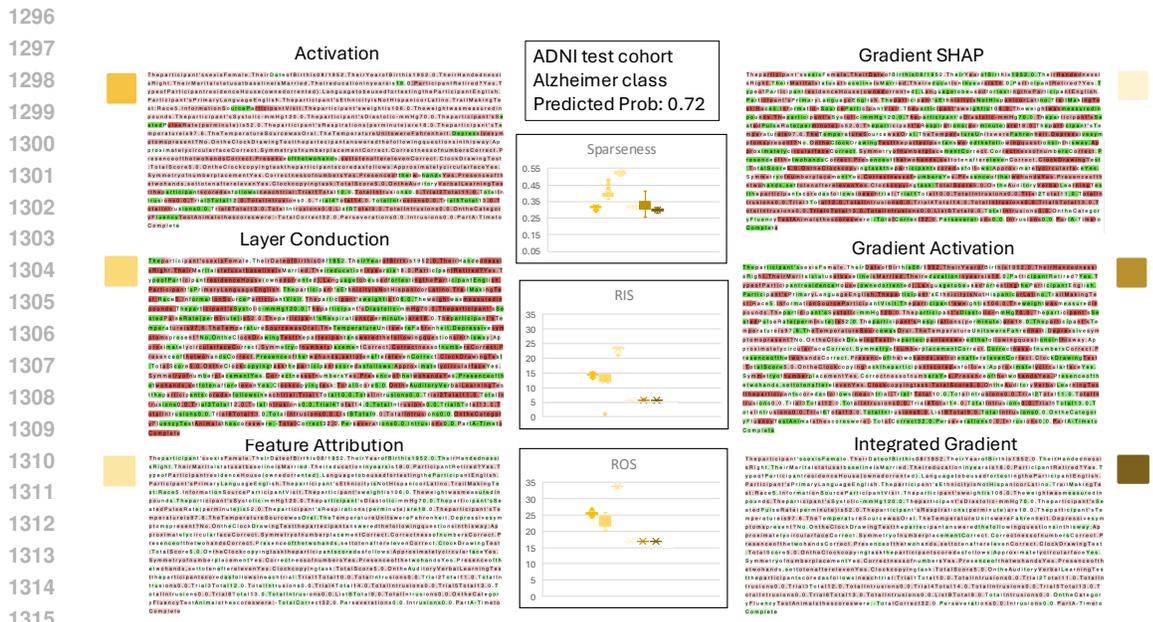


Figure 9: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a binary classification (Alzheimer’s disease vs Control) on the ADNI cohort; the examples shown here are from the Alzheimer class.

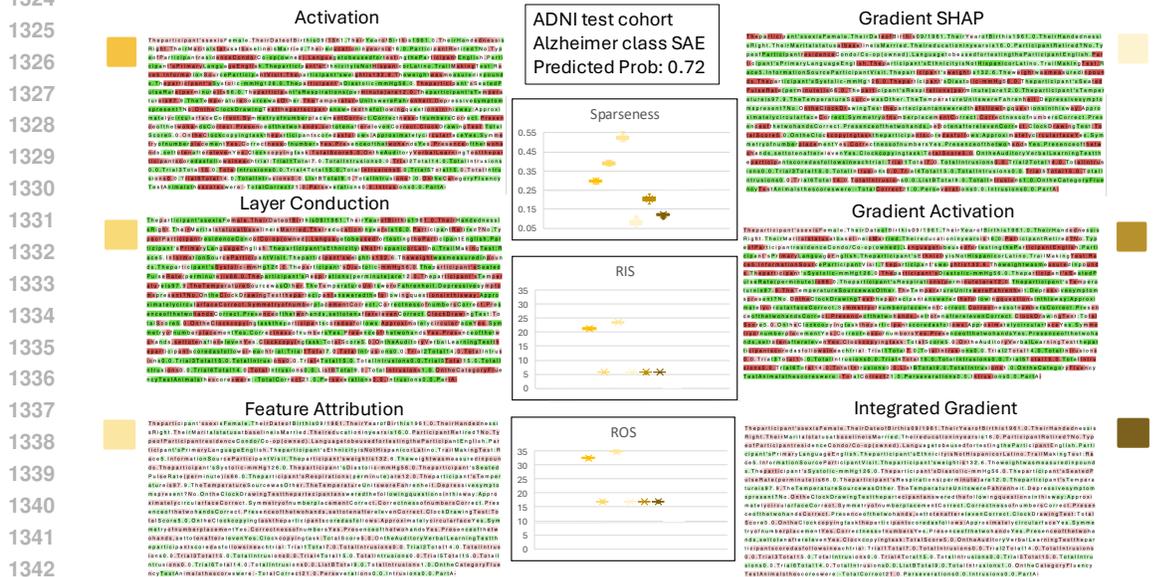


Figure 10: Local attribution examples across different explanation methods with the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a binary classification (Alzheimer’s disease vs Control) on the ADNI cohort; the examples shown here are from the Alzheimer class.

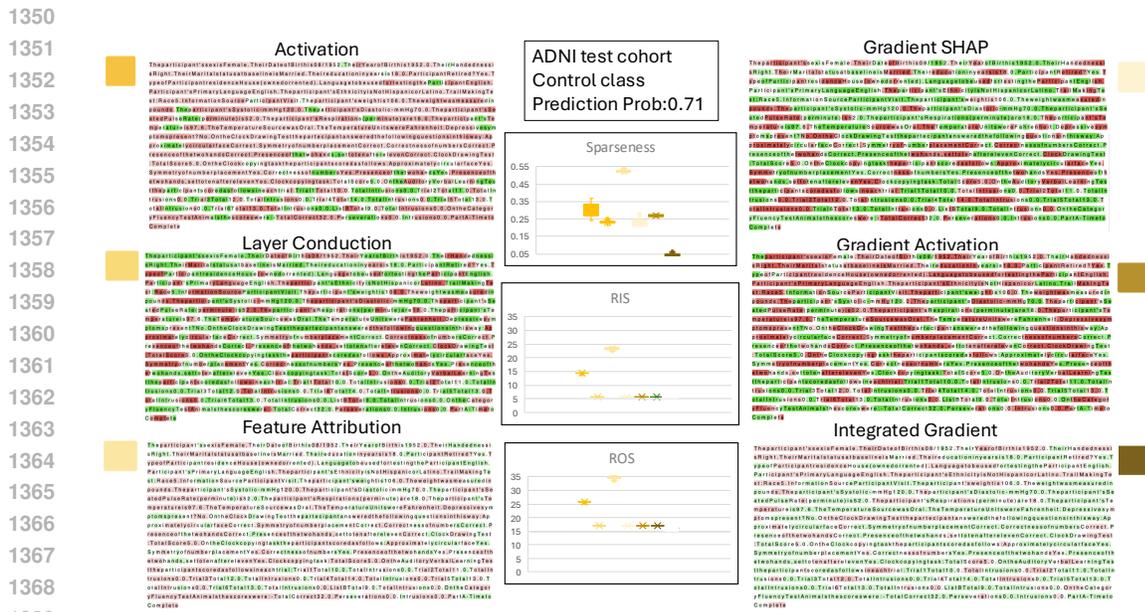


Figure 11: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the Control class.

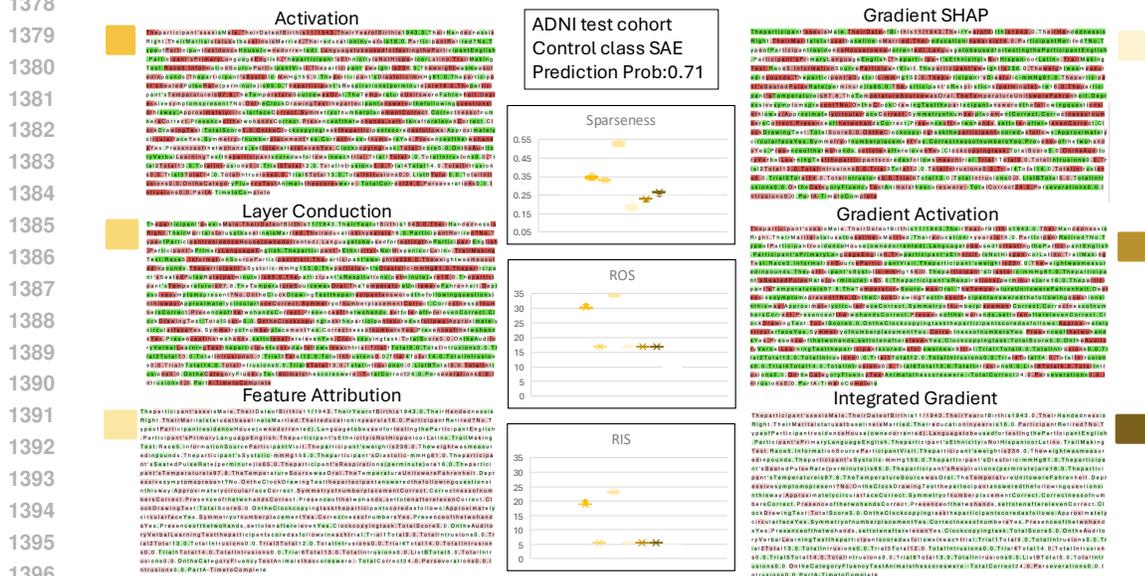


Figure 12: Local attribution examples across different explanation methods with the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the Control class.

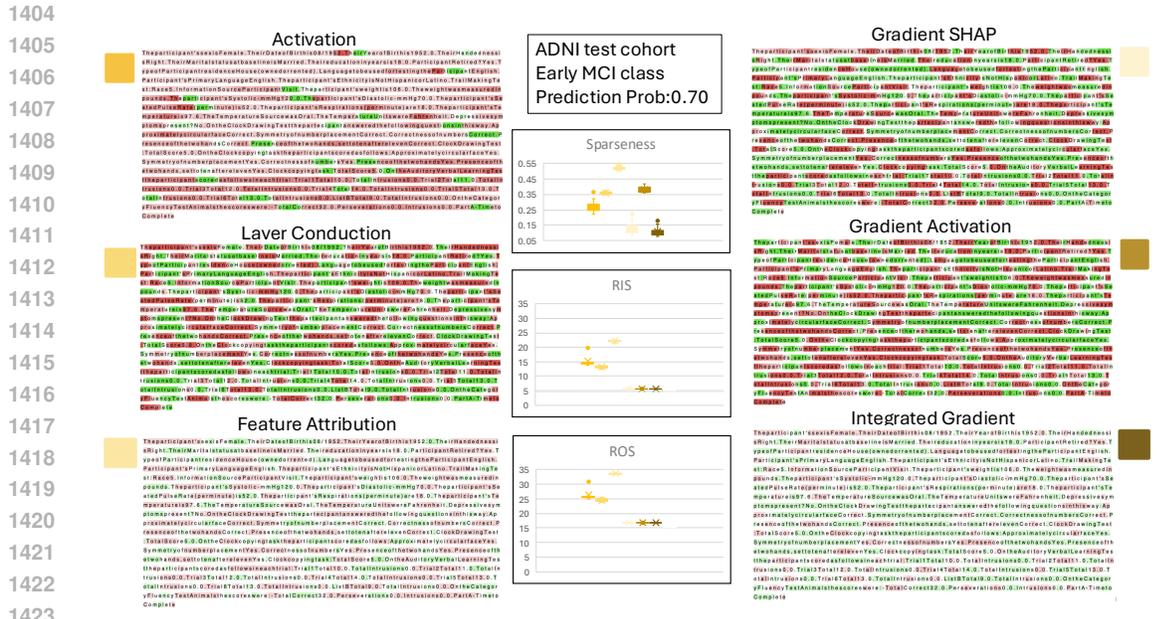


Figure 13: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the LMCI class.

