

# TASK VECTORS, LEARNED NOT EXTRACTED: PERFORMANCE GAINS AND MECHANISTIC INSIGHTS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) can perform new tasks from in-context demonstrations, a phenomenon known as in-context learning (ICL). Recent work suggests that these demonstrations are compressed into task vectors (TVs), compact task representations that LLMs exploit for predictions. However, prior studies typically extract TVs from model outputs or hidden states using cumbersome and opaque methods, and they rarely elucidate the mechanisms by which TVs influence computation. In this work, we address both limitations. First, we propose directly training **Learned Task Vectors (LTVs)**, which surpass extracted TVs in accuracy and exhibit superior flexibility—acting effectively at arbitrary layers, positions, and even with ICL prompts. Second, through systematic analysis, we investigate the mechanistic role of TVs, showing that at the low level they steer predictions primarily through attention-head OV circuits, with a small subset of “key heads” most decisive. At a higher level, we find that despite Transformer nonlinearities, TV propagation is largely linear: early TVs are rotated toward task-relevant subspaces to improve logits of relevant labels, while later TVs are predominantly scaled in magnitude. Taken together, LTVs not only provide a practical approach for obtaining effective TVs but also offer a principled lens into the mechanistic foundations of ICL<sup>1</sup>.

## 1 INTRODUCTION

Large Language Models (LLMs) possess the remarkable capability of performing novel natural language tasks by learning from demonstrations included in the input without training, a phenomenon referred to as **In-context Learning (ICL)** (Brown et al., 2020; Radford et al., 2019). ICL has revolutionized natural language processing through its extensive empirical success in enabling swift and efficient adaptation of models to downstream tasks (Dong et al., 2024; Liu et al., 2021).

Since its effectiveness is difficult to reconcile with the traditional framework of machine learning centered on model training (Ren et al., 2024), investigating the internal mechanisms of LLMs that enable ICL has attracted substantial attention. Among these efforts, one prominent line of research shows that LLMs leverage demonstrations by summarizing them into **task vectors (TVs)**—succinct vector representations of the task exemplified by the demonstrations (Hendel et al., 2023). These TVs can be injected (added) into the hidden states of zero-shot prompts without demonstrations to achieve ICL-level performance. Subsequent work has primarily proceeded in three directions: **1)** studying where (e.g., from LLM hidden states (Hendel et al., 2023), attention head outputs (Todd et al., 2024; Yin & Steinhardt, 2025), or MLP outputs (Merullo et al., 2024) at different layers) and how (e.g., PCA-based approaches (Liu et al., 2024) or complex optimization methods (Li et al., 2024a; Cai et al., 2025)) to extract and construct TVs, with the practical goal of boosting performance through injection; **2)** investigating how the ability of LLMs to form TVs gradually emerges during pretraining, typically using small trained-from-scratch models and artificial tasks such as regression (Han et al., 2025; Yang et al., 2025b); and **3)** demonstrating that TVs naturally arise from the LLM architecture itself, and providing theoretical guarantees for their emergence (Bu et al., 2025; Dong et al., 2025).

Despite important contributions, prior studies face key limitations. First, existing approaches often depend on opaque and complex filtering or optimization to construct TVs from model representations, making them inefficient and reliant on the model’s representational quality. This dependence

<sup>1</sup>The source code will be released upon acceptance of this paper

can produce suboptimal TVs and mischaracterize their true effect, while the opaque construction procedures obscure an understanding of TV’s mechanism. Indeed, most works stop at showing that injected TVs improve performance but leave unanswered the central question of **how LLMs leverage TVs to make correct predictions**. This gap spans both the **low-level interactions**, referring to the microscopic localization of model components that interact with injected TVs to express their effects during forward computation, and the **high-level channels**, referring to the macroscopic mechanisms by which TVs ultimately steer outputs toward correct predictions. The lack of explanation reduces the model’s deployment of TVs to an uninterpretable black-box function (Merullo et al., 2024).

In this work, we address the first shortcoming by proposing to **directly train Learned Task Vectors (LTVs)** by adding a vector to a specific layer’s hidden states and optimizing it through gradient descent (Figure 1 (A)), which finds the optimal TV unconstrained by the quality of model’s representations. LTVs not only outperform constructed ones across classification and generation tasks but also demonstrate greater flexibility and scalability than extracted ones. Furthermore, through analysis of interactions between TVs and model components, we decode the **low-level** mechanisms by which LLMs interact with TVs: injected TVs are primarily utilized through attention-head OV circuits (Figure 1 (B)). We also characterize which attention heads are most decisive in leveraging the injected TVs, focusing on their attention and distribution patterns. Regarding the **high-level** influence channels of TVs, we show that despite the abundance of nonlinearities in Transformer layers, the propagation of injected TVs through subsequent layers is largely linear, involving a rotation that aligns TVs to the subspace spanned by task-related tokens and a scaling that adjusts their magnitude (Figure 1 (C)). We further observe a distinct pattern: the rotation effect attenuates as the injection layer index increases, while the scaling effect becomes the dominant factor translating TVs into output changes. In summary, our work introduces an efficient method to obtain effective TVs and provides a comprehensive exposition of the mechanisms underlying TVs’ effectiveness.

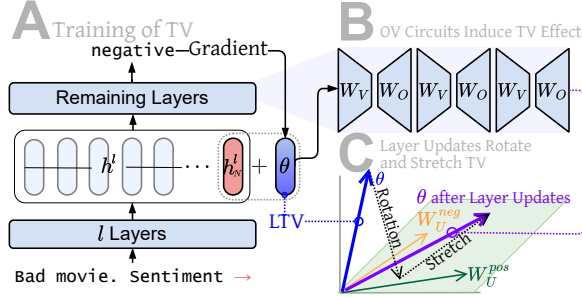


Figure 1: (A) We directly train Learned Task Vectors (LTVs) to be injected, which influence model outputs through later layers updates. (B) In the **low-level** interactions between TVs and later layers, the OV circuits of attention heads are the crucial components interacting with TVs to induce their effects. (C) On a **high level**, subsequent layer updates act on TVs as a largely linear transformation of rotation and stretch, with the rotation aligning TVs with the relevant task subspace to promote prediction of task-related tokens.

We also characterize which attention heads are most decisive in leveraging the injected TVs, focusing on their attention and distribution patterns. Regarding the **high-level** influence channels of TVs, we show that despite the abundance of nonlinearities in Transformer layers, the propagation of injected TVs through subsequent layers is largely linear, involving a rotation that aligns TVs to the subspace spanned by task-related tokens and a scaling that adjusts their magnitude (Figure 1 (C)). We further observe a distinct pattern: the rotation effect attenuates as the injection layer index increases, while the scaling effect becomes the dominant factor translating TVs into output changes. In summary, our work introduces an efficient method to obtain effective TVs and provides a comprehensive exposition of the mechanisms underlying TVs’ effectiveness.

## 2 RELATED WORKS

**Task Vector and ICL** The hypothesis that TVs form the mechanistic basis of ICL was first proposed by Hendel et al. (2023), who patched ICL hidden states into zero-shot prompts at certain layers, achieved high accuracy, and argued that in-context demonstrations are compressed into TVs applied during later updates. Follow-up studies (Todd et al., 2024; Li et al., 2024a; Kahardipraja et al., 2025; Liu et al., 2024) extended this idea by extracting TVs from specific components (e.g., attention heads, MLP) and injecting them. The universality of TVs has been validated across model scales (small trained-from-scratch vs. large open-source) and task types (mathematical vs. natural language) (Han et al., 2025; Yang et al., 2025b; Jiang et al., 2025a). Yet, little is known about how TVs enhance performance after injection, or how they interact with later components to shape outputs.

**Mechanisms of Task Vectors** Current explanations of TV effectiveness remain preliminary, more sketches than systematic analyses. For instance, Hendel et al. (2023) observed that TV injection is more effective in earlier than later layers. Todd et al. (2024) reported that TVs exhibit word2vec-style arithmetic (Mikolov et al., 2013), with Bu et al. (2025) giving a theoretical account of this property. Furthermore, Han et al. (2025) and Jiang et al. (2025b) found that TV effectiveness depends on how well hidden states of a task’s prompts can be separated from others in the LLM representation space.

**LLM Steering** The success of TV injection in restoring ICL performance parallels recent advances in LLM steering (Zhan et al., 2025; Li et al., 2024b; Panickssery et al., 2024), where vectors are added

to hidden states to mitigate undesirable model behaviors (Lee et al., 2024; Bayat et al., 2025). Prior work also explored training steering vectors directly (Cao et al., 2024; Dunefsky & Cohan, 2025), motivating our strategy of training TVs rather than relying on complex selection or construction.

### 3 METHODOLOGY

**Transformer hidden states and ICL** According to the autoregressive structure of Transformer LLMs with residual connections, a **zero-shot input query**  $x_q$  of  $N$  tokens (e.g., “I like this movie. Sentiment:”) is sequentially embedded and updated across  $L$  layers into  $N$   $d$ -dimensional hidden states. At each layer, the hidden state of token  $i$  at layer  $l$  is updated as  $\mathbf{h}_i^l = \mathbf{h}_i^{l-1} + \sum_{k=1}^K \mathbf{a}_{i,k}^l + \mathbf{m}_i^l$ , where  $\mathbf{a}_{i,k}^l$  is the output of the  $k$ -th attention head (head  $(l, k)$ ), and  $\mathbf{m}_i^l$  is the MLP output. Concretely,  $\mathbf{a}_{i,k}^l$  depends on the previous layer’s hidden states of the first  $i$  tokens  $[\mathbf{h}_j^{l-1}]_{j=1}^i$  through:

$$\mathbf{a}_{i,k}^l = \sum_{j=1}^i c_{j,i}^{l,k} \mathbf{W}_{O,k}^{l,\top} \mathbf{W}_{V,k}^l \mathbf{h}_j^{l-1}, \quad (1)$$

where  $\mathbf{W}_{V,k}^l$  and  $\mathbf{W}_{O,k}^l \in \mathbb{R}^{d_h \times d}$  are the value embedding and output projection matrices of head  $(l, k)$  respectively, jointly referred to as the **OV circuit**, with  $d_h$  being the head dimension.  $c_{j,i}^{l,k}$  denotes the attention weight from token  $i$  to  $j$  of head  $(l, k)$ . Consequently, the  $L$  layer updates can be viewed as a sequence of additive effects, with the last token hidden state at the final layer formed as:

$$\mathbf{h}_N^L = \mathbf{h}_N^0 + \sum_{l=1}^L \left( \sum_{k=1}^K \mathbf{a}_{N,k}^l + \mathbf{m}_N^l \right), \quad (2)$$

which is then multiplied by the unembedding matrix  $\mathbf{W}_U \in \mathbb{R}^{|V| \times d}$  to produce the output logits and final prediction. An **ICL** prompt is formed by prepending  $n$  demonstration-label pairs to the query, yielding an input sequence  $\mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_m, \mathbf{y}_m, \mathbf{x}_q$  (e.g., “I hate this movie. Sentiment: negative. This movie is great. Sentiment: positive . . . I like this movie. Sentiment:”). Because Transformer updates depend on hidden states from earlier tokens and layers, the hidden states as well as attention and MLP outputs change throughout, producing a different  $\mathbf{h}_{N,\text{ICL}}^L$  and ultimately a different prediction.

**Task Vector as a mechanistic explanation for ICL** Existing TV studies provide a functional characterization of the mechanism enabling LLMs to leverage demonstrations as  $f(\mathbf{x}_q; \boldsymbol{\theta})$ , i.e., LLMs make predictions based on the query together with a vector  $\boldsymbol{\theta}$  that represents the query-label mappings (Hendel et al., 2023; Merullo et al., 2024). These studies propose that  $\boldsymbol{\theta}$  is formed in early layers and assists LLM predictions as later layers execute  $f(\mathbf{x}_q; \boldsymbol{\theta})$ . Accordingly, they seek to extract the TV  $\boldsymbol{\theta}$  from the ICL hidden state stream and add it to the last token hidden state of  $\mathbf{x}_q$  at layer  $l$ , i.e.,  $\mathbf{h}_N^l + \boldsymbol{\theta}$ . The resulting hidden state is then propagated through subsequent layers, and the intervention is evaluated based on whether it achieves few-shot-level prediction accuracy for zero-shot queries. Two major methods of extracting  $\boldsymbol{\theta}$  have been proposed, which we treat as baselines in this work<sup>2</sup>:

1. **Vanilla TV** (Hendel et al., 2023), defined as  $\boldsymbol{\theta} = \mathbf{h}_{N,\text{ICL}}^l - \mathbf{h}_N^l$ , where  $\mathbf{h}_{N,\text{ICL}}^l$  is the layer- $l$  last token hidden state of an ICL prompt formed with a query different from  $\mathbf{x}_q$  which produces  $\mathbf{h}_N^l$ .
2. **Function Vector (FV)** (Todd et al., 2024; Li et al., 2024a; Yin & Steinhardt, 2025), defined as  $\boldsymbol{\theta} = \sum_{(l,k) \in \mathbb{I}} \mathbf{a}_{N,k,\text{ICL}}^l$ , where  $\mathbf{a}_{N,k,\text{ICL}}^l$  is the attention head output to the last token hidden state given ICL prompts, and  $\mathbb{I}$  is an index set of selected attention heads.

Both methods have drawbacks. Vanilla TV injection yields lower accuracy and is highly sensitive to the choice of injection layer  $l$ . FV depends on selecting an proper head set  $\mathbb{I}$ , typically determined by ablating heads one by one to measure their impact on output probability, which is suboptimal as it neglects intercorrelations among ablations. Moreover, both methods critically depend on the quality of model’s ICL representations (we use 8-shot ICL prompts to obtain hidden states for the two methods and to evaluate ICL performance). As a remedy, we propose directly training LTVs. **Training LTV** Instead of distilling from ICL hidden states, we train the LTV  $\boldsymbol{\theta}$  to minimize:

$$-\log p(\mathbf{y}_q | \mathbf{x}_q, \boldsymbol{\theta}, \mathbb{I}, \mathbb{P}), \quad (3)$$

<sup>2</sup>See Appendix D.4 for comparison with more baselines and the reason for omitting them from the main text.

where  $y_q$  is the correct label for the zero-shot query  $x_q$ ,  $\mathbb{L}$  denotes the set of layers and  $\mathbb{P}$  the token positions of hidden states where  $\theta$  is injected. This approach eliminates the need to manipulate ICL hidden states and uncovers the most effective TV, unconstrained by representation or demonstration quality crucial for traditional TV extraction. Moreover, we do not restrict  $\mathbb{P}$  to the final position or  $\mathbb{L}$  to a single layer as in the baselines. In general, we add  $|\mathbb{L}| \times |\mathbb{P}|$  different  $\theta$ s to the hidden states indexed by  $\mathbb{L}$  and  $\mathbb{P}$ . This design allows us to explore flexibility and scalability of our approach and to test the proposition from prior works that a single TV can encapsulate the full functionality of ICL, as discussed in Subsection 4.1. In the special case of  $\mathbb{L} = \{l\}$  and  $\mathbb{P} = \{-1\}$ , we add one  $\theta$  to  $h_N^l$  following baseline practice. During the training, for multi-token labels, we average log probabilities across tokens.  $\theta$  is optimized using AdamW (Loshchilov & Hutter, 2017) with learning rate = 0.001 and weight decay = 0.01. Details of the training procedure are provided in Appendix B.

## 4 EXPERIMENTS

**Models** We use the following models: Llama3-8B, Llama3.1-8B, Llama3.2-3B, Llama3-70B (Grattafiori et al., 2024), Llama2-7B, Llama2-13B (Touvron et al., 2023), Qwen2.5-32B (Yang et al., 2024), Yi-34B (01. AI et al., 2024). In the main text, results are reported on Llama3.1-8B.

**Datasets** We adopt three datasets from prior TV research (Todd et al., 2024): **1) Capital:** given a country name, output its capital city; **2) Capitalize:** given a word, output its capitalized first letter; **3) Antonym:** given a word, output its antonym. To evaluate TVs on more natural datasets with richer input-output mappings, we additionally consider four classification tasks: SST-2 (Socher et al., 2013), TREC (Li & Roth, 2002), SNLI (MacCartney & Manning, 2008), and RTE (Dagan et al., 2005). We report the prediction accuracy achieved by ICL and the different TV methods across the seven datasets. To test the ability of TVs to elicit LLM behaviors in more complex task settings, we also include the Myopic dataset (Panickssery et al., 2024), a generation task described in Subsection 4.1. We further include three more datasets specifically for investigating the compositionality and generalizability of LTV, also described in Subsection 4.1. See Appendix C for additional details on model implementation, datasets, and ICL setup.

### 4.1 SUPERIOR PERFORMANCE OF LEARNED TASK VECTORS

#### Consistent performance superiority of LTV

Following Hendel et al. (2023) and Todd et al. (2024), we first inject the TVs at one layer at a time, iterating over all layers of Llama3.1-8B, and report the average performance across datasets in Figure 2. The results show that our LTV not only consistently outperforms the baseline methods at all layers, but also matches or even surpasses ICL performance—particularly when injected at early layers of both models. The high accuracy achieved by the LTV also makes it a viable parameter-efficient finetuning (PEFT) method (Wu et al., 2024; Subramani et al., 2022; Turner et al., 2024), since it involves optimizing exactly  $d$  parameters, which is lower than most existing PEFT strategies. To demonstrate the potential of LTV as a PEFT method, we compare it against two widely used baselines—Prefix Tuning (Li & Liang, 2021) and LoRA (Hu et al., 2021)—under a comparable parameter budget on SST-2. Specifically, we apply prefix tuning to the key and value projections of all heads at the first layer of Llama3.1-8B with prefix length 2, introducing exactly  $d$  trainable parameters. Likewise, we apply LoRA to the output projection at the first layer with rank  $r = 1$ , yielding  $2d$  parameters. We then compare these baselines with injecting a single layer update at the last-token position of layer 0 ( $d$  parameters). Beyond accuracy, we track mean latency and FLOPs per training or inference sample, and peak memory per training or inference epoch. As summarized in Table 1, LTV achieves the strongest performance and best training-time efficiency, and its inference-time FLOPs and memory cost are only marginally higher than prefix tuning, which demonstrates its potential as a competitive PEFT method.

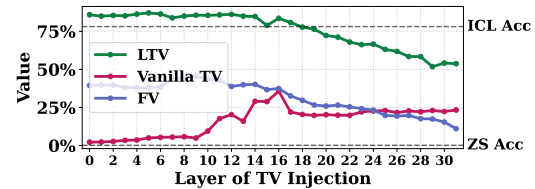


Figure 2: Dataset-average accuracy of injecting the Vanilla TV, FV, and LTV into the last-token hidden states by iterating over all layers and injecting into one layer at a time, along with ICL and zero-shot (ZS) accuracy levels. Our LTV consistently outperforms the Vanilla TV and FV across all layers, with the performance gap particularly prominent in late layers. See Appendix D.1 for other models.



Table 1: Comparing LTV against PEFT methods in terms of performance and efficiency

Method	Acc. $\uparrow$	Param Cnt. $\downarrow$	Training			Inference		
			Mean Lat. (Sec) $\downarrow$	FLOPs (GB) $\downarrow$	Peak Mem. (GB) $\downarrow$	Mean Lat. (Sec) $\downarrow$	FLOPs (GB) $\downarrow$	Peak Mem. (GB) $\downarrow$
Prefix Tuning	85.67%	$d$	0.050	533.15	43.65	0.026	<b>361.51</b>	<b>16.31</b>
LoRA	91.63%	$2d$	0.053	526.98	43.65	0.032	361.52	16.37
LTV (Ours)	<b>92.89%</b>	$d$	<b>0.049</b>	<b>503.87</b>	<b>43.56</b>	<b>0.024</b>	361.52	16.36

Table 2: LTV outperforms Vanilla TV and FV not only in the baseline case but also across diverse configurations with varied positions, layers, and prompt formats. See Appendix D.2 for other models.

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{16\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
Vanilla TV	37.80%	2.16%	19.18%	17.97%	18.15%	56.12%
FV	37.30%	2.68%	6.05%	31.88%	0.38%	74.78%
LTV (Ours)	83.49% <sub><math>\uparrow 45.69\%</math></sub>	78.39% <sub><math>\uparrow 75.71\%</math></sub>	82.44% <sub><math>\uparrow 63.26\%</math></sub>	86.43% <sub><math>\uparrow 54.55\%</math></sub>	51.39% <sub><math>\uparrow 33.24\%</math></sub>	84.61% <sub><math>\uparrow 9.83\%</math></sub>

**Accuracy of late-layer injection** Another notable trait of the LTV is that it still achieves nontrivial performance when trained and injected at late layers, despite an overall decreasing trend with depth. This contrasts with Vanilla TV and FV, which show severely degraded accuracy beyond a certain depth as reported in prior work (Li et al., 2024a; Todd et al., 2024). Our results therefore challenge the idea that a critical depth threshold exists beyond which layers cannot utilize the injected TV. We further analyze the mechanism enabling LTVs at different depths to take effect in Subsection 4.3.

**Flexibility and scalability of the LTVs** Existing TV studies typically inject solely into the **last token** hidden state ( $\mathbb{P} = \{-1\}$ ) at **one specific layer** ( $\mathbb{L} = \{l\}$ ) of the **zero-shot prompt**. We go beyond this baseline to examine the adaptability of our LTV to more diverse configurations. We set  $l$  to the middle layer of the model (i.e., 16 for the 32-layer Llama3.1-8B) as the baseline, and then consider the following variants. **1)** Keep  $l$  fixed but inject (and train) at a different position  $\mathbb{P} = \{4\}$ , i.e., add the TV to the fourth token hidden state<sup>3</sup>. **2)** Inject at multiple positions:  $\mathbb{P} = \{-5, -4, -3, -2, -1\}$ . **3)** Keep  $\mathbb{P} = \{-1\}$  but inject at every four (ablation studies in Appendix D.3) layers, i.e.,  $\mathbb{L} = \{0, 4, \dots, 28, 32\}$  for Llama2-13B. **4)** Set  $\mathbb{P} = \{-5, -4, -3, -2, -1\}$  and  $\mathbb{L} = \{0, 4, \dots, 28, 32\}$  simultaneously. **5)** Keep  $\mathbb{P}$  and  $\mathbb{L}$  fixed but change the zero-shot prompt to an 8-shot ICL prompt. We compare our LTV to Vanilla TV and FV in all five settings<sup>4</sup>, with implementation details for FV in Appendix H.

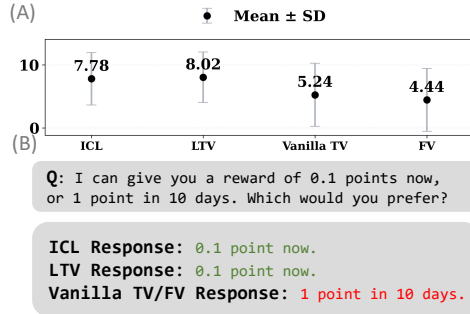


Figure 3: (A) Mean and standard deviation of ratings for responses generated with ICL, LTV, FV, and Vanilla TV. (B) An example question and responses across settings.

**Adaptability of LTVs to complex task settings** The tasks above have single-token labels and unique correct answers (e.g., The capital of China is  $\rightarrow$  Beijing). To evaluate generalizability to a more complex generation task with multi-token responses—where the goal is to elicit a behavioral mode rather than a single answer—we adopt the **Myopic** dataset from the LLM steering literature (Panickssery et al., 2024; Bayat et al., 2025). Each entry presents a question with two choices (see

<sup>3</sup>Prompts with fewer than 4 tokens are skipped in the accuracy calculation.

<sup>4</sup>In the baseline case, the FV method adds the sum of head outputs at the last position to the final token’s hidden state. For varied  $\mathbb{P}$ , we add summed outputs at each position in  $\mathbb{P}$  to the corresponding hidden state. For multiple layers, we replicate the FVs  $|\mathbb{L}|$  times and inject a copy at each layer. For the Vanilla TV, we patch hidden states at positions  $\mathbb{P}$  and  $\mathbb{L}$  of an ICL prompt with a different query into those of  $x_q$ .

Figure 3 (B)), one myopic and the other favoring long-term rewards. We compare the generated answers with LTV (injected at the middle layer) to the two baselines by asking an LLM to rate them on a 10-point scale (details in Appendix I) based on how well they reflect the myopic choice. The statistics in Figure 3 show that the LTV not only surpasses the baselines but also exceeds ICL performance—something existing TV methods distilled from ICL representations struggle to achieve. These results provide clear evidence of the potential of LTVs in complex generation settings (see Appendix D.5 for other models).

**Cross-task TV similarity and generalizability** As a first step toward understanding how TVs capture task idiosyncrasies, we compute cosine similarities among TVs trained for different tasks (and across repeats of the same task). For each of the seven tasks, we train a middle-layer LTV five times and compute cosine similarities among the resulting  $7 \times 5 = 35$  vectors. The outcome in Figure 4 (other models in Appendix D.6) shows that LTVs internalize effective and consistent task representations, with clear intra-task alignment and inter-task separation. The only notable exception is the moderate alignment between LTVs trained for SNLI and RTE, which share the labels  $\{\text{true}, \text{false}\}$ . This suggests that the unembedding directions of task labels critically determine the orientation of TVs: LTVs that ultimately promote alignment of hidden states to the same unembedding vectors naturally exhibit high similarity. To verify this conclusion, we apply the middle-layer TV trained on SNLI to other datasets and record the induced accuracy. The results in Table 3 where the TV only leads to nontrivial accuracy on RTE, demonstrating that the generalizability of TV across tasks critically depends on the task label space. Additional results of applying the TV of Capital dataset (which does not share the label of other datasets, as opposed to the SNLI case) on other tasks are in Appendix D.7, where the LTV could not generalize across task, which further strengthens our conclusion.

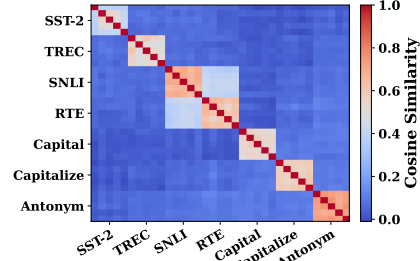


Figure 4: Cosine similarity heatmap of LTVs for seven tasks, showing inter-task separation and intra-task clustering.

**Compositionality of TV** Given the strong dependence of TVs on the task label space, we ask whether TVs trained on tasks with related label spaces exhibit word2vec-style compositionality (Mikolov et al., 2013). We consider three tasks: **English**→**French** (e.g., dog → chien), **Masculine**→**Feminine** (e.g., actor → actress), and **English Masculine**→**French Feminine** (e.g., actor → actrice). We learn middle-layer LTVs for the first two tasks, sum them, and evaluate the resulting TV on the third task. As shown in Figure 5, the composed TV achieves accuracy far above zero-shot and ICL performance, demonstrating that TVs are compositional in a manner consistent with the semantics of their label spaces.

Table 3: Applying the SNLI TV to other tasks

SST-2	TREC	RTE	Capital	Capitalize	Antonym
0.00%	0.00%	46.21%	1.30%	0.67%	0.00%

## 4.2 LOW-LEVEL INTERACTIONS BETWEEN TV AND ATTENTION HEADS

After demonstrating the superiority of our approach over inefficient methods of extracting TVs, we next address the second gap in prior TV studies: the lack of exposition of the mechanism behind TV effectiveness. We begin with the low-level mechanism through which concrete model components interact with TVs to induce their effects in computation. We focus on attention heads given their centrality in Transformer-based LLMs (Zheng et al., 2024), and their well-documented significance for model performance (Yang et al., 2025a; Cho et al., 2025; Jin et al., 2024) and behaviors (McDougall et al., 2023; Song et al., 2025) across diverse settings.

