

# Do Sparse Autoencoders Reveal Faithful Concepts? A Neuron-Level Cross-Validation with QA-Based Probes

Anonymous ACL submission

## Abstract

Sparse Autoencoders (SAEs) have recently been adopted to interpret large language models by decomposing hidden activations into sparse “concept features.” These features are often assumed to provide faithful semantic explanations of individual neurons. However, despite their growing use, it remains unclear whether SAE-discovered concepts genuinely reflect neuron semantics or instead arise from reconstruction-induced artifacts. We propose a framework for cross-validating SAE-based neuron interpretations using an independent semantic probe. Specifically, we leverage question–answering embeddings (QA-Emb), in which each feature corresponds to an explicit yes/no semantic query posed to a language model, as an external reference signal. For each neuron, we derive parallel interpretations from SAE features and QA probes, and quantify their semantic agreement using embedding-based similarity and sign-consistency criteria. Across experiments on large-scale news data and multiple auxiliary domains, we find that while many neurons exhibit consistent interpretations, a substantial fraction shows partial or strong disagreement. Further analysis reveals recurring failure modes in which SAE features appear coherent but still lack support from independent QA evidence. Our results demonstrate that independent semantic cross-validation is essential for assessing the trustworthiness of reconstruction-based neuron interpretations in large language models.

## 1 Introduction

Understanding the internal representations of large language models (LLMs) has become a central challenge in natural language processing, both for scientific insight and for practical concerns such as robustness, controllability, and safety. A growing body of work aims to interpret individual neurons or latent dimensions in LLMs by associating them with human-interpretable concepts, often through

analysis of activation patterns, top-ranked tokens, or synthetic stimuli (Kowalska and Kwaśnicka, 2025; Lin et al., 2025). Recently, Sparse Autoencoders (SAEs) have gained prominence as a promising approach to neuron-level interpretability. By learning a sparse, often overcomplete decomposition of hidden activations, SAEs aim to extract a set of “concept features” that can be interpreted independently and composed to explain neuron behavior (Shu et al., 2025). Qualitative analyses of SAE features frequently reveal coherent token sets or natural-language descriptions, leading to the widespread assumption that these features correspond to meaningful semantic units (Cunningham et al., 2023). However, interpretability does not necessarily imply faithfulness. SAE features are learned under reconstruction and sparsity objectives that do not explicitly enforce alignment with neuron semantics (Heap et al., 2025). As a result, it remains unclear whether the concepts identified by SAEs genuinely reflect the functional role of a neuron, or whether they capture abstractions introduced by the autoencoding process itself (Cho et al., 2025). Crucially, most existing evaluations of SAE interpretability are self-referential, relying on the same features or reconstruction statistics used during training, with limited independent validation (Lin and Bloom, 2023; Paulo et al., 2025).

In this work, we ask: *Can we trust the semantic “concept features” discovered by Sparse Autoencoders as faithful interpretations of individual neurons in large language models?* To address this question, we propose a cross-validation framework that evaluates SAE-based neuron interpretations against an independent semantic probe. Specifically, we leverage question–answering embeddings (QA-Emb), a question-driven semantic representation in which each dimension corresponds to the model’s response to an explicit yes/no semantic question. While QA-Emb was originally introduced in a different context, we repurpose it

here purely as an external, human-interpretable semantic probe that is independent of SAE training objectives. A detailed description of QA-Emb is provided in Appendix A. For each neuron, we derive two parallel interpretations: one from SAE-discovered concept features, and one from QA-based semantic probes. We then quantify their agreement using embedding similarity measures and sign-consistency criteria, and categorize neurons into agreement, partial agreement, and disagreement regimes. This setup allows us to move beyond anecdotal interpretability examples and systematically analyze when SAE interpretations align with an independent semantic signal and when they fail. Our experiments reveal that while SAE and QA interpretations often agree, a substantial fraction of neurons exhibit systematic disagreement, even when SAE features appear highly coherent in isolation. Through detailed case studies and controlled faithfulness checks, we identify recurring failure modes in which SAE features assign misleading semantic labels that are not supported by QA evidence. We further analyze how agreement patterns vary across layers and datasets, highlighting the role of inductive biases introduced by reconstruction-based interpretability. Overall, our work makes three contributions:

- We propose a cross-modal validation framework that evaluates neuron interpretations by contrasting reconstruction-based explanations with independent semantic probes.
- We conduct a systematic study of agreement and disagreement between reconstruction-based and probe-based neuron interpretation methods in large language models.
- We identify systematic failure modes where semantically plausible interpretations diverge across methods, and show that faithfulness-based analysis can distinguish interpretations that are more explainable of neuron behavior.

## 2 Related Work

**Neuron and Feature Interpretability in Large Language Models** Interpreting individual neurons and latent features in neural networks has long been a central goal in representation learning and model analysis. Early work focused on visualizing neuron activations through maximally activating inputs, top-ranked tokens, or manually curated examples, aiming to associate neurons with human-

interpretable concepts (Feldhus and Kopf, 2025; Choi et al., 2024). In the context of large language models, neuron-level interpretability has been explored through activation analysis, probing classifiers, and dataset-based inspection methods. These approaches often reveal that individual neurons respond to recurring lexical, syntactic, or semantic patterns, but also highlight the prevalence of polysemantic neurons whose activations cannot be easily explained by a single concept (Bricken et al., 2023; Bills et al., 2023). While these methods provide valuable qualitative insights, they typically rely on heuristic criteria or manual inspection, making it difficult to systematically evaluate the faithfulness of the resulting interpretations (Bills et al., 2023; Choi et al., 2024).

### Sparse Autoencoders for Concept Discovery

Sparse Autoencoders (SAEs) have recently emerged as a prominent approach for interpreting high-dimensional representations in neural networks. A central motivation for SAE-based interpretability is the widely observed polysemanticity of individual neurons, particularly in large language models, where a single neuron may respond to multiple unrelated concepts (Bricken et al., 2023; Bills et al., 2023). SAEs aim to address this issue by decomposing neuron activations into a set of sparse, often overcomplete latent features that are hypothesized to be more nearly monosemantic and thus easier to interpret in isolation (Elhage et al., 2022; Kim et al., 2020). Applied to large language models, SAEs have been used to extract latent “concept features” from hidden layers, with empirical analyses showing that individual features often correspond to coherent sets of tokens or concise natural-language descriptions (Bricken et al., 2023). This has led to the growing practice of treating SAE features as semantically meaningful units and using them as the basis for neuron- and feature-level interpretation. However, SAE-based representations are learned solely through a reconstruction-driven objective under sparsity constraints, without direct supervision or guarantees of semantic faithfulness (Le et al., 2012; Makhzani and Frey, 2014; Song et al., 2025). As a result, while SAE features are often assumed to correspond to approximately monosemantic concepts, the extent to which they faithfully capture the functional semantics of individual neurons remains an open question. Existing evaluations primarily focus on reconstruction quality or internal feature coherence, and rarely

185 incorporate independent semantic validation.

### 186 **Question-Based and Probe-Based Interpretability**

187 An alternative line of work studies interpretability through explicit probes, in which pre-  
188 defined semantic properties are tested against model  
189 representations. Probing classifiers and diagnostic  
190 tasks have been widely used to assess whether specific  
191 linguistic or semantic information is encoded  
192 in model activations (Alain and Bengio, 2018; Be-  
193 linkov, 2021; Conneau et al., 2018). More recently,  
194 question-based approaches have been proposed to  
195 construct interpretable feature spaces by asking  
196 language models explicit yes/no questions and using  
197 their responses as semantic features (Benara  
198 et al., 2024). These methods offer transparent and  
199 human-readable representations, as each feature  
200 corresponds directly to a semantic query. While  
201 question-based probes provide explicit semantic  
202 grounding, they are not designed to reconstruct hidden  
203 representations or explain individual neurons  
204 in isolation. Instead, they serve as an external semantic  
205 signal with inductive biases distinct from  
206 reconstruction-based methods.  
207

208 **Positioning of the Present Work** Prior work has  
209 largely studied reconstruction-based interpretability  
210 methods, such as sparse autoencoders, and  
211 probe-based semantic representations in isolation.  
212 In particular, SAE features are often implicitly  
213 treated as faithful, approximately monosemantic  
214 decompositions of polysemantic neurons, despite  
215 being optimized without explicit semantic supervision.  
216 In contrast, our work focuses on systematically  
217 cross-validating this assumption by comparing  
218 SAE-based neuron interpretations with independent,  
219 question-based semantic probes. By analyzing when  
220 these two signals agree or diverge, we aim to assess  
221 the reliability of SAE-based interpretations and to  
222 identify recurring failure modes in neuron-level  
223 semantic explanations.