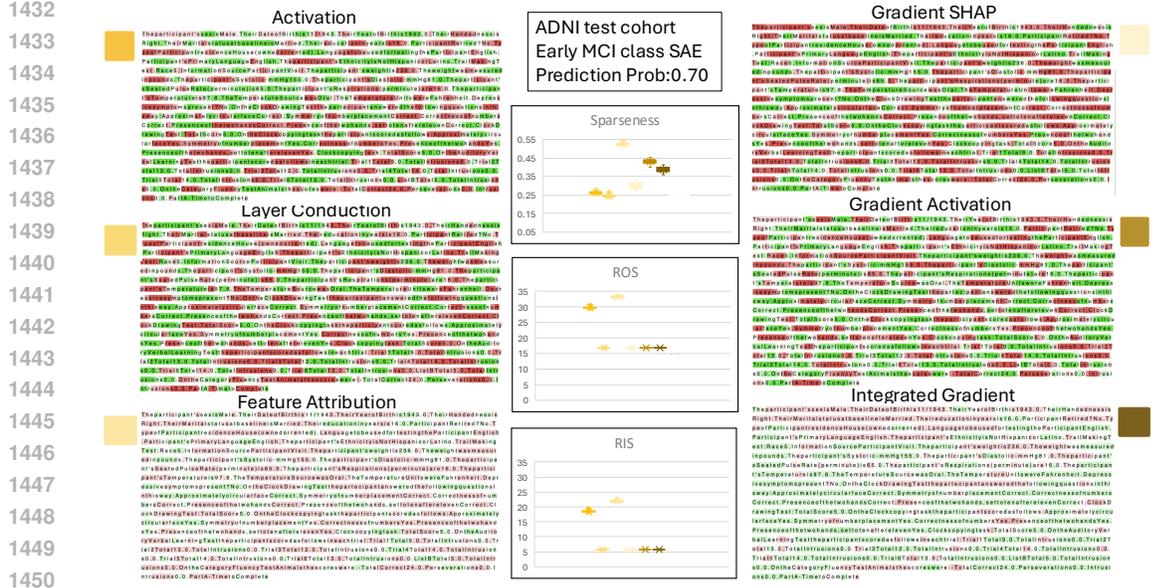


Figure 14: Local attribution examples across different explanation methods with the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the LMCI class.

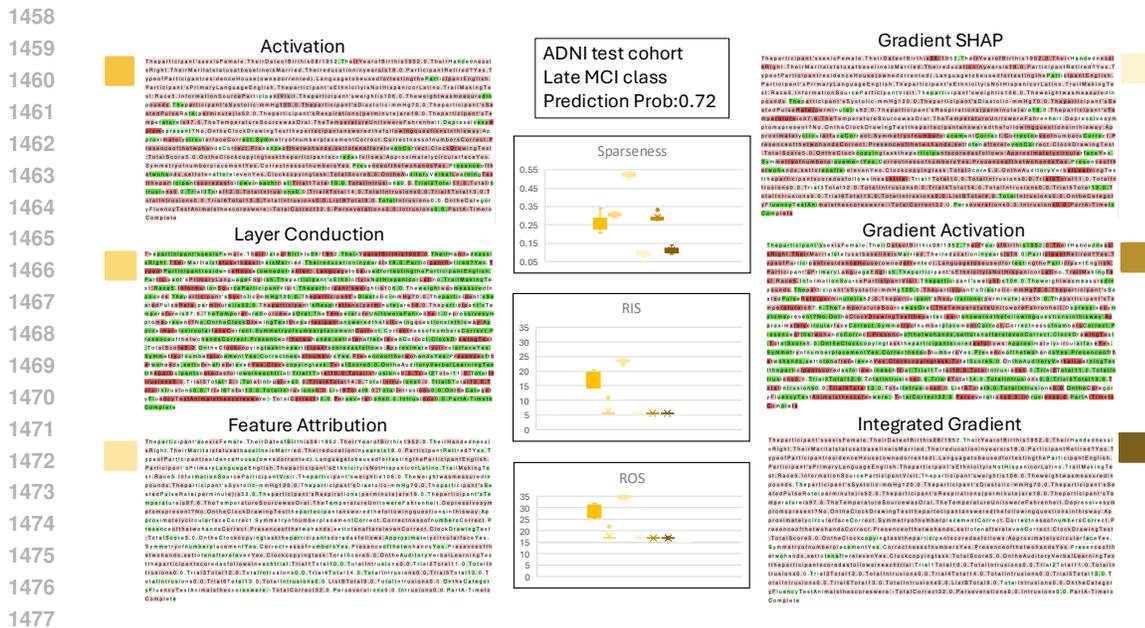


Figure 15: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the MCI class.

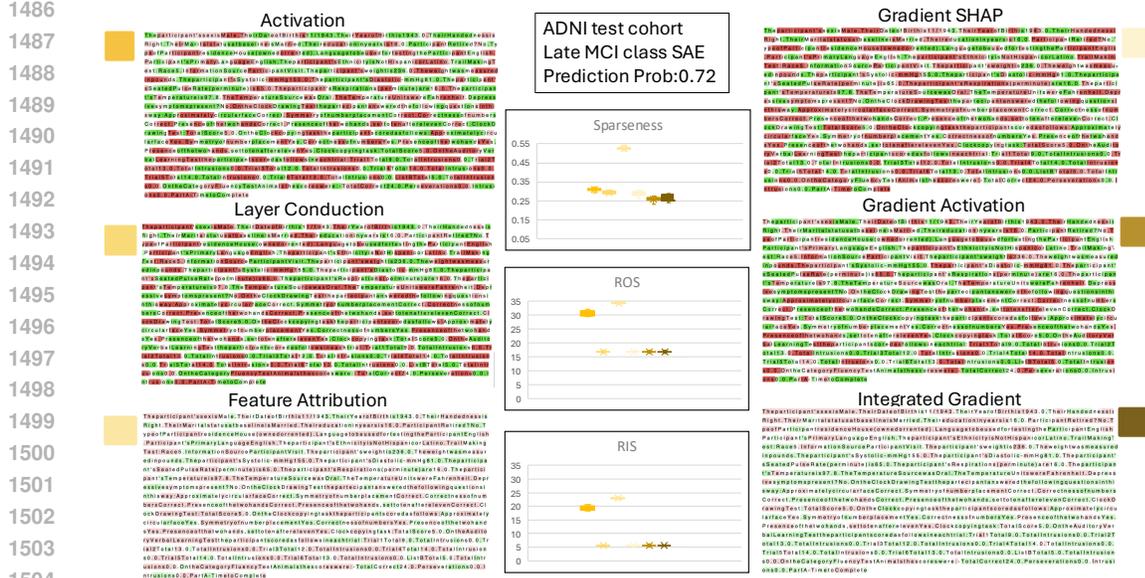


Figure 16: Local attribution examples across different explanation methods with the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to +1 (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) in the ADNI cohort; the examples shown here are from the MCI class.

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

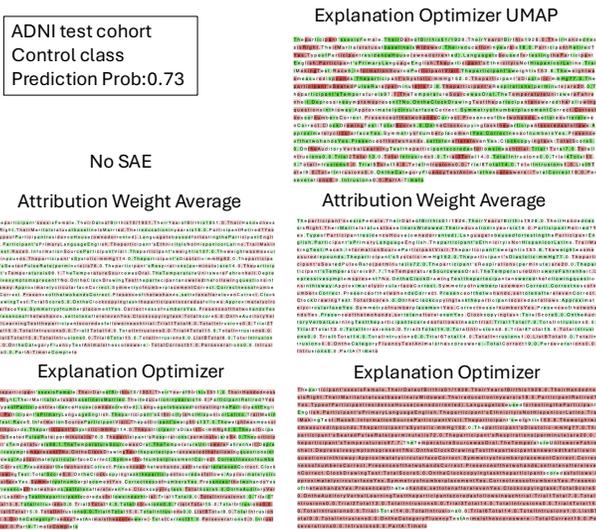


Figure 17: Local attribution examples across different explanation methods with the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to $+1$ (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a binary classification (Alzheimer’s disease vs Control) on the ADNI cohort; the examples shown here are from the Control class.

explanations (strongly negative or positive attributions), whereas higher Sparseness with lower RIS/ROS yields more compact and stable saliency patterns.

Figures 17–21 present qualitative local attribution examples, analogous to Figures 7–17, for the no-SAE analyses of (i) the attributional weighted average (computed from the six base methods), (ii) the Transformer Explanation Optimizer (TEO), and (iii) TEO with a linear UMAP constraint (TEO-UMAP). As shown in the previous subsection, with the SAE layer TEO achieves the best stability—i.e., the lowest RIS and ROS—but at the cost of a marked reduction in Sparseness; this reduction is clearly visible in the binary task (Figures 17–18). Introducing the UMAP constraint yields a more balanced trade-off, producing explanations that are more compact and clinically interpretable; the same behaviour is observed across all classes in the three-class setting (Figures 19–21). By contrast, the weighted-average approach—a linear combination of the six attribution techniques—does not yield superior explanations, consistent with Mamalakis et al. (2025).

B.7 UMAP AND COHORT-LEVEL EXPLANATION AND PATTERNS.

Figures 22–31 present cohort-level attribution examples for both the binary (Control vs Alzheimer’s disease) and three-class (Control, LMCI, MCI) classification tasks on the ADNI test cohort across six explanation methods. Each method is shown without (Figures 22, 24, 26, 28, 30) and with (Figures 23, 25, 27, 29, 31) the Sparse Autoencoder (SAE) layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the 1D panel shows attribution scores along PCA-first component. All plotted values are normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI subgroups (see §B3). In general, moving from the no-SAE to the SAE condition broadens the distribution of features in 2D and increases the density of high-significance points (upper-right boxed region), consistent with a decoder-induced decompression effect and a corresponding reduction in sparsity in the attribution maps.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585



Figure 18: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to $+1$ (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a binary classification (Alzheimer’s disease vs Control) on the ADNI cohort; the examples shown here are from the Alzheimer class.

1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

Figure 19: Local attribution examples across different explanation methods with the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to $+1$ (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the Control class.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

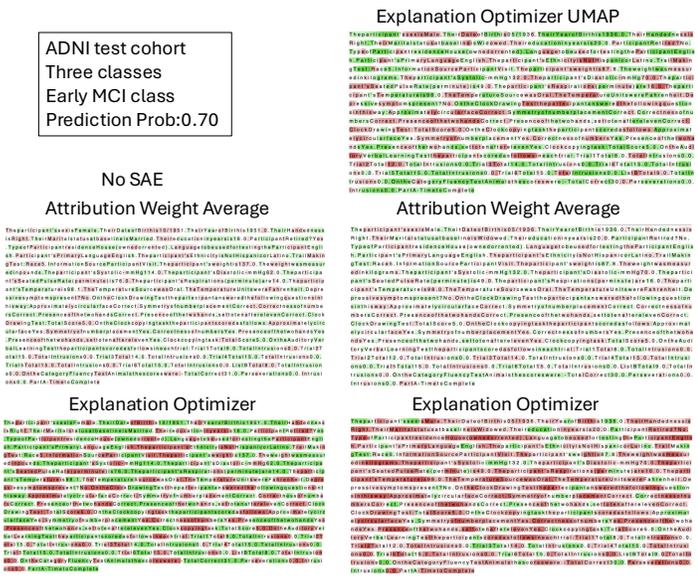


Figure 20: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to $+1$ (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the LMCI class.