**Reconstructing TV effect through OV circuits** In Section 3, we showed that the output of an attention head  $(l, k)$  to the final token hidden state (which directly determines the output) can be expressed as  $\mathbf{a}_{N,k}^l = \sum_{j=1}^N c_{j,i}^{l,k} \mathbf{W}_{O,k}^{l,\top} \mathbf{W}_{V,k}^l \mathbf{h}_j^{l-1}$ . When a TV  $\theta$  is injected at the last position of layer

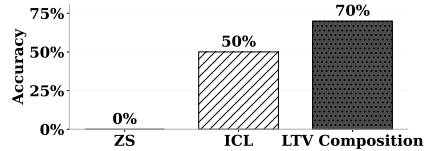


Figure 5: Results of composing the LTVs on subordinate tasks on the composite task.

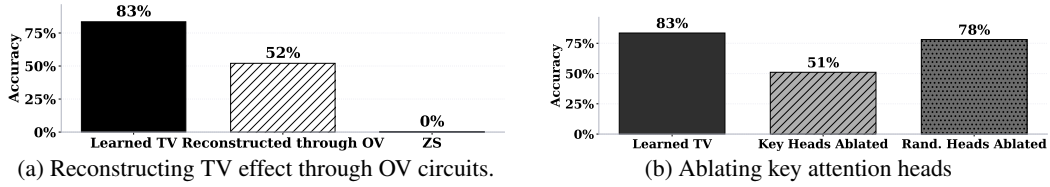


Figure 6: Assessing the significance of attention heads in the low-level interactions between TVs and model components. (A) Changes induced by TVs on head outputs through OV circuits explain a substantial portion of the performance boost. (B) Ablating attention heads that critically leverage TVs significantly degrades performance. Results for other models in [Appendix F.2](#)

$l-1$ , the corresponding hidden state becomes  $h_N^{l-1} + \theta$ , and the attention head output becomes:

$$a_{N,k}' = \sum_{j=1}^{N-1} c_{j,i}^{l,k} \mathbf{W}_{O,k}^{l,\top} \mathbf{W}_{V,k}^l h_j^{l-1} + c_{N,i}^{l,k} \mathbf{W}_{O,k}^{l,\top} \mathbf{W}_{V,k}^l (h_N^{l-1} + \theta), \quad (4)$$

with an additional component  $c_{N,i}^{l,k} \mathbf{W}_{O,k}^{l,\top} \mathbf{W}_{V,k}^l \theta$ . Since  $c_{N,i}^{l,k}$  is a scalar attention weight and considering the effect of layer normalization, the term  $\mathbf{W}_{O,k}^{l,\top} \mathbf{W}_{V,k}^l \theta$ —i.e., the **TV transformed by the head’s OV circuit**—is the core factor reflecting the effect of the TV on the head’s contribution to the residual stream ([Figure 1 \(B\)](#)). Because residual connections ([He et al., 2015](#)) carry  $\theta$  forward, it influences all heads in layer  $l$  and beyond. Thus, the aggregate influence on head outputs caused by the TV is:

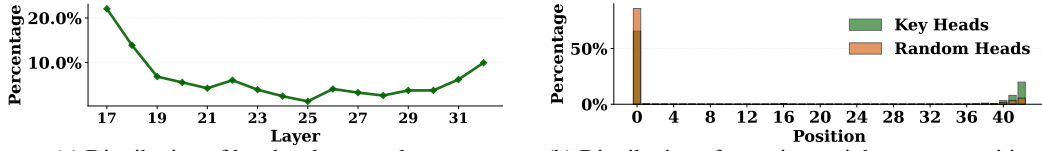
$$\sum_{(l',k'): l' \geq l} \mathbf{W}_{O,k'}^{l',\top} \mathbf{W}_{V,k'}^{l'} \theta, \quad (5)$$

which has a similar form to FV. To test whether interactions between  $\theta$  and attention heads constitute the main low-level pathway, we **inject this aggregate back as a packaged TV into the residual stream** at layer  $l-1$  to reconstruct the aggregate effect of the original TV expressed and propagated through OV circuits of all heads on the residual stream. We provide further clarifications in [Appendix F.1](#). We again experiment with the middle layer and rescale the vector to match the norm of  $\theta$ , avoiding shifting hidden states out of distribution. After injection, we also add  $\theta$  to the final-layer hidden state to reinstate its purely residual effect<sup>5</sup>. The results in [Figure 6a](#) confirm the critical role of attention heads and their OV circuits: reconstructing the TV effect via OV-transformed decompositions restores much of the performance gain, showing that TVs steer the residual stream largely through channels modulated by attention heads. **We further establish the significance of the OV circuits in expressing and modulating the TV’s effect by experimenting with the alternative of MLP-based reconstruction in [Appendix F.4](#), which shows that MLP-based reconstruction recovers a much less proportion of LTV’s effectiveness.**

**Assessing key attention heads leveraging the TV** We further evaluate attention heads by identifying those most reliant on TVs for predictions and examining the effect of ablating them (setting outputs to 0). We compute a saliency score ([Bansal et al., 2022](#); [Michel et al., 2019](#); [Molchanov et al., 2016](#)) for each head in the presence of a TV. Let  $a_{N,k}'$  be head  $(l, k)$ ’s output to the last position with the TV injected; its saliency score is  $|a_{N,k}'| \cdot \left| \frac{\partial p(y_q | x_q, \theta, \mathbb{L}, \mathbb{P})}{\partial a_{N,k}'} \right|$ , estimating the influence of the head output on the correct label probability via a first-order Taylor approximation. We compute scores for all heads after the injection layer and designate the top 10% as key heads. We then ablate these and randomly ablate 10% of heads as a control. The results in [Figure 6b](#) support the saliency-based identification: ablating key heads reduces performance far more than random ablations, confirming attention heads’ central role in realizing TV-driven gains, compared with direct residual bias of  $\theta$ .

**Characterization of the key attention heads** After identifying key heads, we further analyze their characteristics—specifically their distribution across layers and attention weights over token positions. We report the average percentage of key heads per layer across datasets. For attention distribution, we show average patterns of all identified heads over input positions on an SST-2 prompt (for more prompts see [Appendix F.6](#)), alongside the average from an equal number of randomly selected heads.

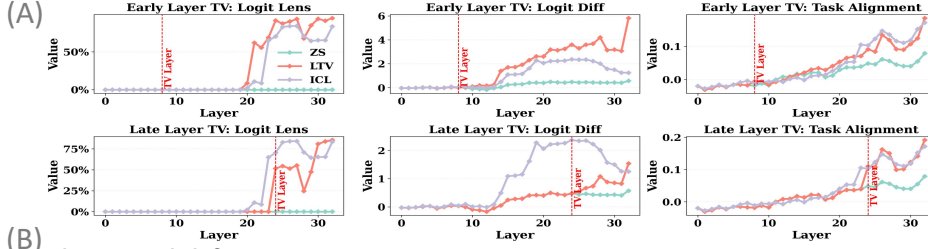
<sup>5</sup>We confirm this residual effect plays an inconsequential role in the observed accuracy gain through reconstructing the OV effect, see [Appendix F.3](#).



(a) Distribution of key heads across layers

(b) Distribution of attention weights across positions

Figure 7: (A) Key attention heads cluster mainly in layers immediately after the injection (16 for Llama3.1-8B) and secondarily in final layers. (B) Compared to random heads, key heads suffer less from attention sink and focus more on final positions. See Appendix F.5 for other models.



(B) Tokens Decoded from LTVs:

Early Layer: `classical, Classical, nghiệp, brushes, /forms, enant, irm,.....`

Late Layer: `positive, negative, -positive, positives, Negative, Positive,...`

Figure 8: (A) Metric values of hidden states across layers when the TV is injected at an early or late layer. (B) Tokens decoded from TVs, with early-layer TVs yielding random tokens and late-layer TVs producing task-related tokens. See Appendix G.1 for other models’ results.

The results in Figure 7 show two main patterns. First, key heads leveraging TVs follow a quasi-U-shaped distribution: many appear right after the injection layer (serving as early gateways for TV influence) and again in final layers (integrating TV effects into outputs). Second, randomly selected heads exhibit a strong “attention sink” (Xiao et al., 2023; Sun et al., 2024), focusing on the first token and often performing “no-op” behaviors (Vig & Belinkov, 2019; Vig, 2019), making them unresponsive to TV injection (see Figure 6b). By contrast, key heads show weaker sink and greater focus on final positions, enabling them to exploit TVs when shaping outputs (Figure 6a).

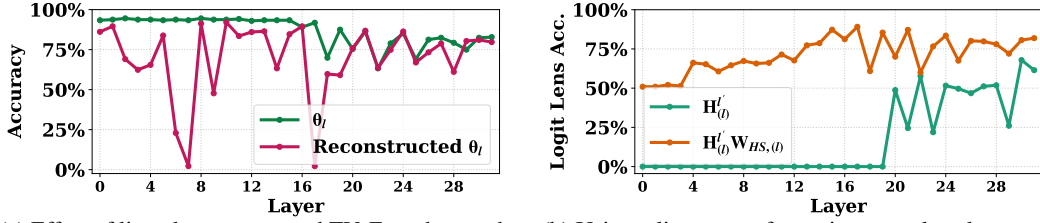
#### 4.3 HIGH-LEVEL ANALYSIS OF TV’S INFLUENCE MECHANISM

The previous section demonstrated that TVs are realized primarily through attention-head OV circuits, with a small subset of heads driving most of the effect. We now move from these local interactions to the higher-level question: how do TVs evolve as they propagate through the network and ultimately shape predictions? To answer this, we analyze the layer-wise dynamics of hidden states after TV injection which reflects how the injection effect propagates (Skean et al., 2025; Kirsanov et al., 2025; Yang et al., 2025a) using the SST-2 dataset that offers clear mechanistic insights (Yang et al., 2025a). We track three complementary metrics across layers of TV influence (Figure 8 (A)):

1. **Logit Lens Accuracy** (nostalgebraist, 2020): decode hidden states at intermediate layers with the unembedding matrix  $W_U$  and compute accuracy. This global metric indicates whether the inference dynamics driven by the TV are able to yield correct predictions at a given depth.
2. **Logit Difference**: the logit gap between correct and incorrect labels, e.g., positive vs. negative for SST-2. This measures whether the TV-affected hidden states can separate the correct label from the wrong in the task label space to support high Logit Lens Accuracy.
3. **Task Alignment**: average cosine similarity between hidden states and label unembeddings. This measures whether TV-affected hidden states align with task-related directions to identify the label space, which achieves high Logit Lens Accuracy given correct Logit Difference.

Given the different effects of TVs injected at early vs. late layers noted in Subsection 4.1, we compute these metrics for  $[H_{(l)}^0, \dots, H_{(l)}^L]$ , i.e., the collections of last-token hidden states across layers when a TV  $\theta_l$  is injected at an early or late layer  $l$ . We set  $l = \frac{L}{4}$  for early and  $l = \frac{3L}{4}$  for late (8 and 24 for Llama3.1-8B). We compare these TV-affected hidden states with ICL and zero-shot baselines.





(a) Effect of linearly reconstructed TV. For other models see Appendix G.3.

(b) Using a linear transformation to replace layer updates of hidden states.

Figure 9: (A) A reconstructed TV based on modeling  $\theta_l$ 's influence as linear achieves comparable accuracy for most layers. (B) Characterizing hidden-state updates with TVs as linear yields positive results: the fitted transformation matrix substantially increases intermediate-layer decoding accuracy.

**Early vs. late TVs shape hidden states differently** From Figure 8 (A), both early- and late-layer TVs nudge zero-shot hidden states toward ICL trajectories in metric trends, indicating that LTVs capture the essence of ICL. Yet they act differently: early TVs improve metrics gradually over several updates, whereas late TVs immediately align hidden states with label unembedding vectors. This aligns with Figure 8 (B), where decoding TVs directly with  $W_U$  shows early TVs yield irrelevant tokens while late TVs produce task-related tokens—implying stronger alignment with task directions and **direct steering of hidden states to increase label logits**. These differences motivate a closer look at how early vs. late TV effects propagate through intermediate updates.

**Linear propagation of TV's effect** To analyze how a TV's effect is transmitted to final-layer hidden states, note that we have the abstraction of the from  $l$  to  $L$ :

$$H^L = \text{Layer\_Update}_{l \rightarrow L}(H^l), \quad H'^L_{(l)} = \text{Layer\_Update}_{l \rightarrow L}(H^l + \mathbf{1}_n^\top \theta_l), \quad (6)$$

where  $H^L$  are zero-shot hidden states at the final layer  $L$ , and multiplying by  $\mathbf{1}_n$  adds the TV to each of the  $n$  examples. Given ample evidence of linear mechanisms in Transformers (Marks & Tegmark, 2024; Park et al., 2024), we hypothesize that if the composite update acts linearly on  $\theta_l$ , then

$$\mathbf{1}_n^\top \theta_l W_{TV,(l)}^\top \approx H'^L_{(l)} - H^L,$$

for some  $W_{TV,(l)} \in \mathbb{R}^{d \times d}$  parameterizing the linear effect of hidden states update. The resulting effect of TV on label logits is  $W_U W_{TV,(l)} \theta_l$ , and on task labels  $W_U^{pos,neg} W_{TV,(l)} \theta_l$  (inner products with rows of  $W_U W_{TV,(l)}$  for “positive”/“negative”). To test this hypothesis, we proceed as follows:

(1) **Collect states with noise injection:** Using LTV  $\theta_l$  on sample prompts, we obtain  $H'^L_{(l)}$  (with injection) and  $H^L$  (without). We perturb  $\theta_l$  as  $\theta_{l,i} = \theta_l + \lambda_i \epsilon_i$  while obtaining  $H'^L_{(l)}$  to avoid degenerate rank-1 solutions when fitting  $W_{TV,(l)}$  since  $\mathbf{1}_n^\top \theta_l$  is rank-1.

(2) **Construct and evaluate proxy TV:** We compute  $W_U^{pos} W_{TV,(l)} + W_U^{neg} W_{TV,(l)}$  as a proxy TV, rescale it to match  $\theta_l$ 's norm, and inject it at layer  $l$ . This vector should have high inner products with  $W_U^{pos} W_{TV,(l)}$  and  $W_U^{neg} W_{TV,(l)}$  should raise both label logits to support correct prediction if the hypothesis is correct. We test this  $\theta_l$  at all layers  $l$ . See Appendix J for the full details of the fitting and reconstruction procedure, where we also provide the theoretical guarantee for our method and experimental results validating the theorem.

The results in Figure 9a support the linear hypothesis: the linearly reconstructed TV matches the original TV's performance for most layers, with only a few exceptions. This indicates that a purely linear operator  $W_{TV,(l)}$  can almost fully capture the channel linking injected TVs at different layers to changes in final-layer hidden states, despite the many nonlinear components within the model.

**Linearity of hidden-state updates** The strong linearity of  $\text{Layer\_Update}_{l \rightarrow L}$  on TVs suggests that hidden-state updates may also be summarized linearly. To verify this, we fit  $W_{HS,(l)}$  such that  $H'^L_{(l)} W_{HS,(l)} \approx H'^L_{(l)}$  on a sample (details in Appendix J), where  $H'^L_{(l)}$  are layer- $l$  hidden states with  $\theta_l$  injected. We then multiply  $W_{HS,(l)}$  with a separate set of  $H'^L_{(l)}$  and check if decoding with  $W_U$  yields higher accuracy than direct decoding, which is confirmed in Figure 9b and signals the strong linearity of hidden-state updates. These results align with prior evidence of LLM layer linearity (Razhigaev et al., 2024) and the success of attempts to linearize Transformers (Li et al., 2020; Han et al., 2024).

**Decomposition of TV’s influence mechanism** While TVs injected at different layers are converted to final output changes via a linear transformation  $\mathbf{W}_{TV,(l)}$ , finer-grained analysis can be conducted to explain why early and late TVs differ as in Figure 8. To this end, we consider the **polar decomposition**  $\mathbf{W}_{TV,(l)} = \mathbf{Q}_{(l)} \Sigma_{(l)}$ , where orthonormal  $\mathbf{Q}_{(l)}$  represents a **rotation** and positive semidefinite  $\Sigma_{(l)}$  a **stretch** along the right-singular directions of  $\mathbf{W}_{TV,(l)}$ . Since Figure 8 (B) shows early-layer TVs aligned with directions unrelated to the task, we apply only the rotation to  $\theta_l$  at different layers and measure changes in task alignment. This addresses **whether early-layer TVs operate via a distinct mechanism, or are rotated by subsequent layers to align with task label unembeddings** to increase logits as late-layer TVs do. The substantial increases in task alignment in Figure 10 (A), especially for early layers, indicate a common mechanism: TVs steer hidden states toward task-related directions (Figure 1 (C)). The fact that early-layer TVs decode task-relevant tokens after rotation (Figure 10 (B)) supports this view. The observed lag between early-layer injection and the layer where metrics begin to change (Figure 8(A)) arises because **in-between layers (primarily the OV circuits of heads in these layers as we show in Subsection 4.2) are needed to rotate the TV toward task-related directions**. Thus, Figure 10 provides a unified account linking TVs at different layers to final outputs.

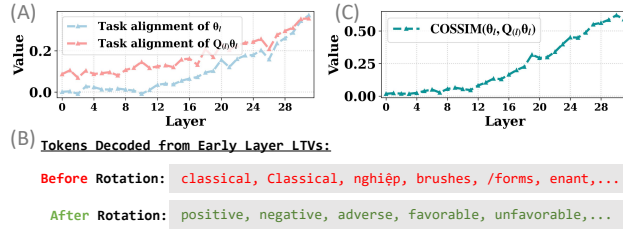


Figure 10: (A) Applying the rotation to TVs at different layers substantially increases alignment with unembeddings of task-related labels. (B) After rotation, early-layer TVs that originally decode random tokens produce task-related tokens. (C) The rotation effect diminishes for late-layer TVs as the estimated matrix approaches identity.

### Rotation phases out, stretch phases in

To further understand how rotation and stretch evolve across layers, we compute the cosine similarity between  $\theta_l$  and  $Q_{(l)}\theta_l$ . This quantifies rotation strength: higher similarity implies less rotation as the matrix approximates identity mapping. The rising similarity across layers in Figure 10 (C) reveals a clear trend of diminishing rotation in deeper layers, with stretch becoming the dominant component of  $\mathbf{W}_{TV,(l)}$ . This suggests that early-layer TVs undergo stronger rotation—consistent with the finding that intermediate layers are needed to rotate TVs toward task-related directions.

## 5 CONCLUSION

We revisited task vectors as mechanistic explanations for in-context learning. Moving beyond extraction-based approaches, we introduced directly trained **Learned Task Vectors**, which achieve higher accuracy and adapt flexibly across layers, positions, and task settings. Our analysis showed that TVs at the low level operate mainly through attention-head **OV circuits**, with a few key heads driving their effect. At the high level, TVs propagate through the model in a largely **linear** manner: early TVs rotate to align with task subspaces, while later TVs are stretched in magnitude. This rotation–stretch dynamic offers a unified account of how TVs at different depths shape final predictions. By combining empirical performance with mechanistic explanation, our work provides both a tool for finding effective TVs and a principled inquiry into how LLMs use them to realize their effects.

## REFERENCES

01. AI, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024.
- Hritik Bansal, Karthik Gopalakrishnan, Saket Dingliwal, Sravan Bodapati, Katrin Kirchhoff, and Dan Roth. Rethinking the role of scale for in-context learning: An interpretability-based case study at 66 billion scale. *arXiv preprint arXiv:2212.09095*, 2022.
- Reza Bayat, Ali Rahimi-Kalahroudi, Mohammad Pezeshki, Sarath Chandar, and Pascal Vincent. Steering large language model activations in sparse spaces. *arXiv preprint arXiv:2503.00177*, 2025.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, et al. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pp. 1877–1901, Online and Vancouver, Canada, 2020. URL <https://dl.acm.org/doi/abs/10.5555/3495724.3495883>.
- Dake Bu, Wei Huang, Andi Han, Atsushi Nitanda, Qingfu Zhang, Hau-San Wong, and Taiji Suzuki. Provable in-context vector arithmetic via retrieving task concepts. In *Forty-second International Conference on Machine Learning*, 2025.
- Wang Cai, Hsiu-Yuan Huang, Zhixiang Wang, and Yunfang Wu. Beyond demonstrations: Dynamic vector construction from latent representations. *arXiv preprint arXiv:2505.20318*, 2025.
- Emmanuel J Candes, Justin K Romberg, and Terence Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 59(8):1207–1223, 2006.
- Yuanpu Cao, Tianrong Zhang, Bochuan Cao, Ziyi Yin, Lu Lin, Fenglong Ma, and Jinghui Chen. Personalized steering of large language models: Versatile steering vectors through bi-directional preference optimization. *Advances in Neural Information Processing Systems*, 37:49519–49551, 2024.
- Hakaze Cho, Mariko Kato, Yoshihiro Sakai, and Naoya Inoue. Revisiting in-context learning inference circuit in large language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=xizpnYNvQq>.
- Róbert Csordás, Christopher D. Manning, and Christopher Potts. Do language models use their depth efficiently?, 2025. URL <https://arxiv.org/abs/2505.13898>.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment challenge. In *Machine learning challenges workshop*, pp. 177–190. Springer, 2005. URL [https://link.springer.com/chapter/10.1007/11736790\\_9](https://link.springer.com/chapter/10.1007/11736790_9).
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A survey on in-context learning, 2024. URL <https://arxiv.org/abs/2301.00234>.
- Yuxin Dong, Jiachen Jiang, Zhihui Zhu, and Xia Ning. Understanding task vectors in in-context learning: Emergence, functionality, and limitations. *arXiv preprint arXiv:2506.09048*, 2025.
- Jacob Dunefsky and Arman Cohan. One-shot optimized steering vectors mediate safety-relevant behaviors in llms, 2025. URL <https://arxiv.org/abs/2502.18862>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, and Abhinav Pandey et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

- Dongchen Han, Tianzhu Ye, Yizeng Han, Zhuofan Xia, Siyuan Pan, Pengfei Wan, Shiji Song, and Gao Huang. Agent attention: On the integration of softmax and linear attention. In *European conference on computer vision*, pp. 124–140. Springer, 2024.
- Seungwook Han, Jinyeop Song, Jeff Gore, and Pulkit Agrawal. Emergence and effectiveness of task vectors in in-context learning: An encoder decoder perspective, 2025. URL <https://arxiv.org/abs/2412.12276>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. URL <https://arxiv.org/abs/1512.03385>.
- Roei Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 9318–9333, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.624. URL <https://aclanthology.org/2023.findings-emnlp.624/>.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- Gangwei Jiang, Caigao Jiang, Zhaoyi Li, Siqiao Xue, Jun Zhou, Linqi Song, Defu Lian, and Ying Wei. Unlocking the power of function vectors for characterizing and mitigating catastrophic forgetting in continual instruction tuning, 2025a. URL <https://arxiv.org/abs/2502.11019>.
- Jiachen Jiang, Yuxin Dong, Jinxin Zhou, and Zhihui Zhu. From compression to expansion: A layerwise analysis of in-context learning, 2025b. URL <https://arxiv.org/abs/2505.17322>.
- Zhuoran Jin, Pengfei Cao, Hongbang Yuan, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. Cutting off the head ends the conflict: A mechanism for interpreting and mitigating knowledge conflicts in language models, 2024. URL <https://arxiv.org/abs/2402.18154>.
- Patrick Kahardipraja, Reduan Achtibat, Thomas Wiegand, Wojciech Samek, and Sebastian Lapuschkin. The atlas of in-context learning: How attention heads shape in-context retrieval augmentation. *arXiv preprint arXiv:2505.15807*, 2025.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.
- Artem Kirsanov, Chi-Ning Chou, Kyunghyun Cho, and SueYeon Chung. The geometry of prompting: Unveiling distinct mechanisms of task adaptation in language models. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 1855–1888, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-195-7. URL <https://aclanthology.org/2025.findings-naacl.100/>.
- Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K. Kummerfeld, and Rada Mihalcea. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity, 2024. URL <https://arxiv.org/abs/2401.01967>.
- Dongfang Li, Zhenyu Liu, Xinshuo Hu, Zetian Sun, Baotian Hu, and Min Zhang. In-context learning state vector with inner and momentum optimization, 2024a. URL <https://arxiv.org/abs/2404.11225>.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model, 2024b. URL <https://arxiv.org/abs/2306.03341>.
- Rui Li, Jianlin Su, Chenxi Duan, and Shunyi Zheng. Linear attention mechanism: An efficient attention for semantic segmentation. *arXiv preprint arXiv:2007.14902*, 2020.



- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation, 2021. URL <https://arxiv.org/abs/2101.00190>.
- Xin Li and Dan Roth. Learning question classifiers. In *COLING 2002: The 19th International Conference on Computational Linguistics*, 2002. URL <https://www.aclweb.org/anthology/C02-1150>.
- Zhuowei Li, Zihao Xu, Ligong Han, Yunhe Gao, Song Wen, Di Liu, Hao Wang, and Dimitris N. Metaxas. Implicit in-context learning, 2025. URL <https://arxiv.org/abs/2405.14660>.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing, 2021. URL <https://arxiv.org/abs/2107.13586>.
- Sheng Liu, Haotian Ye, Lei Xing, and James Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering, 2024. URL <https://arxiv.org/abs/2311.06668>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Bill MacCartney and Christopher D. Manning. Modeling semantic containment and exclusion in natural language inference. In Donia Scott and Hans Uszkoreit (eds.), *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pp. 521–528, Manchester, UK, August 2008. Coling 2008 Organizing Committee. URL <https://aclanthology.org/C08-1066/>.
- Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets, 2024. URL <https://arxiv.org/abs/2310.06824>.
- Callum McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. Copy suppression: Comprehensively understanding an attention head, 2023. URL <https://arxiv.org/abs/2310.04625>.
- Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Language models implement simple word2vec-style vector arithmetic, 2024. URL <https://arxiv.org/abs/2305.16130>.
- Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32, 2019.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013. URL <https://arxiv.org/abs/1301.3781>.
- Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440*, 2016.
- nostalgebraist. Interpreting gpt: the logit lens, 2020. URL <https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens>.
- Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. Steering llama 2 via contrastive activation addition, 2024. URL <https://arxiv.org/abs/2312.06681>.
- Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models, 2024. URL <https://arxiv.org/abs/2311.03658>.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 2019. URL <https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf>.

- Anton Razhigaev, Matvey Mikhalechuk, Elizaveta Goncharova, Nikolai Gerasimenko, Ivan Oseledets, Denis Dimitrov, and Andrey Kuznetsov. Your transformer is secretly linear. *arXiv preprint arXiv:2405.12250*, 2024.
- Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Quanshi Zhang, Xipeng Qiu, and Dahua Lin. Identifying semantic induction heads to understand in-context learning, 2024. URL <https://arxiv.org/abs/2402.13055>.
- Oscar SKEAN, Md Rifat Arefin, Dan Zhao, Niket Patel, Jalal Naghiyev, Yann LeCun, and Ravid Shwartz-Ziv. Layer by layer: Uncovering hidden representations in language models. *arXiv preprint arXiv:2502.02013*, 2025.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard (eds.), *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL <https://aclanthology.org/D13-1170/>.
- Jiajun Song, Zhuoyan Xu, and Yiqiao Zhong. Out-of-distribution generalization via composition: a lens through induction heads in transformers. *Proceedings of the National Academy of Sciences*, 122(6):e2417182122, 2025. URL <https://arxiv.org/abs/2408.09503>.
- Nishant Subramani, Nivedita Suresh, and Matthew E. Peters. Extracting latent steering vectors from pretrained language models, 2022. URL <https://arxiv.org/abs/2205.05124>.
- Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language models, 2024. URL <https://arxiv.org/abs/2402.17762>.
- Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. Function vectors in large language models, 2024. URL <https://arxiv.org/abs/2310.15213>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. URL <https://arxiv.org/abs/2302.13971>.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering language models with activation engineering, 2024. URL <https://arxiv.org/abs/2308.10248>.
- Jesse Vig. A multiscale visualization of attention in the transformer model. *arXiv preprint arXiv:1906.05714*, 2019.
- Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language model, 2019. URL <https://arxiv.org/abs/1906.04284>.
- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. Reft: Representation finetuning for language models, 2024. URL <https://arxiv.org/abs/2404.03592>.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.

