

---

# Geometry-Adaptive Explainer for Faithful Dictionary-Based Interpretability under Distribution Shift

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Mechanistic interpretability aims to explain a model’s behavior by identifying  
2 causally responsible internal structures. Dictionary-based explainers such as sparse  
3 autoencoders and transcoders are a primary tool, but their faithfulness under out-  
4 of-distribution (OOD) shift has received little systematic attention. We show that  
5 distribution shift rotates the subspace that the model actively uses, misaligning  
6 the explainer’s dictionary trained on in-distribution (ID) activations. We formalize  
7 this misalignment as the **faithfulness gap**, a geometric distance between the ID  
8 dictionary and the OOD-active subspace, and show that it controls OOD faith-  
9 fulness degradation. To reduce this gap, we propose the **Geometry-Adaptive**  
10 **Explainer (GAE)**, which realigns the explainer’s dictionary with the OOD-active  
11 subspace while preserving the original feature structure. This requires only unlabeled  
12 OOD activations and no gradient updates. We prove that GAE improves over  
13 the unadapted ID explainer, with excess loss bounded quadratically by the second-  
14 moment shift. Empirically, GAE even matches or surpasses all training-based  
15 baselines in causal faithfulness across multiple models and OOD settings.<sup>1</sup>

## 16 1 Introduction

17 Mechanistic interpretability aims to explain a model’s behavior by identifying internal structures  
18 that are causally responsible for its outputs [1–3]. A primary approach is to train an explainer, a  
19 post-hoc module such as a sparse autoencoder (SAE) [4] or transcoder [5], that decomposes hidden  
20 activations into sparse combinations of learned feature directions (a dictionary). These dictionary-  
21 based explainers have recently been scaled to large language models [6, 7] and used to uncover  
22 interpretable feature circuits [8]. A central requirement for such explanations is faithfulness: they  
23 should accurately reflect the computations the model actually uses [9, 1].

24 When a model encounters out-of-distribution (OOD) inputs, the dictionary learned in-distribution  
25 (ID) can no longer capture the directions the model actively uses [10]. Prior work has addressed  
26 this vulnerability from two angles, each limited in scope. Attribution-level robustness studies [11–  
27 14] focus on input perturbations rather than hidden-state geometry. Dictionary-explainer remedies  
28 such as retraining on the model’s own generations [15], upweighting rare concepts [16, 17], and  
29 adding residual modules [18] remain heuristic, without diagnosing the underlying misalignment.  
30 Consequently, the mechanism driving this failure and a principled correction remain unaddressed.

31 In this work, we identify a geometric mechanism for OOD faithfulness degradation in dictionary-based  
32 explainers. These explainers learn their feature directions from ID activations, so their dictionaries  
33 reflect the geometric structure of ID hidden representations [19]. This structure is captured by

---

<sup>1</sup>Code is available at: <https://anonymous.4open.science/r/GAE/>

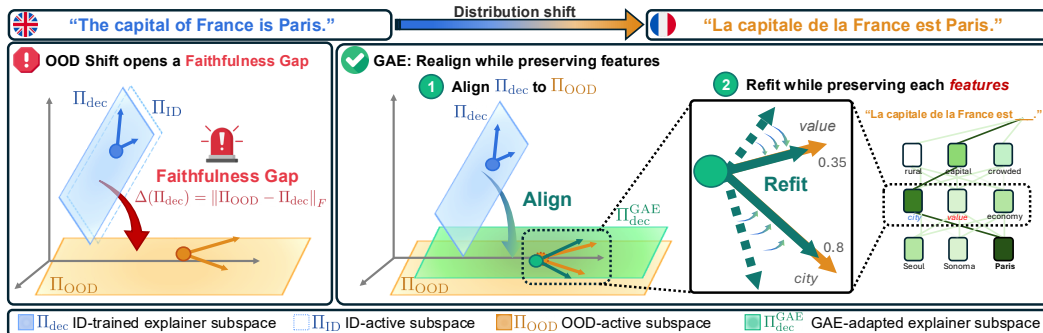


Figure 1: **Faithfulness gap and GAE.** *Left:* distribution shift (illustrated as a language change) rotates the OOD-active subspace  $\Pi_{\text{OOD}}$  away from the ID-trained explainer subspace  $\Pi_{\text{dec}} \approx \Pi_{\text{ID}}$ , opening a faithfulness gap  $\Delta(\Pi_{\text{dec}})$ . *Right:* GAE closes this gap in two steps. Step 1 rotates  $\Pi_{\text{dec}}$  onto  $\Pi_{\text{OOD}}$  via orthogonal Procrustes. Step 2 refits individual feature directions within the aligned subspace to match OOD activations while preserving the original feature structure.

34 the second moment of hidden activations, which distribution shift typically alters [20], leaving the  
 35 ID-trained dictionary misaligned with the OOD-active subspace. As illustrated in Figure 1 (left), we  
 36 call this misalignment the **faithfulness gap**, a geometric distance between the ID dictionary and the  
 37 OOD-active subspace. We prove that this gap controls OOD faithfulness degradation and that it is  
 38 itself upper-bounded by the magnitude of the second-moment shift. Reducing this gap is therefore  
 39 necessary for restoring OOD faithfulness, motivating methods that directly realign the dictionary.

40 We instantiate this idea with the **Geometry-Adaptive Explainer (GAE)**, a closed-form, post-hoc  
 41 method that closes the faithfulness gap using only unlabeled OOD activations. GAE first rotates the  
 42 ID-trained dictionary so that its subspace aligns with the OOD-active subspace, choosing the rotation  
 43 closest to the original dictionary to preserve feature structure (Step 1 in Figure 1). A constrained  
 44 decoder refit then adjusts individual feature directions to match OOD activations while maintaining  
 45 this alignment (Step 2 in Figure 1). The entire pipeline requires no gradient computation, yet we  
 46 prove it is guaranteed to improve over the unadapted explainer (Theorem 1). Empirically, GAE  
 47 matches or surpasses all training-based baselines in causal faithfulness across multiple language  
 48 models and OOD settings, including methods that retrain from scratch on OOD data.

49 Our main contributions are summarized as follows:

- 50 • We identify the **faithfulness gap** as a geometric mechanism for OOD faithfulness degradation and  
 51 prove that excess loss grows at most quadratically with second-moment shift.
- 52 • We propose **GAE**, a closed-form dictionary realignment method that targets the faithfulness gap  
 53 with a theoretical guarantee while preserving feature structure.
- 54 • Without any gradient computation, GAE matches or surpasses all training-based baselines, including  
 55 full OOD retraining, in causal faithfulness across multiple models and diverse OOD settings.

## 56 2 Related Work

### 57 2.1 Faithfulness in Mechanistic Interpretability

58 Dictionary-based explainers such as SAEs [4, 21] and transcoders [5, 22] decompose hidden acti-  
 59 vations into sparse feature directions and have become a primary tool for mechanistic interpretability  
 60 [1, 2, 6–8]. Faithfulness is typically evaluated via causal interventions that ablate features and  
 61 measure the effect on model outputs [9, 23, 24], and is a prerequisite for reliable circuit discovery  
 62 and model editing [8, 25, 26]. These explainers’ dictionaries reflect the geometric structure of ID  
 63 hidden representations [19], so distribution shift can degrade faithfulness by misaligning the learned  
 64 directions with those the model actively uses. Recent empirical reports support this concern: SAE-  
 65 based features underperform dense baselines on OOD downstream tasks [10], yet this vulnerability  
 66 has received little systematic attention, with no formal diagnosis of its cause.

### 67 2.2 Interpretability under Distribution Shift

68 The faithfulness assumption above becomes problematic when the model encounters OOD inputs.  
 69 Early work showed that saliency maps are fragile under input perturbations [11, 12]. Subsequent  
 70 studies examined explanation consistency under shift more broadly [14, 13]. These analyses primarily

71 concern attribution methods and attribute failures to predictive degradation or input-level perturbations,  
 72 rather than to structural changes in hidden representations. Yet evidence from OOD detection  
 73 shows that OOD inputs produce statistically distinguishable activation patterns in hidden layers [20].  
 74 Representation comparison methods [27–29] enable quantifying such geometric changes across  
 75 conditions, but the connection to explainer faithfulness has not been drawn.

76 For dictionary-based explainers, proposed remedies include retraining the explainer on OOD data,  
 77 training on the model’s own generations to avoid external dataset dependence [15], upweighting  
 78 tail samples during training [16, 17], or adding residual capacity via boosting [18]. However, these  
 79 approaches either require full retraining or do not directly address the geometric misalignment  
 80 between the explainer’s learned dictionary and the OOD-active subspace. Post-hoc methods that  
 81 realign the explainer’s dictionary with the OOD-active subspace remain unexplored.

### 82 3 Hidden-Space Geometry Shift and Faithfulness Degradation

83 When a neural network encounters OOD inputs, do its mechanistic explanations remain faithful? We  
 84 show that they generally do not. Distribution shift alters the geometry of hidden activations, creating  
 85 a misalignment between the directions the model actively uses and those the explainer was trained  
 86 to reconstruct. We call this misalignment the **faithfulness gap**, and prove that (i) it grows at most  
 87 proportionally with the second-moment shift for ID-trained explainers (Proposition 1), and (ii) it  
 88 controls the reducible part of OOD faithfulness loss (Proposition 2).

#### 89 3.1 Setup and Faithfulness Gap

90 **Target model and explainer.** The target model is a fixed, pretrained neural network whose internal  
 91 computations we wish to explain. For an input  $X$ , it produces a hidden representation  $h(X) \in \mathbb{R}^d$  at  
 92 a designated layer. An explainer is a post-hoc module that decomposes hidden representations into  
 93 interpretable components. Among various approaches, dictionary-based explainers such as SAEs  
 94 and transcoders have become a primary tool for mechanistic interpretability [4, 30]. These methods  
 95 learn a decoder (dictionary)  $W_{\text{dec}} \in \mathbb{R}^{d \times k}$  and an encoder  $W_{\text{enc}} \in \mathbb{R}^{k \times d}$ , where the  $k$  columns of  
 96  $W_{\text{dec}}$  serve as learned feature directions. SAEs reconstruct the hidden activation  $h(X)$  itself, while  
 97 transcoders reconstruct the MLP output at the same layer; both share the form

$$\hat{h}(X) = W_{\text{dec}} \sigma(W_{\text{enc}} h(X) + b_{\text{enc}}) + b_{\text{dec}}, \quad (1)$$

98 where  $\sigma$  is a sparsifying nonlinearity (e.g., ReLU or TopK) [21]. The explainer represents hidden  
 99 activations through a learned dictionary, so its behavior is tied to hidden-space geometry.

100 **In-distribution, out-of-distribution, and second-moment shift.** We call the data distribution  
 101 on which the explainer was trained ID, denoted by  $P_{\text{ID}}$ ; the dictionary  $W_{\text{dec}}$  is learned from ID  
 102 activations. At deployment, however, the explainer may encounter OOD inputs from a different  
 103 distribution  $P_{\text{OOD}}$ . Our goal is to understand whether the explainer remains faithful under this shift.

104 We characterize the shift via second-moment matrix  $M_e := \mathbb{E}_{X \sim P_e} [h(X)h(X)^\top]$ ,  $e \in \{\text{ID}, \text{OOD}\}$ .  
 105 Dictionary-based explainers minimize reconstruction error whose optimum depends on the eigen-  
 106 structure of  $M_e$  [21], so a shift in  $M_e$  changes the directions the explainer should reconstruct. A  
 107 **second-moment shift** occurs when  $M_{\text{OOD}} \neq M_{\text{ID}}$ ; we measure its magnitude by  $\|M_{\text{OOD}} - M_{\text{ID}}\|_F$ .

108 **Active subspace and faithfulness gap.** Since hidden activations concentrate energy along a  
 109 few directions [31–33], the top eigenspace of  $M_e$  captures most of the model’s representational  
 110 activity [28]. We write  $\Pi_e = U_e U_e^\top$ , where  $U_e \in \mathbb{R}^{d \times r}$  contains the top- $r$  eigenvectors of  $M_e$ ,  
 111 and call this the **active subspace** in environment  $e$ . Every reconstruction lies in the column space  
 112 of  $W_{\text{dec}}$ , but the reconstruction energy concentrates along the top- $r$  left singular directions [21].  
 113 Similarly,  $\Pi_{\text{dec}} = U_{\text{dec}} U_{\text{dec}}^\top$ , where  $U_{\text{dec}} \in \mathbb{R}^{d \times r}$  contains the top- $r$  left singular vectors of  $W_{\text{dec}}$ , is  
 114 the **explainer subspace**. OOD faithfulness depends on how well  $\Pi_{\text{dec}}$  aligns with  $\Pi_{\text{OOD}}$ .

115 Under second-moment shift,  $\Pi_{\text{OOD}}$  may diverge from  $\Pi_{\text{ID}}$ , opening a gap between the explainer  
 116 subspace and OOD-active subspace. We call this misalignment the **faithfulness gap**, which tightly  
 117 controls how much faithfulness degrades under OOD, as we formalize in Proposition 2:

118 **Definition 1** (Faithfulness Gap). *The faithfulness gap of an explainer subspace  $\Pi_{\text{dec}}$  under OOD is*  

$$\Delta(\Pi_{\text{dec}}) := \|\Pi_{\text{OOD}} - \Pi_{\text{dec}}\|_F.$$

119 A large gap means the explainer is reconstructing along directions that the model no longer uses, and  
 120 the second-moment shift directly upper-bounds this gap for ID-trained explainers (Proposition 1).

### 121 3.2 Second-Moment Shift Enlarges the Faithfulness Gap

122 The previous subsection defined the faithfulness gap as a geometric quantity. We now show that  
 123 second-moment shift directly enlarges this gap for ID-trained explainers. In practice, explainers  
 124 are trained on ID activations and deployed without modification [30, 34]. Since a well-trained ID  
 125 explainer satisfies  $\Pi_{\text{dec}} \approx \Pi_{\text{ID}}$  (empirically validated in Appendix B.2, Table 4), its OOD faithfulness  
 126 depends on how far  $\Pi_{\text{ID}}$  lies from  $\Pi_{\text{OOD}}$ . The following result, a consequence of the Davis–Kahan  
 127  $\sin \Theta$  theorem [35], bounds this distance in terms of the second-moment shift.

128 **Proposition 1** (Second-Moment Shift Bounds the Faithfulness Gap of the ID Explainer). *Suppose*  
 129 *that  $M_{\text{ID}}$  has eigengap  $\gamma_{\text{ID}} = \lambda_r(M_{\text{ID}}) - \lambda_{r+1}(M_{\text{ID}}) \geq 0$  at rank  $r$ . Then*

$$\Delta(\Pi_{\text{ID}}) = \|\Pi_{\text{OOD}} - \Pi_{\text{ID}}\|_F \leq \frac{\sqrt{2}}{\gamma_{\text{ID}}} \|M_{\text{OOD}} - M_{\text{ID}}\|_F.$$

130 The faithfulness gap of the ID explainer grows at most proportionally with the second-moment shift  
 131  $\|M_{\text{OOD}} - M_{\text{ID}}\|_F$ , with sensitivity controlled by the inverse eigengap  $1/\gamma_{\text{ID}}$  (proof in Appendix A.1;  
 132 empirical verification in Appendix B.4). This establishes that the faithfulness gap can grow large  
 133 under distribution shift. The next question is whether reducing it improves faithfulness.

### 134 3.3 The Faithfulness Gap Controls OOD Degradation

135 We now formalize the connection between the faithfulness gap and OOD faithfulness loss, showing  
 136 that  $\Delta(\Pi_{\text{dec}})$  is the central quantity an adaptation method should target.

137 **OOD faithfulness objective.** An ideal explainer subspace  $\Pi_{\text{dec}}$  minimizes reconstruction error on  
 138 OOD activations. We formalize this objective as

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}) := \mathbb{E}_{X \sim P_{\text{OOD}}} \|h(X) - \Pi_{\text{dec}} h(X)\|_2^2. \quad (2)$$

139 This measures how much OOD activation is lost when projected onto  $\Pi_{\text{dec}}$ .  $\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}})$  is a valid  
 140 surrogate since hidden-layer reconstruction constrains logit-level faithfulness [23, 24].

