

Multi-Layer Attention is the Amplifier of Demonstration Effectiveness

Anonymous ACL submission

Abstract

Many studies show that not all demonstrations help in-context learning (ICL), limiting performance. Therefore, in this paper, we analyze why demonstrations become ineffective using gradient flow. By setting the gradient flow to zero, we reveal two cases of ineffectiveness: the model has either already learned the information or it is irrelevant to the query. We also prove that in the multi-layer attention model, effectiveness disparities amplify with depth, directing attention toward effective demonstrations. Building on the above discussion, we propose GRADS, which selects demonstrations via gradient-flow signals and explicitly accounts for already assimilated information. We validate our derivation and GRADS on four prominent LLMs across five mainstream datasets. The experiment confirms that the disparity in demonstration effectiveness is magnified as the model layer increases, substantiating our derivations. Moreover, GRADS achieves a relative improvement of 1.3% on average over the strongest baselines, achieving new SOTA results in the demonstration selection.

1 Introduction

In-Context Learning (ICL) is an effective method for enhancing the performance of Large Language Models (LLMs) (Brown et al., 2020; Dong et al., 2024). By providing demonstrations relevant to the user query in the input, ICL guides the reasoning of LLMs, thereby improving inference performance. Recent years have witnessed many efforts in investigating the internal mechanisms of ICL for guiding the design of ICL methods (Zhou et al., 2024). For example, recent research works (Zhang et al., 2024; Mahankali et al., 2024; Smart et al., 2025) discuss the convergence and convergence speed of ICL, while some other works (Olsson et al., 2022; Li et al., 2024b; Chen et al., 2024) study the function of each part and module of the Transformer (Vaswani et al., 2017) during ICL.

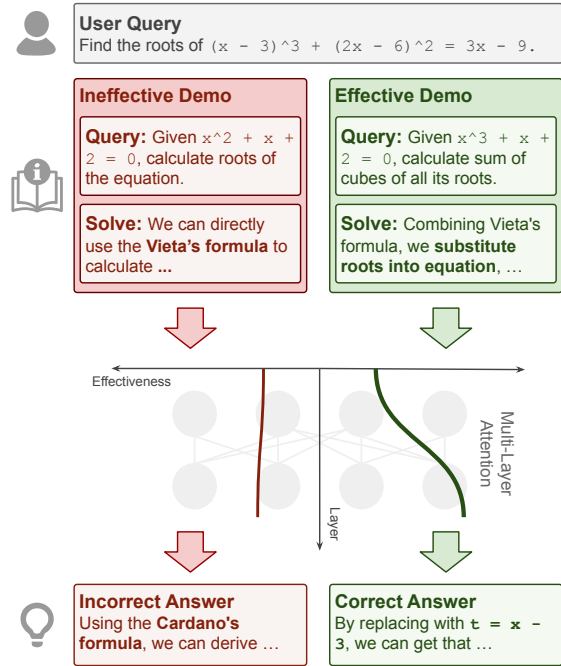


Figure 1: The gradient flow comparison of effective (green) and ineffective demonstrations (red). For the ineffective demonstration, with the increase of model layers, the effectiveness remains at a low level. About the effective ones, the effectiveness is significantly amplified as the layer increases.

As the performance of LLMs continues to improve, a phenomenon emerges that there exist ineffective demonstrations to the given query, which cannot enhance reasoning performance during ICL (DeepSeek-AI, 2025; Wang et al., 2025c). However, existing research on the mechanisms of ICL is predominantly based on the assumption that the given demonstrations are not ineffective (Liu et al., 2025; Cho et al., 2025), which limits the exploration of how to enhance the performance of ICL. Therefore, in this paper, we aim to answer the following research questions (RQs): 1) *What makes the demonstration ineffective?* 2) *How does the effectiveness of the demonstration affect the ICL performance?* and 3) *How to select effective demonstrations to enhance ICL performance?*

Following Wang et al. (2023); Liu and Deng (2025), we employ the gradient flow for our investigation: the greater the gradient flow from the demonstration to the generated answer, the demonstration is more effective. We first study the factors that determine the magnitude of the gradient flow in a single-layer linear self-attention model (LSA), thereby providing the reasons for an ineffective demonstration: the information in the demonstration has already been learned by LLMs, or the demonstration is irrelevant to the user query (RQ1). We then prove that in a multi-layer LSA, **the ratio of the accumulated gradient flow among demonstrations is amplified along the increase of the number of layers** (RQ2), as shown in Figure 1.

Considering that existing demonstration selection methods mainly focus on the relevance between the demonstration and the query, which could select ineffective demonstrations (Wei et al., 2024; Chen et al., 2023a; Wang et al., 2025b). Therefore, based on the above discussion, we propose GRADS, which selects the demonstration that maximizes the gradient flow with the given query (RQ3). Measure whether LLMs utilize the information in a given demonstration through the gradient flow to ensure that the demonstrations provided are effective in reasoning, enhancing performance.

To validate our conclusions and the effectiveness of GRADS, we conduct experiments on five mainstream datasets and four mainstream LLMs. Experimental analysis shows that the ratio of the gradient flow between effective and ineffective demonstrations increases with the model layers, proving that the multi-layer Transformer indeed acts as an amplifier of demonstration effectiveness. Furthermore, the experimental results of GRADS indicate that, compared to the best demonstration selection baselines, GRADS achieves a relative improvement of 1.3% on average, achieving new SOTA results on the demonstration selection.

In summary, our contributions include:

- We prove that a demonstration is ineffective when the information has already been learned by the LLMs or is irrelevant to the given user query.
- We theoretically and empirically analyze that the multi-layer structure amplifies the ratio of the demonstration effectiveness.
- We propose GRADS, a gradient-based demonstration selection method that achieves new SOTA results on the demonstration selection.

2 Analysis of Demonstration Effectiveness

In this section, we discuss *why* and *how* LLMs ignore ineffective demonstrations in ICL from a gradient flow perspective. We first provide definitions for necessary mathematical notations and concepts (§2.1). Then, we discuss the gradient flow in a single-layer linear self-attention network and why LLMs distinguish ineffective demonstrations (§2.2). Subsequently, we discuss the gradient flow in a multi-layer linear self-attention network and how LLMs ignore ineffective demonstrations (§2.3). The main findings of our analysis are summarized in Table 1. All proofs in this section are provided in Appendix A, and we follow Wang et al. (2023) to calculate the gradient flow.

2.1 Preliminary

Without loss of generality, we primarily focus on the 1-shot setting for mechanism analysis. Following Zhang et al. (2024), we denote a demonstration as $d = (d_x \ d_y)$, where $d_x, d_y \in \mathbb{R}^e$ represents the input and output embedding vector of the demonstration, respectively. Here, e is the embedding dimension. We denote a user query as $q = (q_x \ q_y)$, where $q_x \in \mathbb{R}^e$ is the query input embedding vector, and we set $q_y = 0$ to indicate that the corresponding answer to the query is not provided in the input. We denote the complete network input as $E = (d \ q)$. In this paper, we use $\|M\|$ to denote the Frobenius norm (Wu et al., 2022) of a given matrix M .

Following previous work (Zhang et al., 2024), we define the linear self-attention network (LSA):

$$f_{LSA}(E; \theta) = E + W^{PV} E \cdot \frac{E^\top W^{KQ} E}{\rho} \quad (1)$$

$W^{PV} \in \mathbb{R}^{e \times e}$ is the combined projection and value matrix of Attention, and $W^{KQ} \in \mathbb{R}^{e \times e}$ is the combined key and query matrix, where the specific definitions are consistent with Zhang et al. (2024). ρ is the normalization coefficient, where we set $\rho = 1$ in this paper following Zhang et al. (2024). We denote all network parameters as $\theta = \{W^{PV}, W^{KQ}, \rho\}$. It can be observed that Equation 1 is the result of replacing the activation function in a single-layer attention network of a Transformer with a linear function. Specifically, we denote $\hat{q}_y(E; \theta)$ as the predicted answer of the given E and θ , abbreviated as \hat{q}_y , which is the vector corresponding to the last column of the output from the LSA model $f_{LSA}(E; \theta)$.

RQ	Finding	Evidence
RQ1: What makes the demonstration ineffective?	A demonstration is ineffective if the information it contains has already been learned by LLMs or is irrelevant to the user query.	Equation 2
RQ2: How does the effectiveness of the demonstration affect the ICL performance?	With deeper layers, the ratio of the accumulated gradient flow between effective and ineffective demonstrations widens.	Theorem 1

Table 1: The main research questions (RQs), findings, and corresponding evidence of our analysis.

Following previous work (Wang et al., 2023), for a multivariate function f and one of its independent variables x , we define the magnitude of the gradient flow of f with respect to x as the partial derivative of f with respect to x , which is $\nabla_x f = \frac{\partial f}{\partial x}$. The magnitude of the gradient flow measures the influence of the change in x on f . Specifically, in the context of ICL, the gradient flow of the output \hat{q}_y with respect to an input demonstration d reflects the contribution of that demonstration to the answer. Consequently, it can indicate how much information from the demonstration is utilized during generating the answer.

2.2 Single-Layer Linear Self-Attention

We first analyze the gradient flow in a single-layer network. It can be shown that the contribution of the input demonstration to the gradient flow for answer generation is:

$$\nabla_d \hat{q} = (W^{PV} d)(q^\top W^{KQ})^\top + (d^\top W^{KQ} q) W^{PV} \quad (2)$$

Equation 2 formally defines the gradient flow value of the single layer LSA, which indicates how much information from the demonstration is used for answer generation. We can see that, given a certain query q , the factors determining the magnitude of the equation can be divided into two categories: (i) $d^\top W^{KQ} q$: The similarity between the demonstration and the user query. (ii) $W^{PV} d$: Whether the model has already learned the information in the demonstration (Petroni et al., 2019; Roberts et al., 2020). Therefore, the determination of whether to use the information from a given demonstration mainly depends on 1) the similarity between the demonstration and the query, and on 2) whether the information in the demonstration has already been learned by the model.

Based on the discussion above, as the base of the following discussion, we propose to use the parameters in Equation 2 to formally define the effectiveness of demonstrations with the given query and model as follows:

Definition 1 (Demonstration Effectiveness). *Given a query q and model parameters θ , if two demonstrations d_1 and d_2 satisfy that:*

$$\begin{aligned} \|W^{PV} d_1\| &\geq \|W^{PV} d_2\| \\ \|d_1^\top W^{KQ} q\| &\geq \|d_2^\top W^{KQ} q\| \end{aligned}$$

then we say that d_1 is more effective than d_2 with respect to q and θ , denoted as:

$$d_1 \succ_{q;\theta} d_2$$

In Definition 1, we require a demonstration to be more effective than the other if it contains more knowledge that the model has not learned ($\|W^{PV} d_1\| \geq \|W^{PV} d_2\|$), and it is more similar to the query ($\|d_1^\top W^{KQ} q\| \geq \|d_2^\top W^{KQ} q\|$). If only one of these conditions holds, it is difficult to compare the effectiveness of the demonstration because it is not possible to determine whether the knowledge or the similarity to the query has a greater impact on performance. We experimentally evaluate the impact of two conditions in Definition 1 on ICL performance in Figure 2.