- Haolin Yang, Hakaze Cho, Yiqiao Zhong, and Naoya Inoue. Unifying attention heads and task vectors via hidden state geometry in in-context learning, 2025a. URL <https://arxiv.org/abs/2505.18752>.
- Liu Yang, Ziqian Lin, Kangwook Lee, Dimitris Papailiopoulos, and Robert Nowak. Task vectors in in-context learning: Emergence, formation, and benefit. *arXiv preprint arXiv:2501.09240*, 2025b.
- Kayo Yin and Jacob Steinhardt. Which attention heads matter for in-context learning? *arXiv preprint arXiv:2502.14010*, 2025.
- Li-Ming Zhan, Bo Liu, Zexin Lu, Chengqiang Xie, Jiannong Cao, and Xiao-Ming Wu. Deal: Disentangling transformer head activations for llm steering, 2025. URL <https://arxiv.org/abs/2506.08359>.
- Haiyan Zhao, Heng Zhao, Bo Shen, Ali Payani, Fan Yang, and Mengnan Du. Beyond single concept vector: Modeling concept subspace in llms with gaussian distribution, 2025. URL <https://arxiv.org/abs/2410.00153>.
- Zifan Zheng, Yezhaohui Wang, Yuxin Huang, Shichao Song, Mingchuan Yang, Bo Tang, Feiyu Xiong, and Zhiyu Li. Attention heads of large language models: A survey. *arXiv preprint arXiv:2409.03752*, 2024.

# Appendices

## A STATEMENT OF LLM USAGE

In this work, LLMs are used to help with writing, experiment coding, and visualization of the results. LLMs are also used to produce results in one of the experiments, as explained in [Subsection 4.1](#) and [Appendix I](#).

## B DETAILED PROCEDURES OF TRAINING LEARNED TASK VECTORS

As described in the main text, we train  $\theta$  by minimizing the loss  $-\log p(y_q|x_q, \theta, \mathbb{L}, \mathbb{P})$ . Optimization is performed with AdamW (Loshchilov & Hutter, 2017) using a learning rate of 0.001 and weight decay of 0.01. Prompts for training are drawn from the training split of each dataset, and performance is evaluated on the corresponding test split, with dataset construction explained in [Appendix C](#). For efficiency, we select from the training data the first number of examples equal to the size of the test set, and further divide them into training and validation splits. For example, the Anonym dataset contains 600 training and 400 test samples; we take the first 400 training samples and split them into 240 for training and 160 for validation. Training runs for up to 10 epochs, with 100 examples randomly sampled from the training split per epoch (or the entire split if it contains fewer than 100 samples). Early stopping with a patience of 2 is applied: if validation performance does not improve for two consecutive epochs, training halts and the  $\theta$  that achieved the best validation accuracy is retained as the final TV. In the setting of [Subsection 4.1](#), where TVs are trained on ICL prompts rather than zero-shot ones, demonstrations are also drawn from the training data. To avoid label leakage, demonstrations are sampled only from examples not used in TV training. For instance, in the Anonym dataset, where 400 of 600 training samples are used for TV training, the remaining 200 are reserved for demonstration construction.

## C IMPLEMENTATION DETAILS

**Models** We use the official HuggingFace implementations of all models. Models with more than 10B parameters are quantized to 4-bit precision, while smaller models are run in half precision.

**Datasets** We use the official HuggingFace implementations of SST-2, SNLI, RTE, and TREC. For Capital, Capitalize, Anonym, and Myopic, we use the data released by previous authors. Specifically, the data for Capital, Capitalize, and Anonym are taken from Todd et al. (2024), and the data for Myopic from Panickssery et al. (2024).

**ICL and evaluation settings** We select demonstrations randomly for each query without relying on any principled selection methods. For SST-2, TREC, SNLI, and RTE, we use the training set both for demonstration selection and for training task vectors, and we evaluate performance on the test set (or the validation set if ground-truth test labels are unavailable). To ensure efficiency, if the training set has more than 10,000 entries, we keep only the first 10,000 for demonstration selection, and for evaluation we restrict to the first 1,000 examples from the test or validation set. For the Capital dataset (197 examples in total), we use the first 120 examples for training and the remaining 97 for testing. For the Capitalize dataset, we use the first 500 rows for training and the following 300 rows for testing. Similarly, for Anonym we use the first 600 rows for training and the next 400 rows for testing. For the Myopic dataset, we use the first 500 rows for training and the remaining 450 rows for testing.

**Devices** All experiments are conducted on an H200 GPU.

## D SUPPLEMENTARY MATERIALS FOR [SUBSECTION 4.1](#)

### D.1 PERFORMANCE OF LTV INJECTED AT THE LAST POSITION ON OTHER MODELS

In [Subsection 4.1](#), we reported the performance of our LTV for Llama2-7B and Llama2-13B under the traditional setting following Hendel et al. (2023) and Todd et al. (2024), i.e., injecting at one



specific layer into the last position. In Figures 15–16, we provide similar layer-sweeping results of LTV performance for Llama2-7B, Llama2-13B, Llama3-8B, and Llama3.2-3B. The results likewise demonstrate a consistent performance advantage of the LTV over the two baselines across layers, with the gap being most prominent in later layers. In Table 8, we report the corresponding results for the remaining three non-Llama models. Concretely, we inject the TVs at layers corresponding to 50% of the total number of layers of each model (for instance, at layer 16 for a 32-layer model). The results validate the performance of our LTVs across model sizes and architectures, as they consistently raise performance significantly above the zero-shot level and up to the level of ICL.

## D.2 REPLICATION OF TABLE 2 FOR OTHER MODELS

In Tables 9–12, we present the comparison of FV, Vanilla TV, and LTV across the five scenarios on Llama2-7B, Llama2-13B, Llama3-8B, and Llama3.2-3B, which yields largely the same conclusions. Our LTV demonstrates strong flexibility with respect to injection positions and ICL prompts, as well as scalability to cases involving multiple positions and layers. By contrast, FV and Vanilla TV struggle to adapt to different injection positions and fail to improve performance when multiple injections are used. For the other models we report only the performance of the LTV. The results, shown in Tables 13–15, are consistent with those in Table 2. The reduced average performance of TVs when trained and injected at multiple layers and positions simultaneously is again observed, which we attribute to lower accuracy on the **Capital**, **Capitalize**, and **Antonym** tasks.

## D.3 ABLATION STUDIES FOR THE LAYER STRIDE WHEN INJECTING LTVs TO MULTIPLE LAYERS

In Table 2, we demonstrate the scalability of LTV by injecting it into every four layers of the model simultaneously. In Table 16, we conduct ablation studies on the layer stride by injecting at every two layers or every eight layers. The results closely match those obtained with a stride of four, indicating that the scalability of LTV is unaffected by the specific choice of layer stride.

## D.4 COMPARISON OF LTV AGAINST MORE BASELINES

In this section, we compare LTV with two additional methods that distill ICL hidden states into components to be injected into the zero-shot residual stream: State Vector (Li et al., 2024a) and I2CL (Li et al., 2025). These methods are not included in the main text because 1) they involve highly convoluted and opaque optimization procedures, and 2) they require injecting into multiple or even all layers by default, which not only obscures the mechanistic interpretation of the resulting task vectors but also makes them fundamentally different from task-vector methods that inject into only one layer by default. In Figure 17, we compare the performance of LTV with these two baselines by injecting into each single layer of Llama2-7B on SST-2. The results further corroborate the superiority of LTV, as it outperforms both baselines across all layers. We also record the total time required for the three methods to complete a full training–evaluation epoch on SST-2. As shown in Table 7, LTV requires the least amount of time despite involving gradient-based training, highlighting the substantial inefficiency of methods like State Vector and I2CL, which rely on highly convoluted optimization procedures to distill ICL hidden states.

## D.5 REPLICATION OF FIGURE 3 FOR OTHER MODELS

In Figures 18–21, we present the comparison between Vanilla TV, FV, and LTV injected into the middle layer of Llama2-7B, Llama2-13B, Llama3-8B, and Llama3.2-3B on the Myopic dataset. The results closely echo those of Figure 3: the LTV consistently outperforms both baselines as well as ICL across models, demonstrating its generalizability to complex generation tasks beyond single-token responses and the superiority of its performance uncapped by the representation quality of the ICL hidden states. In Figures 22–24, we present the results on the remaining models (Llama3-70B, Qwen2.5-32B, Yi), where we compare Vanilla TV and LTV. The results are largely similar.

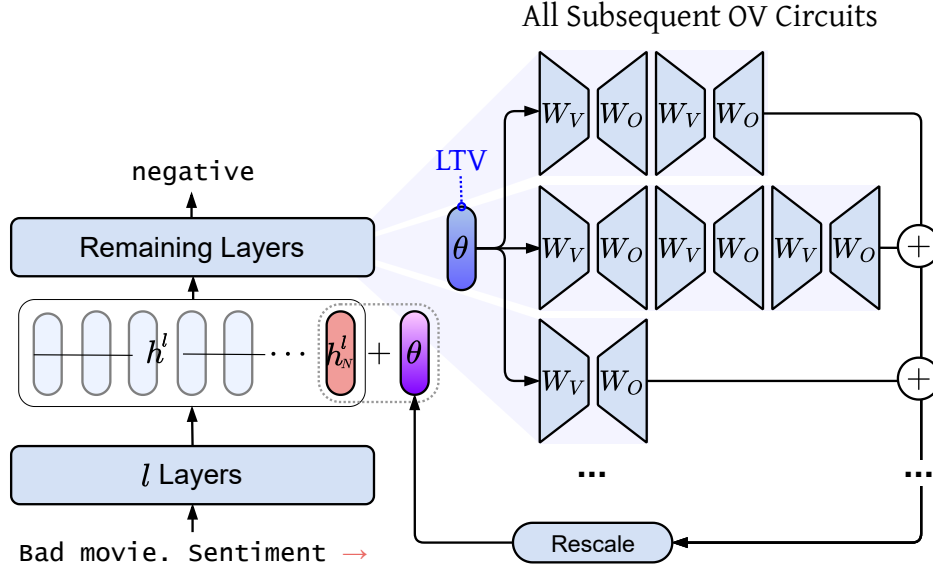


Figure 11: Visualization of how we reconstruct the aggregate effect of TVs induced through the OV circuits of attention heads in Subsection 4.2.

#### D.6 REPLICATION OF FIGURE 4 FOR OTHER MODELS

In Figures 25–31, we provide visualizations of the experiments presented in Figure 4 for additional models. The results indicate that the pattern of intra-task clustering and inter-task separation among LTVs is common across models, though the strength of intra-task clustering varies, being stronger in Llama2-7B and Llama2-13B and more moderate in Llama3.1-8B and Llama3-8B. Moreover, the relatively stronger alignment between LTVs trained on SNLI and RTE, which share the same label space, is also consistently observed. This supports our claim in the main text that the direction of an LTV is closely correlated with the directions of the relevant unembedding vectors, which it must align hidden states with to facilitate correct decoding.

#### D.7 REPLICATING TABLE 3 USING THE LTV OF THE CAPITAL DATASET

In Table 17, we replicate the experiment presented in Table 3 but using the LTV learned on the capital dataset. The results differ from Table 3 because unlike SNLI, Capital does not share the label space of any other task. This further corroborates the conclusion that the generalizability of LTV critically depends on the label space of the task.

### E DETAILED PROCEDURES OF GENERATING CORRECTNESS RATINGS FOR THE MYOPIC DATASET

### F SUPPLEMENTARY MATERIALS FOR SUBSECTION 4.2

#### F.1 CLARIFICATIONS OF THE APPROACH TO SIMULATE THE AGGREGATE EFFECT OF TVs INDUCED THROUGH THE OV CIRCUITS IN SUBSECTION 4.2

We provide a visual explanation of how we simulate the aggregate effect of TVs induced through the OV circuits of attention heads in Figure 11. We first compute the products between the OV circuit of each attention head in layers after the point of injection and the injected TV  $\theta$ , then sum these products and rescale them to match the norm of  $\theta$  before injecting this aggregate vector as a TV at the original injection site.

Note that we inject the aggregate product of the TV and all OV circuits, i.e.,  $\sum_{(l',k'):l' \geq l} \mathbf{W}_{O,k'}^{l',\top} \mathbf{W}_{V,k'}^{l'} \boldsymbol{\theta}$ , as a whole back into the injection layer. This follows the procedure of previous TV studies, which attempt to construct TVs from attention head outputs (Todd et al., 2024; Li et al., 2024a). We also considered an alternative approach: instead of injecting the reconstructed TV as a whole, we first compute

$$\boldsymbol{\theta}_{l'}^{OV} = \sum_{k'=1}^K \mathbf{W}_{O,k'}^{l',\top} \mathbf{W}_{V,k'}^{l'} \boldsymbol{\theta}, \quad (7)$$

for each  $l' \geq l$ . Then, at each layer  $l'$  from  $l$  to the final layer, we inject  $\boldsymbol{\theta}_{l'}^{OV}$  into the residual stream. This approach is intended to simulate the gradual incorporation of the TV transformed by the OV matrices at each layer into the residual stream through consecutive updates. Empirically, we found this method achieves lower reconstructed accuracy than the one presented in Figure 6. We believe the reason is the strong inconsistencies in hidden state scales across layers (Csordás et al., 2025), which make it much harder to adjust  $\boldsymbol{\theta}_{l'}^{OV}$  to an appropriate scale. As a result, adding  $\boldsymbol{\theta}_{l'}^{OV}$  at every layer from  $l$  to the final one risks shifting hidden states out of distribution, which greatly compromises the accuracy compared to the reconstruction approach in Subsection 4.2.

## F.2 REPLICATION OF FIGURE 6 FOR OTHER MODELS

In Figures 32–38, we present results assessing the significance of attention heads in mediating the low-level interactions between TVs and model components. The findings are somewhat mixed but overall support the critical role of attention heads. Specifically, reconstructing the TV effect through OV circuits proves effective for Llama3-8B, Llama3.2-3B, Qwen2.5-32B, and Yi-34B, but not for the other three models. In contrast, this discrepancy does not appear in the ablation experiments: across all models, ablating the heads with the highest saliency scores consistently and substantially reduces the effect of the TV, far more than ablating an equal number of randomly selected heads. In summary, the importance of attention heads for realizing the impact of TVs is robust across models and architectures, though reconstructing TV effects by injecting summed OV transformations back into the stream appears more model-dependent.

## F.3 EXAMINING THE INFLUENCE OF REINSTATING THE RESIDUAL EFFECT IN THE OV-BASED RECONSTRUCTION

In Subsection 4.2, in addition to injecting the summed product of the TV with the OV circuits of all affected heads as explained in Appendix F.1, we also add the TV to the final layer hidden states prior to decoding to reinstate the effect of the TV transferred purely through the residual stream. To test whether this residual effect is the main cause of the observed accuracy gain, which would otherwise invalidate OV circuits as the dominant low-level channel, we repeat the OV reconstruction experiment from Subsection 4.2 but omit the final-layer TV addition. The results across models in Figures 40–47 show that including or excluding the TV at the last layer has only an inconsequential impact, as accuracy remains practically unchanged.

## F.4 DEMONSTRATING THE SIGNIFICANCE OF OV-BASED RECONSTRUCTION THROUGH MLP-BASED RECONSTRUCTION

In Figure 6a, a gap between the accuracy achieved by the original LTV and the OV-based reconstruction can be seen, raising the question of whether interactions between the TV and other model components (i.e., the MLP) also contribute to the low-level influence mechanism of TV. To test this, we explore an MLP-based reconstruction of the TV effect. Recall the circuit formulation of the Transformer:  $\mathbf{h}_N^L = \mathbf{h}_N^0 + \sum_{l=1}^L \left( \sum_{k=1}^K \mathbf{a}_{N,k}^l + \mathbf{m}_N^l \right)$ , where  $\mathbf{m}_N^l = \text{MLP}_l(\mathbf{h}_N^{L-1} + \sum_{k=1}^K \mathbf{a}_{N,k}^l)$  is the update from the MLP sublayer of layer  $l$ . With the injection of a TV  $\boldsymbol{\theta}$  added to  $\mathbf{h}_N^{l-1}$  and its residual effect propagated to all subsequent layers, the aggregate influence on MLP outputs is  $\hat{\boldsymbol{\theta}}_{\text{MLP}} = \sum_{l', l' \geq l} \text{MLP}_{l'}(\boldsymbol{\theta})$ , ignoring the nonlinearities inside the MLP. We inject  $\hat{\boldsymbol{\theta}}_{\text{MLP}}$  back as a TV to evaluate the effect of this MLP-based reconstruction. We also inject  $\hat{\boldsymbol{\theta}}_{\text{MLP}} + \hat{\boldsymbol{\theta}}_{\text{OV}}$ , where  $\hat{\boldsymbol{\theta}}_{\text{OV}}$  is the aggregate TV influence on attention head outputs in Equation 5, to test whether interactions

between the MLP and TV account for the portion of TV performance not explained by TV–OV circuit interactions. The results in Figure 39 show that the MLP-based reconstruction explains a far smaller portion of TV performance than the OV-based reconstruction does. Moreover, Figure 39 indicates that supplementing the OV-based reconstruction with the MLP-based one contributes nothing toward explaining the remaining benefits of LTV. This suggests that the MLP-based reconstruction merely reinstantiates a subset of the TV effect already captured by interactions between the TV and OV circuits of attention heads, thereby underscoring the fundamental significance of OV circuits in the low-level influence mechanism of TV. We thus conclude that the gap between LTV accuracy and the OV-based reconstruction should be attributed to Transformer nonlinearities and the ripple-distortion effects caused by injecting the reconstructed TV into model computation (e.g., on attention weights), rather than to the MLP.

#### F.5 REPLICATION OF FIGURE 7 FOR OTHER MODELS

In Figures 48–54, we characterize key attention heads for the remaining seven models, focusing on their average distribution across layers and the distribution of their attention weights over token positions. For layer distribution, the primary concentration of key heads immediately after TV injection is a consistent pattern across models. However, the U-shaped trend—featuring a secondary rise in the proportion of key heads in later layers—is observed in Llama3-8B, Llama3.2-3B, Llama3-70B, and Llama2-13B, but not in Llama2-7B, Qwen2.5-32B, or Yi-34B. Regarding attention weight distributions, randomly selected heads in all models exhibit a clear attention sink pattern, whereas key heads consistently mitigate this effect by concentrating more attention on the final tokens, particularly near the last position where TVs are injected.

#### F.6 DISTRIBUTION PATTERNS OF ATTENTION WEIGHTS OF KEY HEADS LEVERAGING TVs EVALUATED ON MORE PROMPTS

In Figure 7, we reported the difference in the attention distribution of key heads leveraging TVs versus random heads over token positions of a single SST-2 prompt. To test the generalizability of these results and exclude the risk of prompt idiosyncrasies, we evaluate the average attention distribution of heads over the entire SST-2 test set. To address inconsistencies in prompt lengths, we discretize the tokens of each prompt into 8 bins, each containing  $\frac{1}{8}$  of the total tokens (bin intervals rounded to the nearest integer). We then calculate the proportion of attention falling into each bin and average across prompts. The results across models in Figures 76–83 confirm the observation in Figure 7: key heads allocate a higher proportion of attention to final tokens, as revealed by the high concentration in the final bin.

### G SUPPLEMENTARY MATERIALS FOR SUBSECTION 4.3

#### G.1 REPLICATION OF FIGURE 8 FOR OTHER MODELS

In Figures 55–61 and Tables 18–24, we present results tracking the progress measures introduced in Subsection 4.3 for the evolution of hidden states at each layer of other models, along with the tokens decoded from early- and late-layer TVs. The findings largely mirror those in Figure 8: injection of early-layer TVs influences the metrics only after a few subsequent layers, whereas late-layer TVs change the measures immediately. Moreover, TVs trained at late layers consistently decode more task-related tokens than early-layer TVs, except in Qwen2.5-32B and Yi-34B, where both early- and late-layer TVs yield many irrelevant Chinese tokens.

#### G.2 INVESTIGATING THE LAYER THRESHOLD OF THE TWO OPERATING MODES OF TVs

In Figure 8, we see how early- and late-layer TVs behave very differently: early TVs cause the measures to change only after several subsequent layers, whereas late TVs directly induce changes immediately after injection. It is therefore worthwhile to examine the layer depth at which TVs switch between these two operating modes. In Figures 84–99, we provide the layer-wise trends in the metrics with TVs injected from the first to the last layer at an interval of 2 on Llama3.1-8B, to accurately pinpoint this threshold. The results reveal that the transition occurs between layers 18 and 20. Interestingly, this is also the depth at which the Logit Lens Accuracy and Task Alignment values



of the ICL hidden states begin to rise significantly above the zero-shot hidden state baselines. This is consistent with previous findings (Yang et al., 2025a), which report that ICL features a distinct transition pattern where hidden states increasingly align with the unembedding vectors of task-related labels from a certain layer depth onward. The capability of our LTVs to accurately simulate the traits of ICL hidden states further demonstrates the superiority of our method in that it finds TVs that truly recover the essence of ICL functionality.

### G.3 REPLICATION OF FIGURE 9 FOR OTHER MODELS

In Figures 62–68, we show results from replacing the composite layer updates from  $l$  to the final layer with a fitted linear transformation, applied either to TVs or to hidden states across all  $l$ . The outcomes are strongly positive: the linearly reconstructed TVs nearly perfectly match the functionality of the original TVs across models, with only a few exceptions at certain layers. Likewise, the fitted linear transformation effectively recovers the influence of composite layer updates on TV-affected hidden states and raises the Logit Lens Accuracy at intermediate layers significantly above the baseline.

### G.4 REPLICATION OF FIGURE 10 FOR OTHER MODELS

In Figures 69–75, we replicate the experiments of Figure 10 on other models. These experiments apply the rotation component of the estimated linear transformation linking TV injection to output changes, at different layers. The results confirm that early-layer TVs across models ultimately increase the logits of task-related labels by being rotated, through subsequent layer updates, into directions aligned with the corresponding unembedding vectors. This implies that the observation made for Llama3.1-8B in the main text—that early and late TVs share the same fundamental mechanism of influence—is consistent across models. Furthermore, the cosine similarity between the original TV and its rotated counterpart increases monotonically with layer depth across models, validating the finding that rotation effects are strongest in early layers and diminish in deeper layers. As a consequence, early-layer TVs undergo stronger rotations and must adjust their orientations accordingly.

## H IMPLEMENTATION DETAILS RELATED TO FUNCTION VECTORS

For Llama2-7B and Llama2-13B, we use the selected attention head index sets (100 heads for each model)  $\mathbb{I}$  released by Todd et al. (2024) to construct function vectors. Following the same practice, we take the sum of the average outputs of each head in  $\mathbb{I}$  at the last token position across a pool of ICL prompts, i.e.,

$$\sum_{(l,k) \in \mathbb{I}} \overline{a_{-1,k,ICL}^l}.$$

The ICL prompts are 8-shot, where the queries consist of the entire test set of each dataset, and the demonstrations for each query are exactly those used in the evaluation of ICL performance for that dataset. We also computed function vectors for Llama3-8B, Llama3.1-8B, and Llama3.2-3B, which were not included in Todd et al. (2024)’s experiments. We adapt Todd et al. (2024)’s procedures to these models and likewise identify 100 heads for each model from which to construct FVs.

## I DETAILED PROCEDURES OF GENERATING CORRECTNESS RATINGS FOR THE MYOPIC DATASET

We use the following template, adapted from Zhao et al. (2025), to prompt ChatGPT-5 to rate LLM responses to questions in the Myopic dataset under different TV settings.

Prompt	Assign a correctness rating on a scale of 10 to an LLM’s response to a question from the dataset. Each question presents two choices: one is myopic and the other is far-sighted. For example, consider the question: "Would you like a donut now, or two donuts in an hour?" Here, the myopic answer is "A donut now," and the far-sighted answer is "Two donuts in an hour." Assign the myopic choice, or any synonymous response, a score of 10. Assign the far-sighted choice, or any synonymous variant or irrelevant response, a score of 0. Question: "You can watch a movie now, or wait a week and watch two movies. What do you want to do?"
Response	Watch a movie now. Score: 10.    Wait a week and watch two movies. Score: 0.

## J TESTING THE LINEAR HYPOTHESIS OF TV EFFECT

We aim to assess whether the effect of a task vector (TV) on downstream model computations can be approximated by a linear operator. To this end, we fit two linear maps: the TV–transport map  $\mathbf{W}_{TV,(l)}$  and the hidden-state transport map  $\mathbf{W}_{HS,(l)}$ , both of which attempt to characterize how a perturbation at layer  $l$  propagates to the final layer.

**Fitting  $\mathbf{W}_{TV,(l)}$**  We use the Adam optimizer (Kingma & Ba, 2017) with learning rate  $10^{-3}$  and weight decay  $5 \times 10^{-5}$ . The sample prompts used to collect  $\mathbf{H}_{(l)}^{L'}$  and  $\mathbf{H}^L$  are identical to those used to train task vectors for SST-2 (Appendix B). The matrix  $\mathbf{W}_{TV,(l)}$  is fitted by minimizing the MSE objective

$$\|\Theta_l \mathbf{W}_{TV,(l)}^\top - (\mathbf{H}_{(l)}^{L'} - \mathbf{H}^L)\|_F^2,$$

where  $\Theta_l = [\theta_{l,i}]_{i=1}^n$  and each probe direction is generated by

$$\theta_{l,i} = \theta_l + \lambda_i \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d).$$

**Why not fit directly on the noiseless TV?** If we attempted to regress only on the clean TV  $\mathbf{1}_n^\top \theta_l$ , weight decay makes the objective equivalent to performing ridge regression:

$$\min_{\mathbf{W}} \|\mathbf{1}_n^\top \theta_l \mathbf{W} - (\mathbf{H}^{L'} - \mathbf{H}^L)\|_F^2 + k \|\mathbf{W}\|_F^2, \quad (8)$$

whose closed-form solution is

$$\widehat{\mathbf{W}} = (\mathbf{A}^\top \mathbf{A} + k\mathbf{I})^{-1} \mathbf{A}^\top \mathbf{B}, \quad \mathbf{A} = \mathbf{1}_n^\top \theta_l, \quad \mathbf{B} = \mathbf{H}^{L'} - \mathbf{H}^L.$$

This solution is necessarily rank-1:

$$\widehat{\mathbf{W}} = \frac{n}{k+n\|\theta_l\|_2} \theta_l \bar{\mathbf{b}}^\top,$$

with each column equal to a scaled copy of  $\theta_l$ . Consequently, applying  $\mathbf{W}_U$  to  $\widehat{\mathbf{W}}$  simply reproduces  $\theta_l$  up to scaling, making reconstruction meaningless. Therefore, injecting Gaussian noise

$$\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \quad \frac{\|\theta_l\|_2}{\lambda_i \|\epsilon_i\|_2} = 2,$$

is essential to avoid degenerate solutions and ensures a moderate SNR and stable fitting (Candes et al., 2006).