141 **Decomposition of  $\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}})$ .** To isolate the part of the OOD loss that the explainer can reduce,  
 142 we decompose  $\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}})$  into two terms. Let  $\mathcal{C}_r$  denote the set of rank- $r$  orthogonal projectors in  
 143  $\mathbb{R}^d$ . For any  $\Pi_{\text{dec}} \in \mathcal{C}_r$ ,

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}) = \underbrace{\mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}})}_{\text{irreducible}} + \underbrace{\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}) - \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}})}_{\text{explainer-dependent}}. \quad (3)$$

144 The explainer-dependent component is nonnegative (proof in Appendix A.2), so  $\Pi_{\text{OOD}}$  minimizes  
 145  $\mathcal{L}_{\text{OOD}}$  over  $\mathcal{C}_r$ . The irreducible component depends on the target model and the OOD distribution,  
 146 both fixed at deployment; adapting the explainer cannot reduce it. *The explainer-dependent component*  
 147 *is the only part the explainer can reduce.*

148 **Faithfulness gap as a tight proxy.** The faithfulness gap  $\Delta(\Pi_{\text{dec}})$  measures the distance between  
 149 two subspaces, making it a direct optimization target. The next proposition shows that controlling  
 150  $\Delta(\Pi_{\text{dec}})$  is equivalent to controlling the explainer-dependent component.

151 **Proposition 2** (Faithfulness Gap Controls the Explainer-Dependent Term). *Assume that the OOD*  
 152 *eigengap at rank  $r$ ,  $\gamma_{\text{OOD}} = \lambda_r(M_{\text{OOD}}) - \lambda_{r+1}(M_{\text{OOD}}) > 0$ . Then for any  $\Pi_{\text{dec}} \in \mathcal{C}_r$ ,*

$$\frac{\gamma_{\text{OOD}}}{2} \Delta(\Pi_{\text{dec}})^2 \leq \mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}) - \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) \leq \frac{\lambda_1(M_{\text{OOD}}) - \lambda_d(M_{\text{OOD}})}{2} \Delta(\Pi_{\text{dec}})^2.$$

153 The lower bound shows that any nonzero gap incurs a positive cost (misalignment cannot be free);  
 154 the upper bound shows that reducing  $\Delta(\Pi_{\text{dec}})$  is sufficient. Together, they establish that controlling  
 155  $\Delta(\Pi_{\text{dec}})$  is equivalent to controlling the explainer-dependent faithfulness loss (proof in Appendix A.3;  
 156 empirical verification in Appendix B.3). *Reducing  $\Delta(\Pi_{\text{dec}})$  is therefore both necessary and sufficient*  
 157 *for improving OOD faithfulness.* Section 4 introduces a method that directly targets this quantity.

## 158 4 Geometry-Adaptive Explainer (GAE)

159 Section 3 showed that OOD faithfulness degradation is controlled by the faithfulness gap  $\Delta(\Pi_{\text{dec}})$ .  
 160 We now propose the Geometry-Adaptive Explainer (GAE), which reduces  $\Delta(\Pi_{\text{dec}})$  by realigning the  
 161 explainer’s dictionary with the OOD-active subspace while preserving the original feature structure.  
 162 We present an objective for this adaptation and a closed-form solution.

163 **4.1 Problem Formulation**

164 Section 3.2 showed that an ID-trained explainer’s faithfulness degrades under OOD because its  
 165 subspace diverges from  $\Pi_{\text{OOD}}$ . Minimizing  $\Delta(\Pi_{\text{dec}})$  amounts to choosing a  $W_{\text{dec}}$  whose induced  
 166 subspace aligns with  $\Pi_{\text{OOD}}$ . To restore faithfulness, we adapt the existing dictionary  $W_{\text{dec}}^{\text{ID}} \in \mathbb{R}^{d \times k}$   
 167 using a set of unlabeled OOD activations  $\{h_i\}_{i=1}^N$ , from which we estimate the OOD-active subspace  
 168  $\widehat{\Pi}_{\text{OOD}}$ . We seek  $W_{\text{dec}}$  that closes the faithfulness gap while preserving the original feature structure:

$$\min_{W_{\text{dec}}} \underbrace{\mathcal{L}_{\text{recon}}}_{\text{reconstruction}} + \lambda_{\text{geom}} \underbrace{\|\widehat{\Pi}_{\text{OOD}} - \Pi_{\text{dec}}\|_F^2}_{\text{subspace alignment}} + \lambda_{\text{pres}} \underbrace{\|W_{\text{dec}} - W_{\text{dec}}^{\text{ID}}\|_F^2}_{\text{feature preservation}}, \quad (4)$$

169 where  $\mathcal{L}_{\text{recon}} = \frac{1}{N} \sum_{i=1}^N \|h_i - \hat{h}_i\|_2^2$  is the mean reconstruction error.<sup>2</sup> GAE holds the encoder  
 170 fixed to preserve the learned feature decomposition [4]. The three terms serve complementary roles.  
 171 The *reconstruction* term fits individual OOD activations, since subspace alignment alone does not  
 172 guarantee sample-level reconstruction. The *subspace alignment* term directly targets the faithfulness  
 173 gap  $\Delta(\Pi_{\text{dec}})$ , which Proposition 2 showed controls the explainer-dependent component. The *feature*  
 174 *preservation* term keeps the adapted dictionary close to the original, maintaining the encoder-decoder  
 175 pairing for downstream circuit analyses [8].

176 **4.2 Dictionary Adaptation**

177 Eq. (4) is non-convex in  $W_{\text{dec}}$ , as  $\Pi_{\text{dec}}$  depends on it through a top- $r$  SVD. However, the subspace  
 178 alignment term can be enforced as a hard constraint, and once the subspace is fixed, the remaining  
 179 terms are quadratic with a closed-form solution. GAE exploits this in two steps (Algorithm 1):  
 180 Step 1 enforces *subspace alignment* by rotating the ID-trained dictionary onto  $\widehat{\Pi}_{\text{OOD}}$ , choosing the  
 181 rotation closest to the original. Step 2 refits the decoder via constrained ridge regression, solving the  
 182 *reconstruction* and *feature preservation* terms while preserving this alignment.

---

**Algorithm 1** GAE

---

**Require:**  $W_{\text{dec}}^{\text{ID}} \in \mathbb{R}^{d \times k}$ ;  $W_{\text{enc}}, b_{\text{enc}}$ ; OOD activations  $\{h_i\}_{i=1}^N$ ;  $\lambda_{\text{geom}}, \lambda_{\text{pres}}$

**Ensure:**  $(W_{\text{dec}}^{\text{GAE}}, b_{\text{dec}}^{\text{GAE}})$

- 1:  $\widehat{M}_{\text{OOD}} \leftarrow \frac{1}{N} \sum_{i=1}^N h_i h_i^\top$
  - 2:  $U_{\text{dec}} \leftarrow$  top- $r$  left singular vectors of  $W_{\text{dec}}^{\text{ID}}$ ;  $U_{\text{OOD}}^{(:r)} \leftarrow$  top- $r$  eigenvectors of  $\widehat{M}_{\text{OOD}}$
  - 3:  $G \leftarrow U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}} (W_{\text{dec}}^{\text{ID}})^\top U_{\text{OOD}}^{(:r)}$ ;  $T^* \leftarrow \widetilde{V} \widetilde{U}^\top$  from SVD( $G$ )
  - 4:  $\widetilde{W}_{\text{dec}} \leftarrow U_{\text{OOD}}^{(:r)} T^* U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}}$  ▷ Step 1: Subspace rotation
  - 5:  $z_i \leftarrow \sigma(W_{\text{enc}} h_i + b_{\text{enc}})$  for each  $h_i$  ▷ frozen encoder
  - 6:  $W_{\text{dec}}^{\text{GAE}} \leftarrow$  Eq. (8);  $b_{\text{dec}}^{\text{GAE}} \leftarrow$  Eq. (9) ▷ Step 2: Constrained decoder refit
- 

183 **Step 1: Subspace rotation.** Let  $U_{\text{dec}} \in \mathbb{R}^{d \times r}$  be the top- $r$  left singular vectors of  $W_{\text{dec}}^{\text{ID}}$  (the  
 184 explainer subspace defined in Section 3.1), and let  $U_{\text{OOD}}^{(:r)}$  be the top- $r$  eigenvectors of the empirical  
 185 second-moment matrix  $\widehat{M}_{\text{OOD}} = \frac{1}{N} \sum_{i=1}^N h_i h_i^\top$ . We constrain the rotated dictionary to the form

$$\widetilde{W}_{\text{dec}}(T) = U_{\text{OOD}}^{(:r)} T U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}}, \quad T \in \mathcal{O}_r, \quad (5)$$

186 where  $\mathcal{O}_r = \{T \in \mathbb{R}^{r \times r} : T^\top T = I\}$ . Every column of  $\widetilde{W}_{\text{dec}}(T)$  lies in  $\text{span}(U_{\text{OOD}}^{(:r)})$ , so the  
 187 column space of  $\widetilde{W}_{\text{dec}}(T)$  is contained in this  $r$ -dimensional subspace. Since  $\widetilde{W}_{\text{dec}}(T)$  has rank  $r$ , its  
 188 induced explainer subspace equals  $\widehat{\Pi}_{\text{OOD}}$  exactly, and *the faithfulness gap vanishes* ( $\Delta(\Pi_{\text{dec}}) = 0$ )  
 189 *for any*  $T \in \mathcal{O}_r$ . Among all rotations that achieve this alignment, we select the one that keeps the  
 190 rotated dictionary closest to the original:

$$T^* = \arg \min_{T \in \mathcal{O}_r} \|\widetilde{W}_{\text{dec}}(T) - W_{\text{dec}}^{\text{ID}}\|_F^2. \quad (6)$$

191 This is an orthogonal Procrustes problem [37]. Let  $G = U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}} (W_{\text{dec}}^{\text{ID}})^\top U_{\text{OOD}}^{(:r)} \in \mathbb{R}^{r \times r}$ , with  
 192 SVD  $G = \widetilde{U} \Sigma \widetilde{V}^\top$ . Then  $T^* = \widetilde{V} \widetilde{U}^\top$  (derivation in Appendix A.5). Since  $\Pi_{\text{dec}}^{\text{GAE}} := \widehat{\Pi}_{\text{OOD}}$  by  
 193 construction, the residual faithfulness gap reduces to the eigenspace estimation error. This yields a  
 194 quantitative improvement over the unadapted ID explainer.

<sup>2</sup>The full dictionary-based explainer objective (Eq. (1)) includes a sparsity penalty on  $\mathbf{z}_i = \sigma(W_{\text{enc}} h_i + b_{\text{enc}})$ . GAE holds the encoder fixed, so  $\mathbf{z}_i$  is constant with respect to  $W_{\text{dec}}$  and the sparsity term drops out, leaving only the reconstruction term. This also means GAE applies regardless of sparsity mechanism ( $\ell_1$ , Top-K [7], JumpReLU [36], etc.).

195 **Theorem 1** (Improvement over ID Explainer). Suppose  $\gamma_{\text{OOD}} := \lambda_r(M_{\text{OOD}}) - \lambda_{r+1}(M_{\text{OOD}}) > 0$ ,  
 196  $\Delta(\Pi_{\text{ID}}) > 0$ , and  $\widehat{\Pi}_{\text{OOD}} \approx \Pi_{\text{OOD}}$  (holds with sufficient OOD samples). Then

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}^{\text{GAE}}) \leq \mathcal{L}_{\text{OOD}}(\Pi_{\text{ID}}) - \frac{\gamma_{\text{OOD}}}{2} \Delta(\Pi_{\text{ID}})^2. \quad (7)$$

197 The improvement grows quadratically with the ID explainer’s misalignment, so the more severe the  
 198 shift, the larger the guaranteed gain. The proof is given in Appendix A.7.

199 **Step 2: Constrained decoder refit.** Step 1 aligns the subspace but does not optimize sample-level  
 200 reconstruction. Step 2 refits  $\mathcal{L}_{\text{recon}}$  while preserving the alignment from Step 1. Since the encoder  
 201 is fixed, the feature activations  $z_i = \sigma(W_{\text{enc}} h_i + b_{\text{enc}})$  for each OOD sample  $h_i$  are constants, and  
 202 reconstruction reduces to a linear least-squares problem in  $W_{\text{dec}}$  and  $b_{\text{dec}}$ . To keep the decoder  
 203 geometrically aligned, we penalize decoder mass outside  $\widehat{\Pi}_{\text{OOD}}$  with  $\lambda_{\text{geom}} \|(I - \widehat{\Pi}_{\text{OOD}}) W_{\text{dec}}\|_F^2$ .  
 204 To preserve the feature structure from Step 1, we regularize toward  $\widetilde{W}_{\text{dec}}(T^*)$  with  $\lambda_{\text{pres}} \|W_{\text{dec}} -$   
 205  $\widetilde{W}_{\text{dec}}(T^*)\|_F^2$ . The combined objective is convex and quadratic, yielding the closed-form solution

$$W_{\text{dec}}^{\text{GAE}} = \widehat{\Pi}_{\text{OOD}} C B^{-1} + (I - \widehat{\Pi}_{\text{OOD}}) C (B + \lambda_{\text{geom}} I)^{-1}, \quad (8)$$

$$b_{\text{dec}}^{\text{GAE}} = \frac{1}{N} \sum_i h_i - W_{\text{dec}}^{\text{GAE}} \frac{1}{N} \sum_i z_i, \quad (9)$$

207 where

$$B = \frac{1}{N} \sum_i z_i z_i^\top - \left(\frac{1}{N} \sum_i z_i\right) \left(\frac{1}{N} \sum_i z_i\right)^\top + \lambda_{\text{pres}} I, \quad (10)$$

$$C = \frac{1}{N} \sum_i h_i z_i^\top - \left(\frac{1}{N} \sum_i h_i\right) \left(\frac{1}{N} \sum_i z_i\right)^\top + \lambda_{\text{pres}} \widetilde{W}_{\text{dec}}(T^*). \quad (11)$$

208 To summarize, GAE works by first rotating the ID dictionary so that its column space coincides  
 209 with the OOD-active subspace (Step 1), then refitting the decoder via a closed-form ridge regression  
 210 that matches sample-level reconstruction while preserving this alignment (Step 2). The Step 2  
 211 solution applies regularization strength  $\lambda_{\text{pres}}$  to decoder mass inside the OOD-active subspace and  
 212 the larger strength  $\lambda_{\text{pres}} + \lambda_{\text{geom}}$  outside it, so *any decoder mass that drifts off the Step 1 alignment*  
 213 *is automatically shrunk back*. The full derivation is in Appendix A.6.

214 **Applying GAE at inference.** At inference, the ID-trained encoder remains unchanged: given  
 215 an OOD activation  $h$ , the explainer extracts features  $z = \sigma(W_{\text{enc}} h + b_{\text{enc}})$  and reconstructs  $\hat{h} =$   
 216  $W_{\text{dec}}^{\text{GAE}} z + b_{\text{dec}}^{\text{GAE}}$ . This applies identically to both SAEs and transcoders.

## 217 5 Experiments

218 We evaluate whether GAE restores explanation faithfulness under distribution shift. Section 5.1 tests  
 219 the geometric mechanism in a controlled setting, and Section 5.2 evaluates on language models.

### 220 5.1 Controlled Experiment

221 We first test whether the geometric mech-  
 222 anism from Section 3 holds in a controlled  
 223 setting. We train a 2-layer ReLU MLP with  
 224 hidden dim  $d = 256$  and output dim  $p = 8$ ,  
 225 and a linear-decoder SAE on its ID hid-  
 226 den activations, then continuously increase  
 227 OOD severity by rotating and rescaling the  
 228 input covariance (details in Appendix B.1).  
 229 Figure 2 confirms that as severity increases,  
 230 the faithfulness gap and reconstruction error  
 231 of the Fixed explainer (the ID-trained  
 232 explainer, used without adaptation) grow,  
 233 while GAE closes the gap to near zero and  
 234 keeps reconstruction error nearly flat.

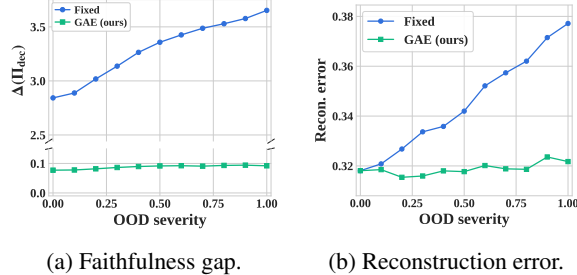


Figure 2: **Controlled experiment on a toy MLP with OOD severity varied from 0 (ID) to 1 (maximum shift).** (a) The Fixed explainer’s faithfulness gap  $\Delta(\Pi_{\text{dec}})$  grows monotonically. (b) Its reconstruction error rises accordingly. GAE maintains near-zero gap and flat error throughout.