2.3 Multi-Layer Linear Self-Attention

Considering that current mainstream LLMs consist of multi-layer modules, in this part, we discuss the gradient flow of the multi-layer linear self-attention network. Let L be the total number of layers in the network. We denote the input to the l -th layer as $E^{(l)} = (d^{(l)} \quad q^{(l)})$ and its parameters as $\theta^{(l)}$, the corresponding predicted answer of the l -th layer is $\hat{q}^{(l)}$. Since in the multi-layer LSA, the output of the $(l-1)$ -th layer is the input of the l -th layer:

$$E^{(l)} = f_{LSA}(E^{(l-1)}; \theta^{(l-1)}) \quad (3)$$

According to the chain rule, we can derive that the gradient flow of the whole model is the product of the gradient flow of each layer:

$$\frac{\partial \hat{q}_y^{(L)}}{\partial d^{(0)}} = \frac{\partial \hat{q}_y^{(L)}}{\partial E^{(L-1)}} \times \frac{\partial E^{(L-1)}}{\partial E^{(L-2)}} \times \cdots \times \frac{\partial E^{(1)}}{\partial d^{(0)}} \quad (4)$$

Due to the complexity of the structure, we only provide a qualitative analysis for the multi-layer model. Based on Equation 2, we know that the magnitude of each term in Equation 4 is positively correlated with the demonstration effectiveness, i.e.:

Lemma 1. *Let an L -layer LSA network be given. For every layer $l = 1, \dots, L$, assume there exist strictly increasing functions*

$$g_l : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}, \quad h_l : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0},$$

such that for every demonstration d and query q , the following equalities hold:

$$\begin{aligned} & \|W^{PV,(l)} d^{(l)}\| \\ &= g_l(\|W^{PV,(l-1)} d^{(l-1)}\|) \end{aligned} \quad (5)$$

$$\begin{aligned} & \|(d^{(l)})^\top W^{KQ,(l)} q^{(l)}\| \\ &= h_l(\|(d^{(l-1)})^\top W^{KQ,(l-1)} q^{(l-1)}\|) \end{aligned} \quad (6)$$

If two demonstrations satisfy $d_1 \succ_{q;\theta^{(0)}} d_2$ at the input layer, then for every $l = 1, \dots, L$, we have $d_1^{(l)} \succ_{q;\theta^{(l)}} d_2^{(l)}$.

The condition in Lemma 1 assumes that each layer of the model has a consistent effect on its input. Based on Lemma 1, we know that for each layer, a more effective input demonstration results in a larger corresponding gradient flow in all layers. **It is worth noting that for the same demonstration d , Lemma 1 does not guaranty that the gradient flow at layer l is always greater than that at layer $l - 1$.** This is because the magnitude of the gradient flow also depends on the model parameters $\theta^{(l)}$ of each layer.

Based on Lemma 1, we can deduce that since the gradient flow is a partial derivative, a larger gradient flow leads to a greater change in the output of subsequent layers. Therefore, as the number of model layers increases, the difference in gradient flow between effective and ineffective demonstrations also becomes larger.

Theorem 1 (Multi-Layer Attention is the Amplifier of Demonstration Effectiveness). *For a given user query q and a model θ , let d_1 and d_2 be two demonstrations with corresponding inputs E_1 and E_2 . If the condition of Lemma 1 holds, for any $L \geq l_1 > l_2 \geq 1$, we can conclude that the following inequalities hold:*

$$\frac{\|\nabla_{d^{(0)}} \hat{q}_y^{(l_1)}(E_1; \theta)\|}{\|\nabla_{d^{(0)}} \hat{q}_y^{(l_1)}(E_2; \theta)\|} \geq \frac{\|\nabla_{d^{(0)}} \hat{q}_y^{(l_2)}(E_1; \theta)\|}{\|\nabla_{d^{(0)}} \hat{q}_y^{(l_2)}(E_2; \theta)\|}$$

Theorem 1 shows that as the layers increase, the **ratio of the accumulated gradient flow** of different demonstrations also becomes increasingly large. This indicates that *the multi-layer LSA acts as an amplifier for the effectiveness.*

3 GRADS

In this section, we introduce a demonstration selection method, named GRADS, based on gradient flow, selecting the demonstration with the largest flow to the given query as the selection result. The prompt we used is shown in Appendix B.1. Considering Theorem 1, in actual selection, we use the gradient flow of the last layer as the selection metric. Because the effectiveness differences between demonstrations are sufficiently amplified, it allows a better distinction between effective and ineffective demonstrations. We discuss the efficiency of GRADS in Appendix C.1 and Appendix D.3.

However, performing a full inference pass on the model to obtain the gradient flow for each user query results in low computational efficiency. From Equation 2, it can be observed that in the calculation of the gradient flow, the computations for the demonstration and the user query are relatively independent. Therefore, we can first compute the encoded vectors for the demonstrations and the user query separately, and then use these results to calculate the magnitude of the gradient flow through matrix operations. This approach significantly reduces the computational overhead, thereby enhancing the efficiency of demonstration selection.

Specifically, for a given demonstration pool, we first input each demonstration individually to extract the encoding result of the final layer as \hat{d} . Then, for each user query, we also input it to extract the encoding result of the final layer as \hat{q} . Subsequently, the computed \hat{d} and \hat{q} are substituted into Equation 2 to obtain $\nabla_{d \hat{q}_y^{(L)}}$. Since $\nabla_{d \hat{q}_y^{(L)}}$ is a $2 \times e$ matrix, containing the gradient flows for both the input and output parts of the demonstration, and considering that both parts are equally important, we use $\|\nabla_{d \hat{q}_y^{(L)}}\|$ as the metric for demonstration selection to balance their importance.

4 Experiment

Our experiments are primarily divided into three parts: (i) Introduction of the experimental settings and baselines. (ii) Verification of Theorem 1 and its related corollaries. (iii) Validation of effectiveness and the impact of different factors on GRADS.

4.1 Experiment Setup

Dataset To thoroughly validate our analytic conclusions and the effectiveness of our proposed method, we conduct experiments on five main-stream datasets that span various tasks and domains, including: (i) math: GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021); (ii) reasoning: ARC-Challenge (Yadav et al., 2019) and MMLU-Pro (Wang et al., 2024b); (iii) domain-specific: Amazon Review (Ni et al., 2019) and FinQA (Tao et al., 2024). Detailed descriptions of these datasets are provided in Appendix B.2. We employ Exact Match (EM) as the evaluation metric across all experimental datasets.

Model We validate our discovery and method on four LLMs: Llama-3.1-8B/70B-Insturct (Llama-3.1-8B/70B) (Aaron Grattafiori, 2024), Deepseek-R1-Distill-Llama-8B (Llama-R1-8B) (DeepSeek-AI, 2025), and Qwen3-8B (An Yang, 2025)¹. Our selection of models encompasses various scales, series, and capabilities for generating long chains of thought (Long-CoT, Llama-R1 and Qwen3) (Chen et al., 2025a), which can fully evaluate the effectiveness of our conclusion.

Baseline We compare GRADS with following demonstration selection baselines: BM25 (Li et al., 2023), MMR (Ye et al., 2023), MoD (Wang et al., 2024a), Influence (Nguyen and Wong, 2023), DICK (Kapuriya et al., 2025), GenICL (Zhang et al., 2025c), TopicK (Kweon et al., 2025), and ICL-Grad (Zhang et al., 2025d). The above covers the most advanced current demonstration selection methods. Detailed introductions and configurations for each baseline are provided in Appendix B.3.

Implementation Detail For the mechanism analysis experiments, we employ a 1-shot setting to align with the setup in our theoretical analysis. For the validation of GRADS, we follow previous works (Wang et al., 2025b) and adopt a 3-shot setting to ensure the comparability of the final performance. Following DeepSeek-AI (2025), we set the maximum generation length to 32768, and for each question, we sample a single answer. Our experiments are performed on a single A100-80G GPU, with the selection and inference for each dataset taking approximately one hour on average. All implementations are based on Transformers (Wolf et al., 2020) and vLLM (Kwon et al., 2023).

¹We utilize the thinking mode of Qwen3 for a comprehensive evaluation.

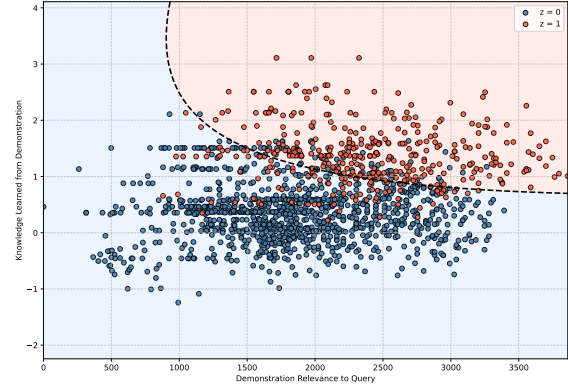


Figure 2: The ICL performance of Llama-3.1-8B across all datasets with different demonstration relevance (X-axis) and unlearned knowledge within the demonstration (Y-axis). Red points denote correct prediction and blue points incorrect predictions. The dashed line represents the decision boundary generated with polynomial logistic regression (Buitinck et al., 2013).

4.2 Experiment of Mechanism Analysis

4.2.1 The Factors Making Demonstrations Ineffective (RQ1)

To verify the condition for low gradient flow in the discussion regarding Equation 2, we record the statistics of the model performance as a function of the changes in demonstration relevance ($d^T W^K Q q$) and learned knowledge ($W^{PV} d$). The experimental results are shown in Figure 2. From the figure, we can observe that: (i) As the relevance of the demonstrations and the knowledge learned from them increase, the number of correctly solved data points also increases, which verifies the conclusion from Equation 2 that the performance of ICL is positively correlated with these two factors. (ii) ICL only begins to correctly solve the given user queries after the relevance of the demonstrations and the knowledge learned from them surpasses a certain threshold, which indicates that a sufficient amount of effective information relevant to the user query is necessary to enable the model to perform correct reasoning process.

4.2.2 Gradient Flow is Amplified as the Number of Layer Increases (RQ2)

To validate the correctness of Lemma 1 and Theorem 1, we compute the gradient flow in different settings. For a given dataset, we first select the data points that are incorrectly predicted under the 0-shot setting. From this subset, we then separate the data into two groups based on the 1-shot performance: those that are correctly predicted (*effective demonstrations*) and those that remain incorrectly

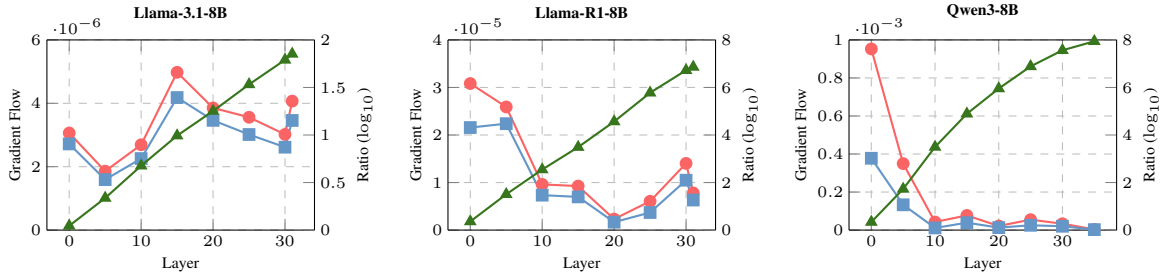


Figure 3: The average gradient flow (left y-axis) and the ratio of the accumulated flow $\frac{\|\nabla_{d(0)} \hat{q}_y^{(l_1)}(E_1; \theta)\|}{\|\nabla_{d(0)} \hat{q}_y^{(l_1)}(E_2; \theta)\|}$ of Theorem 1 (right y-axis) cross different datasets under the i -th layer of each model, where E_1 denotes the input of the effective demonstrations and E_2 denotes the ineffective ones. blue and Red points denote the gradient flow under each layer of the ineffective and effective demonstrations. Green points denote the ratio of the accumulated flow.

389 predicted (*ineffective demonstrations*). We then
 390 compute the average gradient flow for each layer
 391 across these two groups. The number of selected
 392 effective and ineffective data points on each dataset
 393 is detailed in Appendix B.4.

394 The experimental results are presented in Figure
 395 3. From the figure, we can observe the following:
 396 (i) Across all settings, the average gradient flow
 397 of effective demonstrations is stronger than
 398 that of ineffective demonstrations at every layer,
 399 which validates the conclusion of Lemma 1. Fur-
 400 thermore, the ratio of gradient flow between effec-
 401 tive and ineffective demonstrations increases
 402 with the layer depth, which supports Theorem 1.
 403 (ii) The magnitude of the gradient flow, as well
 404 as the accumulated ratio, is significantly more pro-
 405 nounced in the Long-CoT model compared to other
 406 models. This suggests that such models are more
 407 adept at capturing information from demonstrations
 408 and can better distinguish between effective and
 409 ineffective demonstrations, showing a higher sensi-
 410 tivity to the useful information contained within
 411 the demonstrations. (iii) In all models, there are
 412 specific layers where the gradient flow exhibits a
 413 significant increase. This indicates that the learn-
 414 ing from demonstrations is primarily concentrated
 415 in these particular layers. While the specific lay-
 416 ers differ across models, a strong gradient flow is
 417 consistently observed in the final few layers, sug-
 418 gesting that the model pays special attention to the
 419 information in the demonstrations when generat-
 420 ing the output. (iv) The curved boundary shows
 421 that either high relevance with non-trivial learn-
 422 ing or moderate relevance with substantial learning
 423 yields effectiveness during ICL, while high rele-
 424 vance alone with little new knowledge often re-
 425 mains ineffective. Hence, demonstration selection
 426 should combine relevance with a learning signal
 427 rather than rely on the relevance alone.

