LOCALIZING TASK RECOGNITION AND TASK LEARNING IN IN-CONTEXT LEARNING VIA ATTENTION HEAD ANALYSIS

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027 028 029

031

033

034

036

037

040

041

042

043

044

045

046

047

048

049

050

051

052

ABSTRACT

We investigate the mechanistic underpinnings of in-context learning (ICL) in large language models by reconciling two dominant perspectives: the componentlevel analysis of attention heads and the holistic decomposition of ICL into Task Recognition (TR) and Task Learning (TL). We propose a novel framework based on Task Subspace Logit Attribution (TSLA) to identify attention heads specialized in TR and TL, and demonstrate their distinct yet complementary roles. Through correlation analysis, ablation studies, and input perturbations, we show that the identified TR and TL heads independently and effectively capture the TR and TL components of ICL. Via steering experiments with geometric analysis of hidden states, we reveal that TR heads promote task recognition by aligning hidden states with the task subspace, while TL heads rotate hidden states within the subspace toward the correct label to facilitate prediction. We further show how previous findings on ICL's mechanism—including induction heads, task vectors, and more—can be reconciled with our attention-head-level analysis of the TR-TL decomposition. Our framework thus provides a unified and interpretable account of how LLMs execute ICL across diverse tasks and settings¹.

1 Introduction

A key property of Large Language Models (LLMs) is their ability to solve tasks from demonstrations embedded in the input—without further training. This phenomenon, known as In-context Learning (ICL) (Brown et al., 2020; Radford et al., 2019), has reduced the need for large datasets and finetuning, enabling fast adaptation of LLMs to new tasks (Dong et al., 2024; Sun et al., 2022). Since its success cannot be explained by traditional gradient-based paradigms (Ren et al., 2024b), deciphering the mechanism behind ICL has become a central research question of great academic interest.

Two research paradigms dominate this pursuit. (1) The introspective paradigm designates internal model components or representations as critical drivers of ICL functionality. Pioneering works (Elhage et al., 2021; Olsson et al., 2022) formulate the output logits of Transformers as the sum of individual component outputs and highlight the significance of **Induction Heads (IHs)** in toy models, with follow-ups confirming their importance in larger models via ablation (Crosbie & Shutova, 2024; Halawi et al., 2024; Cho et al., 2025a). These studies inspired the concept of task vectors—compact representations distilled from hidden states or attention head outputs that steer zero-shot prompts toward ICL-level predictions (Hendel et al., 2023; Todd et al., 2024; Liu et al., 2024), and spurred further inquiries into the properties, behaviors, and emergence of IHs (Ren et al., 2024b; Singh et al., 2024; Yin & Steinhardt, 2025). (2) The holistic paradigm instead treats the LLM as an entirety and investigates ICL's properties by directly inspecting and probing how different demonstration configurations shape ICL performance. For instance, by perturbing the demonstration labels in context, Pan et al. (2023) factorize ICL into two core components: Task Recognition (TR, recognizing the label space) and Task Learning (TL, learning the text-label mapping), each contributing to part of the ICL functionality (Figure 1 (A)). Min et al. (2022) also systematically explored the effect of the distribution of texts and labels in demonstrations individually, as well as the templates and number of demonstrations.

¹The source code will be released upon acceptance of this paper.

Figure 1: (A) Example of how LLMs deduce the label of a final query through ICL, which consists of two components: task recognition (identifying the label space) and task learning (mapping demonstration texts to labels). (B) The outputs of Task Recognition heads align with the task subspace spanned by candidate label unembeddings, thus they can steer the hidden states to align with the subspace by reducing the angle between the point clouds and the task subspace. (C) Task Learning heads act as rotations within the task subspace, aligning the query's hidden state with the unembedding of the correct label and enabling the correct prediction.

The two research paradigms offer complementary insights but also limitations. The introspective paradigm localizes ICL to individual attention heads, yet its reliance on ablation only measures **how much** performance changes when heads are removed, without explaining **how** these heads realize ICL or behave under varied inputs. The holistic paradigm provides a broad functional view, separating ICL into TR and TL, but cannot trace these roles back to concrete components. A unified framework is needed to combine mechanistic precision with functional clarity.

Therefore, in this paper, we propose such a framework through attention head analysis, illustrated in Figure 1. Using the Task Subspace Logit Attribution (TSLA) method, which analyzes how each head contributes to the movement of hidden states w.r.t. the unembedding vectors of the task-related labels from a geometric perspective, we identify Task Recognition Heads (TR heads) and Task Learning Heads (TL heads) that are critically responsible for the TR and TL components of ICL. Geometric analyses (Kirsanov et al., 2025; Yang et al., 2025; Marks & Tegmark, 2024) reveal that TR heads align hidden states with the subspace spanned by task-related token unembeddings, helping the model recognize them as part of the label space and preventing the prediction of irrelevant tokens (Figure 1 (B)). TL heads then rotate the hidden states within the task subspace toward the correct label unembedding and facilitate correct prediction (Figure 1 (C)). Together, these mechanisms guide LLMs to perform the complete ICL functionality.

We validate these mechanisms through ablation and novel steering experiments across diverse ICL settings, including corrupted demonstrations and free-form generation—scenarios often overlooked in prior work. Our framework also reconciles earlier findings: for **IHs**, correlation analyses indicate they are best understood as a subset of TR heads whose main role is label-space recognition. For **task vectors**, we show that the primary barrier to accurate zero-shot prediction is weak alignment between hidden states and the task subspace, explaining why outputs of TR heads—and thus IHs—naturally serve as effective task vectors (Todd et al., 2024).

2 RELATED WORKS

Early investigations into ICL relied on input perturbation experiments (Min et al., 2022; Pan et al., 2023; Wei et al., 2023). By removing labels with semantic content (e.g., "negative" \rightarrow "0") or scrambling the text-label mapping while still observing non-trivial accuracy, these works concluded that ICL is a composite mechanism comprising two components: Task Recognition (TR) and Task Learning (TL). However, this approach is limited, as it cannot provide mechanistic explanations with finer granularity and mechanistically tie ICL functionality to specific model components.

This gap was narrowed by the circuit formulation of Transformers (Elhage et al., 2021; Olsson et al., 2022), which decomposes output logits into contributions from individual attention heads and MLPs. This allowed precise attribution of LLM behaviors to components (Crosbie & Shutova, 2024) and spotlighted the **Induction Head (IH)** as crucial for ICL (Zheng et al., 2024; Song et al., 2025). In toy copying tasks ([X][Y][X][Y]...[X] \rightarrow [Y]), IHs attend to earlier occurrences of [Y], mimicking label copying. Their importance for ICL has been confirmed in large-scale models through ablation studies that directly replace or remove their outputs and observe a change in final logits (Cho et al., 2025a).

Building on this methodology, later work identified **Function Vector Heads** (Todd et al., 2024; Yin & Steinhardt, 2025)—those exerting the strongest influence on ground-truth label logits and thus conceived as the cause of ICL functionality (Sun et al., 2025). Similar approaches extend to retrieval-augmented generation (distinguishing heads leveraging parametric vs. external knowledge) (Jin et al., 2024; Kahardipraja et al., 2025) and chain-of-thought reasoning (Cabannes et al., 2024). Yet these methods reveal only **how much** a head contributes, not **how** it contributes or **which** ICL component it affects. Existing attempts at explanation, often based on attention heatmaps (Kahardipraja et al., 2025; Ren et al., 2024a), remain too shallow to establish causal significance.

Recent studies move beyond attribution by aggregating head outputs into **task vectors** that steer zero-shot hidden states toward ICL-level accuracy (Hendel et al., 2023), where heads producing effective task vectors are seen as mechanistic origins of ICL (Todd et al., 2024; Yin & Steinhardt, 2025). However, this line of work faces interpretability challenges: 1) candidate heads are selected using ablation-based methods, treating the channel linking attention heads to outputs as a black box; 2) construction of task vectors sometimes involves opaque optimization (Li et al., 2024); and 3) the model's use of injected vectors is described only as a black-box function (Merullo et al., 2024).

To address these interpretability challenges, geometric analysis of layer-wise hidden state evolution which incorporates the effects of attention head outputs and injected task vectors offers a promising alternative (Kirsanov et al., 2025). Jiang et al. (2025) identify a compress-expand pattern in task representations during ICL, while Yang et al. (2025) link IH outputs to alignment between hidden states and unembedding vectors of task-relevant labels. This geometric perspective enables deeper insight into how attention heads influence LLM outputs in complex ICL settings.

3 METHODOLOGY

3.1 BACKGROUND

Circuit Formulation of Transformer In the circuit formulation of Transformer LLMs, an input query q with N tokens $[x_1,...,x_N]$ (e.g., "I like this movie. Sentiment:" for a sentiment analysis task) is first transformed into layer-0 hidden states $h_1^0,...,h_N^0$ via the embedding matrix $W_E \in \mathbb{R}^{|\mathbb{V}| \times d}$, where d is the model dimension and \mathbb{V} the vocabulary. These hidden states then pass through L layers, where the update of the i-th token's hidden state at layer l is:

$$m{h}_i^l = m{h}_i^{l-1} + \sum_{k=1}^K m{a}_{i,k}^l + m{m}_i^l,$$

with $\boldsymbol{a}_{i,k}^l$ the output of the k-th attention head (denoted head (l,k)) in the attention sublayer, and \boldsymbol{m}_i^l the MLP sublayer output. $\boldsymbol{a}_{i,k}^l$ is the weighted sum of the layer-(l-1) hidden states of the first i tokens, $\boldsymbol{H}_{\leq i}^{l-1} = [\boldsymbol{h}_j^{l-1}]_{j=1}^i$, transformed by the embedding matrices of head (l,k). The final hidden state of the last token (":" in the previous example) can thus be written as:

$$h_N^L = h_N^0 + \sum_{l=1}^L \left(\sum_{k=1}^K a_{N,k}^l + m_N^l \right).$$
 (1)

 \boldsymbol{h}_N^L is multiplied by the unembedding matrix $\boldsymbol{W}_U \in \mathbb{R}^{d \times |\mathbb{V}|}$ to form logits. Each head output thus contributes to the logits additively as $\boldsymbol{a}_{N,k}^l \boldsymbol{W}_U$, referred to as **Direct Logit Attribution (DLA)** (Olsson et al., 2022; Chughtai et al., 2024; Yu & Ananiadou, 2024; Lieberum et al., 2023).

ICL and Induction Head In ICL, m text-label demonstration pairs $t_1, y_1, ..., t_m, y_m$ are prepended to the query, forming the sequence $t_1, y_1, ..., t_m, y_m, q$ (e.g., "I hate this movie. Sentiment: negative. This movie is great. Sentiment: positive... I like this movie. Sentiment:"). With these demonstrations, the attention head outputs $a_{N,k}^l$ to the final position, depending on all preceding tokens, producing logits that can lead the LLM to predict y^* , the correct label for q. An Induction Head (IH) is a special attention head that, at each position, searches for earlier occurrences of the current token, attends to the immediately following tokens, and copies their information back to the current position. In the example above, an IH at the final position places its attention on the "positive" and "negative" tokens that follow previous ":" tokens and uses their hidden states to form its output.

3.2 IDENTIFYING TR AND TL HEADS USING TASK SUBSPACE LOGIT ATTRIBUTION

Pan et al. (2023) decomposes ICL into two components. **Task Recognition (TR)** means recognizing the set of candidate task label tokens $\mathbb Y$ from demonstration labels, with $\{y_1,...,y_m\} \in \mathbb Y$, without using the text-label mapping information to deduce the correct token. **Task Learning (TL)**, in contrast, means learning the mapping from demonstration texts to task labels, $f: \mathbb X \to \mathbb Y$, to predict the only correct label for query q. To identify heads contributing to TR and TL, Lieberum et al. (2023) compute $a_{N,k}^l W_U^{\mathbb Y}$ and $a_{N,k}^l W_U^{\mathbb Y}$ for all heads (l,k), where $W_U^{\mathbb Y} \in \mathbb R^{d \times |\mathbb Y|}$ and $W_U^{\mathbb Y} \in \mathbb R^d$ are the unembedding matrix restricted to $\mathbb Y$ and y^* . Heads with the highest element-wise sum $1^T a_{N,k}^l W_U^{\mathbb Y}$ are considered TR heads, and those with the highest $a_{N,k}^l W_U^{\mathbb Y}$ are TL heads.

This approach has two problems. 1) For TR heads, Lieberum et al. (2023) study four-choice tasks where the full label space is "A", "B", "C", "D". In general settings, demonstration labels are arbitrary hyperparameters and may not capture full task semantics. Changing labels from positive/negative to favourable/unfavourable does not alter the task, but heads amplifying logits for positive/negative may not do so for favourable/unfavourable. 2) For TL heads, the method ignores competition among label tokens: heads boosting y^* may also boost incorrect labels \mathbb{Y}/y^* , disqualifying them as true task-mapping heads. A more precise approach must (a) capture task semantics beyond surface tokens and (b) evaluate contributions relative to competing labels.

