

# Data-scarce Behavior Editing of Language Models

Anonymous ACL submission

## Abstract

Large Language Models trained on web-scale text acquire language generation abilities that can solve a wide range of tasks, particularly when task knowledge is refined into the generative prior using in-context examples. However, spurious features learned from noisy data hinder their generalizability. Supervised finetuning can enhance task specificity but may lead to data inefficiency. Prior studies indicate that (i) noisy neural circuitries coexist with generalizable ones within LLMs, and (ii) finetuning typically enhances (or suppresses) existing abilities without introducing newer ones. Building upon these, we propose TaRot, a novel method for task adaptation. TaRot intervenes in the neural circuitries using learnable rotation matrices that are optimized using Bayesian optimization, on labelled samples in the order of standard few-shot prompting examples. Experiments on multiple classification and generation tasks using LLMs of varying sizes reveal the efficacy of TaRot, improving upon both zero- as well as few-shot performance, with average improvements (across models and tasks) of 15.6% and 14%, respectively.

## 1 Introduction

Large Language Models (LLMs) acquire the ability to associate different language concepts presented in a sequential context by optimizing the prediction probability of the next token given a context. Despite its apparent simplicity, when scaled across web-sized text corpora, such a learning strategy introduces the ability to solve a wide range of tasks presented in natural language. However, the web contains almost everything humankind has written, and therefore, it introduces spurious token associations that are irrelevant or even counter-productive to the model to become generalized task-solvers. We observe phenomena like brittle few-shot performance (Sclar et al., 2024), hallucination (Huang et al., 2023), harmful text generation (Wen et al.,

2023), etc. as evidence of learning noisy patterns. Remedial interventions like instruction tuning (Zhang et al., 2024), alignment tuning (Shen et al., 2023), etc. have been proposed. Recent research has shown that such mediation only acts on a superficial level — out-of-distribution inputs can reinforce noisy behavior and break the model (Ghosh et al., 2024). Without an in-depth understanding of the inner workings, remedial strategies become wild goose chase.

Mechanistic disentangling of Transformer-based language models has shed some light on this direction (Elhage et al., 2021; Olsson et al., 2022; Wang et al., 2023). Two recent investigations (Jain et al., 2024; Prakash et al., 2024) on the effects of fine-tuning confirm the inability of supervised fine-tuning to alter fundamental abilities acquired via pretraining. On a tangential investigation, Dutta et al. (2024) recently confirmed the existence of multiple parallel neural pathways of answer processing within LLMs. Bhaskar et al. (2024) echoed similar findings in the case of syntactic generalization while pointing out that different components acquire different generalization behaviors. These findings lead us to the central research question of this work: *is it possible to edit the model behavior by editing internal representations in a generalizable manner?* Prior work in this direction has heavily relied on careful manual effort to localize task-specific neural components and design intervention techniques (Meng et al., 2022; Li et al., 2024a). Two key shortcomings limit the scalability of such methods: (i) Localization complexity grows polynomially with model size, making it difficult to identify task-relevant components and design effective ablations; (ii) Redundant components performing similar neural computations hinder the generalizability of any single intervention.

**Our contributions.** To this end, we propose a novel intervention technique, TaRot – Task-aware

**Rotation** of token-association (see Figure 1 for a representative depiction)<sup>1</sup>. We establish the conceptual prior from Transformers’s implicit gradient descent bias in next token prediction. Specifically, we first show that attention-weighted averaging of value vectors facilitates the memorization of token association from pertaining data in individual attention heads, in the sense that each attention head acts as a mini-language model. Due to the vast number of token associations present in the pretraining corpus compared to the number of attention heads in even the largest of the models, we hypothesize that individual directions of these memorized associations remain in superposition, and removal or downscaling of a head can counteract model performance. Instead, we construct parametrized rotations to align head outputs for task-adaptation. The rotation parameters are then optimized using Bayesian optimization. Furthermore, TaRot is *extremely data- and compute-efficient*: we use 6-20 supervised examples for each task and  $\frac{dL}{4}$  rotation parameters (where  $d$  is the model dimension and  $L$  is the number of layers) for each different task. This renders TaRot at par with standard few-shot prompting in labeled data-efficiency.

We experiment with four different classification tasks and two natural language generation tasks; the choice of tasks seeks to investigate general world knowledge (news topic classification) as well as the ability to generalize beyond imitation (BIG Bench tasks (BIG-bench authors, 2023)). TaRot demonstrates consistent improvements over six different language models of varying sizes: Qwen2-1.5B-Instruct, Phi-3-mini-4k-instruct, Mistral-7B-Instruct-v0.1, Meta-Llama-3-8B-Instruct, Qwen2.5-14B-Instruct, and Qwen2.5-32B-Instruct, in both zero-shot as well as few-shot settings. Furthermore, we analyze the changes in neural representation introduced by TaRot to uncover useful insights.

## 2 Related Work

Our work is primarily relevant to two broad areas of existing literature: adaptation of pretrained language models to downstream tasks, and mechanistic understanding and intervention techniques.

**Task adaptation of pretrained language models.** The *pretrain-finetune* regime for adapting language models to downstream tasks dates back to the

early approaches like BERT (Devlin et al., 2019) — pretrain a language model (LM) on large unstructured text corpora using self-supervised objective, followed by supervised fine-tuning on task-specific, relatively smaller datasets. Despite the apparent simplicity, the pitfalls of this regime have been pointed out in terms of *distribution shift* (Kumar et al., 2022). With the development of large-scale, autoregressive Transformer-based language models and their ability to learn from in-context examples (Brown et al., 2020), a definitive shift has happened in the more recent past. Current practices of using these models for downstream tasks primarily rely on designing suitable prompt templates and labeled example retrieval for in-context learning (ICL) (Liu et al., 2022; Rubin et al., 2022; Tanwar et al., 2023); traditional techniques of fine-tuning have taken a back seat due to the computational cost and catastrophic forgetting introduced by small-scale task-specific data that hurts the pre-trained abilities (Zhai et al., 2024). Instead, fine-tuning to follow task instructions, *aka* instruction-tuning (Zhang et al., 2024), has gained popularity. Instruction-tuning has been shown to introduce zero-shot task adaptation abilities in LLMs (Wei et al., 2022). Additionally, different methods of alignment tuning have been proposed with the primary goal being aligning the generative distribution of the language models with human values and preferences (Shen et al., 2023; Wang et al., 2024b). Despite the popularity of instruction and alignment tuning, their ability to alter fundamental information processing has been put in question in recent literature. Jain et al. (2024) investigated the effects of fine-tuning in toy models trained with formal languages as well as precompiled ones; their findings suggest that supervised fine-tuning does not introduce any new ability into pretrained models but only reinforces (or suppresses) existing ones. Similar concerns have been raised upon investigating entity tracking in the neural representation space (Prakash et al., 2024). Ghosh et al. (2024) identified multiple limitations of instruction tuning, including the inability to introduce new knowledge and deterioration of performance due to over-reliance on pattern matching.

**Mechanistic understanding and interventions.** The umbrella of mechanistic interpretability broadly encompasses methods to disentangle model behavior via reverse engineering the underlying neural algorithm (Elhage et al., 2021; Ferrando et al., 2024). Endeavors to mechanistically under-

<sup>1</sup>The source code of TaRot is attached with the supplementary and will be made public upon acceptance of the paper.

