SAE AS A CRYSTAL BALL : INTERPRETABLE FEATURES PREDICT CROSS-DOMAIN TRANSFERABILITY OF LLMs WITHOUT TRAINING

Anonymous authors

Paper under double-blind review

ABSTRACT

In recent years, pre-trained large language models have achieved remarkable success across diverse tasks. Besides the pivotal role of self-supervised pre-training, their effectiveness in downstream applications also depends critically on the posttraining process, which adapts models to task-specific data and objectives. However, this process inevitably introduces model shifts that can influence performance in different domains, and how such shifts transfer remains poorly understood. To open up the black box, we propose the SAE-based Transferability Score (STS), a new metric that leverages sparse autoencoders (SAEs) to forecast posttraining transferability. Taking supervised fine-tuning as an example, STS identifies shifted dimensions in SAE representations and calculates their correlations with downstream domains, enabling reliable estimation of transferability before fine-tuning. Extensive experiments across multiple models and domains show that STS accurately predicts the transferability of supervised fine-tuning, achieving Pearson correlation coefficients above 0.75 with actual performance changes. Beyond this, we take an initial step toward extending STS to reinforcement learning. We believe that STS can serve as a principled tool for guiding post-training strategies in LLMs.

1 Introduction

Recent advances in large-scale neural networks have demonstrated that pre-training on massive datasets yields models with strong generalization capabilities (Achiam et al., 2023; Dubey et al., 2024; Yang et al., 2025; Liu et al., 2024). However, due to discrepancies between the pretraining objectives and the specific requirements of downstream tasks, pretraining alone is often insufficient to achieve optimal performance on these tasks. Post-training, which includes supervised fine-tuning (Zhang et al., 2023; Luo et al., 2023), and reinforcement learning (Schulman et al., 2017; Shao et al., 2024), plays a critical role in bridging this gap. By selectively adapting the pretrained model, post-training improves performance in target tasks, and allows models to better capture domain-specific characteristics.

However, during the post-training process, it is widely observed that improvements on a target task often come at the expense of performance in other domains (Dong et al., 2023; Kumar et al., 2022). For instance, (Li et al., 2025) state that improvements in the reasoning ability of large language models come at the cost of reduced model robustness. Despite these observations, the mechanisms underlying how model features are linked and transferred during post-training remain largely unexplored. As a result, we currently lack the ability to predict which domain performance is likely to benefit or deteriorate under specific post-training adaptations, limiting both interpretability and principled design of post-training strategies.

In this paper, we analyze the transferability of post-training through the enhanced interpretability provided by sparse autoencoders (Ng et al., 2011). The sparse autoencoder (SAE) is an encoder-decoder architecture that reconstructs the internal activations of models while enforcing sparsity constraints on the hidden layer. Previous works have shown that the SAE encoder features achieve monosemanticity (Cunningham et al., 2023; Gao et al., 2024), where each dimension is only activated by a certain natural concept. Leveraging this property, we observe that post-training only mod-

055

056

060

061

062

063

064

065

066

067

068

069

071

073 074

075

076

077

078 079

081

082

084

085

089

090

092

093

094

095

096

098

102

103

104

105

107

ifies certain SAE dimensions—for example, those associated with mathematical reasoning. This observation motivates a natural approach to predict the transferability of post-training: we can identify the shifted SAE features and examine their correlations with different domains.

Concretely, our analysis consists of two steps: (1) identifying the dimensions that are shifted during post-training, and (2) assessing their correlations with downstream domains. In the first stage, the primary challenge is to identify the shifted dimensions **prior to post-training**. Inspired by the observation that in-context learning exhibits behaviors similar to supervised fine-tuning (Wang et al., 2023; Mosbach et al., 2023), we forecast the shifted dimensions by using the supervised answers as demonstrations for in-context learning, and then identify the dimensions that undergo the largest changes. Empirical results show a clear overlap between the predicted and actual shifted dimensions. In the second stage, leveraging the interpretability of SAE activations, we observe that the activation values of these shifted dimensions in a domain can capture their correlation. We formalize this as the SAE-based transferability score (STS), which quantifies how strongly the shifted dimensions relate to downstream tasks. A higher STS suggests a larger expected performance change after supervised fine-tuning. Empirically, we find that our metric consistently correlates well with actual performance shifts—for instance, the Pearson correlation coefficient exceeds 0.75 when evaluating performance variations across domains in the MMLU-Pro dataset (Wang et al., 2024). At last, we provide a preliminary exploration of extending our metric to reinforcement learning settings. Together, these results allow us to develop a principled framework for predicting the cross-domain transferability without training. We summarize our contributions as follows:

- We propose a method to identify shifted dimensions in supervised fine-tuning without requiring access to the fine-tuned models. We observe that when supervised answers are used as context prompts, the shifted dimensions in in-context learning substantially overlap with those in supervised fine-tuning. This result allows us to predict shifted dimensions without performing any additional training.
- We propose a new metric, the SAE-based Transferability Score (STS) to predict the transferability of LLMs without training. STS builds on the observation that activation values in the SAE feature space reflect correlations between dimensions and domains. By combining with the identified shifted dimensions, STS provides a principled way to predict the transferability of supervised fine-tuning without requiring access to the fine-tuned models.
- We empirically demonstrate that higher STS values strongly correlate with larger performance changes in supervised fine-tuning, which consistently achieve Pearson correlation coefficients above 0.75 across diverse scenarios. These results confirm that STS provides a reliable, finetuning-free metric for predicting the cross-domain transferability of LLMs.

2 RELATED WORK & PRELIMINARY

Post-training. Post-training refers to the stage after large-scale pretraining, where a pretrained model is further adapted to align with specific objectives, user preferences, or downstream applications. Unlike pretraining, which relies on massive unlabeled data (Achiam et al., 2023; Yang et al., 2025), post-training typically leverages smaller but higher-quality datasets. The methods in post-training can be majorly divided into supervised fine-tuning (Zhang et al., 2023; Luo et al., 2023) and reinforcement learning (Schulman et al., 2017; Shao et al., 2024). In this paper, we mainly focus on SFT, and discuss extensions to RL. During the SFT process, given a set of labeled examples $\{x_i, y_i\}$, the model parameters are updated to minimize the discrepancy between the model's predictions and the ground-truth answers via Negative Log-Likelihood (NLL) Loss:

$$\mathcal{L}_{SFT}(\Theta) = -\mathbb{E}_{x_i} \log p(y_i|x_i;\Theta).$$

where Θ denotes the model parameters. By aligning the representations with task-specific supervision, SFT effectively transfers the general knowledge encoded in pretrained models to particular applications.

Previous studies have shown that supervised fine-tuning (SFT) often reduces performance in domains beyond the target task (Dong et al., 2023; Kumar et al., 2022). For example, (Li et al., 2025) find that SFT on reasoning data improves the mathematical ability of LLMs but weakens their robustness against jailbreak attacks. These results suggest that it is important to understand how model

capabilities transfer during post-training. While some works have studied this problem (Sun et al., 2025; Huan et al., 2025), most of them focus on post-hoc analysis after training. However, such approaches are less practical because they cannot help predict transfer effects before the fine-tuning process starts. This limitation motivates our work, where we aim to build a method that can predict transferability without fine-tuning.

Sparse Autoencoders. Although large language models (LLMs) have demonstrated remarkable performance across a wide range of downstream tasks, many of their decisions and internal behaviors remain opaque, which hinders broader deployment in applications. To address this issue, sparse autoencoders (SAEs) have been proposed as a promising framework for improving the mechanistic interpretability of LLMs (Gao et al., 2024; Cunningham et al., 2023). Concretely, given a hidden representation $z \in \mathbb{R}^d$ within the network, an SAE employs an encoder–decoder architecture to project z into a sparse latent representation and reconstruct it back to the original space. For instance, in the case of a top-K SAE (Gao et al., 2024), the encoding-decoding process can be formulated as:

$$h = \text{TopK}(W_e z - b),$$

$$\hat{z} = W_d h + b.$$
(1)

The encoder representation h is computed via a linear transformation defined by $W_e \in \mathbb{R}^{s \times d}$ and a bias $b \in \mathbb{R}^s$, while the decoder reconstructs the input features using $W_d \in \mathbb{R}^{d \times s}$. The SAE is trained by minimizing the reconstruction loss:

$$\mathcal{L}_{SAE}(W_e, W_d, b) = \|\hat{z} - z\|^2.$$

Previous studies (Gao et al., 2024) have shown that when the encoder features are sufficiently sparse (e.g., $K \ll s$), the resulting representations often display monosemanticity. In other words, each feature dimension is only activated by a certain natural concept, such as a mathematical definition, a physical property, or a linguistic pattern.

In-context Learning. In-context learning denotes the ability of large pretrained models to solve tasks by conditioning on demonstrations provided in the input (Kossen et al., 2023; Wang et al., 2025). Formally, given a context consisting of k labeled examples

$$C = \{(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)\},\$$

the model receives a new query input x_{k+1} and generates the output \hat{y}_{k+1} by leveraging the conditional distribution learned during pretraining:

$$\hat{y}_{k+1} \sim p_{\theta}(y \mid x_{k+1}, \mathcal{C}),$$

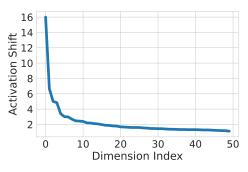
where θ denotes the fixed pretrained parameters. Unlike supervised fine-tuning and reinforcement learning, in-context learning (ICL) adapts to new tasks during inference without requiring additional training. Nevertheless, several studies have shown that models under in-context learning still exhibit many similarities to those trained with supervised fine-tuning and reinforcement learning (Mosbach et al., 2023; Wang et al., 2023).