### 235 5.2 Experiments on Language Models

#### 236 5.2.1 Setup

237 **Target models and explainers.** We evaluate on two frozen pretrained language models: GPT-2  
 238 Small [38] and Pythia-1.4B [39]. For each model, we train transcoders [5] and Top-K SAEs [7], both  
 239 with dictionary size  $k=32d$ , as dictionary-based explainers (Eq. (1)).

240 **OOD settings.** We consider three categories of distribution shift: **temporal** (FineWeb [40], web  
 241 text collected after each model’s pretraining cutoff), **domain** (Edgar [41], financial filings whose  
 242 specialized vocabulary and structure differ from general web text), and **adversarial** (HaluEval [42],  
 243 hallucination-inducing prompts that elicit atypical hidden representations). All three induce measur-  
 244 able second-moment shift in hidden activations, as verified in Appendix B.4.

245 **Baselines.** We compare training-free and training-based approaches. **Fixed** applies the ID-trained  
 246 explainer without adaptation. **TERM** [16, 17] trains the ID explainer with tilted ERM to upweight  
 247 tail samples. Among training-based methods, **Finetune** [43] warm-starts from the ID explainer on  
 248 OOD activations, **Retrain** trains from scratch on OOD data, **SAEBoost** [18] adds a residual booster  
 249 on OOD reconstruction residuals, and **FaithfulSAE** [15] retrains on the model’s own generations.  
 250 **GAE (ours)** is training-free: it uses only unlabeled OOD activations with no gradient computation.  
 251 Detailed descriptions are in Appendix D.2.

252 As shown in Table 1, all training-based base-  
 253 lines require gradient-based optimization over  
 254 millions of tokens. Even the lightest, Finetune,  
 255 processes 5M tokens over several minutes; Re-  
 256 train, SAEBoost, and FaithfulSAE each con-  
 257 sume 100M tokens and take hours on a single  
 258 GPU. *GAE requires no gradient computation at*  
 259 *all: the entire closed-form pipeline completes in*  
 260 *0.5 s for GPT-2 and 2.9 s for Pythia-1.4B, using*  
 261 *only ~2,048 unlabeled OOD activations. This*  
 262 *makes GAE practical for on-the-fly adaptation*  
 263 *whenever the deployment distribution changes.*

Table 1: **Computational cost per method.**  
 Training-based baselines require gradient optimization over millions of tokens. GAE adapts in under 3 s without training.

Method	Tokens	Wall-clock	
		GPT-2	P-1.4B
Finetune	5M	~2 min	~12 min
Retrain	100M	~39 min	~4 hrs
SAEBoost	100M	~39 min	~4 hrs
FaithfulSAE	100M	~39 min	~4 hrs
<b>GAE</b>	<b>2K</b>	<b>0.5 s</b>	<b>2.9 s</b>

264 **Evaluation metrics.** We evaluate causal faithfulness using three metrics. **Normalized AOPC**  
 265 **(nAOPC)** [23] averages the normalized logit drop across multiple feature budgets when top- $m$   
 266 features are removed ( $\uparrow$  is better). **Normalized comprehensiveness (nComp)** [24] measures the  
 267 normalized logit drop at a single budget  $m^*=32$  ( $\uparrow$  is better). Both measure the logit-level effect of  
 268 ablating top features. **Delta cross-entropy ( $\Delta$ CE)** [7] measures reconstruction quality: the cross-  
 269 entropy change when activations are replaced with the explainer’s reconstruction ( $\approx 0$  is better).  
 270 GAE optimizes a geometric objective (the faithfulness gap); improvements on these causal metrics  
 271 confirm that geometric realignment yields faithfulness gains. Formal definitions are in Appendix D.4.

272 **5.2.2 Faithfulness Results**

Table 2: **Faithfulness under distribution shift (Transcoder, two models  $\times$  three OOD settings).**  
 GAE is training-free yet leads in all nine columns on GPT-2 and achieves the best nComp and  $|\Delta$ CE|  
 in 5 of 6 Pythia-1.4B columns. **Bold:** best per column. Underline: second best.

Method	FineWeb (Temporal)			Edgar (Domain)			HaluEval (Adversarial)		
	nAOPC $\uparrow$	nComp $\uparrow$	$ \Delta$ CE  $\downarrow$	nAOPC $\uparrow$	nComp $\uparrow$	$ \Delta$ CE  $\downarrow$	nAOPC $\uparrow$	nComp $\uparrow$	$ \Delta$ CE  $\downarrow$
<i>GPT-2 Small</i>									
Fixed	0.857	1.017	0.0281	0.975	1.025	0.0201	0.735	0.737	0.0473
TERM	0.853	0.993	0.0283	0.964	0.999	0.0198	0.730	0.732	0.0523
Finetune	0.856	0.984	<u>0.0172</u>	0.971	1.117	0.0047	0.579	0.579	<u>0.0105</u>
Retrain	0.895	1.118	0.0218	0.936	1.034	<u>0.0015</u>	0.521	0.528	0.2763
SAEBoost	<u>0.958</u>	<u>1.476</u>	0.0177	<u>0.979</u>	<u>1.542</u>	0.0072	<u>0.840</u>	<u>0.921</u>	0.0212
FaithfulSAE	0.936	<u>1.186</u>	0.0213	0.976	<u>1.134</u>	0.0197	0.738	0.740	0.0586
<b>GAE (ours)</b>	<b>0.960</b>	<b>1.494</b>	<b>0.0167</b>	<b>0.981</b>	<b>1.618</b>	<b>0.0009</b>	<b>0.871</b>	<b>0.963</b>	<b>0.0014</b>
<i>Pythia-1.4B</i>									
Fixed	0.839	1.091	0.0278	0.725	0.828	0.0300	0.899	1.021	0.0354
TERM	0.845	1.043	0.0271	0.895	0.981	0.0298	0.896	1.104	0.0349
Finetune	0.859	1.249	<b>0.0264</b>	0.684	0.860	<u>0.0282</u>	0.925	1.502	<u>0.0280</u>
Retrain	0.894	<u>1.315</u>	0.0405	0.746	0.821	0.0329	0.894	1.393	0.0305
SAEBoost	<u>0.908</u>	1.296	0.0284	<u>0.903</u>	<u>1.297</u>	0.0296	<u>0.965</u>	<u>1.530</u>	0.0283
FaithfulSAE	0.858	1.056	0.0272	0.724	0.822	0.0323	0.899	1.171	0.0307
<b>GAE (ours)</b>	<b>0.915</b>	<b>1.354</b>	<u>0.0269</u>	<b>0.988</b>	<b>1.652</b>	<b>0.0230</b>	<b>0.968</b>	<b>1.693</b>	<b>0.0276</b>

273 Table 2 reports faithfulness across three OOD settings for GPT-2 Small and Pythia-1.4B (Transcoder).  
 274 Despite using no gradient updates, *GAE leads on all three metrics for GPT-2 across every OOD*

275 *setting*, surpassing training-based baselines that consume up to 100M tokens (cf. Table 1). The largest  
 276 gains appear on adversarial shift, where GAE improves nComp over the strongest training-based  
 277 baseline SAEBoost by 4.6% (0.963 vs 0.921) and reduces  $|\Delta\text{CE}|$  by 93% (0.0014 vs 0.0212). On  
 278 Pythia-1.4B, GAE achieves the best nAOPC and nComp on all three settings;  $|\Delta\text{CE}|$  is best on Edgar  
 279 and HaluEval but slightly elevated on FineWeb (0.027 vs Finetune’s 0.026), where the more diffuse  
 280 eigenvalue decay weakens the rank- $r$  approximation.

Table 3: **Faithfulness under distribution shift (SAE, two models  $\times$  three OOD settings)**. GAE leads in all nine columns on GPT-2 and leads on nComp and  $|\Delta\text{CE}|$  in 5 of 6 Pythia-1.4B columns. **Bold**: best. Underline: second best.

Method	FineWeb (Temporal)			Edgar (Domain)			HaluEval (Adversarial)		
	nAOPC $\uparrow$	nComp $\uparrow$	$ \Delta\text{CE} \downarrow$	nAOPC $\uparrow$	nComp $\uparrow$	$ \Delta\text{CE} \downarrow$	nAOPC $\uparrow$	nComp $\uparrow$	$ \Delta\text{CE} \downarrow$
<b>GPT-2 Small</b>									
Fixed	0.735	0.796	0.0185	0.650	0.667	0.0089	0.930	1.134	0.0406
TERM	0.741	0.790	0.0183	0.655	0.682	0.0145	0.898	0.987	0.0401
Finetune	0.725	0.802	<u>0.0015</u>	0.658	0.687	0.0068	<u>0.932</u>	1.155	<u>0.0218</u>
Retrain	<u>0.766</u>	<u>0.856</u>	0.0375	0.715	0.797	<u>0.0065</u>	0.930	1.276	0.0300
SAEBoost	0.704	<u>0.786</u>	0.0120	0.604	0.650	0.0098	0.909	1.243	0.0448
FaithfulSAE	0.725	0.760	0.0262	0.657	0.683	0.0074	0.908	1.082	0.0278
<b>GAE (ours)</b>	<b>0.768</b>	<b>0.871</b>	<b>0.0011</b>	<b>0.723</b>	<b>0.809</b>	<b>0.0037</b>	<b>0.953</b>	<b>1.303</b>	<b>0.0017</b>
<b>Pythia-1.4B</b>									
Fixed	0.962	1.536	0.0216	0.953	1.690	0.0170	0.985	1.425	0.0252
TERM	0.965	1.486	0.0237	0.959	1.787	0.0167	0.984	1.762	0.0292
Finetune	0.963	1.618	<u>0.0216</u>	0.959	1.832	<u>0.0131</u>	<u>0.988</u>	<u>1.880</u>	0.0174
Retrain	<u>0.983</u>	<b>1.681</b>	0.0563	<b>0.970</b>	<u>1.848</u>	0.0261	<b>1.000</b>	1.855	0.0209
SAEBoost	0.971	1.451	0.0211	0.955	1.598	0.0152	<b>1.000</b>	1.830	<u>0.0166</u>
FaithfulSAE	0.982	1.642	0.0307	0.953	1.769	0.0205	<b>1.000</b>	1.678	0.0514
<b>GAE (ours)</b>	<b>0.985</b>	<u>1.677</u>	<b>0.0207</b>	<u>0.968</u>	<b>1.946</b>	<b>0.0098</b>	<b>1.000</b>	<b>1.885</b>	<b>0.0163</b>

281 Table 3 reports the same evaluation for Top-K SAEs. On GPT-2, GAE again surpasses every training-  
 282 based baseline on all nine columns, with the largest margins on Edgar (nComp 0.809 vs Retrain’s  
 283 0.797) and HaluEval ( $|\Delta\text{CE}|$  0.0017 vs Finetune’s 0.0218). On Pythia-1.4B, Retrain is a stronger  
 284 competitor than for transcoders, taking best nComp on FineWeb (1.681 vs 1.677) and best nAOPC on  
 285 Edgar (0.970 vs 0.968). GAE still leads on nComp and  $|\Delta\text{CE}|$  in 5 of 6 columns and achieves the  
 286 best  $|\Delta\text{CE}|$  across all three settings. The pattern is consistent: *GAE matches or surpasses methods  
 287 that require orders of magnitude more computation.*

### 288 5.2.3 Case Study: Circuit Attribution under Distribution Shift

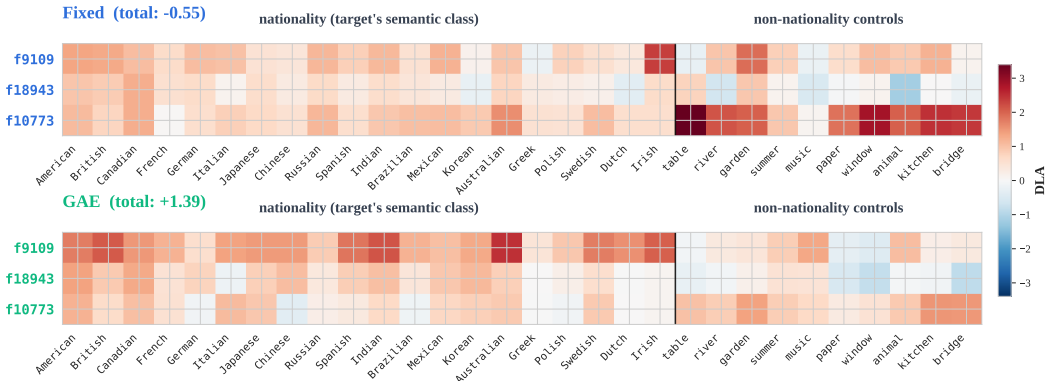


Figure 3: **Per-feature DLA on a prompt predicting ‘American’ (GPT-2, Transcoder)**. Both methods share the same encoder and top-3 features; only the decoder columns differ. Each cell shows a feature’s direct logit attribution (DLA) to nationality tokens (left, 20 tokens) vs. non-nationality controls (right, 10 tokens). Fixed’s total class-specificity is  $-0.55$  (circuit points away from the target class); GAE’s is  $+1.39$  (circuit points toward it).

289 We select a prompt where the model predicts ‘American’ and measure each feature’s direct logit  
 290 attribution (DLA), the dot product of its decoder column with the token’s unembedding vector scaled  
 291 by the feature activation. DLA quantifies how much each feature pushes the model toward a given

292 next token. Since GAE keeps the encoder frozen, both Fixed and GAE extract the same top-3 features  
 293 with the same activations, so any change in DLA isolates the effect of the decoder rotation. For each  
 294 feature, we compute the difference between its mean DLA on 20 nationality tokens (the target’s  
 295 semantic class) and 10 non-nationality controls, then sum across features to obtain a class-specificity  
 296 score that measures whether the identified circuit points toward the correct token class.

297 Fixed scores  $-0.55$ : its circuit points away from the target class on average. GAE scores  $+1.39$ :  
 298 every top feature contributes more to nationalities than to controls, as shown in Figure 3. The  
 299 decoder rotation alone corrects feature-level attribution without altering which features are selected.  
 300 Appendix E repeats this analysis on two further prompts whose target tokens are a male first name  
 301 ( $+1.00 \rightarrow +4.51$ ) and a profession ( $+0.54 \rightarrow +0.99$ ).

### 302 5.2.4 Mechanism Analysis

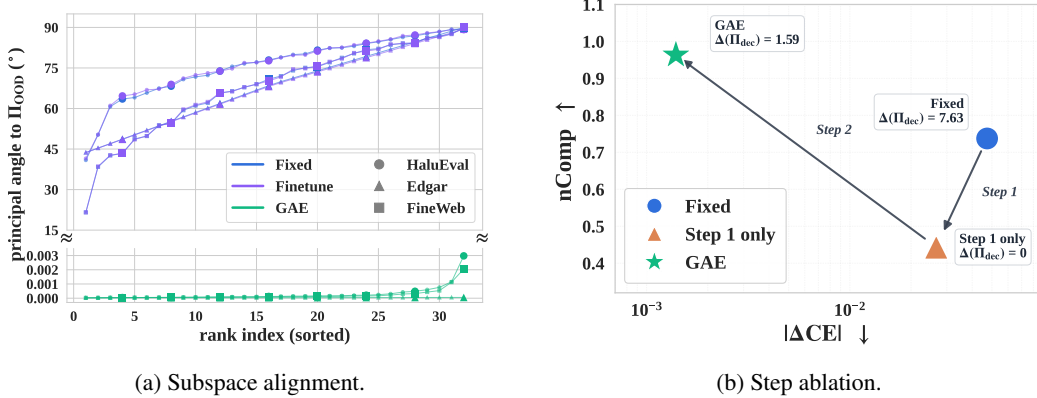


Figure 4: **Mechanism analysis (GPT-2, Transcoder).** (a) Sorted principal angles between each explainer’s top- $r$  subspace and  $\hat{\Pi}_{\text{OOD}}$ . GAE’s subspace aligns with  $\hat{\Pi}_{\text{OOD}}$ , while Fixed and Finetune leave large angular gaps. (b) Step ablation: Step 1 closes the faithfulness gap to 0 yet drops nComp from 0.74 to 0.44. Step 2 restores nComp to 0.96 at the cost of a small gap (1.59).