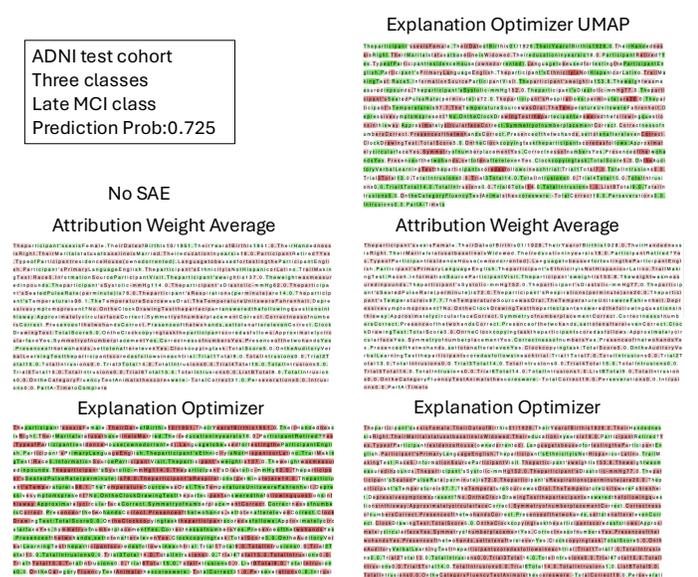


Figure 21: Local attribution examples across different explanation methods without the SAE layer. The colour scale ranges from -1 (dark red; negative attribution), through 0 (white; neutral), to $+1$ (dark green; positive attribution). For each of the six panels, the small colour swatches at the top-left and top-right indicate the colour keys used for the three summary box-plot metrics—Sparseness, RIS, and ROS—for the corresponding attribution technique. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the MCI class.

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

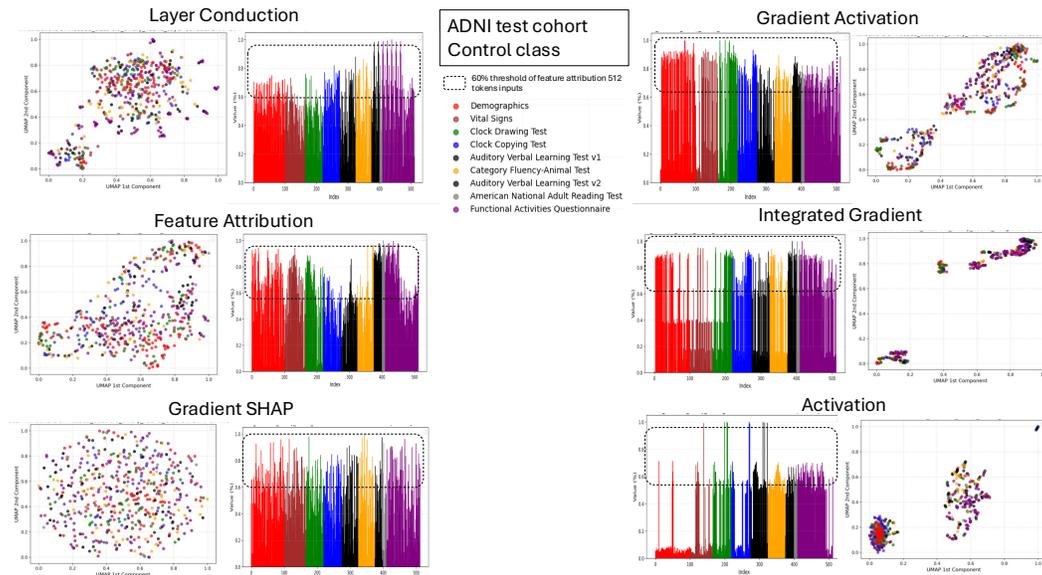
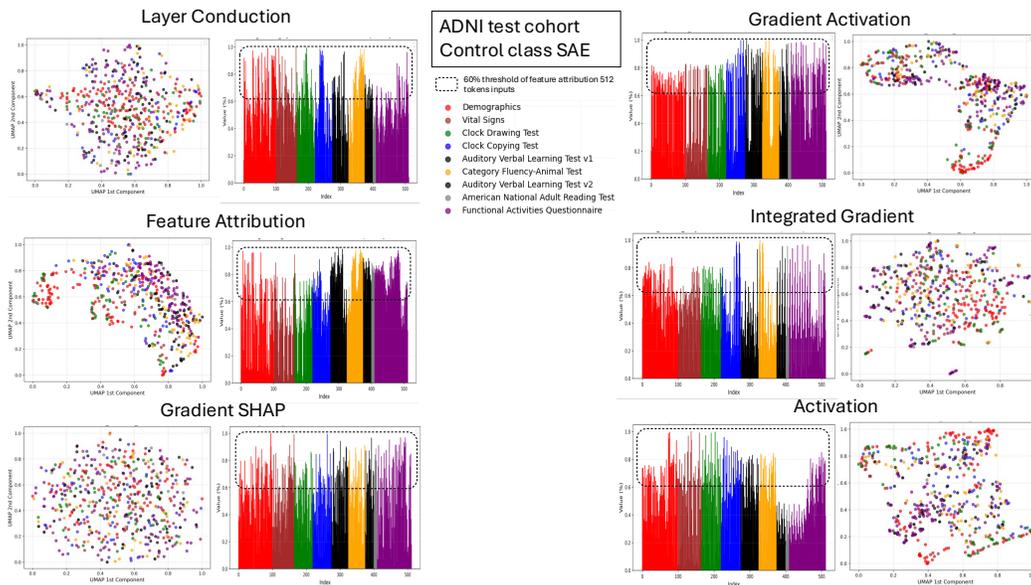


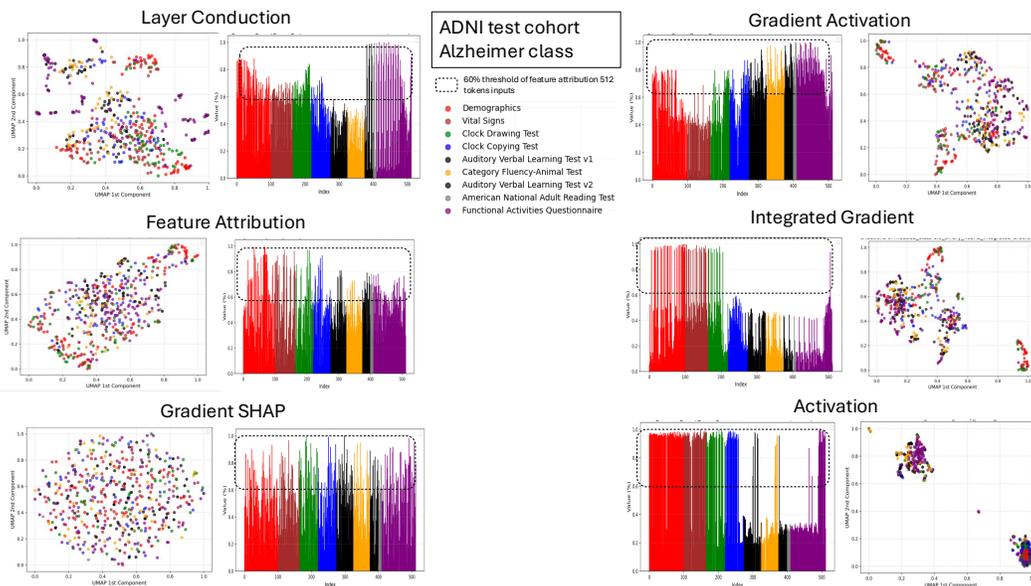
Figure 22: Global (cohort-level) feature attribution across explanation methods without the SAE layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the 1D panel shows attribution scores along PCA first component. All plotted values are normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant tokens in the 1D views. The task is binary classification (Alzheimer’s vs Control) on the ADNI cohort; the examples shown here are from the Control class. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the MCI class.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747



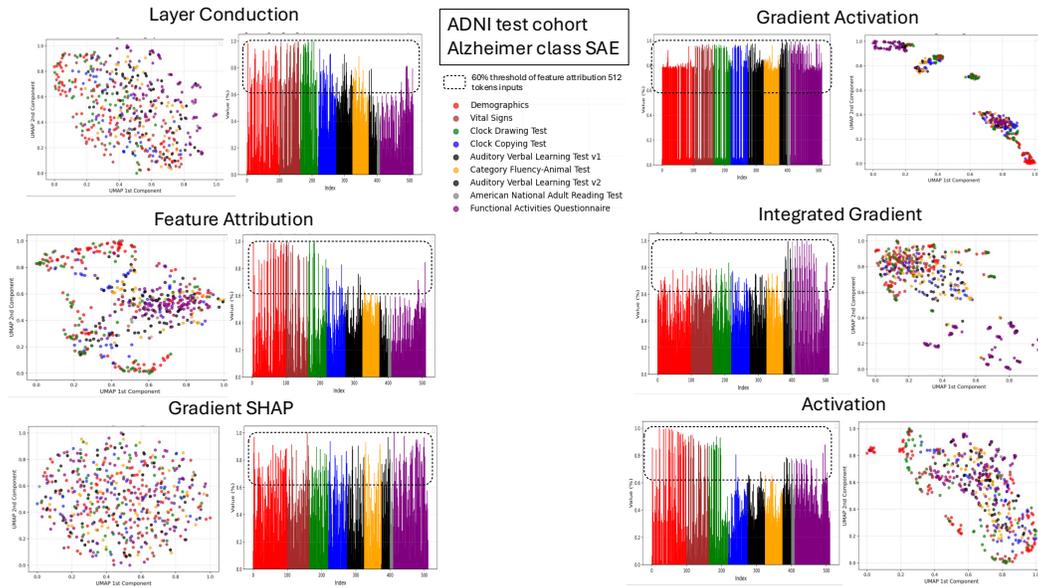
1748 Figure 23: Global (cohort-level) feature attribution across explanation methods with the SAE layer.
1749 The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the
1750 1D panel shows attribution scores along PCA first component. All plotted values are normalised to
1751 [0, 1] and represent positive contributions only. Colours (red→purple) denote the nine ADNI
1752 subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant
1753 tokens in the 1D views. The task is binary classification (Alzheimer’s vs Control) on the ADNI
1754 cohort; the examples shown here are from the Control class.

