

Predict the Retrieval! Test Time Adaptation for Retrieval Augmented Generation

Anonymous ACL submission

Abstract

Retrieval-Augmented Generation (RAG) has emerged as a powerful approach for enhancing large language models' question-answering capabilities through the integration of external knowledge. However, when adapting RAG systems to specialized domains, challenges arise from distribution shifts, resulting in suboptimal generalization performance. In this work, we propose **TTARAG**, a test-time adaptation method that dynamically updates the language model's parameters during inference to improve RAG system performance in specialized domains. Our method introduces a simple yet effective approach where the model learns to predict retrieved content, enabling automatic parameter adjustment to the target domain. Through extensive experiments across six specialized domains, we demonstrate that **TTARAG** achieves substantial performance improvements over baseline RAG systems.

1 Introduction

Retrieval-Augmented Generation (RAG) (Izacard and Grave, 2021; Lewis et al., 2020; Edge et al., 2024) has emerged as a crucial approach for enhancing large language models (LLMs) (Radford et al., 2019; Brown et al., 2020; Bubeck et al., 2023) by addressing their inherent knowledge limitations. Through the integration of external knowledge sources (Pasca, 2019; Bollacker et al., 2008; Jin et al., 2019), RAG systems not only improve the accuracy of LLM responses but also help mitigate hallucination issues while eliminating the need for extensive model retraining.

However, while most current research has focused on the effectiveness of RAG systems for general domains, significant challenges persist in adapting these systems to specialized domains. These systems often struggle with distribution shifts and domain-specific data dependencies (Xu et al., 2025; Shi et al., 2024), frequently failing to

accurately utilize information in domain-specific contexts (Miller et al., 2020; Liu et al., 2022). This limitation is particularly problematic in critical domains such as healthcare (Raja et al., 2024), legal services (Reji et al., 2024), and financial applications (Yepes et al., 2024), where accuracy and reliability are paramount.

To address these challenges, test-time adaptation (TTA) (Sun et al., 2020; Hardt and Sun, 2024; Karmanov et al., 2024) offers a promising solution for enhancing model performance. TTA allows models to dynamically adapt their parameters at inference time through self-supervised learning objectives, without the need for labeled data (Chen et al., 2022; Liang et al., 2024). This approach is particularly valuable when dealing with domain shifts and distribution changes that weren't anticipated during initial training. Building on these insights, we propose a simple yet powerful method for adapting RAG systems during inference: **TTARAG**. Our approach generates self-supervised learning signals by dividing retrieved passages into prefix-suffix pairs and training the model to predict suffix content from prefix context. This technique enables LLMs to perform real-time parameter updates when encountering new domains, effectively leveraging domain knowledge stored within the model parameters.

Through extensive experiments across six specialized domains, we demonstrate that **TTARAG** achieves substantial performance improvements over baseline RAG systems. Our approach consistently outperforms both standard RAG and baselines like Chain-of-Thought and In-Context Learning, achieving the best results in 19 out of 24 experimental settings while maintaining computational efficiency. These results validate the effectiveness of our approach for domain-specific applications.

2 Methodology

Our approach introduces a test-time adaptation mechanism for retrieval-augmented generation that enables model optimization during inference without access to ground truth labels. The key innovation lies in designing a self-supervised learning objective using retrieved passages as supervision signals.

2.1 Overview

Given a test input query q and retrieved passages $\{p_1, \dots, p_k\}$, we formulate a self-supervised adaptation objective by splitting passages into prefix-suffix pairs for prediction:

$$\mathcal{L}_{adapt} = - \sum_{i=1}^k \log P(p_i^{suffix} | p_i^{prefix}, q; \theta) \quad (1)$$

where θ represents the model parameters.

2.2 Context Processing

The adaptation process begins with careful processing of the retrieved passages to create meaningful prefix-suffix pairs for training.

Length Filtering To ensure sufficient context for meaningful adaptation, passages shorter than a configured minimum length threshold are filtered out.