## 224 **3 Method**

225 We propose a neuron interpretation framework  
226 that cross-validates activation-grounded semantic  
227 features with independent question-based semantic  
228 probes. At a high level, our method derives two  
229 complementary interpretations for each neuron—  
230 one based on sparse autoencoder (SAE) features  
231 and one based on QA-based semantic probes—and  
232 evaluates their semantic alignment to assess  
233 interpretability and faithfulness. Figure 1

234 provides an overview of the full pipeline. Given a  
235 collection of input text spans, neuron activations  
236 are processed along two parallel pathways. The upper  
237 pathway extracts SAE features that are directly  
238 grounded in model activations, while the lower  
239 pathway derives QA-based semantic hypotheses by  
240 correlating neuron activations with predefined semantic  
241 questions. These two interpretations are then aligned  
242 in a shared semantic space, cross-filtered to identify  
243 consistently supported concepts, and optionally  
244 synthesized into a refined neuron-level interpretation  
245 using a constrained LLM.

### 246 **3.1 SAE-Based Neuron Interpretation**

247 Let  $h \in \mathbb{R}^d$  denote the hidden activation vector at  
248 a fixed layer of a language model. A sparse autoencoder  
249 (SAE) maps  $h$  to a sparse latent representation  
250  $z \in \mathbb{R}^{d_{\text{sae}}}$ , from which  $h$  can be approximately  
251 reconstructed. For a given neuron  $n$ , we identify  
252 the SAE latent features whose decoder weights  
253 have the strongest influence on that neuron. These  
254 latent features are treated as candidate semantic  
255 explanations for the neuron’s behavior. Each SAE  
256 feature is associated with a set of top-activating  
257 tokens or text spans, and may additionally include  
258 a short natural-language description derived from  
259 these activations. Aggregating the most influential  
260 SAE features yields an activation-grounded, data-  
261 driven semantic interpretation for each neuron.

### 262 **3.2 QA-Based Neuron Interpretation**

263 In parallel, we derive an independent semantic  
264 interpretation using QA-based semantic probes. Intuitively,  
265 QA-based interpretation treats neuron activity as  
266 being explained by explicit semantic questions, such  
267 as whether an input expresses a topic, attribute, or  
268 relation (Benara et al., 2024). More details about  
269 this method can be found in Appendix B. Each input  
270 span is associated with a QA embedding, where each  
271 dimension corresponds to a predefined binary semantic  
272 question. These questions are answered by a language  
273 model, producing a sparse vector of yes/no responses  
274 that serves as an interpretable semantic representation  
275 (see Figure 1, top branch). For a fixed neuron, we  
276 measure the statistical association between the neuron’s  
277 activation values and each QA dimension across inputs.  
278 The QA-based interpretation of a neuron is defined  
279 by the subset of questions whose responses exhibit  
280 the strongest positive or negative association with  
281 that neuron’s activations. These questions serve  
282 as explicit semantic hypotheses about the neuron’s  
283

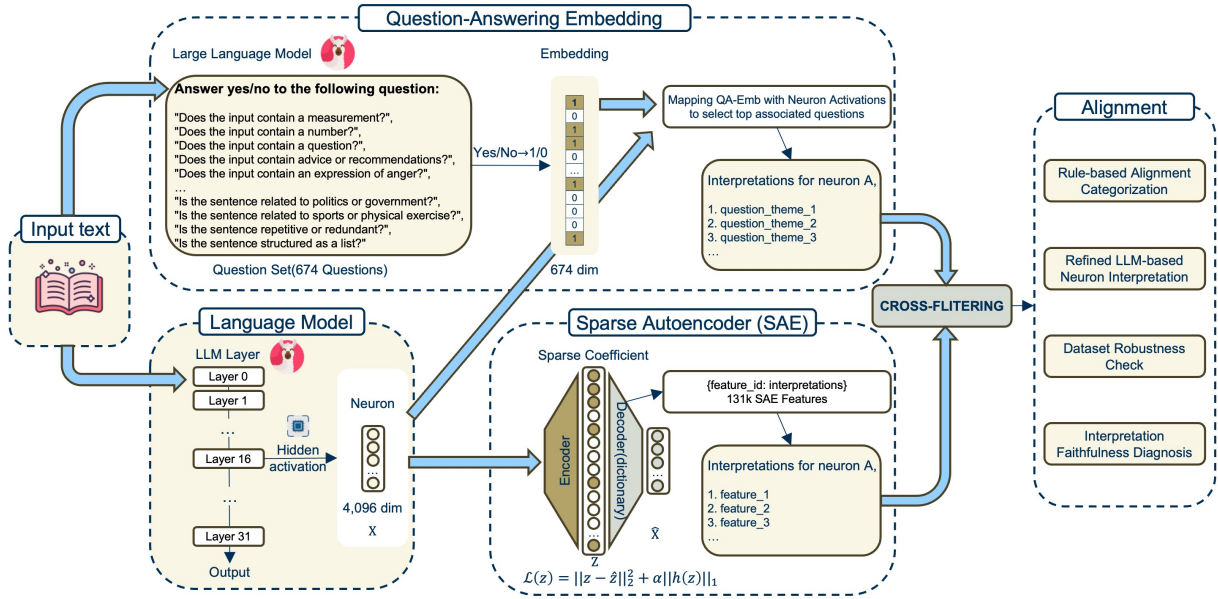


Figure 1: Overview of the proposed SAE-QA cross-validation framework. For each neuron, semantic interpretations are derived through two independent pathways: sparse autoencoder (SAE) features grounded in model activations, and QA-based semantic probes that provide external semantic hypotheses. Their semantic alignment is evaluated to assess agreement, disagreement, and faithfulness of neuron-level interpretations. Optionally, cross-filtered evidence from both pathways can be synthesized into a refined interpretation using a constrained language model.

functional role. Unlike SAE features, QA probes are predefined and independent of neuron activations, providing an external semantic signal for interpretation.

### 3.3 Semantic Alignment Between SAE and QA Interpretations

To assess whether SAE-based and QA-based interpretations describe a common underlying semantic concept, we compare their semantic content in a shared embedding space. Specifically, we compute semantic similarity between textual representations derived from SAE features (e.g., feature descriptions or token sets) and QA question texts. In addition to semantic similarity, we consider the directionality of neuron activations. This allows us to distinguish cases where the semantic content aligns with activation behavior from cases where similar semantics correspond to opposite activation directions. Based on these criteria, neurons can exhibit varying degrees of agreement or disagreement between SAE- and QA-based interpretations.

### 3.4 Cross-Filtering and Interpretation

While SAE and QA interpretations provide complementary semantic evidence, each alone may contain spurious or weakly grounded concepts. We therefore apply a cross-filtering step that retains only semantic elements that are consistently

supported by both methods. This step identifies a shared semantic subspace between activation-grounded SAE features and externally defined QA probes. To synthesize a unified neuron-level interpretation, we employ a constrained LLM-based aggregation step. The language model is provided with: (i) cross-filtered SAE token evidence directly grounded in neuron activations, (ii) cross-filtered QA questions serving as independent semantic confirmation, and (iii) SAE feature descriptions capturing higher-level abstractions. The LLM is explicitly constrained to produce a concise semantic summary that is strictly supported by the provided evidence and may not introduce new semantic content. In this framework, the LLM functions as a semantic aggregator rather than a free-form generator.

### 3.5 Evaluation Overview

To evaluate the quality of refined neuron interpretations, we introduce the **Qualitative Interpretation Score (QIS)**, a composite metric designed to assess semantic coherence, cross-method agreement, grounding in evidence, and semantic focus. QIS is used exclusively for evaluation and does not influence the construction of neuron interpretations. Full definitions and implementation details are provided in the experimental section and Appendix F.