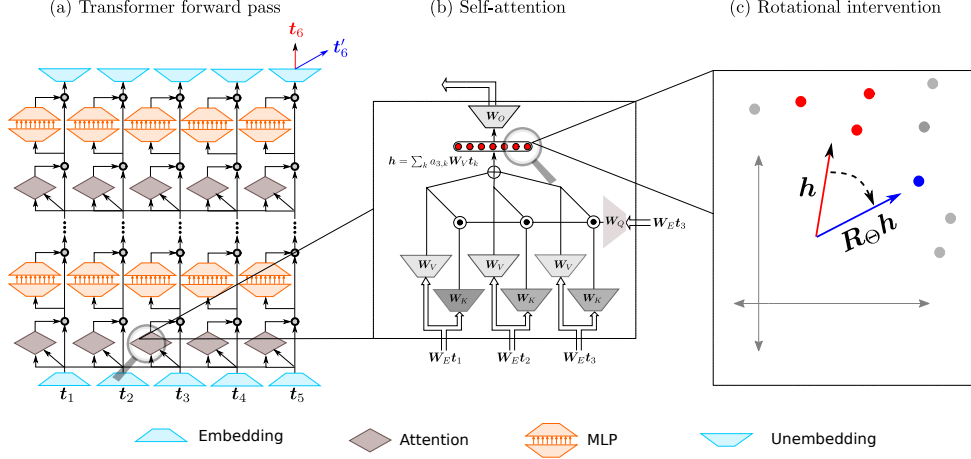


Figure 1: **A conceptual illustration of TaRot.** (a) Model generates an undesired next token  $t_6$  upon an input token sequence. (b) A certain attention head is responsible for associating certain input tokens with the undesired output. (c) TaRot learns a parametrized rotation operator  $R_\Theta$  that rotates  $h$  to the direction of the desired token (red to blue). The intervention results in a change in the forward pass in (a) that outputs the desired token  $t'_6$ .

stand Transformer-based language models trace back to the seminal work by Elhage et al. (2021). Their framework established attention heads as one of the fundamental building blocks of language model interpretation. Subsequent studies have identified the functional roles of different attention heads in pretrained models: induction heads as a primary mechanism of prefix matching (Olsson et al., 2022), circuitries of attention heads responsible for indirect object identification (Wang et al., 2023), neural pathways that implement chain-of-thought reasoning (Dutta et al., 2024), etc. Much relevant to our analysis, Lv et al. (2024) found that certain attention heads memorize the association between country names and their capitals. On a tangential line of investigation, Geiger et al. (2024) introduced the Distributed Alignment Search (DAS) framework for localizing interpretable features in subspaces of the neural representations. Mechanistic methods provide actionable insights that have led to non-traditional techniques to edit model behavior. Elhage et al. (2021) experimented with key propagation to elicit induction heads (and thereby, prefix-matching ability) in single-layer attention-only Transformers. Meng et al. (2022) used causal tracing to locate factual associations in MLP neurons and proposed a gradient-free approach to edit factual recall patterns in pretrained language models. Li et al. (2024a) identified attention head circuitry that elicits toxic text generation in GPT-2; mean-ablation of these circuits is shown to reduce toxicity. Self-detoxification (Leong et al., 2023) identifies toxic generation direction in the internal representation using trigger prompts and then

rewrites in the opposite direction to reduce toxicity. Wang et al. (2024a) formulated toxicity reduction as a knowledge editing task that can permanently alter toxic behaviors instead of suppressive interventions like supervised fine-tuning or RLHF-based alignment. Lamparth and Reuel (2024) localized backdoor mechanisms (i.e., vulnerabilities against adversarial prompt injections) in early-layer MLPs and proposed a low-rank substitution to improve robustness against such injections. Vergara-Browne et al. (2024) employed attribution patching techniques to identify and remove certain singular values in the parameter matrices to improve performance.

In comparison with prior intervention approaches, our work bears two fundamental differences: (i) TaRot does not necessitate task-specific localization of neural behaviors; this significantly reduces intense manual effort and risk of over-localization, eliciting efficient, generalizable interventions; (ii) TaRot is gradient-free, parameter-efficient, and requires supervised samples in the order of standard ICL; this poses TaRot as a practical alternative to intense prompt-engineering.

### 3 Methodology

In this section, we demonstrate the role of attention heads in memorizing token associations. Next, we lay out the working principles of TaRot.

#### 3.1 Attention heads as token-token maps

Inspired by Elhage et al. (2021), we dissect the Transformer-based language models with the following assumptions: (i) each attention head reads

from and writes to the residual stream independently in a linear fashion, and (ii) given that the attention heads utilize hidden representation of dimensionality much smaller than the residual stream (i.e., for a model with 16 attention heads, each attention head uses 1/16-th of the dimension of the residual stream), they typically operate on small subspaces of the residual stream. This way, two attention heads can operate on two distinct subspaces and never interact with each other. These two assumptions allow us to interpret the working of the attention heads meaningfully even while treating each head in isolation. We start with identifying what a single-head attention operation tends to learn in isolation.

Following the standard terminology (Elhage et al., 2021), we represent the embedding and unembedding matrices as  $\mathbf{W}_E \in \mathbb{R}^{d \times V}$  and  $\mathbf{W}_U \in \mathbb{R}^{V \times d}$ , where  $d$  and  $V$  are the dimensionality of the residual stream and the token space, respectively, the query, key, value, and output projection matrices denoted as  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, \mathbf{W}_O \in \mathbb{R}^{d \times d}$ , respectively. Given a sequence of input tokens as one-hot column vectors  $\mathbf{T} = \{\mathbf{t}_1, \dots, \mathbf{t}_n\}$ , the forward pass for single-layer attention-only Transformer can be written as:

$$\hat{\mathbf{t}}_{n+1} = \mathbf{W}_U \left( \mathbf{W}_E \mathbf{t}_n + \mathbf{W}_O \sum_i a_{n,i} \mathbf{W}_V \mathbf{W}_E \mathbf{t}_i \right) \quad (1)$$

where  $a_{n,i} = \frac{\exp(\mathbf{t}_n^\top \mathbf{W}_E^\top \mathbf{W}_Q^\top \mathbf{R}_{\Theta, n-i} \mathbf{W}_K \mathbf{W}_E \mathbf{t}_i)}{\sum_j \exp(\mathbf{t}_n^\top \mathbf{W}_E^\top \mathbf{W}_Q^\top \mathbf{R}_{\Theta, n-j} \mathbf{W}_K \mathbf{W}_E \mathbf{t}_j)}$  is the softmax-attention probability from source token  $\mathbf{t}_i$  to destination token  $\mathbf{t}_n$ , and  $\hat{\mathbf{t}}_{n+1} \in \mathbb{R}^V$  is the logit of the predicted next token. Upon reparametrization of  $\mathbf{W}_U \mathbf{W}_O \mathbf{W}_V \mathbf{W}_E$  as  $\mathbf{W}_{OV}$ , we can rewrite Equation 1 as

$$\hat{\mathbf{t}}_{n+1} = \mathbf{W}_U \mathbf{W}_E \mathbf{t}_n + \sum_i \mathbf{W}_{OV} \mathbf{t}_i \quad (2)$$

Note that  $\mathbf{W}_{OV} \in \mathbb{R}^{V \times V}$ , denoted as OV-circuits by Elhage et al. (2021), maps a distribution over tokens to another distribution over tokens. If the true token is  $\mathbf{t}_{n+1}$  with  $I(\mathbf{t}_{n+1})$  denoting its index (i.e., index of 1 in  $\mathbf{t}_{n+1}$ ), then the typical language modeling loss can be calculated as:

$$\mathcal{L}(\hat{\mathbf{t}}_{n+1}, \mathbf{t}_{n+1}) = -\log \left( \frac{\exp(\hat{\mathbf{t}}_{n+1}^{(I(\mathbf{t}_{n+1}))})}{\sum_k \exp(\hat{\mathbf{t}}_{n+1}^{(k)})} \right) \quad (3)$$

We can compute the gradient dynamics of the OV-circuit (with unit batch size and zero momentum) using Equations 2 and 3 as follows:

$$\mathbf{W}_{OV}^{(s+1)} = \mathbf{W}_{OV}^{(s)} + \eta \mathbf{t}_{n+1} \left( \sum_i a_{n,i} \mathbf{t}_i \right)^\top - \eta \text{SoftMax}(\mathbf{t}_{n+1}) \left( \sum_i a_{n,i} \mathbf{t}_i \right)^\top \quad (4)$$

where  $\mathbf{W}_{OV}^{(s)}$  and  $\mathbf{W}_{OV}^{(s+1)}$  are the OV-circuit parameters before and after the  $s$ -th gradient update step, respectively and  $\eta$  is the learning rate. The positive incremental component in the right-hand side of Equation 4 dictates that, when applied on an attention-weighted linear combination of the context tokens, OV-circuits learn to memorize a linear combination of possible next tokens.