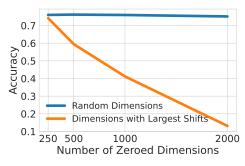
3 CAPTURING SHIFTED FEATURES IN SUPERVISED FINE-TUNING

In this paper, we analyze the transferability of supervised fine-tuning across different domains by understanding how feature representations shift within neural networks. Since sparse autoencoders (SAEs) enhance monosemanticity by disentangling overlapping representations, their features provide a clearer interpretability for tracking representation shifts in different domains. Consequently, in this section, we start by analyzing how the supervised fine-tuning process modifies SAE features and how the shifted features can be predicted in advance of the fine-tuning process.

3.1 SFT-INDUCED CHANGES IN SAE FEATURES

Prior works have shown that SAE encoder features exhibit strong interpretability, with each dimension corresponding to a certain natural concept. Since supervised fine-tuning (SFT) is generally tailored to specific downstream tasks and targeted capabilities, we wonder whether it primarily affects only a small subset of SAE dimensions tied to task-relevant features. Taking the mathematical ability as an example, we investigate how SAE features change when models are fine-tuned on the





- (a) Activation Shifts on SAE Dimensions
- (b) Test Accuracy with Zeroed SAE Dimensions

Figure 1: Analysis of feature shifts induced by supervised fine-tuning (SFT). We fine-tune Qwen2.5-7B-Instruct on the LIMO (a mathematical reasoning dataset) and examine shifts of SAE features on the residual stream at layer 25. Figure (a) shows the distribution of shift magnitudes while Figure (b) shows accuracy on Math-LightEval when progressively zeroing the dimensions with the largest shifts. The results indicate that SFT primarily affects a small subset of SAE dimensions tied to specific model capabilities.

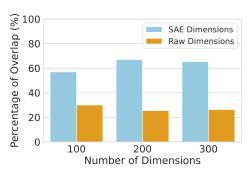
math dataset LIMO (Ye et al., 2025). As noted in (Lieberum et al., 2024), the same dimension in an SAE usually continues to represent similar concepts after fine-tuning. Consequently, we extract monosemantic features before and after the SFT process using the same SAE on the residual streams of 25-th layer in Qwen2.5-7B-Instruct (Team, 2024).

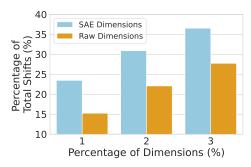
As shown in Figure 1a, we find that changes in SAE features are largely concentrated in a small subset of dimensions. With calculation, we find that the top-100 dimensions account for 25% of the total change, indicating that the SFT process primarily affects only a limited portion of SAE features. Furthermore, to observe the relationship between features shifted during SFT and the mathematical ability of the model, we rank SAE features according to their magnitude of change during fine-tuning and then evaluate model performance on the Math-Lighteval dataset (Hendrycks et al., 2021) by zeroing out different numbers of features. As shown in Figure 1b, performance on math tasks drops rapidly when the shifted features are removed, whereas the model retains strong mathematical abilities when random SAE features are zeroed. These results indicate that the SFT process primarily changes a small subset of SAE features that are closely associated with specific model capabilities.

3.2 IDENTIFYING SHIFTED DIMENSIONS VIA IN-CONTEXT LEARNING

As previously discussed, supervised fine-tuning modifies only a small subset of SAE features that correspond to specific semantics. Intuitively, these features are crucial for studying properties of SFT, such as transferability. However, their identification typically requires examining the model after fine-tuning, which confines the analysis to a post-hoc perspective. This limitation motivates a central question of our work: can such features be identified prior to the fine-tuning process, thereby enabling a predictive understanding of transferability?

To solve this challenge, we draw on the connection between supervised fine-tuning (SFT) and incontext learning (ICL). Previous works demonstrate that ICL can obtain similar performance to SFT in the large language models (Wang et al., 2023; Mosbach et al., 2023). Consequently, this motivates us to investigate whether the SAE features shifted during ICL and SFT are consistent. To verify this hypothesis, we respectively sort the SAE dimensions according to their changes after SFT and ICL. To be specific, we conduct experiments on Qwen2.5-7B-Instruct. For SFT, we employ ground-truth chain-of-thoughts (CoTs) on LIMO as supervision, while for ICL, we use the same CoTs as context prompts. As shown in Figure 2a, we observe substantial overlap between the shifted features in SFT and ICL; for example, 57% of the top 100 most shifted SAE dimensions coincide. These findings confirm that **the shifted dimensions in ICL and SFT are highly consistent**, suggesting that the shifted dimensions can be identified before the SFT process.