We therefore propose the **Task Subspace Logit Attribution** (**TSLA**) method. For TR heads, we compute the TR score:

$$\|\operatorname{Proj}_{\boldsymbol{W}_{U}^{\mathbb{Y}}}\boldsymbol{a}_{N,k}^{l}\|_{2},$$
 (2)

where $\operatorname{Proj}_{\boldsymbol{W}_U^{\mathbb{Y}}} = \boldsymbol{W}_U^{\mathbb{Y}} (\boldsymbol{W}_U^{\mathbb{Y},\top} \boldsymbol{W}_U^{\mathbb{Y}})^{-1} \boldsymbol{W}_U^{\mathbb{Y},\top}$ is the $d \times d$ projection matrix onto $\operatorname{span}(\boldsymbol{W}_U^{\mathbb{Y}})$, the subspace spanned by unembedding vectors of demonstration labels, which also encompasses unembeddings of related tokens since LLMs encode related semantics as subspaces (Saglam et al., 2025; Zhao et al., 2025). The TR score—the projected norm of a head's output onto this subspace—thus captures logit contributions to all task-related semantics regardless of the chosen demonstration labels, alleviating the DLA approach's sensitivity to demonstration-label choice as a hyperparameter. We have the following theoretical guarantee for this metric's effectiveness.

Theorem 1 Let $r = |\mathbb{Y}|$. Assume n distinct r-dimensional subspaces drawn i.i.d. from the Grassmannian Gr(r,d) are spanned by columns of W_U . If head (l,k) has TR score γ , then with probability at least $1 - (n-1)(1 - I_{(\frac{\gamma}{\|\mathbf{a}_{N,k}^l\|_2})^2}(\frac{r}{2},\frac{d-r}{2}))$, $\mathbf{a}_{N,k}^l$ has the largest projected l_2 norm onto $\operatorname{span}(\mathbf{W}_U^{\mathbb{Y}})$ among all such subspaces,

where $I_x(\alpha, \beta) = \frac{B(x; \alpha, \beta)}{B(\alpha, \beta)}$ is the regularized incomplete beta function, increasing in x. This shows that a large TR score implies the head output is best captured by the subspace spanned by demonstration label unembeddings, qualifying it as a TR head. The proof is in Appendix B.

For TL heads, we compute:

$$\frac{\operatorname{Ave}_{y' \in \mathbb{Y}/\{y^*\}}(\boldsymbol{a}_{N,k}^{l,\top}(\boldsymbol{W}_{U}^{y^*} - \boldsymbol{W}_{U}^{y'}))}{\|\operatorname{Proj}_{\boldsymbol{W}_{U}^{\mathbb{Y}}}\boldsymbol{a}_{N,k}^{l}\|_{2}}.$$
(3)

The numerator is the mean inner product of the head output with the difference between the correct label unembedding and each incorrect label, which measures the logit difference a head creates between correct and incorrect labels. The denominator is the TR score. Since $W_U^{y^*} - W_U^{y'} \in \text{span}(W_U^{y'})$ for all $y' \in \mathbb{Y}/\{y^*\}$, the TL score ranges in [-1,1]. Geometrically, it is the proportion of the projected head output that aligns with the unembedding difference between correct and incorrect labels. Heads with high TL scores express logit contributions to task-related labels primarily by increasing the logit gap between correct and incorrect labels. They can steer hidden states to better align with the correct label's unembedding within the task subspace when added to the residual stream. This conforms to task learning, which centers on identifying the correct label and excluding incorrect ones for an input. Moreover, this TL score also mitigates the DLA issue of interference from incorrect labels by disregarding heads that fail to differentiate and instead raise both logits.

For each dataset, we use ICL prompts built from the first 50 queries to calculate the TR and TL scores of each head, with the scores summed across prompts. Heads are ranked by TR and TL scores to identify TR and TL heads. See Appendix C for a comparison with Lieberum et al. (2023)'s method.

EXPERIMENTS

216

217

218

219 220

221 222

224

225

226

227

228

229 230

231 232 233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266 267

268

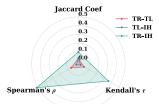
269

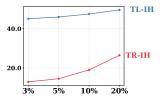
Models We experiment on models with diverse architectures and sizes, including Llama3-8B, Llama3.1-8B, Llama3.2-3B (Grattafiori et al., 2024), Qwen2-7B, Qwen2.5-32B (Yang et al., 2024), and Yi-34B (01. AI et al., 2024). Unless otherwise noted, results are reported on Llama3-8B.

Datasets We evaluate on the following datasets: SUBJ (Wang & Manning, 2012), SST-2 (Socher et al., 2013), TREC (Li & Roth, 2002), MR (Pang & Lee, 2005), SNLI (MacCartney & Manning, 2008), RTE (Dagan et al., 2005), and CB (De Marneffe et al., 2019). We also include an LLMgenerated dataset introduced in Subsection 4.3, with curation details in Appendix D and Appendix J.

ICL setting We use 8-shot demonstrations for ICL. For implementation details (models, datasets, prompt templates, etc.), see Appendix D.

4.1 VALIDATING TR/TL HEAD SPECIALIZATION AND THE ROLE OF INDUCTION HEADS





efficient, Kendall's τ , and Spear- age at four top thresholds for the man's ρ for TR heads, TL heads, TR-IH and TL-IH pairs averand IHs at the top 3% level.

(a) Dataset-averaged Jaccard Co- (b) Conditional Mean Percentaged over datasets.

Figure 2: Overlap, correlation, and consistency of three attention head types averaged across datasets. (A) TR heads exhibit substantially greater overlap and correlation with IHs compared to TR-TL or TL-IH pairs. (B) Top IHs consistently rank higher in the TR ranking than in the TL ranking. Results for other models are in Appendix F.1.

To examine whether the TR and TL heads we identified indeed capture the distinct TR and TL components of ICL, we first analyze overlap, correlation, and consistency between them. We also include IHs in the analysis. Given their significance in mechanistic accounts of ICL (Zheng et al., 2024), we also wish to know whether IHs contribute by improving TR, TL, or both. Following Todd et al. (2024) and Yang et al. (2025), we compute IH scores (i.e., the degree to which a head's attention pattern resembles that of IHs) as described in Appendix E.

Adopting the methodology of Yin & Steinhardt (2025), we report Jaccard

Coefficients² among the top 3% of each head type to measure overlap at the top level. We also compute Kendall's τ and Spearman's ρ among the three rankings to evaluate the global correlations levels among the head types. Finally, we introduce **Conditional Mean Percentage**, which measures the average rank percentile of the top 3%, 5%, 10%, and 20% IHs within the TR and TL rankings. This metric bridges the local and global perspectives and answers the question of how significantly on average does the top IHs exhibit the TR and TL functionality, which is important in subsequent ablation-based experiments in Subsection 4.2.

Strong association between IHs and TR heads Figure 2a shows that TR heads and IHs are highly similar: 1) their overlap at the top 3% is much larger than either TR-TL or TL-IH pairs, and significantly above random baseline³; 2) their rank correlations are consistently higher. By contrast, TR-TL and TL-IH pairs show weaker overlap and correlation. This reveals how IHs, the magnitude of whose influence on ICL has been widely recorded, affect ICL: they influence ICL mainly by enabling recognition of the label space, rather than selectively amplifying the correct label. It reconciles conflicting previous findings with some reporting IHs as recognizing correct labels (Olsson et al., 2022; Cho et al., 2025b), others mentioning "false induction heads" that mislead (Halawi et al., 2024; Yu & Ananiadou, 2024). It also echoes Yin & Steinhardt (2025), who observed IHs overlap strongly with "function vector heads"⁴, reinforcing the centrality of TR heads in ICL.

²For two subsets $\mathbb{A}_1, \mathbb{A}_2 \subseteq \mathbb{A}$, the Jaccard Coefficient is $\frac{|\mathbb{A}_1 \cap \mathbb{A}_2|}{|\mathbb{A}_1 \cup \mathbb{A}_2|}$

³For random subsets of size k%, the expected Jaccard Coefficient is $\frac{k}{200-k}$, which is 0.0152 for k=3.

⁴Heads revealed to have greatest impact on correct label logits through ablations.

Figure 2b further shows that top IHs consistently rank higher within TR rankings than TL rankings. For example, the top 10% IHs correspond to roughly the top 20% TR heads but only the top 50% TL heads. This further justifies the large accuracy implications of ablating top IHs, which would imply ablating fairly high-ranking TR heads and the failure of task recognition.

Layer-wise distribution of special heads Figure 3 shows the per-layer distribution of the top 10% TR, TL, and IH heads for SST-2. We observe: 1) TR heads appear significantly deeper than both TL heads and IHs, while the layer distributions of TL heads and IHs are more similar (see Appendix F.2 for significance tests). This partly echoes but also challenges Yin & Steinhardt (2025), who reported function vector heads as only "slightly deeper" than IHs. 2) The TR-IH overlaps are much greater than TL-IH or TR-TL overlaps, and occur primarily in deeper layers. This indicates that the correlation between TR heads and IHs is systematic rather than haphazard: the overlaps conform to the general trend of TR heads being concentrated in deeper layers, instead of reflecting coincidental matches with scattered TR heads that occasionally appear in early layers.

Cross-dataset correlation To examine whether TR, TL, and IH heads generalize across tasks, we measure pairwise Jaccard, Kendall's τ , and Spearman's ρ among top 3% heads identified on seven datasets,

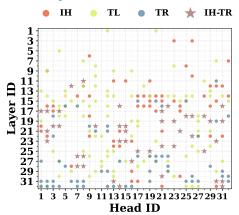


Figure 3: Distribution of the top 10% TR heads, TL heads, and IHs on SST-2. TR heads occur significantly deeper than TL heads and IHs. Overlaps between TR heads and IHs are also more frequent in deeper layers. Results for other models are in Appendix F.2.

averaged across $\binom{7}{2} = 21$ dataset pairs. As shown in Figure 4, TR heads and IHs exhibit much higher cross-task overlap and correlation than task-specific TL heads. This underscores TR heads (and IHs) as task-invariant mechanistic foundations for recognizing label spaces, upon which TL heads specialize to learn dataset-specific mappings.

4.2 TESTING INDEPENDENT CONTRIBUTIONS OF TR AND TL VIA ABLATION

We now show that the TR and TL heads identified by our method indeed independently and effectively capture the respective TR and TL functionalities through ablation experiments. Prior studies have mainly measured the effect of ablation on overall ICL accuracy (Crosbie & Shutova, 2024). However, as argued in Section 2, this provides only a coarse view of **how much** performance drops, without revealing **how** heads contribute via the TR or TL components. To address this, we introduce the **Task Recognition Ratio** (**TR ratio**), defined as the proportion of predictions that fall within the in-context label set. Formally, for m ICL prompts with predicted labels $\hat{y}_1, ..., \hat{y}_m$,

TR ratio =
$$\frac{1}{m} \sum_{i=1}^{m} \mathbf{1}_{\hat{y}_i \in \mathbb{Y}}$$
.

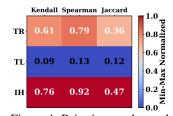
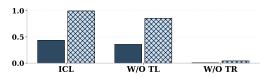
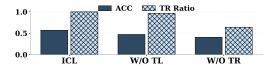


Figure 4: Pairwise overlap and correlation of TR/TL heads, and IHs identified across datasets. TR heads and IHs are consistent across tasks, while TL heads vary greatly. See Appendix F.1 for other models.

Since accuracy is upper-bounded by the TR ratio, the two metrics together let us separately evaluate contributions of TR and TL heads. We conjecture: (1) ablating TR heads should reduce both accuracy and TR ratio, while (2) ablating TL heads should reduce accuracy but leave TR ratio largely intact—causing performance to approximate random guessing over all candidate labels with expected accuracy $\frac{1}{|\mathbb{Y}|}$. As a control, we also ablate 3% of randomly chosen heads disjoint from the identified TR/TL sets. We report the results averaged over the seven datasets.

 $^{^5}For$ the seven datasets with 4 having 2 labels, 2 having 3 labels, 1 having 6 labels, the average random guessing level is 40.48%





(a) Ablating with shuffled demonstration texts.

(b) Ablating with relabeled demonstrations.

Figure 6: Dataset-average effects of ablating TR and TL heads under input perturbations. (A) Shuffling character order of demonstration texts destroys TL, making TL ablation negligible while TR ablation still matters. (B) Relabeling demonstrations alters the label space, thereby greatly reducing the impact of TR head ablation.