1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774



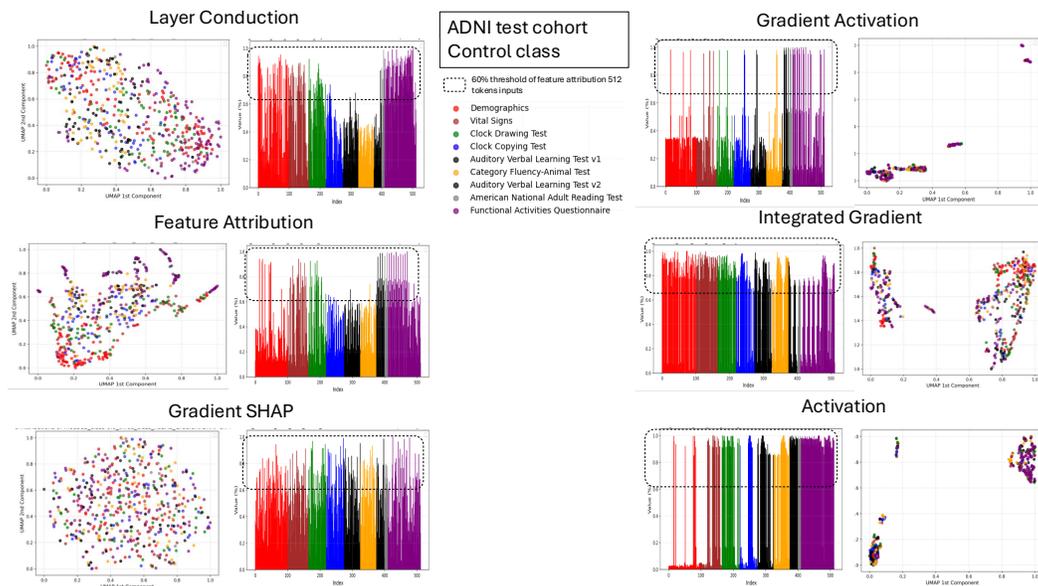
1775 Figure 24: Global (cohort-level) feature attribution across explanation methods without the SAE
1776 layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI
1777 test set; the 1D panel shows attribution scores along PCA first component. All plotted values are
1778 normalised to [0, 1] and represent positive contributions only. Colours (red→purple) denote the
1779 nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most
1780 significant tokens in the 1D views. The task is binary classification (Alzheimer’s vs Control) on the
1781 ADNI cohort; the examples shown here are from the Alzheimer class.

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801



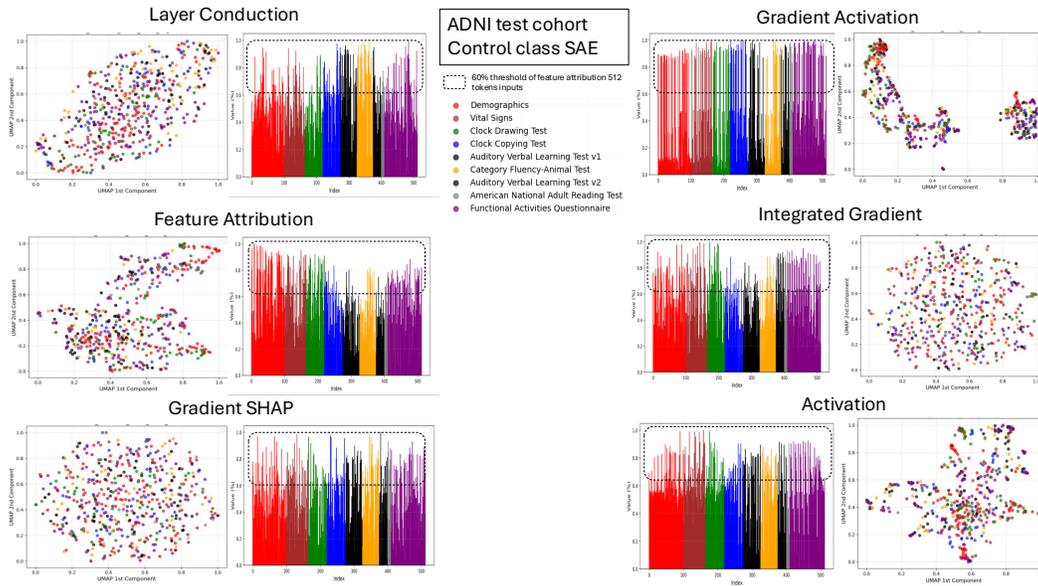
1802 Figure 25: Global (cohort-level) feature attribution across explanation methods with the SAE layer.
1803 The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set
1804 the 1D panel shows attribution scores along PCA first component. All plotted values are normalised
1805 to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI
1806 subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant
1807 tokens in the 1D views. The task is binary classification (Alzheimer’s vs Control) on the ADNI
1808 cohort; the examples shown here are from the Alzheimer class.

1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828



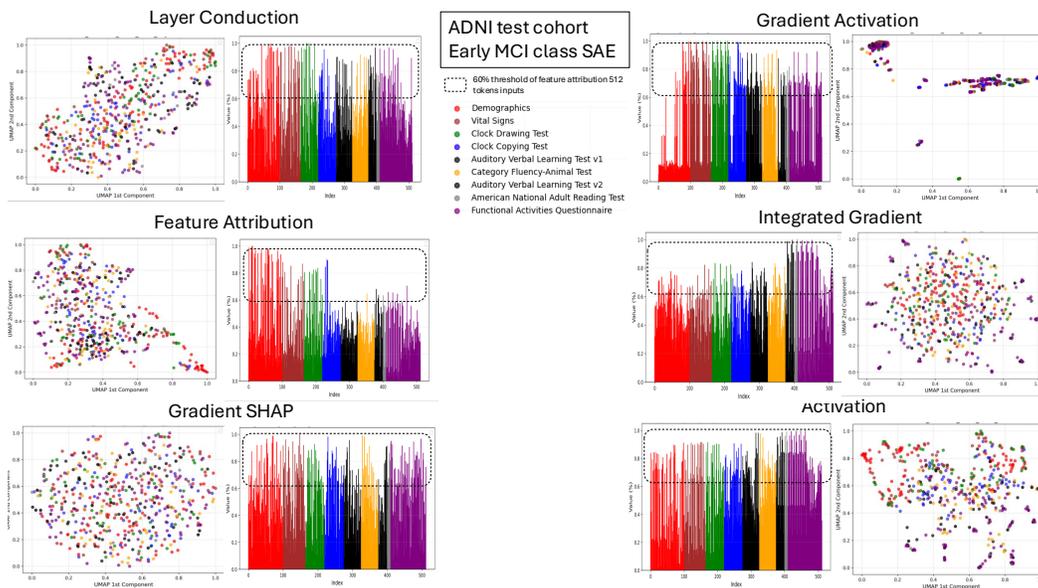
1829 Figure 26: Global (cohort-level) feature attribution across explanation methods without the SAE
1830 layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI
1831 test set the 1D panel shows attribution scores along PCA first component. All plotted values are
1832 normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the
1833 nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most
1834 significant tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs
1835 Control) on the ADNI cohort; the examples shown here are from the Control class.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855



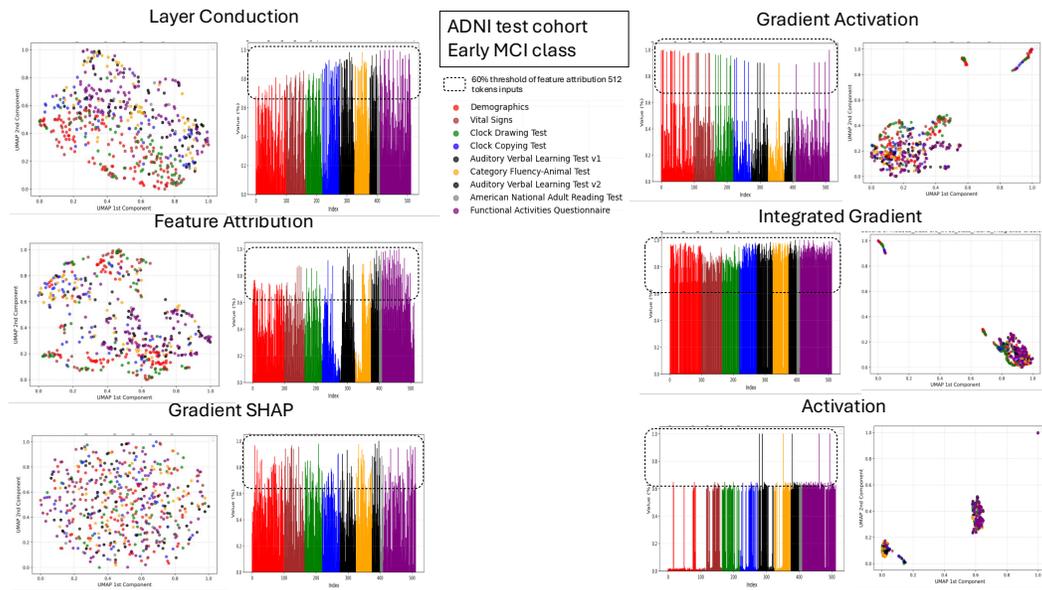
1856 Figure 27: Global (cohort-level) feature attribution across explanation methods with the SAE layer.
1857 The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set
1858 the LMCI class SAE. The 1D panel shows attribution scores along PCA first component. All plotted values are normalised
1859 to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI
1860 subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant
1861 tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs Control) on
1862 the ADNI cohort; the examples shown here are from the Control class.

1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889



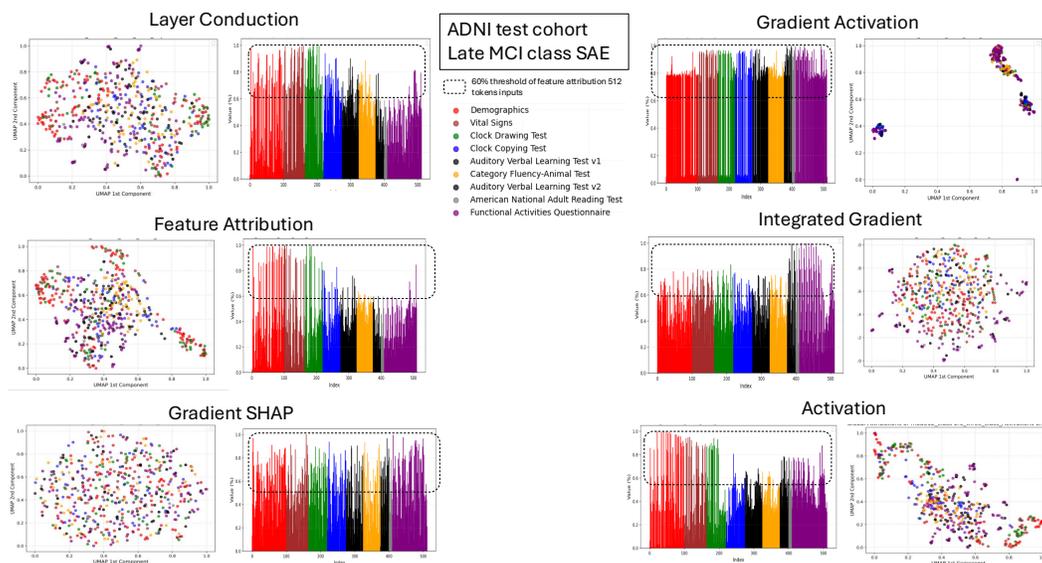
1884 Figure 28: Global (cohort-level) feature attribution across explanation methods with the SAE layer.
1885 The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the
1886 LMCI class SAE. The 1D panel shows attribution scores along PCA first component. All plotted values are normalised
1887 to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI
1888 subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant
1889 tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs Control) on
the ADNI cohort; the examples shown here are from the LMCI class.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909



1910 Figure 29: Global (cohort-level) feature attribution across explanation methods without the SAE
1911 layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI
1912 test set; the 1D panel shows attribution scores along PCA first component. All plotted values are
1913 normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the
1914 nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most
1915 significant tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs
1916 Control) on the ADNI cohort; the examples shown here are from the LMCI class.