303 **Subspace alignment.** Figure 4(a) measures the principal angles between each explainer’s top- $r$   
 304 decoder subspace and the OOD-active subspace  $\hat{\Pi}_{\text{OOD}}$ . Fixed and Finetune leave  $40^\circ$ – $90^\circ$  angular  
 305 gaps across all rank indices and OOD shift types, while GAE drives every angle to  $\sim 10^{-3}^\circ$ . This  
 306 confirms that GAE’s faithfulness gains arise from explicit geometric alignment, validating the  
 307 mechanism of Section 4. The persistence of Finetune’s gap further shows that *gradient-based*  
 308 *reconstruction loss does not, by itself, drive subspace alignment with  $\hat{\Pi}_{\text{OOD}}$* . Appendix C verifies  
 309 that the projection-loss improvement scales quadratically with  $\Delta(\Pi_{\text{ID}})^2$ , as predicted by Theorem 1.

310 **Step ablation.** Figure 4(b) ablates each step of GAE on HaluEval. Step 1 alone closes the faithfulness  
 311 gap from 7.63 to 0 by construction, yet nComp drops from 0.74 to 0.44 since the orthogonal  
 312 rotation diffuses the encoder-decoder feature pairing that the top- $k$  ablation in nComp measures.  
 313 Step 2 refits the decoder within Step 1’s subspace, accepting a small gap (1.59) in exchange for the  
 314 highest nComp (0.96) and the lowest  $|\Delta\text{CE}|$  (0.0014). *The two steps are complementary: Step 1*  
 315 *chooses the subspace, Step 2 makes the dictionary causally coherent within it.* A hyperparameter  
 316 sensitivity analysis is reported in Appendix F.

## 317 6 Conclusion

318 We showed that OOD faithfulness degradation in dictionary-based explainers has a geometric cause:  
 319 the decoder subspace drifts from the directions the model actively uses. The faithfulness gap  $\Delta(\Pi_{\text{dec}})$   
 320 formalizes this misalignment and provably controls the reducible part of OOD faithfulness loss.  
 321 GAE closes the gap with a closed-form subspace rotation and constrained decoder refit, using only  
 322 unlabeled OOD activations. Across two models and three shift types, GAE outperforms all training-  
 323 based baselines on 5 of 6 settings, completing in under 3 seconds without any gradient computation.  
 324 A limitation is that we have not yet evaluated on larger-scale models. GAE also relies on a top- $r$   
 325 SVD truncation of  $W_{\text{dec}}^{\text{ID}}$ , so any feature information carried in the residual  $(d-r)$  singular directions  
 326 is dropped before adaptation. Extending GAE to adaptive rank selection, encoder adaptation, and  
 327 connections to optimal transport on the Grassmannian [44] are promising future directions.

328 **References**

- 329 [1] Atticus Geiger, Duligur Ibeling, Amir Zur, Maheep Chaudhary, Sonakshi Chauhan, Jing Huang,  
330 Aryaman Arora, Zhengxuan Wu, Noah Goodman, Christopher Potts, et al. Causal abstraction: A  
331 theoretical foundation for mechanistic interpretability. Journal of Machine Learning Research,  
332 26(83):1–64, 2025.
- 333 [2] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom  
334 Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning  
335 and induction heads. arXiv preprint arXiv:2209.11895, 2022.
- 336 [3] Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas  
337 Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, et al. Open problems  
338 in mechanistic interpretability. arXiv preprint arXiv:2501.16496, 2025.
- 339 [4] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Con-  
340 erly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu,  
341 Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex  
342 Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter,  
343 Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language  
344 models with dictionary learning. Transformer Circuits Thread, 2023. [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)  
345 [circuits.pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 346 [5] Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable llm feature  
347 circuits. Advances in Neural Information Processing Systems, 37:24375–24410, 2024.
- 348 [6] Adly Templeton. Scaling monosemanticity: Extracting interpretable features from claude 3  
349 sonnet. Anthropic, 2024.
- 350 [7] Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya  
351 Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. arXiv  
352 preprint arXiv:2406.04093, 2024.
- 353 [8] Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.  
354 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models.  
355 arXiv preprint arXiv:2403.19647, 2024.
- 356 [9] Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable nlp systems: How should we  
357 define and evaluate faithfulness? arXiv preprint arXiv:2004.03685, 2020.
- 358 [10] Google DeepMind Safety Research. Negative results for sparse autoencoders on downstream  
359 tasks and deprioritising sae research. DeepMind Safety Research Blog, 2025. Blog post.
- 360 [11] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim.  
361 Sanity checks for saliency maps. Advances in neural information processing systems, 31, 2018.
- 362 [12] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile.  
363 In Proceedings of the AAAI conference on artificial intelligence, volume 33, pages 3681–3688,  
364 2019.
- 365 [13] Chris Lin, Ian Covert, and Su-In Lee. On the robustness of removal-based feature attributions.  
366 Advances in Neural Information Processing Systems, 36:79613–79666, 2023.
- 367 [14] Chiara Balestra, Bin Li, and Emmanuel Müller. On the consistency and robustness of saliency  
368 explanations for time series classification. arXiv preprint arXiv:2309.01457, 2023.
- 369 [15] Seonglae Cho, Harryn Oh, Donghyun Lee, Luis Rodrigues Vieira, Andrew Birmingham, and  
370 Ziad El Sayed. Faithfulsae: Towards capturing faithful features with sparse autoencoders without  
371 external datasets dependency. In Proceedings of the 63rd Annual Meeting of the Association  
372 for Computational Linguistics (Volume 4: Student Research Workshop), pages 297–314, 2025.
- 373 [16] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization.  
374 arXiv preprint arXiv:2007.01162, 2020.

- 375 [17] Aashiq Muhamed, Mona Diab, and Virginia Smith. Decoding dark matter: Specialized sparse  
376 autoencoders for interpreting rare concepts in foundation models. In Findings of the Association  
377 for Computational Linguistics: NAACL 2025, pages 1604–1635, 2025.
- 378 [18] Nikita Koriagin, Yaroslav Aksenov, Daniil Laptev, Gleb Gerasimov, Nikita Balagansky,  
379 and Daniil Gavrilov. Teach old saes new domain tricks with boosting. arXiv preprint  
380 arXiv:2507.12990, 2025.
- 381 [19] Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety—a review.  
382 arXiv preprint arXiv:2404.14082, 2024.
- 383 [20] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for  
384 detecting out-of-distribution samples and adversarial attacks. Advances in neural information  
385 processing systems, 31, 2018.
- 386 [21] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-  
387 coders find highly interpretable features in language models. arXiv preprint arXiv:2309.08600,  
388 2023.
- 389 [22] Gonçalo Paulo, Stepan Shabalín, and Nora Belrose. Transcoders beat sparse autoencoders for  
390 interpretability. arXiv preprint arXiv:2501.18823, 2025.
- 391 [23] Joakim Edin, Andreas Geert Motzfeldt, Casper L Christensen, Tuukka Ruotsalo, Lars Maaløe,  
392 and Maria Maistro. Normalized aopc: Fixing misleading faithfulness metrics for feature  
393 attributions explainability. In Proceedings of the 63rd Annual Meeting of the Association for  
394 Computational Linguistics (Volume 1: Long Papers), pages 1715–1730, 2025.
- 395 [24] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard  
396 Socher, and Byron C Wallace. Eraser: A benchmark to evaluate rationalized nlp models. In  
397 Proceedings of the 58th annual meeting of the association for computational linguistics, pages  
398 4443–4458, 2020.
- 399 [25] Lawrence Chan, Adria Garriga-Alonso, Nicholas Goldowsky-Dill, Ryan Greenblatt, Jenny  
400 Nitishinskaya, Ansh Radhakrishnan, Buck Shlegeris, and Nate Thomas. Causal scrubbing: A  
401 method for rigorously testing interpretability hypotheses. In AI Alignment Forum, volume 2,  
402 2022.
- 403 [26] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
404 associations in gpt. Advances in neural information processing systems, 35:17359–17372,  
405 2022.
- 406 [27] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of  
407 neural network representations revisited. In International conference on machine learning,  
408 pages 3519–3529. PMLR, 2019.
- 409 [28] Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular  
410 vector canonical correlation analysis for deep learning dynamics and interpretability. Advances  
411 in neural information processing systems, 30, 2017.
- 412 [29] Alex H Williams, Erin Kunz, Simon Kornblith, and Scott Linderman. Generalized shape metrics  
413 on neural representations. Advances in neural information processing systems, 34:4738–4750,  
414 2021.
- 415 [30] Tom Lieberum, Senthoooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat,  
416 Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope:  
417 Open sparse autoencoders everywhere all at once on gemma 2. In Proceedings of the 7th  
418 BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP, pages 278–  
419 300, 2024.
- 420 [31] Charles H Martin, Tongsu Peng, and Michael W Mahoney. Predicting trends in the qual-  
421 ity of state-of-the-art neural networks without access to training or testing data. Nature  
422 Communications, 12(1):4122, 2021.

- 423 [32] Alessio Ansuini, Alessandro Laio, Jakob H Macke, and Davide Zoccolan. Intrinsic dimension  
424 of data representations in deep neural networks. Advances in Neural Information Processing  
425 Systems, 32, 2019.
- 426 [33] Armen Aghajanyan, Sonal Gupta, and Luke Zettlemoyer. Intrinsic dimensionality explains the  
427 effectiveness of language model fine-tuning. In Proceedings of the 59th annual meeting of the  
428 association for computational linguistics and the 11th international joint conference on natural  
429 language processing (volume 1: long papers), pages 7319–7328, 2021.
- 430 [34] Callum McDougall, Arthur Conmy, János Kramár, Tom Lieberum, Senthoran Rajamanoharan,  
431 and Neel Nanda. Gemma scope 2: Technical paper. Technical report, Google DeepMind, 2025.
- 432 [35] Chandler Davis and William Morton Kahan. The rotation of eigenvectors by a perturbation. iii.  
433 SIAM Journal on Numerical Analysis, 7(1):1–46, 1970.
- 434 [36] Senthoran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma,  
435 János Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu  
436 sparse autoencoders. arXiv preprint arXiv:2407.14435, 2024.
- 437 [37] Peter H Schönemann. A generalized solution of the orthogonal procrustes problem.  
438 Psychometrika, 31(1):1–10, 1966.
- 439 [38] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.  
440 Language models are unsupervised multitask learners. OpenAI blog, 1(8):9, 2019.
- 441 [39] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien,  
442 Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward  
443 Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In  
444 International Conference on Machine Learning, pages 2397–2430. PMLR, 2023.
- 445 [40] Guilherme Penedo, Hynek Kydlíček, Loubna Ben Allal, and Thomas Wolf. Fineweb: decanting  
446 the web for the finest text data at scale. HuggingFace. Accessed: Jul, 12, 2024.
- 447 [41] Lefteris Loukas, Manos Fergadiotis, Ion Androutsopoulos, and Prodromos Malakasiotis. Edgar-  
448 corpus: Billions of tokens make the world go round. In Proceedings of the Third Workshop on  
449 Economics and Natural Language Processing, pages 13–18, 2021.
- 450 [42] Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. Halueval: A  
451 large-scale hallucination evaluation benchmark for large language models. In Proceedings of  
452 the 2023 conference on empirical methods in natural language processing, pages 6449–6464,  
453 2023.
- 454 [43] Connor Kissane, Robert Krzyzanowski, Arthur Conmy, and Neel Nanda. Saes (usually) transfer between base and chat models. Alignment Forum,  
455 2024. URL <https://www.alignmentforum.org/posts/fmwk6qxrpw8d4jvbd/saes-usually-transfer-between-base-and-chat-models>.  
456  
457
- 458 [44] Alan Edelman, Tomás A Arias, and Steven T Smith. The geometry of algorithms with or-  
459 thogonality constraints. SIAM journal on Matrix Analysis and Applications, 20(2):303–353,  
460 1998.
- 461 [45] Yi Yu, Tengyao Wang, and Richard J Samworth. A useful variant of the davis–kahan theorem  
462 for statisticians. Biometrika, 102(2):315–323, 2015.
- 463 [46] Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. [http://Skyllion007.github.io/](http://Skyllion007.github.io/OpenWebTextCorpus)  
464 [OpenWebTextCorpus](http://Skyllion007.github.io/OpenWebTextCorpus), 2019.
- 465 [47] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason  
466 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse  
467 text for language modeling. arXiv preprint arXiv:2101.00027, 2020.
- 468 [48] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-  
469 tuning can distort pretrained features and underperform out-of-distribution. arXiv preprint  
470 arXiv:2202.10054, 2022.

## 471 A Proofs and Derivations

### 472 A.1 Proof of Proposition 1

473 **Setup.** Write the second-moment shift as  $E = M_{\text{OOD}} - M_{\text{ID}}$ , so that  $M_{\text{OOD}} = M_{\text{ID}} + E$ . The  
 474 projectors  $\Pi_{\text{ID}}$  and  $\Pi_{\text{OOD}}$  correspond to the top- $r$  eigenspaces of  $M_{\text{ID}}$  and  $M_{\text{ID}} + E$ , respectively.  
 475 We wish to bound  $\Delta(\Pi_{\text{ID}}) = \|\Pi_{\text{OOD}} - \Pi_{\text{ID}}\|_F$  in terms of  $\|E\|_F$ .

476 **Step 1: Projector distance via principal angles.** Let  $\theta_1, \dots, \theta_r$  be the principal angles between  
 477 the column spaces of  $\Pi_{\text{ID}}$  and  $\Pi_{\text{OOD}}$ . For two rank- $r$  orthogonal projectors, the Frobenius norm of  
 478 their difference satisfies

$$\|\Pi_{\text{OOD}} - \Pi_{\text{ID}}\|_F^2 = 2 \sum_{i=1}^r \sin^2 \theta_i.$$

479 **Step 2: Applying Davis–Kahan.** The Frobenius-norm form of the Davis–Kahan  $\sin \Theta$  theorem [35,  
 480 45] bounds the sum of squared sines in terms of the perturbation  $E$  and the eigengap:

$$\sum_{i=1}^r \sin^2 \theta_i \leq \frac{\|(I - \Pi_{\text{ID}}) E \Pi_{\text{ID}}\|_F^2}{\gamma_{\text{ID}}^2},$$

481 where  $\gamma_{\text{ID}} = \lambda_r(M_{\text{ID}}) - \lambda_{r+1}(M_{\text{ID}})$  is the eigengap of  $M_{\text{ID}}$  at rank  $r$ .