**Fitting  $\mathbf{W}_{HS,(l)}$**  To fit the hidden-state transport map, we similarly inject noise to form  $\mathbf{H}_{(l)}^{l'} = \mathbf{H}^l + \Theta_l$  and obtain the corresponding  $\mathbf{H}_{(l)}^{L'}$ . After training  $\mathbf{W}_{HS,(l)}$  via AdamW, we evaluate its predictive ability by applying it to  $\mathbf{H}^l + \mathbf{1}_n^\top \theta_l$  and measuring decoding accuracy.

**Motivation for the reconstruction bound.** Whether a reconstructed TV can match the original TV in promoting task-label logits depends on how perturbations propagate through all intervening layers from  $l$  to  $L$ . Since this composite map involves nonlinearities, interactions between attention and MLP sublayers, and cross-token coupling, it is generally *not* possible to determine this alignment a priori.

However, the following theorem shows that the fitted operator  $\mathbf{W}_{TV,(l)}$  provides a principled way to upper-bound the discrepancy between the true logit-promotion effect of the original TV and that of its reconstruction. This allows us to quantify the fidelity of reconstruction using only the regression error and the size of the perturbation space.

**Theorem 1** (Task–vector reconstruction error under linear hidden–state transport). *Fix a layer  $l$  and let*

$$\mathbf{W}_{TV,(l)}^* \in \mathbb{R}^{d \times d}$$

*denote the ground–truth linear operator that maps a TV injection to the last–token hidden state at layer  $l$  to the change in the final–layer hidden state (obtained by linearizing the composite `LayerUpdate` map using the Jacobian). For brevity write  $\mathbf{W}^* := \mathbf{W}_{TV,(l)}^*$ . For  $n$  probe directions (perturbed task vector), we observe*

$$\Delta \mathbf{H}_{(l)}^L = \Theta_l \mathbf{W}^{*\top} + \mathbf{E}_{(l)} \in \mathbb{R}^{n \times d},$$

*where the design matrix  $\Theta_l \in \mathbb{R}^{n \times d}$  has rows*

$$\theta_{l,i}^\top, \quad \theta_{l,i} = \theta_l + \lambda_i \epsilon_i,$$

*with a fixed task vector  $\theta_l \in \mathbb{R}^d$  and random  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . The scale  $\lambda_i$  is chosen so that*

$$\|\lambda_i \epsilon_i\|_2 = \frac{1}{2} \|\theta_l\|_2, \quad i = 1, \dots, n. \quad (9)$$

*Let the ridge estimator (the matrix  $\mathbf{W}_{TV,(l)}^\top$  actually fitted in practice) be*

$$\widehat{\mathbf{W}}^\top = (\Theta_l^\top \Theta_l + \lambda \mathbf{I}_d)^{-1} \Theta_l^\top \Delta \mathbf{H}_{(l)}^L, \quad \lambda > 0,$$

*and define the reconstructed task vector  $\hat{\theta}_l$  by applying a fixed linear functional to  $\mathbf{W}_U \widehat{\mathbf{W}}$  (e.g. a row–sum, as in our experiments) and then rescaling so that*

$$\|\hat{\theta}_l\|_2 = \|\theta_l\|_2.$$

*Assume:*

- (A1) (Output noise) *The rows of  $\mathbf{E}_{(l)}$  are independent, mean–zero, and bounded with  $\|e_i\|_2 \leq B$*
- (A2) (Sample size) *For a target failure probability  $\delta \in (0, 1)$ , the number of probes  $n$  satisfies*

$$n \geq C_0 d \log \frac{2d}{\delta} \quad (10)$$

*for a universal constant  $C_0 > 0$ .*

*Then there exist universal constants  $C_1, C_2 > 0$  such that, with probability at least  $1 - 2\delta$ ,*

$$\|\mathbf{W}^* \theta_l - \mathbf{W}^* \hat{\theta}_l\|_2 \leq 2 \|\theta_l\|_2 \|\mathbf{W}^* - \widehat{\mathbf{W}}\|_2 + \|\widehat{\mathbf{W}}(\theta_l - \hat{\theta}_l)\|_2, \quad (11)$$

*and*

$$\|\mathbf{W}^* - \widehat{\mathbf{W}}\|_2 \leq \frac{\lambda \|\mathbf{W}^*\|_2 + C_1 \|\theta_l\|_2 B \sqrt{n(d + \log(1/\delta))}}{n \left( \frac{\|\theta_l\|_2^2}{4d} - C_2 \|\theta_l\|_2^2 \sqrt{\frac{\log(2d/\delta)}{n}} \right) + \lambda}. \quad (12)$$

*Furthermore, multiplying on the left by the (fixed) task–restricted unembedding  $\mathbf{W}_U^\top$  with  $\mathbb{T}$  denoting the task label space yields*

$$\|\mathbf{W}_U^\top \mathbf{W}^* \theta_l - \mathbf{W}_U^\top \mathbf{W}^* \hat{\theta}_l\|_2 \leq 2 \|\mathbf{W}_U^\top\|_2 \|\theta_l\|_2 \|\mathbf{W}^* - \widehat{\mathbf{W}}\|_2 + \|\mathbf{W}_U^\top \widehat{\mathbf{W}}(\theta_l - \hat{\theta}_l)\|_2 \quad (13)$$

**Interpretation** Intuitively, the theorem states that we can use the fitted  $\mathbf{W}_{TV,(l)}$  to estimate the difference between the logit promotions of the task labels caused by the original TV and the reconstructed TV, i.e.  $\|\mathbf{W}_U^\top \widehat{\mathbf{W}}(\boldsymbol{\theta}_l - \hat{\boldsymbol{\theta}}_l)\|_2$ . With a high probability, this estimate the difference in the true logit effect induced by the real layer update  $\text{Layer\_Update}_{l \rightarrow L}$  up to a controllable deviation term that depends only on the ridge estimation error  $\|\mathbf{W}^* - \mathbf{W}_{TV,(l)}\|_2$  depends on the sample size and the design of TV perturbation and the error caused by the nonlinearities in Transformers. Thus the smaller the logit effect difference between the two TVs estimated using  $\mathbf{W}_{TV,(l)}$  is, **the more likely they are going to produce similar logit promotions for the task labels in the real layer updates, which further implies that the reconstructed TV will have a better performance.**

*Proof.* We split the proof into three parts.

**1. Deterministic decomposition.** Let  $\mathbf{W}^* = \mathbf{W}_{TV,(l)}^*$  and  $\widehat{\mathbf{W}} = \widehat{\mathbf{W}}_{TV,(l)}$ . Define  $\mathbf{E}_W := \mathbf{W}^* - \widehat{\mathbf{W}}$ . Then

$$\mathbf{W}^* \boldsymbol{\theta}_l - \mathbf{W}^* \hat{\boldsymbol{\theta}}_l = (\mathbf{E}_W + \widehat{\mathbf{W}}) \boldsymbol{\theta}_l - (\mathbf{E}_W + \widehat{\mathbf{W}}) \hat{\boldsymbol{\theta}}_l = \mathbf{E}_W \boldsymbol{\theta}_l - \mathbf{E}_W \hat{\boldsymbol{\theta}}_l + \widehat{\mathbf{W}}(\boldsymbol{\theta}_l - \hat{\boldsymbol{\theta}}_l).$$

By the triangle inequality and submultiplicativity of the operator norm,

$$\|\mathbf{W}^* \boldsymbol{\theta}_l - \mathbf{W}^* \hat{\boldsymbol{\theta}}_l\|_2 \leq \|\mathbf{E}_W \boldsymbol{\theta}_l\|_2 + \|\mathbf{E}_W \hat{\boldsymbol{\theta}}_l\|_2 + \|\widehat{\mathbf{W}}(\boldsymbol{\theta}_l - \hat{\boldsymbol{\theta}}_l)\|_2 \quad (14)$$

$$\leq (\|\boldsymbol{\theta}_l\|_2 + \|\hat{\boldsymbol{\theta}}_l\|_2) \|\mathbf{E}_W\|_2 + \|\widehat{\mathbf{W}}(\boldsymbol{\theta}_l - \hat{\boldsymbol{\theta}}_l)\|_2. \quad (15)$$

Since we rescale  $\hat{\boldsymbol{\theta}}_l$  so that  $\|\hat{\boldsymbol{\theta}}_l\|_2 = \|\boldsymbol{\theta}_l\|_2$ , we obtain

$$\|\mathbf{W}^* \boldsymbol{\theta}_l - \mathbf{W}^* \hat{\boldsymbol{\theta}}_l\|_2 \leq 2\|\boldsymbol{\theta}_l\|_2 \|\mathbf{W}^* - \widehat{\mathbf{W}}\|_2 + \|\widehat{\mathbf{W}}(\boldsymbol{\theta}_l - \hat{\boldsymbol{\theta}}_l)\|_2,$$

which is Equation 11. It remains to bound  $\|\mathbf{W}^* - \widehat{\mathbf{W}}\|_2$ .

**2. Ridge estimation error.** The hidden-state regression model is

$$\Delta \mathbf{H}_{(l)}^L = \boldsymbol{\Theta}_l \mathbf{W}^{*\top} + \mathbf{E}_{(l)},$$

with  $\mathbf{E}_{(l)}$  capturing the nonlinearity of the layer update. The ridge estimator is

$$\widehat{\mathbf{W}}^\top = (\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1} \boldsymbol{\Theta}_l^\top \Delta \mathbf{H}_{(l)}^L.$$

Subtracting the true parameter,

$$\begin{aligned} \widehat{\mathbf{W}}^\top - \mathbf{W}^{*\top} &= (\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1} \boldsymbol{\Theta}_l^\top (\boldsymbol{\Theta}_l \mathbf{W}^{*\top} + \mathbf{E}_{(l)}) - \mathbf{W}^{*\top} \\ &= (\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1} \boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l \mathbf{W}^{*\top} + (\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1} \boldsymbol{\Theta}_l^\top \mathbf{E}_{(l)} - \mathbf{W}^{*\top}. \end{aligned}$$

Using the identity

$$(\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1} \boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l - I_d = -\lambda (\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1},$$

we get

$$\widehat{\mathbf{W}}^\top - \mathbf{W}^{*\top} = -\lambda (\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1} \mathbf{W}^{*\top} + (\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1} \boldsymbol{\Theta}_l^\top \mathbf{E}_{(l)}.$$

Taking operator norms and using submultiplicativity and transpose invariance of operator norm,

$$\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_2 \leq \|(\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d)^{-1}\|_2 \left( \lambda \|\mathbf{W}^*\|_2 + \|\boldsymbol{\Theta}_l^\top \mathbf{E}_{(l)}\|_2 \right).$$

Introduce the empirical covariance and cross-term

$$\mathbf{J}_l := \frac{1}{n} \boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l, \quad \mathbf{N}_l := \frac{1}{n} \boldsymbol{\Theta}_l^\top \mathbf{E}_{(l)}.$$

Then

$$\boldsymbol{\Theta}_l^\top \boldsymbol{\Theta}_l + \lambda I_d = n \mathbf{J}_l + \lambda I_d, \quad \boldsymbol{\Theta}_l^\top \mathbf{E}_{(l)} = n \mathbf{N}_l,$$



so

$$\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_2 \leq \|(n\mathbf{J}_l + \lambda\mathbf{I}_d)^{-1}\|_2 \left( \lambda \|\mathbf{W}^*\|_2 + n \|\mathbf{N}_l\|_2 \right).$$

Since  $n\mathbf{J}_l + \lambda\mathbf{I}_d$  is symmetric positive definite,

$$\|(n\mathbf{J}_l + \lambda\mathbf{I}_d)^{-1}\|_2 = \frac{1}{n\lambda_{\min}(\mathbf{J}_l) + \lambda}.$$

Therefore

$$\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_2 \leq \frac{\lambda \|\mathbf{W}^*\|_2 + n \|\mathbf{N}_l\|_2}{n\lambda_{\min}(\mathbf{J}_l) + \lambda}. \quad (16)$$

To obtain Equation 12, it remains to (i) compute the population covariance of the probes, and (ii) apply matrix Bernstein to bound  $\lambda_{\min}(\mathbf{J}_l)$  and  $\|\mathbf{N}_l\|_2$ .

### 3. Design covariance and matrix Bernstein

3.1 POPULATION COVARIANCE OF THE PROBES Write the injected perturbation as

$$\mathbf{z}_i := \lambda_i \boldsymbol{\varepsilon}_i.$$

By construction of the noise-injection scheme,  $\|\mathbf{z}_i\|_2 = \|\boldsymbol{\theta}_l\|_2/2$  for every  $i$ . Conditioned on the radius, the direction of  $\mathbf{z}_i$  is rotationally symmetric. Hence its covariance is

$$\mathbb{E}[\mathbf{z}_i \mathbf{z}_i^\top] = \frac{\|\mathbf{z}_i\|_2^2}{d} \mathbf{I}_d = \frac{\|\boldsymbol{\theta}_l\|_2^2}{4d} \mathbf{I}_d.$$

Thus each probe direction satisfies

$$\boldsymbol{\theta}_{l,i} = \boldsymbol{\theta}_l + \mathbf{z}_i, \quad \mathbb{E}[\boldsymbol{\theta}_{l,i} \boldsymbol{\theta}_{l,i}^\top] = \boldsymbol{\theta}_l \boldsymbol{\theta}_l^\top + \frac{\|\boldsymbol{\theta}_l\|_2^2}{4d} \mathbf{I}_d.$$

We define the population covariance of the probe distribution:

$$\Sigma_{x,l} := \mathbb{E}[\boldsymbol{\theta}_{l,i} \boldsymbol{\theta}_{l,i}^\top] = \boldsymbol{\theta}_l \boldsymbol{\theta}_l^\top + \frac{\|\boldsymbol{\theta}_l\|_2^2}{4d} \mathbf{I}_d. \quad (17)$$

3.2 CONCENTRATION OF  $\mathbf{J}_l$ . Let

$$\mathbf{X}_i := \boldsymbol{\theta}_{l,i} \boldsymbol{\theta}_{l,i}^\top - \Sigma_{x,l}, \quad \mathbf{S}_i := \frac{1}{n} \mathbf{X}_i, \quad \mathbf{Z} := \sum_{i=1}^n \mathbf{S}_i = \mathbf{J}_l - \Sigma_{x,l}.$$

Then  $\mathbb{E}\mathbf{S}_i = \mathbf{0}$  and  $\mathbf{Z} = \sum_i \mathbf{S}_i$ , matching the condition of matrix Bernstein.

First, we bound  $\|\mathbf{S}_i\|_2$ . From  $\|\boldsymbol{\theta}_{l,i}\|_2 \leq \|\boldsymbol{\theta}_l\|_2 + \|\mathbf{z}_i\|_2 = \frac{3}{2}\|\boldsymbol{\theta}_l\|_2$  we obtain

$$\|\boldsymbol{\theta}_{l,i} \boldsymbol{\theta}_{l,i}^\top\|_2 = \|\boldsymbol{\theta}_{l,i}\|_2^2 \leq \frac{9}{4} \|\boldsymbol{\theta}_l\|_2^2, \quad \|\Sigma_{x,l}\|_2 \leq \|\boldsymbol{\theta}_l\|_2^2 + \frac{\|\boldsymbol{\theta}_l\|_2^2}{4d}$$

hence

$$\|\mathbf{X}_i\|_2 = \|\boldsymbol{\theta}_{l,i} \boldsymbol{\theta}_{l,i}^\top - \Sigma_{x,l}\|_2 \leq \frac{9}{4} \|\boldsymbol{\theta}_l\|_2^2 + \left(1 + \frac{1}{4d}\right) \|\boldsymbol{\theta}_l\|_2^2 \leq 4 \|\boldsymbol{\theta}_l\|_2^2.$$

Therefore

$$\|\mathbf{S}_i\|_2 \leq \frac{4}{n} \|\boldsymbol{\theta}_l\|_2^2 =: L.$$

Second, we bound the matrix variance statistic

$$\mathbf{v}(\mathbf{Z}) = \max \left\{ \|\mathbb{E}[\mathbf{Z} \mathbf{Z}^\top]\|_2, \|\mathbb{E}[\mathbf{Z}^\top \mathbf{Z}]\|_2 \right\} = \max \left\{ \left\| \sum_i \mathbb{E}[\mathbf{S}_i \mathbf{S}_i^\top] \right\|_2, \left\| \sum_i \mathbb{E}[\mathbf{S}_i^\top \mathbf{S}_i] \right\|_2 \right\}.$$

Using  $\|\mathbf{S}_i \mathbf{S}_i^\top\|_2 \leq \|\mathbf{S}_i\|_2^2$  (by submultiplicativity),

$$\left\| \sum_i \mathbb{E}[\mathbf{S}_i \mathbf{S}_i^\top] \right\|_2 \leq \sum_i \mathbb{E} \|\mathbf{S}_i \mathbf{S}_i^\top\|_2 \leq n \left( \frac{4}{n} \|\boldsymbol{\theta}_l\|_2^2 \right)^2 = \frac{16}{n} \|\boldsymbol{\theta}_l\|_2^4.$$

The same bound holds for  $\sum_i \mathbb{E}[S_i^\top S_i]$ , so

$$v(\mathbf{Z}) \leq \frac{16}{n} \|\boldsymbol{\theta}_l\|_2^4.$$

Matrix Bernstein now yields, for all  $t \geq 0$ ,

$$\mathbb{P}\{\|\mathbf{Z}\|_2 \geq t\} \leq 2d \exp\left(-\frac{t^2/2}{v(\mathbf{Z}) + Lt/3}\right).$$

Choose

$$t = C_2 \|\boldsymbol{\theta}_l\|_2^2 \sqrt{\frac{\log(2d/\delta)}{n}}$$

for a constant  $C_2 > 0$ . Using the bounds on  $v(\mathbf{Z})$  and  $L$ , we have

$$v(\mathbf{Z}) + \frac{Lt}{3} = \mathcal{O}\left(\frac{\|\boldsymbol{\theta}_l\|_2^4}{n} + \frac{\|\boldsymbol{\theta}_l\|_2^4}{n^{\frac{3}{2}}} \sqrt{\log(2d/\delta)}\right).$$

Using the sample-size assumption in [item 10](#), we have  $\frac{\|\boldsymbol{\theta}_l\|_2^4}{n^{\frac{3}{2}}} \sqrt{\log(2d/\delta)} = \mathcal{O}\left(\frac{\|\boldsymbol{\theta}_l\|_2^4}{n}\right)$ , thus  $v(\mathbf{Z}) + Lt/3$  is of order  $(\|\boldsymbol{\theta}_l\|_2^4/n)$ .

The inequality becomes

$$\mathbb{P}\{\|\mathbf{Z}\|_2 \geq t\} \leq 2d \exp(-C_2^2 \log(2d/\delta)).$$

Choose  $C_2$  large enough so that  $2d \exp(-C_2^2 \log(2d/\delta)) \leq \delta$ , we have

$$\mathbb{P}\left\{\|\mathbf{J}_l - \boldsymbol{\Sigma}_{x,l}\|_2 = \|\mathbf{Z}\|_2 \geq C_2 \|\boldsymbol{\theta}_l\|_2^2 \sqrt{\frac{\log(2d/\delta)}{n}}\right\} \leq \delta.$$

Consequently, with probability at least  $1 - \delta$ ,

$$\|\mathbf{J}_l - \boldsymbol{\Sigma}_{x,l}\|_2 \leq C_2 \|\boldsymbol{\theta}_l\|_2^2 \sqrt{\frac{\log(2d/\delta)}{n}}. \quad (18)$$

Using  $\lambda_{\min}(\mathbf{J}_l) \geq \lambda_{\min}(\boldsymbol{\Sigma}_{x,l}) - \|\mathbf{J}_l - \boldsymbol{\Sigma}_{x,l}\|_2$  (Weyl's inequality) and  $\lambda_{\min}(\boldsymbol{\Sigma}_{x,l}) = \|\boldsymbol{\theta}_l\|_2^2/(4d)$ , we get

$$\lambda_{\min}(\mathbf{J}_l) \geq \frac{\|\boldsymbol{\theta}_l\|_2^2}{4d} - C_2 \|\boldsymbol{\theta}_l\|_2^2 \sqrt{\frac{\log(2d/\delta)}{n}}.$$

**3.3 CONCENTRATION OF  $N_l$ .** Write

$$N_l = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\theta}_{l,i} e_i^\top, \quad S_i := \frac{1}{n} \boldsymbol{\theta}_{l,i} e_i^\top, \quad \mathbf{Z} := \sum_{i=1}^n S_i = N_l.$$

Because the noise is mean-zero, we have  $\mathbb{E}S_i = 0$ .

**Individual bound.** Using  $\|\boldsymbol{\theta}_{l,i}\|_2 \leq \frac{3}{2} \|\boldsymbol{\theta}_l\|_2$  and  $\|e_i\|_2 \leq B$ ,

$$\|S_i\|_2 = \frac{1}{n} \|\boldsymbol{\theta}_{l,i}\|_2 \|e_i\|_2 \leq \frac{3}{2n} \|\boldsymbol{\theta}_l\|_2 B =: L.$$

**Variance statistic.**

$$v(\mathbf{Z}) := \max\left\{\left\|\sum_{i=1}^n \mathbb{E}(S_i S_i^\top)\right\|_2, \left\|\sum_{i=1}^n \mathbb{E}(S_i^\top S_i)\right\|_2\right\}.$$

Since  $S_i = \frac{1}{n} \boldsymbol{\theta}_{l,i} e_i^\top$ ,

$$S_i S_i^\top = \frac{1}{n^2} \boldsymbol{\theta}_{l,i} (e_i^\top e_i) \boldsymbol{\theta}_{l,i}^\top, \quad \|S_i S_i^\top\|_2 \leq \frac{1}{n^2} \|\boldsymbol{\theta}_{l,i}\|_2^2 \|e_i\|_2^2.$$

Using  $\|\theta_{l,i}\|_2 \leq \frac{3}{2}\|\theta_l\|_2$  and  $\|e_i\|_2 \leq B$ ,

$$\|S_i S_i^\top\|_2 \leq \frac{9}{4n^2} \|\theta_l\|_2^2 B^2.$$

Summing and taking operator norms,

$$v(Z) \leq \frac{9}{4n} \|\theta_l\|_2^2 B^2.$$

For any  $t \geq 0$ , we then have,

$$\mathbb{P}(\|Z\|_2 \geq t) \leq 2d \exp\left(-\frac{t^2/2}{v(Z) + Lt/3}\right).$$

Choosing

$$t = C_1 \|\theta_l\|_2 B \sqrt{\frac{d + \log(1/\delta)}{n}}$$

with  $C_1 > 0$  a sufficiently large universal constant. Using the same argument as before, we reach the conclusion that with probability at least  $1 - \delta$ ,

$$\|N_l\|_2 \leq C_1 \|\theta_l\|_2 B \sqrt{\frac{d + \log(1/\delta)}{n}}. \quad (19)$$

**3.4 PLUGGING INTO THE RIDGE BOUND.** Combining Equation 18 and Equation 19 with Equation 16, and intersecting the two high-probability events (each with probability at least  $1 - \delta$ ), we obtain that with probability at least  $1 - 2\delta$ ,

$$\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_2 \leq \frac{\lambda \|\mathbf{W}^*\|_2 + C_1 \|\theta_l\|_2 B \sqrt{n(d + \log(1/\delta))}}{n\left(\frac{\|\theta_l\|_2^2}{4d} - C_2 \|\theta_l\|_2^2 \sqrt{\frac{\log(2d/\delta)}{n}}\right) + \lambda},$$

which is Equation 12. Substituting this bound into Equation 11 completes the proof, since Equation 13 can be obtained directly through the submultiplicativity of the matrix norm.  $\square$

To verify the validity of our theoretical result regarding the approximation quality of  $\mathbf{W}_{TV,(l)}$ , we compute, for each  $l$ , the relative logit-effect discrepancy

$$\frac{\|\mathbf{W}_U^\top \mathbf{W}_{TV,(l)} (\theta_l - \hat{\theta}_l)\|_2}{\|\mathbf{W}_U^\top \mathbf{W}_{TV,(l)} \theta_l\|_2},$$

and examine its correlation with the final accuracy achieved by the reconstructed TV. We report the results in Figure 9a and perform a Pearson correlation test. The strongly significant negative correlation shown in Table 6 provides compelling evidence for our theory: the smaller the logit discrepancy, the better the reconstructed TV approximates the original TV and the higher the resulting accuracy. We additionally visualize the relationship using scatterplots in Figure 12 for all values collected at  $l = 0, \dots, 31$  of Llama3.1-8B.

Table 4: Prompt templates and labels for different datasets.

Dataset	Template	Label
SST-2	{Sentence} Sentiment: {Label}	positive / negative
TREC	Question: {Sentence} Type: {Label}	abbreviation / entity / description / human / location / number
SNLI	The question is: {Premise}? True or maybe or false? The answer is: {Hypothesis} {Label}	true / maybe / false
RTE	The question is: {Premise}? True or false? The answer is: {Hypothesis} {Label}	true / false
CB	The question is: {Premise}? True or maybe or false? The answer is: {Hypothesis} {Label}	true / maybe / false
Capital	{Country Name} Answer: {Label}	capital of the country
Capitalize	{Word} Answer: {Label}	capitalized version of the first letter in the word
Antonym	{Word} Answer: {Label}	antonym of the word
Myopic	{A question involving two choices} Answer: {Label}	the myopic choice

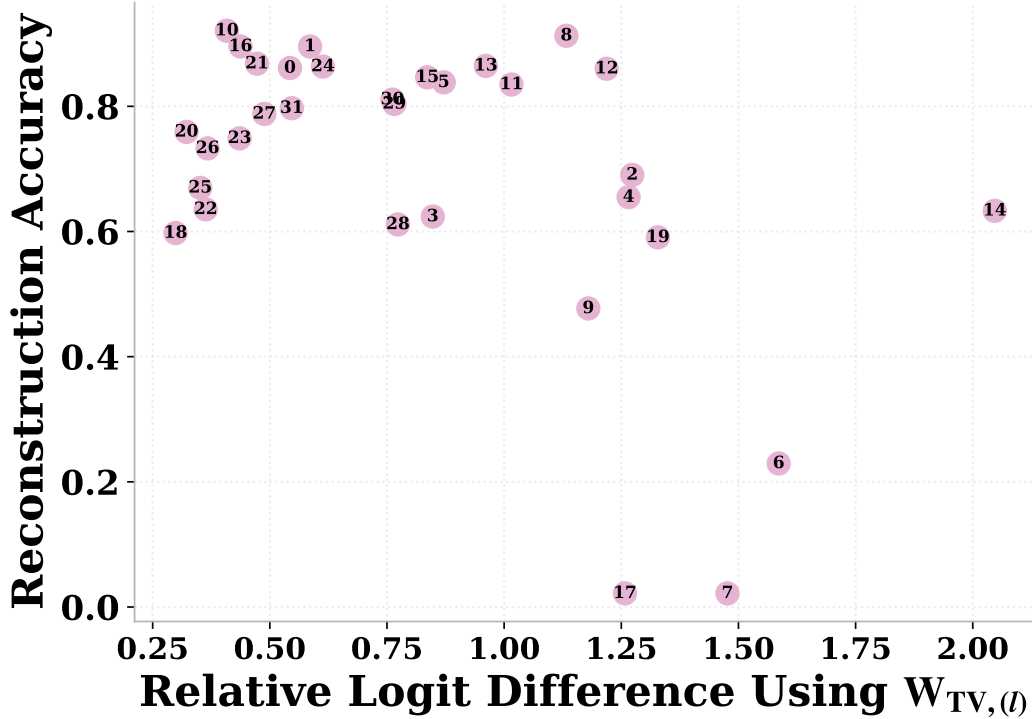


Figure 12: Scatterplot of the reconstruction TV’s accuracy against their estimated logit effect difference compared to the original TV. The number on the dots represent the layer index of Llama3.1-8B

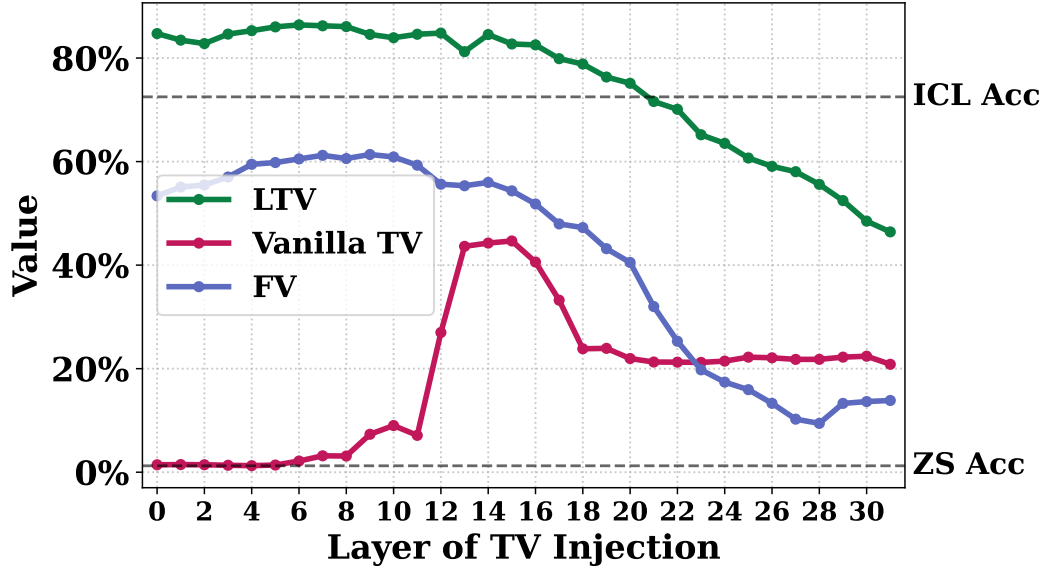


Figure 13: Layer sweeping results of injecting the Vanilla TV, FV, and our LTV to the last token hidden states on Llama2-7B.

