SALVE: SPARSE AUTOENCODER-LATENT VECTOR EDITING FOR MECHANISTIC CONTROL OF NEURAL NETWORKS

Anonymous authors
Paper under double-blind review

000

001

002

004 005 006

007

008 009 010

011 012 013

014

015

016

017

018

019

021

022

025

026

027

028 029 030

031 032 033

034

036

037

038

041

042

043

044

046

047

048 049

051

052

Abstract

Deep neural networks achieve impressive performance but remain difficult to interpret and control. We present **SALVE** (Sparse Autoencoder-Latent Vector Editing), a unified "discover, validate, and control" framework that bridges mechanistic interpretability and model editing. Using an ℓ_1 regularized autoencoder, we learn a sparse, model-native feature basis without supervision. We validate these features with Grad-FAM, a feature-level saliency mapping method that visually grounds latent features in input data. Leveraging the autoencoder's structure, we perform precise and permanent weight-space interventions, enabling continuous modulation of both classdefining and cross-class features. We further derive a critical suppression threshold, $\alpha_{\rm crit}$, quantifying each class's reliance on its dominant feature, supporting fine-grained robustness diagnostics. Our approach is validated on both convolutional (ResNet-18) and transformer-based (ViT-B/16) models, demonstrating consistent, interpretable control over their behavior. This work contributes a principled methodology for turning feature discovery into actionable model edits, advancing the development of transparent and controllable AI systems.

1 Introduction

Understanding the internal mechanisms of deep neural networks remains a central challenge in machine learning. While these models achieve remarkable performance, their opacity hinders our ability to trust, debug, and control their decision-making processes, especially in high-stakes applications where reliability is non-negotiable. The field of Mechanistic interpretability aims to resolve these issues by reverse-engineering how networks compute, identifying internal structures that correspond to meaningful concepts and establishing their influence on outputs (4; 17; 32; 1; 26) However, the bridge between interpretation and intervention remains a critical frontier. While recent advances in model steering successfully use discovered features to guide temporary, inference-time adjustments, a path toward using these insights to perform durable, permanent edits to a model's weights is less established. This paper closes that gap by introducing SALVE (Sparse Autoencoder-Latent Vector Editing), a unified framework that transforms interpretability insights into direct, permanent model control. We build a bridge from unsupervised feature discovery to fine-grained, posthoc weight-space editing. Our core contribution is a "discover, validate, and control" pipeline that uses a sparse autoencoder (SAE) to first learn a model's native feature representations and then leverages that same structure to perform precise, continuous interventions on the model's weights.

Our framework achieves this by first training an ℓ_1 -regularized autoencoder on a model's internal activations to discover a sparse, interpretable feature basis native to the model. We validate these features using visualization techniques, including activation maximization and our proposed Grad-FAM (Gradient-weighted Feature Activation Mapping), to confirm their semantic meaning. This interpretable basis is then used to guide a permanent, continuous weight-space intervention that can suppress or enhance a feature's influence. Finally, to

make this control quantifiable, we derive and validate a critical suppression threshold, $\alpha_{\rm crit}$, providing a measure of a class's reliance on its dominant feature.

2 Related Work

054

055

057

060

061

062

063 064

065

066

067

068

069

071

072

073

074

075

076

077 078

079

081

083

085

086

087

880

089

090

091

092

094

096

098

100

101

102

103 104

105 106

107

Our work is positioned at the intersection of two primary domains: the discovery of interpretable features and the direct editing and control of model behavior. We situate our contributions at the intersection of these fields, focusing on creating a direct pathway from understanding to control.

2.1 Discovering Interpretable Features

A key goal of post-hoc interpretability is to make a trained model's decisions intelligible. Attribution methods, such as Grad-CAM (31), generate saliency maps that highlight influential input regions. While useful for visualization, these methods are correlational and do not expose the internal concepts the model has learned. Other approaches aim to link a model's internal components with human-understandable concepts. For example, Network Dissection (3) quantifies the semantics of individual filters by testing their alignment with a broad set of visual concepts. Similarly, TCAV (14) uses directional derivatives to measure a model's sensitivity to user-defined concepts. These methods are powerful but often rely on a pre-defined library of concepts. Recent work in mechanistic interpretability has used SAEs to discover features in an unsupervised manner, primarily within Transformer-based models. Our work builds on this SAE-based approach to feature discovery by using a SAE to automatically discover semantically meaningful features directly from the model's activations.

2.2 Editing and Controlling Model Behavior

Model editing techniques aim to modify a model's behavior without retraining. Simple interventions like ablation typically zero out neurons, or entire filters, to observe their effect on the output (18; 19; 25). However, while this provides evidence of their importance, these approaches are limited in their ability to perform fine-grained, continuous interventions, and often lack a structured representation in which such interventions can be reasoned about and controlled directly. Most recent work using SAEs has focused on inference-time steering, where interventions involve adding a "steering vector" to a model's activations during a forward pass to influence the output (33; 36; 5; 12). Our approach is fundamentally different in its mechanism: instead of a temporary inference time adjustment to activations, we perform a permanent edit directly on the model's weights, allowing both suppression and enhancement of specific features. This is advantageous for scenarios requiring fixed, verifiable model states, eliminating the computational overhead of steering vectors at inference and ensuring consistent behavior across all uses of the edited model. Our work is also related to specialized, training-free model editing techniques like ROME (22) and MEMIT (23). These methods perform corrective, example-driven edits by calculating a surgical weight update to alter a model's output for a user-provided input. In contrast, our method is feature-driven and diagnostic. Interventions are guided by general, model-native concepts discovered by the SAE. This approach enables the continuous modulation and quantitative analysis of concepts, providing a more transferable mechanism for influencing a model's overall behavior than single-instance correction. Finally, our method differs from other concept-editing paradigms that are more invasive. For instance, Iterative Nullspace Projection (INLP) (30) removes information by projecting representations, but often requires fine-tuning to preserve model performance. Similarly, Concept Bottleneck Models (CBMs) (16) involve substantial architectural changes to incorporate a labeled concept layer. In contrast, our approach performs a localized weight edit guided by discovered features, remaining fully post-hoc and avoiding both retraining and architectural modification

3 Methods

Our framework follows a three-stage "discover, validate, and control" pipeline designed to identify, understand, and intervene on a model's internal features. The stages are: (1)

Table 1: Comparison of related interpretability and model editing methods.

Method Category	Examples	Primary Objective	Key Trait or Limitation
Attribution	Grad-CAM	Interpret	Correlational, not causal
Concept-based	Dissection, TCAV	Interpret	Requires predefined, labeled concepts
Ablation	Filter Pruning	Control	Coarse, not fine-grained
Activation Steering	SAE Steering	Interpret + Control	Temporary (inference-time only)
Factual Editing	ROME, MEMIT	Control	Corrective, not diagnostic; example-dependent
Projection-based	INLP	Control	Often requires fine-tuning
Architectural	CBMs	Control	Invasive; requires retraining
Our Approach	-	${\bf Interpret+Control}$	Permanent edit of model-native concepts

Discover, where we learn a sparse latent representation of a model's activations using an ℓ_1 -regularized autoencoder; (2) **Validate**, where we use visualization techniques to confirm the semantic meaning of the discovered features; and (3) **Control**, where we perform targeted, continuous interventions on the model's weights guided by the autoencoder's structure. We demonstrate this framework on two distinct model architectures. Our primary analysis is conducted on a ResNet-18 model fine-tuned on the Imagenette dataset (11). To demonstrate the generality and robustness across other architectures, the core stages of our methodology are successfully validated on a Vision Transformer (ViT-B/16 (7)). Further details on both models are available in Appendix A.3. The remainder of this section details the procedures for each stage.

3.1 Discovering Interpretable Features

To obtain a sparse and interpretable representation of the model's internal activations, we train a linear autoencoder on the outputs of a semantically rich layer. For ResNet-18, we use the final average pooling layer, while for the Vision Transformer, we extract the representation corresponding to the [CLS] token after the final transformer encoder block. We use a standard reconstruction loss combined with an ℓ_1 penalty on the latent activations to encourage sparsity. The full loss function and further architectural details are provided in Appendix A.4. To identify dominant features associated with specific output classes from this representation, we compute the class-conditional mean of the latent activations, $\mu_k = \frac{1}{|C_k|} \sum_{i \in C_k} Z_i$, where C_k is the set of samples for class k and Z_i is the latent activation vector for a single sample i. Analyzing μ_k reveals which features are strongly associated with a particular class, providing a basis for targeted interventions.