Passage Splitting Each passage is split into prefix-suffix pairs using a two-tier strategy:

- **Primary Strategy** Passages are split at first natural linguistic boundaries marked by punctuation (periods, commas, semicolons, colons, exclamation marks, and question marks)
- **Fallback Strategy** When no suitable punctuation-based split exists, the passage is divided at its midpoint, ensuring each segment contains at least three words.

2.3 Parameter Adaptation Process

The adaptation process employs a gradient-based optimization approach:

2.3.1 Initialization

Prior to the adaptation process, the model parameters are reset to their original pre-trained state to ensure a clean starting point for each adaptation iteration. An AdamW optimizer is then initialized with carefully configured hyperparameters: learning rate α for controlling update step sizes, epsilon ϵ for numerical stability, and weight decay λ for regularization.

2.3.2 Training Loop

For each batch of prefix-suffix pairs:

$$\theta_t = \theta_{t-1} - \alpha \cdot \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \mathcal{L}_{adapt}^i \quad (2)$$

where N is the gradient accumulation steps and \mathcal{L}_{adapt}^i is the loss for the i -th pair.

During training, the complete text (prefix and suffix) is first tokenized. The model then computes the loss on the suffix prediction task, where prefix tokens are masked during label preparation. To ensure stable training, gradients are accumulated over two steps and clipped to a maximum norm threshold. The AdamW optimizer then updates model parameters using these accumulated gradients. Since we only adapt on 1-5 prefix-suffix pairs in our experiments, the computational overhead remains acceptable.

2.4 Response Generation

After parameter adaptation, the model generates the final response using the adapted parameters θ' :

$$y = \arg \max_y P(y|q, \{p_1, \dots, p_k\}; \theta') \quad (3)$$

This approach enables effective domain adaptation through self-supervised learning on retrieved passages, allowing the model to dynamically align with the target domain during inference time without requiring ground truth labels.

3 Experiments

3.1 Datasets

We conduct experiments on CRAG (Yang et al., 2024) as the evaluation benchmark. CRAG is a comprehensive RAG benchmark containing 2,706 question-answer pairs across five domains: Finance, Sports, Music, Movie, and Open domain. The questions are constructed through web content-based creation where annotators formulate real-world questions answerable through web search.

To evaluate the effectiveness of TTARAG in the medical domain, we conduct additional experiments on two specialized datasets: PubMedQA (Jin et al., 2019), which contains 1,000 biomedical research question-answer pairs, and BioASQ (Tsatsaronis et al., 2015), comprising 500 expert-curated question-answer pairs from the biomedical literature.

Table 1: Performance comparison across different domains. Numbers represent accuracy scores (%). Best results for each model group are shown in **bold**.

Model	CRAG					Medical		
	Finance	Sports	Music	Movie	Open	Overall	BioASQ	PubMedQA
<i>Llama-3.1-8b-it</i>								
Base	17.4	27.6	34.9	31.3	42.4	29.8	55.6	46.6
CoT	17.9	30.2	37.6	31.5	45.8	31.6	54.6	50.8
ICL	16.1	24.8	33.5	29.4	40.4	28.0	49.8	53.6
TTARAG	20.1	29.5	37.7	34.6	41.5	31.9	75.0	57.4
Δ vs Base	+2.7	+1.9	+2.8	+3.3	-0.9	+2.1	+19.4	+10.8
<i>Llama-2-7b-chat</i>								
Base	14.7	23.2	36.5	30.4	39.2	27.8	54.1	47.6
CoT	15.7	26.7	34.3	31.4	41.5	29.1	55.1	48.2
ICL	16.0	24.2	36.1	31.2	39.2	28.4	55.6	43.4
TTARAG	16.4	25.8	40.7	33.8	41.1	30.5	71.8	54.0
Δ vs Base	+1.7	+2.6	+4.2	+3.4	+1.9	+2.7	+17.7	+6.4
<i>ChatGLM-3-6b</i>								
Base	9.8	18.7	31.4	22.4	33.4	22.0	51.4	19.8
CoT	12.7	20.6	28.4	25.8	33.9	23.6	44.3	22.4
ICL	9.9	18.2	30.8	22.1	33.0	21.8	50.8	19.2
TTARAG	14.0	22.1	33.5	25.5	38.1	25.7	58.4	44.8
Δ vs Base	+4.2	+3.4	+2.1	+3.1	+4.7	+3.7	+7.0	+25.0