However, in a deep Transformer model with several attention heads, MLP blocks and layer normalization, we can not determine the exact token-token map for the OV-circuits of attention head. Moreover, as Elhage et al. (2021) suggested, multiple attention heads across different layers can construct compositions, where the deeper heads use the output of the shallower heads. Alternatively, we can view each head as memorizing how to write in a specific direction in the residual stream, given a sequence of residual vectors—effectively acting as a *mini-LM*. When pretrained on web-scale corpora, these heads may memorize spurious token-token associations that harm downstream performance or introduce unsafe behaviors.

### 3.2 Editing model behavior via attention rotation

A natural conclusion from the prior discussion would be that, by suppressing undesired associations for certain attention heads, we can improve task performance. However, multiple token associations are expected to be memorized in each attention head in superposition since the number of attention heads is way smaller than the potential token associations present in the pretraining data—one cannot selectively switch off one certain association. Prior research in mechanistic interpretability has shown that, although we can often localize attention heads responsible for particular task, removing the non-dominant attention heads does not deliver the performance of the full model (Wang et al., 2023; Dutta et al., 2024).



Instead, one can *rotate* the output of the attention heads in order to maximize its alignment with rows of  $\mathbf{W}_U$  corresponding to certain tokens while near-orthogonalizing with certain undesired tokens. This way, the model behaviour can be edited without destroying the superposed associations. Defining the complete space of  $d \times d$  rotation matrices and optimizing them can become computationally challenging. Instead, we utilize the fact that any  $d \times d$  orthonormal matrix is similar to a block-diagonal matrix  $\mathbf{R}_\Theta$ , where  $\Theta = \{\theta_1, \dots, \theta_{d/2}\} \subset [0, 2\pi)^{\frac{d}{2}}$ , defined as:

$$\mathbf{R}_\Theta^d = \begin{bmatrix} \mathbf{B}(\theta_1) & 0 & \dots & 0 \\ 0 & \mathbf{B}(\theta_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{B}(\theta_{d/2}) \end{bmatrix} \quad (5)$$

where

$$\mathbf{B}(\theta_i) = \begin{bmatrix} \cos \theta_i & -\sin \theta_i \\ \sin \theta_i & \cos \theta_i \end{bmatrix}$$

Given the multi-head attention with  $H$  heads at layer  $l \in [L]$ , where  $L$  is the total number of layers in the Transformer, defined as:

$$\text{Attn}_l(\mathbf{x}_n^{(l)} \| [\mathbf{x}_1^{(l)}, \dots, \mathbf{x}_n^{(l)}]) = \mathbf{W}_O \parallel \sum_{h=1}^H a_{n,i}^{(h,l)} \mathbf{W}_V^{(h,l)} \mathbf{x}_i^{(l)}$$

where  $\parallel$  is the concatenation operator,  $a_{n,i}^{(h,l)}$  and  $\mathbf{W}_V^{(h,l)}$  denote the attention probability between source and destination residual streams at layer  $l$   $\mathbf{x}_i^{(l)}$  and  $\mathbf{x}_n^{(l)}$  and the value projection matrix corresponding to the attention head with index  $h \in [H]$  at layer  $l$ , respectively; we define the rotated attention as:

$$\text{RotAttn}_l(\mathbf{x}_n^{(l)} \| [\mathbf{x}_1^{(l)}, \dots, \mathbf{x}_n^{(l)}]) = \mathbf{W}_O \mathbf{R}_{\Theta_l}^d \parallel \sum_{h=1}^H a_{n,i}^{(h)} \mathbf{W}_V^{(h)} \mathbf{x}_i^{(l)} \quad (6)$$

Note that the block-diagonal definition of  $\mathbf{R}_\Theta^d$  in Equation 5 implies that applying  $\mathbf{R}_\Theta^d$  on the concatenated head outputs is equivalent to applying  $H$ -distinct  $\mathbf{R}_\Theta^{d/H}$  on each of the head outputs.

Without prior knowledge of which attention heads are responsible for memorizing undesired token associations, we need to apply the intervention defined in Equation 6 on a set of attention

blocks at layers  $l \in \hat{\mathbb{L}}$  (see Section 4 for the choice of the set  $\hat{\mathbb{L}}$ ). Then, the intervened forward pass is denoted as:

$$\hat{\mathbf{t}}_{n+1} = \mathcal{M}_{\text{Rotated}} \left( \{\mathbf{t}_1, \dots, \mathbf{t}_n\} | \Theta_O, \Theta_R \{\Theta_l | l \in \hat{\mathbb{L}}\} \right) \quad (7)$$

where  $\Theta_O$  is the set of pretrained model parameters, and  $\Theta_R$  are the parameters of rotations, and  $\mathcal{M}_{\text{Rotated}}$  denote the function representing the language model upon the designed intervention.

### 3.3 Optimization of rotation parameters

With the rotational interventions defined, all that we are left with is to optimize the rotational parameters. Let  $\mathcal{D} := \{\mathbf{T}_j, \mathbf{Y}_j | j \in [D]\}$  be a set of  $D$  supervised examples for a given task, with  $\mathbf{T}_j$ ,  $\mathbf{Y}_j$  referring to the sequence of tokens corresponding to the input and gold output, respectively. If  $\mathbf{Y}_j = \{\mathbf{y}_j\}$  is a single label token, the cost function to optimize becomes straightforward:

$$\max_{\Theta_R} \sum_j p \left( \mathcal{M}_{\text{Rotated}} \left( \mathbf{T}_j | \Theta_O, \Theta_R \{\Theta_l | l \in \hat{\mathbb{L}}\} \right) = \mathbf{y}_j \right) \quad (8)$$

where  $\Theta \subset [0, 2\pi)$ . For NLG tasks, maximizing the aggregate probability of all the generated tokens can be a solution. However, the goal of our rewiring method is to minimize undesired behaviors. When a model demonstrates such behaviors, depending upon the task, not all tokens equally correspond to the behavior under inspection. The pretrained model is trained using teacher-forcing and is generally able to generate grammatically correct responses. Hence, trying to align the model generation to a single reference response does not make much sense. Instead, we opt for a surrogate scoring function  $s : \{\mathbf{Y}_j\} \rightarrow \mathbb{R}$  that scores the “desirability” of a generated response. We let the model with rotation intervention to generate a complete response given an input, compute the score for the generated response, and seek to minimize the aggregate score across  $\mathcal{D}$ . We implement Bayesian optimization (Snoek et al., 2012) to solve the optimization problem depending upon the task. However, standard Gaussian Process with Matern kernel fails to scale to high dimension input space (Li et al., 2024b). Instead, Infinite-width Bayesian Neural Networks (I-BNN), proposed by Lee et al. (2017), has shown to scale effectively with high-dimensional parameter space<sup>2</sup>. The I-

<sup>2</sup>Here the term “high dimension” is relatively used. Our method seeks to optimize only the rotation configurations that

BNN covariance function does not rely on Euclidean distance, enabling the Gaussian Process to model non-stationary functions, an advantage since rotational effects may vary across the configuration space.

## 4 Experiment Setup

**Training setting.** Dutta et al. (2024) previously found that token associations corresponding to pre-trained knowledge primarily resides in the initial half of the model. Since the rotational intervention designed in Equations 6 and 7 are primarily targeted towards undesired token associations acquired through pretraining, we restrict  $\mathbb{L}$  to the initial half only. Therefore, the total number of parameters to optimise becomes  $\frac{dL}{4}$ . Since we want to optimise the rotation matrix for a particular task, only a small subset of training samples is required, i.e.,  $6 \leq D_{\text{training}} \leq 20$ .

**Models.** Six different instruction-tuned models with varying size are used for all experiments: Qwen2-1.5B-Instruct (Yang et al., 2024), Phi-3-mini-4k-instruct (Abdin et al., 2024) (2.8 billion parameter), Mistral-7B-Instruct-v0.1 (Jiang et al., 2023), Meta-Llama-3-8B-Instruct (Dubey et al., 2024), Qwen2.5-14B-Instruct (Team, 2024), and Qwen2.5-32B-Instruct (Team, 2024); we refer to these models as Qwen2-1.5B, Phi-3-mini, Mistral-7B, Llama-3-8B, Qwen2.5-14B, and Qwen2.5-32B, respectively.