3.2 Validating and Controlling Features

To validate the semantic content of the discovered features, we use two complementary visualization techniques. We employ activation maximization (27), extending it from traditional applications on individual neurons or filters to our discovered latent features to synthesize the abstract concept a feature represents. To ground feature activations in specific input images, we introduce Grad-FAM (Gradient-weighted Feature Activation Mapping), an adaptation of Grad-CAM (Gradient-weighted Class Activation Mapping). While standard Grad-CAM generates saliency maps for a class logit, our method repurposes this logic for a specific latent feature, providing a direct visual link between the feature and the input regions that activate it. See Appendix A.5 for the full derivation and implementation details.

Having validated that the features represent meaningful semantic concepts, we investigate their causal role by using the autoencoder's decoder matrix, $D \in \mathbb{R}^{M \times d}$, to guide a permanent, continuous edit to the weights of the model. Each column of this matrix, D[:, l], represents the "direction" of a latent feature l in the original activation space. We edit the final-layer weights w_{ij} as follows:

$$w'_{ij} = w_{ij} \cdot \max(0, 1 \pm \alpha \cdot |c_j|), \tag{1}$$

where $c_j = D[j, l]$ is the contribution of the latent feature to the original feature j, and α controls the intervention strength. The \pm symbol indicates that we can either enhance (+) or suppress (-) the feature's influence, depending on the desired intervention. This feature-guided modulation differs from traditional ablation as it allows for fine-grained control over

the model's learned concepts while preserving the network's structure. To quantify a class's reliance on a feature, we first define the feature-perturbed logit, $z'_i(\alpha)$, as a function of the intervention strength α :

$$z_i'(\alpha) = \sum_j w_{ij} \cdot \max(0, 1 - \alpha |c_j|) \cdot x_j, \tag{2}$$

where x_j is the j-th component of the activation vector from the penultimate layer. This gives the class logit under feature perturbation, excluding the global class bias term b_i so that we measure the feature's influence independently of baseline class predisposition. We define the critical suppression threshold, $\alpha_{\rm crit}$, as the value of α for which $z_i'(\alpha) = 0$. Intuitively, $\alpha_{\rm crit}$ is the intervention strength required to completely suppress the logit contribution attributable to this feature direction. In the regime of weak perturbations ($\alpha < 1/|c_j|$), a linear approximation yields the analytical estimate:

$$\alpha_{\text{crit,i}}^{(n)} \approx \frac{z_i^{(n)}}{R_i(\mathbf{x}^{(n)})},$$
(3)

where $z_i^{(n)} = \sum_j w_{ij} x_j^{(n)}$ is the original (bias-free) class logit for sample n, and $R_i(\mathbf{x}^{(n)}) = \sum_j |c_j| w_{ij} x_j^{(n)}$ quantifies the logit's sensitivity to suppression along the latent feature direction. While the analytical estimate provides a lower bound, we also compute the exact threshold numerically. The full derivation and further details are provided in Appendix A.6.

4 Results

We validate our framework through two main analyses. First, we identify and visualize the discovered latent features to confirm they represent meaningful semantic concepts. Second, we perform targeted weight-space interventions to demonstrate precise control over the model's output. Together, these analyses demonstrate our framework's ability to link interpretation directly to intervention.

4.1 Latent Features Encode Semantic Concepts

Our first objective is to validate that the features discovered by the SAE correspond to meaningful visual concepts. We find that the SAE successfully learns a class-specific feature basis. As shown in the average feature activations in Figure 1a, the ResNet-18 representations have a sparse, dominant structure where a single feature is strongly associated with a single class. While these dominant features define a class, the less active features often represent finer-grained concepts shared across classes, which we explore further in Section 4.3.

To understand what these features represent abstractly, we first use activation maximization (see Appendix A.5.1 for further details). Figure 1b shows this for the "golf ball" feature, where the optimization reveals objects with the characteristic texture of a golf ball. To further visualize how these features are grounded in specific inputs, we use our proposed Grad-FAM method (see Appendix A.5 for the full derivation and implementation details). Figure 1c shows an example for a golf ball image. The dominant "Feature 1" provides a high-level concept map for the "golf ball" class, whereas less dominant features correspond to granular sub-concepts. For example, "Feature 2" activates on the ball's surface texture, while "Features 3 and 4" highlight different parts of the golf club.

To confirm these findings were not an artifact of the convolutional architecture, we replicated our core validation analyses on the Vision Transformer (ViT). We found that the ViT also learned a sparse, class-specific feature basis and that Grad-FAM successfully localized these features to semantically relevant image regions (see Appendix A.8.1). Together, these results confirm that the discovered features are semantically aligned, providing a valid basis for targeted interventions.

4.2 Controlling Class Predictions via Feature Editing

Having validated the semantic meaning of the features, we now evaluate their causal role by performing permanent model edits. We first focus on the dominant, class-specific features. A

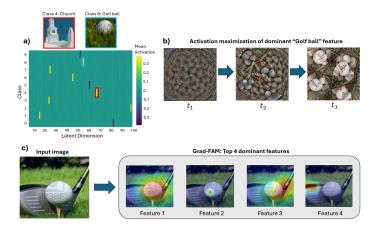


Figure 1: Validation of discovered features for the ResNet-18 model. (a) Average latent feature activations across classes, showing a sparse, class-specific basis. (b) Activation maximization of the "golf ball" feature. The image evolves from random noise (t_1) , through emergence of circular shapes (t_2) , to the final image clearly exhibiting golf ball characteristics (t_3) . (c) Grad-FAM visualizations grounding the top-4 dominant features for a sample "golf ball" image.

qualitative case study on an ambiguous, out-of-distribution image (not part of Imagenette) containing both a "golf ball" and a "church" demonstrates the precision of our method. This experiment is illustrated under "Single prediction" in Figure. 2. The original model predicts "Church", and as a comparative illustration we calculate both the standard Grad-CAM as well as Grad-FAM saliency maps. The Grad-CAM indicates that the model primarily focuses on the tower of the church to make its prediction. Using Grad-FAM, we confirm that the dominantly activated feature corresponds to the church structure (similar to the Grad-CAM), but it also reveals other features activated by the golf ball (Feature 2) and by the church tower (Feature 3).

We then perform two interventions. First, suppressing the dominant "Church" feature predictably flips the classification to "Golf ball". Conversely, enhancing the "Golf ball" feature achieves the same outcome. Post-edit Grad-CAMs confirm the model's attention shifts accordingly in both cases, from focusing primarily on the church tower to the golf ball instead. To validate this effect quantitatively, we evaluate the model on the entire Imagenette test set. The confusion matrices in Figure. 2 for the edited model show two key results. First, suppressing the "Church" feature disables the model's ability to recognize that class, reducing its accuracy to near zero. Second, enhancing the "Golf ball" feature maintains its near-perfect accuracy, indicating the edit is stable and does not degrade performance on well-learned and robust class representations. Critically, for both interventions, performance on the other classes remains relatively unaffected. This confirms that the learned class-dominant features enable precise, modular control, allowing for targeted interventions with minimal off-target effects. While this example targeted a specific class, the same methodology was successfully applied to the other classes in the dataset. To validate the architectural robustness of this control method, we replicate these class-suppression experiments on the Vision Transformer and achieve a similar degree of precise, modular control over its predictions (see Appendix A.8.2).

To contextualize our results, we include a baseline inspired by ROME (22), a well-known method for example-driven factual editing in large language models. We adapt ROME's rank-one update idea to the final classification layer of our model for a comparable "class suppression" task, achieving a similar outcome of reducing the target accuracy to near zero with minimal off-target effects (see Appendix A.7.2). However, this similarity in outcome masks a fundamental difference in approach. ROME is a corrective tool that performs a surgical update to fix a model's output for a specific input. In contrast, our method is a diagnostic tool that intervenes on general, model-native concepts. This feature-driven

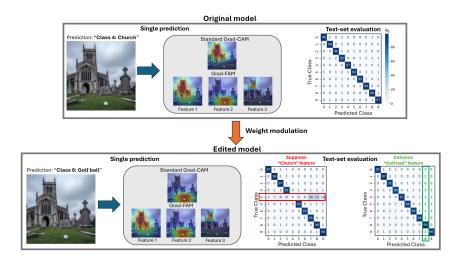


Figure 2: (Left) A qualitative case study where suppressing the "Church" feature or enhancing the "Golf ball" feature successfully flips the model's prediction for an ambiguous image. (Right) Quantitative validation on the test set, showing minimal off-target effects

approach is what enables the capabilities unique to our framework: continuous modulation, quantitative analysis (α_{crit}), and the nuanced, cross-class edits explored in the next sections.