1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943



1937 Figure 30: Global (cohort-level) feature attribution across explanation methods with the SAE
1938 layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI
1939 test set; the 1D panel shows attribution scores along PCA first component. All plotted values are
1940 normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the
1941 nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most
1942 significant tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs
1943 Control) on the ADNI cohort; the examples shown here are from the MCI class.

1944
 1945
 1946
 1947
 1948
 1949
 1950
 1951
 1952
 1953
 1954
 1955
 1956
 1957
 1958
 1959
 1960
 1961
 1962
 1963
 1964
 1965
 1966
 1967
 1968
 1969
 1970
 1971
 1972
 1973
 1974
 1975
 1976
 1977
 1978
 1979
 1980
 1981
 1982
 1983
 1984
 1985
 1986
 1987
 1988
 1989
 1990
 1991
 1992
 1993
 1994
 1995
 1996
 1997

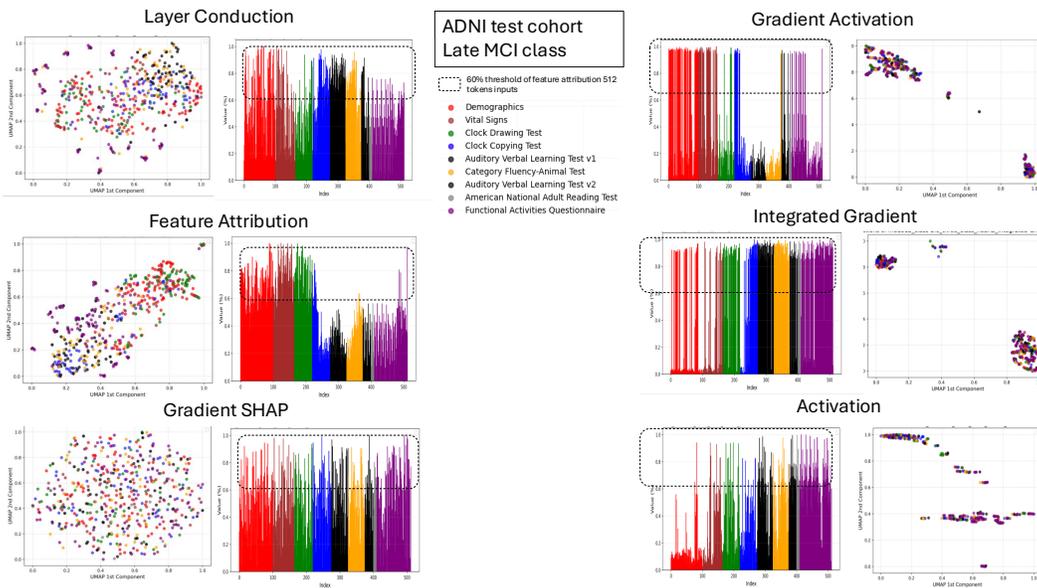


Figure 31: Global (cohort-level) feature attribution across explanation methods without the SAE layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the 1D panel shows attribution scores along PCA first component. All plotted values are normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the MCI class.

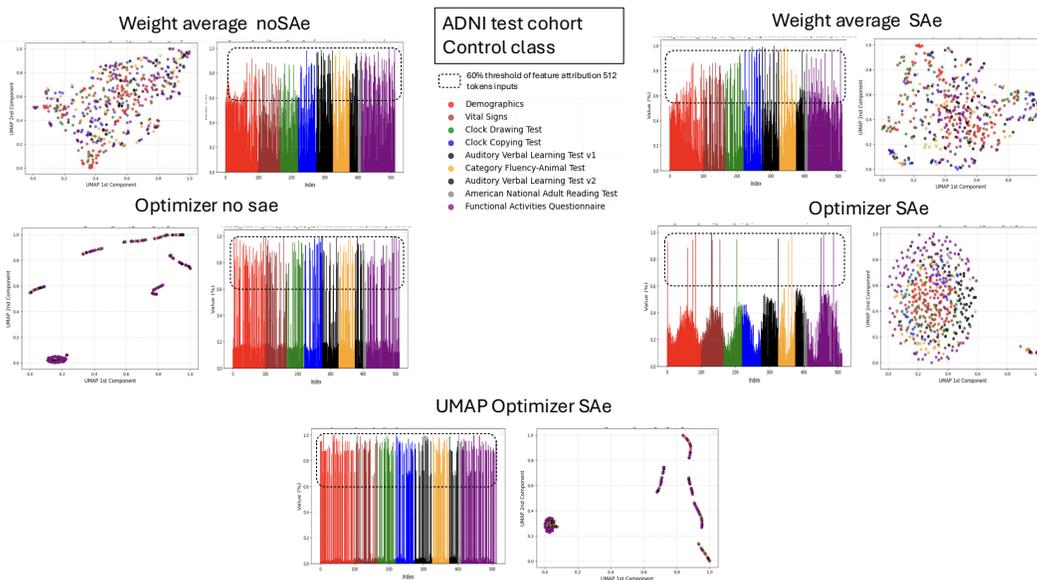
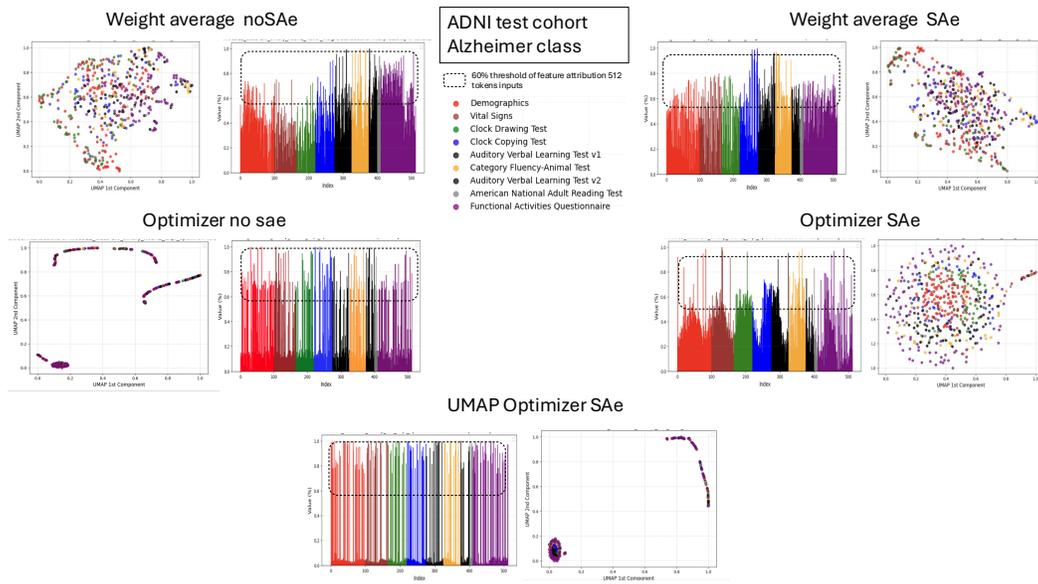


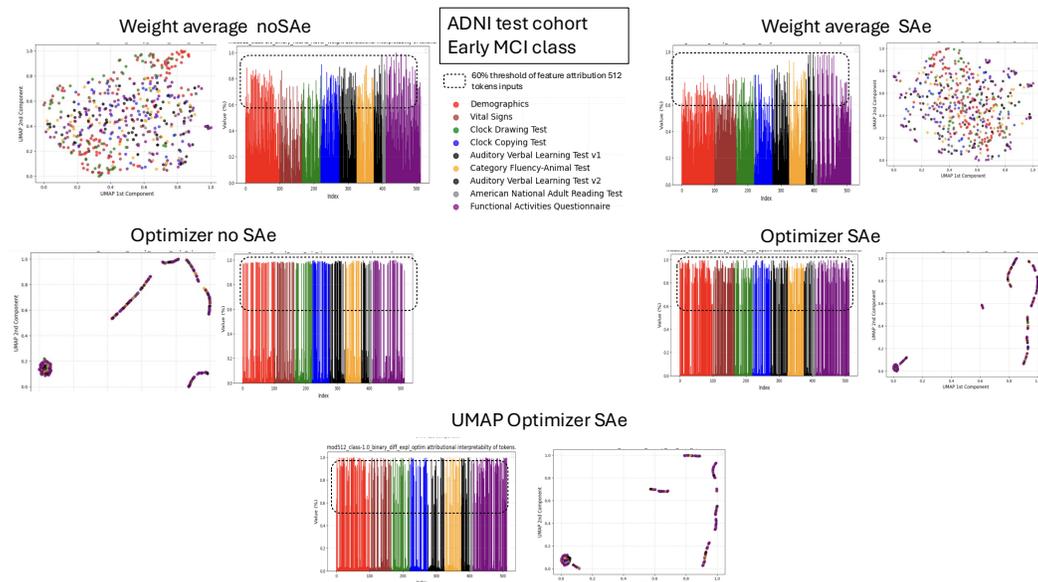
Figure 32: Global (cohort-level) feature attribution across explanation methods with the SAE layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the 1D panel shows attribution scores along PCA first component. All plotted values are normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant tokens in the 1D views. The task is a binary classification (Alzheimer vs Control) on the ADNI cohort; the examples shown here are from the Control class.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017



2018 Figure 33: Global (cohort-level) feature attribution across explanation methods with the SAE layer.
2019 The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the
2020 1D panel shows attribution scores along PCA first component. All plotted values are normalised to
2021 $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI
2022 subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant
2023 tokens in the 1D views. The task is a binary classification (Alzheimer vs Control) on the ADNI
2024 cohort; the examples shown here are from the Alzheimer class.

2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044



2045 Figure 34: Global (cohort-level) feature attribution across explanation methods without the SAE
2046 layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI
2047 test set; the 1D panel shows attribution scores along PCA first component. All plotted values are
2048 normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the
2049 nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most
2050 significant tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs
2051 Control) on the ADNI cohort; the examples shown here are from the LMCI class.