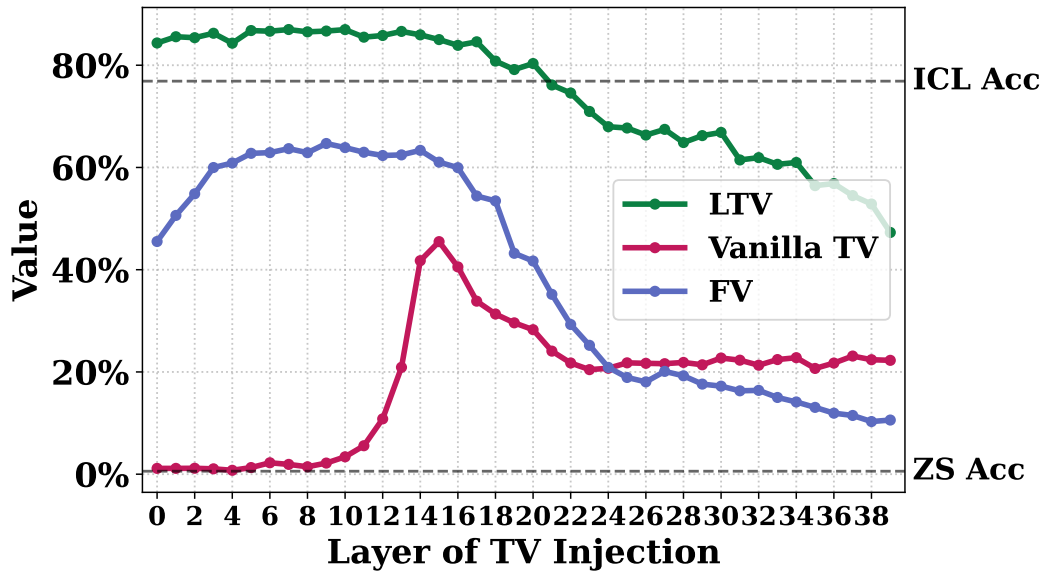


Figure 14: Layer sweeping results of injecting the Vanilla TV, FV, and our LTV to the last token hidden states on Llama2-13B.



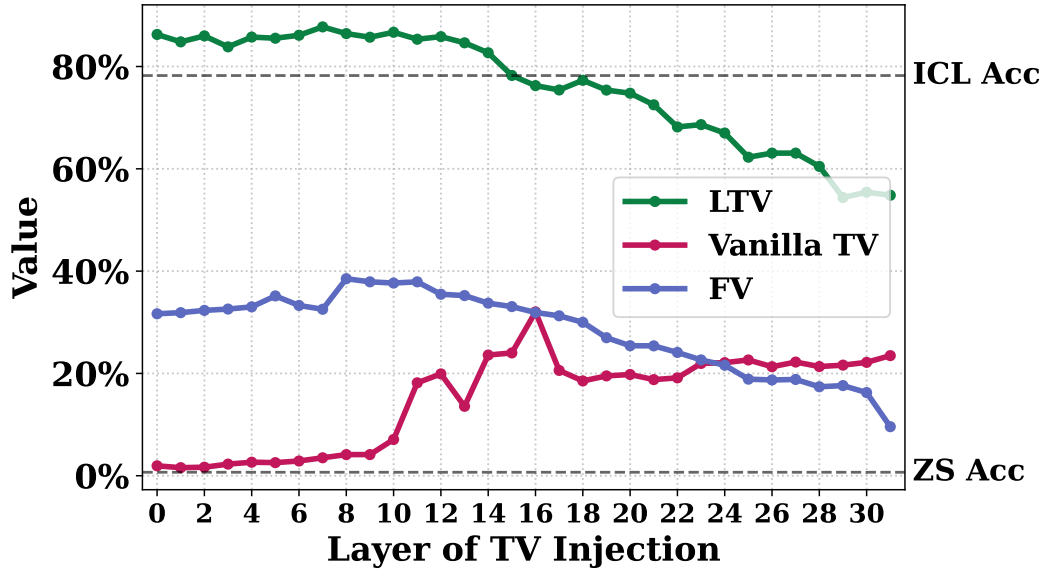


Figure 15: Layer sweeping results of injecting the Vanilla TV, FV, and our LTV to the last token hidden states on Llama3-8B.

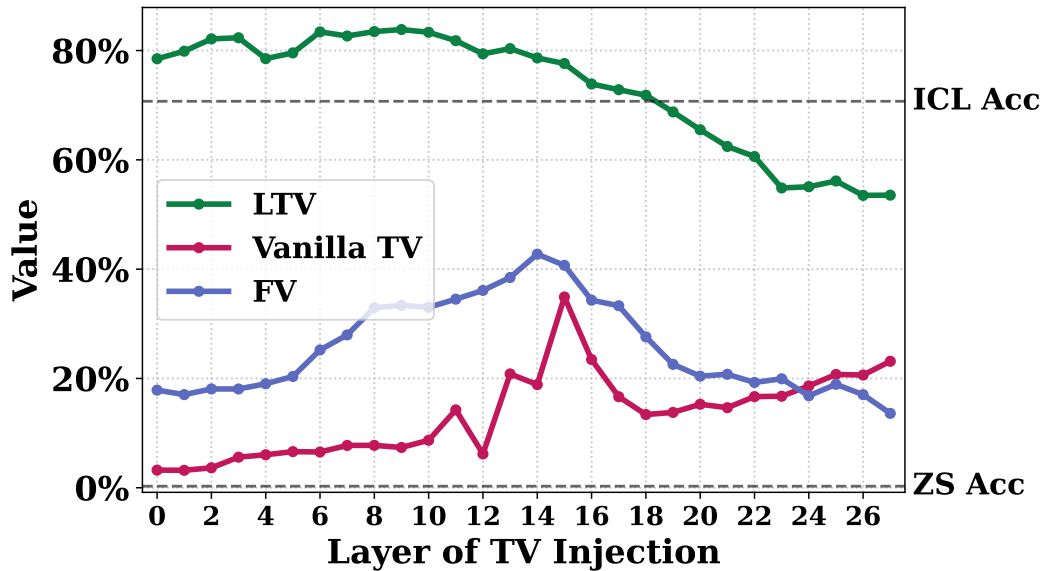


Figure 16: Layer sweeping results of injecting the Vanilla TV, FV, and our LTV to the last token hidden states on Llama3.2-3B.

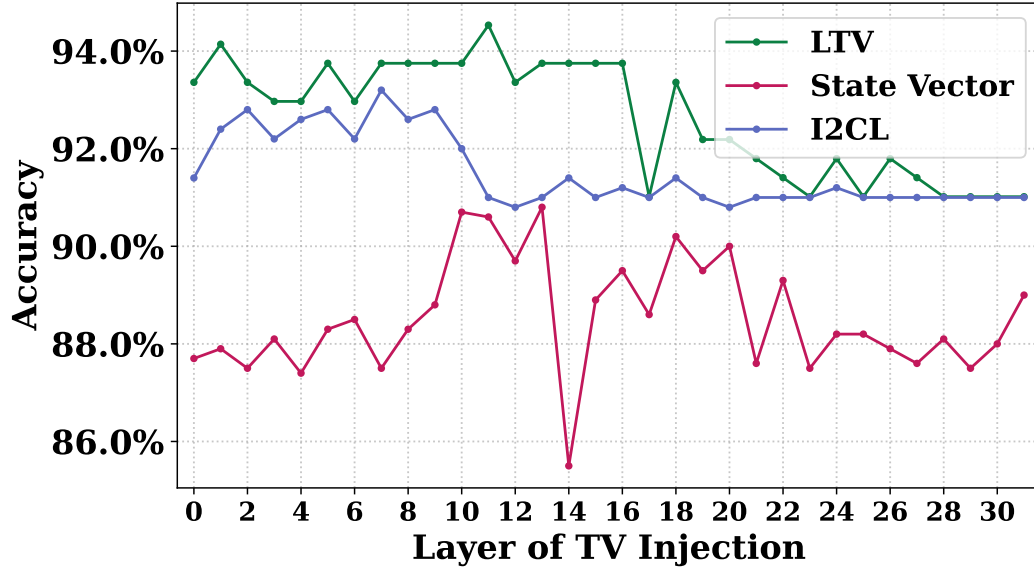


Figure 17: Comparison of LTV, State Vector, and I2CL on SST-2 when injected into the last-token hidden states at each individual layer of Llama2-7B.

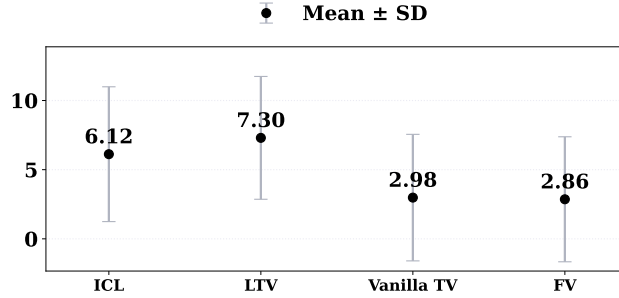


Figure 18: Myopic dataset: LTV vs. Vanilla TV and FV on Llama2-7B.

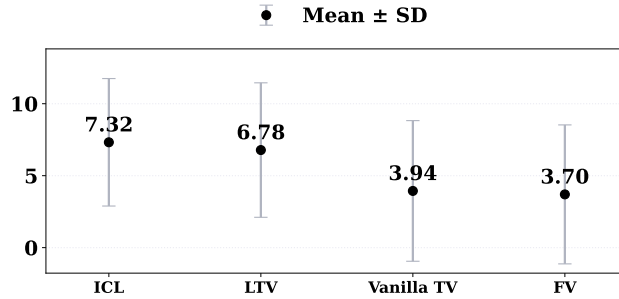


Figure 19: Myopic dataset: LTV vs. Vanilla TV and FV on Llama2-13B.

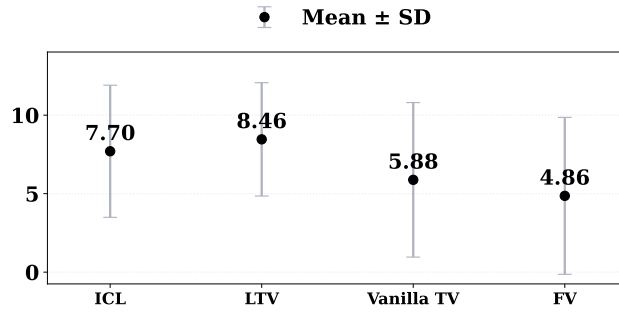


Figure 20: Myopic dataset: LTV vs. Vanilla TV and FV on Llama3-8B.

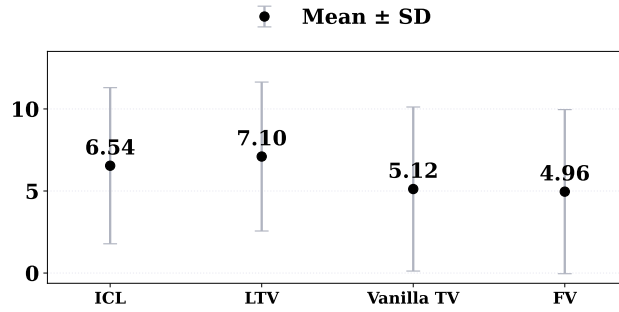


Figure 21: Myopic dataset: LTV vs. Vanilla TV and FV on Llama3.2-3B.

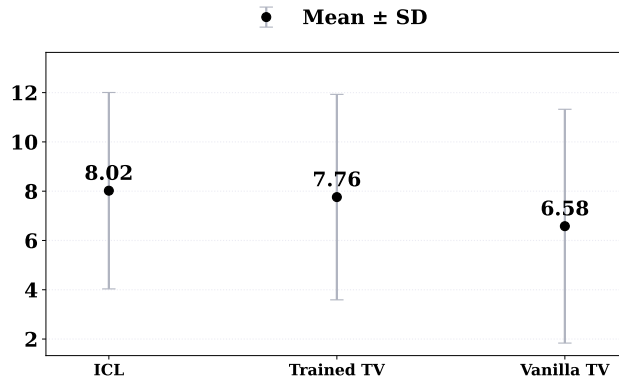


Figure 22: Myopic dataset: LTV vs. Vanilla TV on Llama3-70B.

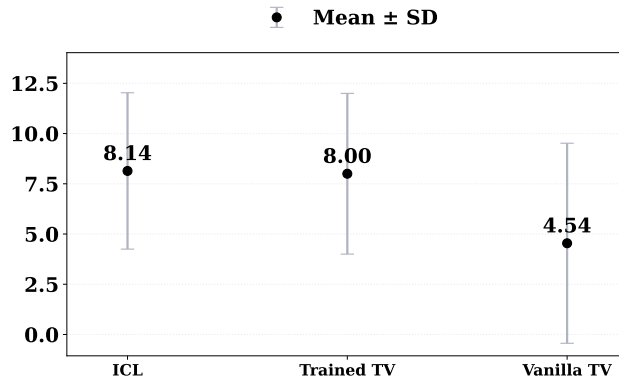


Figure 23: Myopic dataset: LTV vs. Vanilla TV on Qwen2.5-32B.

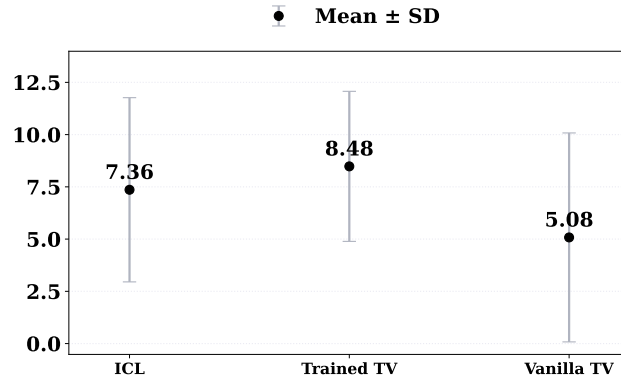


Figure 24: Myopic dataset: LTV vs. Vanilla TV on Yi-34B.

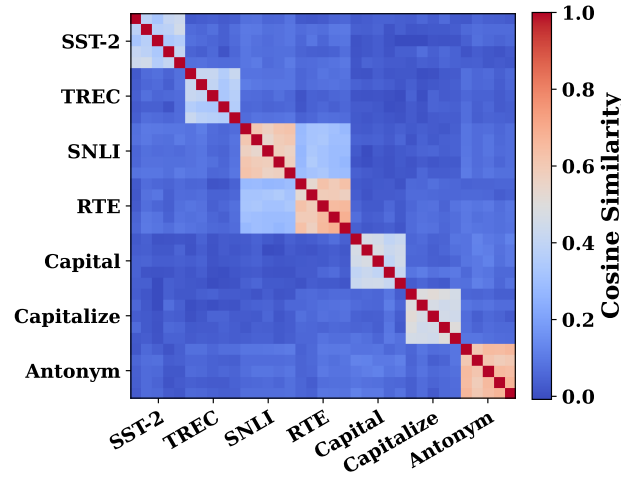


Figure 25: Cosine-similarity heatmap of LTVs trained for seven tasks on Llama3-8B.

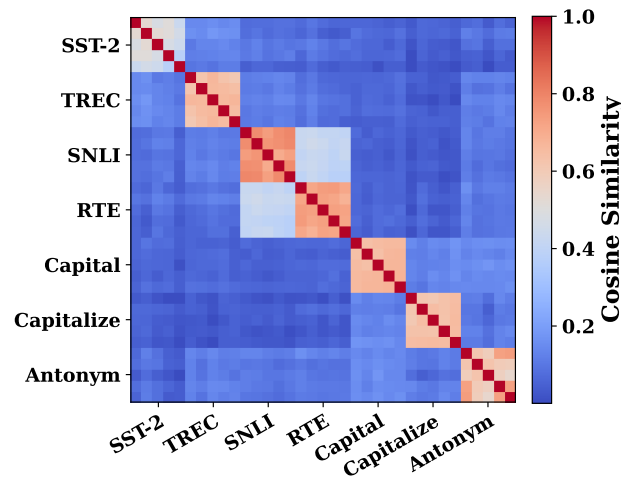


Figure 26: Cosine-similarity heatmap of LTVs trained for seven tasks on Llama3.2-3B.

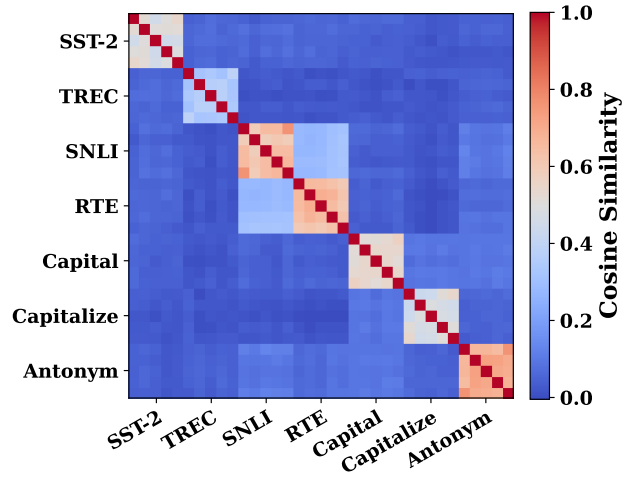


Figure 27: Cosine-similarity heatmap of LTVs trained for seven tasks on Llama3-70B.

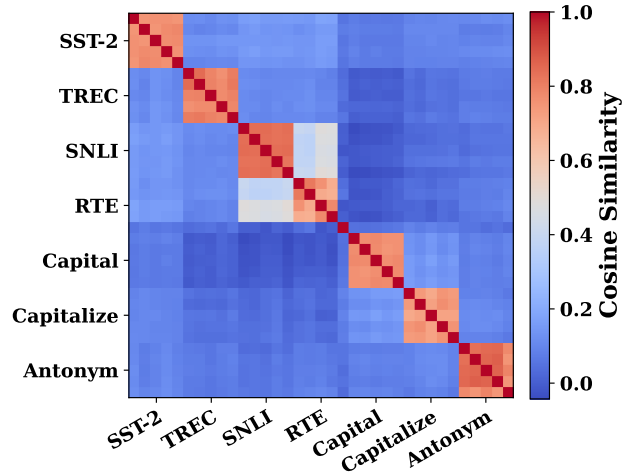


Figure 28: Cosine-similarity heatmap of LTVs trained for seven tasks on Llama2-7B.



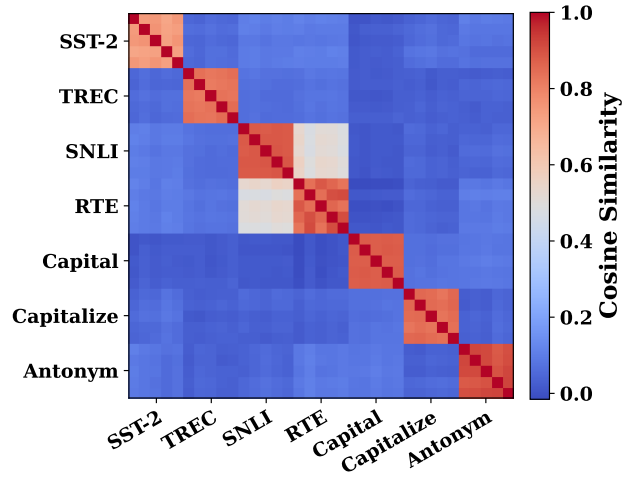


Figure 29: Cosine-similarity heatmap of LTVs trained for seven tasks on Llama2-13B.

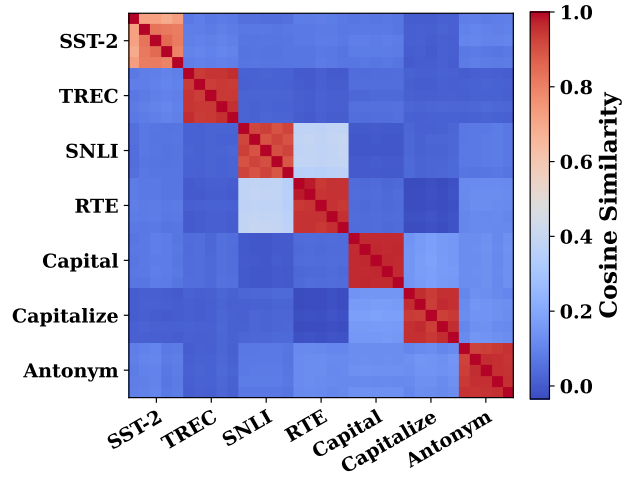


Figure 30: Cosine-similarity heatmap of LTVs trained for seven tasks on Qwen2.5-32B.

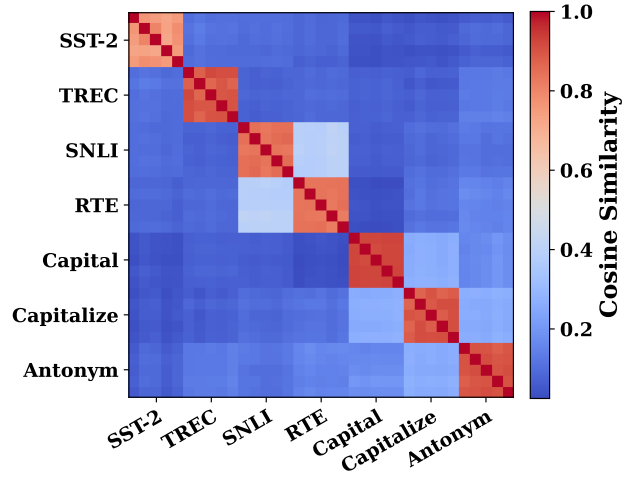


Figure 31: Cosine-similarity heatmap of LTVs trained for seven tasks on Yi-34B.

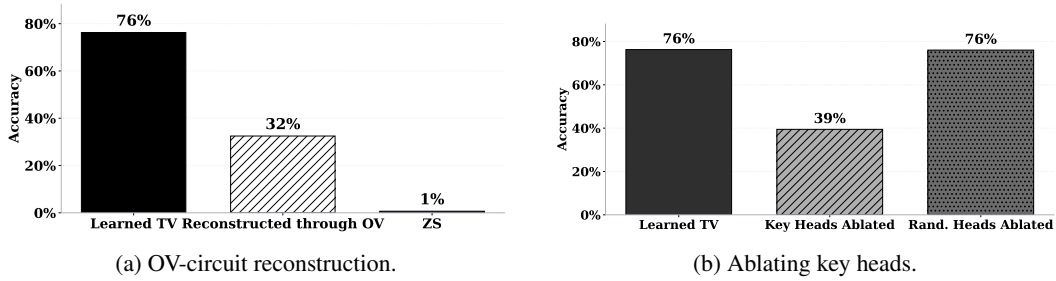


Figure 32: Attention heads and TV on Llama3-8B: OV-circuit reconstruction (left) and ablation of key heads (right).

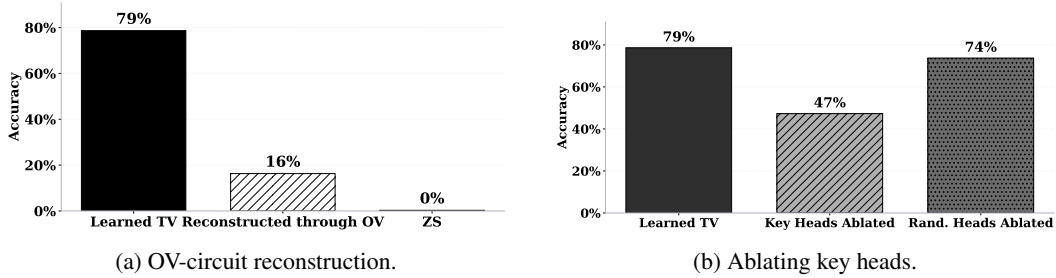


Figure 33: Attention heads and TV on Llama3.2-3B: OV-circuit reconstruction (left) and ablation of key heads (right).

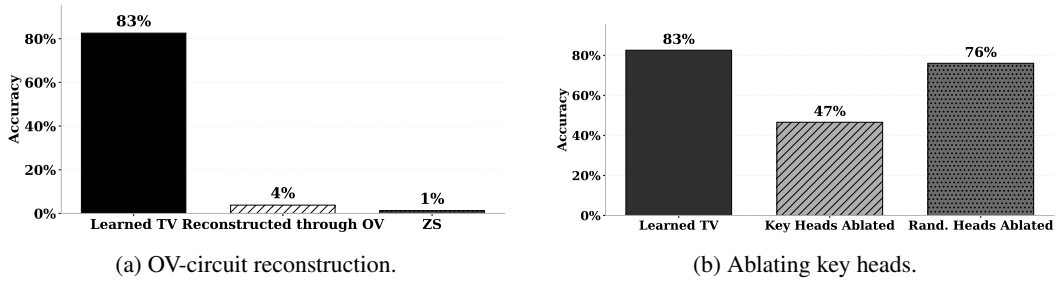


Figure 34: Attention heads and TV on Llama2-7B: OV-circuit reconstruction (left) and ablation of key heads (right).

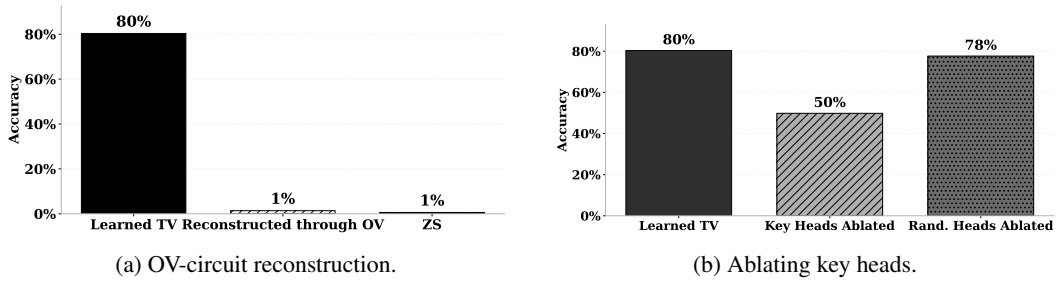


Figure 35: Attention heads and TV on Llama2-13B: OV-circuit reconstruction (left) and ablation of key heads (right).