- (a) Degrees of Overlaps between Estimated and Actual Shifted Dimensions
- (b) Percentage of Overall Activation Change from Top-Shifted Dimensions

Figure 2: Overlap between estimated and actual shifted dimensions. Figure (a) demonstrates that the SAE shifted dimensions predicted by ICL substantially overlap with the actual shifted dimensions identified after SFT, whereas applying the same method directly on raw dimensions is less effective. Figure (b) further shows that raw model dimensions, prior to applying SAE, are influenced more uniformly by the SFT process, thereby limiting the ability to identify crucial shifted dimensions.

We have shown that, by leveraging the connection between ICL and SFT, the shifted SAE dimensions can be identified prior to the SFT process. A natural question then arises: are SAEs necessary for this method? Intuitively, as dimensions in the raw feature space are highly polysemantic and entangled (Elhage et al., 2022), specific abilities of the model (e.g., mathematical reasoning) are distributed across multiple dimensions. As a result, the SFT process tends to affect a broader set of features, which makes it more difficult to identify crucial shifted dimensions. To verify the analysis, we further conduct experiments on the raw features before applying SAEs. We again sort the features by their changes during SFT and ICL. As shown in Figure 2b, the shifted features before SAE are more uniformly distributed. We then investigate the influence on the accuracy of identifying shifted dimensions. With the same empirical settings, Figure 2a shows that the overlap between shifted features in ICL and SFT is much reduced in the raw feature space. These findings indicate that the enhanced monosemanticity introduced by SAEs is crucial for identifying shifted features prior to supervised fine-tuning.

4 PREDICTING THE TRANSFERABILITY OF SUPERVISED FINE-TUNING ACROSS DOMAINS

In Section 3, we introduced a method for predicting shifted SAE dimensions prior to the SFT process. Intuitively, when applying fine-tuned models to downstream tasks, if the shifted dimensions are closely related to a given domain, the SFT influence on that domain will be stronger. Thus, understanding the transferability of SFT across domains requires analyzing the correlations between shifted dimensions and different domains. In this section, we present a metric for evaluating this correlation and predict transferability based on that. Specifically, Section 4.1 introduces the proposed metric, while Section 4.2 validates it across diverse scenarios.

4.1 METRICS FOR MEASURING CROSS-DOMAIN CORRELATIONS

Due to the enhanced monosemanticity of SAEs (Cunningham et al., 2023), the activations in the SAE feature space become interpretable, meaning that the top-activated sequences within a given dimension usually share similar semantics. This property allows us to associate each dimension with specific semantic concepts. Consequently, if the sequences from a particular domain exhibit higher activation values in a certain dimension, it implies that this dimension is more strongly correlated with that domain. Building on this intuition, we can use the degrees of activations across dimensions to quantify domain–feature correlations and further analyze the transferability of supervised fine-tuning.

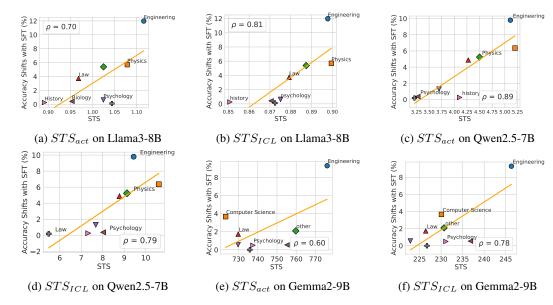


Figure 3: The Pearson correlation (ρ) between STS and actual performance shifts on MMLU-Pro induced by SFT on LIMO. We extract SAE features from Llama3-8B-Instruct, Qwen2.5-7B-Instruct, and Gemma2-9B-Instruct. During the evaluation process, we select four MMLU-Pro domains with the largest and smallest performance shifts under SFT. The detailed performance shifts can be found in Appendix A. The results demonstrate that STS consistently maintains a strong correlation with the actual performance changes.