**Tasks.** We experiment with four different classification (i.e., single token generation) tasks and two NLG tasks. Classification tasks used are as follows: AG News (Zhang et al., 2015), Entailed Polarity (Srivastava et al., 2022), Navigate (Srivastava et al., 2022), and Winowhy (BIG-bench authors, 2023). The generation tasks used include Imdb Positive Review (Maas et al., 2011), and Detoxify (Gehman et al., 2020) Further details and examples of tasks are available in Appendix A.1

**Baselines.** We compare TaRot with four different baselines: *Base model*, Eigen Pruning (Vergara-Browne et al., 2024), *RED (Representation EDiting)* (Wu et al., 2024), and Rescaling (additional Details in Appendix A.2).

**Evaluation metrics.** For NLG tasks, Imdb and Detoxify, two different types of reward models are

scales as  $\mathcal{O}(Ld)$ , which is substantially low-dimensional if compared to the parameter space of the LM itself.

used. To calculate the fluency of the generated text, GPT4 (Achiam et al., 2023) is used as an oracle. For both the tasks, the average of fluency and the score from the reward models are reported. Further details are present in Appendix A.3

## 5 Results

Tables 1 and 2 summarize the performance of various methods across classification tasks in zero- and 6-shot settings, respectively. Eigen Pruning is included only in zero-shot comparisons, per its original design. Table 3 presents results for NLG tasks.

**Consistent improvement with TaRot.** Across LLMs of varying sizes, TaRot consistently ranks as the best or second-best method across all tasks. Notably, it achieves relative gains in task-wise average F1 scores of 13.7%, 1.1%, 8.9%, 13%, 3.2%, and 1.3% over the base versions of Qwen2-1.5B, Phi-3-mini, Mistral-7B, Llama-3-8B, Qwen2.5-14B, and Qwen2.5-32B, respectively, in the zero-shot setting (see Table 1). The only exceptions are the Entailed polarity task with Qwen2-1.5B and Winowhy with Qwen2.5-32B, where TaRot slightly underperforms (e.g., 0.98 F1 vs. perfect score). In contrast, baseline methods like Eigen Pruning and Rescaling suffer from inconsistency—while they may improve performance in some cases, they often cause severe drops without any clear task- or model-specific patterns. For instance, Eigen Pruning improves Qwen2-1.5B on all but one task, yet fails on all tasks with Phi-3-mini.

**In-context examples vs. TaRot.** Unlike Eigen Pruning (or even, traditional fine-tuning), TaRot is optimized with a mixture of M-shot inference to avoid zero-shot bias. Consequently, we can observe the improvement over the base model achieved via TaRot while provided with in-context examples, except with Mistral-7B on AG News and Navigate (c.f. Table 2). Specifically, we observe improvements with Qwen2-1.5B on AG News, Entailed Polarity, and Winowhy; and with Phi-3-mini, Llama-3-8B, Qwen 2.5 14B, and Qwen 2.5 32B across all tasks. For Mistral-7B, gains are limited to Entailed Polarity and Winowhy.

**Importance of rotation over rescaling attention heads.** Comparing TaRot with the rotation-free Rescaling approach highlights key differences in intervention effectiveness. Rescaling is often brittle, with no consistent performance pattern. For

Method		AG News	Entailed polarity	Navigate	Winowhy	Avg.
Qwen2-1.5B	Base	0.691	<b>1.000</b>	0.173	0.389	0.563
	Eigen Pruning	0.720	0.919	0.290	0.415	0.586
	Rescaling	<b>0.796</b>	0.719	0.214	0.458	0.547
	TaRot	0.778	0.980	<b>0.515</b>	<b>0.547</b>	<b>0.705</b>
Phi-3-mini	Base	0.729	<b>1.000</b>	0.470	0.588	0.697
	Eigen Pruning	0.519	0.878	0.392	0.099	0.472
	Rescaling	0.739	0.921	0.273	<b>0.629</b>	0.641
	TaRot	<b>0.740</b>	<b>1.000</b>	<b>0.491</b>	0.600	<b>0.708</b>
Mistral-7B	Base	0.653	0.762	0.140	0.618	0.543
	Rescaling	0.437	<b>0.896</b>	<b>0.550</b>	0.683	<b>0.642</b>
	TaRot	<b>0.721</b>	0.823	0.216	<b>0.767</b>	0.632
Llama-3-8B	Base	0.662	0.980	0.155	0.568	0.591
	RED	0.688	<b>0.980</b>	<b>0.957</b>	0.236	0.715
	Rescaling	0.636	0.544	0.550	0.255	0.496
	TaRot	<b>0.718</b>	<b>1.000</b>	0.464	<b>0.701</b>	<b>0.721</b>
Qwen 2.5 14B	Base	0.753	0.763	0.424	0.723	0.666
	Rescaling	0.738	0.517	0.463	0.506	0.556
	TaRot	<b>0.754</b>	<b>0.826</b>	<b>0.480</b>	<b>0.732</b>	<b>0.698</b>
Qwen 2.5 32B	Base	0.808	0.901	0.717	<b>0.788</b>	0.803
	Rescaling	0.803	0.892	0.625	0.593	0.728
	TaRot	<b>0.824</b>	<b>0.927</b>	<b>0.734</b>	0.767	<b>0.813</b>

Table 1: **Overall performance in zero-shot regime.** Performance of methods with different LLMs in terms of F1 scores are presented across different tasks and on average. **Bold-faced** number denote the best method. For Mistral-7B, Llama-3-8B, Qwen 2.5 14B and Qwen 2.5 32B, Eigen Pruning resulted in OOM and RED codebase is only compatible with LLaMA architecture. Further details in Appendix A.5

Method		AG News	Entailed polarity	Navigate	Winowhy	Avg.
Qwen2-1.5B	Base	0.680	<b>0.902</b>	0.173	0.393	0.537
	Rescaling	0.662	0.765	0.314	<b>0.576</b>	0.579
	TaRot	<b>0.695</b>	<b>0.902</b>	<b>0.494</b>	0.544	<b>0.659</b>
Phi-3-mini	Base	0.745	0.974	0.440	0.604	0.691
	Rescaling	0.732	0.980	0.196	0.562	0.618
	TaRot	<b>0.764</b>	<b>0.991</b>	<b>0.494</b>	<b>0.647</b>	<b>0.724</b>
Mistral-7B	Base	0.691	0.921	<b>0.236</b>	<b>0.790</b>	0.660
	Rescaling	<b>0.746</b>	0.698	0.196	0.580	0.555
	TaRot	0.684	<b>0.960</b>	0.196	<b>0.790</b>	<b>0.658</b>
Llama-3-8B	Base	0.524	0.950	0.645	0.651	0.693
	Rescaling	0.444	0.702	0.196	0.577	0.480
	TaRot	<b>0.638</b>	<b>1.000</b>	<b>0.727</b>	<b>0.761</b>	<b>0.782</b>
Qwen 2.5 14B	Base	0.749	0.868	0.527	0.691	0.709
	Rescaling	0.739	0.807	0.362	0.422	0.583
	TaRot	<b>0.752</b>	<b>0.888</b>	<b>0.605</b>	<b>0.759</b>	<b>0.751</b>
Qwen 2.5 32B	Base	0.877	0.950	0.791	0.647	0.816
	Rescaling	0.844	0.941	0.715	0.674	0.793
	TaRot	<b>0.882</b>	<b>0.966</b>	<b>0.802</b>	<b>0.688</b>	<b>0.835</b>

Table 2: **Overall performance in few-shot regime.** Performance of methods with different LLMs in terms of F1 scores are presented across different tasks (and on average). **Bold-faced** numbers denote the best methods.