338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384

## 4 Experiments

This section describes the experimental setup used to instantiate and evaluate the proposed SAE–QA cross-validation framework, including model selection, datasets, implementation details, and evaluation protocols. Code will be released upon acceptance.

### 4.1 Model and Layers

All experiments are conducted on meta-llama/LLaMA-3.1-8B, which has a hidden dimensionality of  $d_{in} = 4096$  (AI@Meta, 2024). Unless otherwise specified, our primary analysis focuses on a middle layer (layer 16), which prior work has suggested to contain rich semantic information (Skean et al., 2024, 2025). To assess layer-wise variation and faithfulness, we additionally analyze an early layer (layer 4) and a late layer (layer 28) in selected experiments.

### 4.2 Datasets

Our primary dataset consists of 9,800 text segments sampled from CC-News, a large-scale news corpus (Mackenzie et al., 2020). To evaluate robustness across domains and languages, we additionally include auxiliary datasets covering financial text, scientific abstracts, programming-related text, and non-English summaries. Dataset details are provided in Appendix G. For each dataset, we extract fixed-length token snippets of 64 tokens and segment each snippet into overlapping spans of 16 tokens using a sliding window, resulting in 49 spans per snippet. Each span is treated as an independent input for neuron activation analysis.

### 4.3 Representations

**Hidden State Extraction** For each input span, we extract hidden activations from the selected layer(s) using SAELens (Bloom et al., 2024). This yields an activation matrix of size  $[N, d_{in}]$ , where  $N$  denotes the number of input spans.

**Sparse Autoencoder Configuration** We use pre-trained Sparse Autoencoders (SAEs) from *Llama Scope* to analyze hidden activations of the selected model and layer (He et al., 2024). Each SAE maps residual stream activations to a higher-dimensional sparse latent space. For each neuron, we select the top-ranked SAE features based on the absolute decoder weight magnitude connecting latent features to that neuron.

**QA-Based Semantic Probes** Each input span is associated with a QA embedding derived from a fixed set of binary yes/no questions (Benara et al., 2024). For each neuron, we compute Pearson correlation coefficients between its activation values and each QA dimension across inputs. QA dimensions with the strongest positive and negative correlations are selected to form the neuron’s QA-based interpretation.

### 4.4 Implementation Details

To compare SAE- and QA-based interpretations, we embed textual representations using all-mpnet-base-v2, a sentence embedding model (Reimers and Gurevych, 2019) and compute cosine similarity in the resulting semantic space. All similarity thresholds, correlation settings, and selection hyperparameters are fixed across experiments. Exact values are reported below.

### 4.5 Evaluation Protocol

**Agreement Categories** Neurons are assigned to one of four categories based on semantic similarity and sign consistency between SAE and QA interpretations: *agreement*, *partial agreement*, *weak agreement*, or *disagreement*. Neurons with cosine similarity  $\geq 0.45$  are labeled as *agreement*, those with similarity in  $[0.30, 0.45)$  as *partial agreement*, and those with similarity  $< 0.30$  as *weak agreement*. Cases with high semantic similarity but opposite activation signs are labeled as *disagreement*. These thresholds are selected empirically based on preliminary analysis, balancing semantic coherence and coverage across neurons. We observe that the overall agreement trends and qualitative conclusions are robust to moderate variations of these thresholds. Additional details on threshold selection are provided in Appendix C.

**Qualitative Interpretation Score (QIS)** We evaluate refined neuron interpretations using QIS, which measures: (i) internal semantic consistency, (ii) semantic overlap between SAE and QA evidence, (iii) grounding of the refined summary in the provided evidence, and (iv) semantic focus. QIS is computed as a fixed weighted combination of these components and is used solely for evaluation.

**Visualization** We visualize SAE–QA alignment using similarity matrices, token-level UMAP overlap plots, and neuron-wise agreement distributions.

385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

## 5 Results and Analysis

### 5.1 Overall SAE-QA Alignment

We begin by measuring semantic agreement between SAE-based and QA-based interpretations on the main dataset at the primary analysis layer (Layer 16). Each neuron is assigned to one of four agreement categories based on semantic similarity and sign consistency, as defined in Section 4.5. Table 1 summarizes the distribution of neurons across agreement categories in this primary setting. A

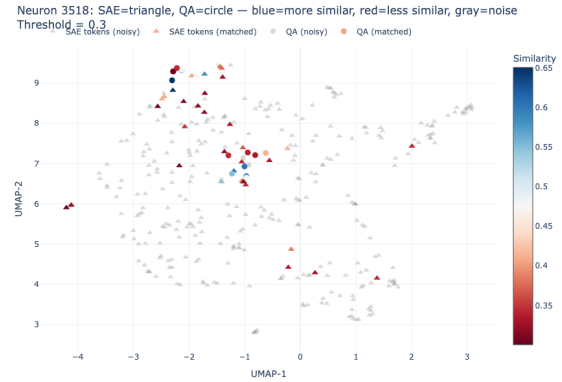
Agreement Category	Proportion of Neurons (%)
Agreement	29.3
Partial Agreement	60.8
Weak Agreement	9.9
Disagreement	< 0.1

Table 1: Distribution of neurons across agreement categories for the main dataset at the middle layer.

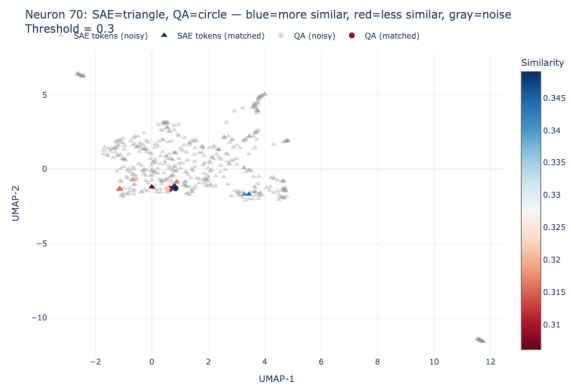
substantial majority of neurons fall into the agreement or partial agreement regimes, indicating that SAE features often capture semantic information consistent with independent QA probes. A smaller but non-negligible fraction of neurons exhibit weak agreement, while explicit disagreement cases are exceedingly rare. Interestingly, in the few disagreement cases we observe, the associated SAE features often appear semantically coherent when viewed on their own, for example by inspecting their top activating tokens. However, these apparent interpretations are not supported by the neuron’s activation patterns as measured by QA probes. This suggests that some SAE features can be intuitively interpretable at the surface level, which is patterns tied to the literal form of the input rather than a single underlying semantic function—into a single abstraction, yet fail to faithfully reflect the semantic factors that actually drive neuron behavior.

**Qualitative Alignment Patterns** To better understand what semantic agreement entails, we visualize the alignment between SAE-derived semantic representations and QA question texts using a cross-filtered UMAP projection computed per neuron. In each view, points below the similarity threshold are rendered in gray, while above-threshold points are highlighted to reveal whether SAE- and QA-derived semantics co-localize. Figure 2 shows two representative neurons: one with strong co-localization (agreement) and one with limited overlap (weak alignment). An interactive version of these visualizations (clickable points with linked

SAE/QA matches) is provided in the supplementary material. Neurons in the agreement regime



(a) **Well-aligned neuron.** SAE-highlighted tokens and QA questions overlap strongly in the high-similarity region.



(b) **Poorly-aligned neuron.** Limited overlap between SAE and QA; highlighted regions are largely method-specific.

Figure 2: Cross-filtered UMAP views for two representative neurons on the main dataset (Layer 16). Points below the similarity threshold are gray; above-threshold points are highlighted, allowing visual inspection of whether SAE and QA derived semantics co-localize.

typically exhibit consistent semantic themes across activation-grounded SAE tokens, SAE feature descriptions, and QA questions. These cases suggest that reconstruction-based and question-based interpretability can converge on a shared semantic explanation when neuron behavior is well captured by both signals. To better understand why such convergence fails for other neurons, we examine cases of weak alignment and disagreement in more detail. We find that weakly aligned neurons often display method-specific semantic clusters: SAE features may form coherent and interpretable descriptions when examined in isolation, yet these descriptions are not supported by the neuron’s activation patterns as measured by QA probes. This observation indicates that a coherent SAE feature

description does not necessarily imply behavioral faithfulness, and motivates a systematic analysis of failure modes enabled by our cross-validation framework.