Separability of TR and TL functionality Figure 5 confirms our conjecture. Removing top TR heads collapses the TR ratio from nearly 100% to $\sim\!20\%$, leading to a drastic accuracy drop. In contrast, removing top TL heads lowers accuracy by $\sim\!30\%$ but only slightly reduces the TR ratio (by $\sim\!10\%$). This highlights a key property: separability, i.e., TR and TL can be independently controlled and intervened upon, consistent with the conclusions in Pan et al. (2023)

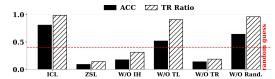


Figure 5: Effects of ablating the top 3% of different heads, averaged across datasets. TR head ablation severely reduces both accuracy and TR ratio, while TL head ablation primarily reduces accuracy.

achieved through input perturbations. (see Appendix G.1 for other models).

TR heads, IHs, and implications for zero-shot Ablating IHs produces a pattern closely resembling TR head ablation: large accuracy losses primarily due to failed task recognition. This supports the conclusion that IHs influence ICL mainly by strengthening TR. Likewise, the root cause of poor zero-shot performance is insufficient task recognition. Thus, restoring ICL-level accuracy in zero-shot settings hinges on activating the TR functionality—a question we revisit in Subsection 4.3.

Testing independence via input perturbations To further probe the separability of TR and TL, we perturb ICL inputs following Wei et al. (2023); Pan et al. (2023). Specifically: **Case 1:** Keep labels unchanged but randomize the character order of demonstration texts (e.g., "I like it: positive" \rightarrow "tkl iieI: positive"). This destroys TL, since no meaningful mapping from such nonsensical texts to labels remains. **Case 2:** Keep texts unchanged but replace demonstration labels with arbitrary tokens (e.g., "negative" \rightarrow "0", "positive" \rightarrow "1"), thereby altering the label space recognized by TR heads.

We hypothesize that: if the TR heads and TL heads maintain sufficient independence, then in Case 1, ablating TL heads should have little effect (TL is already disabled), while in Case 2, ablating TR heads should matter less (the original TR functionality is nullified).

As shown in Figure 6a, when texts are shuffled, TL ablation barely matters while TR ablation remains devastating. Conversely, in Figure 6b, TR ablation has little effect, compared to the significant impact shown in Figure 5, since the recognized label space has shifted, while TL ablation behaves as in the standard case. These results confirm the robustness of TR/TL independence across diverse ICL settings (see Appendix G.2 for other models).

4.3 STEERING WITH TR/TL HEADS: FUNCTIONAL AND GEOMETRIC INSIGHTS

Ablation experiments, while informative, are insufficient for mechanistic explanations: they show what happens when heads are removed, but not what happens when added. To address this limitation, we complement ablation with steering experiments that examine how TR and TL heads contribute when actively injected. Specifically, we test their suitability as task vectors (TVs) (Todd et al., 2024; Hendel et al., 2023) by extracting their outputs at the final token position from ICL prompts, summing them across prompts, and injecting them into the residual stream of zero-shot inputs to test whether accuracy improves toward ICL levels. We use the top 3% TR/TL heads and compare against a baseline of 3% random heads and follow procedures in Appendix I.1 to construct task vectors.

Task recognition as the key to zero-shot failure Figure 7 mirror our ablation findings (Figure 5): poor zero-shot performance stems primarily from weak task recognition. Injecting TR-based TVs

restores this functionality and improves performance. TL-based vectors are less effective, reinforcing that TL operates based on TR (see Appendix H.1 for other models).

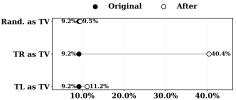


Figure 7: Zero-shot accuracy gains from steering with task vectors constructed from TR, TL, or random heads. TR-based task vectors consistently recover ICL-level accuracy, while TL-based vectors have weaker effects.

Task-type dependence of steering effectiveness Note that the relative ineffectiveness of TL heads as task vectors can be partly attributed to the classification datasets we use, where performance is tightly linked to task recognition and effectively upper-bounded by the TR ratio. In contrast, generation tasks differ fundamentally in that no fixed label space exists—the label space is indefinite and, in principle, infinite. As a result, model success in such tasks is less constrained by recognizing a closed set of labels, and instead depends more on learning and applying the correct input—output mapping. To

examine this scenario, we consider a sentiment-controlled review generation task with prompts such as: "Write a positive/negative review of a movie within 30 words." Labels are coherent reviews with the desired sentiment⁶. We identify TR and TL heads on ICL-styled prompts from this task following Appendix I.2, and use their outputs as task vectors to influence the zero-shot generations. An LLM evaluator is then used to rate generations from 0 to 10 based on sentiment adherence and language coherence. As shown in Figure 8, TL-based vectors significantly outperform TR-based and random vectors, consistent with TL heads capturing mappings from demonstrations to the sentiment values and semantic coherence of the labels (see Appendix H.2 for other models).

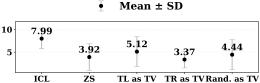


Figure 8: Ratings in the review generation task when steering with TR, TL, or random TVs. TL vectors yield the largest improvements, reflecting their strength in capturing in-context mappings.

Geometric effects of TR and TL outputs To understand the significance of TR/TL heads in ICL at a finer level than task vector experiments, we invoke the geometric analysis of hidden states (Kirsanov et al., 2025; Yang et al., 2025), which analyzes the evolution of ICL hidden states and the role of different components. Concretely, given an ICL prompt, we extract the summed outputs of the top 3% TR or TL heads, revert to the hidden state at an earlier layer, and steer it with these outputs. This mimics how

head outputs are added to the residual stream during layer progression inside the model. We measure two geometric metrics before and after steering: (1) **Logit Difference**: inner product of the hidden state with the mean unembedding difference between correct and incorrect labels, reflecting label discrimination. (2) **Subspace Alignment**: cosine similarity between the hidden state and the subspace spanned by label unembeddings, reflecting alignment with task-related semantics ⁷.

The results in Figure 9 demonstrate the specialized geometric effects TR and TL heads have in the evolution of hidden states (for other models, see Appendix H.3). Steering with TR outputs causes hidden states to align significantly more with the task subspace. In contrast, TL outputs adjust the hidden state to align better with the unembedding direction of the correct label in the task space but not with the task subspace overall. This leads us to conjecture that TR outputs are well-aligned with the task subspace, thus increase hidden-state alignment with the subspace after addition by decreasing the angle in between. By contrast, TL heads create pure rotation toward the correct label unembedding direction.

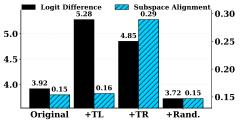
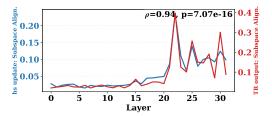


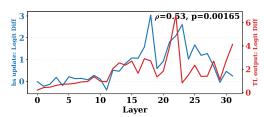
Figure 9: Geometric effects of TR and TL steering. TR outputs enhance alignment with the task subspace, while TL outputs rotate hidden states toward the correct label unembedding within the subspace.

toward the correct label unembedding direction, fine-tuning hidden-state orientation toward the correct label without boosting subspace alignment. See Appendix H.3 for other models.

⁶Example positive review: "Bold experimental narrative structure defies genre conventions delightfully. Socioeconomic themes challenge viewers' perceptions thoughtfully and respectfully."

⁷See Appendix I.3 for full definitions.





(a) Correlation between hidden state updates and TR head outputs in subspace alignment. Strong correlation confirms TR heads as main drivers of alignment.

(b) Correlation between hidden state updates and TL head outputs in logit difference. TL heads consistently drive discrimination toward correct labels.

Figure 10: Layerwise verification of TR and TL geometric effects. TR heads enforce alignment with the task subspace, while TL heads enforce rotations toward correct label directions.

Layerwise verification of geometric influence To validate that TR and TL heads indeed primarily drive these geometric dynamics, we examine hidden state updates under ICL across layers. At each layer, we compute the mean subspace alignment of top-3 TR heads (i.e. heads with top-3 TR scores at the layer) outputs and the mean logit difference of top-3 TL heads outputs, then correlate them with the same metrics computed on the full hidden state updates. Since head outputs contribute directly to layer updates, their correlations with hidden-state geometry across layers indicate how strongly TR and TL heads drive the dynamics. As shown in Figure 10, the correlations are strong, confirming that TR and TL heads dominate layerwise geometric shaping of hidden states. For other models see Appendix H.4. Additional ablation-based verification is provided in Appendix H.5.

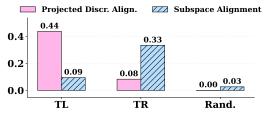


Figure 11: Decomposed geometric effects of TR and TL outputs. TR heads align hidden states to the task subspace; TL heads rotate states within the subspace toward correct label directions.

TR Heads align to task space, TL heads rotate within it To support our geometric intuition from Figure 9 that TR heads foster alignment while TL heads perform rotation, we consider two geometric measures of TR and TL head outputs. (1) Subspace Alignment — cosine similarity with the task subspace, and (2) Projected Discriminant Alignment — cosine similarity with the mean unembedding difference between the correct and incorrect labels after projection onto the task subspace. These measures dissect the geometric effects of head outputs into

steering towards the task space and steering within the task space, enabling more fine-grained verification of the heads' distinct effects (Figure 1 (B), (C)). Figure 11 shows that TL head outputs have high cosine similarity with the mean unembedding difference after projection, confirming that TL heads, when restricted to the task subspace, propel rotation from wrong-label to correct-label unembedding directions. The high cosine similarity between TR heads and the task subspace itself strongly evidences their capability to steer hidden states towards the task subspace and support prediction of task-related labels (see Appendix H.6 for more models).

5 Conclusion

We presented a unified framework that reconciles component-level and holistic views of in-context learning (ICL) by identifying attention heads specialized for **Task Recognition** (**TR**) and **Task Learning** (**TL**). Using Task Subspace Logit Attribution (TSLA), we showed that TR heads align hidden states with the task subspace, enabling recognition of candidate labels, while TL heads rotate states within this subspace toward the correct label. Ablation experiments confirmed their separable roles: removing TR heads collapses task recognition, whereas removing TL heads primarily reduces accuracy. Steering experiments further highlighted task dependence: TR-based vectors are crucial for classification tasks with fixed label spaces, while TL-based vectors dominate in open-ended generation. Geometric analyses corroborated these findings, attributing alignment effects to TR heads and discriminative rotations to TL heads. Our results also clarify the roles of induction heads and task vectors, positioning them as manifestations of TR functionality. Together, this work establishes TR and TL heads as mechanistic foundations of ICL.

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.
- 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.
- Vivien Cabannes, Charles Arnal, Wassim Bouaziz, Xingyu Yang, Francois Charton, and Julia Kempe. Iteration head: A mechanistic study of chain-of-thought. *Advances in Neural Information Processing Systems*, 37:109101–109122, 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*, 2025a. URL https://openreview.net/forum?id=xizpnYNvQq.
- Hakaze Cho, Yoshihiro Sakai, Mariko Kato, Kenshiro Tanaka, Akira Ishii, and Naoya Inoue. Token-based decision criteria are suboptimal in in-context learning. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 5378–5401, Albuquerque, New Mexico, April 2025b. Association for Computational Linguistics. ISBN 979-8-89176-189-6. URL https://aclanthology.org/2025.naacl-long.278/.
- Bilal Chughtai, Alan Cooney, and Neel Nanda. Summing up the facts: Additive mechanisms behind factual recall in llms. *arXiv preprint arXiv:2402.07321*, 2024.
- Joy Crosbie and Ekaterina Shutova. Induction heads as an essential mechanism for pattern matching in in-context learning. *arXiv preprint arXiv:2407.07011*, 2024.
- 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.
- Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank: Investigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*, volume 23, pp. 107–124, 2019. URL https://ojs.ub.uni-konstanz.de/sub/index.php/sub/article/view/601.
- 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.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1(1):12, 2021. URL https://transformer-circuits.pub/2021/framework/index.html.
- 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.
- Danny Halawi, Jean-Stanislas Denain, and Jacob Steinhardt. Overthinking the truth: Understanding how language models process false demonstrations, 2024. URL https://arxiv.org/abs/2307.09476.