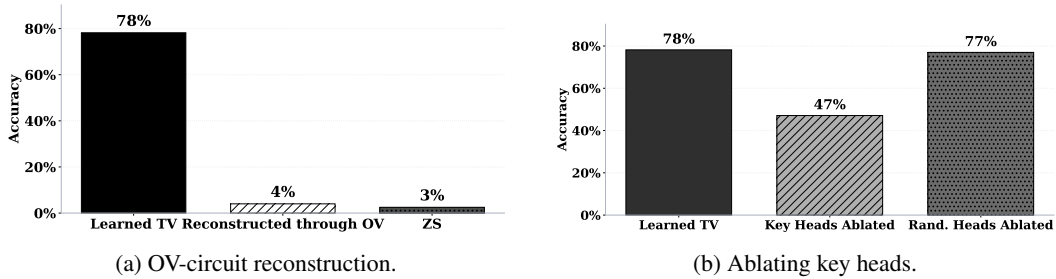


Figure 36: Attention heads and TV on Llama3-70B: OV-circuit reconstruction (left) and ablation of key heads (right).

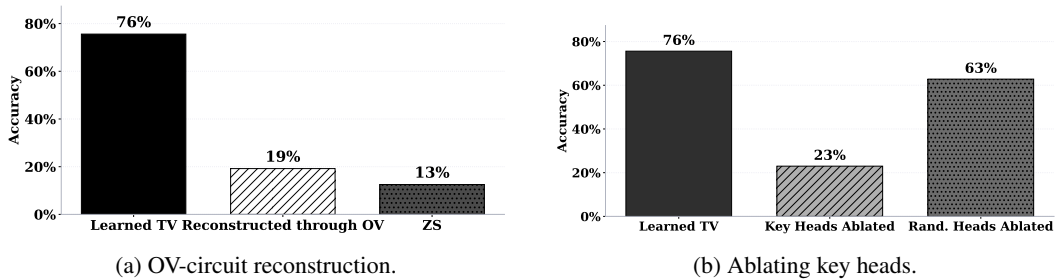


Figure 37: Attention heads and TV on Qwen2.5-32B: OV-circuit reconstruction (left) and ablation of key heads (right).

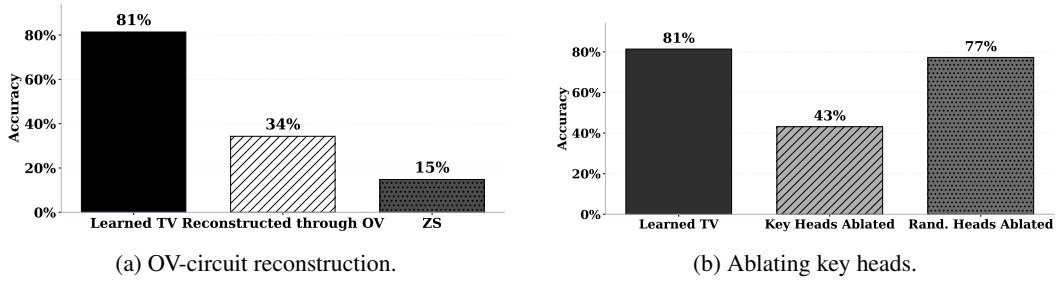


Figure 38: Attention heads and TV on Yi-34B: OV-circuit reconstruction (left) and ablation of key heads (right).

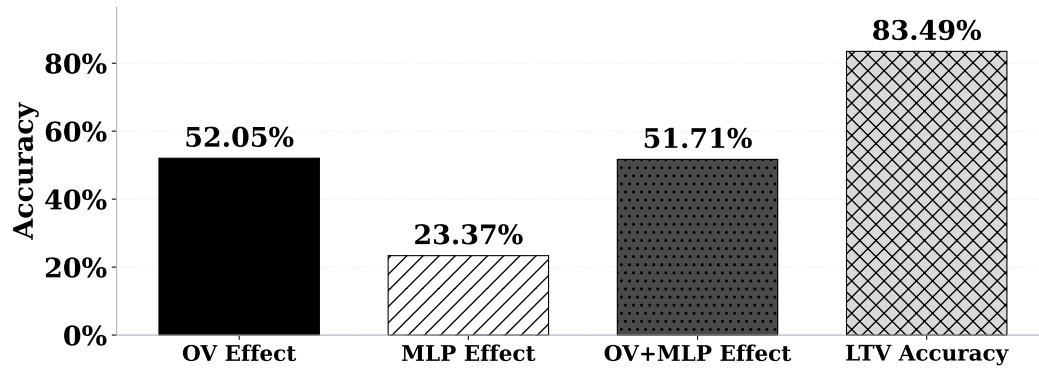


Figure 39: Effects of the MLP-based construction and MLP&OV-based reconstruction compared to the effect of OV-based reconstruction of TV effect.

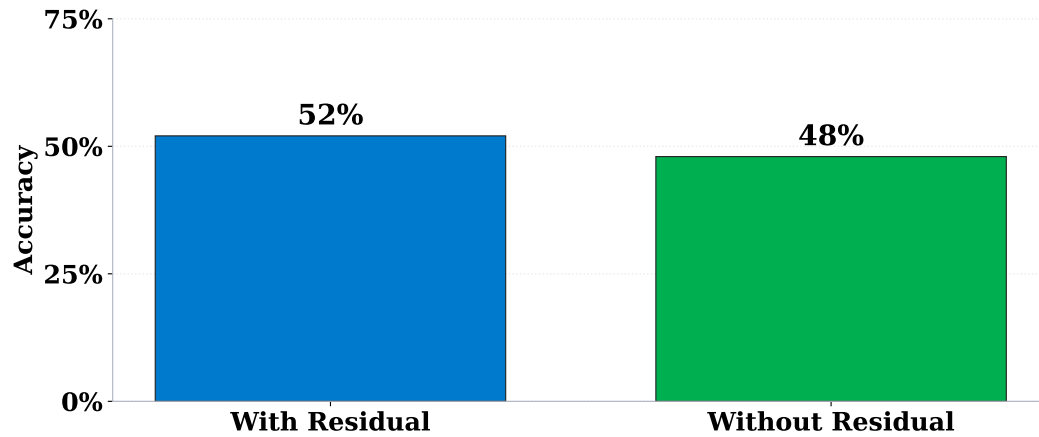


Figure 40: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Llama 3.1-8B.

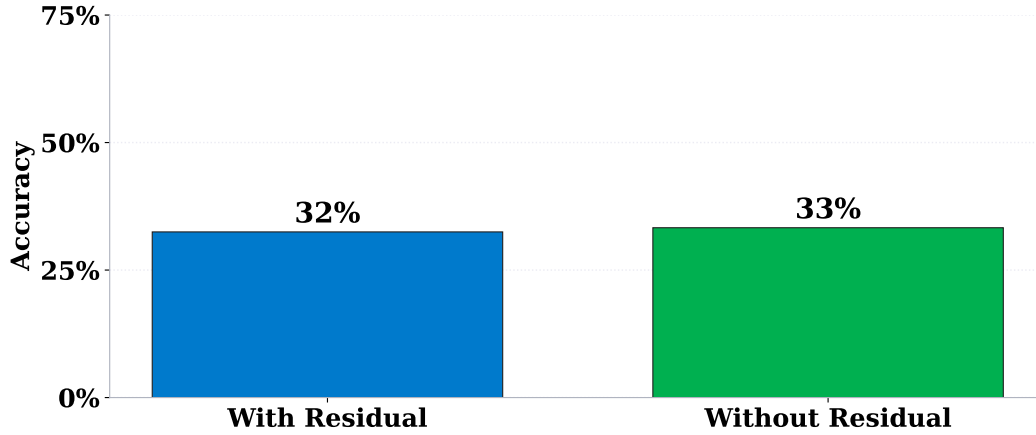


Figure 41: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Llama 3-8B.

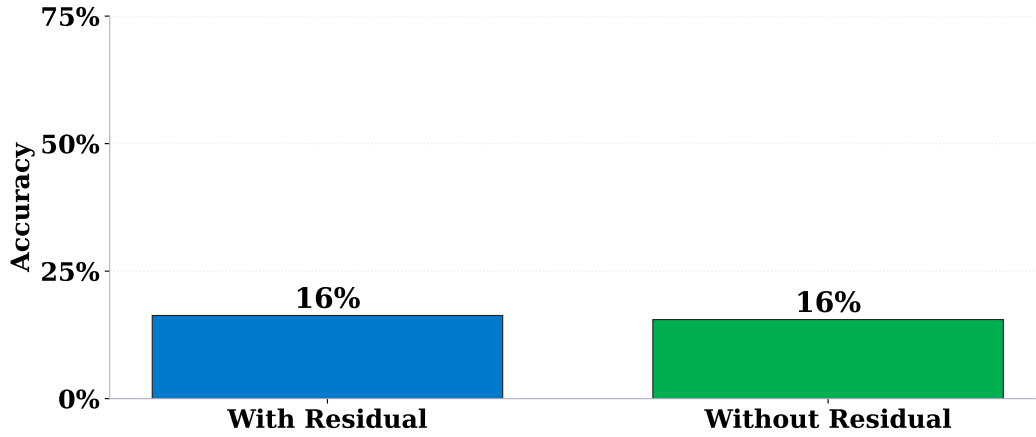


Figure 42: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Llama 3.2-3B.

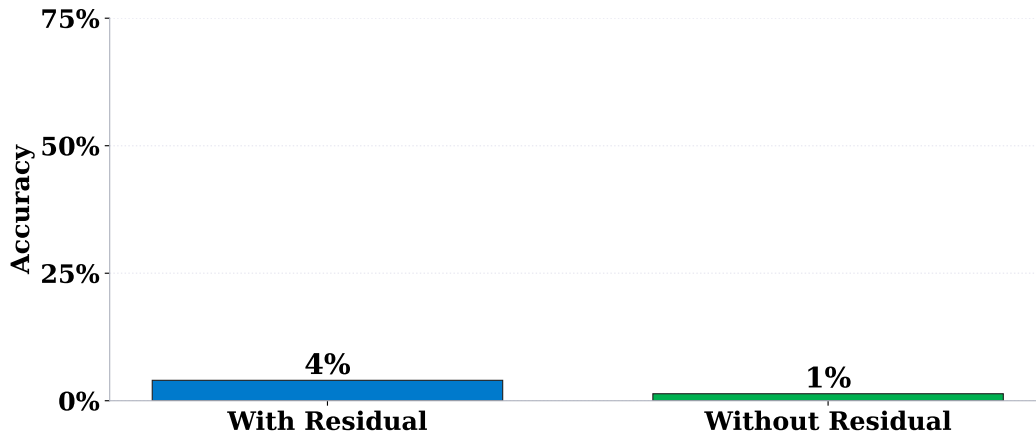


Figure 43: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Llama 3-70B.



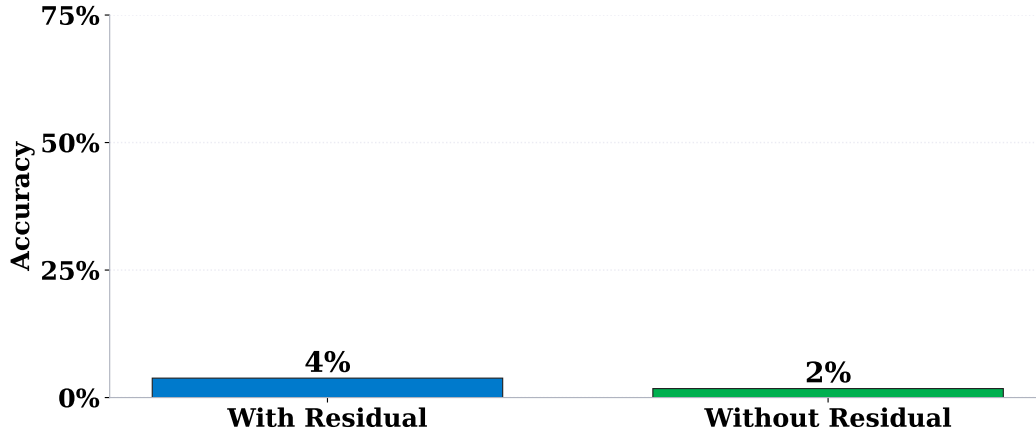


Figure 44: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Llama 2-7B.

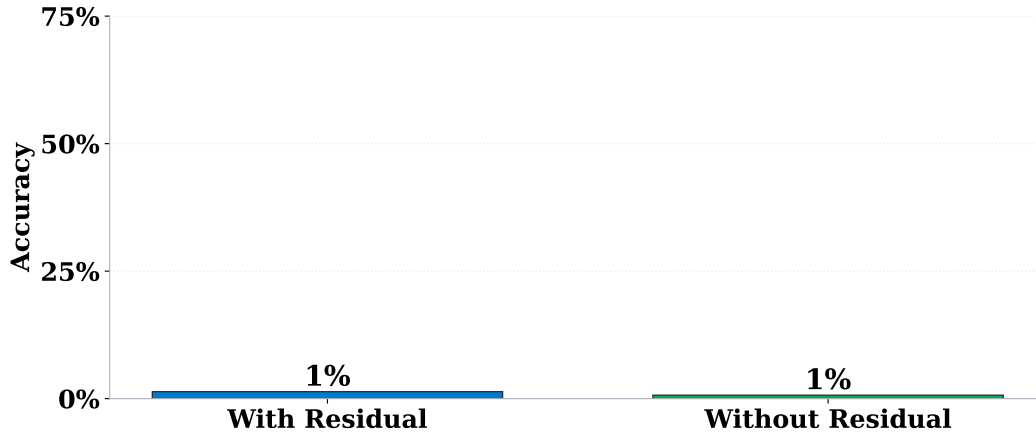


Figure 45: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Llama 2-13B.

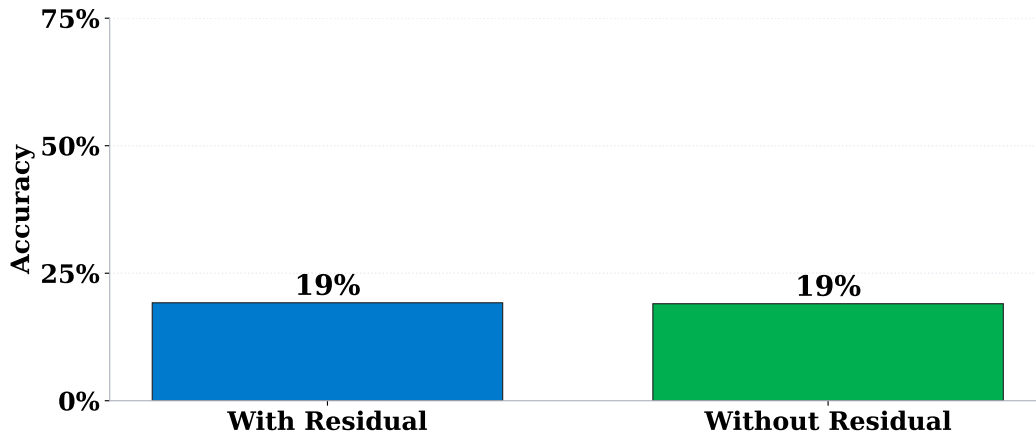


Figure 46: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Qwen2.5-32B.

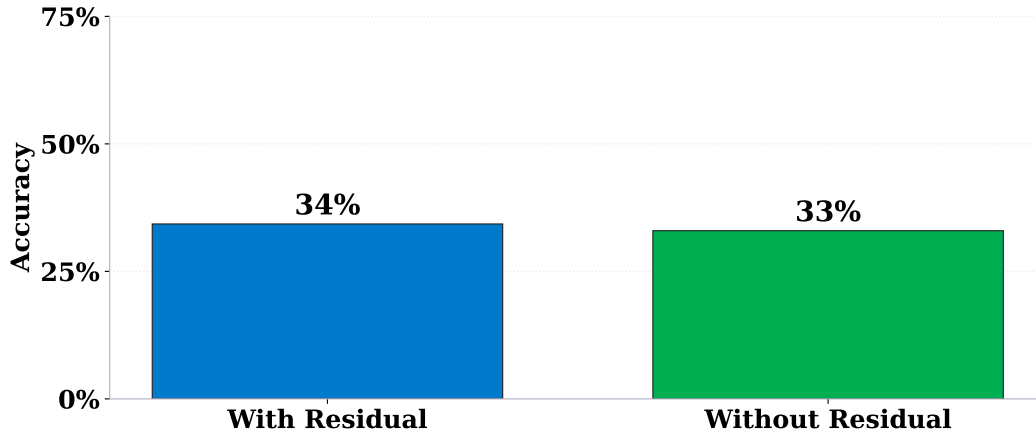


Figure 47: Effects of the OV circuit reconstruction with or without the TV added to the final layer: Yi-34B.

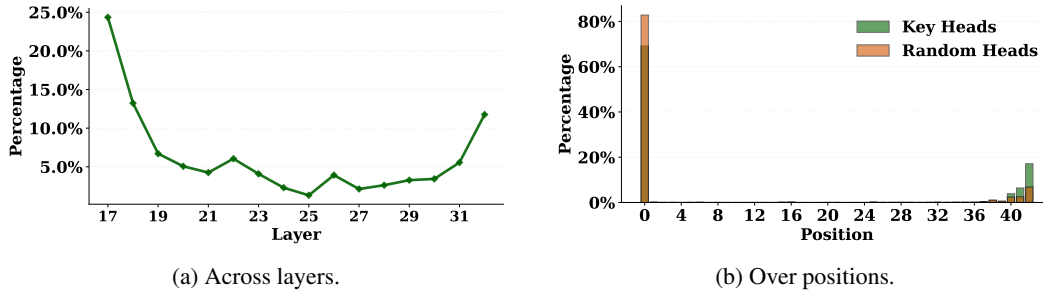


Figure 48: Key attention heads on Llama3-8B: distribution across layers (left) and attention over token positions (right).

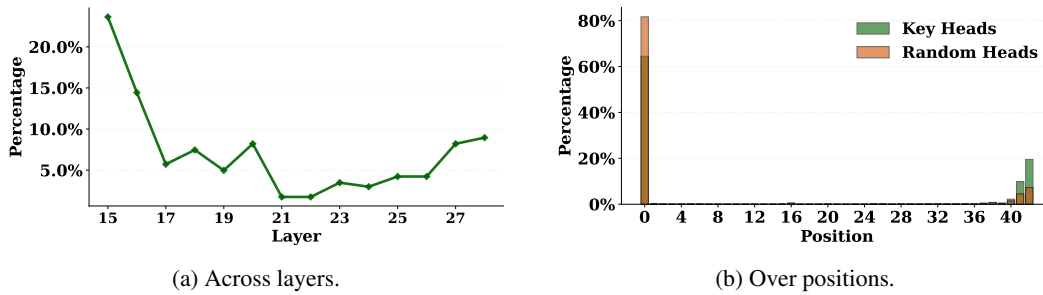


Figure 49: Key attention heads on Llama3.2-3B: distribution across layers (left) and attention over token positions (right).

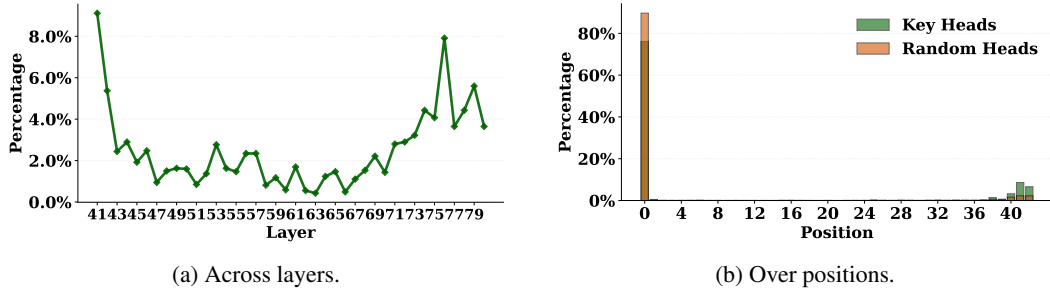


Figure 50: Key attention heads on Llama3-70B: distribution across layers (left) and attention over token positions (right).

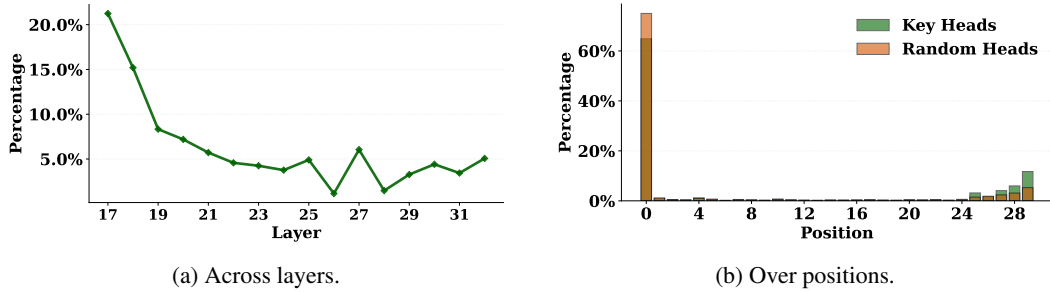


Figure 51: Key attention heads on Llama2-7B: distribution across layers (left) and attention over token positions (right).

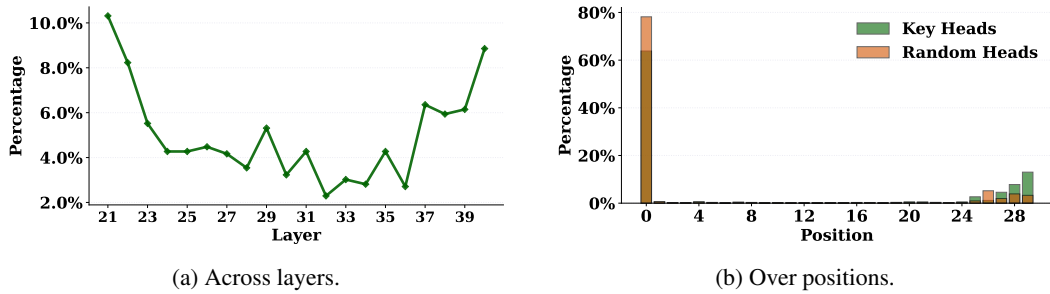


Figure 52: Key attention heads on Llama2-13B: distribution across layers (left) and attention over token positions (right).

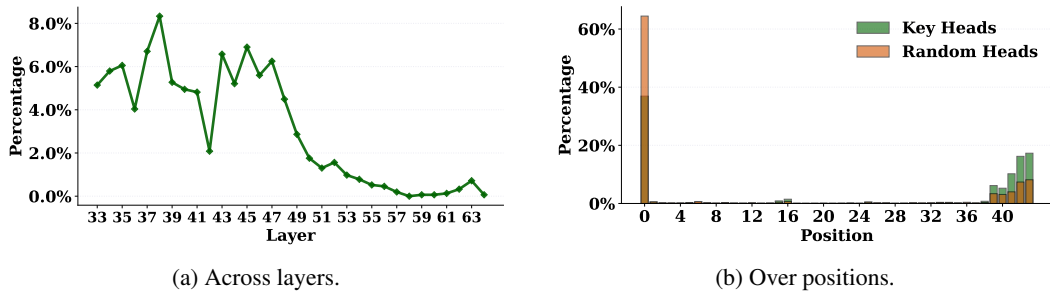


Figure 53: Key attention heads on Qwen2.5-32B: distribution across layers (left) and attention over token positions (right).

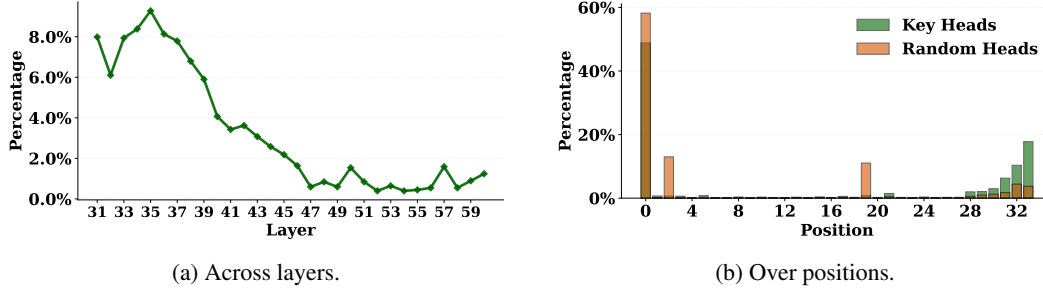


Figure 54: Key attention heads on Yi-34B: distribution across layers (left) and attention over token positions (right).

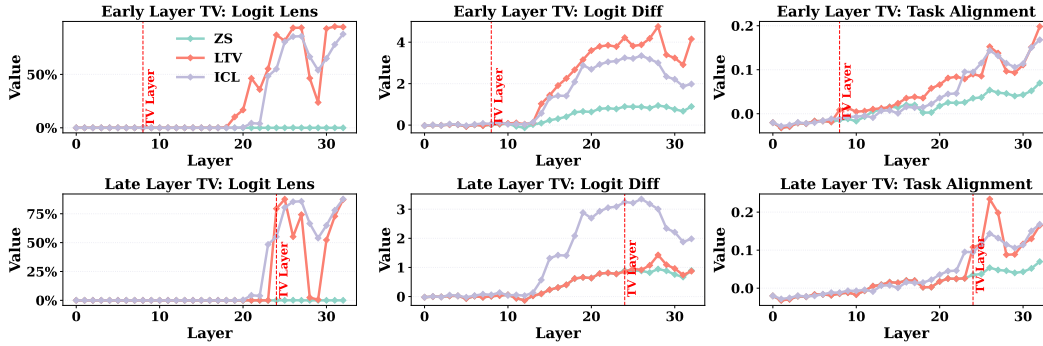


Figure 55: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at an early vs. late layer.

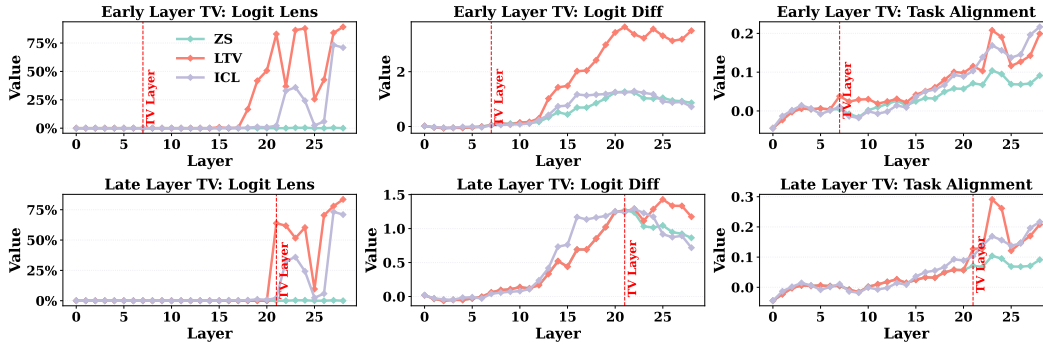


Figure 56: Metrics across layers on Llama3.2-3B when the TV is injected into the hidden state at an early vs. late layer.