### Failure Modes of SAE-Based Interpretability

A central contribution of our framework is its ability to expose systematic failure modes of SAE-based neuron interpretation. Figure 3 provides an annotated example that illustrates how such failures manifest in practice. The highlighted neu-

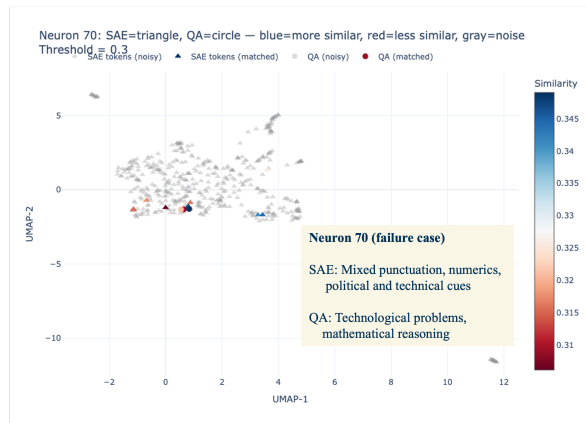


Figure 3: Annotated failure case illustrating a representative SAE-QA disagreement. The circled point corresponds to Neuron 70, previously shown in the interactive visualization (Figure 2(b)). Despite a semantically rich SAE feature description, the SAE-highlighted region shows little overlap with QA-derived semantics, which instead indicate a narrower functional role centered on technological and mathematical reasoning.

ron (Neuron 70) exhibits a semantically rich and internally coherent SAE feature description, yet shows minimal alignment with QA-based evidence or activation behavior. Closer inspection reveals that the SAE interpretation aggregates heterogeneous surface-level regularities. These regularities include punctuation usage, numerical formatting patterns, and loosely related topical cues that frequently co-occur in the data but do not correspond to a unified semantic role. In contrast, QA probes consistently indicate a narrower functional role related to technological problem-solving and mathematical reasoning. This discrepancy suggests that reconstruction objectives can induce abstractions that appear meaningful but are not faithful to the neuron’s causal contribution. An interactive version of this example, allowing inspection of individual points and their corresponding SAE and QA matches, is provided in the supplementary material. Additional failure cases and detailed token-level

analyses are provided in the appendix I.

### Constrained LLM-Based Neuron Interpretation

To synthesize a unified neuron-level explanation from cross-validated evidence, we apply a constrained LLM-based semantic aggregation step. The LLM integrates cross-filtered SAE tokens, QA questions, and SAE feature descriptions into a concise interpretation supported by activation-grounded and independently validated evidence. Figure 4 shows the distribution of QIS values for

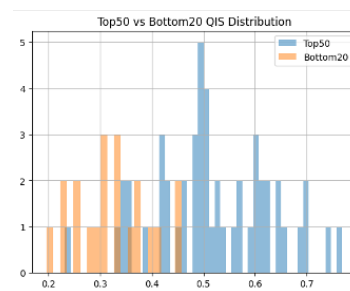


Figure 4: Distribution of Qualitative Interpretation Score (QIS) values for neurons grouped by SAE-QA similarity ranking. Higher QIS values indicate more coherent and well-grounded neuron interpretations.

neurons ranked by SAE-QA semantic similarity. Neurons in the top similarity group exhibit substantially higher mean QIS values (approximately  $\sim 0.52$ ) compared to neurons in the bottom group (approximately  $\sim 0.32$ ). This result indicates that neurons with stronger SAE-QA alignment tend to admit clearer and more coherent interpretations after refinement. Importantly, we also observe a non-trivial overlap in QIS distributions between the two groups. We find that LLM-aggregated interpretations are consistently more specific and faithful than interpretations derived from either SAE or QA alone. An illustrative example is provided in Appendix E. In agreement cases, aggregation reinforces shared semantics. In disagreement cases, aggregation resolves ambiguities by prioritizing activation-grounded and cross-validated evidence, often discarding misleading abstractions introduced by SAE features.

### 5.2 Dataset Robustness Analysis

Finally, we assess whether the observed agreement patterns generalize beyond the primary dataset. We repeat the SAE-QA cross-validation procedure on several auxiliary datasets spanning different domains and languages. Across datasets, we observe qualitatively similar distributions of agreement and

disagreement, which means over 90% of neurons in each dataset show at least one aligned SAE-QA concept, indicating that the SAE-QA alignment is not dependent on dataset-specific lexical statistics, but instead captures higher-level neuron behaviors that generalize across domains. We further analyze the overlap of aligned neurons across datasets to assess whether the same neurons tend to express similar semantic concepts in different corpora. When comparing each auxiliary dataset against CCNews, we observe that approximately 4,000 neurons (out of the full layer population) exhibit overlapping aligned SAE-QA concepts across datasets. This indicates that the overlap is not driven by a single dominant dataset, but reflects a broadly shared set of neuron-level semantic units that persist across domains. Taken together, these results provides strong evidence that the proposed interpretation pipeline captures intrinsic neuron semantics rather than dataset-specific patterns.

## 6 Faithfulness-Oriented Diagnostics

To assess faithfulness more directly, we design a controlled diagnostic that evaluates whether competing interpretations better predict neuron activation behavior under targeted semantic interventions. Rather than assuming either SAE- or QA-based interpretations to be correct, we test each interpretation as a falsifiable hypothesis.

**Diagnostic Setup** For neurons exhibit failure modes, we prompt a language model to propose candidate semantic hypotheses corresponding to the SAE and QA interpretations, respectively. These hypotheses serve as competing explanations for the neuron’s activation behavior. For each hypothesis pair, we construct controlled text subsets that isolate the corresponding semantic dimensions. In particular, we generate  $2 \times 2$  contrastive text conditions that independently vary the two hypothesized factors (e.g., language form and religious content), yielding four conditions that differ along exactly one dimension at a time. Neuron activation is measured using a simple detector defined as the maximum activation over tokens in the input text. We evaluate each hypothesis by computing the area under the ROC curve (AUC), where activation scores are used to distinguish positive and negative contrast conditions. An AUC of 1 indicates that the hypothesis fully explains the neuron’s behavior, an AUC of 0.5 indicates no predictive power, and an AUC of 0 indicates an incorrect explanation.

**Layer-wise Faithfulness Outcomes** Applying this diagnostic to neurons with SAE-QA disagreement reveals clear layer-dependent patterns. Among middle-layer neurons, a clear majority (~60%) are more faithfully explained by QA-based hypotheses under controlled diagnostics. In contrast, late-layer neurons more frequently favor SAE-based interpretations, with approximately two-thirds of neurons exhibiting higher AUC scores for SAE-based hypotheses. Early-layer neurons often yield ambiguous outcomes, with nearly half of neurons showing no clear predictive advantage for either interpretation. Representative neuron-level case studies illustrating each outcome are provided in Appendix J.

**Implications** While individual faithfulness diagnostics yield mixed outcomes, a consistent pattern emerges when results are grouped by layer. For neurons in the middle layer, failure cases identified via SAE-QA disagreement are more often resolved in favor of QA-based interpretations under controlled diagnostics. In contrast, early-layer neurons frequently yield ambiguous results, with neither method demonstrating clear predictive superiority. For late-layer neurons, SAE-based interpretations more often align with activation behavior under contrastive testing. These observations suggest that faithfulness is not a property of an interpretation method in isolation, but rather a function of its interaction with layer-specific representational characteristics.

## 7 Conclusion

Overall, our results demonstrate the importance of independent semantic cross-validation for neuron-level interpretability. By comparing reconstruction-based interpretations with external semantic probes, we show that apparent interpretability does not necessarily imply faithfulness, and that systematic failure modes can be revealed through cross-modal discrepancies. We further introduce a constrained LLM-based aggregation procedure and a quantitative metric, QIS, to assess the coherence, grounding, and focus of neuron interpretations. Taken together, our findings suggest that interpretability quality depends on more than raw semantic similarity, highlighting the need for principled diagnostic tools when evaluating neuron explanations in large language models.