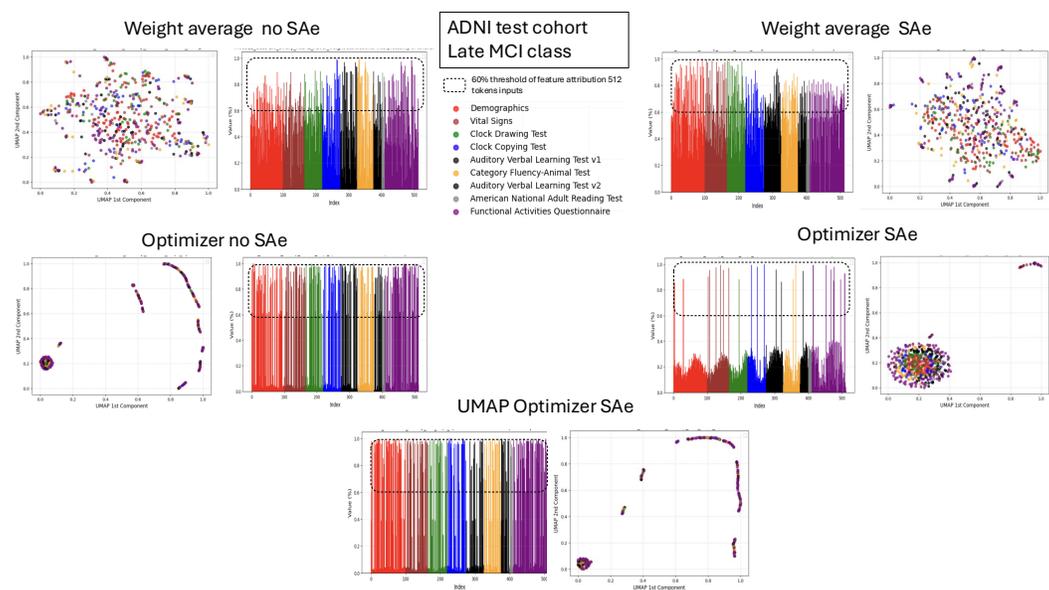


Figure 35: Global (cohort-level) feature attribution across explanation methods with the SAE layer. The 2D panel shows a UMAP embedding (UMAP-1 vs UMAP-2) computed on the ADNI test set; the 1D panel shows attribution scores along PCA first component. All plotted values are normalised to $[0, 1]$ and represent positive contributions only. Colours (red→purple) denote the nine ADNI subgroups (see §B3). Square boxes mark the 0.6–1.0 interval, highlighting the most significant tokens in the 1D views. The task is a three-class classification (LMCI, MCI disease vs Control) on the ADNI cohort; the examples shown here are from the MCI class.

Figures 32–35 present cohort-level attribution examples for both the binary (Control vs Alzheimer’s disease) and three-class (Control, LMCI, MCI) classification tasks on the ADNI test cohort, analogous to Figures 22–31, for the no-SAE analyses of (i) the attributional weighted average (computed from the six base methods), (ii) the Transformer Explanation Optimizer (TEO), and (iii) TEO with a linear UMAP constraint (UMAP Optimizer). As shown in the previous subsection, with the SAE layer TEO achieves the best stability—i.e., the lowest RIS and ROS—but at the cost of a marked reduction in Sparseness; this reduction is clearly visible in the binary task (Figures 32–35), where a spreading of tokens in 2D is observed when moving from no-SAE to SAE, as with the other methods. TEO with SAE reorganises the space, yielding a more homogeneous low-to-high attribution gradient. The drawback is that, without appropriate guidance, there may be too few features in the squares denoting significant contribution, and not all subgroups in the global observations are represented (e.g., Figure 32). However, this can be mitigated by constraining the 2D manifold in the attribution space. To that end, we proposed a linear constraint to further smooth the regrouping of tokens in the attribution manifold. Introducing the UMAP linear constraint yields an even more balanced trade-off compared with unconstrained TEO with SAE, producing explanations that share similar significant traits across the different subgroups (colours) and are more homogeneous (very clear in Figures 32, 33, and 35, less so in 34). Consequently, the maps are more compact and clinically interpretable; the same behaviour is observed across all classes in the three-class setting (Figures 33–35). By contrast, at both cohort and local levels, the weighted-average approach—a linear combination of the six attribution techniques—does not yield superior explanations, consistent with Mamalakis et al. (2025).

B.8 CLINICAL RELEVANCE IN ALZHEIMER’S DIAGNOSIS PROGRESSION: EVIDENCE THAT SAE-GUIDED ATTRIBUTION YIELDS MORE RELIABLE EXPLANATIONS THAN TRADITIONAL ATTRIBUTION.

To further validate the role of the SAE layer in shaping the attribution space into a more monosemantic and clinically coherent feature representation, we conducted an auxiliary evaluation. Specifically, for each test input, we extracted the top 50% most influential token attributions

2106 produced by our attribution framework—TEO without SAE, TEO with SAE, and TEO-UMAP. We
2107 then generated a CSV file for each classification class, in which the highlighted characters for
2108 the three attribution methods were organized column-wise. The complete character sequence
2109 for each sample, beginning with the CLS token, was included in the first column to ensure clear
2110 sample-level distinction. These CSV files were subsequently provided as input to a large language
2111 model (ChatGPT-5.1 OpenAI (2024)) using a fixed prompt to obtain an external, model-agnostic
2112 assessment of the interpretability structure encoded by each explanation space.

2113 Three experiments were performed:

- 2114 1. **Binary ADNI** (Control vs. Alzheimer’s disease), with each class provided as a separate
2115 CSV file.
- 2116 2. **Binary BrainLAT** (Control vs. Alzheimer’s disease), also split into two class-specific CSV
2117 files.
- 2118 3. **Three-class ADNI** (Control, MCI, LMCI), with each diagnostic category represented in its
2119 own CSV file.

2122 We evaluated two primary criteria: (i) whether the language model could distinguish, based solely
2123 on the highlighted features, which CSV corresponded to the pathological versus the healthy control
2124 class; and (ii) whether the model could identify meaningful pathology-related biomarkers.

2125 For the first two experiments, the model was prompted with:

2126
2127 *Given the two CSV files, and recognizing that medical biases exist in the `char`
2128 column with each sample beginning with the character sequence [CLS], determine
2129 how each of the three attribution methods (`attr1`, `attr2`, `attr3`) highlight features
2130 associated with healthy or unhealthy interpretations, and analyze the reasons
2131 for these differences. Predict the pathology and specify which of the two CSV files
2132 corresponds to the pathological case for each attribution technique.*

2133 For the three-class ADNI experiment, we used:

2134
2135 *Given the three CSV files, and recognizing that medical biases exist in the `char`
2136 column with each sample beginning with the character sequence [CLS], determine
2137 how each attribution method (`attr1`, `attr2`, `attr3`) highlights features associated
2138 with healthy or pathological interpretations. Predict the pathology or pathologies,
2139 identify which CSVs correspond to the pathological and healthy groups for each
2140 attribution method, and specify the associated conditions.*

2141
2142 The resulting GPT-generated interpretations, shown in Figures 36, provide an external linguistic
2143 lens on each explanation space.

2144 In the binary ADNI experiment, all the three attribution frameworks, TEO without SAE, TEO
2145 with SAE, and TEO-UMAP, correctly identified the pathological file (CSV1) and associated it
2146 with Alzheimer’s disease, noting that the signal was more consistent with late-stage rather than
2147 early cognitive impairment. A similar pattern emerged for the binary BrainLAT experiment: the
2148 pathological CSV was attributed to Alzheimer’s disease rather than MCI, with clear differentiation
2149 from the healthy control (Figures 36a,b).

2150 Across both binary tasks, TEO without SAE exhibited erratic and clinically uninformative behavior,
2151 frequently attending to task labels, instruction counts, or other artefactual patterns rather than
2152 neurocognitive biomarkers (Figures 36a,b). In contrast, TEO-SAE and TEO-UMAP consistently
2153 highlighted clinically meaningful domains, including demographic risk factors, processing-speed
2154 impairments, and neurophysiological indicators.

2155 In the more complex three-class ADNI experiment, the advantages of the SAE-induced monose-
2156 mantic structure became even clearer (Figure 36c). TEO without SAE failed to correctly identify
2157 pathological classes and did not surface meaningful biomarkers. Conversely, TEO-SAE achieved
2158 clearer diagnostic separation and more coherent feature concentrations, while TEO-UMAP fur-
2159 ther emphasized structured biomarkers, particularly within demographics and vitals, and also
provided correct class-level predictions.

Group	Dem	VS	CDT	CCT	AVLT1	CFA	AVLT2	ANART	FAQ
Control TEO	0.00	0.00	0.00	0.00	0.02	0.00	0.00	0.00	0.11
Control TEO-UMAP	0.33	0.21	0.29	0.29	0.12	0.30	0.32	0.30	0.36
Alzheimer TEO	0.05	0.00	0.09	0.01	0.12	0.12	0.00	0.00	0.13
Alzheimer TEO-UMAP	0.33	0.29	0.30	0.32	0.22	0.30	0.32	0.50	0.35
Control TEO	0.74	0.55	0.87	0.63	0.88	0.68	0.20	0.80	0.38
Control TEO-UMAP	0.51	0.61	0.67	0.60	0.72	0.66	0.56	0.60	0.45
MCI TEO	0.23	0.30	0.29	0.38	0.22	0.30	0.36	0.10	0.31
MCI TEO-UMAP	0.31	0.26	0.21	0.22	0.22	0.36	0.40	0.40	0.23
LMCI TEO	0.19	0.19	0.20	0.10	0.22	0.16	0.24	0.10	0.22
LMCI TEO-UMAP	0.32	0.27	0.23	0.25	0.24	0.28	0.40	0.40	0.43

Table 6: Abbreviations: Dem = Demographics; VS = Vital Signs; CDT = Clock Drawing Test; CCT = Clock Copying Test; AVLT1/2 = Auditory Verbal Learning Test (v1/v2); CFA = Category Fluency (Animals); ANART = American National Adult Reading Test; FAQ = Functional Activities Questionnaire.

Collectively, these findings demonstrate that incorporating the SAE layer—thereby enforcing a more monosemantic, disentangled attribution representation—substantially enhances the clinical meaningfulness, stability, and diagnostic alignment of the resulting explanations.

B.9 THE CLINICAL IMPACT AND OUTCOME IN THE DIAGNOSIS OF ALZHEIMER, EARLY MCI AND LATE MCI.

This study shows that the Transformer Explanation Optimizer (TEO) with a Sparse Autoencoder (SAE) and TEO-UMAP provide the most reliable identification of informative sources across nine multimodal subgroups: Demographics (DEM), Vital Signs (VS), Clock Drawing Test (CDT), Clock Copying Test (CCT), Auditory Verbal Learning Test v1 (AVLT1), Category Fluency—Animals (CFA), Auditory Verbal Learning Test v2 (AVLT2), American National Adult Reading Test (ANART), and Functional Activities Questionnaire (FAQ). Using a significance threshold of 0.6 on PCA principal components PC1, we observe in the binary task that, for Control, TEO-SAE is dominated by FAQ, whereas TEO-UMAP emphasises DEM, AVLT2, and FAQ; for Alzheimer’s, TEO prioritises FAQ, AVLT1, and CFA, while TEO-UMAP highlights ANART, FAQ, and DEM. In the three-class task, for Control the main contributors are AVLT1, CDT, and ANART under TEO, and AVLT1, CDT, and CFA under TEO-UMAP; for MCI, TEO favours CCT, AVLT2, and FAQ, whereas TEO-UMAP favours AVLT2, ANART, and CFA; and for LMCI, TEO elevates AVLT1, FAQ, and CDT, while TEO-UMAP elevates FAQ, ANART, and AVLT2. These patterns, summarised in Table 6, support the clinical interpretability of the proposed optimisers.