- Roee 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/.
 - Jiachen Jiang, Yuxin Dong, Jinxin Zhou, and Zhihui Zhu. From compression to expansion: A layerwise analysis of in-context learning, 2025. 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.
 - 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/.
 - Dongfang Li, Zhenyu Liu, Xinshuo Hu, Zetian Sun, Baotian Hu, and Min Zhang. In-context learning state vector with inner and momentum optimization, 2024. URL https://arxiv.org/abs/2404.11225.
 - 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.
 - Tom Lieberum, Matthew Rahtz, János Kramár, Neel Nanda, Geoffrey Irving, Rohin Shah, and Vladimir Mikulik. Does circuit analysis interpretability scale? evidence from multiple choice capabilities in chinchilla, 2023. URL https://arxiv.org/abs/2307.09458.
 - 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.
 - 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.
 - Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Language models implement simple word2vecstyle vector arithmetic, 2024. URL https://arxiv.org/abs/2305.16130.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv* preprint arXiv:2202.12837, 2022.
 - Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,

- and Chris Olah. In-context learning and induction heads, 2022. URL https://arxiv.org/abs/2209.11895.
 - OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, et al. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.
 - Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What in-context learning "learns" in-context: Disentangling task recognition and task learning. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 8298–8319, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.527. URL https://aclanthology.org/2023.findings-acl.527/.
 - Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *arXiv* preprint cs/0506075, 2005.
 - 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.
 - Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Quanshi Zhang, Xipeng Qiu, and Dahua Lin. Identifying semantic induction heads to understand in-context learning. *arXiv preprint arXiv:2402.13055*, 2024a.
 - Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Quanshi Zhang, Xipeng Qiu, and Dahua Lin. Identifying semantic induction heads to understand in-context learning, 2024b. URL https://arxiv.org/abs/2402.13055.
 - Baturay Saglam, Paul Kassianik, Blaine Nelson, Sajana Weerawardhena, Yaron Singer, and Amin Karbasi. Large language models encode semantics in low-dimensional linear subspaces, 2025. URL https://arxiv.org/abs/2507.09709.
 - Aaditya K Singh, Ted Moskovitz, Felix Hill, Stephanie CY Chan, and Andrew M Saxe. What needs to go right for an induction head? a mechanistic study of in-context learning circuits and their formation. In *Forty-first International Conference on Machine Learning*, 2024. URL https://arxiv.org/abs/2404.07129.
 - 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.
 - Chenghao Sun, Zhen Huang, Yonggang Zhang, Le Lu, Houqiang Li, Xinmei Tian, Xu Shen, and Jieping Ye. Interpret and improve in-context learning via the lens of input-label mappings. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3873–3895, 2025.
 - JTianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning for language-model-as-a-service. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 20841–20855, Baltimore, Maryland, USA, 2022. ACM. URL https://proceedings.mlr.press/v162/sun22e/sun22e.pdf.

- 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.
- Sida Wang and Christopher Manning. Baselines and bigrams: Simple, good sentiment and topic classification. In Haizhou Li, Chin-Yew Lin, Miles Osborne, Gary Geunbae Lee, and Jong C. Park (eds.), *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics* (*Volume 2: Short Papers*), pp. 90–94, Jeju Island, Korea, July 2012. Association for Computational Linguistics. URL https://aclanthology.org/P12-2018/.
- Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning differently, 2023. URL https://arxiv.org/abs/2303.03846.
- 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, 2025. URL https://arxiv.org/abs/2505.18752.
- Kayo Yin and Jacob Steinhardt. Which attention heads matter for in-context learning? *arXiv preprint arXiv:2502.14010*, 2025.
- Zeping Yu and Sophia Ananiadou. How do large language models learn in-context? query and key matrices of in-context heads are two towers for metric learning. *arXiv preprint arXiv:2402.02872*, 2024.
- 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

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 4.3 and Appendix J.

PROOF OF THEOREM 1

Let $S \in \mathbb{R}^{d \times r}$ be one of the *n* distinct *r*-dimensional subspaces in span(W_U) drawn uniformly i.i.d. from the Grassmannian Gr(r,d). Denote P_S as the projection matrix of S. Let $c = \frac{\|P_S a_{N,k}^l\|_2}{\|a_{N,k}^l\|_2}$ be

the projected norm of the normalized head output $\frac{a_{N,k}^l}{\|a_{N,k}^l\|_2}$ onto S. Since the uniform distribution over Gr(r,d) is induced by the Haar measure over the orthogonal group O(d), the distribution is rotation-invariant; i.e., multiplying S by an orthogonal matrix $U \in \mathbb{R}^{d \times d}$ does not change its distribution. Because orthogonal transformations preserve angles, we also have

$$\frac{\|\boldsymbol{P_{US}Ua_{N,k}^l}\|_2}{\|\boldsymbol{a}_{N,k}^l\|_2} = \frac{\|\boldsymbol{P_{S}a_{N,k}^l}\|_2}{\|\boldsymbol{a}_{N,k}^l\|_2}.$$

Hence, without loss of generality, we pick an U such that $U \frac{a'_{N,k}}{\|a'_{N,k}\|_2} = e_1$, the unit vector in the first coordinate, with $c' = \|P_{US}e_1\|_2$ having the same distribution as c.

Since $US = \text{span}(v_1, ..., v_r)$, where $v_1, ..., v_r$ are the first r columns of a Haar orthogonal matrix V, let $V_r = [v_1,...,v_r]$. Then $P_{US} = V_r V_r^{ op}$, and we have

$$c'^2 = e_1^{\top} V_r V_r^{\top} e_1 = \|V_r^{\top} e_1\|_2^2 = \sum_{i=1}^r \langle e_1, v_i \rangle^2 = \sum_{i=1}^r V_{1,i}^2,$$

where $V_{1,:}$ denotes the first row of V_r . Since V_r is Haar orthogonal, $V_{1,:}$ is uniformly distributed on \mathbb{S}^{d-1} and has the same distribution as $\frac{g}{\|g\|_2}$ with $g \sim \mathcal{N}(0, I)$. Therefore, $\sum_{i=1}^r V_{1,i}^2$ has the same distribution as

$$\frac{\sum_{i=1}^r g_i^2}{\|g\|_2^2} = \frac{\sum_{i=1}^r g_i^2}{\sum_{i=1}^d g_i^2} = \frac{\chi_r^2}{\chi_r^2 + \chi_{d-r}^2},$$
 since $g_i \sim \mathcal{N}(0,1)$ for all i . Because $\chi_r^2 \perp \!\!\! \perp \!\!\! \chi_{d-r}^2$, we have

$$\frac{\chi_r^2}{\chi_r^2 + \chi_{d-r}^2} \sim \operatorname{Beta}\left(\frac{r}{2}, \frac{d-r}{2}\right).$$

If $c^2 \sim \text{Beta}(\frac{r}{2}, \frac{d-r}{2})$, then the tail probability is

$$\Pr(c \ge x) = 1 - I_{x^2}\left(\frac{r}{2}, \frac{d-r}{2}\right) = 1 - \frac{B(x^2; \frac{r}{2}, \frac{d-r}{2})}{B(\frac{r}{2}, \frac{d-r}{2})},$$

where B is the Beta function. Since the TR score of the head is γ , the probability that $\|P_S a_{N,k}^l\|_2 \ge \gamma$

$$1 - I_{\left(\frac{\gamma}{\|\boldsymbol{a}_{N,k}^I\|_2}\right)^2\left(\frac{r}{2},\frac{d-r}{2}\right).}$$

Because there are n-1 subspaces alongside $W_U^{\mathbb{Y}}$, the probability that the head output has the largest projected norm on $W_U^{\mathbb{Y}}$ is

$$1-(n-1)\left(1-I_{\left(\frac{\gamma}{\|\boldsymbol{a}_{N,k}^l\|_2}\right)^2\left(\frac{r}{2},\frac{d-r}{2}\right)}\right)$$

via the union bound.

C ABLATION EXPERIMENTS REGARDING THE IDENTIFICATION OF TR AND TL HEADS

In this section, we demonstrate the advantage of our Task Subspace Logit Attribution (TSLA) method over the naive approach of selecting TR and TL heads based on $a_{N,k}^l W_U^{\mathbb{Y}}$ and $a_{N,k}^l W_U^{\mathbb{Y}^*}$, i.e., Direct Logit Attribution (DLA) to the demonstration labels and the correct label. Specifically, Figure 12 shows the consequences of ablating the top 3% TR and TL heads identified via DLA, averaged across datasets on Llama3-8B. While ablating the identified TR heads achieves the intended effect of disabling task recognition by reducing the TR ratio, ablating the identified TL heads fails to induce the expected outcome of driving ICL toward random guessing over the label space. This indicates that the DLA approach cannot isolate distinct mechanistic causes for the TR and TL components of ICL, and reflects its inability to correct identify the real TL heads. Instead, it largely identifies heads that broadly amplify the logits of all demonstration label tokens, which may also increase the correct label logits but still primarily function through task recognition rather than true label differentiation.

To validate this statement, and following the setup of Figure 2, we report the dataset-averaged Jaccard Coefficient, Kendall's τ , and Spearman's ρ values between TR/TL heads identified by DLA and those identified by TSLA, as well as their relationship with IHs. As shown in Figure 13, both TL and TR heads selected via DLA strongly overlap with the TR heads identified by TSLA at the 3% level. This corroborates our conclusion that DLA fails to effectively recover genuine TL heads. Furthermore, the weak correlation between the TR/TL sets obtained from the two methods is reinforced by Figure 14, which displays overlap, correlation, and consistency analyses between DLA TR/TL heads and IHs. The strikingly high consistency between DLA TR and TL heads across all three metrics demonstrates the lack of a meaningful distinction between them. Meanwhile, the low correlation between DLA heads and IHs highlights another limitation of DLA: it cannot provide mechanistic explanations for the well-documented importance of IHs in ICL.

Finally, to justify our second critique of the DLA approach in Subsection 3.1 regarding its sensitivity to the concrete set of demonstration labels as a hyperparameter and its inability to comprehensively capture the task semantics, we consider the following experiment on SST-2. We replace "positive" and "negative", i.e. the default demonstration labels used to create ICL prompts from the dataset, with "unfavourable" and "favourable", which do not alter the essence of the task. Then we test how the ablation of TR heads identified with DLA and our TSLA using the ICL prompts with the original labels will impact the ICL accuracy and TR ratio with the new labels. The results in Figures 15–20 confirm the robustness of our TSLA method against demonstration label shifts in identifying TR heads. On all models except Qwen2.5-32B, ablating the TR heads causes a significantly larger impact on ICL performance and TR ratio with the new demonstration labels on the SST-2 dataset, with the gap being most prominent for the three Llama family models.

D IMPLEMENTATION DETAILS

Models We use the official HuggingFace implementations of all models. Models with more than 10B parameters are quantized to 4-bit for efficiency.

Datasets We use the official HuggingFace implementations of all datasets, except for the Review dataset, which we curated ourselves. The Review dataset was generated using ChatGPT-40 (OpenAI et al., 2024) and contains 200 datapoints. Each datapoint consists of a prompt instructing the model to generate a movie review in the format: 'Write a positive review for a movie. The positive review should be within 30 words.'' The 30-word limit was chosen to set the max_new_tokens parameter (set to 45) when calling the generation function. Labels are ChatGPT-40-generated reviews that comprehensively assess a movie from multiple aspects in the requested positive/negative tone. For example: "Bold experimental narrative structure defies genre conventions delightfully. Rich orchestral score enhances every pivotal moment. Progressive messages inspire reflection on equality and justice. Raw vulnerability on screen fosters sincere emotional investment." as a positive review. Details of dataset curation are provided in Appendix J. The dataset is balanced, with 100 positive and 100 negative reviews.

ICL setting For each dataset (except the Review dataset), we select demonstrations from the training set and queries from the test set, or the validation set if ground-truth test labels are unavailable. For

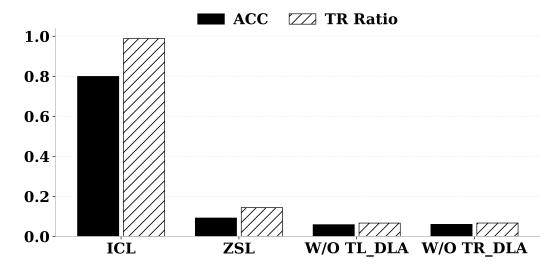


Figure 12: Effects of ablating the top 3% TR and TL heads identified using DLA, averaged across datasets on Llama3-8B. While TR heads reduce task recognition as expected, TL heads do not replicate the behavior predicted for task-learning components.

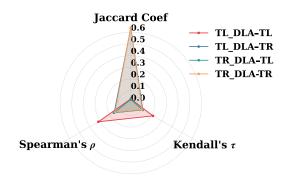
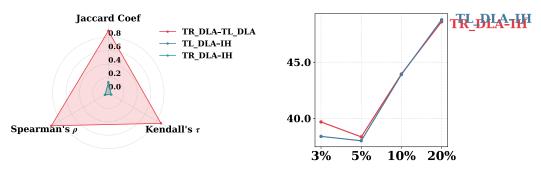


Figure 13: Dataset-averaged Jaccard Coefficient, Kendall's τ , and Spearman's ρ between TR and TL heads identified using DLA and TSLA at the top 3% level. DLA heads overlap substantially with TR heads, confirming their inability to recover distinct TL heads.



- (a) Dataset-averaged Jaccard Coefficient, Kendall's τ , and Spearman's ρ values for TR heads, TL heads, and IHs at the top 3% level.
- (b) Conditional Average Percentage at the top 3%, 5%, 10%, and 20% levels for TR-IH and TL-IH pairs, averaged across datasets.

Figure 14: Dataset-averaged overlap, correlation, and consistency analyses of TR and TL heads identified using DLA and their relationship with IHs. Results show high redundancy between DLA TR and TL heads and weak association with IHs, underscoring the limitations of DLA in separating TR and TL mechanisms or explaining IH significance.

demonstration selection, we retain at most the first 10,000 training examples. For evaluation, we use the first 1,000 test or validation examples. For the Review dataset, we use the first 50 examples for demonstration selection and the remaining 150 for testing. Prompt templates used to construct ICL prompts are listed in Table 1.