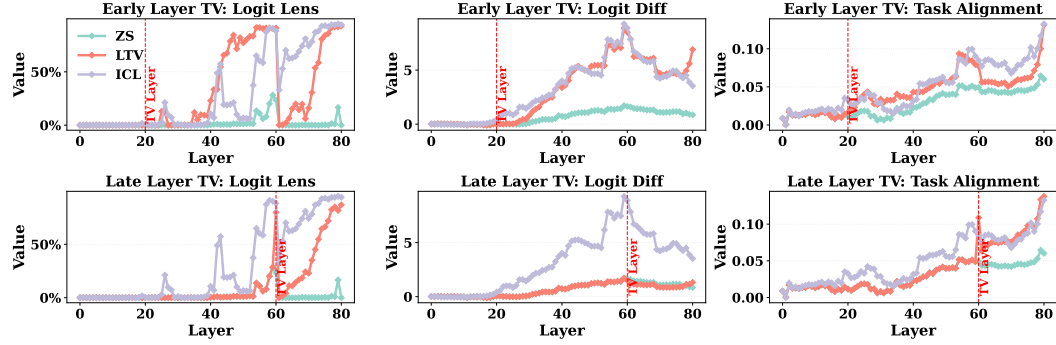


Figure 57: Metrics across layers on Llama3-70B when the TV is injected into the hidden state at an early vs. late layer.

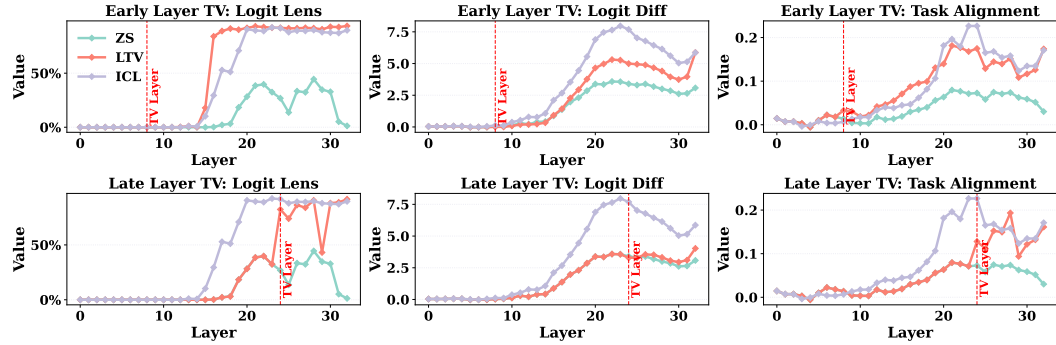


Figure 58: Metrics across layers on Llama2-7B when the TV is injected into the hidden state at an early vs. late layer.

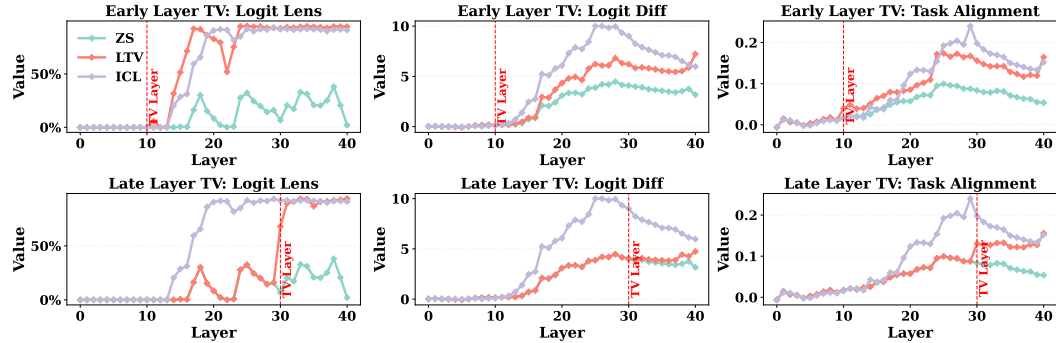


Figure 59: Metrics across layers on Llama2-13B when the TV is injected into the hidden state at an early vs. late layer.



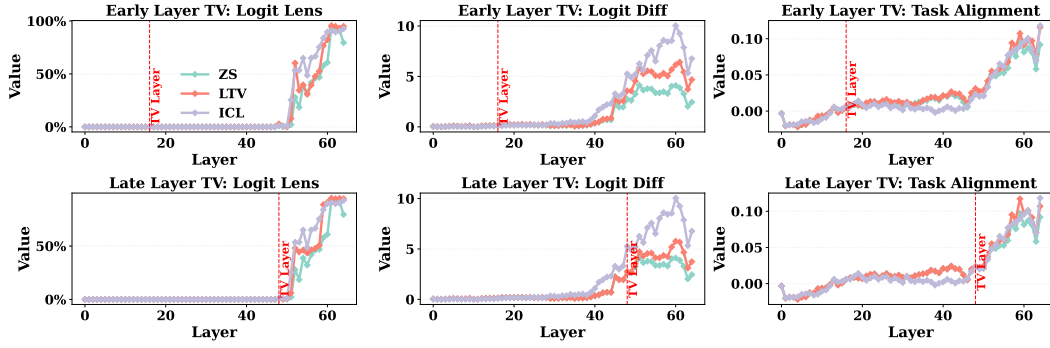


Figure 60: Metrics across layers on Qwen2.5-32B when the TV is injected into the hidden state at an early vs. late layer.

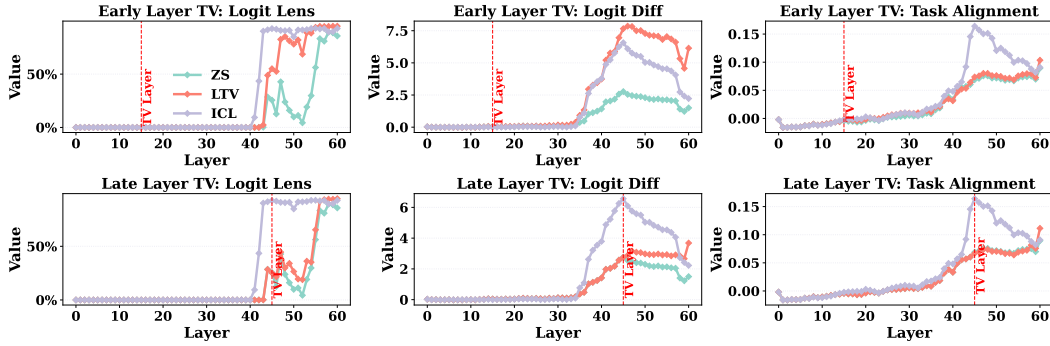
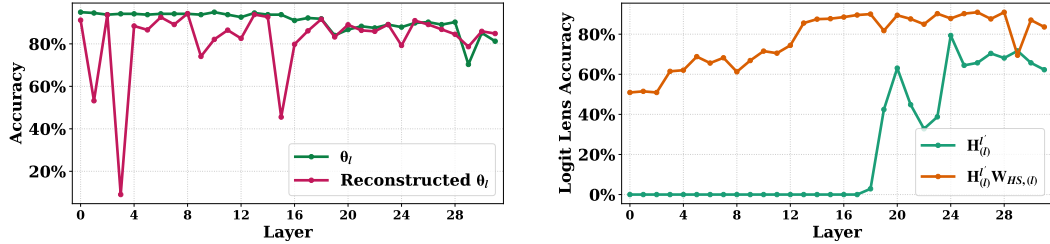


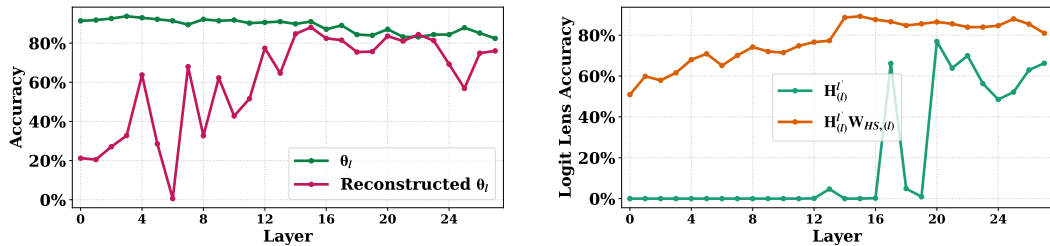
Figure 61: Metrics across layers on Yi-34B when the TV is injected into the hidden state at an early vs. late layer.



(a) Reconstructed TV.

(b) Hidden-state surrogate.

Figure 62: Linear hypothesis on Llama3-8B: linearly reconstructed TV (left) and linear surrogate for hidden-state updates (right).



(a) Reconstructed TV.

(b) Hidden-state surrogate.

Figure 63: Linear hypothesis on Llama3.2-3B: linearly reconstructed TV (left) and linear surrogate for hidden-state updates (right).

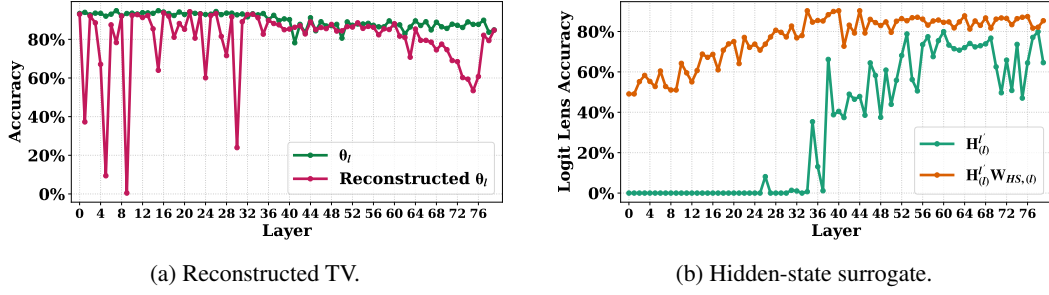


Figure 64: Linear hypothesis on Llama3-70B: linearly reconstructed TV (left) and linear surrogate for hidden-state updates (right).

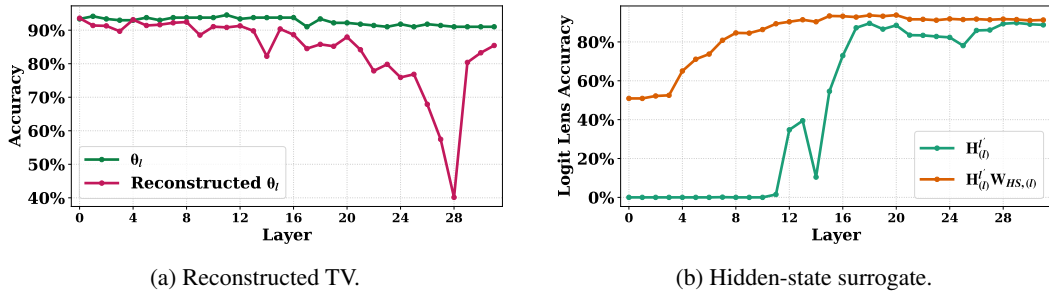


Figure 65: Linear hypothesis on Llama2-7B: linearly reconstructed TV (left) and linear surrogate for hidden-state updates (right).

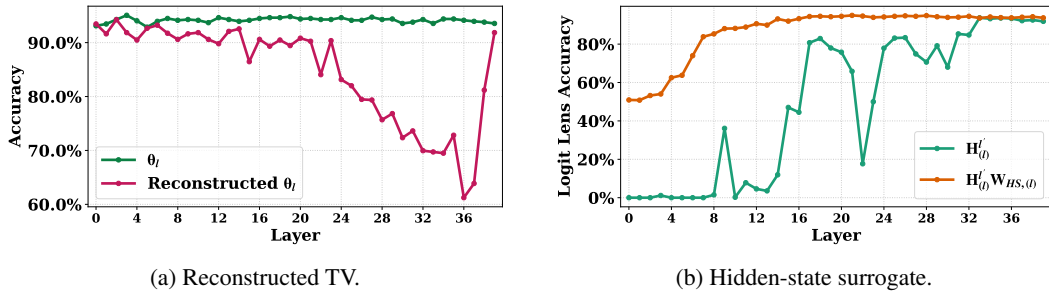


Figure 66: Linear hypothesis on Llama2-13B: linearly reconstructed TV (left) and linear surrogate for hidden-state updates (right).

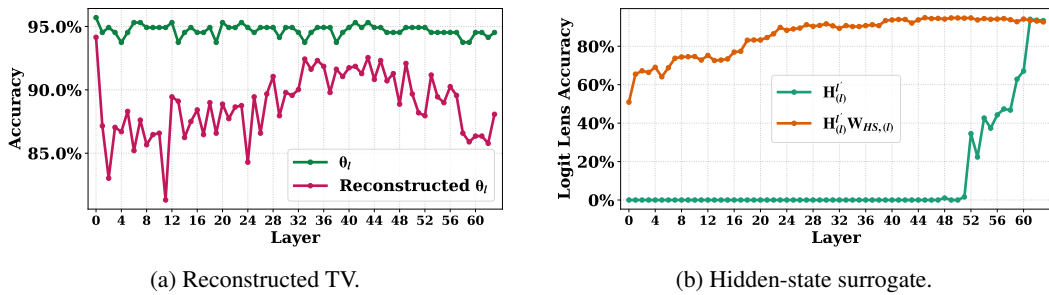


Figure 67: Linear hypothesis on Qwen2.5-32B: linearly reconstructed TV (left) and linear surrogate for hidden-state updates (right).

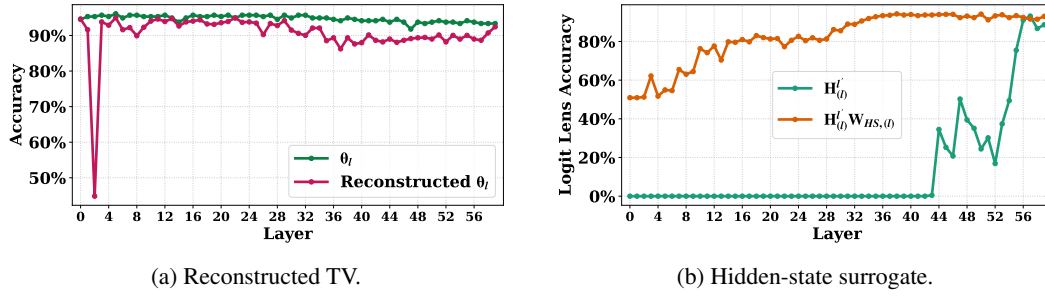


Figure 68: Linear hypothesis on Yi-34B: linearly reconstructed TV (left) and linear surrogate for hidden-state updates (right).

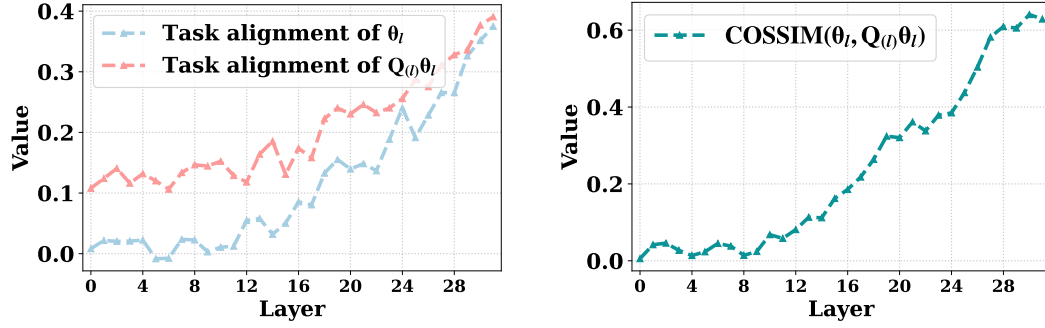


Figure 69: Rotation analysis on Llama3-8B: applying the fitted rotation  $Q_{(l)}$  to the TV increases task alignment (left); rotation strength vs. layer depth (right).

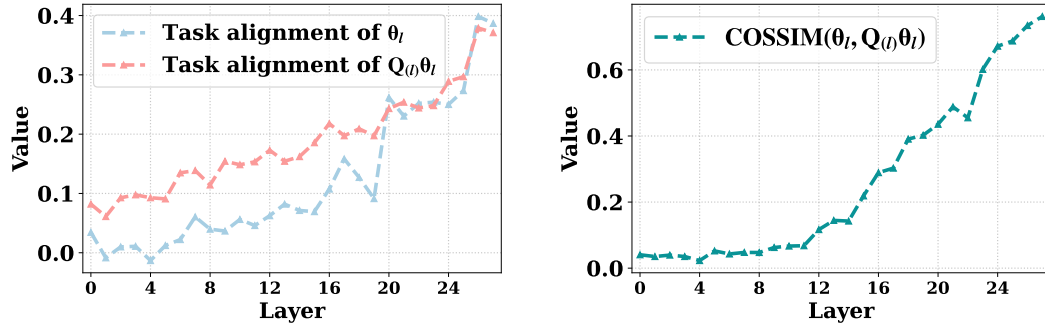


Figure 70: Rotation analysis on Llama3.2-3B: applying the fitted rotation  $Q_{(l)}$  to the TV increases task alignment (left); rotation strength vs. layer depth (right).

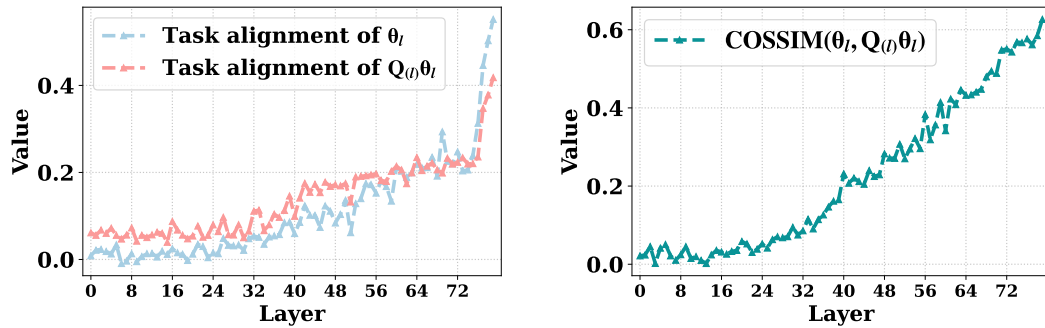


Figure 71: Rotation analysis on Llama3-70B: applying the fitted rotation  $Q_{(l)}$  to the TV increases task alignment (left); rotation strength vs. layer depth (right).

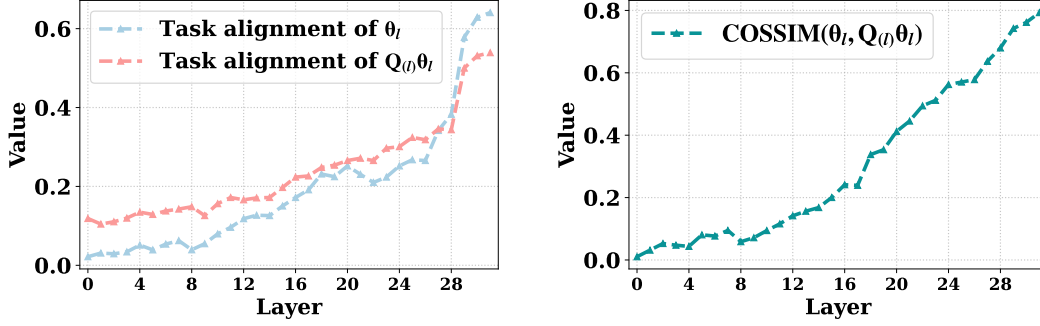


Figure 72: Rotation analysis on Llama2-7B: applying the fitted rotation  $Q_{(l)}$  to the TV increases task alignment (left); rotation strength vs. layer depth (right).

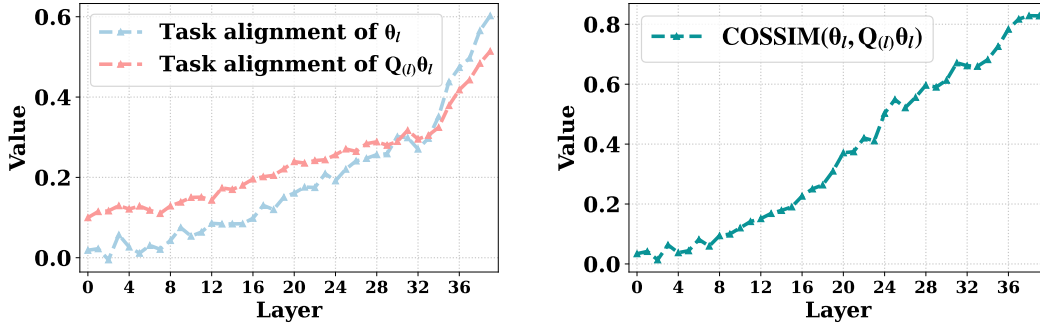


Figure 73: Rotation analysis on Llama2-13B: applying the fitted rotation  $Q_{(l)}$  to the TV increases task alignment (left); rotation strength vs. layer depth (right).

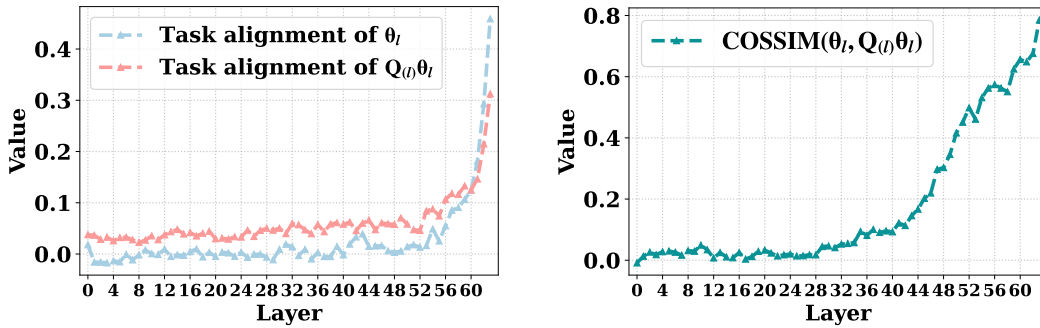


Figure 74: Rotation analysis on Qwen2.5-32B: applying the fitted rotation  $Q_{(l)}$  to the TV increases task alignment (left); rotation strength vs. layer depth (right).

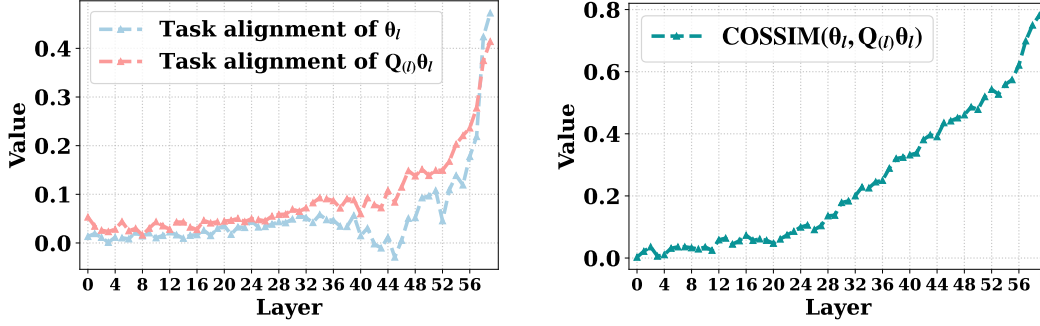


Figure 75: Rotation analysis on Yi-34B: applying the fitted rotation  $Q_{(l)}$  to the TV increases task alignment (left); rotation strength vs. layer depth (right).

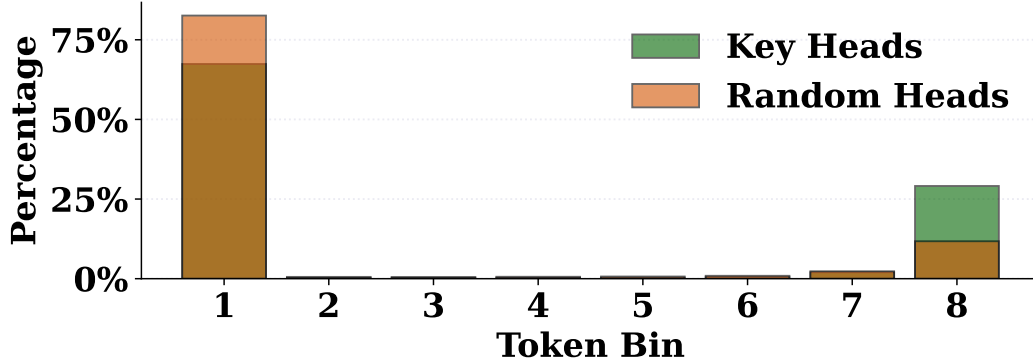


Figure 76: Average attention distribution of Llama3.1-8B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.

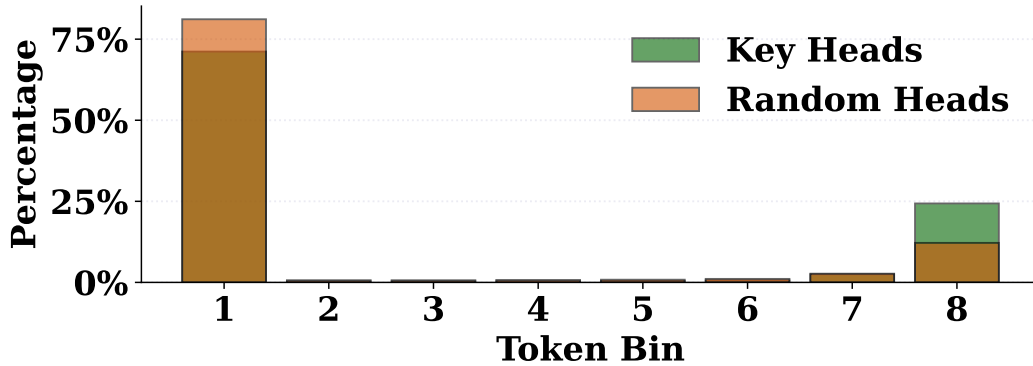


Figure 77: Average attention distribution of Llama3-8B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.

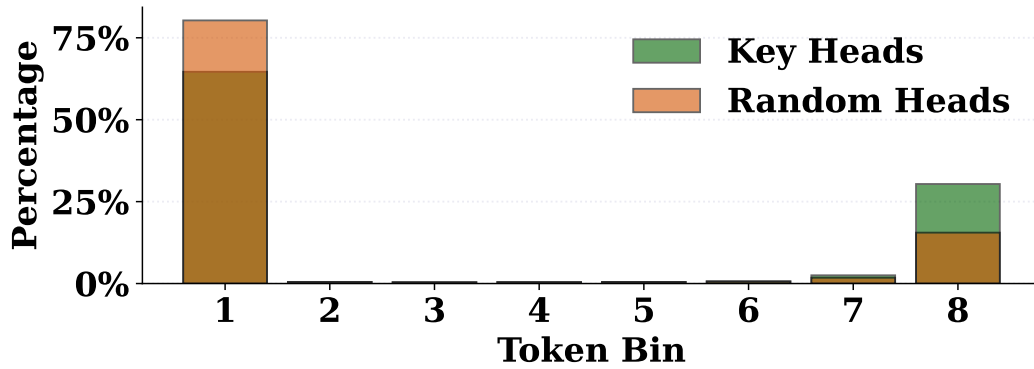


Figure 78: Average attention distribution of Llama3.2-3B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.