3.2 Baselines

We evaluate TTARAG against several strong baselines, including prompting techniques (**Chain-of-Thought** (Wei et al., 2022), **In-Context Learning** (Brown et al., 2020)) and state-of-the-art pretrained RAG models (**Ret-Robust** (Yoran et al., 2024), **RAAT** (Fang et al., 2024), **Self-RAG** (Asai et al., 2023)). Detailed descriptions of each baseline are provided in Appendix B.

3.3 Experimental Results

Table 1 presents comprehensive evaluation results across different domains and model architectures. Several key observations emerge from our experiments: TTARAG demonstrates consistent improvements across specialized domains, with Llama-3.1-8b-it showing notable gains in Finance (+2.7%), Music (+2.8%), and Movie (+3.3%) domains, and particularly strong performance in medical domains (BioASQ +19.4%, PubMedQA +10.8%). All three model architectures benefit from our approach: Llama-3.1-8b-it achieves the highest overall accuracy (31.9%), Llama-2-7b-chat shows remarkable adaptation capability in medical domains

(+17.7% on BioASQ), and ChatGLM-3-6b demonstrates significant improvements in PubMedQA (+25.0%) and consistent gains across CRAG domains (+3.7% overall). While both CoT and ICL show some improvements over the base models, TTARAG consistently outperforms these baselines in specialized domains, with the only exception being Open domain tasks where CoT occasionally shows stronger performance, particularly with Llama-3.1-8b-it (45.8% vs 41.5%).

Table 2 presents a performance comparison between different RAG models across various domains. Notably, three of the models (Ret-robust, RAAT, and Self-rag) are pre-trained RAG models based on Llama-2. Despite Ret-robust using the larger Llama-2-13b as its base, and RAAT and Self-rag using Llama-2-7b, all three pre-trained RAG models perform worse than the Llama-2-7b-chat model (which achieves 27.8% overall accuracy). This underperformance is consistent across most domains, with only RAAT showing strength in the BioASQ medical domain (64.9%). The results suggest that current RAG pre-training methods have limited generalization capabilities, as they

Table 2: Performance comparison with state-of-the-art pretrained RAG models.

Model	CRAG						Medical	
	Finance	Sports	Music	Movie	Open	Overall	BioASQ	PubMedQA
Base	14.7	23.2	36.5	30.4	39.2	27.8	54.1	47.6
Ret-Robust	14.6	20.6	33.2	32.4	33.5	26.1	24.7	28.4
RAAT	13.4	18.1	28.6	25.2	31.7	22.7	64.9	46.6
Self-rag	11.4	19.8	22.5	20.9	26.7	19.8	57.1	43.4
TTARAG	16.4	25.8	40.7	33.8	41.1	30.5	71.8	54.0

fail to match or exceed the performance of the base model, even when using larger model architectures. TTARAG outperforms all other models across all domains, demonstrating the effectiveness of its approach compared to existing RAG pre-training methods.

The effectiveness of segment-based adaptation

We compare our segment-based approach (splitting passages into prefix-suffix pairs) with a baseline that does not segment the passage, where we perform next-token prediction on the entire passage without segmentation. The results in Table 3 demonstrate that the segmentation strategy yields consistent performance gains across all model architectures: +1.1% for Llama-3.1-8b-it, +0.4% for Llama-2-7b-chat, and +0.7% for ChatGLM-3-6b. We attribute these improvements to the front-to-back prediction task better aligning with natural language understanding compared to token-by-token prediction, enabling more effective parameter updates. The larger improvement observed with Llama-3.1-8b-it (+1.1%) suggests that higher-capacity models may particularly benefit from structured adaptation approaches.

Table 3: The effectiveness of segmentation.

Strategy	Llama-3.1-8b-it	Llama-2-7b-chat	ChatGLM-3-6b
TTARAG	31.9	30.5	25.7
wo seg	30.8	30.1	25.0

We also conduct hyper-parameter analysis about the number of adaptation pairs and learning rate in Section C.