Across ADNI cohorts, the most stable signals for clinical stratification are functional status (FAQ) and memory measures (AVLT1/AVLT2), with visuospatial performance (CDT) recurrent in Control/LMCI. TEO+SAE preferentially elevates neuropsychological performance features (AVLT1/2, CDT, CCT), while TEO-UMAP surfaces complementary contextual/language markers (DEM, ANART, CFA), yielding class-specific, interpretable profiles: Control—FAQ/AVLT1/CDT; Alzheimer’s—FAQ with AVLT1/CFA (TEO) or ANART/DEM (TEO-UMAP); MCI—AVLT2 with CCT/FAQ (TEO) or ANART/CFA (TEO-UMAP); LMCI—FAQ with AVLT1/CDT (TEO) or ANART/AVLT2 (TEO-UMAP). Using a simple PC1 0.6 significance rule, these optimisers provide actionable attribution maps that can prioritise assessments, reduce testing burden, support trial enrichment, and guide personalised monitoring. Together, they offer a practically deployable, transparent framework for clinically meaningful multimodal reasoning in neurodegenerative disease.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

Binary ADNI Control (0) **Binary ADNI Alzheimer (1)**

Given the two CSV files, and recognizing that medical biases exist in the char column with each of samples beginning with the character sequence [C L S], determine how each of the three attribution methods (attr1, attr2, attr3) highlights features that may be associated with healthy or unhealthy interpretations, and analyze the reasons for these differences. Try to predict the pathology. Mention which of the two csv files highlight the pathological case for each of the three attr techniques and try to specify the pathology.

Method	What It Highlights Most	Which CSV Shows the Pathological Pattern	Why This CSV Looks Pathological	Likely Pathology Suggested
attr1	Score terms (Total, Intrusions, Correct), demographic fields (Birth year, education, Race), test-related error wording	CSV 1	Stronger emphasis on demographic risk indicators and on error-related cognitive test features	Amnesic cognitive disorder (likely early Alzheimer's / amnesic MCI)
attr2	Task names, specific scores, intrusion counts, slow Trails B areas	CSV 1	More attention to the worst-performing cognitive segments and to demographic correlates	Amnesic cognitive disorder (likely early Alzheimer's / amnesic MCI)
attr3	Whole neuropsych blocks (Clock, AVLT, Trails), demographic segments, sections with intrusions or slow processing	CSV 1	Highlights both demographic risk factors and impaired memory/executive performance more strongly	Amnesic cognitive disorder (likely early Alzheimer's / amnesic MCI)

(a) Two class ADNI classification experiment

Binary BrainLAT Control (BL_0) **Binary BrainLAT Alzheimer (BL_1)**

Given the two CSV files, and recognizing that medical biases exist in the char column with each of samples beginning with the character sequence [C L S], determine how each of the three attribution methods (attr1, attr2, attr3) highlights features that may be associated with healthy or unhealthy interpretations, and analyze the reasons for these differences. Try to predict the pathology. Mention which of the two csv files highlight the pathological case for each of the three attr techniques and try to specify the pathology.

Method	What It Highlights Most	Pathological CSV	Why This CSV Appears Pathological	Suggested Pathology
attr1	Notable emphasis on key cognitive-test terms, scoring markers, and error-related text segments	BL_1	Shows greater sensitivity to segments containing performance weakness indicators and demographic risk markers	Memory-dominant cognitive deficit (amnesic MCI / early Alzheimer's)
attr2	Focus on specific task labels and parts of the text associated with reduced performance or irregular responses	BL_1	Highlights sections where performance is most impaired, such as memory intrusions or slowed sequences	Amnesic pattern consistent with early neurodegenerative disease
attr3	Broad blocks of neuropsych content, especially where errors, slow timing, or intrusions appear	BL_1	Shows widespread emphasis across impaired-performance regions and demographic variables	Early Alzheimer's-type cognitive impairment

(b) Two class BrainLAT classification experiment

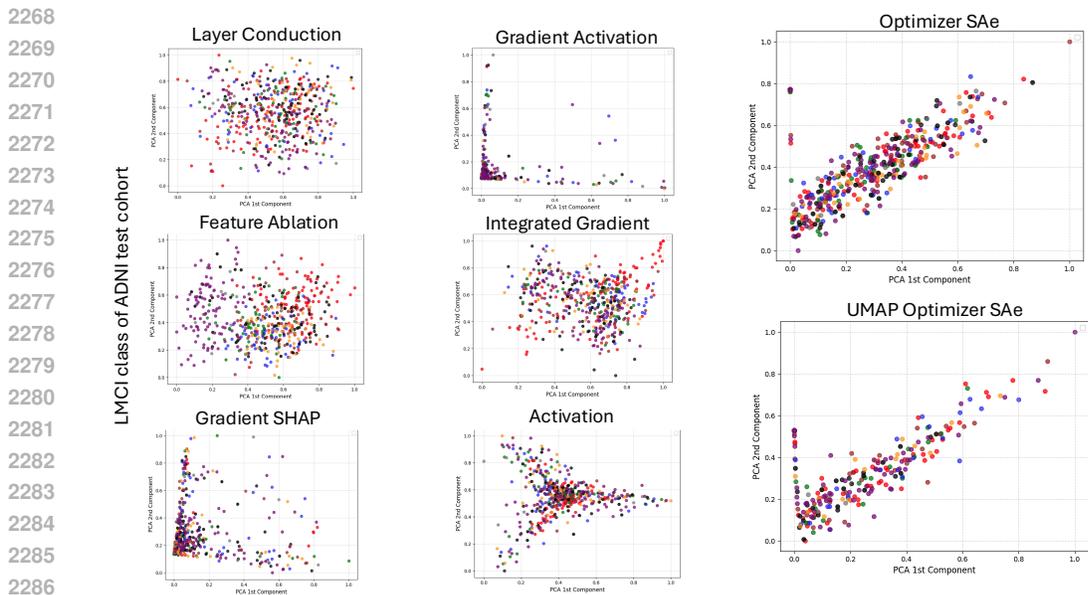
Three class ADNI Control (0) **Three class ADNI MCI (1)** **Three class ADNI LMCI (2)**

Given the three CSV files, and recognizing that medical biases exist in the char column with each of samples beginning with the character sequence [C L S], determine how each of the three attribution methods (attr1, attr2, attr3) highlights features that may be associated with healthy or unhealthy interpretations, and analyze the reasons for these differences. Try to predict the pathology or pathologies and the healthy if exist. Mention which of the three csv files highlight the pathological case/s for each of the three attr techniques and try to specify the pathology/ies and the same for the healthy (control if exist).

Attribution Pattern Type	Behavior Observed in CSV File 0	Behavior Observed in CSV File 1	Behavior Observed in CSV File 2	Interpretation (Pathology / Healthy) ¹
Pattern Type A (very sparse, near-zero emphasis overall)	Minimal highlighting of any medically relevant fields; almost all characters left unmarked.	Slight but still very limited emphasis on demographic fields and test labels.	Noticeably more highlighting on demographics, test headers, and structural elements of the note compared to CSV 0 and 1.	CSV 0 corresponds best to healthy / control class . CSV 1 corresponds to milder pathology , CSV 2 to more severe pathology .
Pattern Type B (extremely flat, low-information pattern)	Very low number of highlighted characters; almost uniformly neutral.	Similarly low highlighting; no meaningful differentiation.	Same low-density pattern; virtually no medically decisive areas marked.	This pattern does not meaningfully distinguish healthy vs pathological classes across any CSV files.
Pattern Type C (dense, diffuse, broad coverage)	Broad but slightly lower coverage of demographic fields, vitals, and test labels.	Broad coverage similar to CSV 2 but with slightly less emphasis.	Highest density; demographics, vitals, and test labels consistently highlighted.	CSV 0 again aligns with healthy / control . CSV 1 and CSV 2 align with pathology classes , with CSV 2 indicating more severe pathology .

(c) Three class ADNI classification experiment.

Figure 36: GPT-5.1 generated global explanations characterizing attribution-score performance in biomarker identification for TEO, both with and without the SAE layer, and for the TEO-UMAP model. (a) Binary ADNI: Control vs. Alzheimer's disease, with each class provided as a separate CSV file for external evaluation. (b) Binary BrainLAT: Control vs. Alzheimer's disease, similarly split into two class-specific CSV files. (c) Three-class ADNI: Control, MCI, and LMCI categories, each represented in its own CSV file.



2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299

Figure 37: Principal Component Analysis (PCA) of token-level attribution representations across methods. Each scatter plot shows the first two principal components extracted from the top eight PCA components of the attribution matrix for the LMCI class of the ADNI test cohort. PCA is used here to assess whether the proposed UMAP linear constraint encourages the attribution space to adopt an approximately linear structure. Traditional attribution methods (Layer Conduction, Feature Ablation, Gradient SHAP, Gradient Activation, Integrated Gradients, and Activation) exhibit more dispersed or irregular PCA trajectories, indicating nonlinear or highly variable attribution topologies. In contrast, the Optimizer SAE and UMAP Optimizer SAE methods show a clear linear trend along the first two principal components, demonstrating that the constrained UMAP formulation produces a substantially more linear and stable representation of attribution scores across the shared tokenizer feature space. These results support the hypothesis that enforcing a linear constraint in the UMAP embedding enhances structural consistency and robustness of feature-level attributions.

2300

2301 B.10 VALIDATION OF THE UMAP LINEAR CONSTRAINT VIA PCA STRUCTURE ANALYSIS

2302

2303 To verify the claim that the proposed UMAP linear constraint effectively linearizes the majority
2304 of the attribution space and yields robust attribution scores within the same tokenized feature,
2305 we performed an additional PCA analysis. Specifically, we extracted the top eight principal
2306 components from the attribution matrix and visualized the first two components, which capture
2307 the highest proportion of variance in the data. This allows us to assess whether the tokenized
2308 features exhibit an approximately linear structure in their dominant statistical directions. The
2309 results (Figure 37) demonstrate that the proposed method indeed produces an embedding that
2310 approaches a linear configuration, thereby supporting our hypothesis that the UMAP linear
2311 constraint leads to more stable and structurally consistent attribution representations.

2312

2313 B.11 BROAD IMPACT STATEMENT

2314 The clinical deployment of large language models (LLMs) in high-stakes neurodegenerative disease
2315 diagnosis, such as Alzheimer’s Disease (AD), is hindered by the inherent polysemanticity of their
2316 representations, which renders traditional attribution methods (e.g., gradients, SHAP) unreliable
2317 due to ambiguous or inconsistent explanations. By aligning LLM explanations with clinical
2318 reasoning and enforcing statistical fidelity, this work establishes a foundation for trustworthy,
2319 deployable AI systems in medicine, transforming complex models into transparent partners for life-
2320 critical decision-making and paving the way for safer, ethically sound integration of advanced AI
2321 into cognitive health applications. Critically, this framework’s adaptability and rigorous validation
position it for immediate real-world deployment in healthcare settings, enabling clinicians to

2322 harness LLMs' diagnostic power without compromising transparency, thereby accelerating the
2323 translation of AI research into measurable improvements in patient care.
2324

2325 REFERENCES
2326

2327 Chirag Agarwal, Nari Johnson, Martin Pawelczyk, Satyapriya Krishna, Eshika Saxena, Marinka
2328 Zitnik, and Himabindu Lakkaraju. Rethinking stability for attribution-based explanations, 2022.
2329 URL <https://arxiv.org/abs/2203.06877>.