4.3 Experiment of GRADS (RQ3) 428

4.3.1 GRADS is More Effective than Baselines 429

430 To validate the effectiveness of GRADS, we con-
 431 duct a comparative analysis against several base-
 432 lines, with the experimental results presented in
 433 Table 2. The results show that GRADS achieve
 434 a relative improvement of 1.3% on average over
 435 the best-performing baseline across various mod-
 436 els and datasets, achieving new demonstration selec-
 437 tion SOTA, which substantiates its efficacy and
 438 generalizability. Furthermore, a deeper analysis of
 439 the results reveals several key observations:

440 **Baseline** In some settings, the performance of
 441 other baselines does not exceed that of simpler
 442 methods like BM25, and is even inferior to the zero-
 443 shot approach in some cases. This suggests that
 444 methods based on the similarity between demon-
 445 strations and the user query do not guarantee an
 446 enhancement in ICL performance due to the fact
 447 that the model could have already been exposed to
 448 the information present in the demonstrations, and
 449 irrelevant information within these demonstrations
 450 could consequently mislead the model’s reasoning
 451 process. In contrast, GRADS ensures that the infor-
 452 mation in the selected demonstrations is effectively
 453 utilized during inference, securing the effectiveness
 454 of the selected demonstrations.

455 **Model** The most significant performance im-
 456 provement from GRADS is observed in models
 457 without Long-CoT in most settings (e.g., Llama-
 458 3.1). This is because such models do not employ
 459 Long-CoT, rendering them less capable of effec-
 460 tively leveraging the useful information within the
 461 demonstrations. Consequently, their performance
 462 is more dependent on the quality of the demon-
 463 strations compared to Llama-R1-8B and Qwen3-8B.

Model	Method	GSM8K	MATH	ARC-C	MMLU-Pro	Amazon	FinQA
Llama-3.1-8B	Zero	83.7	47.0	82.0	50.4	63.5	48.1
	BM25 (Li et al., 2023)	84.1	44.6	84.7	54.0	68.4	50.5
	MMR (Ye et al., 2023)	83.9	45.2	85.1	53.6	67.1	50.2
	MoD (Wang et al., 2024a)	84.0	47.0	84.9	53.8	68.7	50.9
	Influence (Nguyen and Wong, 2023)	84.1	47.4	84.6	54.3	68.2	50.8
	DICL (Kapuriya et al., 2025)	84.3	47.4	85.0	54.7	68.9	51.3
	GenICL (Zhang et al., 2025c)	84.5	47.6	85.2	55.0	69.1	51.6
	ICL-Grad (Zhang et al., 2025d)	84.1	47.0	84.8	53.2	67.5	50.4
	TopicK (Kweon et al., 2025)	83.9	47.0	84.1	52.6	66.8	49.8
	GRADS	85.6	48.2	86.3	56.0	70.0	52.3
Llama-3.1-70B	Zero	87.3	63.6	85.2	56.3	69.0	58.3
	BM25 (Li et al., 2023)	86.9	63.2	85.9	57.7	70.8	63.4
	MMR (Ye et al., 2023)	87.2	63.8	86.3	57.4	70.4	63.2
	MoD (Wang et al., 2024a)	87.4	64.2	86.2	57.6	71.0	63.6
	Influence (Nguyen and Wong, 2023)	87.5	64.0	86.6	58.0	71.2	64.0
	DICL (Kapuriya et al., 2025)	87.6	64.0	86.4	57.9	71.1	63.8
	GenICL (Zhang et al., 2025c)	87.8	64.4	86.7	58.3	71.5	64.3
	ICL-Grad (Zhang et al., 2025d)	87.4	64.0	86.0	57.0	70.1	62.0
	TopicK (Kweon et al., 2025)	87.3	63.6	85.7	56.6	69.6	60.8
	GRADS	88.9	65.2	87.8	59.0	72.4	65.1
Llama-R1-8B	Zero	86.0	74.6	83.5	58.2	61.5	60.2
	BM25 (Li et al., 2023)	80.2	74.0	84.0	52.9	65.2	62.9
	MMR (Ye et al., 2023)	85.5	73.0	84.2	58.4	65.9	63.2
	MoD (Wang et al., 2024a)	85.7	74.4	84.0	57.9	65.8	62.6
	Influence (Nguyen and Wong, 2023)	85.8	74.4	84.1	58.6	66.1	63.0
	DICL (Kapuriya et al., 2025)	85.9	74.4	84.3	58.5	66.0	63.4
	GenICL (Zhang et al., 2025c)	86.1	74.6	84.6	58.7	66.2	63.6
	ICL-Grad (Zhang et al., 2025d)	86.0	74.6	83.8	58.3	65.4	62.0
	TopicK (Kweon et al., 2025)	86.0	74.6	83.6	58.2	64.7	61.4
	GRADS	87.2	75.6	85.7	59.5	67.0	64.4
Qwen3-8B	Zero	92.9	76.2	89.2	64.9	62.5	59.0
	BM25 (Li et al., 2023)	92.2	77.2	91.0	65.3	71.5	65.5
	MMR (Ye et al., 2023)	92.5	77.6	90.7	66.5	74.0	66.6
	MoD (Wang et al., 2024a)	92.6	78.2	90.8	66.1	72.0	66.0
	Influence (Nguyen and Wong, 2023)	92.8	78.0	91.2	66.6	73.0	66.5
	DICL (Kapuriya et al., 2025)	92.7	78.0	91.1	67.0	73.8	67.1
	GenICL (Zhang et al., 2025c)	93.0	78.4	91.5	67.6	74.1	67.3
	ICL-Grad (Zhang et al., 2025d)	92.9	76.4	90.2	65.2	71.8	65.0
	TopicK (Kweon et al., 2025)	92.9	76.2	89.7	65.0	70.9	64.2
	GRADS	94.2	79.4	92.6	68.5	75.0	68.2

Table 2: The performance of GRADS compared with different baselines. ARC-C denotes ARC-Challenge, Amazon denotes Amazon Review. Zero denotes the zero-shot setting. The best result under each setting is marked in **bold**.

Dataset GRADS yielded more substantial performance gains on the MATH and MMLU-Pro datasets. This is because such two datasets demand a higher degree of specialized knowledge, which could not be contained in LLMs. GRADS, being more effective at selecting demonstrations that contain the requisite knowledge for the model, achieves a more significant performance improvement in these contexts.

4.3.2 The Performance of GRADS is Positively Correlated to Model Layer

To validate the conclusion that the disparity in gradient flow between effective and ineffective demonstrations increases with network layer in Theo-

Dataset	Layer							
	0	5	10	15	20	25	30	31
GSM8K	83.2	84.4	84.9	85.1	85.2	85.6	85.2	85.6
MATH	44.4	44.8	45.8	45.2	46.2	47.6	46.8	48.2
ARC-Challenge	82.1	82.3	82.1	82.8	83.4	84.0	85.1	86.3
MMLU-Pro	50.8	51.5	53.6	54.8	54.8	55.6	56.2	56.0

Table 3: The performance of GRADS using the gradient flow of different layer on Llama-3.1-8B. The best result under each setting is marked in **bold**.

rem 1, thereby enhancing the demonstration selection capability, we conduct the experiment. We propose to employ GRADS to select demonstrations using gradient flows from different layers for demonstration selection. Considering the experimental cost, we evaluate under each 5 layer.

The experimental results are presented in Table 3. From the table, we can observe the following: (i) Across all datasets, there is a general upward trend in model performance as the number of layers increases, which confirms the conclusion of Theorem 1 that as the network layers increase, the difference in effectiveness among demonstrations is amplified, which enables a better selection of effective demonstrations. (ii) However, the performance of GRADS does not monotonically increase with the number of layers, since the models used in our experiments are more complex than those in our theoretical analysis. For instance, the presence of residual streams (He et al., 2016) can diminish the amplifying effect of multi-layer transformers on the gradient flow, which could lead the model to erroneously select ineffective demonstrations, resulting in a decline in performance.

5 Related Works

In-Context Learning ICL can effectively improve LLM reasoning performance by adding query-relevant demonstrations to the prompt with additional training (Brown et al., 2020; Dong et al., 2024). Prior work falls into three axes: synthesizing, selecting, and utilizing demonstrations. To cut labeling cost, LLM-based synthesis draws on existing examples, related information, and similar-task data (Long et al., 2024a; Su et al., 2024; He et al., 2024; Chen et al., 2023b; Wang et al., 2025b; Chen et al., 2025b; Wang et al., 2025a). Selection has progressed from gram-based heuristics to model-aware relevance estimation (Rubin et al., 2022; Li et al., 2023; Peng et al., 2024; Wang et al., 2025e; Yang et al., 2023; Ye et al., 2023; Wang et al., 2024a; Zhang et al., 2025c). Utilization improves effectiveness and reduces inference cost via demonstration reordering or encoding-and-injection (Lu et al., 2022; Pham et al., 2025; Li et al., 2024a; Wang et al., 2025d; Li et al., 2025). Across these strands, a unifying theme is maximizing task alignment under tight context budgets—balancing coverage, diversity, and cost during inference.

However, existing demonstration selection methods mainly consider query–demonstration relevance, overlooking that demonstrations may be ineffective if the model has already learned their information. Thus, we propose GRADS, a gradient-flow–based approach that selects demonstrations contributing significant new information during inference, thereby improving ICL performance.

Mechanism of In-Context Learning Understanding ICL’s mechanism guides performance improvements and clarifies model reasoning (Zhou et al., 2024). Prior work falls into four strands: theory, architecture, data, and inference. Theoretical studies establish ICL’s effectiveness, like convergence and rates (Wies et al. (2023); Huang et al. (2024); Yang et al. (2024); Smart et al. (2025); Fu et al. (2024); Huang et al. (2025); Vladymyrov et al. (2024)—often via LSA-based analyses (Zhang et al., 2024; Mahankali et al., 2024; Lu et al., 2024). Architectural analyses show attention chiefly drives ICL (Olsson et al., 2022; Chen et al., 2024; Oko et al., 2024), while MLPs play a supporting role (Li et al., 2024b; Nguyen and Reddy, 2025). Data-centric studies attribute ICL emergence to task diversity and a training-phase shift from in-weight to in-context learning (Raventos et al., 2023; Wibisono and Wang, 2024; Zhang et al., 2025a; Goddard et al., 2025; Zhang et al., 2025b; Singh et al., 2025; Park et al., 2025; de Wynter, 2025). Inference-time analyses examine observed behaviors, including factor influences and changes in internal computations (Bigelow et al., 2024; Shi et al., 2024; Lin and Lee, 2024; Long et al., 2024b; Sia et al., 2024; Zhao et al., 2024).

However, most studies of the ICL mechanism assume that provided demonstrations are effective. In practice, however, demonstrations can be ineffective, yielding no performance gain (DeepSeek-AI, 2025; Wang et al., 2025c). In this work, we examine this phenomenon by identifying two key factors, including irrelevant content or information already learned by the model, and reveal that multi-layer transformers amplify demonstration effectiveness.

6 Conclusion

In this paper, we first discuss that a demonstration is ineffective because the its information has been learned by the model or is irrelevant to the query. We then demonstrate that as the model layer increases, the ratio of the accumulated gradient flow among demonstrations is amplified. Analytical experiments show that as the model layer increases, the difference in effectiveness among demonstrations is magnified, which corroborates our theoretical derivations. Second, we present GRADS, which achieves an average relative performance improvement of 1.3% compared to existing demonstration selection methods, achieving new SOTA results, which proves the effectiveness of our method.