Devices All experiments were conducted on an H200 GPU.

Random label mappings For the experiment in Appendix F, where demonstration labels are replaced with numbers, we use the mappings in Table 2.

When demonstration labels are replaced with numeric symbols, we also modify the prompt templates in Table 1. Specifically, for SNLI and CB, "True or maybe or false" is replaced with "0 or 1 or 2," and for RTE, "True or false" is replaced with "0 or 1."

Flipped label mappings For the experiment in Subsection 4.2, where demonstration labels are flipped, we use the mappings in Table 3.

E EXPERIMENT DETAILS CONCERNING THE IDENTIFICATION OF IHS

For each dataset, we use the first 50 queries to identify the top IHs. Let the queries be q_1, \ldots, q_{50} , where q_i has token length $s(q_i)$. For each q_i , the LLM outputs an attention tensor $Attn_i \in \mathbb{R}^{L \times N_h \times s(q_i) \times s(q_i)}$, with L being the number of layers and N_h the number of attention heads per layer. The n_h -th head in layer l has an attention matrix $Attn(l, n_h)_i \in \mathbb{R}^{s(q_i) \times s(q_i)}$, where $Attn(l, n_h)_{i,j,k}$ denotes the attention from the k-th token to the j-th token in x_i .

Identification of IHs For each q_i , we randomly sample 8 demonstrations and prepend them to q_i . The resulting ICL prompt, Q_i (length $s(Q_i)$), follows the format $\langle t_{i,1} \rangle : \langle y_{i,1} \rangle, \ldots, \langle t_{i,8} \rangle : \langle y_{i,8} \rangle, \langle q_i \rangle$: where $\langle t_{i,k} \rangle$ is the sentence part of demonstration k (e.g., "I like this movie. Sentiment"), and $\langle y_{i,k} \rangle$ is the label (e.g., "positive"), separated by a colon. $\langle q_i \rangle$ is the sentence for the query. At the position of the final colon, an IH is expected to attend to tokens after previous colons—that is, the label tokens $\langle y_{i,1} \rangle, \ldots, \langle y_{i,8} \rangle$. Let \mathbb{I}_i be the set of label token indices in Q_i . The IH score for head (l,n_h) over the 50 queries is defined as $\sum_{i=1}^{50} \sum_{k \in \mathbb{I}_i} Attn(l,n_h)_{i,k,s(Q_i)}$, i.e., the total attention a head assigns at the final ":" position to the positions of all the label tokens, summed over all 50 queries. We calculate the IH scores for all (l,n_h) pairs and choose the top 3% attention heads as the identified **IH**s.

F SUPPLEMENTARY MATERIALS FOR SUBSECTION 4.1

F.1 REPLICATION OF FIGURE 2 FOR OTHER MODELS

Figures 21–25 replicate the experiments from Figure 2, demonstrating the robustness of our findings across different models. In every case, TR heads show markedly stronger overlap, correlation, and consistency with IHs than the other head pairs. The consistently higher values of the TR–IH pair in terms of Jaccard Coefficient, Kendall's τ , Spearman's ρ , and Conditional Average Percentage across all levels and architectures confirm our conclusion.

F.2 REPLICATION OF FIGURE 3 FOR OTHER MODELS

Figures 26–30 replicate the experiments from Figure 3 on additional models. These visualizations support our claims in Subsection 4.1 that: 1) TR heads generally reside in deeper layers than TL heads and IHs; 2) The overlap between TR heads and IHs is larger and predominantly occurs in deeper layers.

To complement these figures, Tables 4–9 report the mean layer indices of the top 3%, 5%, and 10% TR heads, TL heads, and IHs averaged across datasets. We also conduct Mann–Whitney U tests to assess whether the differences in layer distributions between TR heads and IHs, and between IHs and TL heads, are statistically significant. The results show that the distributional differences between IHs and TL heads are often not significant ($p \ge 0.05$). Even when significant, the p-values are much larger than those observed between TR heads and IHs, indicating that the TR–IH distinction is far more robust.

F.3 REPLICATION OF FIGURE 4 FOR OTHER MODELS

Figures 21–25 extend the analysis of Figure 4 to the remaining models. The results reinforce our conclusion that TR heads and IHs identified across datasets or tasks are largely consistent, whereas TL heads vary substantially.

To further test this, we evaluate the cross-dataset transferability of TR heads. Specifically, we ablate top 3% TR heads identified using SST-2 prompts and measure their impact on accuracy and TR ratio for the six remaining datasets. The results for all models in Figures 36-41 in general confirm the transferrability of TR heads across datasets, but some interesting variations among datasets and models are also worth noting. First, on the three Llama family models, ablating the TR heads identified on the SST-2 dataset can effectively drive both the accuracy and TR ratio on RTE, CB, and MR datasets to near zero, and to a lesser extent impact the two metrics on TREC and SNLI. On Yi-34B, the ablation instead significantly impact the model performance on SNLI rather than MR. For the remaining two Qwen family models the consequence of the ablation over datasets is similar to the case of Llama models but to a lesser degree overall.

G SUPPLEMENTARY MATERIALS FOR SUBSECTION 4.2

G.1 REPLICATION OF FIGURE 5 FOR OTHER MODELS

In Figures 31–35, we replicate the ablation experiments on additional models and examine their effects on dataset-average accuracy and TR ratio. The results echo our observations in Subsection 4.2:

1) Ablation of TR heads and TL heads impacts the TR and TL components of ICL separately. 2) Ablating IHs produces effects similar to ablating TR heads. 3) The primary cause of low accuracy in the zero-shot case—as well as in cases where TR heads or IHs are ablated—is the failure to adequately activate the TR functionality.

G.2 REPLICATION OF FIGURE 6 FOR OTHER MODELS

In Figures 42–46, we repeat the experiments from Figure 6 on other models, focusing on the ablation of TR and TL heads when ICL inputs are subjected to perturbations. The results closely mirror those in Figure 6: when the in-context text–label mapping is destroyed or reversed, ablating TL heads has no effect—or even a positive effect—on accuracy. By contrast, since these perturbations do not alter the demonstration label space, the TR component of ICL remains unaffected.

G.3 ASSESSING THE INDEPENDENCE OF TR AND TL WITH FLIPPED DEMONSTRATION LABELS

To further validate the independence of TR and TL and the mechanisms of their associated heads, we analyze the effect of ablating TR/TL heads under flipped demonstration labels. Specifically, we apply a mapping $g: \mathbb{Y}' \to \mathbb{Y}'$ that reverses the demonstration labels (e.g., "negative" \to "positive," "positive" \to "negative"), as listed in Table 3. Since label flipping invalidates the original text–label mapping captured by TL heads, we conjecture that ablating top TL heads will *increase* accuracy, while the effect of ablating TR heads will remain unchanged because the label space itself is preserved.

The dataset-average results in Figures 47–52 confirm this conjecture: ablating top TL heads indeed raises accuracy, whereas ablating top TR heads still drives accuracy close to zero, as observed in the standard setting of Figure 5.

H SUPPLEMENTARY MATERIALS FOR SUBSECTION 4.3

H.1 REPLICATION OF FIGURE 7 FOR OTHER MODELS

Figures 53–57 replicate the steering experiments from Figure 7, evaluating the effectiveness of task vectors constructed from special attention head outputs in other models. For all models except the two Qwen-family models, the results are consistent with Subsection 4.2: task vectors built from TR heads are substantially more effective than those from TL heads. In the Qwen models, however,

TL heads match TR heads as task vectors. This deviation can be explained by the high zero-shot accuracy of the Qwen models (Figure 56, Figure 55), which exceeds 20%—considerably higher than the other models. Because these models already achieve strong task recognition in the zero-shot setting, injecting TR-head-based task vectors (which primarily encode recognition) provides less additional benefit.

H.2 REPLICATION OF FIGURE 8 FOR OTHER MODELS

Figures 58–62 evaluate how task vectors built from different types of heads affect the quality of generated reviews across models. Consistently, TL heads outperform TR heads and random heads as task vectors. An exception is Yi-34B, where steering reduces the average rating below the original zero-shot level. For Qwen2-7B and Qwen2.5-32B, TL-head task vectors even push ratings above the ICL-level baseline. Interestingly, the zero-shot reviews of these models score higher than their ICL reviews. Closer inspection reveals why: ICL reviews, though coherent and stylistically faithful, sometimes contradict the sentiment required in the query. TL heads appear to filter out such inconsistencies by correctly capturing the text–label mapping and discarding misleading signals, thereby boosting zero-shot review quality beyond ICL.

In addition to cross-model replication, Table 10 presents sampled outputs under ICL, zero-shot, and steering with different task vectors. These examples highlight the TL heads' ability to extract the correct text—label mapping and use it for generation. In contrast, zero-shot or TR-head steering often yields generic, off-topic sentences loosely related to the concept of "review."

H.3 REPLICATION OF FIGURE 9 FOR OTHER MODELS

Figures 63–67 extend the geometric analysis of Figure 9 to other models. The results largely confirm our earlier observation: TL heads tend to align hidden states with label-unembedding difference directions, while TR heads align hidden states with the broader task subspace.

H.4 REPLICATION OF FIGURE 10 FOR OTHER MODELS

Figures 73–77 report layer-wise correlations between mean TR/TL head outputs and full layer updates, measured by logit difference and subspace alignment. Across models, we observe clear and consistent correlation patterns, reinforcing that TR and TL heads are the primary drivers of the geometric shaping of hidden states in layer updates.

${ m H.5}$ Ablating top TR and TL heads per layer to verify their geometric significance

In Figure 10, we validated the geometric importance of TR/TL heads by correlating their outputs with full layer updates. Here, we provide an alternative perspective. Specifically, we ablate the top three TR/TL heads per layer and then remeasure layer-wise hidden state updates under the same two metrics. Figures 78–83 show that TR and TL heads are indeed crucial: without top TR heads, hidden states fail to gradually align with the task subspace, crippling task recognition; without top TL heads, logit differences collapse, preventing hidden states from rotating toward the correct label's unembedding direction. By contrast, ablating three random heads per layer has negligible impact.

H.6 REPLICATION OF FIGURE 11 FOR OTHER MODELS

Figures 68–72 replicate the analysis of Figure 11 across models. The results are consistent: TL heads excel in projected discriminant alignment, rotating hidden states toward the correct label unembedding and away from incorrect ones. TR heads, conversely, excel in subspace alignment, keeping hidden states well-positioned within the task subspace. Both substantially outperform randomly chosen heads on their respective strengths.

I EXPERIMENT DETAILS RELATED TO TASK VECTORS

I.1 CONSTRUCTION AND APPLICATION OF TASK VECTORS

For each dataset, we first construct 8-shot ICL prompts using the last 50 queries. The demonstrations are identical to those used when evaluating the 8-shot ICL accuracy for each dataset. Following the procedure of Todd et al. (2024), we compute the average output (across the 50 prompts) of each identified top 3% TR, TL, or random head at the final token position. We then sum these average outputs across heads to form the task vector.

In the steering experiment, the task vector is added to the hidden state of the final token of each zero-shot query at the midpoint layer (e.g., layer 16 for the 32-layer Llama3-8B). The modified hidden states are then propagated through the subsequent layers, and accuracy as well as TR ratio are measured at the final layer.

I.2 IDENTIFYING TR AND TL HEADS ON THE MOVIE REVIEW DATASET

A key difficulty in identifying TR and TL heads for free-form generation tasks is the unbounded label space, since labels are not restricted to a finite set of tokens. To address this, we define the relevant label tokens as "positive" and "negative," reflecting the sentiment nature of the review-generation task. Specifically, the TR score of a head is defined as the projection norm of its output onto the span of the unembedding vectors of "positive" and "negative" when processing ICL prompts from the review dataset. The TL score is defined as the inner product between the head output and the difference between the unembeddings of "positive" and "negative," normalized by its TR score. After identifying TR and TL heads, we construct task vectors from their outputs following the procedure in Appendix I.1.

- I.3 MATHEMATICAL DETAILS OF THE MEASURES IN SUBSECTION 4.3 AND CALCULATION PROCEDURE
- 1. **Logit Difference** Given a hidden state h, we compute $\text{Ave}_{y' \in \mathbb{Y}/\{y^*\}}(h^\top (W_U^{y^*} W_U^{y'}))$, where W_U is the unembedding matrix, y^* is the correct label, and \mathbb{Y} is the demonstration label space.
- 2. **Subspace Alignment** We compute $\frac{\boldsymbol{h}^{\top}\operatorname{Proj}_{\boldsymbol{W}_{U}^{\mathbb{Y}}}^{\top}\boldsymbol{h}}{\|\operatorname{Proj}_{\boldsymbol{W}_{U}^{\mathbb{Y}}}^{\top}\boldsymbol{h}\|_{2}\|\boldsymbol{h}\|_{2}}$, which is the cosine similarity between $\operatorname{Proj}_{\boldsymbol{W}_{U}^{\mathbb{Y}}}^{\top}\boldsymbol{h}$ and \boldsymbol{h} .