Method		Imdb	Toxicity
Qwen2-1.5B	Base	0.677	0.566
	Rescale	0.252	0.161
	TaRot	0.708	0.581
Phi-3-mini	Base	0.707	0.536
	Rescale	0.686	0.416
	TaRot	0.749	0.564
Llama-3-8B	Base	0.708	0.571
	Rescale	0.669	0.566
	TaRot	0.729	0.579

Table 3: **Performance comparison on NLG tasks.** Performance of Imdb review and toxicity task. The reported score are the average of the fluency and reward scores. A higher score indicates better performance on both NLG tasks.

instance, in zero-shot Entailed Polarity prediction, Rescaling significantly outperforms both the base model and TaRot on Mistral-7B (Table 1), but fails to scale in the few-shot setting (Table 2) and deteriorates performance across most other models. Two factors explain this: (1) As discussed in Section 3.2, attention head token associations exist in superposed states, making direct scaling or ablation unreliable; (2) large fluctuations introduced by Rescaling hinder optimization. While Rescaling requires fewer parameters— $H$  per layer vs.  $\frac{d}{2}$  in TaRot—the difficulty arises from the polysemantic nature of OV-circuits. In some cases, downscaling all associations in a head helps, likely

Method		AG News	Average
Qwen2-1.5B	SFT	0.603	0.362
	TaRot	0.655	0.447
Phi-3-mini	SFT	0.677	0.745
	TaRot	0.738	0.614
Llama-3-8B	SFT	0.693	0.661
	TaRot	0.744	0.520

Table 4: **Generalizability of TaRot.** Performance of supervised fine-tuning (SFT) and TaRot when trained on the AG News dataset and evaluated on both AG News and the average of two other tasks (Winowhy and Navigate).

due to non-interacting associations, but this varies unpredictably across tasks and models. In contrast, TaRot’s rotational alignment offers fine-grained control and robust, consistent performance. Future work can develop a formal theoretical framework to directly compare rotation-based (TaRot) and rescaling-based interventions, potentially by analyzing their effects on the residual stream. In case of NLG tasks, the combined score of the individual task specific reward model and fluency, is higher for both the tasks across the models. We believe that combining reward model and fluency scores provides a more comprehensive evaluation – the reward model captures task alignment, and fluency ensures the outputs remain coherent and natural. This combination better reflects overall performance (details in Appendix A.4). Table 3 presents TaRot’s results on IMDB and toxicity tasks, where it consistently outperforms both the base model and the rescaling approach. The reported scores reflect the combined metric, with higher values indicating better performance. On average, the performance of TaRot is improved on IMDB by 3.1% over the base model and 1.7% on toxicity tasks. However, on observing fluency and reward score separately, we see that solely in terms of reward values, Rescaling performs better than TaRot, and both interventions perform better than the original model. However, TaRot delivers more fluent response in terms of evaluation by GPT-4, pointing towards the more drastic edits of rescaling as compared to TaRot (see Appendix A.4 for complete results).

**Generalizability of TaRot.** TaRot applies fine grained intervention to the model attention heads, without altering the performance on remaining tasks. To show this, we perform supervised fine tuning(SFT), keeping the size of the train set similar to that of TaRot. We choose AG news as the train task as this is the only multi-class classification problem. Table 4 compares the performance of SFT

and TaRot trained on the AG News dataset across three different models: Qwen2-1.5B, Phi-3-mini, and Llama-3-8B. The results indicate that TaRot outperforms SFT on the trained tasks. We observe strong generalization in smaller models but weaker gains in larger ones (e.g., LLaMA-3-8B), due to two main factors: (1) TaRot removes spurious features but cannot inject new task-specific or syntactic knowledge; (2) the high-dimensional rotation space in larger models makes optimization harder. Moreover, since all tasks are classification-based, SFT, being explicitly task-driven, offers stronger supervision, reducing TaRot’s relative impact at scale. This explains the performance boost from SFT even on unseen tasks. Future work can focus on enhancing TaRot’s robustness for larger models.

**Ablation studies.** To assess the robustness of our approach, we conduct two ablation studies. **(1) Hyperparameter  $\mathbb{L}$ :** Using the Qwen 2.5-14B model, we apply the rotation transformation on varying layers of the model by changing the hyperparameter  $\mathbb{L}$ . Results show that applying Bayesian optimization to the initial layers yields the best performance with minimal parameters (see Appendix A.6.1 for details). **(2) System Prompt:** To guide the model towards accurate outputs, we use a fixed system prompt per task. We evaluate TaRot’s robustness to prompt variation on AG News using Qwen 2.5-14B with three different prompts. Despite fluctuations in base model performance, TaRot consistently outperforms it and reduces variance. Full results are in Appendix A.6.2.

## 6 Conclusion

In this work, we proposed TaRot, a novel, gradient-free, mechanistic intervention method for editing language models. TaRot builds on observations from implicit gradient descent bias of causal attention and applies parametrized rotation on the attention output to minimize the effects of undesired memorizations, doing away with effort-intensive localization steps and task-specificity of prior intervention techniques. Using Bayesian optimization of the rotational parameters, TaRot renders as data-efficient as in-context learning; yet, across a variety of tasks and language models of different sizes and families, robust improvement is observed. In a nutshell, TaRot can pave the path for general-purpose model editing methods in the future beyond supervised fine-tuning.



## 7 Limitations and Ethical Considerations

TaRot is designed to perform when the model has a generalization ability that is suppressed by noisy memorization. In that sense, it is limited by the boundaries of pretraining and cannot be used for domain adaptation. Fundamentally, it is not applicable to proprietary models. Finally, similar to any intervention technique, TaRot can be used in reverse to bypass alignment tuning and reinforce undesired behaviors. Future work can be to address the potential misuse of TaRot for bypassing alignment. One can develop detection mechanisms to identify whether TaRot or similar transformations have been applied to manipulate a model’s behavior. Incorporating regularization strategies that penalize rotations leading to toxic, biased, or otherwise misaligned generations would further ensure that the optimization process remains consistent with ethical AI principles. These directions can help mitigate potential misuse of TaRot and similar model editing techniques.

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Adithya Bhaskar, Dan Friedman, and Danqi Chen. 2024. [The heuristic core: Understanding subnetwork generalization in pretrained language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14351–14368, Bangkok, Thailand. Association for Computational Linguistics.

BIG-bench authors. 2023. [Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models](#). *Transactions on Machine Learning Research*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack

Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Subhabrata Dutta, Joykirat Singh, Soumen Chakrabarti, and Tanmoy Chakraborty. 2024. How to think step-by-step: A mechanistic understanding of chain-of-thought reasoning. *arXiv preprint arXiv:2402.18312*.

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1(1):12.

Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R. Costa-jussà. 2024. [A primer on the inner workings of transformer-based language models](#). *Preprint*, arXiv:2405.00208.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*.

Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. 2024. [Finding alignments between interpretable causal variables and distributed neural representations](#). In *Proceedings of the Third Conference on Causal Learning and Reasoning*, volume 236 of *Proceedings of Machine Learning Research*, pages 160–187. PMLR.

Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Ramaneswaran S, Deepali Aneja, Zeyu Jin, Ramani Duraiswami, and Dinesh Manocha. 2024. [A closer look at the limitations of instruction tuning](#). In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 15559–15589. PMLR.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *Preprint*, arXiv:2311.05232.