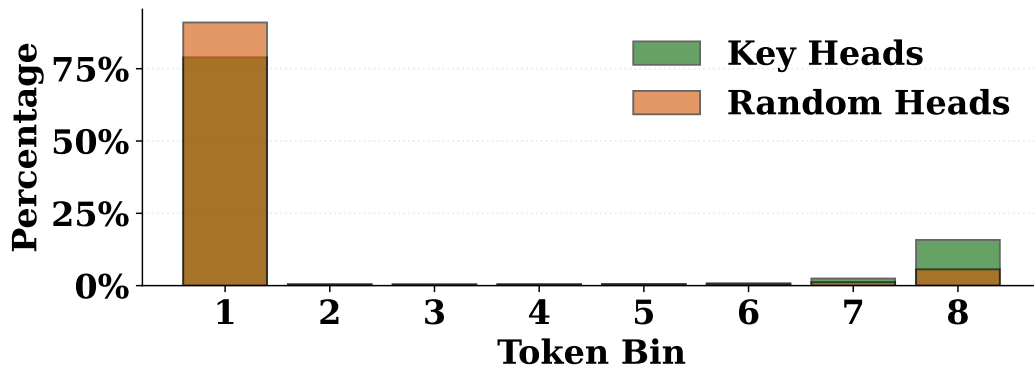


Figure 79: Average attention distribution of Llama3-70B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.

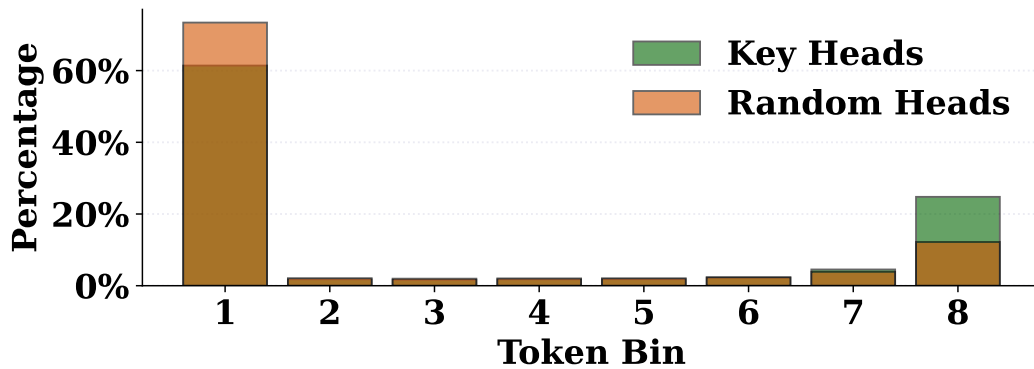


Figure 80: Average attention distribution of Llama2-7B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.



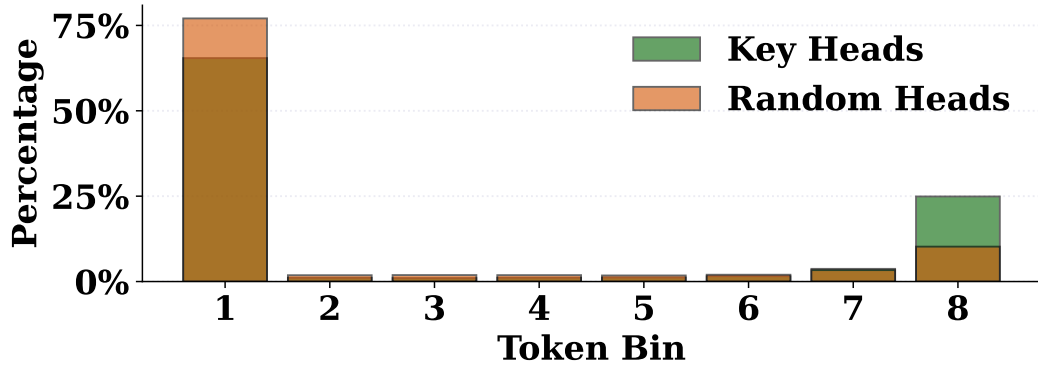


Figure 81: Average attention distribution of Llama2-13B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.

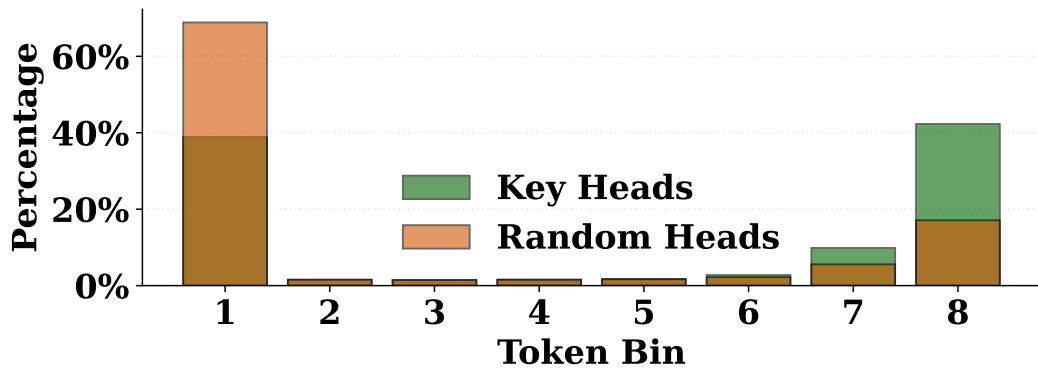


Figure 82: Average attention distribution of Qwen2.5-32B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.

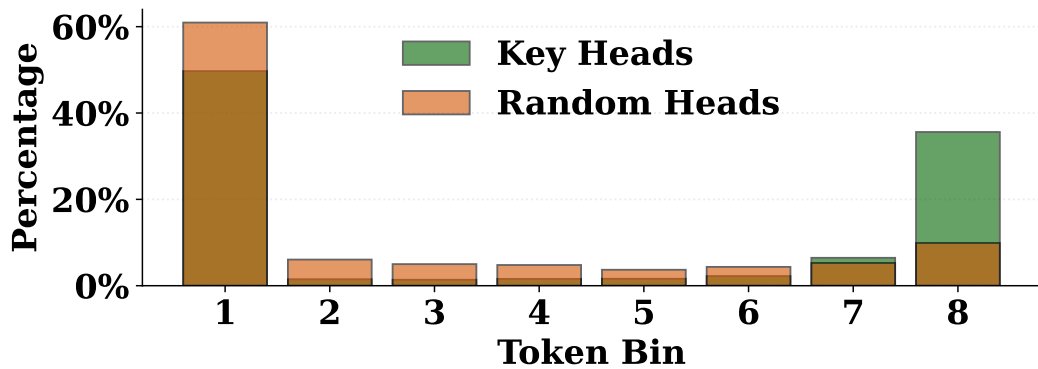


Figure 83: Average attention distribution of Yi-34B on SST-2: proportions of attention weights assigned to 8 tokens intervals each comprising  $\frac{1}{8}$  of all tokens.

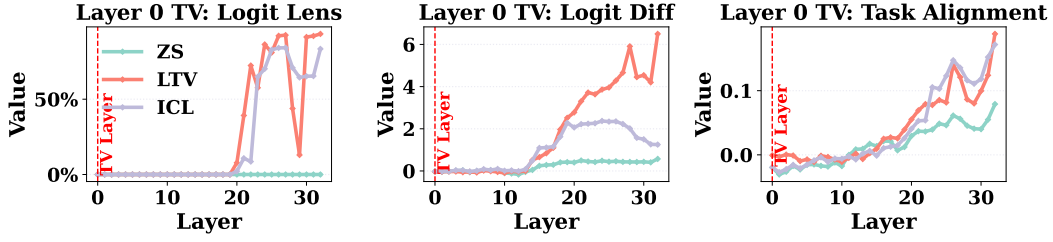


Figure 84: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 0.

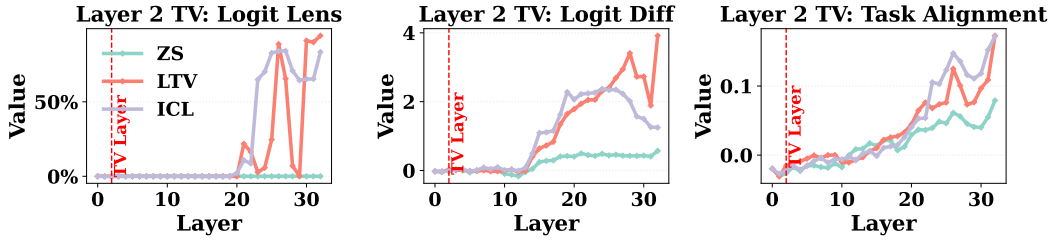


Figure 85: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 2.

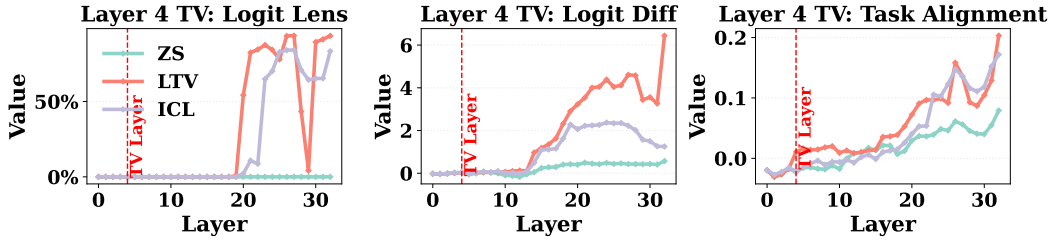


Figure 86: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 4.

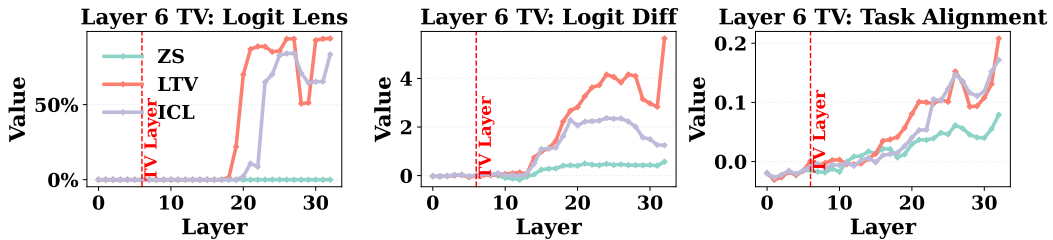


Figure 87: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 6.

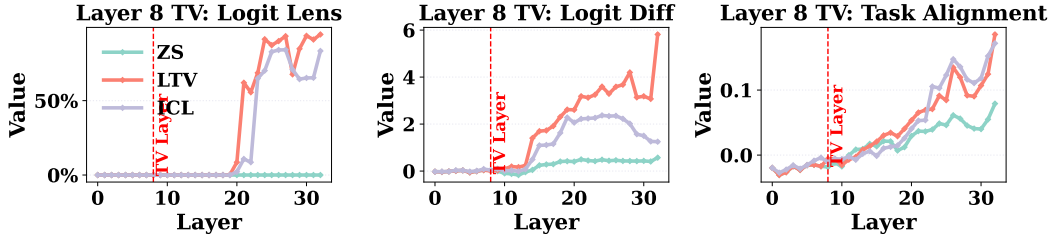


Figure 88: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 8.

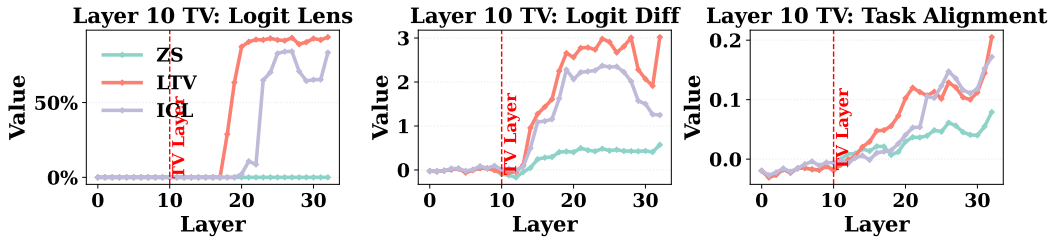


Figure 89: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 10.

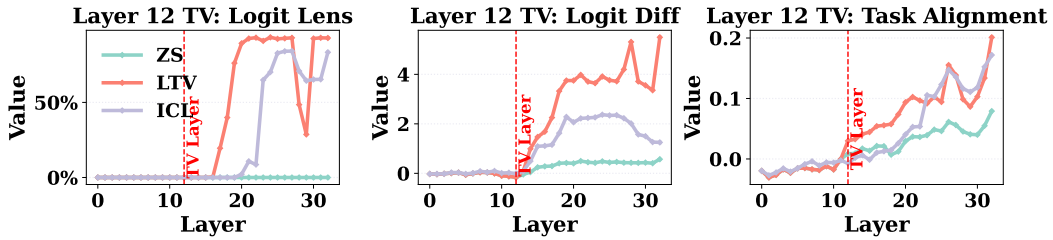


Figure 90: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 12.

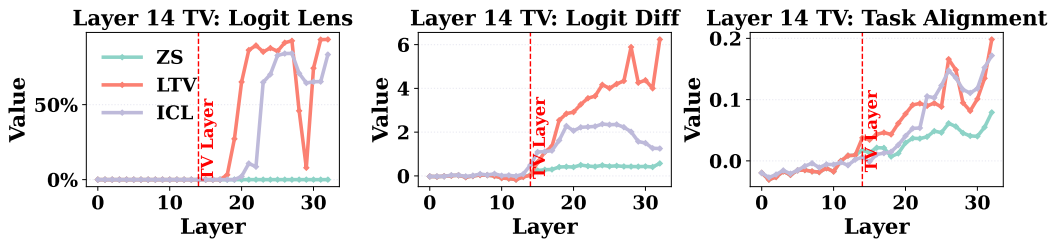


Figure 91: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 14.

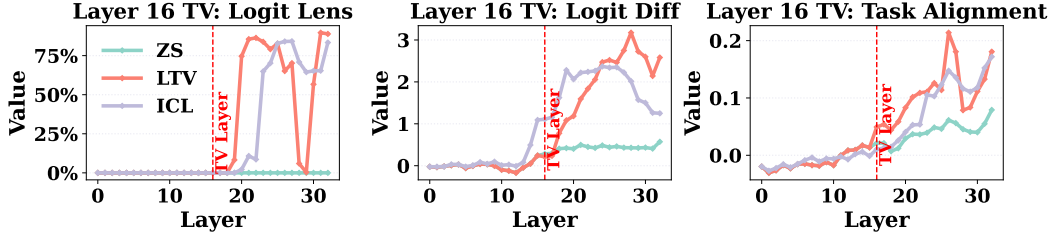


Figure 92: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 16.

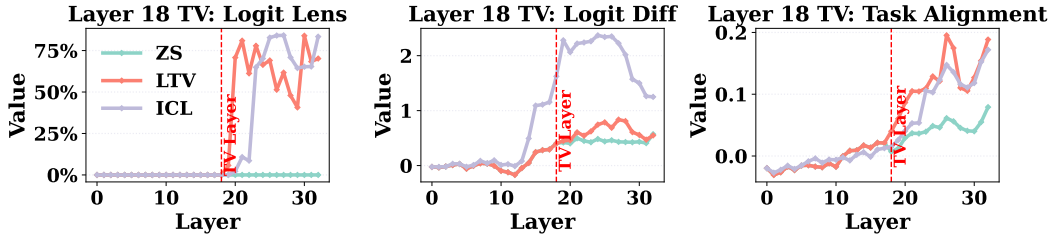


Figure 93: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 18.

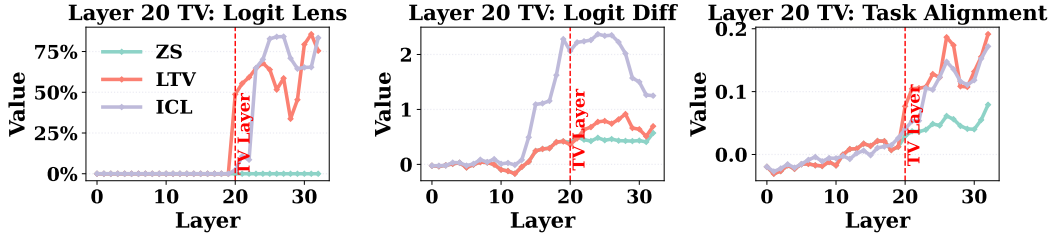


Figure 94: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 20.

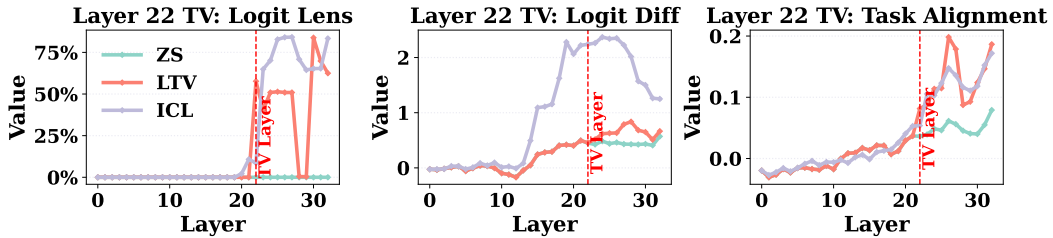


Figure 95: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 22.

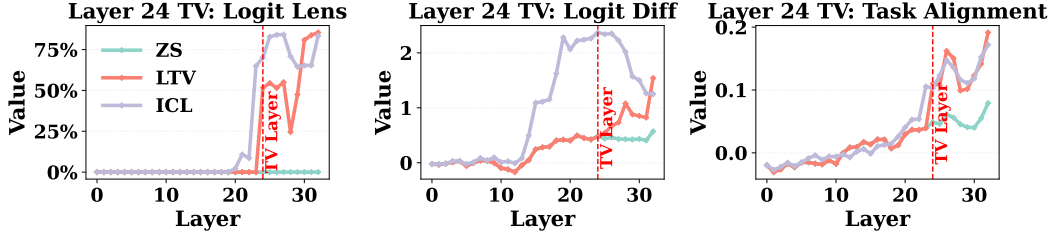


Figure 96: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 24.

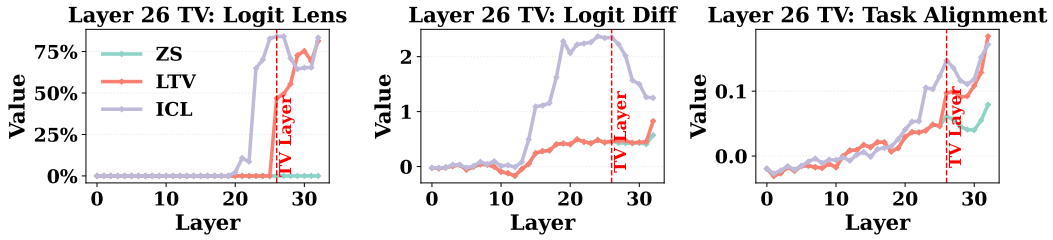


Figure 97: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 26.

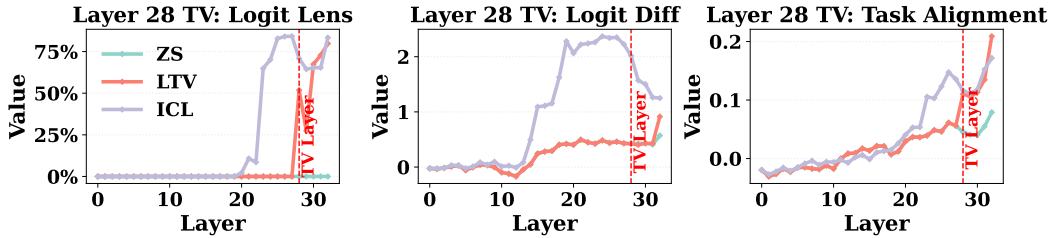


Figure 98: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 28.

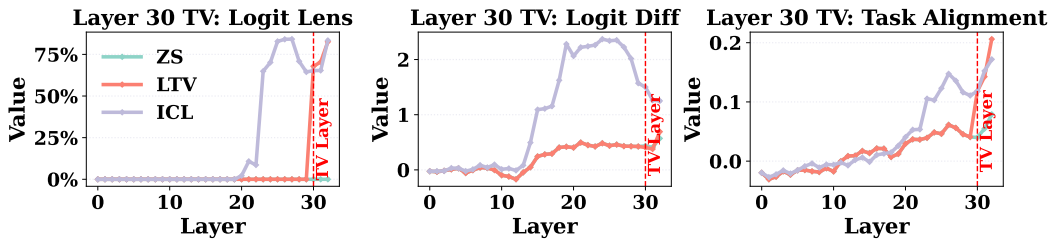


Figure 99: Metrics across layers on Llama3-8B when the TV is injected into the hidden state at layer 30.

Table 5: Comparison of Representation Methods

State Vector	I2CL	LTV (Ours)
91.31	66.04	<b>58.07</b>
Pearson correlation coefficient		p-value
-0.4914		0.0042

Table 6: Correlation strength between accuracy of reconstructed TV and the relative estimated logit effect difference

Table 7: Comparison of LTV, State Vector, and I2CL in terms of the time (seconds) required to complete the entire training and evaluation procedures.

State Vector	I2CL	LTV (Ours)
91.31	66.04	<b>58.07</b>

Model	ZS Accuracy	ICL Accuracy	Accuracy with LTV
Llama3-70B	2.51%	81.93%	78.18%
Qwen2.5-32B	12.52%	85.44%	75.59%
Yi-34B	14.82%	81.33%	81.37%

Table 8: Performance of LTVs under the traditional setting (injecting into the last-token hidden state at a single layer). Injection layers correspond to 50% of each model’s total depth.

Table 9: Comparison of LTV vs. FV and Vanilla TV across five scenarios on Llama2-7B.

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{16\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
Vanilla TV	38.26%	1.96%	14.16%	18.85%	13.30%	52.82%
FV	51.81%	1.40%	28.60%	47.14%	20.44%	73.23%
LTV	82.54% $\uparrow_{30.73\%}$	79.34% $\uparrow_{77.38\%}$	84.60% $\uparrow_{56.00\%}$	82.24% $\uparrow_{35.10\%}$	51.60% $\uparrow_{31.16\%}$	85.16% $\uparrow_{11.93\%}$

Table 10: Comparison of LTV vs. FV and Vanilla TV across five scenarios on Llama2-13B.

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{20\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
Vanilla TV	27.67%	1.84%	16.42%	20.46%	16.07%	43.84%
FV	41.59%	1.22%	42.25%	36.97%	24.74%	77.51%
LTV	80.33% $\uparrow_{38.74\%}$	71.53% $\uparrow_{69.69\%}$	87.69% $\uparrow_{45.44\%}$	82.25% $\uparrow_{45.28\%}$	51.46% $\uparrow_{26.72\%}$	84.99% $\uparrow_{7.48\%}$

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{16\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
Vanilla TV	31.69%	2.02%	1.05%	26.68%	0.33%	75.83%
FV	33.28%	2.93%	18.38%	16.95%	17.72%	53.93%
LTV	76.26% $\uparrow_{42.98\%}$	76.22% $\uparrow_{73.29\%}$	77.93% $\uparrow_{59.55\%}$	83.48% $\uparrow_{56.80\%}$	44.82% $\uparrow_{27.10\%}$	84.51% $\uparrow_{18.68\%}$

Table 11: Comparison of LTV vs. FV and Vanilla TV across five scenarios on Llama3-8B.

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{14\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
Vanilla TV	42.61%	3.07%	18.73%	37.05%	11.33%	65.38%
FV	19.54%	3.53%	15.07%	4.69%	13.26%	62.12%
LTV	78.65% $\uparrow_{36.04\%}$	74.10% $\uparrow_{70.57\%}$	78.18% $\uparrow_{59.45\%}$	80.43% $\uparrow_{43.38\%}$	46.38% $\uparrow_{33.12\%}$	82.80% $\uparrow_{17.42\%}$

Table 12: Comparison of LTV vs. FV and Vanilla TV across five scenarios on Llama3.2-3B.

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{40\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
LTV	78.18%	75.34%	76.13%	75.59%	48.75%	88.40%

Table 13: Performance of LTV across settings on Llama3-70B.

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{32\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
LTV	75.59%	36.04%	75.20%	87.24%	53.30%	87.08%

Table 14: Performance of LTV across settings on Qwen2.5-32B.

Method	Baseline $\mathbb{P} = \{-1\}, \mathbb{L} = \{30\}$	1) Diff. Pos. $\mathbb{P} = \{4\}$	2) More Pos. $\mathbb{P} = \{-5, \dots, -1\}$	3) More layers $\mathbb{L} = \{0, 4, 8, \dots\}$	4) More layers & Pos. $\mathbb{P} = \{-5, \dots\}, \mathbb{L} = \{0, 4, \dots\}$	5) ICL prompts
LTV	81.37%	73.53%	84.39%	82.47%	51.29%	89.69%

Table 15: Performance of LTV across settings on Yi-34B.

Table 16: Performance of LTV while injecting to multiple layers and positions simultaneously with different layer strides

	$\mathbb{P} = \{-1\}$	$\mathbb{P} = \{-5, \dots\}$
<b>Layer Stride = 2</b> $\mathbb{L} = \{0, 2, 4, \dots\}$	82.40%	51.08%
<b>Layer Stride = 4</b> $\mathbb{L} = \{0, 4, 8, \dots\}$	86.43%	51.39%
<b>Layer Stride = 8</b> $\mathbb{L} = \{0, 8, 16, \dots\}$	88.50%	50.47%

Table 17: Applying the Capital LTV to other tasks. The LTV yields no substantial accuracy improvements in any case because the Capital dataset does not share its label space with any of the other datasets.

SST-2	TREC	RTE	SNLI	Capitalize	Antonym
0.00%	0.00%	9.39%	5.06%	2.33%	0.00%

Table 18: Top-10 tokens decoded from early- and late-layer TVs on Llama3-8B.

Layer	Decoded Tokens
Early Layer (8)	tring, CCA, erk, bart, uge, ensor, , テル, a3a, emer
Late Layer (24)	positive, negative, positive, Positive, Negative, negative, Negative, Positive, _positive, -negative

Table 19: Top-10 tokens decoded from early- and late-layer TVs on Llama3.2-3B.

Layer	Decoded Tokens
Early Layer (7)	ync, flip, stress, hope, haven, Lor, negative, ugi, stressed, hab
Late Layer (21)	positive, positive, negative, -positive, Positive, Positive, negative, -positives, negative, Negative

Table 20: Top-10 tokens decoded from early- and late-layer TVs on Llama3-70B.

Layer	Decoded Tokens
Early Layer (20)	EventData, esteem, Я, 𐄂, spath, hores, raya, idth, , _priv
Late Layer (60)	negative, negative, Negative, positive, Negative, -negative, positive, Positive, Positive, _negative



Table 21: Top-10 tokens decoded from early- and late-layer TVs on Llama2-7B.

Layer	Decoded Tokens
Early Layer (8)	bah, arith, arna, revers, feder, HOST, BIT, Pat, orr, IP
Late Layer (24)	positive, negative, negative, posit, pos, Pos, neg, Pos, Neg, poz

Table 22: Top-10 tokens decoded from early- and late-layer TVs on Llama2-13B.

Layer	Decoded Tokens
Early Layer (8)	negative, bin, ed, agg, electric, myself, eda, hed, isser, positive
Late Layer (24)	negative, negative, positive, Neg, neg, neg, отри, pos, Pos, negro

Table 23: Top-10 tokens decoded from early- and late-layer TVs on Qwen2.5-32B.

Layer	Decoded Tokens
Early Layer (16)	fd, Reverse, inverted, Trait, ocale, Hack, ic, Traits, Aware, 逆转
Late Layer (48)	. constraint, registrations, 魏, 传奇, 看点, (SE, ApplicationContext, Offensive, 产量, 浓缩

Table 24: Top-10 tokens decoded from early- and late-layer TVs on Yi-34B.

Layer	Decoded Tokens
Early Layer (15)	一分, iency, , shit, oc, , orating, 正能量, Gap, unbiased
Late Layer (45)	Mpc, elf, izza, Parish, 炳, 莫, nexper, 流行的, 增长率, rst