4.3 Intervening on Cross-Class Features

To demonstrate our framework's ability to edit fine grained concepts shared across classes, we identify a "Tower Feature", which is frequently activated by images of both churches (Class 4) as well as petrol pumps (Class 7) containing tower-like structures. A selection of the top-activating images from the test set, shown in Figure 3a, confirm this cross-class association. We then perform symmetrical interventions on this feature (Figure 3b). Suppressing the feature reduces accuracy for "Petrol Pump" while leaving the "Church" class largely unaffected. Conversely, enhancing the feature increases the accuracy for "Petrol Pump". This differential impact suggests that the model's classification of certain older, column-shaped petrol pumps is highly reliant on the "Tower Feature". In contrast, the "Church" class appears more robust due to a richer set of redundant features (e.g., stained glass, steeples, etc.). Notably, these interventions also revealed a subtle feature entanglement.

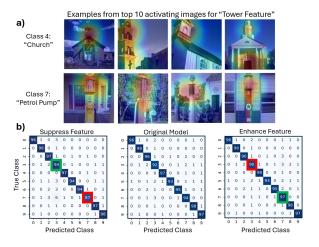


Figure 3: a) Validation of the "Tower Feature", showing example top-activating images from the test set. b) Effect of suppressing and enhancing this feature on model predictions, where "red" and "green" corresponds to decreasing or increasing class accuracy, respectively

Suppressing the "Tower Feature" slightly increased accuracy for "Chain Saw" (Class 3), while enhancing it caused a corresponding decrease. This consistent, inverse effect, suggests that the model has learned a spurious negative correlation, causing the "Tower Feature" to act as an inhibitor for the "Chain Saw" classification. This illustrates a nuanced relationship that would be difficult to uncover without such targeted interventions.

4.4 Quantifying Intervention Sensitivity

Having established that interventions are effective, we now quantify their sensitivity to the scaling factor α . Systematically varying α allows us to probe the causal importance of each feature and measure a class's reliance on it. Figure 4a shows a characteristic suppression curve for the "Church" class. As α increases, its accuracy remains stable before dropping sharply past a critical threshold, while performance on other classes remains high, confirming the targeted nature of the edit. This drop in accuracy corresponds to a reallocation of the model's predictions. As the dominant feature is suppressed, the model's confidence is redistributed among the other classes (Figure 4b)

To evaluate the sensitivity of the learned latent basis on the class intervention, we performed experiments across multiple initializations and training runs (see Appendix A.9.1). Crucially, the results are consistent across multiple SAE training runs, as indicated by the narrow shaded regions in the curves of Figure 4, representing the standard deviation across n=10 different realizations. This demonstrates that while the specific basis vectors learned by the SAE may vary, the effect of our intervention on the model's underlying concepts is stable.

We then extend this analysis by computing suppression curves for the dominant feature of all classes. As these results proved robust to different autoencoder initializations, we simplify the subsequent analysis by presenting results from a single realization. We observe a similar suppression behavior across all classes for the ResNet-18 model (Fig. 5a), with variations in the threshold suggesting that each class relies on its dominant feature to a different degree. To confirm that this behavior was not an artifact of the convolutional architecture, we replicate the same analysis on the Vision Transformer. The ViT exhibited the same characteristic suppression curves, indicating that this is a general property of the intervention method (see Appendix A.8.3). As the results are robust to the stochasticity of the SAE training process, this indicates that the observed differences in class sensitivity stem mainly from the base model's intrinsic feature representations. In addition, the properties of the discovered feature basis are also dependent on the SAE's training objective (e.g., the sparsity coefficient λ_1). Therefore, the observed heterogeneity between classes is a product of both the backbone's intrinsic structure and the specific feature decomposition learned by the SAE. Further details assessing the robustness of our experiments are provided in Appendix A.9.

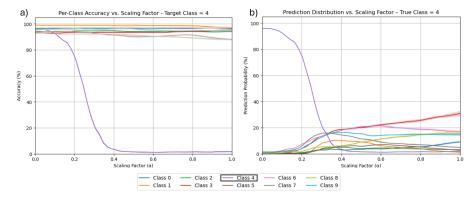


Figure 4: Suppression sensitivity for "Church" (Class 4). (a) Per-class accuracy vs. intervention strength α . (b) Distribution of predictions for images of the target class, showing how confidence is reallocated as the feature is suppressed. Shaded regions indicate the standard deviation across 10 SAE initializations.

4.4.1 Critical Suppression Strength: $\alpha_{\rm crit}$

To quantify this class-level feature dependency, we calculate the critical suppression threshold, $\alpha_{\rm crit}$ (derivation and implementation details are found in Appendix A.6). This metric estimates, on a per-sample basis, the precise intervention strength required to reduce the feature contribution to the logit for the correct class to zero, providing a granular measure of how strongly the model relies on a specific feature for a given input. We compute $\alpha_{\rm crit}$ in two ways: an analytical estimate based on the linear approximation (given by Eq.3), and a numerical calculation (given by Eq.2). To validate both, we compare their distributions to an empirical threshold, $\alpha_{50\%}$. It is important to note the conceptual distinction between these metrics. The analytical and numerical $\alpha_{\rm crit}$ are per-sample measures that identify the intervention strength required to drive the correct class's feature contribution to the logit to zero, representing a point where the feature's evidence for that class is entirely suppressed. In contrast, the empirical $\alpha_{50\%}$ is a population-level metric that marks the point where accuracy on the class falls to 50%. This happens when, for a typical sample, the logit for the correct class has dropped just enough to be surpassed by a competing logit. Figure 5b presents this comparison for the ResNet-18 model. As expected, the analytical estimate provides a lower bound on the critical intervention strength. The numerically calculated $\alpha_{\rm crit}$, in turn, aligns well with both this lower bound and the empirical estimate, and its distribution also captures a wider range of values for samples where the linear approximation is less valid.

To confirm the architectural generality of the $\alpha_{\rm crit}$ metric, we replicate this analysis also on the Vision Transformer. Here, we find two key differences compared to the ResNet-18 model. First, the analytical estimate of $\alpha_{\rm crit}$ represents a much more conservative lower bound. Second, we observe a larger discrepancy between the numerically calculated $\alpha_{\rm crit}$ (representing total evidence loss) and the empirical $\alpha_{50\%}$ (representing a typical flip of predictions). We hypothesize that these observations can be attributed to the ViT's architectural properties. Recent research has shown that ViTs learn a "curved" and nonlinear representation space (15). This effect, combined with the ViT's tendency to produce less confident, more distributed logits (24), means a smaller intervention is required to be surpassed by a competitor. This widens the gap between the point of prediction failure ($\alpha_{50\%}$) and total evidence loss ($\alpha_{\rm crit}$), an effect less pronounced in the more linear representations of the ResNet model. (See Appendix A.8.3 for results and further discussions).

Analyzing the distribution of $\alpha_{\rm crit}$ enables both class-level and sample-level diagnostics, allowing us to pinpoint fragile representations that may be susceptible to adversarial perturbations and to guide targeted strategies for improving robustness.

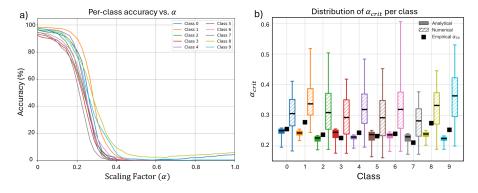


Figure 5: (a) Per-class accuracy vs. the intervention strength α . (b) Comparison of the critical suppression threshold estimates. The distributions of the per-sample analytical (filled box) and numerical (hatched box) $\alpha_{\rm crit}$ are shown as boxplots. The central line indicates the median, the box spans the interquartile range (25th to 75th percentile), and the whiskers extend to the 5th and 95th percentiles. The empirical threshold ($\alpha_{50\%}$) is overlaid as a square marker.

5 Limitations and Future Work

While our framework provides a complete pipeline for discovering, validating, and controlling a model's internal features, this work is presented as a mechanistic case study on a basic dataset and well-understood architectures. This approach establishes the bounds of our current study and motivates the future directions discussed below.

A primary direction for future work is to scale our validation across a broader range of models and datasets. The results of our study, supported by recent findings that sparse autoencoders effectively uncover semantic concepts in Vision Transformers (21), provide a strong foundation for this expansion. Investigating how properties like feature modularity and intervention efficacy vary between different architectures, and extending this comparative analysis to larger models, more complex datasets, and other data modalities like natural language processing, represents crucial next steps.