726	Samyak Jain, Robert Kirk, Ekdeep Singh Lubana,	2024. <a href="#">Interpreting key mechanisms of factual re-</a>	781
727	Robert P. Dick, Hidenori Tanaka, Edward Grefen-	<a href="#">call in transformer-based language models</a> . <i>Preprint</i> ,	782
728	stette, Tim Rocktäschel, and David Scott Krueger.	arXiv:2403.19521.	783
729	2024. <a href="#">Mechanistically analyzing the effects of fine-</a>		
730	<a href="#">tuning on procedurally defined tasks</a> . <i>Preprint</i> ,	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,	784
731	arXiv:2311.12786.	Dan Huang, Andrew Y. Ng, and Christopher Potts.	785
		2011. <a href="#">Learning word vectors for sentiment analysis</a> .	786
732	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	In <i>Proceedings of the 49th Annual Meeting of the</i>	787
733	sch, Chris Bamford, Devendra Singh Chaplot, Diego	<i>Association for Computational Linguistics: Human</i>	788
734	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	<i>Language Technologies</i> , pages 142–150, Portland,	789
735	laume Lample, Lucile Saulnier, et al. 2023. Mistral	Oregon, USA. Association for Computational Lin-	790
736	7b. <i>arXiv preprint arXiv:2310.06825</i> .	guistics.	791
737	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina	Kevin Meng, David Bau, Alex Andonian, and Yonatan	792
738	Toutanova. 2019. Bert: Pre-training of deep bidirec-	Belinkov. 2022. <a href="#">Locating and editing factual asso-</a>	793
739	tional transformers for language understanding. In	<a href="#">ciations in gpt</a> . In <i>Advances in Neural Information</i>	794
740	<i>Proceedings of naacL-HLT</i> , volume 1, page 2. Min-	<i>Processing Systems</i> , volume 35, pages 17359–17372.	795
741	neapolis, Minnesota.	Curran Associates, Inc.	796
742	Ananya Kumar, Aditi Raghunathan, Robbie Matthew	Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas	797
743	Jones, Tengyu Ma, and Percy Liang. 2022. <a href="#">Fine-</a>	Joseph, Nova DasSarma, Tom Henighan, Ben Mann,	798
744	<a href="#">tuning can distort pretrained features and underper-</a>	Amanda Askell, Yuntao Bai, Anna Chen, Tom Con-	799
745	<a href="#">form out-of-distribution</a> . In <i>International Conference</i>	erly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds,	800
746	<i>on Learning Representations</i> .	Danny Hernandez, Scott Johnston, Andy Jones,	801
		Jackson Kernion, Liane Lovitt, Kamal Ndousse,	802
747	Max Lamparth and Anka Reuel. 2024. <a href="#">Analyzing and</a>	Dario Amodei, Tom Brown, Jack Clark, Jared Ka-	803
748	<a href="#">editing inner mechanisms of backdoored language</a>	plan, Sam McCandlish, and Chris Olah. 2022. <a href="#">In-</a>	804
749	<a href="#">models</a> . In <i>Proceedings of the 2024 ACM Confer-</i>	<a href="#">context learning and induction heads</a> . <i>Preprint</i> ,	805
750	<i>ence on Fairness, Accountability, and Transparency</i> ,	arXiv:2209.11895.	806
751	FAccT '24, page 2362–2373, New York, NY, USA.		
752	Association for Computing Machinery.	Nikhil Prakash, Tamar Rott Shaham, Tal Haklay,	807
		Yonatan Belinkov, and David Bau. 2024. <a href="#">Fine-tuning</a>	808
753	Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S	<a href="#">enhances existing mechanisms: A case study on en-</a>	809
754	Schoenholz, Jeffrey Pennington, and Jascha Sohl-	<a href="#">tity tracking</a> . In <i>The Twelfth International Confer-</i>	810
755	Dickstein. 2017. Deep neural networks as gaussian	<i>ence on Learning Representations</i> .	811
756	processes. <i>arXiv preprint arXiv:1711.00165</i> .		
757	Chak Tou Leong, Yi Cheng, Jiashuo Wang, Jian Wang,	Ohad Rubin, Jonathan Herzig, and Jonathan Berant.	812
758	and Wenjie Li. 2023. <a href="#">Self-detoxifying language mod-</a>	2022. <a href="#">Learning to retrieve prompts for in-context</a>	813
759	<a href="#">els via toxification reversal</a> . In <i>Proceedings of the</i>	<a href="#">learning</a> . In <i>Proceedings of the 2022 Conference of</i>	814
760	<i>2023 Conference on Empirical Methods in Natural</i>	<i>the North American Chapter of the Association for</i>	815
761	<i>Language Processing</i> , pages 4433–4449, Singapore.	<i>Computational Linguistics: Human Language Tech-</i>	816
762	Association for Computational Linguistics.	<i>nologies</i> , pages 2655–2671, Seattle, United States.	817
		Association for Computational Linguistics.	818
763	Maximilian Li, Xander Davies, and Max Nadeau. 2024a.	V Sanh. 2019. Distilbert, a distilled version of bert:	819
764	<a href="#">Circuit breaking: Removing model behaviors with</a>	Smaller, faster, cheaper and lighter. <i>arXiv preprint</i>	820
765	<a href="#">targeted ablation</a> . <i>Preprint</i> , arXiv:2309.05973.	arXiv:1910.01108.	821
766	Yucen Lily Li, Tim G. J. Rudner, and Andrew Gordon	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane	822
767	Wilson. 2024b. <a href="#">A study of bayesian neural network</a>	Suhr. 2024. <a href="#">Quantifying language models' sensitiv-</a>	823
768	<a href="#">surrogates for bayesian optimization</a> . In <i>The Twelfth</i>	<a href="#">ity to spurious features in prompt design or: How i</a>	824
769	<i>International Conference on Learning Representa-</i>	<a href="#">learned to start worrying about prompt formatting</a> .	825
770	<i>tions</i> .	In <i>The Twelfth International Conference on Learning</i>	826
		<i>Representations</i> .	827
771	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan,	Tianhao Shen, Renren Jin, Yufei Huang, Chuang Liu,	828
772	Lawrence Carin, and Weizhu Chen. 2022. <a href="#">What</a>	Weilong Dong, Zishan Guo, Xinwei Wu, Yan Liu,	829
773	<a href="#">makes good in-context examples for GPT-3?</a> In	and Deyi Xiong. 2023. <a href="#">Large language model align-</a>	830
774	<i>Proceedings of Deep Learning Inside Out (DeeLIO</i>	<a href="#">ment: A survey</a> . <i>Preprint</i> , arXiv:2309.15025.	831
775	<i>2022): The 3rd Workshop on Knowledge Extrac-</i>		
776	<i>tion and Integration for Deep Learning Architectures</i> ,	Jasper Snoek, Hugo Larochelle, and Ryan P Adams.	832
777	pages 100–114, Dublin, Ireland and Online. Associa-	2012. Practical bayesian optimization of machine	833
778	tion for Computational Linguistics.	learning algorithms. <i>Advances in neural information</i>	834
		<i>processing systems</i> , 25.	835
779	Ang Lv, Yuhan Chen, Kaiyi Zhang, Yulong Wang,		
780	Lifeng Liu, Ji-Rong Wen, Jian Xie, and Rui Yan.		