482 **Step 3: Bounding the cross term.** The matrix  $(I - \Pi_{\text{ID}}) E \Pi_{\text{ID}}$  is the component of  $E$  that maps  
 483 the ID-active subspace into its orthogonal complement. Since  $I - \Pi_{\text{ID}}$  and  $\Pi_{\text{ID}}$  are orthogonal  
 484 projectors, each has operator norm 1, so by the submultiplicativity of the Frobenius norm under  
 485 operator-norm factors,

$$\|(I - \Pi_{\text{ID}}) E \Pi_{\text{ID}}\|_F \leq \|I - \Pi_{\text{ID}}\|_2 \cdot \|E\|_F \cdot \|\Pi_{\text{ID}}\|_2 = 1 \cdot \|E\|_F \cdot 1 = \|E\|_F.$$

486 **Step 4: Combining.** Substituting back,

$$\Delta(\Pi_{\text{ID}}) = \|\Pi_{\text{OOD}} - \Pi_{\text{ID}}\|_F = \sqrt{2 \sum_{i=1}^r \sin^2 \theta_i} \leq \frac{\sqrt{2}}{\gamma_{\text{ID}}} \|E\|_F = \frac{\sqrt{2}}{\gamma_{\text{ID}}} \|M_{\text{OOD}} - M_{\text{ID}}\|_F,$$

487 where the first equality is Step 1, the inequality combines Steps 2 and 3, and the last equality uses  
 488  $E = M_{\text{OOD}} - M_{\text{ID}}$  from the Setup.  $\square$

### 489 A.2 Decomposition Eq. (3)

490 **Step 1: Expressing  $\mathcal{L}_{\text{OOD}}$  in terms of  $M_{\text{OOD}}$ .** For any rank- $r$  orthogonal  $\Pi_{\text{dec}} \in \mathcal{C}_r$ , the recon-  
 491 struction error under  $P_{\text{OOD}}$  is

$$\begin{aligned} \mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}) &= \mathbb{E}_{X \sim P_{\text{OOD}}} \|h(X) - \Pi_{\text{dec}} h(X)\|_2^2 \\ &= \mathbb{E} \|(I - \Pi_{\text{dec}}) h(X)\|_2^2 \\ &= \mathbb{E} \text{tr}[(I - \Pi_{\text{dec}}) h(X) h(X)^\top (I - \Pi_{\text{dec}})^\top] \\ &= \text{tr}[(I - \Pi_{\text{dec}}) M_{\text{OOD}} (I - \Pi_{\text{dec}})], \end{aligned} \quad (12)$$

492 where the third step uses  $\|a\|_2^2 = \text{tr}(aa^\top)$  and the fourth step swaps expectation and trace, with  
 493  $M_{\text{OOD}} = \mathbb{E}[h(X)h(X)^\top]$ . Since  $I - \Pi_{\text{dec}}$  is an orthogonal projector, it is idempotent:  $(I - \Pi_{\text{dec}})^2 =$   
 494  $I - \Pi_{\text{dec}}$ . Using the cyclic property of trace,

$$\text{tr}[(I - \Pi_{\text{dec}}) M_{\text{OOD}} (I - \Pi_{\text{dec}})] = \text{tr}[(I - \Pi_{\text{dec}})^2 M_{\text{OOD}}] = \text{tr}[(I - \Pi_{\text{dec}}) M_{\text{OOD}}]. \quad (13)$$

495 **Step 2: Deriving the decomposition.** Applying the same identity with  $\Pi_{\text{OOD}}$  gives  
 496  $\mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) = \text{tr}[(I - \Pi_{\text{OOD}}) M_{\text{OOD}}]$ . Taking the difference,

$$\begin{aligned} \mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}) - \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) &= \text{tr}[(I - \Pi_{\text{dec}}) M_{\text{OOD}}] - \text{tr}[(I - \Pi_{\text{OOD}}) M_{\text{OOD}}] \\ &= \text{tr}[(\Pi_{\text{OOD}} - \Pi_{\text{dec}}) M_{\text{OOD}}]. \end{aligned} \quad (14)$$

497 Rearranging gives the claimed decomposition:

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}) = \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) + \text{tr}[(\Pi_{\text{OOD}} - \Pi_{\text{dec}}) M_{\text{OOD}}].$$

498 **Step 3: Nonnegativity and optimality.** It remains to show that the second term is nonnegative. Let  
 499  $M_{\text{OOD}} = \sum_{i=1}^d \lambda_i u_i u_i^\top$  be the eigendecomposition with  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ . Since  $\Pi_{\text{OOD}}$  projects  
 500 onto the top- $r$  eigenspace,

$$\text{tr}(\Pi_{\text{OOD}} M_{\text{OOD}}) = \sum_{i=1}^r \lambda_i.$$

501 By Ky Fan's maximum principle,  $\sum_{i=1}^r \lambda_i = \max_{\Pi_{\text{dec}} \in \mathcal{C}_r} \text{tr}(\Pi_{\text{dec}} M_{\text{OOD}})$ . Therefore, for any  
 502  $\Pi_{\text{dec}} \in \mathcal{C}_r$ ,

$$\text{tr}[(\Pi_{\text{OOD}} - \Pi_{\text{dec}}) M_{\text{OOD}}] = \text{tr}(\Pi_{\text{OOD}} M_{\text{OOD}}) - \text{tr}(\Pi_{\text{dec}} M_{\text{OOD}}) \geq 0,$$

503 with equality if and only if  $\Pi_{\text{dec}}$  also projects onto a top- $r$  eigenspace of  $M_{\text{OOD}}$ . This proves both  
 504 the nonnegativity and the optimality  $\Pi_{\text{OOD}} \in \arg \min_{\Pi_{\text{dec}} \in \mathcal{C}_r} \mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}})$ .  $\square$

### 505 A.3 Proof of Proposition 2

506 From the decomposition (3), the explainer-dependent component equals  $\text{tr}[(\Pi_{\text{OOD}} - \Pi_{\text{dec}}) M_{\text{OOD}}]$ .  
 507 We derive both bounds by expanding this trace in the eigenbasis of  $M_{\text{OOD}}$ .

508 **Setup.** Let  $M_{\text{OOD}} = \sum_{i=1}^d \lambda_i u_i u_i^\top$  be the eigendecomposition with  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ . The  
 509 OOD-active subspace is  $\Pi_{\text{OOD}} = \sum_{i=1}^r u_i u_i^\top$ , so  $u_i^\top \Pi_{\text{OOD}} u_i = \mathbf{1}_{i \leq r}$ . For the explainer subspace  
 510  $\Pi_{\text{dec}} \in \mathcal{C}_r$ , define  $p_i := u_i^\top \Pi_{\text{dec}} u_i \in [0, 1]$ , the fraction of the  $i$ -th OOD eigendirection captured by  
 511  $\Pi_{\text{dec}}$ . Expanding the trace in this eigenbasis gives

$$\begin{aligned} \text{tr}[(\Pi_{\text{OOD}} - \Pi_{\text{dec}}) M_{\text{OOD}}] &= \sum_{i=1}^d \lambda_i (u_i^\top \Pi_{\text{OOD}} u_i - p_i) \\ &= \sum_{i=1}^r \lambda_i (1 - p_i) - \sum_{i=r+1}^d \lambda_i p_i. \end{aligned} \quad (15)$$

512 The first sum captures the OOD energy in the top- $r$  directions that the explainer misses (since  $1 - p_i$   
 513 is the fraction lost). The second sum captures the OOD energy in the bottom directions that the  
 514 explainer unnecessarily covers.

515 **Connecting to the faithfulness gap.** For two rank- $r$  subspaces, the Frobenius norm of their  
 516 difference satisfies  $\|\Pi_{\text{OOD}} - \Pi_{\text{dec}}\|_F^2 = 2 \sum_{i=1}^r (1 - p_i)$ . Moreover, since both have the same rank,  
 517  $\sum_{i=1}^r (1 - p_i) = \sum_{i=r+1}^d p_i$ . Denoting  $S := \sum_{i=1}^r (1 - p_i)$ , we have

$$\Delta(\Pi_{\text{dec}})^2 = \|\Pi_{\text{OOD}} - \Pi_{\text{dec}}\|_F^2 = 2S, \quad \text{and} \quad \sum_{i=r+1}^d p_i = S. \quad (16)$$

518 **Upper bound.** Using  $\lambda_i \leq \lambda_1$  for  $i \leq r$  and  $\lambda_i \geq \lambda_d$  for  $i > r$  in Eq. (15),

$$\begin{aligned} \sum_{i=1}^r \lambda_i (1 - p_i) - \sum_{i=r+1}^d \lambda_i p_i &\leq \lambda_1 \sum_{i=1}^r (1 - p_i) - \lambda_d \sum_{i=r+1}^d p_i \\ &= \lambda_1 S - \lambda_d S = (\lambda_1 - \lambda_d) S = \frac{\lambda_1(M_{\text{OOD}}) - \lambda_d(M_{\text{OOD}})}{2} \Delta(\Pi_{\text{dec}})^2. \end{aligned} \quad (17)$$

519 **Lower bound.** Using  $\lambda_i \geq \lambda_r$  for  $i \leq r$  and  $\lambda_i \leq \lambda_{r+1}$  for  $i > r$ ,

$$\begin{aligned} \sum_{i=1}^r \lambda_i (1 - p_i) - \sum_{i=r+1}^d \lambda_i p_i &\geq \lambda_r \sum_{i=1}^r (1 - p_i) - \lambda_{r+1} \sum_{i=r+1}^d p_i \\ &= \lambda_r S - \lambda_{r+1} S = (\lambda_r - \lambda_{r+1}) S = \frac{\gamma_{\text{OOD}}}{2} \Delta(\Pi_{\text{dec}})^2. \end{aligned} \quad (18)$$

520  $\square$

521 **A.4 Corollary: Second-Moment Shift Upper-Bounds OOD Faithfulness Degradation**

522 Combining Propositions 1 and 2 yields an explicit bound on how much faithfulness the ID explainer  
523 loses under OOD.

524 **Corollary 1** (Second-Moment Shift Upper-Bounds OOD Faithfulness Degradation for the ID Ex-  
525 plainer). *Suppose  $\gamma_{\text{ID}} := \lambda_r(M_{\text{ID}}) - \lambda_{r+1}(M_{\text{ID}}) > 0$ . Then*

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{ID}}) - \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) \leq \frac{\lambda_1(M_{\text{OOD}}) - \lambda_d(M_{\text{OOD}})}{2} \left( \frac{\sqrt{2}}{\gamma_{\text{ID}}} \|M_{\text{OOD}} - M_{\text{ID}}\|_F \right)^2.$$

526 *Proof.* Setting  $\Pi_{\text{dec}} = \Pi_{\text{ID}}$  in the upper bound of Proposition 2,

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{ID}}) - \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) \leq \frac{\lambda_1(M_{\text{OOD}}) - \lambda_d(M_{\text{OOD}})}{2} \Delta(\Pi_{\text{ID}})^2.$$

527 Proposition 1 gives  $\Delta(\Pi_{\text{ID}}) \leq \frac{\sqrt{2}}{\gamma_{\text{ID}}} \|M_{\text{OOD}} - M_{\text{ID}}\|_F$ . Substituting yields the result. The upper  
528 bound of Proposition 2 does not require  $\gamma_{\text{OOD}} > 0$ .  $\square$

529 The explainer-dependent component grows at most quadratically with the second-moment shift.

530 **A.5 Derivation of the GAE Procrustes Solution (Step 1)**

531 We derive the closed-form solution to the feature preservation problem in Eq. (6). Substituting  
532  $\widetilde{W}_{\text{dec}}(T) = U_{\text{OOD}}^{(:r)} T U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}}$  and expanding:

$$\|\widetilde{W}_{\text{dec}}(T) - W_{\text{dec}}^{\text{ID}}\|_F^2 = \|W_{\text{dec}}^{\text{ID}}\|_F^2 + \|U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}}\|_F^2 - 2 \text{tr}[(W_{\text{dec}}^{\text{ID}})^\top U_{\text{OOD}}^{(:r)} T U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}}], \quad (19)$$

533 where we used  $T^\top T = I$  and  $(U_{\text{OOD}}^{(:r)})^\top U_{\text{OOD}}^{(:r)} = I$ . The first two terms are independent of  $T$ , so  
534 minimization reduces to

$$T^* = \arg \max_{T \in \mathcal{O}_r} \text{tr}(GT), \quad G = U_{\text{dec}}^\top W_{\text{dec}}^{\text{ID}} (W_{\text{dec}}^{\text{ID}})^\top U_{\text{OOD}}^{(:r)}.$$

535 Let  $G = \widetilde{U} \Sigma \widetilde{V}^\top$  be the SVD of  $G$ . Setting  $R = \widetilde{V}^\top T \widetilde{U}$ , we have  $\text{tr}(GT) = \text{tr}(\Sigma R) = \sum_i \sigma_i R_{ii}$ .  
536 Since  $R \in \mathcal{O}_r$ , each  $|R_{ii}| \leq 1$ , so the maximum is achieved at  $R = I_r$ , giving  $T^* = \widetilde{V} \widetilde{U}^\top$ .  $\square$

537 **A.6 Derivation of the Step 2 Closed-Form Solution**

538 With the encoder fixed, Step 2 solves the following convex quadratic objective over  $W_{\text{dec}} \in \mathbb{R}^{d \times k}$   
539 and  $b \in \mathbb{R}^d$ :

$$\min_{W_{\text{dec}}, b} \frac{1}{N} \sum_{i=1}^N \|h_i - W_{\text{dec}} z_i - b\|^2 + \lambda_{\text{geom}} \|(I - \widehat{\Pi}_{\text{OOD}}) W_{\text{dec}}\|_F^2 + \lambda_{\text{pres}} \|W_{\text{dec}} - \widetilde{W}_{\text{dec}}(T^*)\|_F^2.$$

540 **Bias.** Setting  $\partial/\partial b = 0$  gives  $b^* = \frac{1}{N} \sum_i h_i - W_{\text{dec}}^* \frac{1}{N} \sum_i z_i$ , i.e., Eq. (9).

541 **Decoder.** Substituting  $b^*$  centers the data: define  $h_i^c = h_i - \frac{1}{N} \sum_j h_j$  and  $z_i^c = z_i - \frac{1}{N} \sum_j z_j$ .  
542 The problem reduces to

$$\min_{W_{\text{dec}}} \frac{1}{N} \sum_{i=1}^N \|h_i^c - W_{\text{dec}} z_i^c\|^2 + \lambda_{\text{geom}} \|(I - \widehat{\Pi}_{\text{OOD}}) W_{\text{dec}}\|_F^2 + \lambda_{\text{pres}} \|W_{\text{dec}} - \widetilde{W}_{\text{dec}}(T^*)\|_F^2.$$

543 Setting  $\partial/\partial W_{\text{dec}} = 0$ :

$$[\lambda_{\text{pres}} I + \lambda_{\text{geom}}(I - \widehat{\Pi}_{\text{OOD}})] W_{\text{dec}} + W_{\text{dec}} \frac{1}{N} \sum_i z_i^c z_i^{c\top} = \frac{1}{N} \sum_i h_i^c z_i^{c\top} + \lambda_{\text{pres}} \widetilde{W}_{\text{dec}}(T^*).$$

544 The left-hand coefficient  $\Lambda := \lambda_{\text{pres}} I + \lambda_{\text{geom}}(I - \widehat{\Pi}_{\text{OOD}})$  is diagonal in the OOD basis,  
545 with eigenvalue  $\lambda_{\text{pres}}$  on  $\text{span}(U_{\text{OOD}}^{(:r)})$  and  $\lambda_{\text{pres}} + \lambda_{\text{geom}}$  on its complement. Denoting  $B =$   
546  $\frac{1}{N} \sum_i z_i^c z_i^{c\top} + \lambda_{\text{pres}} I$  and  $C = \frac{1}{N} \sum_i h_i^c z_i^{c\top} + \lambda_{\text{pres}} \widetilde{W}_{\text{dec}}(T^*)$ , the system decouples row-wise in  
547 the OOD basis into two standard ridge regressions:

- 548 • **Inside**  $\widehat{\Pi}_{\text{OOD}}$  (first  $r$  rows):  $W_{\text{in}} = C_{\text{in}} B^{-1}$ , with ridge level  $\lambda_{\text{pres}}$ .
- 549 • **Outside**  $\widehat{\Pi}_{\text{OOD}}$  (remaining  $d-r$  rows):  $W_{\text{out}} = C_{\text{out}} (B + \lambda_{\text{geom}} I)^{-1}$ , with ridge level  $\lambda_{\text{pres}} +$
- 550  $\lambda_{\text{geom}}$ .

551 Combining via the projector yields Eq. (8):  $W_{\text{dec}}^{\text{GAE}} = \widehat{\Pi}_{\text{OOD}} C B^{-1} + (I - \widehat{\Pi}_{\text{OOD}}) C (B +$   
 552  $\lambda_{\text{geom}} I)^{-1}$ .  $\square$

## 553 A.7 Proof of Theorem 1

554 *Proof.* By Eq. (5), every column of  $\widetilde{W}_{\text{dec}}(T^*)$  lies in  $\text{span}(U_{\text{OOD}}^{(:,r)})$ , so

$$\Pi_{\text{dec}}^{\text{GAE}} = U_{\text{OOD}}^{(:,r)} (U_{\text{OOD}}^{(:,r)})^\top = \widehat{\Pi}_{\text{OOD}}. \quad (20)$$

555 Under the condition  $\widehat{\Pi}_{\text{OOD}} = \Pi_{\text{OOD}}$ , this gives  $\Delta(\Pi_{\text{dec}}^{\text{GAE}}) = \|\Pi_{\text{OOD}} - \Pi_{\text{dec}}^{\text{GAE}}\|_F = 0$  and therefore

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}^{\text{GAE}}) = \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}). \quad (21)$$

556 Applying the lower bound of Proposition 2 to  $\Pi_{\text{dec}} = \Pi_{\text{ID}}$ :

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{ID}}) - \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) \geq \frac{\gamma_{\text{OOD}}}{2} \Delta(\Pi_{\text{ID}})^2. \quad (22)$$

557 Substituting  $\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}^{\text{GAE}}) = \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}})$  and rearranging:

$$\mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}^{\text{GAE}}) = \mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) \leq \mathcal{L}_{\text{OOD}}(\Pi_{\text{ID}}) - \frac{\gamma_{\text{OOD}}}{2} \Delta(\Pi_{\text{ID}})^2. \quad (23)$$

558  $\square$

## 559 B Empirical Evidence for Section 3

560 This appendix provides empirical support for the theoretical results in Section 3. We verify three  
 561 claims: (i) the explainer subspace aligns with the ID-active subspace, (ii) the explainer-dependent  
 562 component in the decomposition (3) accounts for a meaningful fraction of the total OOD error, and  
 563 (iii) second-moment shift enlarges the faithfulness gap as predicted by Proposition 1.