Furthermore, we found that intervention effectiveness depends on both the backbone model's training dynamics (e.g., batch size) and the parameters used to train the SAE (e.g., ℓ_1 regularization strength) (see Appendix A.9.2 for a detailed discussion). This indicates a crucial link between training procedures and post-hoc controllability. Future work should investigate how to co-design training and intervention methods to produce models that are not only performant but also inherently more editable. Our work also opens several avenues for a deeper feature-level analysis. While we identified concrete instances of feature entanglement (e.g., the "Tower Feature"), future work could incorporate formal disentanglement metrics for a more rigorous evaluation. Additionally, a powerful extension would be to adapt our weight modulation technique to deeper layers of the model. Such a method could enable more fundamental edits to a model's core feature representations, moving from modifying how concepts are combined to changing how they are formed.

Finally, our $\alpha_{\rm crit}$ metric suggests a link between feature reliance and model robustness. Systematically investigating if low- $\alpha_{\rm crit}$ features and samples correlate with known adversarial vulnerabilities represents another promising direction for future work.

6 Conclusions

This work introduced a complete framework for interpreting and controlling neural networks through a unified "discover, validate, and control" pipeline. By training an SAE on a model's internal activations, we extracted a basis of "model native" features that represent not only dominant, class-defining concepts but also fine-grained representations shared across multiple classes. We validated these discovered concepts using activation maximization to visualize their abstract form and our proposed Grad-FAM method to ground their presence in specific input images.

Our primary contribution is demonstrating that this interpretable feature basis can guide precise, permanent, and post-hoc edits to the model's weights. Our experiments showed we can not only switch class-level predictions but also perform nuanced interventions on subtle, cross-class features. To make this control quantifiable, we introduced the critical suppression threshold, $\alpha_{\rm crit}$, a per-sample metric that measures a class's reliance on its dominant feature. This provides a powerful tool for diagnosing and probing robustness of class and feature representations in models. Our methodology is validated on both convolutional (ResNet-18) and transformer-based (ViT-B/16) models, demonstrating consistent, interpretable control over their behavior. By providing a principled method to discover, validate, quantify, and ultimately control a model's internal concepts, our framework contributes to the development of more transparent, robust, and reliable AI systems.

7 Reproducibility statement

The methodology, derivations, parameters, and implementation details are described in the Appendix A). A comprehensive and fully documented GitHub repository will be made publicly available upon publication.

REFERENCES

- [1] Anthropic. Mapping the mind of a large language model. 2024.
- [2] Yu Bai, Yao Tang, Tianle Zhai, Haoran Wang, Xiyu Zhao, Mingyu Zhou, Boqing Gong, Dahua Lin, and Zhouyuan Lin. Beyond pruning criteria: Fine-tuning is the key to robustness. arXiv preprint arXiv:2410.15176, 2024.
- [3] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6541–6549, 2017.
- [4] L. Bereska and E. Gavves. Mechanistic interpretability for ai safety a review. arXiv, 2024.
- [5] Sviatoslav Chalnev, Matthew Siu, and Arthur Conmy. Improving steering vectors by targeting sparse autoencoder features, 2024. *URL https://arxiv. org/abs/2411.02193*.
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009.
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- [8] Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3):1, 2009.
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [10] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. arXiv preprint arXiv:1705.08741, 2017.
- [11] Jeremy Howard. imagenette.
- [12] Sonia Joseph, Praneet Suresh, Ethan Goldfarb, Lorenz Hufe, Yossi Gandelsman, Robert Graham, Danilo Bzdok, Wojciech Samek, and Blake Aaron Richards. Steering clip's vision transformer with sparse autoencoders. arXiv preprint arXiv:2504.08729, 2025.
- [13] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *Proceedings of the 5th International Conference on Learning Representations (ICLR)*, 2017.
- [14] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, pages 2668–2677. PMLR, 2018.
- [15] Juyeop Kim, Junha Park, Songkuk Kim, and Jong-Seok Lee. Curved representation space of vision transformers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 13142–13150, 2024.
- [16] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *International conference on machine learning*, pages 5338–5348. PMLR, 2020.
- [17] L. Kästner and B. Crook. Explaining ai through mechanistic interpretability. European Journal for Philosophy of Science, 2024.

- [18] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. arXiv preprint arXiv:1608.08710, 2016.
- [19] M. Li and L. Janson. Optimal ablation for interpretability. NeurIPS Spotlight, 2024.

- [20] Sihui Li, Haotian Wang, and Zhouyuan Lin. Characterizing the robustness–accuracy trade-off in fine-tuning. arXiv preprint arXiv:2503.14836, 2025.
- [21] Hyesu Lim, Jinho Choi, Jaegul Choo, and Steffen Schneider. Sparse autoencoders reveal selective remapping of visual concepts during adaptation. arXiv preprint arXiv:2412.05276, 2024.
- [22] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. Advances in neural information processing systems, 35:17359– 17372, 2022.
- [23] Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. arXiv preprint arXiv:2210.07229, 2022.
- [24] Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. Advances in neural information processing systems, 34:15682–15694, 2021.
- [25] Ari S. Morcos, David G.T. Barrett, Neil C. Rabinowitz, and Matthew Botvinick. On the importance of single directions for generalization. In *International Conference on Learning Representations*, 2018.
- [26] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- [27] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. Distill, 2(11):e7, 2017.
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems (NeurIPS), 32, 2019.
- [29] Gonçalo Paulo and Nora Belrose. Sparse autoencoders trained on the same data learn different features. arXiv preprint arXiv:2501.16615, 2025.
- [30] Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. Null it out: Guarding protected attributes by iterative nullspace projection. ACL, 2020.
- [31] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [32] Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lindsey, Jeff Wu, Lucius Bushnaq, Nicholas Goldowsky-Dill, Stefan Heimersheim, Alejandro Ortega, Joseph Bloom, et al. Open problems in mechanistic interpretability. arXiv preprint arXiv:2501.16496, 2025.
- [33] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. Transformer Circuits Thread, 2024.
- [34] Kaidi Xu, Sijia Liu, Gaoyuan Zhang, Mengshu Sun, Pu Zhao, Quanfu Fan, Chuang Gan, and Xue Lin. Interpreting adversarial examples by activation promotion and suppression. arXiv preprint arXiv:1904.02057, 2019.

[35] Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. arXiv preprint arXiv:1506.06579, 2015.

[36] Yu Zhao, Alessio Devoto, Giwon Hong, Xiaotang Du, Aryo Pradipta Gema, Hongru Wang, Xuanli He, Kam-Fai Wong, and Pasquale Minervini. Steering knowledge selection behaviours in llms via sae-based representation engineering. arXiv preprint arXiv:2410.15999, 2024.

A Appendix

A.1 LLM USAGE

LLM-based tools were utilized for enhancing the writing quality and clarity of the paper, and for providing support in writing parts of the analysis code. Specifically, Gemini 2.5 Pro was used for grammar and style text polishing and coding support.

The authors maintain full responsibility for the content, accuracy, and conclusions presented in this paper.

A.2 Compute Resources

All experiments in this paper were performed on a single local workstation, and can be reproduced on a standard modern laptop with GPU acceleration.

A.3 FINE-TUNING OF PRE-TRAINED BACKBONE MODELS

We initialized our feature extraction backbone using standard, pre-trained architectures. We used a ResNet-18 model (9) and a Vision Transformer (ViT-B/16 (7)) with a patch size of 16x16. Both models were pre-trained on the ImageNet dataset (6) and loaded via established PyTorch libraries (28).

As the Imagenette dataset (11) contains only 10 target classes, the final classification layer of each model was replaced to match this number of outputs. For ResNet-18, this involved replacing the fc layer, and for the ViT, the head layer.

Both models were then fine-tuned on Imagenette using the Adam optimizer, categorical cross-entropy loss and a batch size of 16. ViT with a learning rate of 10^{-5} for 5 epochs, and ResNet-18 10^{-4} for 10 epochs. We performed additional sensitivity experiments with different batch sizes, as further described in Section A.9.2

A.4 Autoencoder Definition and Training

To learn a sparse, interpretable representation of the backbone models' features, we trained an autoencoder to reconstruct high-dimensional activation vectors extracted from a specific layer in the backbone. The activations were extracted using forward hooks in PyTorch. The target layer was chosen to be a late-stage, semantically rich layer just before the final classification head.

- For ResNet-18, we extracted the 512-dimensional feature vector from the output of the final average pooling layer (avgpool).
- For the ViT, we extracted the 768-dimensional representation corresponding to the [CLS] token after the final transformer encoder block.

This process yielded a matrix of activations $A \in \mathbb{R}^{N \times M}$, where N is the number of samples and M is the feature dimensionality of the chosen layer. The autoencoder's encoder maps these activations to a latent space $Z \in \mathbb{R}^{N \times d}$, and the decoder reconstructs them as $\hat{A} \in \mathbb{R}^{N \times M}$.