836	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,	893
837	Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,	894
838	Adam R Brown, Adam Santoro, Aditya Gupta,	895
839	Adrià Garriga-Alonso, et al. 2022. Beyond the	896
840	imitation game: Quantifying and extrapolating the	897
841	capabilities of language models. <i>arXiv preprint</i>	898
842	<i>arXiv:2206.04615</i> .	
843	Eshaan Tanwar, Subhabrata Dutta, Manish Borthakur,	899
844	and Tanmoy Chakraborty. 2023. <a href="#">Multilingual LLMs</a>	900
845	<a href="#">are better cross-lingual in-context learners with align-</a>	901
846	<a href="#">ment</a> . In <i>Proceedings of the 61st Annual Meeting of</i>	902
847	<i>the Association for Computational Linguistics (Vol-</i>	903
848	<i>ume 1: Long Papers)</i> , pages 6292–6307, Toronto,	904
849	Canada. Association for Computational Linguistics.	905
850	Qwen Team. 2024. <a href="#">Qwen2.5: A party of foundation</a>	906
851	<a href="#">models</a> .	907
852	Tomás Vergara-Browne, Álvaro Soto, and Akiko	908
853	Aizawa. 2024. <a href="#">Eigenpruning: an interpretability-</a>	909
854	<a href="#">inspired peft method</a> . <i>Preprint</i> , arXiv:2404.03147.	910
855	Kevin Ro Wang, Alexandre Variengien, Arthur Conmy,	911
856	Buck Shlegeris, and Jacob Steinhardt. 2023. <a href="#">Inter-</a>	912
857	<a href="#">pretability in the wild: a circuit for indirect object</a>	913
858	<a href="#">identification in GPT-2 small</a> . In <i>The Eleventh Inter-</i>	914
859	<i>national Conference on Learning Representations</i> .	
860	Mengru Wang, Ningyu Zhang, Ziwun Xu, Zekun Xi,	915
861	Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi	916
862	Yang, Jindong Wang, and Huajun Chen. 2024a. <a href="#">Detoxifying large language models via knowledge</a>	917
863	<a href="#">editing</a> . In <i>Proceedings of the 62nd Annual Meeting</i>	918
864	<i>of the Association for Computational Linguistics (Vol-</i>	919
865	<i>ume 1: Long Papers)</i> , pages 3093–3118, Bangkok,	920
866	Thailand. Association for Computational Linguistics.	
867		
868	Zhichao Wang, Bin Bi, Shiva Kumar Pentiyala, Kiran	921
869	Ramnath, Sougata Chaudhuri, Shubham Mehrotra,	922
870	Zixu, Zhu, Xiang-Bo Mao, Sitaram Asur, Na, and	923
871	Cheng. 2024b. <a href="#">A comprehensive survey of llm align-</a>	
872	<a href="#">ment techniques: Rlhf, rlai, ppo, dpo and more</a> .	
873	<i>Preprint</i> , arXiv:2407.16216.	
874	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	924
875	Adams Wei Yu, Brian Lester, Nan Du, Andrew M.	925
876	Dai, and Quoc V Le. 2022. <a href="#">Finetuned language mod-</a>	
877	<a href="#">els are zero-shot learners</a> . In <i>International Confer-</i>	
878	<i>ence on Learning Representations</i> .	
879	Jiaxin Wen, Pei Ke, Hao Sun, Zhexin Zhang, Chengfei	926
880	Li, Jinfeng Bai, and Minlie Huang. 2023. <a href="#">Unveiling</a>	927
881	<a href="#">the implicit toxicity in large language models</a> . In <i>The</i>	928
882	<i>2023 Conference on Empirical Methods in Natural</i>	
883	<i>Language Processing</i> .	
884	Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li,	929
885	Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan	930
886	Zhang, Xiaoqing Zheng, and Xuanjing Huang. 2024.	931
887	<a href="#">Advancing parameter efficiency in fine-tuning via</a>	
888	<a href="#">representation editing</a> . <i>Preprint</i> , arXiv:2402.15179.	
889	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	932
890	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	933
891	Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2	934
892	technical report. <i>arXiv preprint arXiv:2407.10671</i> .	
	Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing	935
	Qu, Yong Jae Lee, and Yi Ma. 2024. <a href="#">Investigating the</a>	936
	<a href="#">catastrophic forgetting in multimodal large language</a>	937
	<a href="#">model fine-tuning</a> . In <i>Conference on Parsimony and</i>	938
	<i>Learning</i> , volume 234 of <i>Proceedings of Machine</i>	939
	<i>Learning Research</i> , pages 202–227. PMLR.	940
	Hongming Zhang, Xinran Zhao, and Yangqiu Song.	
	2020. <a href="#">WinoWhy: A deep diagnosis of essential</a>	
	<a href="#">commonsense knowledge for answering Winograd</a>	
	<a href="#">schema challenge</a> . In <i>Proceedings of the 58th An-</i>	
	<i>ual Meeting of the Association for Computational</i>	
	<i>Linguistics</i> , pages 5736–5745, Online. Association	
	for Computational Linguistics.	
	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang,	
	Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tian-	
	wei Zhang, Fei Wu, and Guoyin Wang. 2024. <a href="#">In-</a>	
	<a href="#">struction tuning for large language models: A survey</a> .	
	<i>Preprint</i> , arXiv:2308.10792.	
	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	
	Character-level convolutional networks for text classi-	
	fication. <i>Advances in neural information processing</i>	
	<i>systems</i> , 28.	
	<b>A Appendix</b>	
	<b>A.1 Task Details</b>	
	We experimented with five different classification	
	(i.e., single token generation) tasks and two NLG	
	tasks. Below are the details of the tasks with their	
	prompt templates used:	
	<b>AG News:</b> The goal of the task is to categories	
	new articles into one of the four predefined cate-	
	gories.	
	• World – News about global events, interna-	
	tional politics, and worldwide issues.	
	• Sports – News related to sporting events, ath-	
	letes, competitions, and sports industry devel-	
	opments.	
	• Business - News focusing on the economy,	
	financial markets, companies, and business	
	trends.	
	• Science & Technology – News about techno-	
	logical advancements, scientific discoveries,	
	and research.	
	<b>System prompt used for AG News task:</b> <i>You</i>	
	<i>are a news classification model. Your task is to</i>	
	<i>classify news articles into one of the following four</i>	
	<i>categories: World, Sports, Business, or Science.</i>	
	<i>You should respond with only the category name</i>	
	<i>and no other characters.</i>	



**Entailed polarity:** The Entailed Polarity task is a yes/no question-answering task (Srivastava et al., 2022). Given a fact and a question, the goal is to determine whether the fact entails a yes or no answer to the question. The task tests the model’s ability to infer whether the factual statement logically supports the answer in terms of polarity (positive or negative). Example:

- Fact: “Ed remembered to go.”
- Question: “Did Ed go?”
- Answer: “Yes”

**System prompt used for Entailed Polarity task:** *Follow the instructions below and answer with Yes / No.*

**Navigate:** The objective is to follow a set of directional or spatial instructions and determine if, after following those steps, the entity returns to the starting point. The answer is either True or False, depending on whether the instructions guide the entity back to where they started. Example:

- Instruction: “If you follow these instructions, do you return to the starting point?”
- Steps: “Always face forward.”, “Take 7 steps left.”, “Take 2 steps backward.”, “Take 7 steps backward.”, “Take 7 steps backward.”, “Take 3 steps forward.”
- Question: “Do you return to the starting point?”
- Answer: False

**System prompt used for the task:** *Answer the following question and output only True/False.*

**Winowhy:** This task (Srivastava et al., 2022) requires models to identify the correct reasons behind the answers to the Winograd Schema Challenges (Zhang et al., 2020).

This task is based on the original Winograd Schema Challenge (WSC) dataset and 4095 WinoWhy reasons (15 for each WSC question) that could justify the pronoun coreference choices in WSC. The model is presented with a passage that contains a pronoun and an explanation of which word or entity the pronoun refers to. The model’s job is to assess whether the explanation given is correct or incorrect based on the context of the passage.

- Text: “Fred is the only man alive who still remembers my father as an infant. When Fred first saw my father, he was twelve years old. The ‘he’ refers to Fred because, in his own words, he is ‘a very odd man’.”
- Question: “The above reasoning is:”
- Answer: “Incorrect”.

**System prompt used for Winowhy task:** *Follow the instructions and output Correct/Incorrect.*

**Imdb:** Tune model to generate positive movie reviews using a BERT (Kenton and Toutanova, 2019) sentiment classifier as a reward function. The reward model evaluates the sentiment of the generated reviews, and the goal is to maximize the likelihood of generating reviews classified as positive.

- Dataset Used: imdb (Maas et al., 2011)
- Reward Model: lwwerra/distilbert-imdb, a fine-tuned version of distilbert-base-uncased (Sanh, 2019) on the imdb dataset.

**Detoxify:** Involves reducing the toxicity of language model outputs. The toxicity evaluation is done using a classifier, such as facebook/roberta-hate-speech-dynabench-r4-target, which distinguishes between “neutral” and “toxic” text. The classifier provides feedback (reward or penalty) based on the toxicity of the model’s output, guiding the model to produce less toxic text. The dataset used is allenai/real-toxicity-prompts (Gehman et al., 2020).

## A.2 Experimental Setup Details

**Bayesian optimization.** We use I-BNN with 12 hidden layers, and LogExpectedImprovement as the acquisition function. We use a mixture of  $M$ -shots generation to avoid biasing the intervention, with  $M$  chosen randomly from 0 to 6. Each task was optimized for 150 iterations.

**Baselines.** We compare TaRot with four different baselines: (1) *Base model* denotes the pretrained LLM (zero-shot or few-shot) without any interventions. (2) *Eigen Pruning* (Vergara-Browne et al., 2024) removes singular values from weight matrices in an LLM to improve its performance in a particular task. (3) *RED (Representation Editing)* (Wu et al., 2024), which modifies the representations generated at some layers through the application of scaling and biasing operations. To



have a fair comparison, we also use a maximum of 20 prompts in its training phase. (4) *Rescaling* ablates attention heads by scaling their output in the unit interval instead of rotating their outputs; we use the same optimization technique to figure out the optimal scaling configuration.

**Evaluation metrics.** For Imdb positive review tasks, a sentiment analysis reward model, `lvwerra/distilbert-imdb`<sup>3</sup> is used. Roberta-hate-speech-dynabench-r4-target<sup>4</sup> is used for detoxification. For fluency GPT4 (Achiam et al., 2023) is used as an oracle to assign a value between 1 and 5, 1 being the least and 5 being the highest.