### 564 B.1 Controlled Toy Setting

565 All experiments in this appendix use a controlled toy setting that allows us to vary OOD severity  
 566 continuously while keeping the model and explainer fixed.

567 **Target model.** The target model is a 2-layer ReLU MLP with input dimension  $d_{\text{in}}=128$ , hidden  
 568 dimension  $d=256$ , and output dimension  $p \in \{4, 8, 16\}$ :

$$h(x) = \text{ReLU}(W_1 x + b_1) \in \mathbb{R}^d, \quad o(x) = W_2 h(x) + b_2 \in \mathbb{R}^p.$$

569 The hidden activations  $h(x) \in \mathbb{R}^d$  correspond to the hidden representations analyzed in the main text.  
 570 The explainer operates on these  $d$ -dimensional activations.

571 **Explainer.** We train both a transcoder and an SAE on ID hidden activations using the standard  
 572 reconstruction-plus-sparsity objective (ERM), with dictionary sizes  $k \in \{d/2, 1d, 2d, 4d, 8d, 32d\}$ .  
 573 The subspace rank is set to  $r = p$ , matching the rank of the output weight matrix  $W_2 \in \mathbb{R}^{p \times d}$ :  
 574 since  $o(x) = W_2 h(x) + b_2$ , the model’s output depends on  $h(x)$  only through its projection onto  
 575  $\text{span}(W_2^\top)$ , which has dimension  $p$ . This makes  $r = p$  the natural rank at which the active subspace  
 576 captures all output-relevant directions.

577 **OOD generation.** ID inputs are drawn from  $x \sim \mathcal{N}(0, I_{d_{\text{in}}})$ . OOD inputs are generated as  
 578  $x = A_s z$  where  $z \sim \mathcal{N}(0, I_{d_{\text{in}}})$  and  $A_s$  is a severity-dependent transformation matrix. Specifically,  
 579 let  $Q \in \mathbb{R}^{d_{\text{in}} \times d_{\text{in}}}$  be a fixed random orthogonal matrix and  $\mathbf{s} \in \mathbb{R}^{d_{\text{in}}}$  be fixed slopes linearly spaced  
 580 in  $[-S, S]$  with  $S=6$ . The base input covariance under severity  $s$  is

$$\Sigma_{\text{base}}(s) = Q \text{diag}(e^{s \cdot \mathbf{s}}) Q^\top (I + s\rho(1 + \frac{s}{2}) VV^\top),$$

581 where  $V \in \mathbb{R}^{d_{\text{in}} \times r_V}$  is a fixed random orthonormal matrix ( $r_V=32$ ) and  $\rho=10$ . The base covariance  
 582 is then rescaled directionally: variance in the output-relevant subspace  $\text{span}(W_2^\top)$  is reduced by  
 583 factor  $(1 - 0.6s)$ , and variance in its orthogonal complement is amplified by factor  $(1 + 2s^2)$ . Finally,  
 584 the covariance is globally normalized so that  $\text{tr}(\Sigma(s)) = d_{\text{in}}$  for all  $s$ , ensuring that reconstruction  
 585 error differences are not driven by trivial scale changes. The transformation matrix  $A_s$  is the matrix  
 586 square root of the resulting  $\Sigma(s)$ . At  $s=0$ ,  $\Sigma(0) = I$  (ID); as  $s$  increases toward 1, the input  
 587 covariance undergoes progressive rotation and anisotropic rescaling, which propagates through the  
 588 ReLU layer to induce second-moment shift in hidden space. We use  $N=20,000$  samples for all  
 589 geometric computations.

590 **Validation on real models.** We also validate the results on pretrained language models (GPT-2  
 591 Small, Pythia-1.4B) under temporal, domain, and adversarial distribution shifts, using the experimen-  
 592 tal setup described in Section 5. Unlike the toy setting, OOD severity is not varied continuously;  
 593 instead, we compare ID and pure OOD for each shift type.

## 594 B.2 Explainer Subspace Alignment with the ID Dominant Subspace

595 Section 3.1 claims that the explainer subspace  $\Pi$  aligns closely with  $\Pi_{\text{ID}}$  for a well-trained ID  
 596 explainer. We verify this by measuring the subspace overlap

$$\text{overlap}(U_{\text{dec}}, U_{\text{ID}}) := \frac{1}{r} \|U_{\text{dec}}^\top U_{\text{ID}}\|_F^2 \in [0, 1],$$

597 where  $U_{\text{ID}} \in \mathbb{R}^{d \times r}$  contains the top- $r$  eigenvectors of  $M_{\text{ID}}$ . A value of 1 indicates perfect alignment;  
 598 0 indicates orthogonality.

### 599 B.2.1 Toy Setting

600 We sweep dictionary sizes  $k \in \{d/2, 1d, 2d, 4d, 8d, 32d\}$  for both transcoders and SAEs to test  
 601 whether the claim holds across different levels of overcompleteness. Figure 5 reports the overlap as a  
 602 function of OOD severity for each dictionary size. In both panels, solid lines show the explainer–ID  
 603 overlap and dashed lines show the explainer–OOD overlap.

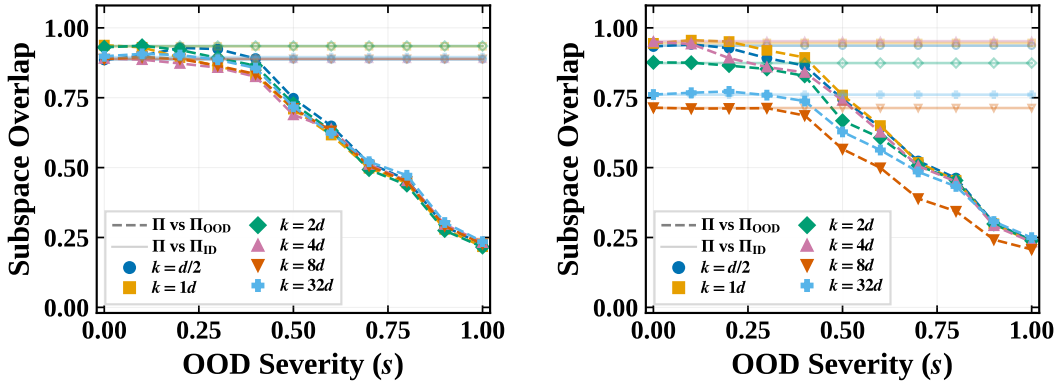


Figure 5: **Explainer subspace overlap between the explainer and the ID-active subspace (solid) vs. the OOD-active subspace (dashed), as a function of OOD severity  $s$ , for dictionary sizes  $k \in \{d/2, 1d, 2d, 4d, 8d, 32d\}$ . Left: Transcoder. Right: SAE. For  $k \geq 4d$ , both explainer types maintain high ID overlap ( $> 0.89$ ) regardless of severity, while OOD overlap degrades monotonically.**

604 **Transcoder (left).** The explainer–ID overlap (solid) remains above 0.89 for all dictionary sizes and  
 605 all severity levels. The overlap is largely insensitive to  $k$ : even an undercomplete dictionary ( $k=d/2$ )  
 606 captures the ID-active subspace well. The explainer–OOD overlap (dashed) degrades monotonically  
 607 with severity, dropping below 0.3 at  $s=1.0$  regardless of  $k$ .

608 **SAE (right).** The explainer–ID overlap depends more on  $k$ . For  $k \geq 4d$ , the overlap exceeds  
 609 0.93, comparable to the transcoder. For smaller  $k$  ( $k=d/2$  or  $k=1d$ ), the overlap drops to 0.70–  
 610 0.77, indicating that undercomplete SAEs do not fully capture the ID-active subspace. As with the  
 611 transcoder, the explainer–OOD overlap degrades with severity for all  $k$ .

612 **Interpretation.** For sufficiently overcomplete dictionaries ( $k \geq 4d$ , the standard setting in practice),  
 613 both transcoders and SAEs align closely with the ID-active subspace regardless of OOD severity,  
 614 confirming the claim in Section 3.1. The divergence between the ID overlap (flat) and the OOD  
 615 overlap (decreasing) is precisely the faithfulness gap: the explainer remains anchored to the ID  
 616 geometry while the model’s active subspace rotates away under OOD shift.

## 617 B.2.2 Real-Data Setting

618 Figure 6 reports the subspace overlap for ID-trained explainers on GPT-2 Small and Pythia-1.4B  
 619 under temporal, domain, and adversarial shifts. The rank is  $r=64$  for both models.

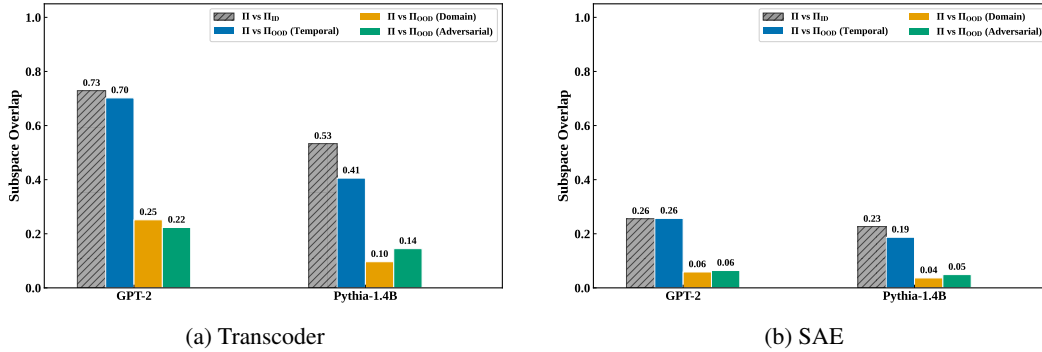


Figure 6: **Subspace overlap on pretrained language models.** Hatched bars: explainer vs. ID-active subspace. Colored bars: explainer vs. OOD-active subspace under temporal, domain, and adversarial shifts. The explainer–ID overlap consistently exceeds the OOD overlaps across both explainer types and all shift types.

620 **Transcoder.** Across both models, the explainer–ID overlap (hatched bars, 0.53–0.73) substantially  
 621 exceeds the explainer–OOD overlap, which drops to 0.10–0.25 for domain and adversarial shifts. The  
 622 gap is larger for domain and adversarial shifts than for temporal shift, consistent with the stronger  
 623 geometric distortion induced by these shift types.

624 **SAE.** SAE explainers show a similar pattern. The explainer–ID overlap (0.23–0.26) exceeds the  
 625 OOD overlap under domain and adversarial shifts (0.04–0.06), while temporal shift retains a larger  
 626 portion of the ID overlap (0.19–0.26). The overall overlap values are lower than transcoders because  
 627 SAEs reconstruct residual-stream activations rather than MLP outputs, spreading energy across more  
 628 directions.

629 **Interpretation.** The ID overlap values are lower than in the toy setting (0.89+). This is expected:  
 630 real language models have higher effective dimensionality, and the rank  $r=64$  captures a smaller  
 631 fraction of the total hidden dimension ( $d=768$  for GPT-2,  $d=2048$  for Pythia-1.4B) than  $r=p$  in  
 632 the toy setting. Despite this, the relative pattern (ID overlap above OOD overlap under domain and  
 633 adversarial shifts) holds consistently for both transcoders and SAEs, confirming the alignment claim  
 634 in Section 3.1 on real models. Temporal shift produces a milder gap, consistent with its smaller  
 635 second-moment perturbation.

636 **Faithfulness gap on  $M_{ID}$  vs. on  $M_{OOD}$ .** Definition 1 measures the gap of  $\Pi_{dec}$  against  $\Pi_{OOD}$ . To  
 637 assess the residual ID-side misalignment, we additionally report the analogous quantity against  $\Pi_{ID}$ ,

$$\Delta_{ID}(\Pi_{dec}) := \|\Pi_{ID} - \Pi_{dec}\|_F,$$

638 which directly quantifies how well the ID-trained explainer captures  $\Pi_{ID}$ . Frobenius gap and overlap  
 639 encode the same information through  $\|\Pi_A - \Pi_B\|_F = \sqrt{2r(1 - \text{overlap}(U_A, U_B))}$ , so Table 4  
 640 reports the gaps converted from the same measurements as Figure 6. The maximum value at  $r=64$  is  
 641  $\sqrt{2r} = \sqrt{128} \approx 11.31$ .

642 For transcoders,  $\Delta_{ID}(\Pi_{dec})$  is consistently smaller than  $\Delta(\Pi_{dec})$ , with the OOD-side gap exceeding  
 643 the ID-side gap by 1.35–1.69 $\times$  under domain and adversarial shifts, where the second-moment shift

Table 4: **Frobenius faithfulness gap of the ID-trained explainer at  $r=64$ .** Columns 4–6 report  $\Delta_{\text{ID}}(\Pi_{\text{dec}}) = \|\Pi_{\text{ID}} - \Pi_{\text{dec}}\|_F$ ,  $\Delta(\Pi_{\text{dec}}) = \|\Pi_{\text{OOD}} - \Pi_{\text{dec}}\|_F$ , and  $\|\Pi_{\text{ID}} - \Pi_{\text{OOD}}\|_F$ . The last column is the ratio  $\Delta(\Pi_{\text{dec}})/\Delta_{\text{ID}}(\Pi_{\text{dec}})$ . Maximum possible value at  $r=64$  is  $\sqrt{128} \approx 11.31$ . All values are derived from the overlap measurements behind Figure 6.

Explainer	Model	Shift	$\Delta_{\text{ID}}(\Pi_{\text{dec}})$	$\Delta(\Pi_{\text{dec}})$	$\ \Pi_{\text{ID}} - \Pi_{\text{OOD}}\ _F$	Ratio
Transcoder	GPT-2 Small	Temporal	5.89	6.17	5.65	1.05×
	GPT-2 Small	Domain	5.89	9.79	9.83	1.66×
	GPT-2 Small	Adversarial	5.89	9.97	9.83	1.69×
	Pythia-1.4B	Temporal	7.73	8.72	7.68	1.13×
	Pythia-1.4B	Domain	7.73	10.75	10.78	1.39×
	Pythia-1.4B	Adversarial	7.73	10.46	10.49	1.35×
SAE	GPT-2 Small	Temporal	9.76	9.75	4.83	1.00×
	GPT-2 Small	Domain	9.76	10.98	10.33	1.12×
	GPT-2 Small	Adversarial	9.76	10.95	10.30	1.12×
	Pythia-1.4B	Temporal	9.95	10.20	7.65	1.03×
	Pythia-1.4B	Domain	9.95	11.10	10.88	1.12×
	Pythia-1.4B	Adversarial	9.95	11.03	10.71	1.11×

644 is largest. Moreover,  $\|\Pi_{\text{ID}} - \Pi_{\text{OOD}}\|_F$  tracks  $\Delta(\Pi_{\text{dec}})$  to within 3% under these shifts, empirically  
 645 supporting the substitution  $\Pi_{\text{dec}} \approx \Pi_{\text{ID}}$  used in Proposition 1. SAE explainers exhibit a larger  
 646  $\Delta_{\text{ID}}(\Pi_{\text{dec}})$  (close to the upper bound at  $r=64$ ), reflecting the lower per-rank overlap of residual-  
 647 stream dictionaries; the ordering  $\Delta(\Pi_{\text{dec}}) \geq \Delta_{\text{ID}}(\Pi_{\text{dec}})$  still holds under domain and adversarial  
 648 shifts but with a smaller margin (1.11–1.12×). Temporal shift produces nearly equal ID- and  
 649 OOD-side gaps for both explainer types, consistent with its milder second-moment perturbation.