The architecture of the autoencoder is defined as follows:

```
694
       class SparseAutoencoder(nn.Module):
           def __init__(self, input_dim, hidden_dim):
               super(SparseAutoencoder, self).__init__()
696
               self.encoder = nn.Linear(input_dim, hidden_dim)
697
               self.decoder = nn.Linear(hidden_dim, input_dim)
698
               nn.init.xavier_uniform_(self.encoder.weight)
699
               nn.init.xavier_uniform_(self.decoder.weight)
700
701
           def forward(self, x):
               encoded = self.encoder(x)
```

decoded = self.decoder(encoded) return encoded, decoded

The model was trained to minimize the reconstruction loss combined with an ℓ_1 -regularization term to promote sparsity:

$$\mathcal{L} = \|A - \hat{A}\|_F^2 + \lambda_1 \|Z\|_1. \tag{4}$$

This sparsity encourages modularity, such that any given input activates only a few latent dimensions.

The autoencoders for both models were trained for 1000 epochs using an Adam optimizer and a step decay learning rate scheduler (reducing LR by a factor of 0.8 every 200 epochs). The ResNet-18 SAE used a learning rate of 10^{-3} , batch size of 32, and $\lambda_1 = 10^{-3}$, while the ViT-B/16 SAE used a learning rate of 10^{-4} , batch size of 16, and $\lambda_1 = 10^{-2}$.

The training loss curves and examples of the reconstructed activations are illustrated in Figure 6.

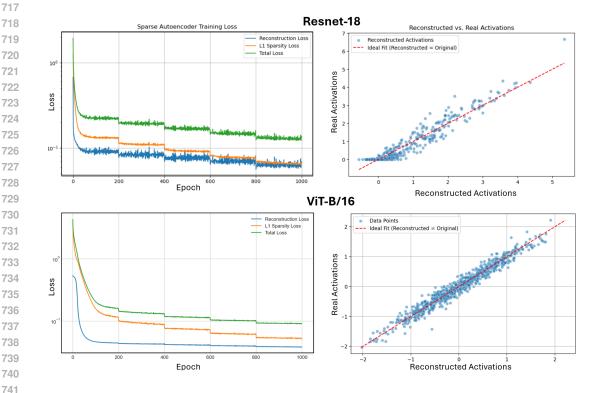


Figure 6: Autoencoder training and reconstruction performance for the ResNet-18 and ViT-B/16 backbones. Left: Training curves for the Sparse Autoencoders, showing total loss, reconstruction loss, and L1 sparsity loss (log-scaled) per epoch. Right: A comparison of the autoencoder's reconstructed feature activations with the original ones, where the red dashed line indicates perfect reconstruction (A = A).

VISUALIZATION OF LATENT FEATURES A.5

A.5.1ACTIVATION MAXIMIZATION

To visualize the abstract concept encoded by a feature, we employ activation maximization (27), extending it from traditional applications on individual neurons or filters to the discovered latent features. The method synthesizes an image from random noise by optimizing its pixels to maximally activate a chosen feature, revealing the visual pattern it has learned to detect.

In this procedure, we optimize an input image $x \in \mathbb{R}^{C \times H \times W}$, where C denotes the number of channels and H, W represent the image dimensions. The objective is to maximize the activation $\phi_l(x)$ of a target latent feature l, previously identified as dominant for a particular class. The optimization objective is defined as:

$$\underset{x}{\text{maximize}} \quad \phi_l(x) - \lambda_{L2} \cdot ||x||_2^2 - \lambda_{TV} \cdot TV(x), \tag{5}$$

where $\phi_l(x)$ is the activation of the *l*-th latent feature given input x, $||x||_2^2$ penalizes large pixel values, and TV(x) is the total variation loss. This encourages smoother and more natural visual structures by discouraging high-frequency noise and sharp changes between neighboring pixels, and is formally defined as:

$$TV(x) = \sum_{c=1}^{C} \sum_{i=1}^{H-1} \sum_{j=1}^{W-1} \sqrt{(x_{c,i+1,j} - x_{c,i,j})^2 + (x_{c,i,j+1} - x_{c,i,j})^2},$$
 (6)

The regularization parameters λ_{L2} and λ_{TV} control the strength of the respective penalties, helping ensure that the optimized image is visually coherent, rather than adversarial or noisy (27).

To further increase the robustness of the visualization, we apply a sequence of stochastic image transformations at each optimization step. Specifically, we apply constant padding (16 pixels, value 0.5), random jittering (8 pixels), random scaling $[-\delta\%, +\delta\%]$, and random rotation $[-\theta, +\theta]$. After augmentation, a center crop restores the input to a fixed resolution of 224×224 pixels. These transformations, inspired by prior work on robust feature visualization (8; 35), act as implicit regularizers that prevent overfitting to high-frequency artifacts and encourage the emergence of semantically meaningful patterns.

The complete procedure is given by the following steps:

- 1. Initialize the input image x with random noise.
- 2. For a fixed number of iterations:

- (a) Apply the sequence of image transformations (padding, jittering, scaling, rotation, center cropping).
- (b) Forward the transformed image through the network and encoder to obtain the latent activations.
- (c) Evaluate the objective function combining target activation, total variation loss, and L_2 regularization.
- (d) Update x using gradient ascent to maximize the objective function.

A.5.2 Grad-FAM Implementation

While standard Grad-CAM is designed to produce saliency maps that highlight image regions influencing the final prediction for a specific *class*, Grad-FAM (Gradient-weighted Feature Activation Map), repurposes this logic to visualize which parts of an image are responsible for activating a specific *latent feature* learned by the SAE.

The core modification lies in the target of the gradient calculation. In standard Grad-CAM, the heatmap is generated by weighting the feature maps F^k from a target layer by importance weights β_k^c . These weights are calculated by global average pooling the gradients of the score for a class c, y^c , with respect to the feature maps. Denoting the number of pixels in the feature map as P, the importance weight for a feature map k is given by:

$$\beta_k^c = \frac{1}{P} \sum_i \sum_j \frac{\partial y^c}{\partial F_{ij}^k} \tag{7}$$

In Grad-FAM, we replace the class score y^c with the activation value of a latent feature, ϕ_l , from our learned sparse representation. The importance weight β_k^l is therefore calculated based on the gradient of this specific latent feature's activation:

$$\beta_k^l = \frac{1}{P} \sum_i \sum_j \frac{\partial \phi_l}{\partial F_{ij}^k} \tag{8}$$

Here, F^k represents the k-th channel of the activation tensor from the target layer. While the formulation is model-agnostic, its application requires handling architectural differences. For a CNN like ResNet, F^k is an inherently spatial 2D feature map. A Vision Transformer, however, outputs a 1D sequence of tokens. To apply the same spatially-based pooling, these tokens must be reshaped back into a 2D grid. This transformation is achieved by first discarding the non-spatial [CLS] token from the sequence, then reshaping the remaining N patch tokens into their original $\sqrt{N} \times \sqrt{N}$ grid. This process correctly restores the visual layout because the token sequence preserves the original raster-scan order of the image patches.

The final heatmap, L^l , is produced by taking the absolute value of the weighted combination of feature maps. This differs from the conventional Grad-CAM method, which applies a ReLU. This choice is motivated by the design of our latent feature representation (ϕ_l) , which can take on both positive and negative values to represent concepts. Since a strong influence in either direction is a meaningful signal of importance, we take the absolute value to ensure our heatmap highlights all regions that significantly contribute to the feature's final activation.

$$L^{l} = \left| \sum_{k} \beta_{k}^{l} F^{k} \right| \tag{9}$$

To implement this, the forward pass must be explicitly chained from the base model through the autoencoder's encoder to establish a continuous computational graph for backpropagation. We make use of PyTorch hooks to capture the necessary intermediate data during this single forward and backward pass. The full procedure is summarized in the following pseudocode.

```
function Generate-Grad-FAM(model, autoencoder, image, target_layer, feature_index):
    # 1. Register hooks on target_layer to capture feature maps and gradients
    register_hooks(target_layer, forward_hook, backward_hook)
    # 2. Forward pass to get the latent score
    internal_activations = model.get_internal_activations(image)
    latent_vector = autoencoder.encoder(internal_activations)
    target_score = latent_vector[feature_index]
    # 3. Backward pass from the latent score
    target_score.backward()
    # 4. Generate the heatmap
    feature_maps = captured_feature_maps
    gradients = captured_gradients
    # For ViT, reshape token-based maps/gradients to a spatial grid
    if is_transformer(model):
        feature_maps = reshape_transform(feature_maps)
        gradients = reshape_transform(gradients)
    pooled_gradients = global_average_pool(gradients)
    weighted_maps = weight_feature_maps(feature_maps, pooled_gradients)
    # Take absolute value to get magnitude of influence
    heatmap = abs(sum(weighted_maps))
    normalized_heatmap = normalize(heatmap)
    return normalized_heatmap
```

By backpropagating from an intermediate latent feature instead of a final class logit, Grad-FAM shifts the focus from standard input attribution to the spatial localization of internal

concepts. The resulting saliency maps thus provide a direct, visual link between the abstract features learned by the autoencoder and their manifestation in the input data.