### A.3 Fluency

To evaluate the fluency of a given text, the following prompt was used with GPT4 (Achiam et al., 2023):

**System prompt used:** *Please rate the fluency of the following text on a scale of 1 to 5, where 1 is least fluent and 5 is most fluent: text. Provide only the number.*

where text is the output from the model.

### A.4 NLG tasks performance

Table 3 presents the performance of TaRot on IMDB sentiment classification and toxicity detection, where it consistently outperforms both the base model and rescaling-based methods. The table reports a combined score of fluency and reward model outputs, where a higher score indicates better performance for both tasks.

#### A.4.1 Combing score of reward mode and fluency

Evaluation Setup of the two NLG tasks and fluency is described below:

- **Fluency:** Assessed using GPT-4, which assigns a score from 1 to 5 (where 1 = least fluent and 5 = most fluent).
- **IMDB Sentiment Reward Model:** We use `lvwerra/distilbert-imdb`, where higher scores indicate better sentiment classification.
- **Toxicity Reward Model:** We use `Roberta-hate-speech-dynabench-r4-target`, where higher scores indicate higher toxicity.

<sup>3</sup><https://huggingface.co/lvwerra/distilbert-imdb>

<sup>4</sup><https://huggingface.co/facebook/roberta-hate-speech-dynabench-r4-target>

Method	Imdb		Detoxify	
	Reward	Fluency	Reward	Fluency
Qwen2-1.5B				
Base	-0.80	—	4.50	—
Rescaling	<b>0.72</b>	1.25	<b>2.29</b>	1.26
TaRot	-0.25	<b>2.24</b>	4.01	<b>4.56</b>
Mistral-7B				
Base	-0.05	—	4.31	—
Rescaling	<b>0.19</b>	2.12	<b>3.18</b>	4.12
TaRot	0.16	<b>2.5</b>	4.01	<b>4.30</b>
Llama-3-8B				
Base	-0.31	—	4.05	—
Rescaling	<b>0.28</b>	<b>2.56</b>	<b>3.18</b>	<b>4.76</b>
TaRot	.002	2.38	3.90	4.24

Table 5: Performance comparison on NLG tasks.

The scoring methodology of the NLG tasks and fluency combined is described below:

- **IMDB + Fluency:** Both scores were normalized to [0,1] and summed to obtain the final score.
- **Toxicity + Fluency:** The toxicity score was normalized and inverted (so that lower toxicity results in a higher score), then combined with fluency.

Thus, in both cases, a higher final score reflects improved overall performance (i.e., better fluency and alignment with task objectives). The complete breakdown of the toxicity and fluency of the NLG tasks is shown in Table 5.

### A.5 Additional baselines

Eigenpruning and RED were used as baselines in our study. Below, we outline the key challenges that prevented us from incorporating additional baselines:

#### Eigenpruning:

- Eigenpruning requires fine-tuning the model before identifying circuits, which is a computationally expensive process.
- Given our resource constraints, we were unable to perform the necessary fine-tuning steps required for a fair comparison.

#### RED:

- The methodology and code for RED have only been demonstrated on GPT-2 and LLaMA models.
- The publicly available codebase lacks implementation details for extending RED to other model architectures.

Method	Layer	AG News	Navigate	Winowhy
Base	NA	0.753	0.424	0.723
TaRot	0 - 6	0.752	0.45	0.723
TaRot	0 - 12	<b>0.757</b>	0.432	0.715
TaRot	0 - 24	0.754	<b>0.480</b>	<b>0.732</b>
TaRot	0 - 32	0.743	0.439	0.712

Table 6: Zero Shot Performance with TaRot applied on different layers.

Method	Layer	AG News	Navigate	Winowhy
Base	NA	0.749	0.527	0.691
TaRot	0 - 6	0.755	0.549	0.733
TaRot	0 - 12	0.75	0.556	0.73
TaRot	0 - 24	0.752	<b>0.605</b>	<b>0.759</b>
TaRot	0 - 36	<b>0.757</b>	0.601	0.735

Table 7: Few Shot Performance with TaRot applied on different layers.

As a result, we were only able to run RED on the LLaMA model for comparison.

## A.6 Ablation Study

### A.6.1 Hyperparameter $\mathbb{L}$

Previous studies indicate that token associations related to pretrained knowledge primarily reside in the first half of the model. Based on this insight, we applied the rotation transformation only to the first half of the model. However, we acknowledge that an ablation study on this hyperparameter is necessary to fully assess the robustness of our approach.

To address this, we conducted an ablation study on the Qwen 2.5-14B model, which has 48 layers. We tested different layer ranges for applying the rotation matrix: 0-6, 0-12, 0-24, and 0-36. The tables below report zero-shot and few-shot F1 scores across three tasks: AG News, Navigate, and Winowhy. Table 6, 7 shows the performance zero shot and few shot performance respectively.

### Key Observations

- For Navigate and Winowhy, the best performance was achieved when applying TaRot to the first 24 layers (0-24).
- Ideally, 0-32 layers should also perform well, but the increased parameter space dimensionality makes it harder for Bayesian optimization to converge effectively.
- For task AG news we see comparable performance of TaRot optimized on the first half of the model with the best performing setting.

Method	Prompt 1	Prompt 2	Prompt 3
Base (Zero Shot)	0.724	0.815	0.73
TaRot (Zero Shot)	<b>0.772</b>	<b>0.832</b>	<b>0.768</b>
Base (Few Shot)	0.67	0.792	0.697
TaRot (Few Shot)	<b>0.7</b>	0.792	<b>0.715</b>

Table 8: Ablation of different prompt used for Qwen/Qwen2-1.5B-Instruct on AG News Tasks.

Therefore we see that applying optimization on the first half provides with us with the best performance.

### A.6.2 System Prompt

We tested the model with three semantically equivalent but syntactically different prompts. TaRot was optimized on each of these prompts separately to evaluate its effectiveness in mitigating performance fluctuations.

#### Prompt 1:

System prompt: You specialize in classifying news articles into distinct categories. Given a news article, determine whether it belongs to World, Sports, Business, or Science. Only provide the category name as a response.

Question: News Content: <review>

Query: What is the most suitable category for this news piece?

#### Prompt 2:

System prompt: You are an expert in news topic classification. Your role is to analyze articles and assign them to one of these four categories: Sports, World, Business, or Science. Do not add extra text—respond with just the category name.

Question: Text: <review>

Question: Which of the four categories (World, Sports, Business, or Science) does this article belong to?

#### Prompt 3:

System prompt: Your function is to categorize news articles into one of four groups: World, Sports, Business, or Science. Given a news article, determine its category and respond using only the category name.

Question: News Article: <review>

Task: Identify the correct category for this article.

Table 8 shows the performance of TaRot in zero and few shot settings compared with the base model. The model used is Qwen/Qwen2-1.5B-Instruct and the dataset is Ag News. The results demonstrate that the base model exhibits considerable fluctuation across different prompts, indicating a high sensitivity to prompt phrasing. In

contrast, TaRot consistently outperforms the base model across all prompt settings, showcasing its reliability and robustness. This consistency highlights TaRot’s ability to generalize better and remain stable despite variations in input structure. Notably, while the base model suffers a significant absolute drop of 0.145 in few-shot performance between prompts (from 0.815 to 0.67), TaRot substantially minimizes such performance degradation. This suggests that TaRot enhances the model’s resilience to prompt perturbations, reducing brittleness and improving reliability. Furthermore, by demonstrating improved performance across varied prompt templates, the experiment effectively addresses the reviewer’s concern—confirming that TaRot’s improvements are not limited to a single prompt instance but instead generalize across different prompt structures.

## A.7 System Prompt Used

For each system custom system prompts were used to help guide the model to output the final answer directly. We ensured that system prompts were only used when necessary—for instance, they were not applied in tasks like Entailed Polarity, where the model naturally follows the desired output structure, i.e output the final answer directly. The system prompt used for each tasks are as follows:

- **Task Navigate:** Prompt: “If you follow these instructions, do you return to the starting point?”
- **Task Entailed Polarity:** Prompt: “Given a fact, answer the following question with a yes or a no.”
- **Task Winowhy:** Prompt: “Please answer the following questions about which words certain pronouns refer to.”
- **Task AG News:** Prompt: “News Article: review: Question: What category does this news article belong to?”