For evaluation, we take the hidden state of the final position at the layer corresponding to 75% of the model depth (e.g., layer 24 in Llama3-8B). Reported metric values are averaged across the first 30 ICL prompts of each dataset.

J CURATION DETAILS OF THE REVIEW DATASET

We use the following template, adapted from Zhao et al. (2025), to prompt ChatGPT-40 to generate movie reviews.

rompt

Compose a concise 30-word movie review that addresses the following four aspects: plot, sound and music, cultural impact, and emotional resonance. Use a positive tone throughout the review. For the plot, comment on its structure or originality. For sound and music, describe how they enhance the storytelling. For cultural impact, mention any relevant social commentary. Finally, highlight how the film resonates emotionally. Ensure the positive tone is consistent throughout and include positive descriptions of the movie.

Samples

Inventive non-linear storyline weaves intrigue with clever twists. Soaring vocal melodies heighten the film's emotional arcs. Relevant socioeconomic themes challenge viewers' perceptions thoughtfully and respectfully. Joyful humor interwoven with drama creates comforting resonance.

We use the following template to ask ChatGPT-40 to rate the movie reviews. Rate the following movie review on a scale of 10. Your rating should be based on two criteria: (1) whether the text is indeed a movie review, and (2) whether it conveys the posi-Prompt tive or negative sentiment indicated by the label. Review: Inventive non-linear storyline weaves intrigue with clever twists. Soaring vocal melodies heighten the film's emotional arcs. Relevant socioeconomic themes challenge viewers' perceptions thoughtfully and respectfully. Joyful humor interwoven with drama creates comforting resonance. Sentiment: Positive Response

Table 1: Prompt templates and labels for different datasets.

Dataset	Template	Labels
SST-2	{Sentence} Sentiment: {Label}	positive / negative
SUBJ	{Sentence} Type: {Label}	subjective / objective
TREC	Question: {Sentence} Type: {Label}	abbreviation / entity / description / human / location / number
MR	{Sentence} Sentiment: {Label}	positive / negative
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

Table 2: Mappings used to replace ground-truth labels with numeric symbols.

Dataset	Label Mapping
SST-2	negative/positive \rightarrow 0/1
SUBJ	objective/subjective \rightarrow 0/1
MR	negative/positive $\rightarrow 0/1$
TREC	abbreviation/entity/description/person/number/location \rightarrow 0/1/2/3/4/5
SNLI	true/maybe/false \rightarrow 0/1/2
RTE	true/false \rightarrow 0/1
CB	true/maybe/false \rightarrow 0/1/2

Table 3: Mappings used to flip the demonstration labels for each dataset.

Dataset	Label Mapping
SST-2	negative/positive → positive/negative
SUBJ	objective/subjective \rightarrow subjective/objective
TREC	abbreviation/entity/description/person/number/location \rightarrow entity/description/person/number/location/abbreviation
MR	negative/positive → positive/negative
SNLI	true/maybe/false → maybe/false/true
RTE	true/false → false/true
CB	$true/maybe/false \rightarrow maybe/false/true$

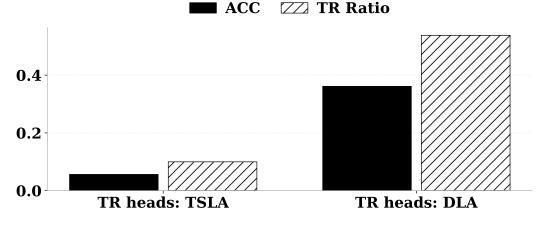


Figure 15: Results on LLama3-8B: Effects of ablating top 10% TR heads identified using TSLA or DLA when the SST-2 demonstration labels are shifted from positive/negative to favourable/unfavourable.

-	TL heads mean layer	IHs mean layer	TR heads mean layer	TL < IHs?	IHs < TR heads?
0.03	16.89	16.69	26.22	0.84985	5.6052e-45
0.05	16.94	16.74	25.08	0.80371	0.0000e+00
0.10	16.75	16.88	23.27	0.59518	0.0000e+00

Table 4: Mean layer index of TL heads, IHs, and TR heads across datasets on Llama3-8B, with p-values for distribution differences.

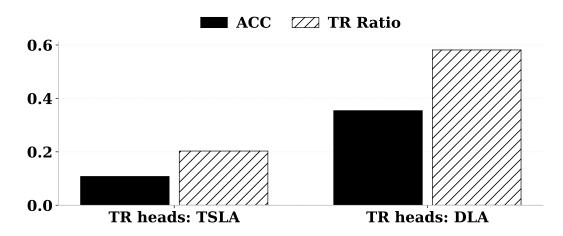


Figure 16: Results on Llama3.1-8B: Effects of ablating top 10% TR heads identified using TSLA or DLA when the SST-2 demonstration labels are shifted from positive/negative to favourable/unfavourable.

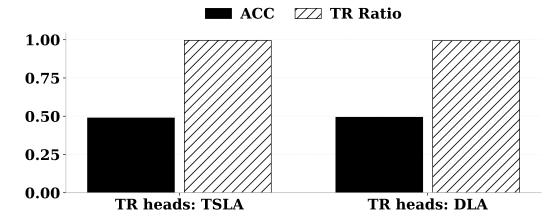


Figure 17: Results on Llama3.2-3B: Effects of ablating top 10% TR heads identified using TSLA or DLA when the SST-2 demonstration labels are shifted from positive/negative to favourable/unfavourable.

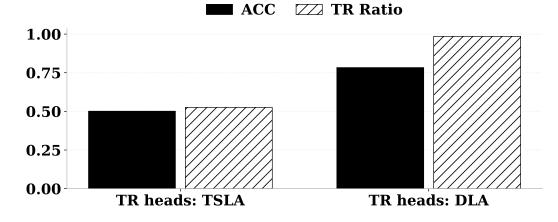


Figure 18: Results on Qwen2-7B: Effects of ablating top 10% TR heads identified using TSLA or DLA when the SST-2 demonstration labels are shifted from positive/negative to favourable/unfavourable.

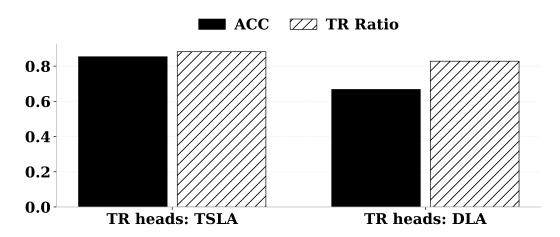


Figure 19: Results on Qwen2.5-32B: Effects of ablating top 10% TR heads identified using TSLA or DLA when the SST-2 demonstration labels are shifted from positive/negative to favourable/unfavourable.

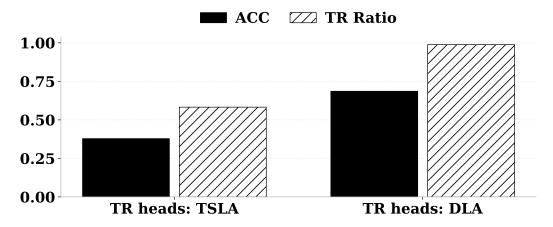
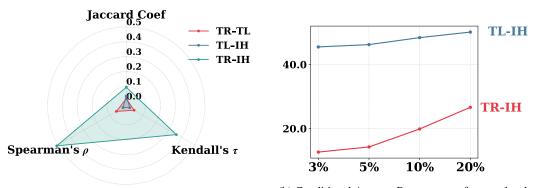


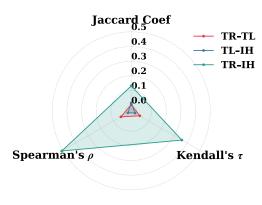
Figure 20: Results on Yi-34B: Effects of ablating top 10% TR heads identified using TSLA or DLA when the SST-2 demonstration labels are shifted from positive/negative to favourable/unfavourable.

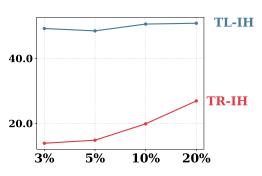


(a) Jaccard Coefficient, Kendall's τ , and Spearman's ρ values for TR heads, TL heads, and IHs.

(b) Conditional Average Percentage at four top levels for TR-IH and TL-IH pairs.

Figure 21: Results of overlap, correlation, and consistency analysis of attention head types averaged across datasets on Llama3.1-8B.

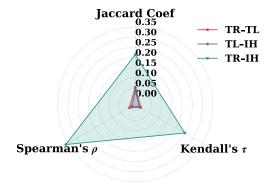


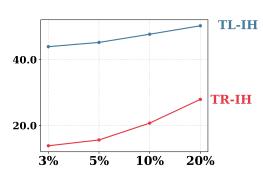


(a) Jaccard Coefficient, Kendall's τ , and Spearman's ρ values for TR heads, TL heads, and IHs.

(b) Conditional Average Percentage at four top levels for TR-IH and TL-IH pairs.

Figure 22: Results of overlap, correlation, and consistency analysis of attention head types averaged across datasets on Llama 3.2-3B.

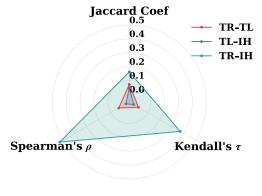


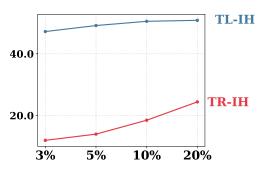


(a) Jaccard Coefficient, Kendall's τ , and Spearman's ρ values for TR heads, TL heads, and IHs.

(b) Conditional Average Percentage at four top levels for TR-IH and TL-IH pairs.

Figure 23: Results of overlap, correlation, and consistency analysis of attention head types averaged across datasets on Qwen2-7B.

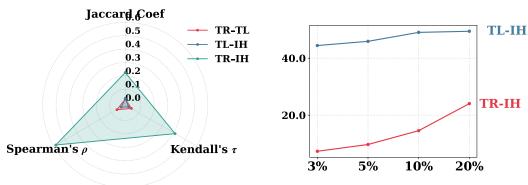




(a) Jaccard Coefficient, Kendall's τ , and Spearman's ρ values for TR heads, TL heads, and IHs.

(b) Conditional Average Percentage at four top levels for TR-IH and TL-IH pairs.

Figure 24: Results of overlap, correlation, and consistency analysis of attention head types averaged across datasets on Qwen2.5-32B.



(a) Jaccard Coefficient, Kendall's τ , and Spearman's ρ values for TR heads, TL heads, and IHs.

(b) Conditional Average Percentage at four top levels for TR-IH and TL-IH pairs.

Figure 25: Results of overlap, correlation, and consistency analysis of attention head types averaged across datasets on Yi-34B.

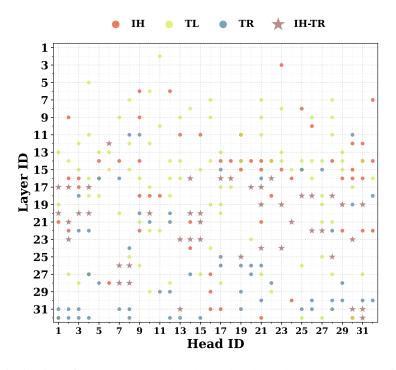


Figure 26: Distribution of the top 10% TR heads, TL heads, and IHs across layers for the SST-2 dataset on Llama3.1-8B.

	TL heads mean layer	IHs mean layer	TR heads mean layer	TL < IHs?	IHs < TR heads?
0.03	16.85	16.67	26.07	0.71998	8.4078e-45
0.05	16.36	16.77	25.33	0.22671	0.0000e+00
0.10	16.12	16.79	23.25	0.049241	0.0000e+00

Table 5: Mean layer index of TL heads, IHs, and TR heads across datasets on Llama3.1-8B, with p-values for distribution differences.

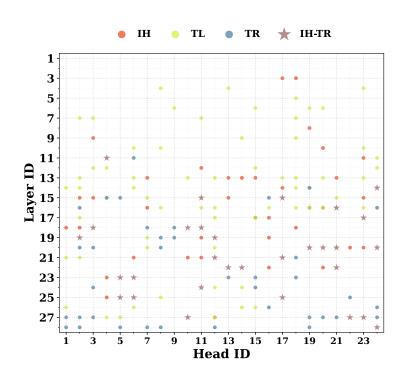


Figure 27: Distribution of the top 10% TR heads, TL heads, and IHs across layers for the SST-2 dataset on Llama3.2-3B.