A.6 Critical Suppression Threshold

A.6.1 Derivation and Procedure: Analytical Method

Here, we provide the full derivation for the critical suppression threshold, $\alpha_{\rm crit}$, introduced in Section 3.2. The goal is to analyze how our weight modulation affects the feature contribution to the logit for a class i.

Let **x** be the activation vector from the penultimate layer. The feature-perturbed logit, $z'_i(\alpha)$, after applying weight suppression with scaling factor α , is then given by:

$$z_i'(\alpha) = \sum_j w_{ij}' x_j$$

$$= \sum_j w_{ij} \cdot \max(0, 1 - \alpha |c_j|) \cdot x_j,$$
(10)

where w_{ij} is the original weight and c_j is the contribution of the target latent feature to the original feature j. This definition isolates the direct contribution to the logit from the model features, excluding the global class bias term.

For small values of α , we can use the linear approximation: $\max(0, 1 - y) \approx 1 - y$. This approximation is valid under the condition that $\alpha < 1/|c_j|$ for all relevant components j (see A.6.3 for validation of this approximation).

Applying this, we get:

$$z_i'(\alpha) \approx \sum_j w_{ij} (1 - \alpha |c_j|) x_j$$

$$= \sum_j w_{ij} x_j - \alpha \sum_j |c_j| w_{ij} x_j$$

$$= z_i - \alpha R_i(\mathbf{x}),$$
(11)

where $z_i = \sum_j w_{ij} x_j$ is the original feature contribution to the logit, and we define the sensitivity term as:

$$R_i(\mathbf{x}) = \sum_j |c_j| w_{ij} x_j \tag{12}$$

This term, $R_i(\mathbf{x})$, quantifies how sensitive the logit z_i is to suppression along the direction of the chosen latent feature. We define the critical threshold as the value of α where the feature contribution of the logit is driven to zero, which gives the per-sample threshold for sample n:

$$z_i^{(n)} - \alpha_{\text{crit}}^{(n)} R_i(\mathbf{x}^{(n)}) = 0 \quad \Longrightarrow \quad \alpha_{\text{crit}}^{(n)} = \frac{z_i^{(n)}}{R_i(\mathbf{x}^{(n)})}$$
 (13)

Our model is trained for multi-class classification using a standard cross-entropy loss function. Because this loss function in PyTorch internally applies a log-softmax operation, the model's final layer is correctly designed to output raw, unbounded logits. For the purpose of our $\alpha_{\rm crit}$ derivation, we thus approximate the effect of suppression by treating these logits independently, providing a practical and interpretable proxy for class-sensitivity.

To ensure this computed factor is well-defined, we restrict the analysis to test samples $\{\mathbf{x}^{(n)}\}$ from class i for which the following conditions hold:

$$z_i^{(n)} > 0$$
 and $R_i^{(n)} > 0$. (14)

Positive logit requirement $(z_i^{(n)} > 0)$ This ensures the sample is initially classified in favor of class *i*. Since we aim to compute the scaling required to suppress a confident prediction, it only makes sense to include samples where the logit is already positive.

Positive relevance requirement $(R_i^{(n)} > 0)$ This avoids division by zero and ensures that the latent features are, on balance, contributing positively to the logit. Including samples with non-positive relevance would contradict the assumption that scaling down their contributions leads to suppression.

To compute $\alpha_{\rm crit}$ using this method, we follow this procedure:

1. Compute per-image thresholds.

For each qualifying sample $\mathbf{x}^{(n)}$, compute $z_i^{(n)}$ and $R_i^{(n)}$ and then calculate:

$$\alpha_{\text{crit}}^{(n)} = \frac{z_i^{(n)}}{R_i^{(n)}}.$$

2. Aggregate across images.

918

919 920

921

922

923 924

925

926 927

928 929 930

931

932 933

934

935 936

937 938

939

940 941

942 943 944

945

946 947

948

949

950 951

952

953

954

955

956

957 958

959 960

961

962

963

964

965

966

967

968

969 970

971

Collect the set $\{\alpha_{\text{crit}}^{(n)}\}$ for all qualifying images n for class i.

3. Calculate summary statistics.

Calculate summary statistics, e.g., the median $\tilde{\alpha} = \text{median}_n(\alpha_{\text{crit}}^{(n)})$.

Procedure: Numerical Method

To obtain a more precise estimate of the critical suppression threshold free from the inaccuracies of the linear approximation, we compute α_{crit} numerically. This approach directly finds the root of the exact modified logit equation:

$$z_i'(\alpha) = \sum_j w_{ij} \cdot \max(0, 1 - \alpha |c_j|) \cdot x_j = 0.$$
 (15)

We restrict the analysis to a subset of samples that meet specific criteria to ensure the results are well-defined and interpretable.

Positive logit requirement $(z_i^{(n)} > 0)$ This requirement is identical to that of the analytical method and is maintained for the same reason: the concept of suppressing a logit to zero is only meaningful for samples that are already classified in favor of the target class.

Existence of a Zero-Crossing While this method does not involve division, we must ensure that increasing α leads to suppression. A sufficient condition is that the function $z_i'(\alpha)$ is positive at $\alpha=0$ (covered by the first requirement) and becomes negative for some sufficiently large α . This guarantees that a root exists. Samples where the logit does not cross zero are excluded.

To compute α_{crit} numerically, we follow this procedure:

1. Find the per-image root. For each qualifying sample $\mathbf{x}^{(n)}$, we solve the equation $z_i'(\alpha^{(n)}) = 0$ for $\alpha^{(n)}$ using a simple and robust bracketing method. Specifically, we evaluate $z_i'(\alpha)$ over a discrete range of α values to find the interval where the function's sign changes from positive to negative. The precise root within this interval is then estimated using linear interpolation for higher precision.

2. Aggregate across images.

Collect the set $\{\alpha_{\text{crit}}^{(n)}\}$ for all images n of class i for which a root was found.

3. Calculate summary statistics.

Calculate relevant summary statistics from the collected set, e.g., the median $\tilde{\alpha}_{crit} = \text{median}_n(\alpha_{crit}^{(n)}).$

This numerical procedure yields the true critical suppression threshold for each sample by avoiding the limitations of the linear approximation.

A.6.3 VALIDATION OF THE ANALYTICAL APPROXIMATION

The derivation of the analytical $\alpha_{\rm crit}$ relies on a linear approximation that is valid when $1-\alpha\cdot|c_j|>0$ for all components j. To rigorously assess how frequently this condition holds, we calculated the full distribution of the term $1-\alpha_{\rm crit}\cdot|c_j|$. Specifically, for each per-sample numerical $\alpha_{\rm crit}$, we calculated its interaction with every component of the corresponding $|c_j|$ vector, yielding a the distribution of all pairwise suppression terms.

Figure 7 shows a clear architectural divergence. For the ResNet-18, the distribution is concentrated above zero, indicating that the linear approximation condition holds for the majority of samples. For the Vision Transformer, however, a significant mass of the distribution is in the negative region. This confirms that the linear approximation is a reasonable

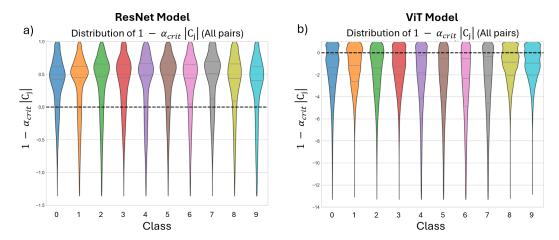


Figure 7: Validation of the linear approximation used to derive the analytical $\alpha_{\rm crit}$. The plots show the distribution of the suppression term $1 - \alpha_{\rm crit} \cdot |c_j|$, where negative values indicate the approximation is invalid. (a) For the **ResNet-18 model**, the distribution is concentrated above zero, confirming the approximation holds for the majority of samples. (b) For the **Vision Transformer**, a significant mass of the distribution is in the negative region, demonstrating that the linear approximation is invalid for the majority of samples.

approximation for the ResNet-18 but a poor one for the ViT. Consequently, we expect the analytical estimate of $\alpha_{\rm crit}$ to represent a reliable lower bound for ResNet-18 but a significantly more conservative lower bound for the ViT. This is consistent with the results presented in the main text for ResNet-18 (Figure 5) and in the Appendix for the ViT (Figure 12)

A.7 ROME METHOD FOR CLASS SUPPRESSION

To provide a point of comparison for permanent weight editing, we implemented a ROME-inspired rank-one update adapted for vision models. The original ROME method (22) was designed for factual editing in large language models by modifying mid-layer feed-forward networks. Our adaptation applies the same rank-one update principle to the final classification layer of a ResNet-18 model for a class suppression task. This section introduces the mathematical formulation and implementation workflow.