584 Limitations

585 (i) We have not validated GRADS on additional
586 datasets such as Question Answering and Code
587 Generation, and we will include results on these
588 datasets in future work. (ii) GRADS cannot be ap-
589 plied to closed-source models, where the gradient
590 flow is not available.

591 Ethics Statement

592 All datasets and models used in this paper are pub-
593 licly available, and our usage follows their licenses
594 and terms. We have employed the LLM tools for
595 coding and polishing.

596 References

597 et al Aaron Grattafiori, Abhimanyu Dubey. 2024. [The](#)
598 [llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

599 et al An Yang, Anfeng Li. 2025. [Qwen3 technical report](#).
600 *Preprint*, arXiv:2505.09388.

601 Eric J Bigelow, Ekdeep Singh Lubana, Robert P. Dick,
602 Hidenori Tanaka, and Tomer Ullman. 2024. [In-](#)
603 [context learning dynamics with random binary se-](#)
604 [quences](#). In *The Twelfth International Conference on*
605 *Learning Representations*.

606 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
607 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
608 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
609 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
610 Gretchen Krueger, Tom Henighan, Rewon Child,
611 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
612 Clemens Winter, and 12 others. 2020. Language
613 models are few-shot learners. In *Proceedings of the*
614 *34th International Conference on Neural Information*
615 *Processing Systems, NIPS '20*, Red Hook, NY, USA.
616 Curran Associates Inc.

617 Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian
618 Pedregosa, Andreas Mueller, Olivier Grisel, Vlad
619 Niculae, Peter Prettenhofer, Alexandre Gramfort,
620 Jaques Grobler, Robert Layton, Jake VanderPlas, Ar-
621 naud Joly, Brian Holt, and Gaël Varoquaux. 2013.
622 API design for machine learning software: experi-
623 ences from the scikit-learn project. In *ECML PKDD*
624 *Workshop: Languages for Data Mining and Machine*
625 *Learning*, pages 108–122.

626 Jiuhai Chen, Lichang Chen, Chen Zhu, and Tianyi Zhou.
627 2023a. [How many demonstrations do you need for](#)
628 [in-context learning?](#) In *Findings of the Association*
629 *for Computational Linguistics: EMNLP 2023*, pages
630 11149–11159, Singapore. Association for Computa-
631 tional Linguistics.

632 Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng,
633 Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang
634 Zhou, Te Gao, and Wanxiang Che. 2025a. [Towards](#)

[reasoning era: A survey of long chain-of-thought](#)
635 [for reasoning large language models](#). *Preprint*,
636 arXiv:2503.09567. 637

Wei-Lin Chen, Cheng-Kuang Wu, Yun-Nung Chen,
638 and Hsin-Hsi Chen. 2023b. [Self-ICL: Zero-shot in-](#)
639 [context learning with self-generated demonstrations](#).
640 In *Proceedings of the 2023 Conference on Empiri-*
641 *cal Methods in Natural Language Processing*, pages
642 15651–15662, Singapore. Association for Computa-
643 tional Linguistics. 644

Xingwu Chen, Lei Zhao, and Difan Zou. 2024. [How](#)
645 [transformers utilize multi-head attention in in-context](#)
646 [learning? a case study on sparse linear regression](#).
647 In *The Thirty-eighth Annual Conference on Neural*
648 *Information Processing Systems*. 649

Zihan Chen, Song Wang, Zhen Tan, Jundong Li, and
650 Cong Shen. 2025b. [MAPLE: Many-shot adaptive](#)
651 [pseudo-labeling for in-context learning](#). In *Forty-*
652 *second International Conference on Machine Learn-*
653 *ing*. 654

Hakaze Cho, Mariko Kato, Yoshihiro Sakai, and Naoya
655 Inoue. 2025. [Revisiting in-context learning inference](#)
656 [circuit in large language models](#). In *The Thirteenth*
657 *International Conference on Learning Representa-*
658 *tions*. 659

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
660 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
661 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
662 Nakano, Christopher Hesse, and John Schulman.
663 2021. [Training verifiers to solve math word prob-](#)
664 [lems](#). *Preprint*, arXiv:2110.14168. 665

Adrian de Wynter. 2025. [Is in-context learning learning?](#)
666 *Preprint*, arXiv:2509.10414. 667

DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing rea-](#)
668 [soning capability in llms via reinforcement learning](#).
669 *Preprint*, arXiv:2501.12948. 670

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan
671 Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu,
672 Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui.
673 2024. [A survey on in-context learning](#). In *Proceed-*
674 *ings of the 2024 Conference on Empirical Methods*
675 *in Natural Language Processing*, pages 1107–1128,
676 Miami, Florida, USA. Association for Computational
677 Linguistics. 678

Deqing Fu, Tian qi Chen, Robin Jia, and Vatsal Sharan.
679 2024. [Transformers learn to achieve second-order](#)
680 [convergence rates for in-context linear regression](#).
681 In *The Thirty-eighth Annual Conference on Neural*
682 *Information Processing Systems*. 683

Chase Goddard, Lindsay M. Smith, Vudtiwat Ngam-
684 pruetikorn, and David J. Schwab. 2025. [When can](#)
685 [in-context learning generalize out of task distribu-](#)
686 [tion?](#) In *Forty-second International Conference on*
687 *Machine Learning*. 688

689	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition . In <i>2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 770–778.	
690		
691		
692		
693	Wei He, Shichun Liu, Jun Zhao, Yiwen Ding, Yi Lu, Zhiheng Xi, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. Self-demos: Eliciting out-of-demonstration generalizability in large language models . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 3829–3845, Mexico City, Mexico. Association for Computational Linguistics.	
694		
695		
696		
697		
698		
699		
700	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset . In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)</i> .	
701		
702		
703		
704		
705		
706	Jianhao Huang, Zixuan Wang, and Jason D. Lee. 2025. Transformers learn to implement multi-step gradient descent with chain of thought . In <i>The Thirteenth International Conference on Learning Representations</i> .	
707		
708		
709		
710	Yu Huang, Yuan Cheng, and Yingbin Liang. 2024. In-context convergence of transformers .	
711		
712	Janak Kapuriya, Manit Kaushik, Debasis Ganguly, and Sumit Bhatia. 2025. Exploring the role of diversity in example selection for in-context learning . In <i>Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '25</i> , page 2962–2966, New York, NY, USA. Association for Computing Machinery.	
713		
714		
715		
716		
717		
718		
719	Wonbin Kweon, SeongKu Kang, Runchu Tian, Pengcheng Jiang, Jiawei Han, and Hwanjo Yu. 2025. Topic coverage-based demonstration retrieval for in-context learning . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 19911–19923, Suzhou, China. Association for Computational Linguistics.	
720		
721		
722		
723		
724		
725		
726	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention . In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	
727		
728		
729		
730		
731		
732		
733	Dongfang Li, zhenyu liu, Xinshuo Hu, Zetian Sun, Baotian Hu, and Min Zhang. 2024a. In-context learning state vector with inner and momentum optimization . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	
734		
735		
736		
737		
738	Hongkang Li, Meng Wang, Songtao Lu, Xiaodong Cui, and Pin-Yu Chen. 2024b. How do nonlinear transformers learn and generalize in in-context learning? In <i>Forty-first International Conference on Machine Learning</i> .	
739		
740		
741		
742		
743	Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng Qiu. 2023. Unified demonstration retriever for in-context learning . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4644–4668, Toronto, Canada. Association for Computational Linguistics.	745
		746
		747
		748
		749
		750
	Zhuowei Li, Zihao Xu, Ligong Han, Yunhe Gao, Song Wen, Di Liu, Hao Wang, and Dimitris N. Metaxas. 2025. Implicit in-context learning . In <i>The Thirteenth International Conference on Learning Representations</i> .	751
		752
		753
		754
		755
	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let’s verify step by step . In <i>The Twelfth International Conference on Learning Representations</i> .	756
		757
		758
		759
		760
	Ziqian Lin and Kangwook Lee. 2024. Dual operating modes of in-context learning . In <i>Forty-first International Conference on Machine Learning</i> .	761
		762
		763
	Hui Liu, Wenya Wang, Hao Sun, Chris Xing Tian, Chenqi Kong, Xin Dong, and Haoliang Li. 2025. Unraveling the mechanics of learning-based demonstration selection for in-context learning . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2623–2641, Vienna, Austria. Association for Computational Linguistics.	764
		765
		766
		767
		768
		769
		770
		771
	Yiting Liu and Zhi-Hong Deng. 2025. Iterative vectors: In-context gradient steering without backpropagation . In <i>Forty-second International Conference on Machine Learning</i> .	772
		773
		774
		775
	Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. 2024a. On LLMs-driven synthetic data generation, curation, and evaluation: A survey . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 11065–11082, Bangkok, Thailand. Association for Computational Linguistics.	776
		777
		778
		779
		780
		781
		782
	Quanyu Long, Yin Wu, Wenya Wang, and Sinno Jialin Pan. 2024b. Does in-context learning really learn? rethinking how large language models respond and solve tasks via in-context learning . In <i>First Conference on Language Modeling</i> .	783
		784
		785
		786
		787
	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.	788
		789
		790
		791
		792
		793
		794
		795
	Yue Lu, Mary Letey, Jacob A Zavatone-Veth, Anindita Maiti, and Cengiz Pehlevan. 2024. In-context learning by linear attention: Exact asymptotics and experiments . In <i>NeurIPS 2024 Workshop on Mathematics of Modern Machine Learning</i> .	796
		797
		798
		799
		800

801	Arvind V. Mahankali, Tatsunori Hashimoto, and Tengyu Ma. 2024. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. In <i>The Twelfth International Conference on Learning Representations</i> .	856
802		857
803		858
804		859
805		860
806	Alex Nguyen and Gautam Reddy. 2025. Differential learning kinetics govern the transition from memorization to generalization during in-context learning. In <i>The Thirteenth International Conference on Learning Representations</i> .	861
807		862
808		863
809		864
810		865
811	Tai Nguyen and Eric Wong. 2023. In-context example selection with influences. <i>Preprint</i> , arXiv:2302.11042.	866
812		867
813		868
814	Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 188–197, Hong Kong, China. Association for Computational Linguistics.	869
815		870
816		871
817		872
818		873
819		874
820		875
821		876
822		877
823	Kazusato Oko, Yujin Song, Taiji Suzuki, and Denny Wu. 2024. Pretrained transformer efficiently learns low-dimensional target functions in-context. In <i>NeurIPS</i> .	878
824		880
825		881
826	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, and 7 others. 2022. In-context learning and induction heads. <i>Preprint</i> , arXiv:2209.11895.	882
827		883
828		884
829		885
830		886
831		887
832		888
833		889
834	Core Francisco Park, Ekdeep Singh Lubana, and Hide-nori Tanaka. 2025. Competition dynamics shape algorithmic phases of in-context learning. In <i>The Thirteenth International Conference on Learning Representations</i> .	890
835		891
836		892
837		893
838		894
839	Keqin Peng, Liang Ding, Yancheng Yuan, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2024. Revisiting demonstration selection strategies in in-context learning. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9090–9101, Bangkok, Thailand. Association for Computational Linguistics.	895
840		896
841		897
842		898
843		899
844		900
845		901
846		902
847	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.	903
848		904
849		905
850		906
851		907
852		908
853		909
854		910
855		911
		912
		913
	Kha Pham, Hung Le, Man Ngo, and Truyen Tran. 2025. Rapid selection and ordering of in-context demonstrations via prompt embedding clustering. In <i>The Thirteenth International Conference on Learning Representations</i> .	914
		915
	Allan Raventos, Mansheej Paul, Feng Chen, and Surya Ganguli. 2023. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	916
		917
	Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5418–5426, Online. Association for Computational Linguistics.	918
		919
	Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2655–2671, Seattle, United States. Association for Computational Linguistics.	920
		921
	Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang. 2024. Why larger language models do in-context learning differently? In <i>Forty-first International Conference on Machine Learning</i> .	922
		923
	Suzanna Sia, David Mueller, and Kevin Duh. 2024. Where does in-context learning \(\backslash\) happen in large language models? In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	924
		925
	Aaditya K Singh, Ted Moskovitz, Sara Dragutinović, Felix Hill, Stephanie C.Y. Chan, and Andrew M Saxe. 2025. Strategy coopeitition explains the emergence and transience of in-context learning. In <i>Forty-second International Conference on Machine Learning</i> .	926
		927
	Matthew Smart, Alberto Bietti, and Anirvan M. Sengupta. 2025. In-context denoising with one-layer transformers: Connections between attention and associative memory retrieval. In <i>Forty-second International Conference on Machine Learning</i> .	928
		929
	Yi Su, Yunpeng Tai, Yixin Ji, Juntao Li, Yan Bowen, and Min Zhang. 2024. Demonstration augmentation for zero-shot in-context learning. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 14232–14244, Bangkok, Thailand. Association for Computational Linguistics.	930
		931
		932
		933
	Wenbiao Tao, Hanlun Zhu, Keren Tan, Jiani Wang, Yuanyuan Liang, Huihui Jiang, Pengcheng Yuan, and Yunshi Lan. 2024. Finqa: A training-free dynamic knowledge graph question answering system in finance with llm-based revision. In <i>Machine Learning and Knowledge Discovery in Databases. Research Track and Demo Track: European Conference, ECML PKDD 2024, Vilnius, Lithuania, September 9–13, 2024, Proceedings, Part VIII</i> , page 418–423, Berlin, Heidelberg. Springer-Verlag.	934
		935