## 656 Limitations

657 Overall, our results demonstrate the importance of  
658 independent semantic cross-validation for neuron-  
659 level interpretability. By comparing reconstruction-  
660 based interpretations with external semantic probes,  
661 we show that apparent interpretability does not nec-  
662 essarily imply faithfulness, and that systematic fail-  
663 ure modes can be revealed through cross-modal  
664 discrepancies. We further introduce a constrained  
665 LLM-based semantic aggregation procedure and  
666 a quantitative metric, QIS, to evaluate the coher-  
667 ence, grounding, and focus of neuron interpreta-  
668 tions. Together, these findings suggest that inter-  
669 pretability quality depends on more than raw se-  
670 mantic similarity, and that neuron explanations can  
671 be assessed in a principled and quantitative man-  
672 ner. Rather than identifying a universally superior  
673 method, our results highlight the need for diag-  
674 nostic tools that account for interactions between  
675 interpretation methods and model representations,  
676 especially as interpretability claims are applied to  
677 larger and higher-stakes models.

## 678 References

679 AI@Meta. 2024. [Llama 3 model card](#).

680 Guillaume Alain and Yoshua Bengio. 2018. [Under-  
681 standing intermediate layers using linear classifier  
682 probes](#). *Preprint*, arXiv:1610.01644.

683 Yonatan Belinkov. 2021. [Probing classifiers:  
684 Promises, shortcomings, and advances](#). *Preprint*,  
685 arXiv:2102.12452.

686 Vinamra Benara, Chandan Singh, John X. Morris,  
687 Richard Antonello, Ion Stoica, Alexander G. Huth,  
688 and Jianfeng Gao. 2024. [Crafting interpretable  
689 embeddings by asking llms questions](#). *Preprint*,  
690 arXiv:2405.16714.

691 Steven Bills, Nick Cammarata, Dan Moss-  
692 ing, Henk Tillman, Leo Gao, Gabriel Goh,  
693 Ilya Sutskever, Jan Leike, Jeff Wu, and  
694 William Saunders. 2023. Language mod-  
695 els can explain neurons in language models.  
696 [https://openai-public.blob.core.windows.  
697 net/neuron-explainer/paper/index.html](https://openai-public.blob.core.windows.net/neuron-explainer/paper/index.html).

698 Joseph Bloom, Curt Tigges, Anthony Duong, and David  
699 Chanin. 2024. Saelens. [https://github.com/  
700 decoderresearch/SAELens](https://github.com/decoderresearch/SAELens).

701 Trenton Bricken, Adly Templeton, Joshua Batson,  
702 Brian Chen, Adam Jermy, Tom Conerly, Nick  
703 Turner, Cem Anil, Carson Denison, Amanda Askell,  
704 Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas  
705 Schiefer, Tim Maxwell, Nicholas Joseph, Zac

Hatfield-Dodds, Alex Tamkin, Karina Nguyen, and  
6 others. 2023. Towards monosemanticity: Decom-  
posing language models with dictionary learning.  
*Transformer Circuits Thread*. [https://transformer-  
circuits.pub/2023/monosemantic-  
features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).

Seonglae Cho, Harryn Oh, Donghyun Lee, Luis Ed-  
uardo Rodrigues Vieira, Andrew Bermingham, and  
Ziad El Sayed. 2025. [Faithfulsae: Towards cap-  
turing faithful features with sparse autoencoders  
without external dataset dependencies](#). *Preprint*,  
arXiv:2506.17673.

Dami Choi, Vincent Huang, Kevin Meng, Daniel D  
Johnson, Jacob Steinhardt, and Sarah Schwettmann.  
2024. Scaling automatic neuron description. [https:  
//transluce.org/neuron-descriptions](https://transluce.org/neuron-descriptions).

Alexis Conneau, German Kruszewski, Guillaume Lam-  
ple, Loïc Barrault, and Marco Baroni. 2018. [What  
you can cram into a single vector: Probing sen-  
tence embeddings for linguistic properties](#). *Preprint*,  
arXiv:1805.01070.

Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert  
Huben, and Lee Sharkey. 2023. [Sparse autoencoders  
find highly interpretable features in language models](#).  
*Preprint*, arXiv:2309.08600.

Nelson Elhage, Tristan Hume, Catherine Olsson,  
Nicholas Schiefer, Tom Henighan, Shauna Kravec,  
Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain,  
Carol Chen, Roger Grosse, Sam McCandlish, Jared  
Kaplan, Dario Amodei, Martin Wattenberg, and  
Christopher Olah. 2022. Toy models of superposition.  
*Transformer Circuits Thread*. [https://transformer-  
circuits.pub/2022/toy\\_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html).

Nils Feldhus and Laura Kopf. 2025. [Interpreting lan-  
guage models through concept descriptions: A survey](#).  
*Preprint*, arXiv:2510.01048.

Zhengfu He, Wentao Shu, Xuyang Ge, Lingjie Chen, Junx-  
uan Wang, Yunhua Zhou, Frances Liu, Qipeng Guo, Xu-  
anjing Huang, Zuxuan Wu, Yu-Gang Jiang, and Xipeng  
Qiu. 2024. [Llama scope: Extracting millions of features  
from llama-3.1-8b with sparse autoencoders](#). *Preprint*,  
arXiv:2410.20526.

Thomas Heap, Tim Lawson, Lucy Farnik, and Lau-  
rence Aitchison. 2025. [Sparse autoencoders can in-  
terpret randomly initialized transformers](#). *Preprint*,  
arXiv:2501.17727.

Baotian Hu, Qingcai Chen, and Fangze Zhu. 2015. [LCSTS:  
A large scale Chinese short text summarization dataset](#).  
In *Proceedings of the 2015 Conference on Empirical  
Methods in Natural Language Processing*, pages 1967–  
1972, Lisbon, Portugal. Association for Computational  
Linguistics.

Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis  
Allamanis, and Marc Brockschmidt. 2019. CodeSearch-  
Net challenge: Evaluating the state of semantic code  
search. *arXiv preprint arXiv:1909.09436*.