2330 Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller,
2331 and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-
2332 wise relevance propagation. *PLOS ONE*, 10(7):1–46, 07 2015. doi: 10.1371/journal.pone.0130140.
2333 URL <https://doi.org/10.1371/journal.pone.0130140>.

2334 Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety – a review, 2024.
2335 URL <https://arxiv.org/abs/2404.14082>.

2337 Peter Bills, Jyothi Guntupalli, et al. Language models represent space and time. *Nature Neuro-*
2338 *science*, 26(5):707–717, 2023a.

2339 Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever,
2340 Jan Leike, Jeff Wu, and William Saunders. Language models can explain neurons in language
2341 models. *arXiv preprint arXiv:2306.00604*, 2023b. URL [https://arxiv.org/abs/2306.](https://arxiv.org/abs/2306.00604)
2342 [00604](https://arxiv.org/abs/2306.00604).

2343 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
2344 Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decom-
2345 posing language models with dictionary learning. *Transformer Circuits Thread*, 2023. URL
2346 <https://transformer-circuits.pub/2023/monosemanticity/index.html>.

2347 Prasad Chalasani, Jiefeng Chen, Amrita Roy Chowdhury, Xi Wu, and Somesh Jha. Concise expla-
2348 nations of neural networks using adversarial training. In Hal Daumé III and Aarti Singh (eds.),
2349 *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of
2350 *Proceedings of Machine Learning Research*, pp. 1383–1391, Virtual Event, online, July 13–18 2020.
2351 PMLR. Originally released as arXiv:1810.06583 (2018).
2352

2353 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-
2354 coders find highly interpretable features in language models, 2023. URL [https://arxiv.](https://arxiv.org/abs/2309.08600)
2355 [org/abs/2309.08600](https://arxiv.org/abs/2309.08600).

2356 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,
2357 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCand-
2358 lish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of
2359 superposition, 2022a. URL <https://arxiv.org/abs/2209.10652>.

2360 Nelson Elhage, Neel Nanda, et al. A mechanistic interpretability analysis of superposition in neural
2361 networks. *Transformer Circuits Thread*, 2022b. URL [https://transformer-circuits.](https://transformer-circuits.pub/2022/superposition/)
2362 [pub/2022/superposition/](https://transformer-circuits.pub/2022/superposition/).

2363 Maria Luisa Gorno-Tempini, Argye E Hillis, Sandra Weintraub, Andrew Kertesz, Mario Mendez,
2364 Stefano F Cappa, Jennifer M Ogar, Jonathan D Rohrer, Sandra Black, Bradley F Boeve, et al.
2365 Classification of primary progressive aphasia and its variants. *Neurology*, 76(11):1006–1014,
2366 2011.
2367

2368 Liv Gorton. The missing curve detectors of inceptionv1: Applying sparse autoencoders to incep-
2369 tionv1 early vision. *arXiv preprint arXiv:2406.03662*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2406.03662)
2370 [2406.03662](https://arxiv.org/abs/2406.03662).
2371

2372 Anna Hedström, Leander Weber, Daniel Krakowczyk, Dilyara Bareeva, Franz Motzkus, Wojciech
2373 Samek, Sebastian Lapuschkin, and Marina Marina M.-C. Höhne. Quantus: An explainable
2374 ai toolkit for responsible evaluation of neural network explanations and beyond. *Journal of*
2375 *Machine Learning Research*, 24(34):1–11, 2023. URL [http://jmlr.org/papers/v24/](http://jmlr.org/papers/v24/22-0142.html)
[22-0142.html](http://jmlr.org/papers/v24/22-0142.html).

- 2376 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint*
2377 *arxiv:2006.11239*, 2020.
2378
- 2379 Clifford R Jack, David A Bennett, Kaj Blennow, Maria C Carrillo, Bruce Dunn, Susan B Haeblerlein,
2380 David M Holtzman, William Jagust, Frank Jessen, Jason Karlawish, et al. Nia-aa research
2381 framework: Toward a biological definition of alzheimer’s disease. *Alzheimer’s & Dementia*, 14(4):
2382 535–562, 2018.
- 2383 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014. URL
2384 <https://arxiv.org/abs/1412.6980>.
2385
- 2386 Patrick Kolpaczki, Viktor Bengs, Maximilian Muschalik, and Eyke Hüllermeier. Approximating the
2387 shapley value without marginal contributions, 2024.
- 2388 Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *CoRR*,
2389 <abs/1705.07874>, 2017. URL <http://arxiv.org/abs/1705.07874>.
2390
- 2391 Michail Mamalakis, Antonios Mamalakis, Ingrid Agartz, Lynn Egeland Mørch-Johnsen, Graham K.
2392 Murray, John Suckling, and Pietro Lio. Solving the enigma: Enhancing faithfulness and com-
2393 prehensibility in explanations of deep networks. *AI Open*, 6:70–81, 2025. ISSN 2666-6510. doi:
2394 [10.1016/j.aiopen.2025.02.001](https://doi.org/10.1016/j.aiopen.2025.02.001). URL [http://dx.doi.org/10.1016/j.aiopen.2025.](http://dx.doi.org/10.1016/j.aiopen.2025.02.001)
2395 [02.001](http://dx.doi.org/10.1016/j.aiopen.2025.02.001).
- 2396 Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse
2397 feature circuits: Discovering and editing interpretable causal graphs in language models. *arXiv*
2398 *preprint arXiv:2403.19647*, 2024. URL <https://arxiv.org/abs/2403.19647>.
2399
- 2400 Callum McDougall. Sae visualizer, 2024. URL [https://github.com/callumcdougall/](https://github.com/callumcdougall/SAE-Visualizer)
2401 [SAE-Visualizer](https://github.com/callumcdougall/SAE-Visualizer).
- 2402 Susanne G. Mueller, Michael W. Weiner, Leon J. Thal, Ronald C. Petersen, Clifford Jack, William
2403 Jagust, John Q. Trojanowski, Arthur W. Toga, and Laurel Beckett. The alzheimer’s disease
2404 neuroimaging initiative. *Neuroimaging Clinics of North America*, 15(4):869–877, 2005. ISSN 1052-
2405 5149. doi: <https://doi.org/10.1016/j.nic.2005.09.008>. URL [https://www.sciencedirect.](https://www.sciencedirect.com/science/article/pii/S1052514905001024)
2406 [com/science/article/pii/S1052514905001024](https://www.sciencedirect.com/science/article/pii/S1052514905001024). Alzheimer’s Disease: 100 Years of
2407 Progress.
- 2408 Scott C. Neu, Karen L. Crawford, and Arthur W. Toga. The image and data archive at the laboratory
2409 of neuro imaging. *Frontiers in Neuroinformatics*, Volume 17 - 2023, 2023. ISSN 1662-5196.
2410 doi: [10.3389/fninf.2023.1173623](https://doi.org/10.3389/fninf.2023.1173623). URL [https://www.frontiersin.org/journals/](https://www.frontiersin.org/journals/neuroinformatics/articles/10.3389/fninf.2023.1173623)
2411 [neuroinformatics/articles/10.3389/fninf.2023.1173623](https://www.frontiersin.org/journals/neuroinformatics/articles/10.3389/fninf.2023.1173623).
2412
- 2413 Toan Q. Nguyen and Julian Salazar. Transformers without tears: Improving the normalization of
2414 self-attention. 2019. doi: [10.5281/zenodo.3525484](https://doi.org/10.5281/zenodo.3525484).
- 2415 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.
2416 Zoom in: An introduction to circuits. *Distill*, 5(3):e00024.001, 2020a. doi: [10.23915/distill.00024.](https://doi.org/10.23915/distill.00024.001)
2417 [001](https://doi.org/10.23915/distill.00024.001). URL <https://distill.pub/2020/circuits/zoom-in>.
2418
- 2419 Chris Olah, Arvind Satyanarayan, Ludwig Schubert Wusser, and Shan Carter. Zoom in: An intro-
2420 duction to circuits. *Distill*, 2020b. doi: [10.23915/distill.00024.001](https://doi.org/10.23915/distill.00024.001).
2421
- 2422 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
2423 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction
2424 heads. *Transformer Circuits Thread*, 2022. URL [https://transformer-circuits.pub/](https://transformer-circuits.pub/2022/in-context-learning/index.html)
2425 [2022/in-context-learning/index.html](https://transformer-circuits.pub/2022/in-context-learning/index.html).
- 2426 OpenAI. ChatGPT-4o. <https://openai.com/chatgpt>, 2024. Accessed May 2025.
2427
- 2428 Ronald C. Petersen, Glenn E. Smith, Susan C. Waring, Robert J. Ivnik, Eric G. Tangalos, and Emre
2429 Kokmen. Mild cognitive impairment: clinical characterization and outcome. *Archives of*
Neurology, 56(3):303–308, 1999.

- 2430 Pavel Prado, Vicente Medel, Agustín Sainz-Ballesteros, Hernando Santamaría-García, Sebastián
2431 Moguilner, Jhony Mejía, Raúl González-Gómez, Andrea Slachevsky, María Isabel Behrens, David
2432 Aguillón, Francisco Lopera, Mario A. Parra, Diana Matallana, Marcelo Adrián Maito, Adolfo M.
2433 García, Nilton Custodio, Alberto Ávila Funes, Stefanie Piña-Escudero, Agustina Birba, Sol Fitti-
2434 paldi, Agustina Legaz, and Agustín Ibáñez. The brainlat project: a multimodal neuroimaging
2435 dataset of neurodegeneration from underrepresented backgrounds. *Scientific Data*, 10:889,
2436 2023. doi: 10.1038/s41597-023-02806-8. URL <https://www.nature.com/articles/s41597-023-02806-8>.
2437
- 2438 Rodrigo Quiroga et al. Invariant visual representation by single neurons in the human brain.
2439 *Nature*, 435(7045):1102–1107, 2005.
2440
- 2441 Senthoooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János
2442 Kramár, Rohin Shah, and Neel Nanda. Improving dictionary learning with gated sparse autoen-
2443 coders, 2024. URL <https://arxiv.org/abs/2404.16014>.
- 2444 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the
2445 predictions of any classifier, 2016. URL <https://arxiv.org/abs/1602.04938>.
2446
- 2447 Mattia Rigotti, Omri Barak, Melissa Warden, et al. The importance of mixed selectivity in complex
2448 cognitive tasks. *Nature*, 497(7451):585–590, 2013.
- 2449 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for
2450 biomedical image segmentation. In *International Conference on Medical Image Computing
2451 and Computer-Assisted Intervention (MICCAI)*, pp. 234–241. Springer, 2015. doi: 10.1007/
2452 978-3-319-24574-4_28.
- 2453 Amitojdeep Singh, Sourya Sengupta, and Vasudevan Lakshminarayanan. Explainable deep learn-
2454 ing models in medical image analysis, 2020.
2455
- 2456 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks, 2017.
2457
- 2458 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam
2459 Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, et al. Scaling monosemanticity: Extracting
2460 interpretable features from claude 3. *Transformer Circuits Thread*, 2024. URL [https://
2461 transformer-circuits.pub/2024/scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 2462 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
2463 Kaiser, and Illia Polosukhin. Attention is all you need, 2017.
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483