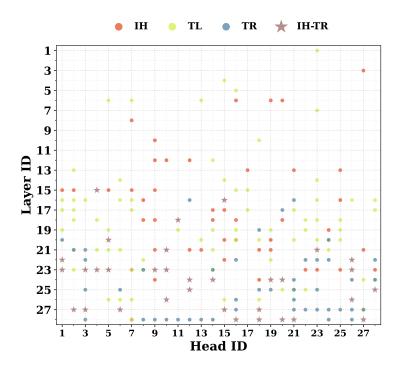


Figure 28: Distribution of the top 10% TR heads, TL heads, and IHs across layers for the SST-2 dataset on Qwen2-7B.

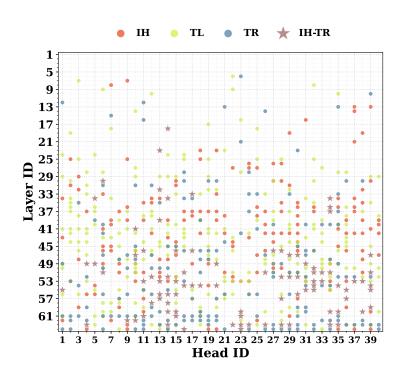


Figure 29: Distribution of the top 10% TR heads, TL heads, and IHs across layers for the SST-2 dataset on Qwen2.5-32B.

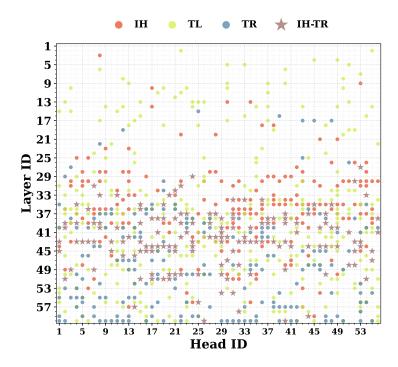


Figure 30: Distribution of the top 10% TR heads, TL heads, and IHs across layers for the SST-2 dataset on Yi-34B.

	TL heads mean layer	IHs mean layer	TR heads mean layer	TL < IHs?	IHs < TR heads?
0.03	14.51	15.69	23.04	0.040803	1.3459e-26
0.05	14.25	15.78	22.59	0.0024930	3.0489e-35
0.10	14.38	15.75	20.96	0.00044495	1.0738e-36

Table 6: Mean layer index of TL heads, IHs, and TR heads across datasets on Llama3.2-3B, with p-values for distribution differences.

	TL heads mean layer	IHs mean layer	TR heads mean layer	TL < IHs?	IHs < TR heads?
0.03	19.06	18.57	24.87	0.31254	2.6148e-28
0.05	18.01	18.62	24.64	0.37007	1.7432e-42
0.10	16.22	18.97	23.38	2.4055e-09	1.8049e-41

Table 7: Mean layer index of TL heads, IHs, and TR heads across datasets on Qwen2-7B, with p-values for distribution differences.

	TL heads mean layer	IHs mean layer	TR heads mean layer	TL < IHs?	IHs < TR heads?
0.03	46.26	44.08	56.80	0.0012338	0.0000e+00
0.05	43.19	44.53	54.89	0.35802	0.0000e+00
0.10	39.42	45.11	51.36	1.7597e-20	0.0000e+00

Table 8: Mean layer index of TL heads, IHs, and TR heads across datasets on Qwen2.5-32B, with p-values for distribution differences.

	TL heads mean layer	IHs mean layer	TR heads mean layer	TL < IHs?	IHs < TR heads?
0.03	36.64	38.16	47.21	0.071766	0.0000e+00
0.05	35.57	38.20	46.90	0.0011017	0.0000e+00
0.10	34.09	38.25	45.20	1.0199e-12	0.0000e+00

Table 9: Mean layer index of TL heads, IHs, and TR heads across datasets on Yi-34B, with p-values for distribution differences.

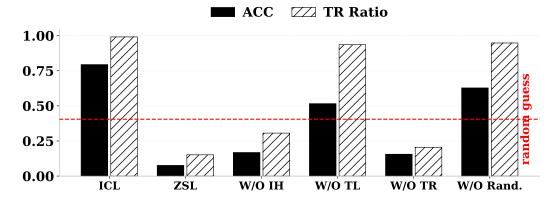


Figure 31: Effects of ablating the top 3% of TR, TL, and IH heads across datasets on Llama3.1-8B.

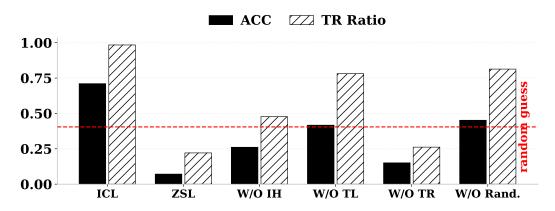


Figure 32: Effects of ablating the top 3% of TR, TL, and IH heads across datasets on Llama3.2-3B.

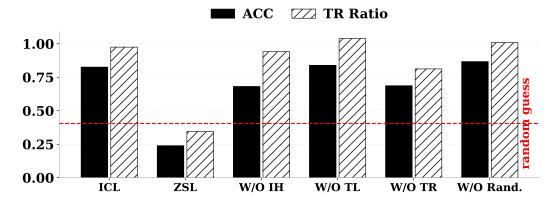


Figure 33: Effects of ablating the top 3% of TR, TL, and IH heads across datasets on Qwen2-7B.

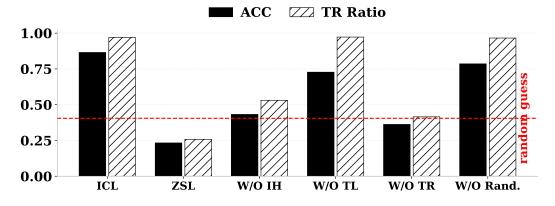


Figure 34: Effects of ablating the top 3% of TR, TL, and IH heads across datasets on Qwen2.5-32B.

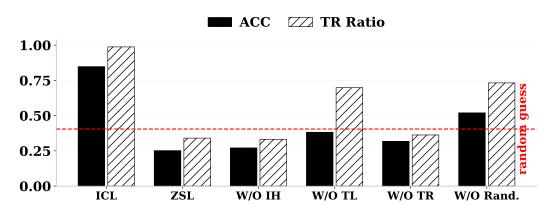


Figure 35: Effects of ablating the top 3% of TR, TL, and IH heads across datasets on Yi-34B.

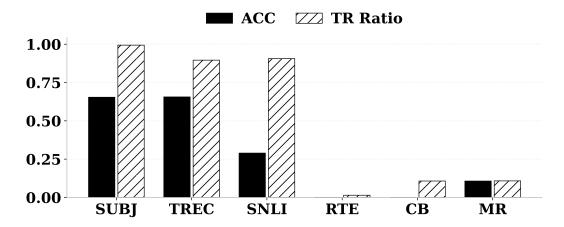


Figure 36: Effects of ablating TR heads identified on SST-2 when transferred to other datasets using Llama3-8B.

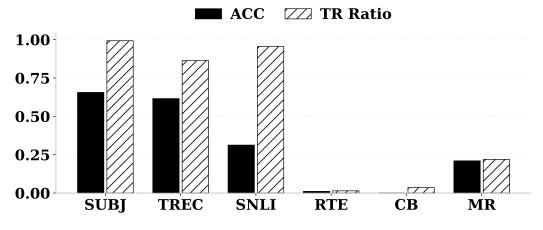


Figure 37: Effects of ablating TR heads identified on SST-2 when transferred to other datasets using Llama3.1-8B.

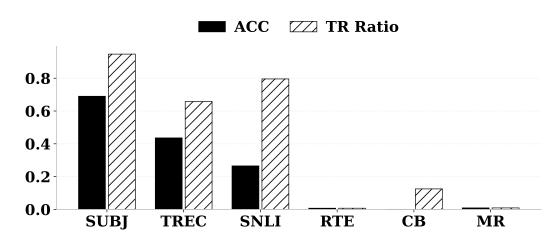


Figure 38: Effects of ablating TR heads identified on SST-2 when transferred to other datasets using Llama3.2-3B.

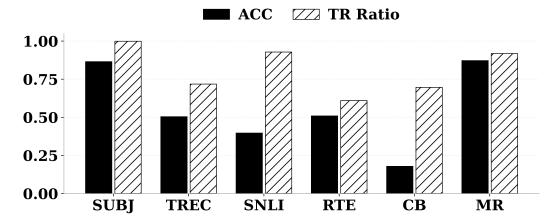


Figure 39: Effects of ablating TR heads identified on SST-2 when transferred to other datasets using Qwen2-7B.

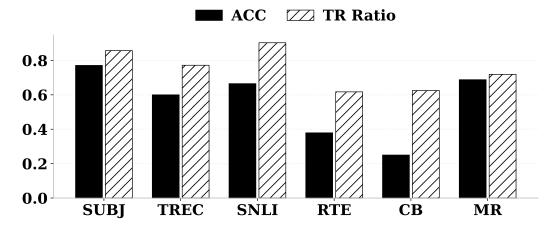


Figure 40: Effects of ablating TR heads identified on SST-2 when transferred to other datasets using Qwen2.5-32B.

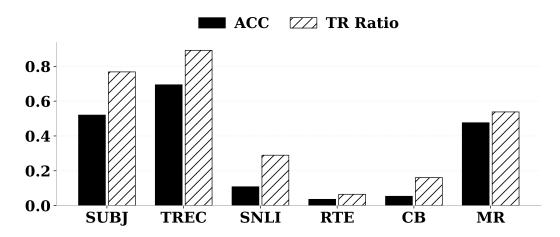
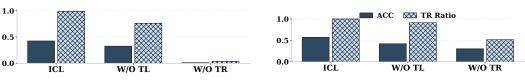


Figure 41: Effects of ablating TR heads identified on SST-2 when transferred to other datasets using Yi-34B.



(a) Ablating TR and TL heads with shuffled demonstration text order.

(b) Ablating TR and TL heads with labels replaced by numbers.

Figure 42: Effects of ablating TR and TL heads under perturbed ICL inputs on Llama3.1-8B.

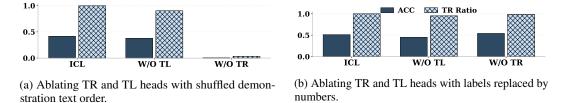


Figure 43: Effects of ablating TR and TL heads under perturbed ICL inputs on Llama3.2-3B.

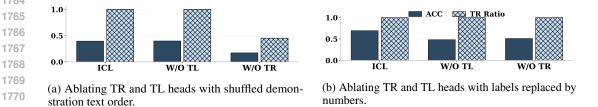


Figure 44: Effects of ablating TR and TL heads under perturbed ICL inputs on Qwen2-7B.

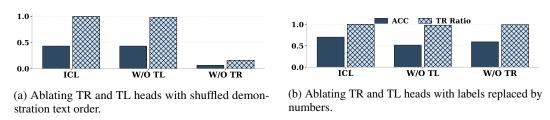


Figure 45: Effects of ablating TR and TL heads under perturbed ICL inputs on Qwen2.5-32B.

stration text order.

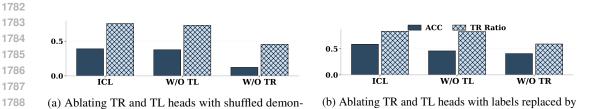


Figure 46: Effects of ablating TR and TL heads under perturbed ICL inputs on Yi-34B.

numbers.

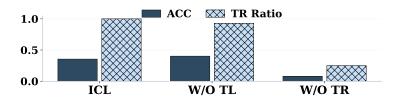


Figure 47: Effects of ablating TR and TL heads on Llama3-8B when demonstration labels are flipped.

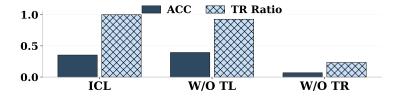


Figure 48: Effects of ablating TR and TL heads on Llama3.1-8B when demonstration labels are flipped.

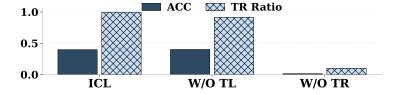


Figure 49: Effects of ablating TR and TL heads on Llama3.2-3B when demonstration labels are flipped.

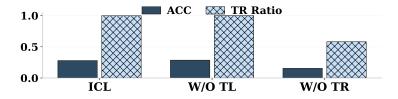


Figure 50: Effects of ablating TR and TL heads on Qwen2-7B when demonstration labels are flipped.

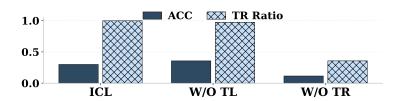


Figure 51: Effects of ablating TR and TL heads on Qwen2.5-32B when demonstration labels are flipped.

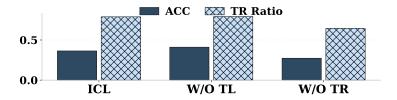


Figure 52: Effects of ablating TR and TL heads on Yi-34B when demonstration labels are flipped.