Formally, we define our metric as the SAE-based transferability score (STS). In the first step, given an SFT dataset $\mathcal{T}=\{x_i,y_i\}$, we extract SAE features for each sample. Let the SAE features be denoted as $h(x_i;\Theta)$, where Θ represents the parameters of the pretrained model, and h is the SAE encoder trained on top of the pretrained model. Similarly, with in-context learning, the features are denoted as $h(x_0,y_0,\cdots,x_t,y_t,x_i;\Theta)$, where $\{x_0,y_0,\cdots,x_t,y_t\}$ are the context prompts. We then identify the N dimensions with the largest changes, i.e.,

$$D_N = \text{TopN}(\mathbb{E}_{x_i} || h_j(x_i; \Theta) - h_j(x_0, y_0, \cdots, x_t, y_t, x_i; \Theta) ||^2),$$

where h_j denotes the j-th dimension of the SAE features.

In the second step, given a downstream domain dataset $\tilde{T} = \{\tilde{x}_i\}$, we compute the activation values on the shifted dimensions identified in the first step:

$$\mathrm{STS}_{\mathrm{act}}(\tilde{\mathcal{T}}) = \mathbb{E}_{\tilde{x}_i} \sum_{j \in D_N} h_j(\tilde{x}_i; \Theta).$$

It is important to note that **we do not use the model after SFT in this estimation**, which means that our metric serves as a predictive measure rather than a post-hoc analysis.

Besides directly computing the average activation values, we introduce an alternative method based on in-context learning (ICL) to capture the correlation between the downstream domain and the shifted dimensions. We know that ICL leverages multiple demonstrations to guide the model in performing downstream tasks, effectively injecting domain-specific signals into the representation. Intuitively, by comparing the SAE features extracted with and without ICL demonstrations, we can isolate the effect of domain knowledge on the shifted dimensions. This difference reflects how strongly the features are modulated by task-relevant information, thereby offering a reliable estimate of the correlation between a domain and the shifted dimensions.

Consequently, we estimate the correlation by measuring the difference between features extracted with and without in-context demonstrations. Formally, let $\{\tilde{x}_i, \tilde{y}_i\}_{i=1}^m$ denote m ground-truth question-answer pairs in the downstream domain. We define the metric as

$$\mathrm{STS}_{\mathrm{ICL}}(\tilde{\mathcal{T}}) = \mathbb{E}_{\tilde{x}_i} \sum_{j \in D_N} \|h_j(\tilde{x}_0, \tilde{y}_0, \cdots, \tilde{x}_m, \tilde{y}_m, \tilde{x}_i; \Theta) - h_j(\tilde{x}_i; \Theta)\|^2.$$

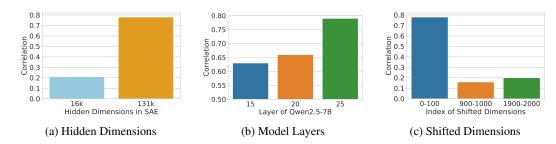


Figure 4: Ablation studies on the implementation of our metric. We evaluate (a) SAEs with varying hidden dimensions in the representation space, (b) SAEs trained on different layers of the pre-trained model, and (c) different ranges of top-shifted dimensions.

This metric captures how much the presence of domain-relevant context (the demonstrations) influences the shifted dimensions, providing another reliable estimate of domain-feature correlation besides using maximum activation values.

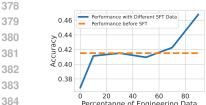
4.2 Empirical Verifications on Predicting the Transferability

Based on the proposed metric, we now empirically evaluate the correlation between the STS score and the actual transferability across different downstream domains. Specifically, we use the LIMO dataset (Ye et al., 2025) as the SFT training set, which contains 817 high-quality mathematical examples. We fine-tune three models (Qwen2.5-7B-Instruct, Llama3-8B-Instruct, and Gemma2-9B-Instruct) on LIMO, and compare their performances before and after SFT on different domains of MMLU-Pro (Wang et al., 2024). To extract SAE features, we apply SAEs to the residual streams of the models prior to SFT. For computing the STS metric, we employ two ground-truth CoTs from LIMO as in-context demonstrations to identify the top-100 shifted dimensions. When estimating the correlations between the domains and the predicted shifted dimensions, we use five ground-truth CoTs from the domain of MMLU-Pro as prompts to calculate $\rm STS_{ICL}$. More details of the experiments can be found in the Appendix.

As shown in Figure 3, the correlation between STS and performance changes across different domains remains consistently high for all three models, with Pearson correlation coefficients exceeding 60%. These findings validate that our metric provides a reliable estimation of the transferability of the SFT process. Besides, we note that the performance of $STS_{\rm ICL}$ is more stable than $STS_{\rm act}$ (the coefficients keeps above 75%), which suggests that leveraging in-context learning yields a more accurate estimation of the correlation between domains and SAE dimensions than relying solely on activation values.