### 650 B.3 Relative Magnitude of the Explainer-Dependent Term

651 The decomposition of Eq. (3) separates OOD faithfulness as

$$\mathcal{L}_{\text{OOD}}(\Pi) = \underbrace{\mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}})}_{\text{irreducible}} + \underbrace{\text{tr}[(\Pi_{\text{OOD}} - \Pi) M_{\text{OOD}}]}_{\text{explainer-dependent}}.$$

652 Adaptation can only reduce the explainer-dependent component. If this component is negligible rela-  
 653 tive to the irreducible component, no adaptation strategy can meaningfully improve OOD faithfulness.  
 654 We measure the explainer-dependent ratio

$$\eta := \frac{\text{tr}[(\Pi_{\text{OOD}} - \Pi) M_{\text{OOD}}]}{\mathcal{L}_{\text{OOD}}(\Pi)}$$

655 to assess whether the explainer-dependent component is an actionable target.

#### 656 B.3.1 Toy Setting

657 Figure 7 plots  $\eta$  as a function of OOD severity  $s$  for each dictionary size  $k$ . At  $s=0$ ,  $\eta < 0.05$  for  
 658 most  $k$ : when ID and OOD coincide, the explainer subspace is near-optimal. As severity increases,  $\eta$   
 659 grows steadily to  $\eta \approx 0.31$  at  $s=1.0$  for both transcoders and SAEs, with little variation across  $k$ .

660 The moderate value of  $\eta$  at maximum severity is a consequence of the toy setting’s low rank ratio:  
 661  $r = p = 8$  out of  $d = 256$  dimensions, so only 3.1% of the hidden space is retained. The irreducible  
 662 term  $\mathcal{L}_{\text{OOD}}(\Pi_{\text{OOD}}) = \sum_{i=r+1}^d \lambda_i(M_{\text{OOD}})$  sums 248 discarded eigenvalues and dominates by  
 663 construction. In real models where  $r$  is chosen to capture a larger fraction of the activation energy,  
 664  $\eta$  is substantially higher (Section B.3.2). The key observation in this toy setting is that  $\eta$  increases  
 665 monotonically with severity, confirming that distribution shift enlarges the explainer-dependent  
 666 component relative to the total error.

#### 667 B.3.2 Real-Data Setting

668 Figure 8 reports  $\eta$  at pure ID and pure OOD for GPT-2 Small and Pythia-1.4B under temporal,  
 669 domain, and adversarial shifts.

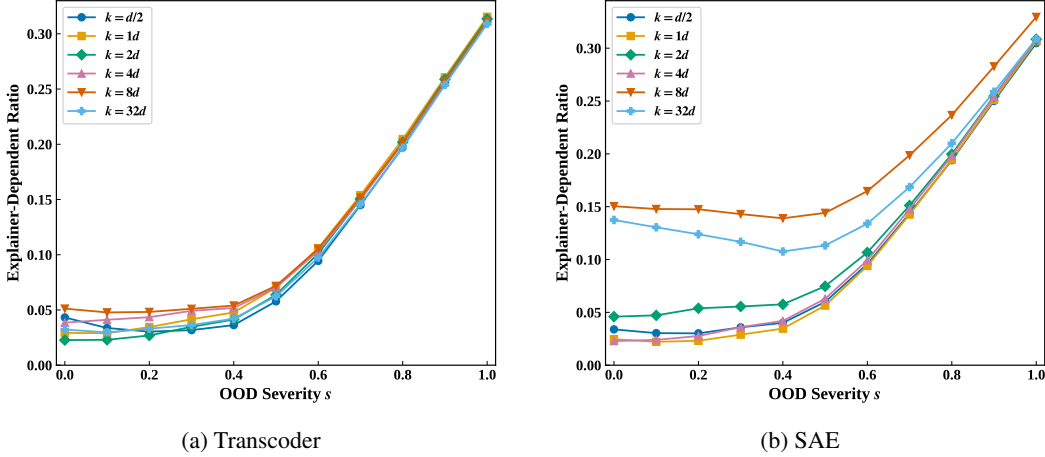


Figure 7: **Explainer-dependent ratio  $\eta$  as a function of OOD severity  $s$  for dictionary sizes  $k \in \{d/2, 1d, 2d, 4d, 8d, 32d\}$ .** At  $s=1.0$ ,  $\eta \approx 0.31$  for both explainer types, independent of  $k$ .

670 **Transcoder (left).** At pure OOD, domain and adversarial shifts reach  $\eta > 0.99$  across both  
 671 models: the explainer-dependent component dominates the total error almost entirely. Temporal  
 672 shift yields  $\eta \approx 0.66$ – $0.99$  depending on the model, reflecting its milder geometric distortion. These  
 673 values are substantially higher than in the toy setting ( $\eta \approx 0.31$ ), because the toy setting uses  
 674  $r/d = 8/256 = 3.1\%$  so the irreducible term dominates by construction.

675 **SAE (right).** SAE explainers exhibit the same trend. Under domain and adversarial shifts,  $\eta > 0.99$   
 676 for both models, confirming that the explainer-dependent component dominates. Temporal shift  
 677 yields  $\eta \approx 0.68$ – $0.99$  depending on the model, consistent with its milder geometric distortion. At  
 678 pure ID,  $\eta \approx 0.51$  (GPT-2) and  $\eta \approx 0.99$  (Pythia-1.4B), reflecting the larger model’s higher effective  
 679 dimensionality relative to  $r=64$ .

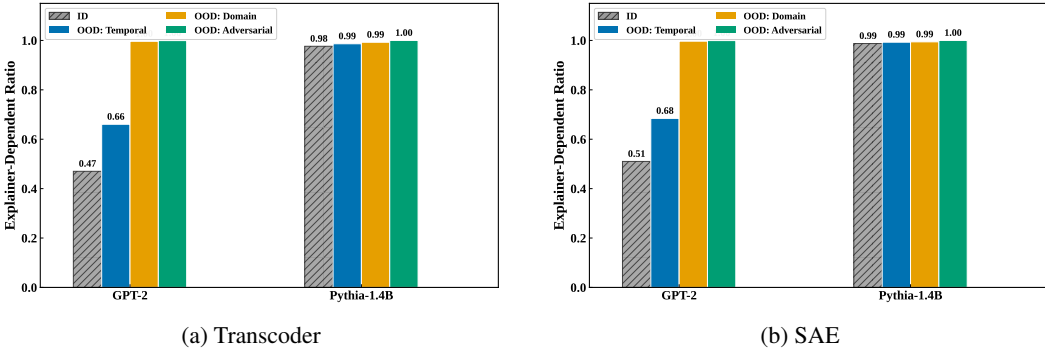


Figure 8: **Explainer-dependent ratio  $\eta$  at pure ID (hatched) and pure OOD (colored) ( $r=64$ ).** Under domain and adversarial shifts,  $\eta > 0.99$  for both explainer types.

#### 680 B.4 Empirical Verification of Proposition 1

681 Proposition 1 predicts that second-moment shift controls the faithfulness gap via

$$\Delta(\Pi_{\text{ID}}) \leq \frac{\sqrt{2}}{\gamma_{\text{ID}}} \|M_{\text{OOD}} - M_{\text{ID}}\|_F.$$

682 Since this bound depends only on  $M_{\text{ID}}$  and  $M_{\text{OOD}}$ , it is independent of the explainer architecture  
 683 and dictionary size. We verify it empirically on both the toy setting and real models.

684 **B.4.1 Toy Setting**

685 Figure 9 plots the normalized second-moment shift against  $\Delta(\Pi_{ID})$ , with color indicating OOD  
 686 severity  $s$ . The two quantities are near-perfectly correlated (Pearson  $r=0.993$ , Spearman  $\rho=1.000$ ),  
 687 consistent with the linear upper bound.

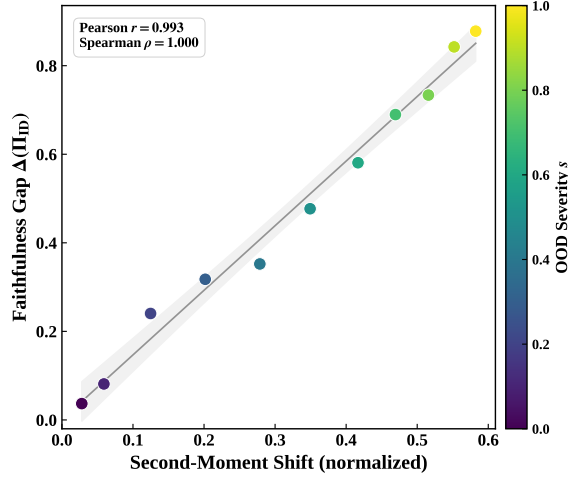


Figure 9: **Proposition 1 verification (toy)**. X: normalized second-moment shift. Y: faithfulness gap  $\Delta(\Pi_{ID})$ . Color: OOD severity  $s$ . The result is independent of explainer type and dictionary size.

688 **B.4.2 Real-Data Setting**

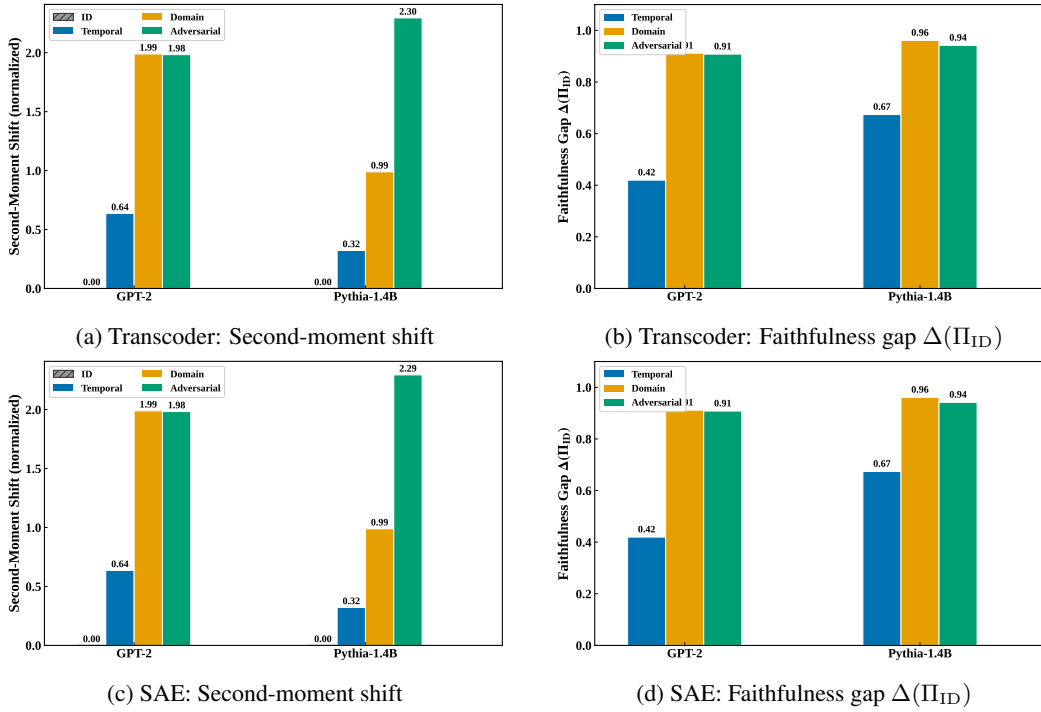


Figure 10: **Proposition 1 verification ( $r=64$ ) at pure OOD**. Top row: Transcoder. Bottom row: SAE. Within each model, larger shifts correspond to larger gaps for both explainer types.

689 Figure 10 shows the second-moment shift and faithfulness gap for GPT-2 Small and Pythia-1.4B  
 690 ( $r=64$ ) at pure OOD for both transcoders (top row) and SAEs (bottom row). Within each the

691 ordering is consistent: temporal shift produces the smallest second-moment shift and the smallest  
 692 faithfulness gap, while domain and adversarial shifts produce larger shifts and correspondingly larger  
 693 gaps.

694 **SAE.** The same ordering holds for SAE explainers: temporal shift produces the smallest second-  
 695 moment shift and faithfulness gap, while domain and adversarial shifts produce larger values. The  
 696 magnitudes are comparable to transcoders, confirming that the proposition is independent of the  
 697 explainer architecture.

698 **Summary.** The proposition holds across all models and both explainer types under diverse real-  
 699 world distribution shifts, confirming that the theoretical predictions generalize beyond the controlled  
 700 toy setting.

## 701 C Empirical Verification of Theorem 1

702 Theorem 1 predicts that GAE’s projection-loss improvement over the ID explainer grows at least  
 703 quadratically with the faithfulness gap:  $\mathcal{L}_{\text{OOD}}(\Pi_{\text{ID}}) - \mathcal{L}_{\text{OOD}}(\Pi_{\text{dec}}^{\text{GAE}}) \geq \frac{1}{2}\gamma_{\text{OOD}} \Delta(\Pi_{\text{ID}})^2$ . We verify  
 704 this in the controlled toy setting (Section 5.1) by fixing  $r = p = 8$  and sweeping OOD severity from  
 705  $s=0$  (ID) to  $s=1$  (maximum shift). At each severity, we compute the Step 1 GAE projector  $\Pi_{\text{dec}}^{\text{GAE}} =$   
 706  $\hat{\Pi}_{\text{OOD}}$  and measure the projection-loss improvement  $I(s) = \mathcal{L}_{\text{OOD},s}(\Pi_{\text{ID}}) - \mathcal{L}_{\text{OOD},s}(\hat{\Pi}_{\text{OOD}})$   
 707 against the squared faithfulness gap  $\Delta(\Pi_{\text{ID}})^2$ .

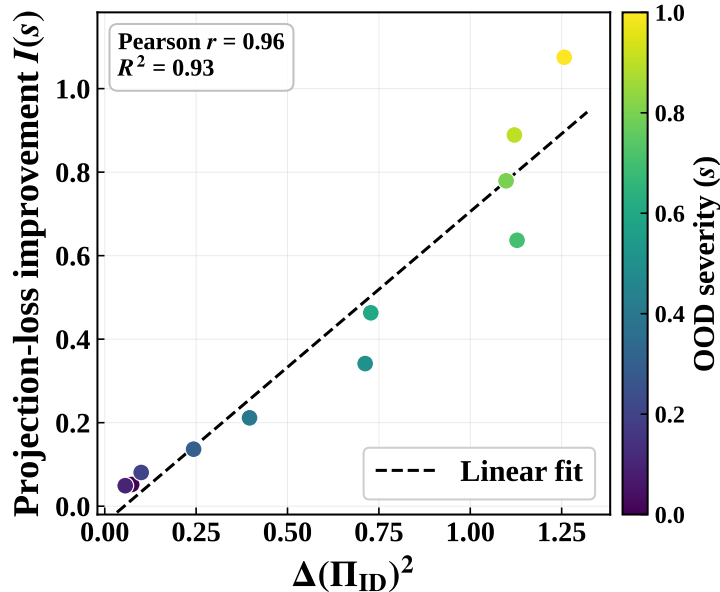


Figure 11: **Empirical verification of Theorem 1 on the controlled toy setting.** Projection-loss improvement  $I(s)$  versus the squared faithfulness gap  $\Delta(\Pi_{\text{ID}})^2$ , swept across OOD severity  $s \in [0, 1]$ . The dashed line is a linear fit ( $R^2 = 0.93$ , Pearson  $r = 0.96$ ), supporting the quadratic dependence predicted by Theorem 1.

708 Figure 11 shows a strong linear relationship between  $I(s)$  and  $\Delta(\Pi_{\text{ID}})^2$  (Pearson  $r = 0.96$ ,  $R^2 =$   
 709  $0.93$ ), supporting the quadratic dependence predicted by Theorem 1. The empirical improvement  
 710 exceeds the guaranteed lower bound at every severity (0 violations out of 11), confirming that the  
 711 bound uses the worst-case eigengap  $\gamma_{\text{OOD}}/2$  as its constant, which is  
 712 conservative relative to the effective improvement rate, as expected for a spectral-gap-based guarantee.

713 **D Experimental Details**

714 **D.1 Model and Explainer Details**

715 Table 5 summarizes the model and explainer configurations. All models are frozen pretrained  
716 checkpoints; only explainer components are adapted.

Table 5: **Model and explainer configurations.**

Model	$d$	Layer	$k$ ( $32d$ )	ID Corpus
GPT-2 Small	768	8	24,576	OpenWebText [46]
Pythia-1.4B	2,048	15	65,536	The Pile [47]

717 For each model, we train two explainer types: Top-K SAEs [7], which reconstruct residual-stream  
718 activations using Top-K sparsity, and transcoders [5], which reconstruct MLP outputs from MLP  
719 inputs. All explainers use dictionary size  $k=32d$  [4] and are trained on in-distribution activations  
720 with the standard reconstruction-plus-sparsity objective.