914	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert	970
915	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,	971
916	Kaiser, and Illia Polosukhin. 2017. Attention is all	Da Huang, Denny Zhou, and Tengyu Ma. 2024.	972
917	you need . In <i>Advances in Neural Information Pro-</i>	Larger language models do in-context learning dif-	973
918	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	ferently .	974
919	Max Vladymyrov, Johannes von Oswald, Mark Sandler,	Kevin Christian Wibisono and Yixin Wang. 2024. From	975
920	and Rong Ge. 2024. Linear transformers are versatile	unstructured data to in-context learning: Exploring	976
921	in-context learners . In <i>ICML 2024 Workshop on In-</i>	what tasks can be learned and when . In <i>The Thirty-</i>	977
922	<i>Context Learning</i> .	eighth Annual Conference on Neural Information	978
923	Dingzirui Wang, Xuanliang Zhang, Qiguang Chen,	<i>Processing Systems</i> .	979
924	Longxu Dou, Xiao Xu, Rongyu Cao, YINGWEI MA,	Noam Wies, Yoav Levine, and Amnon Shashua. 2023.	980
925	Qingfu Zhu, Wanxiang Che, Binhua Li, Fei Huang,	The learnability of in-context learning . In <i>Thirty-</i>	981
926	and Yongbin Li. 2025a. In-context transfer learn-	seventh Conference on Neural Information Process-	982
927	ing: Demonstration synthesis by transferring similar	ing Systems .	983
928	tasks .	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	984
929	Dingzirui Wang, Xuanliang Zhang, Keyan Xu, Qingfu	Chaumond, Clement Delangue, Anthony Moi, Pier-	985
930	Zhu, Wanxiang Che, and Yang Deng. 2025b. V-	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,	986
931	synthesis: Task-agnostic synthesis of consistent and	Joe Davison, Sam Shleifer, Patrick von Platen, Clara	987
932	diverse in-context demonstrations from scratch via	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven	988
933	v-entropy . <i>Preprint</i> , arXiv:2506.23149.	Le Scao, Sylvain Gugger, and 3 others. 2020. Trans-	989
934	Dingzirui Wang, Xuanliang Zhang, Keyan Xu, Qingfu	formers: State-of-the-art natural language processing .	990
935	Zhu, Wanxiang Che, and Yang Deng. 2025c.	In <i>Proceedings of the 2020 Conference on Empirical</i>	991
936	Learning-to-context slope: Evaluating in-context	<i>Methods in Natural Language Processing: System</i>	992
937	learning effectiveness beyond performance illusions .	<i>Demonstrations</i> , pages 38–45, Online. Association	993
938	<i>Preprint</i> , arXiv:2506.23146.	for Computational Linguistics.	994
939	Futing Wang, Jianhao Yan, Yue Zhang, and Tao Lin.	Yuan Wu, Diana Inkpen, and Ahmed El-Roby. 2022.	995
940	2025d. ELICIT: LLM augmentation via external in-	Maximum batch frobenius norm for multi-domain	996
941	context capability . In <i>The Thirteenth International</i>	text classification . In <i>IEEE International Conference</i>	997
942	<i>Conference on Learning Representations</i> .	<i>on Acoustics, Speech and Signal Processing, ICASSP</i>	998
943	Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou,	2022, <i>Virtual and Singapore, 23-27 May 2022</i> , pages	999
944	Fandong Meng, Jie Zhou, and Xu Sun. 2023. Label	3763–3767. IEEE.	1000
945	words are anchors: An information flow perspective	Vikas Yadav, Steven Bethard, and Mihai Surdeanu.	1001
946	for understanding in-context learning . In <i>Proceed-</i>	2019. Quick and (not so) dirty: Unsupervised se-	1002
947	<i>ings of the 2023 Conference on Empirical Methods</i>	lection of justification sentences for multi-hop ques-	1003
948	<i>in Natural Language Processing</i> , pages 9840–9855,	tion answering . In <i>Proceedings of the 2019 Confer-</i>	1004
949	Singapore. Association for Computational Linguis-	<i>ence on Empirical Methods in Natural Language Pro-</i>	1005
950	tics.	<i>cessing and the 9th International Joint Conference</i>	1006
951	Song Wang, Zihan Chen, Chengshuai Shi, Cong Shen,	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	1007
952	and Jundong Li. 2024a. Mixture of demonstrations	pages 2578–2589, Hong Kong, China. Association	1008
953	for in-context learning . In <i>The Thirty-eighth Annual</i>	for Computational Linguistics.	1009
954	<i>Conference on Neural Information Processing Sys-</i>	Tong Yang, Yu Huang, Yingbin Liang, and Yuejie Chi.	1010
955	<i>tems</i> .	2024. In-context learning with representations: Con-	1011
956	Xubin Wang, Jianfei Wu, Yuan Yichen, Deyu Cai,	textual generalization of trained transformers . In	1012
957	Mingzhe Li, and Weijia Jia. 2025e. Demonstration	<i>ICML 2024 Workshop on Theoretical Foundations of</i>	1013
958	selection for in-context learning via reinforcement	<i>Foundation Models</i> .	1014
959	learning . In <i>Forty-second International Conference</i>	Zhao Yang, Yuanzhe Zhang, Dianbo Sui, Cao Liu, Jun	1015
960	<i>on Machine Learning</i> .	Zhao, and Kang Liu. 2023. Representative demon-	1016
961	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni,	stration selection for in-context learning with two-	1017
962	Abhranil Chandra, Shiguang Guo, Weiming Ren,	stage determinantal point process . In <i>Proceedings of</i>	1018
963	Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max	<i>the 2023 Conference on Empirical Methods in Natu-</i>	1019
964	Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang	<i>ral Language Processing</i> , pages 5443–5456, Singa-	1020
965	Yue, and Wenhui Chen. 2024b. MMLU-pro: A more	pore. Association for Computational Linguistics.	1021
966	robust and challenging multi-task language under-	Xi Ye, Srinivasan Iyer, Asli Celikyilmaz, Veselin Stoy-	1022
967	standing benchmark . In <i>The Thirty-eight Conference</i>	anov, Greg Durrett, and Ramakanth Pasunuru. 2023.	1023
968	<i>on Neural Information Processing Systems Datasets</i>	Complementary explanations for effective in-context	1024
969	<i>and Benchmarks Track</i> .	learning . In <i>Findings of the Association for Compu-</i>	1025
		<i>tational Linguistics: ACL 2023</i> , pages 4469–4484,	1026

- 1027 Toronto, Canada. Association for Computational Lin-
1028 guistics.
- 1029 Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. 2024.
1030 [Trained transformers learn linear models in-context.](#)
1031 *Journal of Machine Learning Research*, 25(49):1–55.
- 1032 Xingxuan Zhang, Haoran Wang, Jiansheng Li, Yuan
1033 Xue, Shikai Guan, Renzhe Xu, Hao Zou, Han Yu, and
1034 Peng Cui. 2025a. [Understanding the generalization
1035 of in-context learning in transformers: An empirical
1036 study.](#) In *The Thirteenth International Conference
1037 on Learning Representations*.
- 1038 Yedi Zhang, Aaditya K Singh, Peter E. Latham, and
1039 Andrew M Saxe. 2025b. [Training dynamics of in-
1040 context learning in linear attention.](#) In *Forty-second
1041 International Conference on Machine Learning*.
- 1042 Zheng Zhang, Shaocheng Lan, Lei Song, Jiang Bian,
1043 Yexin Li, and Kan Ren. 2025c. [Learning to select
1044 in-context demonstration preferred by large language
1045 model.](#) In *Findings of the Association for Computa-
1046 tional Linguistics: ACL 2025*, pages 11345–11360,
1047 Vienna, Austria. Association for Computational Lin-
1048 guistics.
- 1049 Ziniu Zhang, Zhenshuo Zhang, Dongyue Li, Lu Wang,
1050 Jennifer Dy, and Hongyang R. Zhang. 2025d. [Linear-
1051 time demonstration selection for in-context learn-
1052 ing via gradient estimation.](#) In *Proceedings of the
1053 2025 Conference on Empirical Methods in Natural
1054 Language Processing*, pages 16470–16488, Suzhou,
1055 China. Association for Computational Linguistics.
- 1056 Siyan Zhao, Tung Nguyen, and Aditya Grover. 2024.
1057 [Probing the decision boundaries of in-context learn-
1058 ing in large language models.](#) In *The Thirty-eighth
1059 Annual Conference on Neural Information Process-
1060 ing Systems*.
- 1061 Yuxiang Zhou, Jiazheng Li, Yanzheng Xiang, Hanqi
1062 Yan, Lin Gui, and Yulan He. 2024. [The mystery
1063 of in-context learning: A comprehensive survey on
1064 interpretation and analysis.](#) In *Proceedings of the
1065 2024 Conference on Empirical Methods in Natural
1066 Language Processing*, pages 14365–14378, Miami,
1067 Florida, USA. Association for Computational Lin-
1068 guistics.

A Proof

A.1 Proof of Equation 2

Proof. Based on Equation 1 and Zhang et al. (2024), the predicted answer of LSA is:

$$\hat{y}_{query} = ((w_{21}^{PV})^\top \quad w_{22}^{PV}) \cdot \left(\frac{EE^\top}{N} \right) \cdot \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} q_x, \quad (7)$$

where N is the number of demonstrations. For $N = 1$, the matrix product EE^\top can be expressed compactly as $EE^\top = dd^\top + qq^\top$. Substituting this into the equation gives:

$$\hat{y}_{query} = ((w_{21}^{PV})^\top \quad w_{22}^{PV})(dd^\top + qq^\top) \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} q_x \quad (8)$$

We can expand this expression and separate the terms that depend on d :

$$\hat{y}_{query} = \underbrace{((w_{21}^{PV})^\top \quad w_{22}^{PV})(dd^\top) \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} q_x}_{\text{Term depending on } d} + \underbrace{((w_{21}^{PV})^\top \quad w_{22}^{PV})(qq^\top) \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} q_x}_{\text{Constant w.r.t. } d}$$

The term depending on d can be rewritten using the property of scalar products:

$$\begin{aligned} & ((w_{21}^{PV})^\top \quad w_{22}^{PV}) dd^\top \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} q_x \\ &= \left(((w_{21}^{PV})^\top \quad w_{22}^{PV}) d \right) \left(d^\top \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} q_x \right) \end{aligned}$$

This is a quadratic form of the vector d . To compute its gradient, we use the vector calculus identity $\nabla_z((a^\top z)(b^\top z)) = a(b^\top z) + b(a^\top z)$. We set:

$$\begin{aligned} a &= \begin{pmatrix} w_{21}^{PV} \\ w_{22}^{PV} \end{pmatrix} \\ b &= \begin{pmatrix} W_{11}^{KQ} q_x \\ (w_{21}^{KQ})^\top q_x \end{pmatrix} \end{aligned}$$

The second term in the expression for \hat{y}_{query} is constant with respect to d , so its gradient is zero.