In the following, we conduct a series of ablation studies to evaluate the performance of STS across different scenarios. We start by investigating the importance of enhanced monosemanticity provided by SAEs. We note that when identifying shifted dimensions, the monosemanticity enables us to find the small subsets of dimensions related to the SFT dataset. While for estimating the dimension-domain correlation, the monosemanticity makes the activation values interpretable and enables us to find the correlation. To empirically verify the importance of monosemanticity, we calculate STS with different SAEs. Since larger hidden dimensions in SAE feature spaces is believed to enhance monosemanticity (Gao et al., 2024), we extract features using SAEs on the residual streams of models with different hidden dimensions (16k and 131k). Figure 4a clearly shows that weaker monosemanticity leads to a substantial drop in prediction accuracy, demonstrating that strong monosemanticity is essential for the success of STS.

We further evaluate the effectiveness of STS by comparing SAEs applied to different model layers and using different SAE implementations. Specifically, we extract features from the residual streams of layers 15, 20, and 25 in Qwen2.5-7B-Instruct, and we consider two SAE variants: one based on topK activations and the other using JumpReLU activations. As shown in Figure 4b, STS maintains a consistently strong correlation with performance changes across domains, regardless of the layer or SAE variant used. The results further verify the effectiveness of our methods in different scenarios.



387

388

389

390

391

392

393

394 395

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

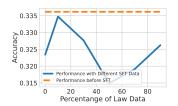
425

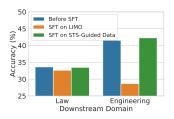
426

427 428

429 430

431





(a) Performance with Different Amounts of Engineering Data

Percentange of Engineering Data

- (b) Performance with Different Amounts of Law Data
- (c) Performance Shifts with Different SFT Data

Figure 5: Comparison of data mixture strategies in the SFT process. We focus on the domains with the largest (engineering) and smallest (law) performance shifts induced by SFT of Qwen2.5-7B-Instruct on LIMO. In total, 220 extra examples from a mixture of engineering and law data are added. Figure (a) reports engineering performance with varying amounts of engineering data, while Figure (b) reports law performance with varying amounts of law data. Figure (c) compares the downstream performance without additional data and with additional data mixed according to the ratio of their corresponding STS values.

In addition to evaluating different SAEs, we examine another critical hyperparameter of our metric: the range of estimated shifted dimensions. Selecting dimensions with small shifts risks including those unaffected by the SFT process. Since correlations between unaffected dimensions and downstream domains can introduce noise, the accuracy of STS is consequently undermined. As shown in Figure 4c, the correlation between STS and downstream performance decreases when smallershifted dimensions are selected, supporting this analysis. These findings underscore that the choice of shifted dimensions directly impacts metric reliability, highlighting the need for an appropriate selection strategy that only identify dimensions with largest shifts.

5 APPLICATIONS: A DATA MIXTURE PRINCIPLE

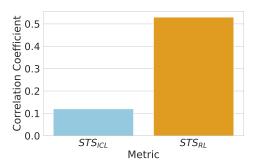
In Section 4, we demonstrated that our proposed metric, the SAE-based transferability score (STS), exhibits a strong correlation with actual performance changes. Building on this result, we now explore a practical application of STS. Specifically, we leverage predicted transferability to optimize data mixture strategies during post-training. Using STS, we can identify the domains most likely to be affected by supervised fine-tuning (SFT). A common approach to mitigate performance degradation in such domains is to introduce additional data. Intuitively, with a predicted ranking of performance changes across domains, we can allocate more data to those at higher risk of degradation. To validate this idea, we examine two domains in MMLU-Pro: engineering (with the largest performance drop when fine-tuning Qwen2.5-7B-Instruct on LIMO) and law (with the smallest). In our experiments, we fine-tune Qwen2.5-7B-Instruct on LIMO with additional data in MMLU-Pro. We split MMLU-Pro into training and testing sets with a 1:1 ratio. For training, we use Qwen2.5-7B-Instruct's outputs prior to SFT as supervised answers. In addition to the original LIMO data, we augment the training set with 220 extra examples sampled from the mixture of engineering and law domains.

As shown in Figure 5, we observe that domains with larger performance changes require more additional data. For example, allocating more data to the engineering domain leads to substantial improvements (Figure 5a), whereas allocating extra data to the law domain yields only marginal gains (Figure 5b). Moreover, when the data mixture ratio is adjusted to align with the ratio of their corresponding STS values, the resulting model achieves balanced performance across both engineering and law (Figure 5c). These findings suggest that STS can serve as a principled guide for designing data mixture ratios, enabling more effective post-training while mitigating uneven performance shifts across domains.