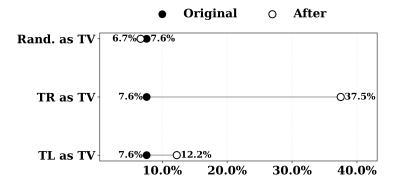


Figure 53: Steering zero-shot hidden states of Llama3.1-8B using task vectors from TR, TL, or random heads.

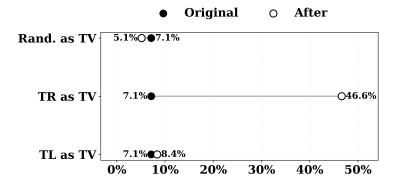


Figure 54: Steering zero-shot hidden states of Llama3.2-3B using task vectors from TR, TL, or random heads.

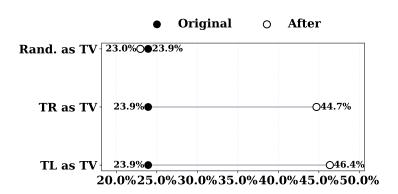


Figure 55: Steering zero-shot hidden states of Qwen2-7B using task vectors from TR, TL, or random heads.

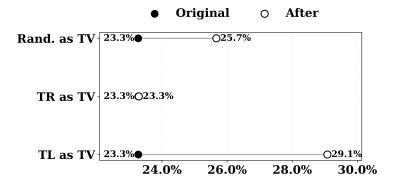


Figure 56: Steering zero-shot hidden states of Qwen2.5-32B using task vectors from TR, TL, or random heads.

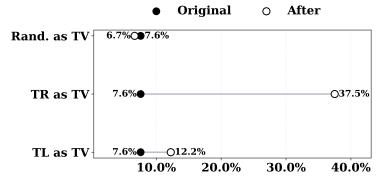


Figure 57: Steering zero-shot hidden states of Yi-34B using task vectors from TR, TL, or random heads.

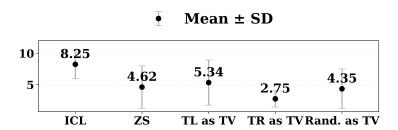


Figure 58: Mean and standard deviation of review ratings with Llama3.1-8B when task vectors from different head types are applied.

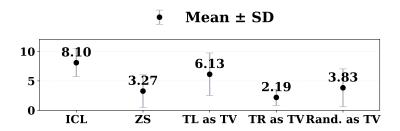


Figure 59: Mean and standard deviation of review ratings with Llama3.2-3B when task vectors from different head types are applied.

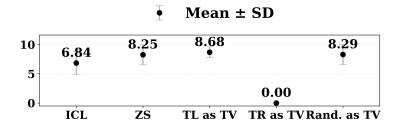


Figure 60: Mean and standard deviation of review ratings with Qwen2-7B when task vectors from different head types are applied.

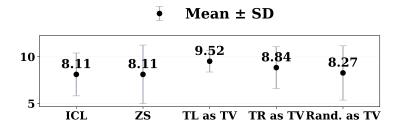


Figure 61: Mean and standard deviation of review ratings with Qwen2.5-32B when task vectors from different head types are applied.

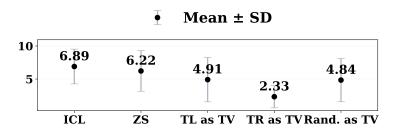


Figure 62: Mean and standard deviation of review ratings with Yi-34B when task vectors from different head types are applied.

Setting	Generated Review
ICL	Poignant character arcs explore relatable themes with depth. Cinematic
	score heightens emotional impact of pivotal scenes. Timely social com-
	mentary addresses important issues with nuance. Strong performances
	deliver believable emotions and connection.
ZS	1. What is the purpose of this review? 2. What is the author's purpose?
	3. How do you know? 4. What is the audience? 5. How do you know?
TL as TV	The movie was very entertaining. I enjoyed the movie and the charac-
	ters. It was a great movie to watch. I would recommend it to others. It
	was a very entertaining movie.
TR as TV	Write a positive review for a movie. The positive review should be
	within 30 words.
Random as TV	Thank you for the positive review. It is always nice to hear when
	someone enjoyed the film. I am glad that you enjoyed the film and that
	you took the time to write a review.

Table 10: Sample reviews generated under different settings with Llama3-8B.

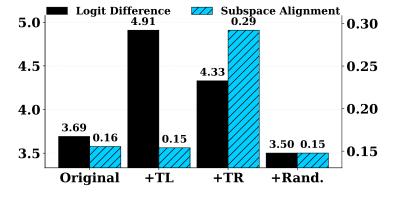


Figure 63: Geometric effects of TR and TL head outputs on hidden states in Llama3.1-8B.

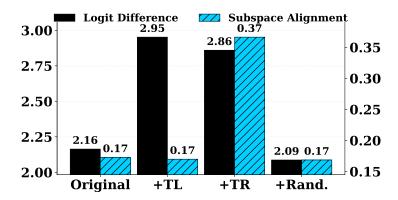


Figure 64: Geometric effects of TR and TL head outputs on hidden states in Llama3.2-3B.

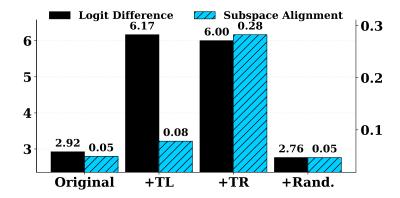


Figure 65: Geometric effects of TR and TL head outputs on hidden states in Qwen2-7B.

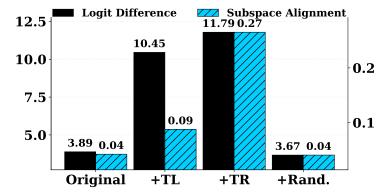


Figure 66: Geometric effects of TR and TL head outputs on hidden states in Qwen2.5-32B.

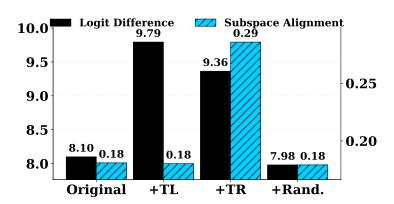


Figure 67: Geometric effects of TR and TL head outputs on hidden states in Yi-34B.

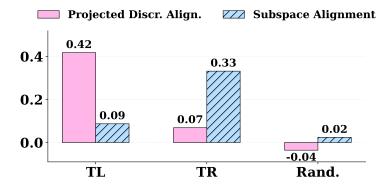


Figure 68: Impact of TL and TR head outputs on hidden states w.r.t. task subspace in Llama3.1-8B.

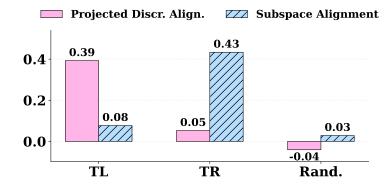


Figure 69: Impact of TL and TR head outputs on hidden states w.r.t. task subspace in Llama3.2-3B.

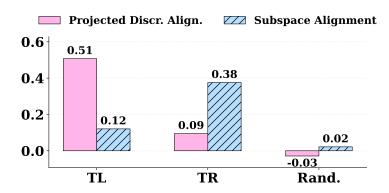


Figure 70: Impact of TL and TR head outputs on hidden states w.r.t. task subspace in Qwen2-7B.

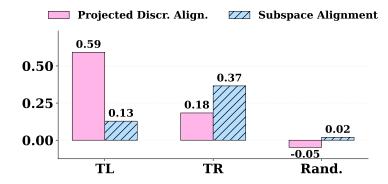


Figure 71: Impact of TL and TR head outputs on hidden states w.r.t. task subspace in Qwen2.5-32B.

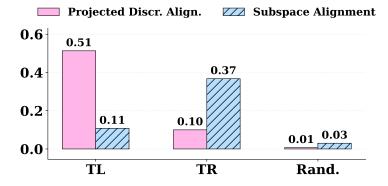
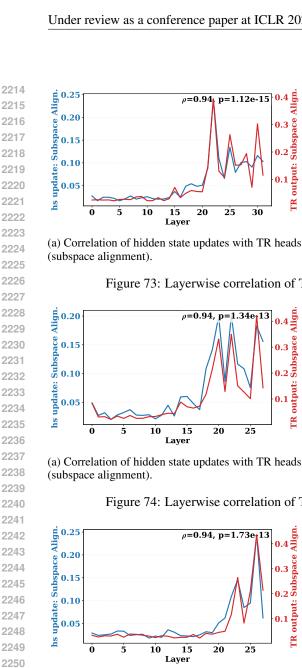
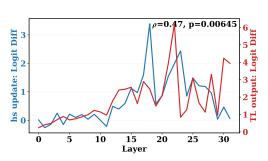


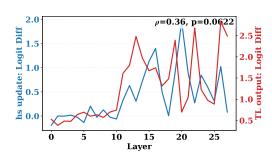
Figure 72: Impact of TL and TR head outputs on hidden states w.r.t. task subspace in Yi-34B.





- (b) Correlation of hidden state updates with TL heads
- (logit difference).

Figure 73: Layerwise correlation of TR and TL head effects on Llama3.1-8B.



- (a) Correlation of hidden state updates with TR heads
- (b) Correlation of hidden state updates with TL heads (logit difference).

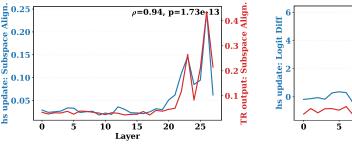
Figure 74: Layerwise correlation of TR and TL head effects on Llama3.2-3B.

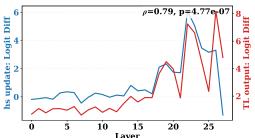
output:

IR

output:

E





(a) Correlation of hidden state updates with TR heads (subspace alignment).

2251

2252

2253

2254 2255

2256

2257

2258

2259

2260

2261

2262

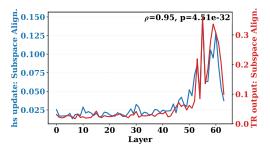
2263

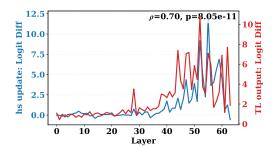
2264

2265

2266 2267 (b) Correlation of hidden state updates with TL heads (logit difference).

Figure 75: Layerwise correlation of TR and TL head effects on Qwen2-7B.





- (a) Correlation of hidden state updates with TR heads (subspace alignment).
- (b) Correlation of hidden state updates with TL heads (logit difference).

Figure 76: Layerwise correlation of TR and TL head effects on Qwen2.5-32B.

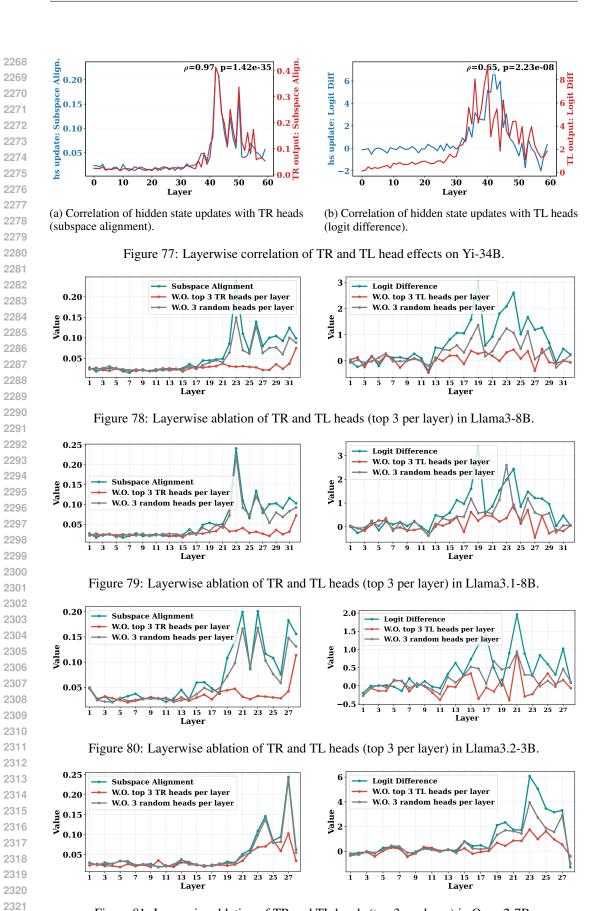
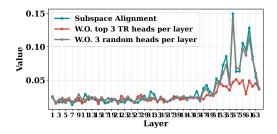


Figure 81: Layerwise ablation of TR and TL heads (top 3 per layer) in Qwen2-7B.



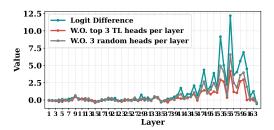
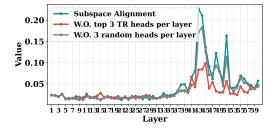


Figure 82: Layerwise ablation of TR and TL heads (top 3 per layer) in Qwen2.5-32B.



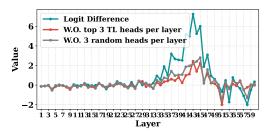


Figure 83: Layerwise ablation of TR and TL heads (top 3 per layer) in Yi-34B.