Applying the identity to the first term yields the gradient of \hat{y}_{query} with respect to d :

$$\begin{aligned} \nabla_d \hat{y}_{query} &= \begin{pmatrix} w_{21}^{PV} \\ w_{22}^{PV} \end{pmatrix} \left(\begin{pmatrix} d_x \\ d_y \end{pmatrix}^\top \begin{pmatrix} W_{11}^{KQ} q_x \\ (w_{21}^{KQ})^\top q_x \end{pmatrix} \right) \\ &+ \begin{pmatrix} W_{11}^{KQ} q_x \\ (w_{21}^{KQ})^\top q_x \end{pmatrix} \left(\begin{pmatrix} w_{21}^{PV} \\ w_{22}^{PV} \end{pmatrix}^\top \begin{pmatrix} d_x \\ d_y \end{pmatrix} \right) \\ &= \begin{pmatrix} w_{21}^{PV} \\ w_{22}^{PV} \end{pmatrix} \left(d_x^\top W_{11}^{KQ} q_x + d_y (w_{21}^{KQ})^\top q_x \right) \\ &+ \begin{pmatrix} W_{11}^{KQ} q_x \\ (w_{21}^{KQ})^\top q_x \end{pmatrix} \left((w_{21}^{PV})^\top d_x + d_y w_{22}^{PV} \right) \end{aligned}$$

This completes the derivation. \square

A.2 Proof of Lemma 1

Proof. We prove by induction on the layer index l .

Base case ($l = 0$). The two inequalities in the statement are satisfied by assumption.

Induction step. Assume the inequalities hold for some layer $l - 1$ ($1 \leq l \leq L$), i.e.

$$\begin{aligned} \|W^{PV,(l-1)} d_1^{(l-1)}\| &\geq \|W^{PV,(l-1)} d_2^{(l-1)}\| \\ \|(d_1^{(l-1)})^\top W^{KQ,(l-1)} q^{(l-1)}\| &\geq \\ \|(d_2^{(l-1)})^\top W^{KQ,(l-1)} q^{(l-1)}\|. \end{aligned}$$

Apply (5) to both demonstrations and use the strict monotonicity of g_l :

$$\begin{aligned} & \|W^{PV,(l)} d_1^{(l)}\| \\ &= g_l(\|W^{PV,(l-1)} d_1^{(l-1)}\|) \\ &\geq g_l(\|W^{PV,(l-1)} d_2^{(l-1)}\|) \\ &= \|W^{PV,(l)} d_2^{(l)}\|. \end{aligned}$$

Similarly, apply (6) and the strict monotonicity of h_l :

$$\begin{aligned} & \|(d_1^{(l)})^\top W^{KQ,(l)} q^{(l)}\| \\ &= h_l(\|(d_1^{(l-1)})^\top W^{KQ,(l-1)} q^{(l-1)}\|) \\ &\geq h_l(\|(d_2^{(l-1)})^\top W^{KQ,(l-1)} q^{(l-1)}\|) \\ &= \|(d_2^{(l)})^\top W^{KQ,(l)} q^{(l)}\|. \end{aligned}$$

Thus the claim holds for layer l . By induction, it holds for all layers $l = 0, \dots, L$. \square

Prompt of Inference
{task}
Below are some examples
—
{demo}
—
Based on the above instruction and examples, solve the following problem.
{question}

Table 4: The prompt of the inference.

A.3 Proof of Theorem 1

Proof. Based on Equation 4, we can derive:

$$\begin{aligned}
& \frac{\partial \hat{q}_y^{(l_1)}}{\partial d^{(0)}} / \frac{\partial \hat{q}_y^{(l_2)}}{\partial d^{(0)}} \\
&= \left(\frac{\partial \hat{q}_y^{(l_1)}}{\partial E^{(l_1-1)}} \times \frac{\partial E^{(1)}}{\partial d^{(0)}} \times \prod_{i=2}^{l_1} \frac{\partial E^{(i)}}{\partial E^{(i-1)}} \right) / \\
& \left(\frac{\partial \hat{q}_y^{(l_2)}}{\partial E^{(l_2-1)}} \times \frac{\partial E^{(1)}}{\partial d^{(0)}} \times \prod_{i=2}^{l_2} \frac{\partial E^{(i)}}{\partial E^{(i-1)}} \right) \\
&= \left(\frac{\partial \hat{q}_y^{(l_1)}}{\partial E^{(l_1-1)}} / \frac{\partial \hat{q}_y^{(l_2)}}{\partial E^{(l_2-1)}} \right) \times \prod_{i=l_2+1}^{l_1} \frac{\partial E^{(i)}}{\partial E^{(i-1)}}
\end{aligned}$$

Consider the total derivative of f_{LSA} :

$$\begin{aligned}
df_{\text{LSA}} &= dE + W^{PV} dE E^\top W^{KQ} E + \\
& W^{PV} E (E^\top W^{KQ} dE + dE^\top W^{KQ} E)
\end{aligned}$$

It can be observed that the coefficients of the differential terms are consistent with those in Definition 1. Therefore, from $d_1^{(l)} \succ_{q^{(l)}; \theta^{(l)}} d_2^{(l)}$, we know that:

$$\frac{\partial E^{(l)}}{\partial E^{(l-1)}}(E_1; \theta) \geq \frac{\partial E^{(l)}}{\partial E^{(l-1)}}(E_2; \theta),$$

where $l \in \{l_2 + 1, \dots, l_1\}$. Therefore, we can conclude that:

$$\frac{\|\partial_{d^{(0)}} \hat{q}_y^{(l_1)}(E_1; \theta)\|}{\|\partial_{d^{(0)}} \hat{q}_y^{(l_1)}(E_2; \theta)\|} \geq \frac{\|\partial_{d^{(0)}} \hat{q}_y^{(l_2)}(E_1; \theta)\|}{\|\partial_{d^{(0)}} \hat{q}_y^{(l_2)}(E_2; \theta)\|}$$

□

B Additional Information

B.1 Prompt

In this section, we present the inference prompt of our main experiment, as shown in Table 4. The task definition we used is the same as the previous works (DeepSeek-AI, 2025; Aaron Grattafiori, 2024).

Dataset	Test Set	Demonstration
GSM8K	1319	7473
MATH	500	7496
ARC-Challenge	1172	1119
MMLU-Pro	1000	70
Amazon Review	200	1800
FinQA	6251	1147

Table 5: The scales of test set and demonstrations of each dataset.

B.2 Dataset

In this section, we detail the datasets used in our study. Table 5 summarizes the scale of the test set and demonstrations for each.

GSM8K GSM8K (Cobbe et al., 2021) is a high-quality collection of elementary school-level math problems. We utilize its training set directly as the demonstration pool.

MATH The MATH dataset (Hendrycks et al., 2021) consists of challenging high school competition-level math problems in fields like algebra, probability, and geometry. Following the approach of Lightman et al. (2024), we evaluate GRADS on a random sample of 500 problems. The demonstrations are drawn from the official training set.

ARC-Challenge The ARC-Challenge (Yadav et al., 2019) is a question-answering dataset with difficult, science-focused questions. For this dataset, the training set is used as our demonstration pool.

MMLU-Pro MMLU-Pro (Wang et al., 2024b) serves as a multi-task benchmark for the comprehensive evaluation of LLMs on professional domain knowledge and complex reasoning. As the dataset is only divided into validation and test sets, we use the validation set as our demonstration pool and conduct evaluations on the test set.

Amazon Review The Amazon Review dataset (Ni et al., 2019), containing a vast amount of user ratings and reviews, is widely used for research in sentiment analysis and recommender systems. Due to the immense size, we select the *Health and Personal Care* category for testing, while using the *All Beauty*, *Digital Music*, and *Software* categories to form the demonstrations.

FinQA FinQA (Tao et al., 2024) is a large-scale benchmark for numerical reasoning in fi-

1179	nance. It contains 8,281 expert-written question-answer pairs built from S&P 500 companies' earnings/annual reports, where evidence can come from both narrative text and structured tables. For each question, FinQA provides supporting facts and an annotated reasoning program that can be executed to obtain the answer, enabling fully explainable evaluation.	1229
1180		1230
1181		
1182		
1183		
1184		
1185		
1186		
1187	B.3 Baseline	
1188	BM25 BM25 is a classic sparse retrieval method based on the probabilistic relevance framework, serving as an extension of TF-IDF. In the context of ICL, it treats the test query as a search query and the candidate demonstrations as documents. It ranks demonstrations by scoring them based on the frequency and distribution of query terms, without considering word order or deep semantics. The top-K scored demonstrations are then selected as the in-context demonstrations. This "bag-of-words" approach is known for its computational efficiency and serves as a strong baseline.	
1189		
1190		
1191		
1192		
1193		
1194		
1195		
1196		
1197		
1198		
1199		
1200	MMR Maximal Marginal Relevance (MMR) is a selection strategy designed to balance the relevance of demonstrations to the query with the diversity within the selected set. The rationale is that selecting only the nearest neighbors (most relevant examples) can result in a set of overly similar demonstrations, which may limit the variety of reasoning processes shown to the model. MMR addresses this by iteratively selecting demonstrations that maximize a combined score of relevance to the query and dissimilarity from the examples already chosen. This approach aims to create a set of demonstrations that are not only relevant but also complementary, using diversity as a proxy for complementarity to improve ICL performance.	
1201		
1202		
1203		
1204		
1205		
1206		
1207		
1208		
1209		
1210		
1211		
1212		
1213		
1214		
1215	MoD Mixture of Demonstrations (MoD) is a framework designed to overcome the challenges of a large search space and suboptimal retriever optimization in ICL demonstration selection. The core idea is to partition the entire demonstration pool into distinct groups, typically using K-means clustering on sentence embeddings. Each group is governed by a dedicated "expert", a unique retriever model trained specifically for that partition. During inference, these experts collaboratively retrieve demonstrations for a given query, with the final set being an aggregation of examples selected by the most relevant experts. This "mixture of experts" approach reduces the search complexity	
1216		
1217		
1218		
1219		
1220		
1221		
1222		
1223		
1224		
1225		
1226		
1227		
1228		
	while ensuring the selected demonstrations are diverse and effective.	1229
		1230
	Influence Influence (Nguyen and Wong, 2023) treats demonstration quality as the causal contribution an example makes to in context learning performance. It repeatedly samples many random subsets of size k from the training set, places each subset in the prompt, and evaluates the model on a validation set to obtain a score. For a candidate example i , it estimates in context influence as the average score of subsets that include i minus the average score of subsets that exclude i . It then ranks examples by this influence and selects the top k positive ones, while also revealing harmful examples and order effects such as recency bias.	1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
	DICL DICL (Kapuriya et al., 2025) aims to balance relevance and diversity when choosing demonstrations. It first retrieves a candidate pool larger than k using sparse or dense similarity, such as TF-IDF or sentence embeddings, so there is room to diversify within a locally relevant neighborhood. It then applies a greedy Maximum Marginal Relevance style reranking that iteratively picks examples that are highly similar to the input while remaining dissimilar to already selected demonstrations. A weight parameter λ controls the trade off between relevance and diversity. This reduces topical redundancy from pure similarity retrieval and tends to yield more stable gains across context sizes and similarity functions.	1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
	GenICL GenICL (Zhang et al., 2025c) frames demonstration selection as learning which examples an LLM <i>prefers</i> to see in the prompt, rather than optimizing a surrogate similarity objective. For each training query, it uses LLM feedback to assign a preference score to candidate demonstrations (e.g., how well the LLM predicts the gold output when conditioned on that single demonstration), and forms preferred vs. non-preferred demonstration pairs. It then trains a generative preference model with a latent variable that bridges selection and inference, optimizing a preference-learning objective (with regularization to avoid drifting too far from a reference behavior). At inference time, GenICL first narrows the search space with an off-the-shelf retriever, and then ranks candidates by the learned probability of generating the latent "preferred demonstration" variable for the test input. The top- K candidates under this learned preference score are used as the in-context demonstra-	1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278

1279 tions, aiming to select examples that are directly
1280 aligned with downstream ICL utility rather than
1281 surface similarity.