ROME treats a feed-forward network layer as a key-value memory. For our use case, we target the final fully connected layer (model.fc). The goal is to modify its weight matrix, W, such that a specific input activation ("key") produces a new desired output logit vector ("value").

1. **Key** (k): The key is the activation vector from the penultimate layer for a single, representative sample of the class to be suppressed. Let $k \in \mathbb{R}^M$ be this activation vector.

- 2. Value (v^*) : The desired new value is a target logit vector, $v^* \in \mathbb{R}^C$. To suppress class i, we define v^* by taking the original output logits, z^0 , and setting the logit for the target class to a large negative value (e.g., -10.0).
- 3. **Update Calculation:** The ROME update is a rank-one modification, Δ , to the weight matrix W. The error vector is defined as $\delta = v^* Wk$. Note that the term Wk recalculates the logits without the layer's bias term, k. Consequently, for the target class k, the error k is large (e.g., $-10.0 (z_i^0 b_i)$). For any other class k, the error is simply the bias term, k is large (e.g., k). The weight update is then given by:

$$\Delta = \frac{\delta k^T}{\|k\|_2^2} \tag{16}$$

The new weight matrix is then $W' = W + \Delta$. Since the bias terms are typically small, the error vector δ is dominated by its target component. This means that while the update matrix Δ is dense, its values are overwhelmingly concentrated in the row corresponding to the target class. The edit therefore primarily affects the weights connected to the target logit, with only minor adjustments to the weights for the other classes.

A.7.1 PSEUDOCODE FOR IMPLEMENTATION

return model

The general workflow for applying ROME to suppress a single class is outlined below:

```
function Apply-ROME-Edit(model, sample_image, target_class):
    # 1. Get the activation "key" from the penultimate layer
    key = get_penultimate_activation(model, sample_image)

# 2. Define the new target output "value"
    original_logits = model(sample_image)
    target_logits = original_logits
    target_logits[target_class_idx] = -10.0 # Suppress class by setting negative value

# 3. Calculate the rank-one weight update
    weights = model.fc.weight
    error_vector = target_logits - (weights @ key)
    key_norm_squared = dot(key, key)
    delta = outer_product(error_vector, key) / key_norm_squared

# 4. Apply the update to the model's weights
    model.fc.weight = weights + delta
```

A.7.2 Benchmarking ROME for Class Suppression

To contextualize the performance of our feature-based editing method, we benchmarked it against our ROME implementation on the same class suppression task detailed in Section 4.2.

As shown in Figure 8, the results are highly comparable. The confusion matrices illustrate that for the class suppression task, ROME achieves a similar outcome to our method. Both techniques effectively reduce the accuracy for the target class to near zero while leaving accuracy on other classes largely unaffected. This benchmark indicates that our latent feature-guided approach is competitive with established, state-of-the-art model editing methods for this type of intervention, while also offering the additional benefits of mechanistic interpretability and fine-grained feature control. We observed similar performance parity across the other classes.

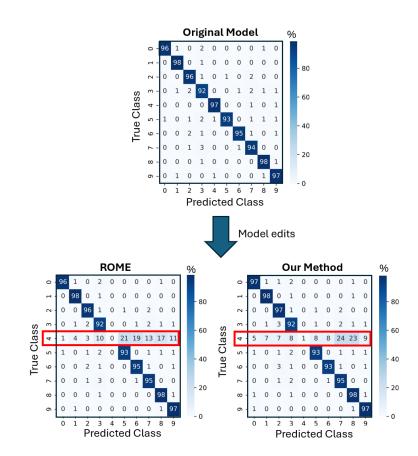


Figure 8: Comparison of class suppression performance. Confusion matrices showing the model's predictions on the Test set before editing (Top), and after editing (Bottom) with ROME and our proposed weight modulation method. Both methods successfully ablate the target class prediction (here shown for class 4)

A.8 VISION TRANSFORMER RESULTS

This section provides the results for the complementary validation analysis performed on the Vision Transformer (ViT-B/16) model.

A.8.1 LATENT FEATURES ENCODE SEMANTIC CONCEPTS

To confirm the generality of our feature discovery method, we first replicate the analysis of the learned feature basis on the ViT. As illustrated in Figure 9a, the SAE successfully learns a sparse and class-specific latent basis, analogous to the one found for ResNet-18. We then validate the semantic meaning of these discovered features using Grad-FAM to visualize their grounding in specific input images. Figure 9b shows examples for two class-dominant features, where the saliency maps correctly highlight the class-relevant objects in the images (the church building and the English springer spaniel). This provides strong evidence that the SAE uncovers semantically meaningful and spatially localized concepts for the Vision Transformer, just as it does for the CNN.

A.8.2 Controlling Class Predictions via Feature Editing

To validate the architectural robustness of our control method, we replicate the class suppression experiments on the Vision Transformer. The results, shown in Figure 10, confirm that our feature-guided intervention is effective also on the transformer-based model.

Specifically, suppressing the dominant feature for "Church" (Class 4) reduces its accuracy on the test set to near zero. Critically, the performance on other classes remains largely

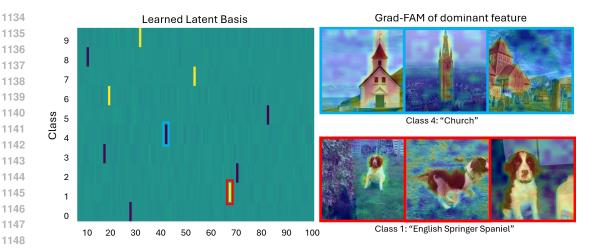


Figure 9: Validation of discovered features for the ViT-B/16 model. (a) Average latent feature activations across classes, highlighting a strong class-specific structure. (b) Grad-FAM visualizations grounding the dominant features for Class 4 ("Church") and Class 1 ("English springer") in representative images.

unaffected, as shown in the confusion matrix. This provides strong evidence that the features learned by the ViT also act as modular switches, allowing for precise, targeted interventions with minimal off-target effects. We observe similar successful suppression across the other classes in the dataset.

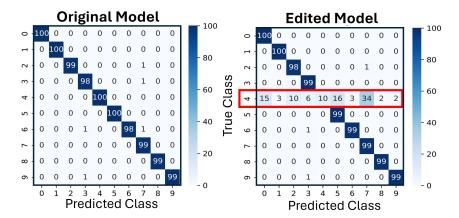


Figure 10: Quantitative validation of feature-guided control in the ViT-B/16 model. The confusion matrices show performance on the test set before (left) and after (right) suppressing the dominant feature for the "Church" class. The intervention effectively ablates the target class while preserving accuracy on the other classes.

A.8.3 QUANTIFYING INTERVENTION SENSITIVITY

To investigate the sensitivity of the interventions also for the Vision Transformer, we replicate the suppression analysis from the main text.

First, we analyze the characteristic suppression curve for a single class. As shown in Figure 11, systematically increasing the intervention strength α results in a stable accuracy that drops sharply past a critical threshold, while the model's predictions are reallocated to other classes. This shows that the intervention exhibits the same targeted behavior for the ViT as it did for the ResNet-18.

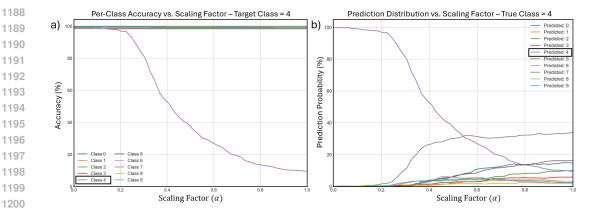


Figure 11: Suppression sensitivity for a single class (Class 4, "Church") in the ViT model. (a) Per-class accuracy as a function of the scaling factor α . (b) Distribution of predictions for images of the target class, showing how confidence is reallocated as the feature is suppressed.

Next, we extend this analysis to all classes and compute the analytical, numerical, and empirical estimates for $\alpha_{\rm crit}$. Figure 12a shows the full set of per-class suppression curves, which exhibit a wider spread in thresholds compared to the ResNet-18, though all classes could be suppressed to near-zero accuracy. Figure 12b presents the comparison of the $\alpha_{\rm crit}$ estimates. This plot confirms two key findings also discussed in the main text.