754	Edward Kim, Connor Onweller, Andrew O'Brien, and Kathleen McCoy. 2020. <a href="#">The interpretable dictionary in sparse coding</a> . <i>Preprint</i> , arXiv:2011.11805.	Oscar Skean, Md Rifat Arefin, Dan Zhao, Niket Patel, Jalal Naghiyev, Yann LeCun, and Ravid Shwartz-Ziv. 2025. <a href="#">Layer by layer: Uncovering hidden representations in language models</a> . <i>Preprint</i> , arXiv:2502.02013.	806
755			807
756			808
757	Bianka Kowalska and Halina Kwaśnicka. 2025. <a href="#">Unboxing the black box: Mechanistic interpretability for algorithmic understanding of neural networks</a> . <i>Preprint</i> , arXiv:2511.19265.	Xiangchen Song, Aashiq Muhamed, Yujia Zheng, Lingjing Kong, Zeyu Tang, Mona T. Diab, Virginia Smith, and Kun Zhang. 2025. <a href="#">Position: Mechanistic interpretability should prioritize feature consistency in saes</a> . <i>Preprint</i> , arXiv:2505.20254.	810
758			811
759			812
760			813
761	Quoc V. Le, Marc'Aurelio Ranzato, Rajat Monga, Matthieu Devin, Kai Chen, Greg S. Corrado, Jeff Dean, and Andrew Y. Ng. 2012. <a href="#">Building high-level features using large scale unsupervised learning</a> . <i>Preprint</i> , arXiv:1112.6209.	Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. <a href="#">Benchmarking retrieval-augmented generation for medicine</a> . <i>arXiv preprint arXiv:2402.13178</i> .	814
762			815
763			816
764			817
765			
766	Johnny Lin and Joseph Bloom. 2023. <a href="#">Neuronpedia: Interactive reference and tooling for analyzing neural networks</a> . <i>Software available from neuronpedia.org</i> .	<b>A Background on Question–Answering Embeddings</b>	818
767			819
768			
769	Zihao Lin, Samyadeep Basu, Mohammad Beigi, Varun Manjunatha, Ryan A. Rossi, Zichao Wang, Yufan Zhou, Sriram Balasubramanian, Arman Zarei, Keivan Rezaei, Ying Shen, Barry Menglong Yao, Zhiyang Xu, Qin Liu, Yuxiang Zhang, Yan Sun, Shilong Liu, Li Shen, Hongxuan Li, and 2 others. 2025. <a href="#">A survey on mechanistic interpretability for multi-modal foundation models</a> . <i>Preprint</i> , arXiv:2502.17516.	Question–answering embeddings (QA-Emb) were introduced by <a href="#">Benara et al. (2024)</a> as a framework for constructing explicitly interpretable semantic representations of language model behavior. Rather than learning latent dimensions whose meanings must be inferred post hoc, QA-Emb defines each representational dimension a priori via a human-readable semantic question.	820
770			821
771			822
772			823
773			824
774			825
775			826
776			827
777	Joel Mackenzie, Rodger Benham, Matthias Petri, Johanne R. Trippas, J. Shane Culpepper, and Alistair Moffat. 2020. <a href="#">Cc-news-en: A large english news corpus</a> . In <i>Proceedings of the 29th ACM International Conference on Information &amp; Knowledge Management, CIKM '20</i> , page 3077–3084, New York, NY, USA. Association for Computing Machinery.	At a high level, QA-Emb shifts the locus of interpretability from representation geometry to question design. Each dimension corresponds to a fixed yes/no query (e.g., whether an input concerns a person, an event, or an affective property), and the embedding value reflects the language model's response to that query. This design makes the semantic meaning of each coordinate transparent by construction, enabling direct inspection without auxiliary probing or interpretation models.	828
778			829
779			830
780			831
781			832
782			833
783			834
784	Alireza Makhzani and Brendan Frey. 2014. <a href="#">k-sparse autoencoders</a> . <i>Preprint</i> , arXiv:1312.5663.	QA-Emb was originally proposed as a general-purpose semantic representation for analyzing language model behavior across tasks. Importantly, it does not assume that the resulting question set is exhaustive, complete, or cognitively grounded. Instead, its primary goal is transparency: each dimension encodes a clearly specified semantic hypothesis that can be individually examined, validated, or rejected.	835
785			836
786	P. Malo, A. Sinha, P. Korhonen, J. Wallenius, and P. Takala. 2014. <a href="#">Good debt or bad debt: Detecting semantic orientations in economic texts</a> . <i>Journal of the Association for Information Science and Technology</i> , 65.	Because QA-Emb does not rely on reconstruction objectives, sparsity constraints, or representation learning, it embodies inductive biases that differ fundamentally from feature-discovery methods such as Sparse Autoencoders. This conceptual separation makes QA-Emb well suited as an external semantic reference for diagnostic analysis, even when it is not treated as an authoritative or ground-truth semantic decomposition.	837
787			838
788			839
789			840
790	Gonçalo Paulo, Alex Mallen, Caden Juang, and Nora Belrose. 2025. <a href="#">Automatically interpreting millions of features in large language models</a> . <i>Preprint</i> , arXiv:2410.13928.		841
791			842
792			843
793			844
794	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentencebert: Sentence embeddings using siamese bert-networks</a> . <i>Preprint</i> , arXiv:1908.10084.		845
795			846
796			847
797	Dong Shu, Xuansheng Wu, Haiyan Zhao, Daking Rai, Ziyu Yao, Ninghao Liu, and Mengnan Du. 2025. <a href="#">A survey on sparse autoencoders: Interpreting the internal mechanisms of large language models</a> . <i>Preprint</i> , arXiv:2503.05613.		848
798			849
799			850
800			851
801			852
802	Oscar Skean, Md Rifat Arefin, Yann LeCun, and Ravid Shwartz-Ziv. 2024. <a href="#">Does representation matter? exploring intermediate layers in large language models</a> . <i>Preprint</i> , arXiv:2412.09563.		853
803			854
804			855
805			856
			857

858	of human semantics or neuron-level ground truth,	effects on category distributions and qualitative in-	906
859	but as a transparent, question-driven semantic lens	terpretability.	907
860	that can be contrasted with reconstruction-based	We find that a threshold of 0.45 for strong agree-	908
861	interpretations to expose differences in inductive	ment yields a stable set of neurons with high se-	909
862	bias.	semantic coherence between SAE- and QA-based	910
863	<b>B QA-Based Semantic Embeddings</b>	interpretations, while avoiding overly restrictive fil-	911
864	We use question–answering embeddings (QA-	tering. Lower thresholds increase coverage but ad-	912
865	Emb) as an independent semantic probe for neuron-	mit semantically ambiguous cases, whereas higher	913
866	level interpretation. Unlike reconstruction-based	thresholds reduce the number of interpretable neu-	914
867	methods, QA-Emb provides a fixed, externally	rons without substantially improving semantic con-	915
868	specified semantic feature space whose dimensions	sistency.	916
869	are defined by explicit yes/no questions, and are	Importantly, we observe that the main agreement	917
870	therefore independent of model training or repre-	trends and qualitative failure modes discussed in	918
871	sentation learning objectives.	the paper remain consistent across a reasonable	919
872	In our framework, QA-Emb is not used as a task	range of threshold values, indicating that our con-	920
873	representation or predictive embedding. Instead, it	clusions are not sensitive to the exact choice of	921
874	serves as a source of semantic hypotheses against	thresholds.	922
875	which neuron behavior can be evaluated. For each	<b>D LLM-Based Semantic Interpretation</b>	923
876	input span, we obtain a binary-valued QA embed-	<b>Prompt</b>	924
877	ding by prompting a language model to answer	<b>Example Neuron Interpretation Prompt</b>	925
878	a predefined set of semantic questions. Each di-	Below we show an example of the prompt used to ge-	926
879	dimension reflects the model’s response to a specific	nerate a refined semantic interpretation for a single	927
880	semantic query.	neuron. The example illustrates how cross-filtered	928
881	To derive a QA-based interpretation for a neu-	tokens and QA questions are treated as primary ev-	929
882	ron, we measure the statistical association between	idence, while SAE feature descriptions serve only	930
883	the neuron’s activation values and each QA dimen-	as secondary context.	931
884	sion across the dataset. Dimensions exhibiting	You are an interpretability assistant	932
885	strong positive or negative association are inter-	for language models.	933
886	preted as candidate semantic factors influencing	Your task is to summarize the semantic	934
887	the neuron’s behavior. This procedure yields a	role of a neuron based on evidence	935
888	neuron-level semantic profile expressed entirely in	from two independent sources.	936
889	question-defined coordinates.	IMPORTANT:	937
890	Crucially, QA-Emb is treated as a non-	- Cross-filtered TOKENS and QA QUESTIONS	938
891	authoritative reference rather than a ground-truth	are the highest-trust evidence.	939
892	explanation. Agreement between QA-based and	- SAE FEATURE DESCRIPTIONS are secondary	940
893	SAE-based interpretations suggests convergent ev-	and must not override shared	941
894	idence for a shared semantic factor, whereas dis-	evidence.	942
895	agreement highlights differences in inductive bias	NEURON_ID: 3518	943
896	between question-driven and reconstruction-driven	CROSS-FILTERED QA QUESTIONS:	944
897	interpretability methods. Our analysis focuses on	1. an expression of anger	945
898	these agreement and divergence patterns rather than	2. body of water	946
899	privileging either representation a priori.	3. technical or specialized terminology	947
900	<b>C Threshold Selection for Semantic</b>	4. technological problem or bug	948
901	<b>Agreement</b>	CROSS-FILTERED TOKENS:	949
902	The cosine similarity thresholds used to define	Engineers, River, country, haven, island	950
903	agreement categories are chosen empirically. We	, islands, ire, land	951
904	evaluated a range of threshold values for high and	SAE FEATURE DESCRIPTIONS:	952
905	moderate semantic similarity and inspected their	1. specific locations and addresses	953
		related to events or activities	954
		2. terms related to functional	955
		programming concepts and operations	956
		3. references to military ranks and	957
		units	958
			959
			960
			961
			962
			963