1282 **TopicK** TopicK (Kweon et al., 2025) selects
1283 demonstrations by explicitly maximizing *topic cov-*
1284 *erage* of the fine-grained knowledge required by
1285 the test input, instead of relying only on embed-
1286 ding similarity. It builds a global topic set from
1287 the demonstration pool and trains a lightweight
1288 topic predictor to estimate (i) the required topic
1289 distribution of a test input and (ii) the covered
1290 topic distribution of each candidate demonstration.
1291 In addition, TopicK estimates the model’s topical
1292 knowledge (a prior over topics the LLM already
1293 “knows well”) using aggregated zero-shot perfor-
1294 mance, and down-weights such topics during re-
1295 trieval. Demonstrations are then chosen iteratively
1296 with a greedy coverage strategy: each new exam-
1297 ple is selected to introduce previously uncovered
1298 required topics, especially those where the model
1299 exhibits low topical knowledge. This yields a set
1300 that is both relevant and non-redundant at the topic
1301 level, improving complementarity beyond nearest-
1302 neighbor retrieval.

1303 **ICL-Grad** ICL-Grad (Zhang et al., 2025d) accel-
1304 erates demonstration selection by estimating how
1305 candidate examples affect ICL performance us-
1306 ing gradients in the *input embedding space*. The
1307 method pre-computes model outputs and gradients
1308 for an anchor prompt, and then applies a first-
1309 order (Taylor) approximation to predict the loss
1310 of many different demonstration subsets *without*
1311 running full LLM inference for each subset. By
1312 sampling many random candidate subsets and es-
1313 timating their losses cheaply, it aggregates these
1314 subset outcomes into an influence-style score for
1315 each demonstration (e.g., the average estimated
1316 loss over subsets that include that example). The fi-
1317 nal k demonstrations are selected as those with
1318 the most beneficial scores (typically the lowest
1319 estimated losses), enabling scalable subset selec-
1320 tion with near-linear cost after one gradient pass.
1321 Compared with similarity-only retrieval, this ap-
1322 proach better reflects the model’s behavior under
1323 ICL while remaining computationally efficient for
1324 large candidate pools.

1325 B.4 Scale of Effective and Ineffective Data

1326 The number of effective and ineffective demonstra-
1327 tions we sampled is shown in Table 6.

1328 C Additional Discussion

1329 C.1 Efficiency of GRADS

1330 The computational cost of GRADS is mainly di-
1331 vided into two parts: the offline computation of the
1332 encoding result for each demonstration and the on-
1333 line retrieval of the demonstration and subsequent
1334 inference for a query. Let $D = \{d_1, \dots, d_n\}$ repre-
1335 sent the entire demonstration pool, and let $\mathcal{M}_\theta(x)$
1336 denote the computational cost of the model \mathcal{M}_θ
1337 for a given input x . For a query q , let the retrieved
1338 demonstration be d_q .

1339 The time complexity of the offline processing
1340 for GRADS is: $O(\sum_{i=1}^n \mathcal{M}_\theta(d_i))$. Although the
1341 offline cost is relatively high, the pre-computation
1342 of demonstration encodings is done offline and thus
1343 does not affect the online user query process.

1344 The time complexity for online processing is:

$$1345 O(\mathcal{M}_\theta(q) + 4e^2 + \mathcal{M}_\theta(q + d_q)) \quad (9)$$

1346 In Equation 9, the first term represents encoding the
1347 user query, the second term corresponds to calculat-
1348 ing Equation 2, and the third term is for generating
1349 the answer to the user query based on the retrieved
1350 demonstration. Considering that \mathcal{M}_θ is positively
1351 correlated with the input length and the complex-
1352 ity of calculating Equation 2 is significantly lower
1353 than that of a full model inference \mathcal{M}_θ , the online
1354 processing time complexity simplifies to:

$$1355 O(\mathcal{M}_\theta(q + d_q)) \quad (10)$$

1356 This is equivalent to the time complexity of direct
1357 1-shot inference, which demonstrates the high effi-
1358 ciency of GRADS. We also compare the run time
1359 with the existing demonstration selection methods
1360 in Appendix D.3.

1361 C.2 Effect of Non-Linear Factor of 1362 Transformer

1363 In this part, we discuss how different non-linear
1364 factors of the Transformer affect the conclusions
1365 of this paper.

1366 **Softmax attention.** While monotone in the align-
1367 ment score $d^\top W^K Q$, softmax introduces compet-
1368 itive normalization that sharpens high-logit tokens
1369 and compresses marginal gradients for medium/low
1370 scores. Consequently, demonstrations with high
1371 alignment but little novel value information are
1372 further down-weighted; with depth, saturation can
1373 dampen inter-layer amplification unless the distri-
1374 bution is sufficiently diffuse (e.g., higher tempera-
1375 ture).

Model	Type	GSM8K	MATH	ARC-Challenge	MMLU-Pro	Amazon	FinQA
Llama-3.1-8B	Effective	84	52	86	174	26	92
	Ineffective	210	280	212	466	63	267
Llama-R1-8B	Effective	41	27	77	203	20	83
	Ineffective	519	146	198	550	102	202
Qwen3-8B	Effective	33	31	28	251	23	37
	Ineffective	99	111	110	368	77	107

Table 6: The number of effective and ineffective demonstrations under each setting.

FFN activations and LayerNorm. Nonlinear gates (e.g., GELU/SwiGLU) truncate or saturate the propagated value magnitude $\|W^{PV}d\|$, while LayerNorm projects representations onto a sphere, partially washing out norm-based cues. Hence, advantages due to “unassimilated information” are reduced unless coupled with strong alignment; amplification persists when layers operate in near-linear regimes but degrades under saturation.

Residual connections. Residual pathways bypass nonlinear transforms and carry forward prior states, diluting the relative gain of newly injected information along the main branch. This weakens strictly monotonic amplification with depth and explains the non-monotonic, layer-dependent optima observed in practice.

C.3 Why Gradient Flow Can Reflect the Effectiveness

As discussed in Wang et al. (2023), for irrelevant demonstrations, LLMs actively ignore their information, which leads to smaller gradient flow. In other words, the model only produces large gradient flow for effective information, which is consistent with the way attention weights are computed. The experimental results in Table 2 are also consistent with this observation, where the effectiveness of one demonstration is positively correlated with the magnitude of the gradient flow.

D Additional Experiment

D.1 Performance under Multiple Runs

To verify that GRADS consistently enhances ICL performance, we conduct experiments with multiple runs. We use Llama-3.1-8B as the experimental model and perform 5 trials for each setting. As shown in Table 7, GRADS outperforms the baselines across all settings, demonstrating stable improvements in ICL performance.

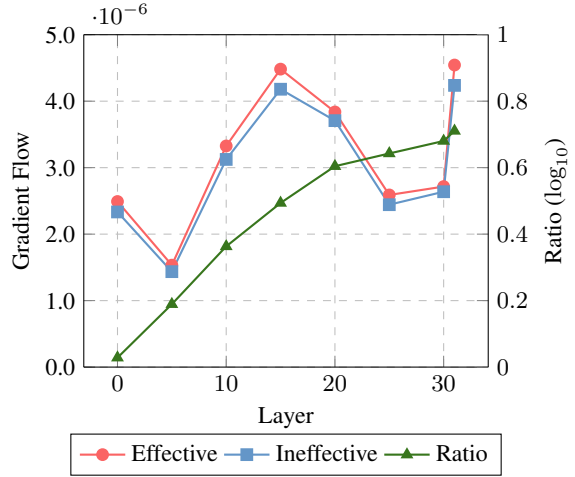


Figure 4: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-3.1-8B on GSM8K.

D.2 Detailed Gradient Flow under Each Setting

In this part, we present the gradient flow under each setting from Figure 4 to Figure 14.

D.3 Run Time Comparison Cross Different Methods

In this section, we compare the run-time efficiency of different demonstration selection methods. Using the same device, we calculate the average time required for demonstration selection per question on the MATH dataset using Llama 3.1-8B. The results are shown in Table 8. We observe that GRADS achieves superior performance with comparable computational efficiency, demonstrating the high efficiency of our method.

D.4 Ablation Study

In this section, we present ablation studies of GRADS to validate the effectiveness of each component in our approach. The results are shown in Table 9, from which we observe that: (i) Removing Similarity causes a noticeable drop on MMLU-Pro (-4.0) and ARC-C/FinQA (-2.2 each),

GSM8K	MATH	ARC-C	MMLU-Pro	Amazon	FinQA
85.6 ± 1.4	48.2 ± 2.2	86.3 ± 1.3	56.0 ± 2.0	70.0 ± 1.8	51.1 ± 1.3

Table 7: The performance of GRADS using Llama-3.1-8B on different datasets with five runs.

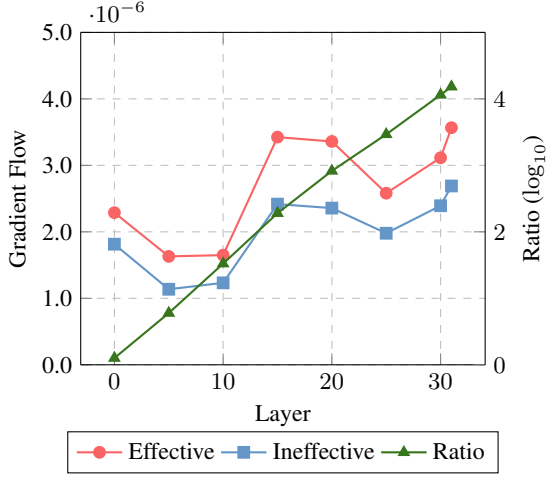


Figure 5: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-3.1-8B on MATH.

Method	EM	T (s)
Zero	47.0	—
BM25	44.6	0.02
MMR	45.2	0.13
MoD	47.0	0.21
Influence	47.4	0.27
DICL	47.4	0.43
GenICL	47.6	0.96
TopicK	47.0	0.83
ICL-Grad	47.0	0.59
GRADS	48.2	0.32

Table 8: The performance and run time of each method on MATH using Llama-3.1-8B. T denotes the average run time of each question.

1435 suggesting that similarity-based signals are par-
1436 ticularly important for knowledge-intensive rea-
1437 soning and multi-step decision tasks, while hav-
1438 ing only marginal impact on GSM8K (-0.6) and
1439 MATH (-0.2). (ii) Removing Information leads to
1440 the largest degradation on Amazon (-3.5) and ARC-
1441 C (-3.4), indicating that the information component
1442 contributes more to domain- and context-sensitive
1443 understanding, whereas its effect is relatively small
1444 on MATH/FinQA (-0.6 each) compared with the
1445 similarity component.

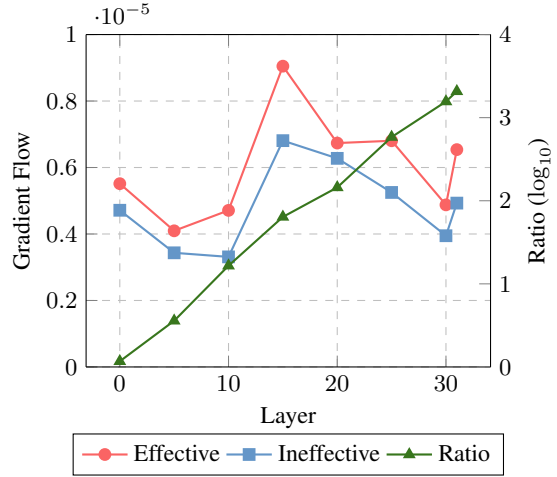


Figure 6: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-3.1-8B on ARC-Challenge.