721 **D.2 Baseline Details**

- 722 • **Fixed (ERM).** The ID-trained explainer applied to OOD inputs without any adaptation. This is  
723 the default deployment setting for existing dictionary-based explainers.
- 724 • **TERM.** An ID explainer trained with tilted empirical risk minimization [16, 17], which upweights  
725 high-loss (rare/tail) samples during training to improve coverage of infrequent concepts. This is  
726 an alternative ID training strategy, not an OOD adaptation method.
- 727 • **Finetune [43].** The ID-trained explainer finetuned on OOD activations with a warm start. This  
728 adapts the existing dictionary to OOD data via gradient-based training.
- 729 • **Retrain.** The explainer retrained from scratch on OOD activations with the same architecture  
730 and hyperparameters. This baseline provides a reference point but is not an oracle upper bound:  
731 retraining on OOD data can distort pretrained feature structure [48].
- 732 • **SAEBoost.** A residual boosting approach [18]: a secondary explainer is trained on the OOD  
733 reconstruction residuals of the ID-trained base explainer, and the two outputs are summed at  
734 inference ( $\hat{h} = \hat{h}_{\text{base}} + \hat{h}_{\text{resid}}$ ). This adds OOD-specific capacity while retaining the base  
735 dictionary, but requires OOD training data.
- 736 • **FaithfulSAE.** The explainer retrained on the target model’s own unconditional generations [15],  
737 avoiding dependence on external datasets. Requires full retraining but no OOD data.
- 738 • **GAE (ours).** Training-free geometric adaptation (Algorithm 1). Step 1 rotates the ID dictionary’s  
739 subspace to align with the OOD-active subspace via orthogonal Procrustes. Step 2 refits the  
740 decoder via constrained ridge regression with geometry and preservation regularization. The  
741 entire pipeline is closed-form; no gradient computation or iterative training is required.

742 **D.3 GAE Implementation Details**

743 **Step 2 regularization.** The closed-form decoder refit (Section 4.2, Step 2) regularizes the decoder  
744 toward the Step 1 output  $\widetilde{W}_{\text{dec}}(T^*)$  with weight  $\lambda_{\text{pres}}$ , following Eq. (8).

745 **Decoder interpolation.** The Step 2 closed-form solution  $W_{\text{dec}}^{\text{GAE}}$  optimizes sample-level recon-  
746 struction under geometry constraints. With limited OOD samples, this solution can overfit to the  
747 estimation noise in  $\{z_i, h_i\}$ . To mitigate this, we interpolate the Step 2 output with the Step 1 rotated  
748 dictionary:

$$W_{\text{final}} = (1 - \alpha) W_{\text{dec}}^{\text{GAE}} + \alpha \widetilde{W}_{\text{dec}}(T^*), \tag{24}$$

749 where  $\alpha \in [0, 1]$  controls the interpolation. When  $\alpha = 0$ , the output equals the closed-form solution  
750 in Section 4.2. When  $\alpha = 1$ , the output equals the Step 1 rotation without reconstruction refinement.  
751 We treat  $\alpha$  as a hyperparameter selected per OOD setting.

752 **Hyperparameter selection.** GAE requires no gradient computation or iterative optimization. The  
 753 hyperparameters ( $r$ ,  $\lambda_{\text{geom}}$ ,  $\lambda_{\text{pres}}$ ,  $\alpha$ ) are selected per OOD setting using a small held-out portion  
 754 of unlabeled OOD activations, monitoring reconstruction quality ( $|\Delta\text{CE}|$ ) and causal faithfulness  
 755 (nComp). No OOD labels are required. Once selected, the same hyperparameters are used for all  
 756 evaluation prompts.

757 **Hyperparameter summary.** Tables 6 and 7 list the GAE hyperparameters for transcoders and  
 758 SAEs, respectively. All settings use the second-moment matrix (not centered covariance) for OOD  
 759 subspace estimation, as prescribed in Algorithm 1.

Table 6: GAE hyperparameters for transcoder experiments.

Model	OOD Setting	$r$	$\lambda_{\text{geom}}$	$\lambda_{\text{pres}}$	$\alpha$	$N_{\text{fit}}$
GPT-2	HaluEval (Adversarial)	32	0.1	0.2	0	2,048
GPT-2	Edgar (Domain)	3	0.1	0.2	0	2,048
GPT-2	FineWeb (Temporal)	6	0.1	0.04	0	2,048
Pythia-1.4B	HaluEval (Adversarial)	64	15	2	0	2,048
Pythia-1.4B	Edgar (Domain)	64	0.1	0.2	0	2,048
Pythia-1.4B	FineWeb (Temporal)	64	20	1	0	2,048

Table 7: GAE hyperparameters for SAE experiments. When only Step 1 (Procrustes rotation) is applied, all Step 2 hyperparameters are set to zero.

Model	OOD Setting	$r$	$\lambda_{\text{geom}}$	$\lambda_{\text{pres}}$	$\alpha$	$N_{\text{fit}}$
GPT-2	HaluEval (Adversarial)	639	0	0	1	0
GPT-2	Edgar (Domain)	462	0	0	1	0
GPT-2	FineWeb (Temporal)	700	0	0	1	0
Pythia-1.4B	HaluEval (Adversarial)	3	0.1	0.2	0	2,048
Pythia-1.4B	Edgar (Domain)	1,750	0	0	1	0
Pythia-1.4B	FineWeb (Temporal)	3	0	0.2	0	2,048

#### 760 D.4 Evaluation Details

761 **Normalized comprehensiveness (nComp).** We measure causal faithfulness via logit-level feature  
 762 ablation [4, 21]. Given a prompt, let  $\ell_0$  denote the target-token logit under the explainer’s full  
 763 reconstruction, and  $\ell_{\emptyset}$  the logit when all features are ablated to zero. For a feature budget  $m^*$ , let  
 764  $\ell_{\setminus m^*}$  be the logit after removing the top- $m^*$  features. We define

$$\text{nComp} = \frac{\ell_0 - \ell_{\setminus m^*}}{|\ell_0 - \ell_{\emptyset}|}, \quad (25)$$

765 where  $m^* = 32$ . Higher values indicate that the top features are causally important for the model’s  
 766 output.

767 **Delta cross-entropy ( $\Delta\text{CE}$ ).** We measure reconstruction quality by the cross-entropy increase  
 768 when original activations are replaced with the explainer’s reconstruction [7]:

$$\Delta\text{CE} = \text{CE}(\hat{h}) - \text{CE}(h), \quad (26)$$

769 where  $\text{CE}(h)$  is the loss with original activations and  $\text{CE}(\hat{h})$  is the loss with reconstructed activations.  
 770 Lower values indicate better preservation of the model’s predictive behavior.

771 **Normalized AOPC (nAOPC).** nAOPC [23] averages the normalized logit drop across multiple  
 772 feature budgets when top- $m$  features are removed:

$$\text{nAOPC} = \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \frac{\ell_0 - \ell_{\setminus m}}{|\ell_0 - \ell_{\emptyset}|}, \quad (27)$$

773 where  $\mathcal{M} = \{1, 2, 4, 8, 16, 32, 64, 128\}$ . Higher values indicate that the identified features are  
 774 causally important across a range of budgets.

775 **Evaluation protocol.** We evaluate at the last token position using zero-residual ablation (replacing  
 776 ablated features with zero). The denominator  $|\ell_0 - \ell_\emptyset|$  normalizes each example by the logit range  
 777 between full and empty reconstruction, enabling cross-example comparison regardless of the absolute  
 778 logit scale. We exclude examples where  $|\ell_0 - \ell_\emptyset| < 0.1$  to avoid unstable normalization. Feature  
 779 budgets are  $\mathcal{M} = \{1, 2, 4, 8, 16, 32, 64, 128\}$  with  $m^* = 32$ . We use  $N_{\text{eval}} = 1,000$  evaluation  
 780 prompts per setting and seed = 2026 throughout.

## 781 D.5 Compute Resources

782 All experiments run on a single NVIDIA RTX A6000 GPU (48 GiB VRAM) with an Intel Xeon  
 783 Gold 6326 CPU and 252 GiB of system RAM. No experiment requires multi-GPU or model-parallel  
 784 execution. GPT-2 runs use peak GPU memory under 8 GiB. Pythia-1.4B runs with batch size 64 use  
 785 peak GPU memory under 24 GiB.

786 **Per-method wall-clock.** Table 1 reports the cost of a single (model, OOD setting) run for each  
 787 adaptation method. Finetune processes 5M tokens, taking about 2 minutes on GPT-2 and 12 minutes  
 788 on Pythia-1.4B. Retrain, SAEBoost, and FaithfulSAE each process 100M tokens, taking about  
 789 39 minutes on GPT-2 and 4 hours on Pythia-1.4B. GAE finishes in 0.5 s on GPT-2 and 2.9 s on  
 790 Pythia-1.4B using a single forward pass over  $\sim 2,000$  OOD activations and no gradient computation.  
 791 Faithfulness evaluation (nAOPC, nComp,  $\Delta\text{CE}$  on 1,000 prompts) adds about 1 minute per (model,  
 792 OOD setting, baseline).

793 **Total compute.** The full result table (two models, three OOD settings, six adaptation baselines)  
 794 requires roughly 50 GPU-hours on a single RTX A6000, dominated by the Retrain-style baselines on  
 795 Pythia-1.4B. The GAE rows themselves contribute under 1 GPU-minute to this total. Pretraining of  
 796 the ID dictionaries (transcoders and SAEs on OpenWebText / The Pile) is a one-time cost that we  
 797 treat as external to the adaptation experiments.

## 798 E Additional Case Studies on Other Semantic Classes

799 This appendix repeats the body case-study protocol on two further HaluEval prompts whose target  
 800 tokens fall in distinct semantic classes: male first names and professions. The protocol is unchanged  
 801 from Section 5.2.3: we keep the encoder frozen, take the top-3 features by GAE causal effect on  
 802 the target token, and report each feature’s direct logit attribution to 20 class-member tokens and 10  
 803 unrelated noun controls. Fixed and GAE share the same encoder and the same top-3 features, so any  
 804 difference in attribution comes entirely from the decoder rotation.

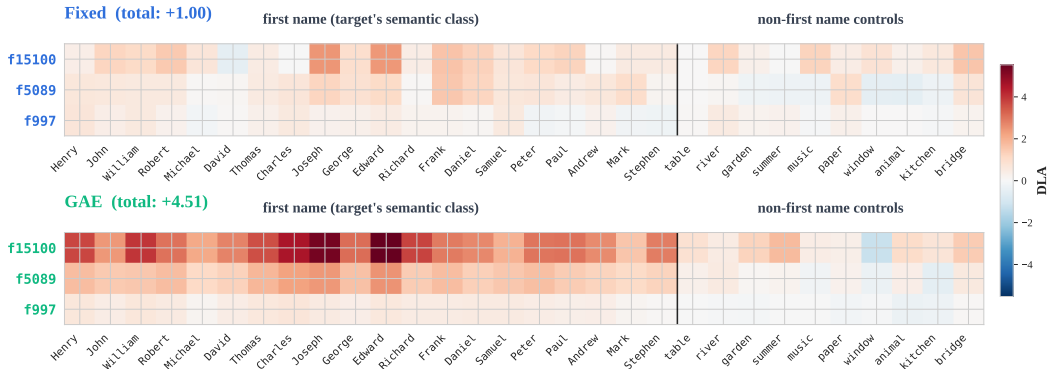


Figure 12: **Per-feature DLA on a prompt predicting ‘Henry’ (GPT-2, Transcoder).** The truncated input is “Question: What nationality was James”; the next token is a male first name. Each cell reports the feature’s direct logit attribution to a candidate token, with 20 male first names on the left and 10 unrelated noun controls on the right. Fixed’s total class-specificity is +1.00 and GAE’s is +4.51, a  $4.5\times$  amplification of the same encoder-selected features’ pull toward the first-name class.

805 In both cases the decoder rotation alone reproduces the body-case finding: the same top-3 features,  
 806 with the same activations, contribute more to their target’s semantic class under GAE than under  
 807 Fixed. The feature selection itself is identical because the encoder is shared.

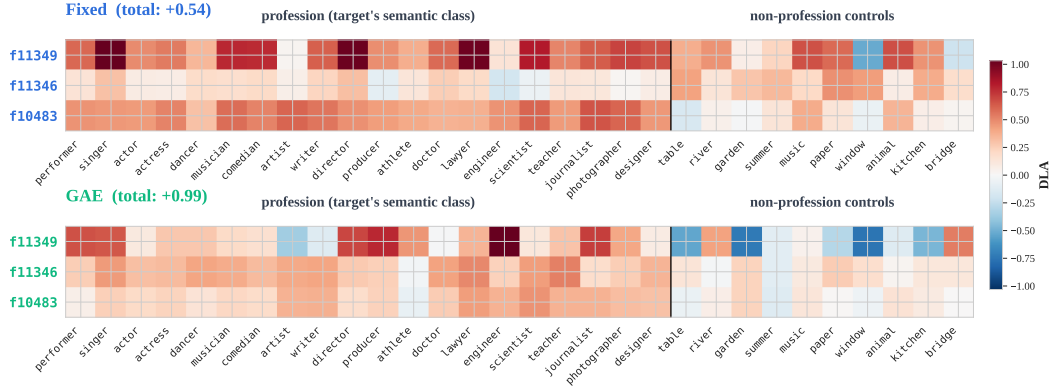


Figure 13: **Per-feature DLA on a prompt predicting ‘politician’ (GPT-2, Transcoder).** The truncated input is “Question: Which American”; the next token is a profession. The 20 class-member tokens are common professions and the 10 controls are unrelated nouns. Fixed’s total class-specificity is +0.54 and GAE’s is +0.99. The GAE row drives several control cells negative (blue), where Fixed leaves them positive, sharpening the contrast between the class and its controls without changing which features were selected.

## 808 F Hyperparameter Sensitivity

809 We sweep GAE’s three hyperparameters on HaluEval (GPT-2 + Transcoder), holding the other two at  
 810 the defaults  $r = 32$ ,  $N_{\text{OOD}} = 2000$ ,  $\lambda_{\text{pres}} = 0.2$  (Figure 14).

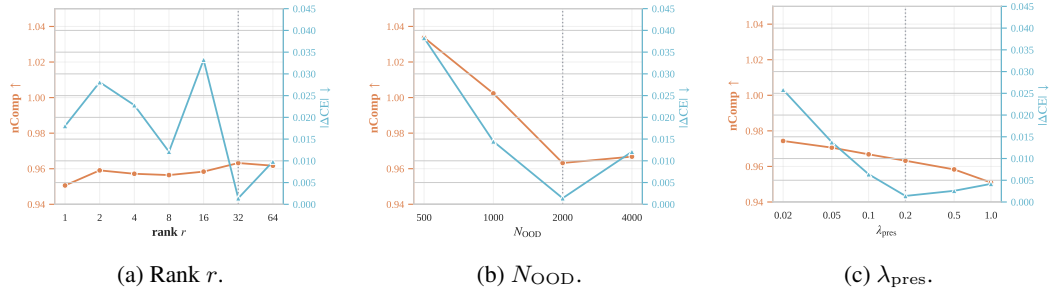


Figure 14: **Hyperparameter sweeps on HaluEval (GPT-2, Transcoder):** nComp (orange, left axis) and  $|\Delta\text{CE}|$  (cyan, right axis) are stable across rank  $r$ , OOD sample size  $N_{\text{OOD}}$ , and preservation weight  $\lambda_{\text{pres}}$ .

811 **Rank  $r$ :** nComp stays above 0.95 for every  $r \in \{1, \dots, 64\}$ ; rank-1 already gives 0.951, confirming  
 812 that the ID-to-OOD drift concentrates in a few directions. **OOD sample size  $N_{\text{OOD}}$ :**  $|\Delta\text{CE}|$  improves  
 813 from 0.038 at  $N = 500$  to 0.001 at  $N \geq 2000$  as the covariance estimate stabilizes. **Preservation**  
 814 **weight  $\lambda_{\text{pres}}$ :** increasing  $\lambda_{\text{pres}}$  trades a small nComp decrease (0.02) for a large  $|\Delta\text{CE}|$  improvement  
 815 (0.026  $\rightarrow$  0.001), with the default near the elbow.