964	4. references to historical events or	“This neuron is related to geographical	1009
965	items related to mail delivery	locations, particularly bodies of water,	1010
966		and technical or specialized terminology,	1011
967	Return ONLY a JSON object with a single	possibly referring to engineers or experts	1012
968	field “summary”.	working on or near rivers, islands, or	1013
		countries, potentially addressing techno-	1014
		logical problems or bugs.”	1015
969	<b>LLM Output</b> Here is the output for refined inter-		
970	pretation for this neuron.		
971	“summary”: “This neuron is associated		
972	with geographical locations,	This aggregated interpretation is more specific	1016
973	particularly bodies of water and land-	than either the SAE-based or QA-based inter-	1017
974	related terms, often appearing in	pretation alone. It retains activation-grounded se-	1018
975	technical or specialized contexts such	semantic elements (e.g., geographical references	1019
976	as engineering or problem-solving	such as <i>river</i> and <i>country</i> ) while prioritizing func-	1020
977	discussions.”	tional cues supported by QA probes (e.g., technical	1021
		problem-solving). At the same time, the aggrega-	1022
978	<b>E Example of Aggregated Interpretation</b>	tion discards misleading or weakly supported	1023
979	<b>Refinement</b>	abstractions that arise when SAE features or QA	1024
980	To illustrate how LLM-based aggregation improves	probes are considered in isolation. This example	1025
981	specificity and faithfulness of neuron interpreta-	illustrates how LLM-based aggregation leverages	1026
982	tions, we present a representative example from our	cross-validated evidence to produce neuron inter-	1027
983	analysis. We consider Neuron 3518(128), which	pretations that are both semantically coherent and	1028
984	exhibits partial semantic alignment between SAE-	more faithful to observed activation behavior.	1029
985	based and QA-based interpretations. After cross-		
986	filtering, the primary QA evidence associated with	<b>F Qualitative Interpretation Score (QIS)</b>	1030
987	this neuron includes the following semantic ques-		
988	tions:	While semantic similarity between SAE-based and	1031
989	• an expression of anger	QA-based interpretations provides a coarse mea-	1032
990	• body of water	sure of cross-modal convergence, it does not cap-	1033
991	• technical or specialized terminology	ture the quality of the <i>refined neuron-level inter-</i>	1034
992	• technological problem or bug	<i>pretation</i> produced after cross-filtering and LLM-	1035
993		based semantic aggregation. In particular, simple	1036
994	The corresponding cross-filtered SAE evidence	similarity scores do not distinguish between fo-	1037
995	consists of the following top activating tokens:	ocused, coherent interpretations and vague or frag-	1038
996	• Engineers	mented ones.	1039
997	• River	To address this limitation, we introduce the	1040
998	• country	<b>Qualitative Interpretation Score (QIS)</b> , a com-	1041
999	• haven	posite metric designed to quantify how well a re-	1042
1000		efined neuron summary captures the shared semantic	1043
1001	When considered independently, neither inter-	structure supported by both SAE and QA evidence.	1044
1002	pretation provides a fully satisfactory explanation	Importantly, QIS does not assume that SAE and	1045
1003	of the neuron’s functional role. The SAE tokens	QA interpretations must match exactly. Instead, it	1046
1004	mix geographical and occupational cues without in-	evaluates whether the refined summary accurately	1047
1005	dicating a clear semantic focus, while the QA ques-	represents the <i>shared semantic subspace</i> between	1048
1006	tions span both topical (e.g., bodies of water) and	the two methods while maintaining internal consis-	1049
1007	functional (e.g., technological problems) dimen-	tency and semantic focus.	1050
1008	sions. We apply LLM-based aggregation to jointly	<b>Overview</b> Given a neuron with: (i) a set of SAE-	1051
	summarize the shared and activation-grounded evi-	derived semantic descriptions, (ii) a set of QA-	1052
	dence. The resulting refined interpretation is:	derived semantic questions, and (iii) a refined neu-	1053
		ron summary produced via constrained LLM-based	1054
		aggregation, QIS combines four complementary	1055
		components:	1056

- 1057 • **Intra-set coherence:** whether semantic ev- 1101  
1058 idence within each method (SAE or QA) is 1102  
1059 internally consistent. 1103
- 1060 • **Cross-set overlap:** the degree of semantic 1104  
1061 overlap between SAE and QA evidence. 1105
- 1062 • **Summary grounding:** how well the refined 1106  
1063 summary aligns with evidence from both SAE 1107  
1064 and QA. 1108
- 1065 • **Semantic focus:** a dispersion-based penalty 1109  
1066 that captures whether the supporting evidence 1110  
1067 forms a tight semantic cluster around the sum- 1111  
1068 mary.

1069 Together, these components reward neuron inter-  
1070 pretations that are coherent, cross-validated, and  
1071 semantically focused, while penalizing interpreta-  
1072 tions that mix unrelated concepts or rely on overly  
1073 generic abstractions.

1074 **Formal Definition** Let  $S$  denote the set of SAE  
1075 semantic descriptions for a neuron,  $Q$  denote the set  
1076 of QA semantic questions, and  $r$  denote the refined  
1077 neuron summary. All elements are embedded into  
1078 a shared semantic space.

1079 We define the four components as follows:

1080 (a) **Intra-set coherence.** For each set  $X \in$   
1081  $\{S, Q\}$ , intra-set coherence is computed as the av-  
1082 erage pairwise cosine similarity within  $X$ .

1083 (b) **Cross-set overlap.** Cross-set overlap is de-  
1084 fined as the average cosine similarity between ele-  
1085 ments of  $S$  and  $Q$ .

1086 (c) **Summary grounding.** Summary grounding  
1087 is defined as the average cosine similarity between  
1088  $r$  and elements in  $S \cup Q$ , measuring whether the  
1089 summary lies near the semantic center of the sup-  
1090 porting evidence.

1091 (d) **Semantic focus.** Semantic focus is defined  
1092 as a dispersion-based penalty measuring the aver-  
1093 age distance between evidence elements and the  
1094 summary  $r$ . Higher dispersion corresponds to  
1095 lower focus.

1096 We define the final quality score as a weighted  
1097 sum of grounding, overlap, intra-set consistency  
1098 (SAE and QA), and a dispersion penalty:

$$1099 \text{QIS} = w_g G + w_o O + w_s I_S + w_q I_Q - w_d D,$$

1100 with all weights fixed across experiments.

**Weighting Rationale** The weights are assigned  
as: Grounding 0.4, Overlap 0.3, SAE coherence  
0.15, QA coherence 0.15, Dispersion 0.1. We  
assign higher weight to summary grounding and  
cross-set overlap, as these directly reflect whether  
the refined interpretation captures the shared se-  
mantic structure supported by both methods. Intra-  
set coherence terms capture internal consistency  
within each method and are assigned lower weight.  
The dispersion term serves as a mild penalty to dis-  
courage overly diffuse or generic interpretations.

## 1112 G Dataset Details

1113 Table 2 summarizes all datasets used in our experi-  
1114 ments, including their domains, languages, and the  
1115 text fields from which token snippets are extracted.

1116 For each dataset, we use the primary  
1117 free-text field as the source for snippet ex-  
1118 traction (e.g., article, sentence, content,  
1119 func\_code\_string, or text), depending on  
1120 the dataset schema. All datasets are processed  
1121 using the same tokenization and segmentation  
1122 procedure: fixed-length snippets of 64 tokens are  
1123 extracted and further segmented into overlapping  
1124 contiguous spans of 16 tokens with a sliding  
1125 window, resulting in 49 spans per snippet. No  
1126 dataset-specific preprocessing or filtering is applied  
1127 unless otherwise stated.

## 1128 H Interactive UMAP Visualizations

1129 We include interactive cross-filtered UMAP visu-  
1130 alizations as an offline HTML artifact. To view,  
1131 open `interactive/umap_viewer.html` in a mod-  
1132 ern browser (Chrome/Firefox). Each point is click-  
1133 able and reveals the corresponding neuron, along  
1134 with the top-matching SAE tokens/feature descrip-  
1135 tions and QA questions. No external network ac-  
1136 cess is required.