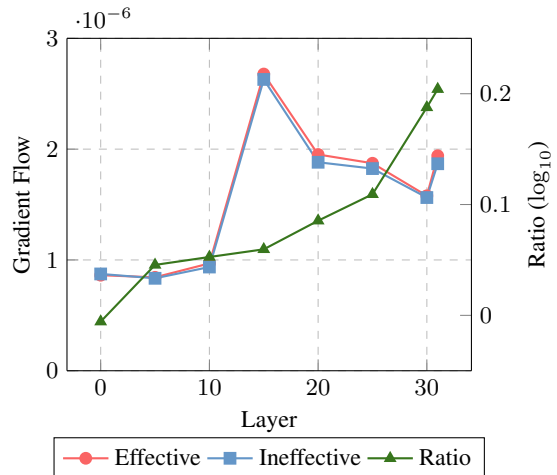


Figure 7: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-3.1-8B on MMLU-Pro.

Method	GSM8K	MATH	ARC-C	MMLU-Pro	Amazon	FinQA
GRADS	85.6	48.2	86.3	56.0	70.0	52.3
- Similarity	85.0 _(-0.6)	48.0 _(-0.2)	84.1 _(-2.2)	52.0 _(-4.0)	68.0 _(-2.0)	50.1 _(-2.2)
- Information	84.2 _(-1.4)	47.6 _(-0.6)	82.9 _(-3.4)	53.5 _(-2.5)	66.5 _(-3.5)	51.7 _(-0.6)

Table 9: The ablation study of GRADS using Llama-3.1-8B. “- Similarity” denotes removing $d^\top W^{KQ} q$. “- Information” denotes removing $W^{(PV)} d$.

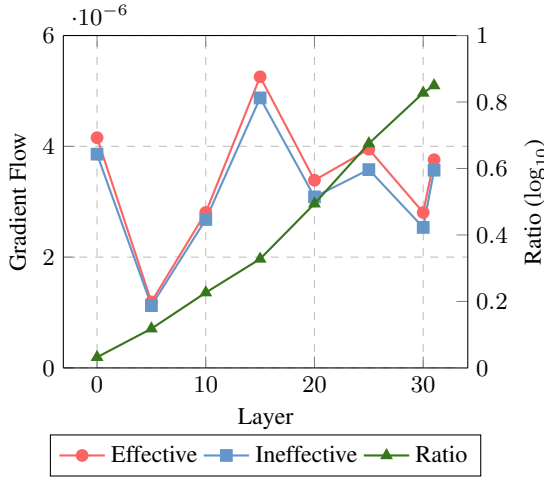


Figure 8: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-3.1-8B on Amazon Review.

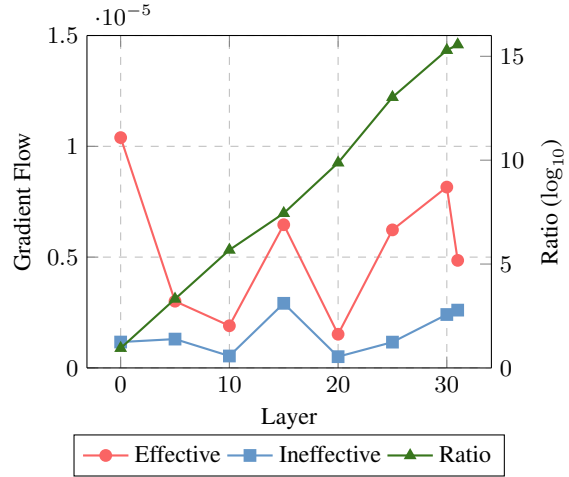


Figure 10: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-R1-8B on MATH.

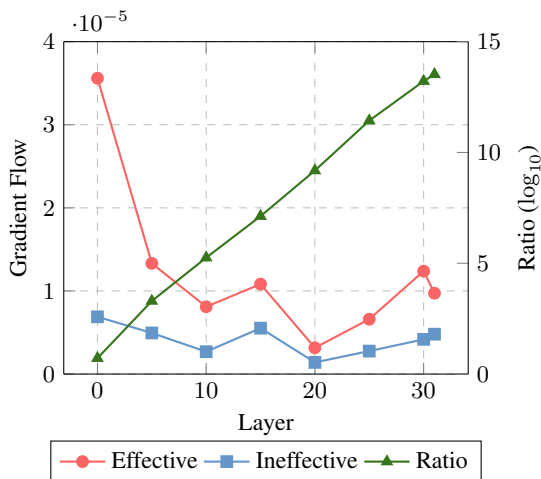


Figure 9: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-R1-8B on GSM8K.

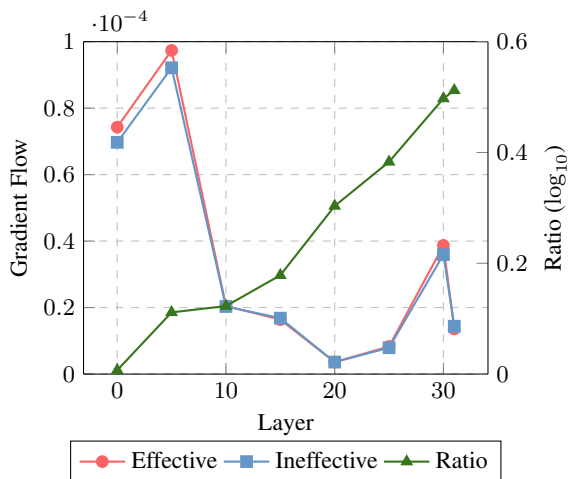


Figure 11: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-R1-8B on ARC-Challenge.

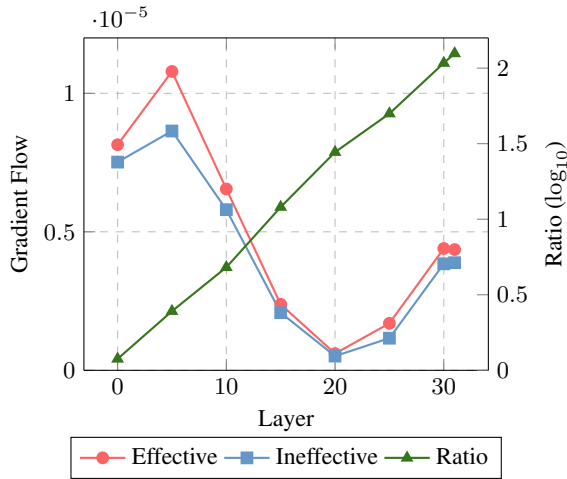


Figure 12: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-R1-8B on MMLU-Pro.

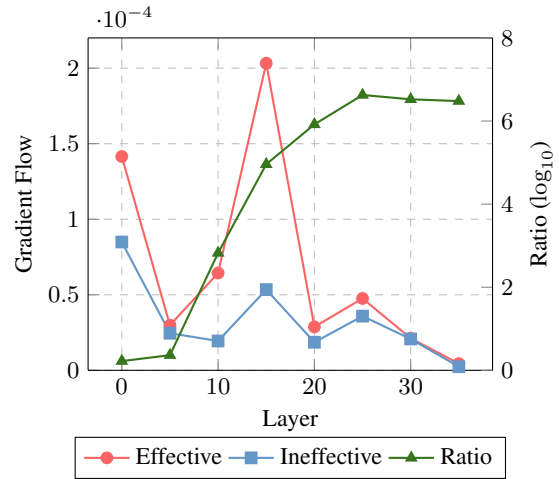


Figure 15: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Qwen3-8B on MATH.

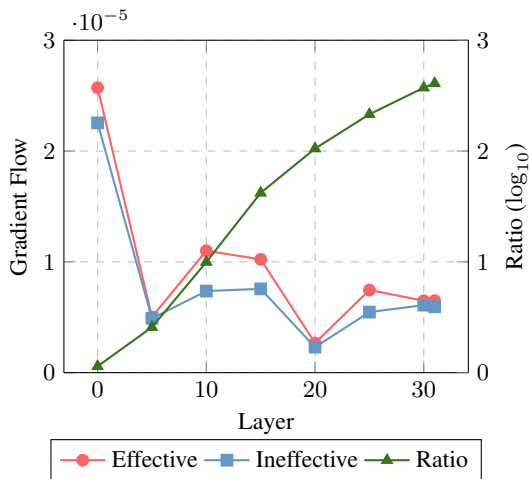


Figure 13: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Llama-R1-8B on Amazon Review.

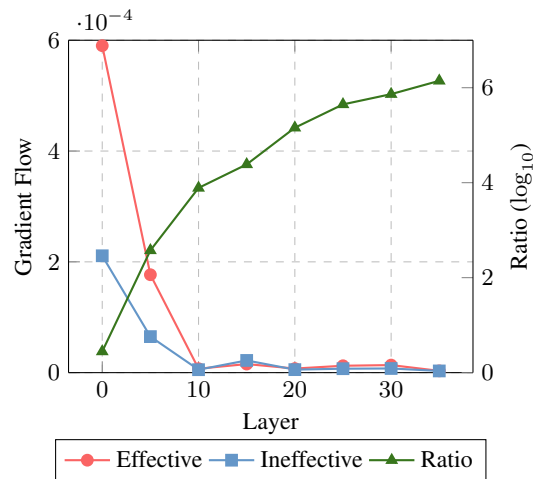


Figure 16: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Qwen3-8B on ARC-Challenge.

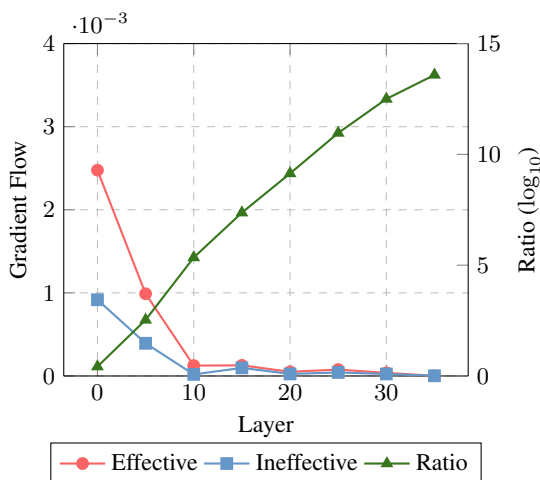


Figure 14: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Qwen3-8B on GSM8K.

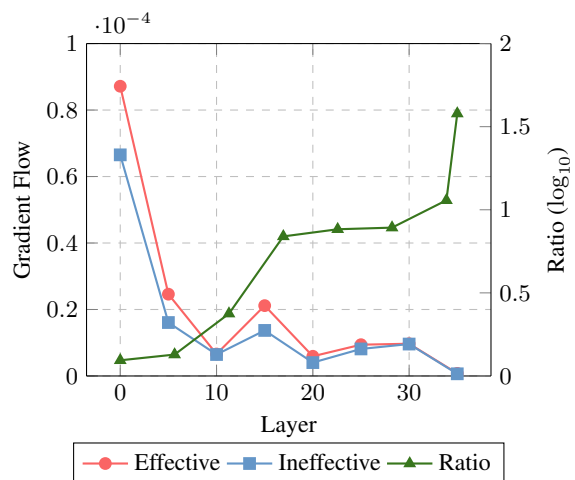


Figure 17: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Qwen3-8B on MMLU-Pro.

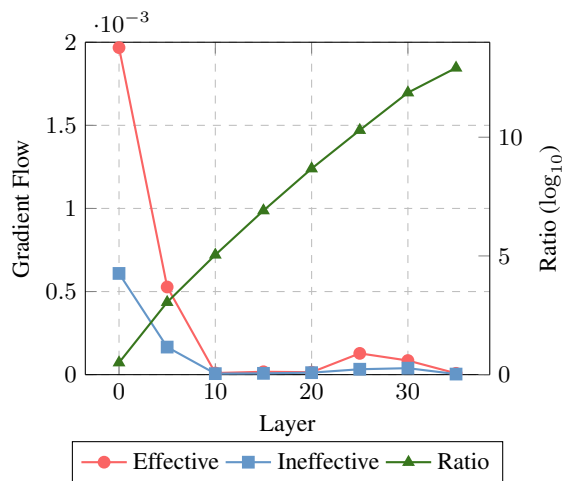


Figure 18: The average gradient flow (left y-axis) and the ratio between the effective and the ineffective (right y-axis) of Qwen3-8B on Amazon Review.