EXPLORATIONS ON REINFORCEMENT LEARNING

In this paper, we primarily investigate the impact of supervised fine-tuning across different domains. Nevertheless, reinforcement learning (RL) represents another non-negligible post-training paradigm,





- (a) Overlaps between Estimated and Actual Shifted Dimensions Induced by SFT and RL
- (b) Different Implementations of STS Metric

Figure 6: The explorations on extending STS metric to the RL scenarios. We optimize Qwen2.5-7B-Instruct on Math-LightEval using GRPO. Figure (a) shows that the identified dimensions in RL show smaller overlap with the actual shifted dimensions compared to SFT, which limits the accuracy of the metric. Figure (b) shows that when the estimated shifted dimensions are replaced with the actual ones, the correlation between STS and downstream performance changes becomes much stronger.

motivating us to explore whether our method can be extended to RL. We begin by applying the STS metric directly, as in the supervised fine-tuning setting. Specifically, we train Qwen2.5-7B-Instruct using the GRPO framework (Shao et al., 2024) on the Math-LightEval dataset and evaluate performance changes across domains in MMLU-Pro. When computing STS, we follow the same procedure as in supervised fine-tuning: ground-truth CoTs from Math-LightEval serve as demonstrations to estimate shifted dimensions, and correlations between these dimensions and downstream domains are calculated based on in-context learning. However, as shown in Figure 6b, STS exhibits low correlations with performance changes in the RL setting. In the following, we try to find the reasons behind this discrepancy.

We note that a key distinction between SFT and RL is that RL lacks access to ground-truth answers, making it challenging to select appropriate demonstrations for in-context learning. As a result, the estimation of shifted features may be inaccurate. To test this hypothesis, we compare the overlaps between the actual and predicted shifted dimensions. As shown in Figure 6a, the overlap is substantially lower in RL than in SFT. To further validate this observation, we replace the estimated top-100 dimensions with the ground-truth dimensions with the largest changes after RL and recompute the STS metric. Figure 6b demonstrates that STS calculated with ground-truth shifted dimensions shows a strong correlation with performance changes, indicating that the main challenge lies in accurately estimating shifted dimensions in RL prior to training. And this will be the future direction of our explorations. For now, STS in RL can serve as a metric to predict transferability without evaluating RL models on downstream tasks. While its effectiveness does not yet match that in SFT, it can still reduce computational overhead, particularly when assessing multiple abilities across various domains.

7 Conclusion

In this work, we introduce a metric for predicting the transferability of post-training in large-scale neural networks, leveraging the interpretability of sparse autoencoders (SAEs). By identifying the SAE dimensions that are shifted during post-training and measuring their correlations with downstream domains, we propose the SAE-based transferability score (STS) as a predictive indicator of performance changes. Our experiments show that STS reliably forecasts performance variations across multiple domains, providing new insights into the internal mechanisms of post-training. Beyond supervised fine-tuning, we further demonstrate the applicability of our approach to reinforcement learning settings. Overall, we believe our work establishes a principled framework for understanding and anticipating post-training effects, paving the way for more targeted and effective post-training strategies.

ETHICS STATEMENT

This work makes use of publicly available datasets and models. No private or sensitive data are involved, and no harmful content is included. Therefore, we believe this paper does not raise any ethical concerns.

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our results, we provide comprehensive details of our experiments in both the main paper and the appendix. In Section 3, we describe the details of the methodology for identifying shifted dimensions, while Section 4 presents the procedure for computing correlations between these dimensions and downstream domains. Additionally, Appendix A reports the concrete performance changes across different MMLU-Pro domains, and Appendix B offers further details and complementary discussions. The code will be released upon the publication of this paper.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.
- Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are affected by supervised fine-tuning data composition. *arXiv* preprint arXiv:2310.05492, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.
- Maggie Huan, Yuetai Li, Tuney Zheng, Xiaoyu Xu, Seungone Kim, Minxin Du, Radha Poovendran, Graham Neubig, and Xiang Yue. Does math reasoning improve general llm capabilities? understanding transferability of llm reasoning. *arXiv preprint arXiv:2507.00432*, 2025.
- Jannik Kossen, Yarin Gal, and Tom Rainforth. In-context learning learns label relationships but is not conventional learning. *arXiv preprint arXiv:2307.12375*, 2023.
- Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv* preprint *arXiv*:2202.10054, 2022.
- Ang Li, Yichuan Mo, Mingjie Li, Yifei Wang, and Yisen Wang. Are smarter llms safer? exploring safety-reasoning trade-offs in prompting and fine-tuning. *arXiv* preprint arXiv:2502.09673, 2025.
- Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*, 2024.

- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. arXiv preprint
 arXiv:2412.19437, 2024.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
 - Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. Fewshot fine-tuning vs. in-context learning: A fair comparison and evaluation. *arXiv preprint arXiv:2305.16938*, 2023.
 - Andrew Ng et al. Sparse autoencoder. CS294A Lecture notes, 72(2011):1–19, 2011.
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
 - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
 - Yiyou Sun, Shawn Hu, Georgia Zhou, Ken Zheng, Hannaneh Hajishirzi, Nouha Dziri, and Dawn Song. Omega: Can llms reason outside the box in math? evaluating exploratory, compositional, and transformative generalization. *arXiv preprint arXiv:2506.18880*, 2025.
 - Qwen Team. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2, 2024.
 - Qixun Wang, Yifei Wang, Xianghua Ying, and Yisen Wang. Can in-context learning really generalize to out-of-distribution tasks? In *ICLR*, 2025.
 - Xindi Wang, Yufei Wang, Can Xu, Xiubo Geng, Bowen Zhang, Chongyang Tao, Frank Rudzicz, Robert E Mercer, and Daxin Jiang. Investigating the learning behaviour of in-context learning: a comparison with supervised learning. *arXiv preprint arXiv:2307.15411*, 2023.
 - Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In *NeurIPS*, 2024.
 - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
 - Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning. *arXiv preprint arXiv:2502.03387*, 2025.
 - Yifan Zhang, Jingqin Yang, Yang Yuan, and Andrew Chi-Chih Yao. Cumulative reasoning with large language models. *arXiv preprint arXiv:2308.04371*, 2023.

A PERFORMANCE CHANGES AFTER SFT

Table 1: Performance change on different domains of MMLU-Pro (Llama3-8B-Instruct fine-tuned on LIMO).

Performance Change (%)	engineering 11.97	physics 5.7	chemistry 5.39	law 3.72
Performance Change (%)	philosophy 2.61	other 2.49	health	computer science 0.97
Performance Change (%)	economics 0.95	math 0.89	psychology 0.63	biology 0.42
Performance Change (%)	history 0.26	business 0.12		

Table 2: Performance change on different domains of MMLU-Pro (Qwen2.5-7B-Instruct fine-tuned on LIMO).

Performance Change (%)	engineering 9.8	chemistry 6.36	physics 5.24	computer science 4.88
Performance Change (%)	business 4.43	health 2.93	math 2.89	economics 2.61
Performance Change (%)	philosophy 2.4	biology 1.53	other 1.3	psychology 0.37
Performance Change (%)	history 0.26	law 0.19		

Table 3: Performance change on different domains of MMLU-Pro (Gemma2-9B-Instruct fine-tuned on LIMO).

Performance Change (%)	engineering 9.29	computer science 3.66	other 2.05	law 1.73
Performance Change (%)	health	math	business	history
	1.47	1.41	1.27	0.79
Performance Change (%)	economics	chemistrcy	biology	physics
	0.71	0.62	0.56	0.54
Performance Change (%)	psychology 0.05	philosophy 0		

We show the concrete performance change of models after SFT on LIMO in Table 1,2, 3.

B EXPERIMENTS DETAILS

When fine-tuning the pre-trained models Qwen2.5-7B-Instruct, Llama3-8B-Instruct, and Gemma2-9B-Instruct on the LIMO dataset, we adopt a unified experimental setup across models. Specifically, the maximum prompt length is set to 2048 tokens, with sequences truncated from the left to fit within this constraint. All models are fine-tuned for 10 epochs using four H20 GPUs, where each GPU runs with a micro-batch size of 2.

For feature extraction with sparse autoencoders (SAEs), we use one SAE for each backbone model on a specific layer of the transformer. Concretely, we use an SAE with 16,384 hidden dimensions on the 25th layer of Qwen2.5-7B-Instruct, an SAE with 65,536 hidden dimensions on the 25th layer of

Llama3-8B-Instruct, and an SAE with 131,072 hidden dimensions on the 31st layer of Gemma2-9B-Instruct. These configurations strike a balance between reconstruction fidelity and interpretability.

To identify shifted SAE dimensions induced by post-training, we first sample activations from 20,000 tokens before and after in-context learning (ICL), where the ICL demonstrations are constructed using supervised answers from LIMO. We then compute the activation differences across SAE dimensions and select those exhibiting the largest shifts. Finally, to quantify the correlations between shifted dimensions and downstream domains, we sample activations from 10,000 tokens for each domain and compute their correlations with the identified dimensions

C THE USE OF LARGE LANGUAGE MODELS (LLMS)

In this work, the use of LLMs was limited to minor language editing to improve readability. All conceptual development, theoretical analysis, experimental design, and result interpretation were conducted independently by the authors. Thus, the use of LLMs was purely auxiliary and had no impact on the scientific contributions of this paper.