## 1137 I Details on SAE–QA Disagreement Cases

1138 We present additional neuron-level disagreement  
1139 cases to illustrate recurring failure modes of SAE-  
1140 based interpretations. In these cases, SAE features  
1141 yield semantically rich but internally diffuse de-  
1142 scriptions, while QA-based probes reveal a more  
1143 constrained functional role that is not captured by  
1144 the SAE interpretation. Such cases highlight the  
1145 limits of reconstruction-driven feature discovery  
1146 when neurons respond to heterogeneous surface  
1147 patterns rather than a coherent semantic concept.  
1148 For example, neuron 70 exhibits a representative

Dataset	Domain	Lang.	Text Field Used	#Segments	Usage
CC-News (Mackenzie et al., 2020)	News	EN	article	9,800	Main
takala/financial_phrasebank (Malo et al., 2014)	Finance	EN	sentence	4,312	Robustness
MedRAG/pubmed (Xiong et al., 2024)	Scientific	EN	content	4,312	Robustness
code-search-net/code_search_net (Husain et al., 2019)	Programming	EN	func_code_string	4,312	Robustness
hugcyp/LCSTS (Hu et al., 2015)	Non-English Summaries	ZH	text	4,312	Robustness

Table 2: Dataset details and text fields used for snippet extraction.

disagreement pattern in which SAE-derived features appear to capture diverse high-level semantic domains, while QA-based evidence indicates a narrower functional sensitivity.

**SAE-Based Interpretation.** The top SAE features associated with this neuron are characterized by highly heterogeneous descriptions, including references to political and military events, sports statistics, mathematical expressions, programming syntax, browser functionality, historical events, and customer service language. Correspondingly, the top-activating tokens span a wide range of surface forms, such as punctuation marks, numerical sequences, dates, programming-related symbols (e.g., parentheses, operators, code fragments), and common function words.

While individual SAE feature descriptions are semantically plausible in isolation, their aggregate lacks a coherent unifying concept. Instead, the features collectively emphasize structural properties of text, including heavy punctuation, numerical formatting, code-like syntax, and structured or templated expressions. This results in an interpretation that is semantically rich but diffuse, making it difficult to ascribe a clear functional role to the neuron based on SAE evidence alone.

**QA-Based Interpretation.** In contrast, the QA-based probes for Neuron 70 show strong positive associations with questions related to technological concepts, technological problems or bugs, and mathematical reasoning. These questions consistently point to technical or formal contexts in which structured symbols, numerical values, and code-like patterns are prevalent.

**Disagreement Analysis.** The disagreement arises because SAE-based interpretations surface a broad set of domain-level labels (e.g., politics, sports, military history) that are not supported by the QA-based evidence. The QA probes suggest that the neuron’s activation is driven less by semantic content in these domains and more by the presence of structurally dense, technical, or formally expressed text. In this case, SAE features

appear to conflate multiple surface-correlated patterns into semantically labeled features, obscuring the neuron’s underlying functional sensitivity.

**Failure Mode Characterization.** This example illustrates a recurring failure mode in which SAE features assign semantic labels to neurons that are primarily sensitive to formatting, symbolic structure, or syntactic regularities. Although such features reconstruct activations effectively, they do not faithfully capture the neuron’s semantic role. QA-based probes provide an external constraint that exposes this mismatch by highlighting the neuron’s consistent association with technical and formal contexts rather than the diverse semantic categories suggested by SAE features.

## J Representative Neuron Case Studies

In this appendix, we present representative neuron-level case studies to illustrate the qualitative behavior underlying the layer-wise faithfulness trends reported in the main text. Rather than relying solely on semantic plausibility, each case study evaluates faithfulness by testing whether a proposed interpretation can *predict neuron activation under controlled semantic interventions*. Neurons are selected to exemplify QA-favored, SAE-favored, and ambiguous diagnostic outcomes.

### J.1 QA-Favored Middle-Layer Neuron

We first examine a representative neuron from a middle layer (layer 16) for which QA-based hypotheses provide a more faithful explanation than SAE-based interpretations. This neuron is associated with QA concepts related to technical or bug-oriented content. Controlled contrasts demonstrate that QA-derived semantic conditions reliably predict activation behavior, while alternative explanations fail to generalize.

=== Neuron 70 @ Layer 16 ===

TECH vs CONTROL:

AUC : 0.8007  
Cliff delta : 0.6014  
t-test pval : 4.02e-22  
means : 0.5235 vs 0.4020

1235  
1236 MATH vs CONTROL:  
1237 AUC : 0.6661  
1238 Cliff delta : 0.3322  
1239 t-test pval : 3.05e-14  
1240 means : 0.4680 vs 0.4020  
1241  
1242 TECH vs MATH:  
1243 AUC : 0.7204  
1244 Cliff delta : 0.4408  
1245 t-test pval : 1.30e-08  
1246  
1247 Conclusion:  
1248 Neuron activation is more controllable  
1249 by TECH/BUG-style text,  
1250 supporting the QA-based hypothesis.

1251 The consistently high AUC values and large ef-  
1252 fect sizes indicate that the QA-based interpretation  
1253 captures a functionally relevant trigger for this neu-  
1254 ron. In contrast, SAE-based interpretations for this  
1255 neuron do not yield comparable predictive power  
1256 under matched controls, suggesting lower faithful-  
1257 ness despite their semantic richness.

## 1258 J.2 SAE-Favored Late-Layer Neuron

1259 We next analyze a representative late-layer neuron  
1260 (layer 28) for which SAE-based interpretations pro-  
1261 vide a more faithful explanation. Initial QA-based  
1262 probes suggest sensitivity to religious content; how-  
1263 ever, controlled conditional tests reveal that this  
1264 apparent effect does not generalize once language  
1265 form is accounted for. In contrast, SAE-derived  
1266 hypotheses related to language form remain predic-  
1267 tive across conditions.

1268 === Neuron 1216 @ Layer 28 ===  
1269  
1270 Language main effect (Non-religious):  
1271 A1B0 vs A0B0 AUC : 0.848  
1272  
1273 Religion main effect (Non-foreign):  
1274 A0B1 vs A0B0 AUC : 0.889  
1275  
1276 Religion within FOREIGN language:  
1277 A1B1 vs A1B0 AUC : 0.472  
1278  
1279 Language within RELIGION:  
1280 A1B1 vs A0B1 AUC : 0.630  
1281  
1282 Conclusion:  
1283 Apparent religion effects disappear  
1284 under language control,  
1285 while language-form effects persist.

1286 These results falsify the QA hypothesis that reli-  
1287 gion is the primary semantic driver for this neuron.  
1288 Although religion-related content initially appears  
1289 predictive, its effect does not generalize across lan-  
1290 guage conditions. In contrast, the SAE-based inter-  
1291 pretation—sensitivity to language form or encod-  
1292 ing—remains stable under targeted interventions,

indicating higher faithfulness. 1293

## J.3 Ambiguous Early-Layer Neuron 1294

1295 Finally, we consider an early-layer neuron (layer 4)  
1296 that exhibits weak and inconsistent predictive  
1297 power for both SAE- and QA-based interpretations. 1298  
1299 For this neuron, neither hypothesis achieves robust  
1300 discrimination, and observed effects are sensitive  
1301 to experimental conditioning. 1302

1303 === Neuron 162 @ Layer 4 === 1304  
1305  
1306 Encoding effect (CY vs EN, Non-financial  
1307 ): 1308  
1309 AUC : 0.540 1310  
1311  
1312 Finance effect (FIN vs Non-FIN, English)  
1313 : 1314  
1315 AUC : 0.465 1316  
1317  
1318 Conclusion: 1319  
1320 Mixed or weak effects; neither  
1321 hypothesis generalizes reliably. 1322  
1323  
1324 Both QA-based semantic probes and SAE-based  
1325 feature descriptions yield only marginal improve-  
1326 ments over chance and fail to generalize across  
1327 subconditions. This suggests that the neuron may  
1328 encode low-level, entangled, or highly context-  
1329 dependent features that are not well captured by  
1330 either semantic abstraction. 1331  
1332 Together, these case studies demonstrate that  
1333 faithfulness depends on whether an interpretation  
1334 can predict neuron behavior under controlled in-  
1335 terventions, not merely on semantic plausibility.  
1336 QA-based and SAE-based methods each succeed  
1337 in different regimes, while some neurons admit no  
1338 clean semantic explanation. These findings moti-  
1339 vate a comparative, falsification-oriented approach  
1340 to neuron interpretability. 1341