First, the analytical estimate of $\alpha_{\rm crit}$ represents a much more conservative lower bound. Second, we observe a larger discrepancy between the numerically calculated $\alpha_{\rm crit}$ (representing total evidence loss) and the empirical $\alpha_{50\%}$ (representing a prediction flip) for the ViT. We hypothesize that this discrepancy can be attributed to the ViT's architectural properties. Recent research has shown that ViTs learn a "curved" and non-linear representation space (15). This effect, combined with the ViT's tendency to produce less confident, more distributed logits (24), means a smaller intervention is required to be surpassed by a competitor. This widens the gap between the point of prediction failure ($\alpha_{50\%}$) and total evidence loss ($\alpha_{\rm crit}$), an effect less pronounced in the more linear representations of the ResNet model.

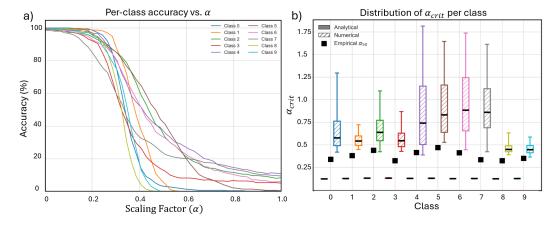


Figure 12: Full class-level sensitivity analysis for the ViT model. (a) Per-class accuracy vs. intervention strength α for all classes. (b) Comparison of analytical (filled), numerical (hatched), and empirical (square marker) $\alpha_{\rm crit}$ estimates, illustrating the inaccuracy of the linear approximation for the ViT.

Another factor that can contribute to the observed wider spread in α_{crit} thresholds and the increased difficulty in driving some classes to near-zero accuracy, is that of feature

entanglement. When we suppress a single dominant latent feature associated with a class, other entangled features might still contribute to the same class's logit. This makes it intrinsically harder to drive the total feature contribution for that class's logit completely to zero. This effect would result in larger required $\alpha_{\rm crit}$ values and can explain instances where suppression curves, as seen in Figure 12a, level off above zero accuracy even at high α . The combined characteristics of the model backbone's feature representation and the SAE's ability to extract a sparse decomposition can thus critically influence intervention efficacy. Further discussion on these intertwined effects is provided in Section A.9.2.

A.9 Additional Experiments on Class-Robustness

This section presents supplementary experiments assessing the sensitivity of the class-robustness results.

A.9.1 Robustness Across Autoencoder Initializations

We trained multiple autoencoders (n = 10) with identical architectures and loss parameters, but different random seeds, weight initializations, and data shuffling during training. This led each autoencoder to learn a distinct latent basis, as illustrated in Figure 13a.

For each learned representation, we repeated the experiments from Section 4.4, computing accuracy—vs— α curves (Figure 13b) and analyzed how model confidence redistributed across the classes (Figure 13c). Shaded regions in these plots indicate the standard deviations across the different realizations.

Despite learning distinct latent bases, the resulting suppression curves were nearly identical, and the same subset of classes remained most resistant to suppression. This observation aligns with findings on the stability of features learned by sparse autoencoders trained on the same data (29), suggesting that the heterogeneity in class robustness is not an artifact of a particular latent decomposition due to random initialization.

This indicates that for a fixed SAE training objective, the observed differences in class sensitivity stem mainly from the base model's intrinsic feature representations. However, the training objective itself, particularly the sparsity coefficient λ_1 , fundamentally shapes the feature decomposition. Therefore, the observed heterogeneity is a product of both the backbone's intrinsic structure and the specific features learned by the SAE. While a detailed analysis of the impact of the SAE training objective is a key topic for future work, the role of the backbone model is addressed further in the next section.

A.9.2 Robustness Across Backbone Models

To examine whether the backbone's learned feature space contributes to the observed class heterogeneity, we fine-tune several ResNet-18 variants on Imagenette with different random initializations and training parameters.

Prior work has reported that minor changes in training can alter class robustness profiles (2; 20). Consistent with these findings, we observe that the "hard-to-suppress" classes varies across backbones (Figure 14), suggesting that the base network's feature space is the dominant factor determining class heterogeneity.

In our experiments, we find that the optimizer's batch size during fine-tuning of the backbone is a key parameter that strongly affects class suppression sensitivity. When using small batches (e.g. 8–16), all classes can be driven to near-zero accuracy with moderate α , and the suppression curves remain tightly aligned in both low- and high- α regimes. In contrast, large batches (e.g. 64–128) substantially increases the heterogeneity of suppression, where $\alpha_{\rm crit}$ grows larger for many classes and only a subset of classes could be fully suppressed.

Some experiments also exhibit more complex suppression curves, where certain classes are only partially suppressed even under high attenuation, and a few show non-monotonic behavior, where accuracy initially drop, then partially recover at higher α . These suppression behaviors resemble the promotion-suppression effects observed in previous work on adversarial

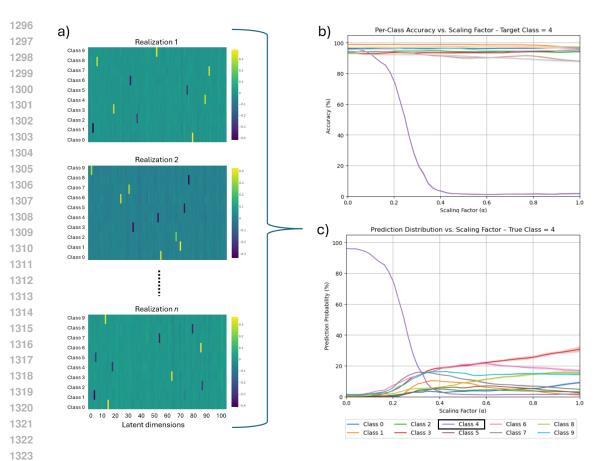


Figure 13: a) Learned latent bases across the n=10 realizations. b) Per-class prediction accuracy as a function of the scaling factor α , when suppressing the dominant feature associated with class 4. c) Distribution of predictions among the different classes, where the "true" class is 4, showing how the models confidence becomes re-distributed among the other classes. The shaded region in these plots represent the standard deviation across n=10 different realizations.

perturbations, where certain features are either suppressed or promoted, leading to changes in model predictions (34).

We hypothesize that these effects may rise from the well-known impact of batch size on the geometry of the loss landscape (13). Large-batch training tends to converge to sharper minima, which are known to produce more entangled internal representations. These entangled features are naturally harder to isolate and suppress with our targeted method. In contrast, the stochasticity of small-batch training acts as an implicit regularizer, guiding the model to flatter minima associated with more disentangled and modular features (10). In these solutions, concepts are more cleanly separated, allowing our interventions to suppress class predictions more effectively.

Beyond backbone training, our experiments also indicate that the parameters used for training the Sparse Autoencoder (SAE) significantly influence the learned latent feature representations and, consequently, their downstream controllability. Specifically, the ℓ_1 regularization loss plays a critical role: too low ℓ_1 values can lead to a less sparse, entangled latent space where multiple features activate for a single class, making targeted suppression difficult. Conversely, overly high ℓ_1 values may result in some class-specific features failing to activate at all. Achieving an effective balance between reconstruction accuracy and ℓ_1 sparsity is crucial, and the optimal parameter values depends on the specific backbone

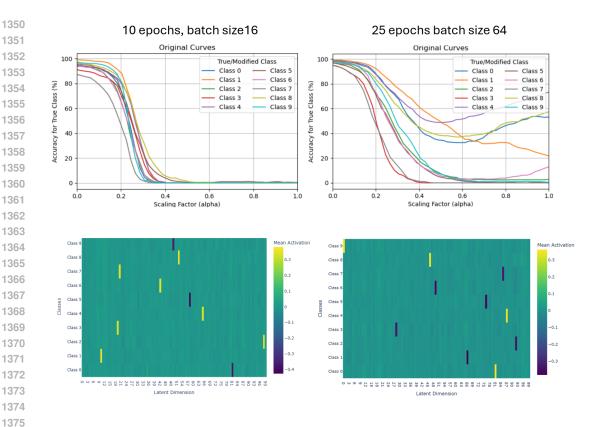


Figure 14: Accuracy vs. scaling factor α for ResNet-18 backbones trained with different initialization and fine-tuning parameters.

architecture and its training dynamics, necessitating tailored values for models like both ResNet and ViT.

Together, these findings highlight a critical link between training dynamics, architectural choices, and downstream controllability. They suggest that co-designing training regimes and architectures to favor flatter, more modular solutions may be a key strategy for developing models that are not only performant but also inherently more interpretable and editable. While our initial analysis provides evidence for this connection, a more comprehensive exploration of batch size, SAE parameters (especially ℓ_1 regularization), and other training parameters on feature modularity and intervention efficacy represents a crucial direction for future work.