On the computation efficiency To evaluate the computational overhead of our approach, we measure the total inference time across different configurations and compare it with baseline methods.

Table 4 shows the total and average inference times for different numbers of adaptation pairs (1-5), compared against Chain-of-Thought (CoT) and the original model without adaptation. The results are based on processing 2,706 queries from the CRAG dataset.

Table 4: Computation time analysis

Metric	1pair	2pair	3pair	4pair	5pair	CoT	Vanilla
Total	4,740	5,723	6,621	7,001	7,023	11,688	961
Avg	1.75	2.11	2.45	2.59	2.60	4.32	0.36

While our method does introduce additional computational overhead compared to the original model, it remains significantly more efficient than CoT. The average processing time per query ranges from 1.75s (1-pair) to 2.60s (5-pair), which is substantially lower than CoT’s 4.32s. This demonstrates that TTARAG achieves its performance improvements with reasonable computational cost, making it practical for real-world applications.

4 Conclusion

In this paper, we present TTARAG, a test-time adaptation approach for retrieval-augmented generation that enables dynamic model optimization during inference. Our method introduces a simple yet effective self-supervised learning objective where the model learns to predict retrieved content, allowing automatic parameter adjustment to target domains without requiring labeled data. Through extensive experiments across six specialized domain, we demonstrate that TTARAG achieves consistent improvements over the base RAG system, suggesting that test-time adaptation is a promising direction for improving RAG systems’ performance in specialized domains while maintaining computational efficiency.

275 **Limitations**

276 While TTARAG demonstrates strong performance
277 improvements across various domains, there are
278 several important limitations to consider:

279 The test-time adaptation process introduces ad-
280 ditional computational overhead during inference.
281 As shown in our experiments, the adaptation step
282 increases the average inference time by 1.75-2.60
283 seconds per query compared to the base model,
284 depending on the number of adaptation pairs
285 used. This additional latency may impact real-
286 time applications where response speed is criti-
287 cal. What’s more, our approach requires additional
288 GPU memory during inference for adaptation train-
289 ing compared to standard RAG systems. For larger
290 models, this increased memory requirement may
291 limit deployment options, particularly in resource-
292 constrained environments.

293 **Ethical Considerations**

294 Test-time adaptation may potentially affect the
295 model’s safety alignment due to parameter updates.
296 However, since our method only updates paramet-
297 ers for a limited number of iterations, the model’s
298 safety alignment likely remains largely intact, with
299 minimal risk of disruption. Nevertheless, we be-
300 lieve it is important to investigate the extent to
301 which gradient updates on domain-specific data
302 can impact a model’s established safety alignment
303 without compromising it. This represents an im-
304 portant direction for future research to better un-
305 derstand the relationship between adaptation and
306 safety preservation.

307 **References**

308 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
309 Hannaneh Hajishirzi. 2023. Self-rag: Learning to
310 retrieve, generate, and critique through self-reflection.
311 *arXiv preprint arXiv:2310.11511*.

312 Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim
313 Sturge, and Jamie Taylor. 2008. Freebase: a collabo-
314 ratively created graph database for structuring human
315 knowledge. In *Proceedings of the 2008 ACM SIG-
316 MOD international conference on Management of
317 data*, pages 1247–1250.

318 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann,
319 Trevor Cai, Eliza Rutherford, Katie Millican, George
320 van den Driessche, Jean-Baptiste Lespiau, Bogdan
321 Damoc, Aidan Clark, Diego de Las Casas, Aurelia
322 Guy, Jacob Menick, Roman Ring, Tom Hennigan,
323 Saffron Huang, Loren Maggiore, Chris Jones, Albin
324 Cassirer, Andy Brock, Michela Paganini, Geoffrey

Irving, Oriol Vinyals, Simon Osindero, Karen Si- 325
monyuan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 326
2022. [Improving language models by retrieving from 327](#)
[trillions of tokens](#). In *International Conference on 328*
Machine Learning, ICML 2022, 17-23 July 2022, Bal- 329
timore, Maryland, USA, volume 162 of *Proceedings 330*
of Machine Learning Research, pages 2206–2240. 331
PMLR. 332

Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and 333
Luca Bertinetto. 2022. Parameter-free online test- 334
time adaptation. In *IEEE Conference on Computer 335*
Vision and Pattern Recognition, pages 8344–8353. 336

Tom Brown, Benjamin Mann, Nick Ryder, Melanie 337
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 338
Neelakantan, Pranav Shyam, Girish Sastry, Amanda 339
Askell, et al. 2020. Language models are few-shot 340
learners. *Advances in neural information processing 341*
systems, 33:1877–1901. 342

Sébastien Bubeck, Varun Chandrasekaran, Ronen El- 343
dan, Johannes Gehrke, Eric Horvitz, Ece Kamar, 344
Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lund- 345
berg, et al. 2023. Sparks of artificial general intelli- 346
gence: Early experiments with gpt-4. *arXiv preprint 347*
arXiv:2303.12712. 348

Dian Chen, Dequan Wang, Trevor Darrell, and Sayna 349
Ebrahimi. 2022. Contrastive test-time adaptation. In 350
IEEE Conference on Computer Vision and Pattern 351
Recognition, pages 295–305. 352

Darren Edge, Ha Trinh, Newman Cheng, Joshua 353
Bradley, Alex Chao, Apurva Mody, Steven Truitt, 354
and Jonathan Larson. 2024. From local to global: A 355
graph rag approach to query-focused summarization. 356
arXiv preprint arXiv:2404.16130. 357

Feiteng Fang, Yuelin Bai, Shiwen Ni, Min Yang, Xiao- 358
jun Chen, and Ruifeng Xu. 2024. Enhancing noise 359
robustness of retrieval-augmented language models 360
with adaptive adversarial training. *arXiv preprint 361*
arXiv:2405.20978. 362

Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chen- 363
hui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Han- 364
lin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Ji- 365
adai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie 366
Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, 367
Lucen Zhong, Mingdao Liu, Minlie Huang, Peng 368
Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shu- 369
dan Zhang, Shulin Cao, Shuxun Yang, Weng Lam 370
Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan 371
Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, 372
Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan 373
An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, 374
Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, 375
Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan 376
Wang. 2024. [Chatglm: A family of large language 377](#)
[models from glm-130b to glm-4 all tools](#). *Preprint,* 378
arXiv:2406.12793. 379

Moritz Hardt and Yu Sun. 2024. [Test-time training 380](#)
[on nearest neighbors for large language models](#). 381
Preprint, arXiv:2305.18466. 382

599 adapted to target domains with limited labeled data.
600 Recent advances have extended this to fully un-
601 supervised scenarios, leveraging techniques like
602 entropy minimization (Wang et al., 2021), self-
603 training (Sun et al., 2020), and test-time normaliza-
604 tion statistics calibration (Schneider et al., 2020) to
605 adapt models using only unlabeled test samples.

606 Building on these advances, TTARAG introduces
607 a simple yet effective approach for test-time adapta-
608 tion in retrieval-augmented generation. By learning
609 to predict subsequent tokens in retrieved passages,
610 our method enables fully unsupervised adaptation
611 without requiring access to source domain data or
612 labeled examples.

613 **B Baseline Details**

614 We evaluate TTARAG using three state-of-the-art
615 instruction-tuned LLMs: **Llama-2-7b-chat** (Tou-
616 vron et al., 2023), **Llama-3.1-8b-it** (Meta-AI,
617 2024), and **ChatGLM3-6b** (GLM et al., 2024).
618 We compare against two widely adopted baselines:

619 **Chain-of-Thought (CoT)** A prompting tech-
620 nique that guides the model to generate step-by-
621 step reasoning before producing the final answer.

622 **In-Context Learning (ICL)** A method that pro-
623 vides relevant examples in the input prompt to
624 demonstrate the desired task behavior.

625 We also compare TTARAG with the three state-of-
626 the-art general domain pretrained RAG models:

627 **Ret-Robust** An approach focused on improving
628 retrieval robustness through strategic passage selec-
629 tion during training. The model learns to discrimi-
630 nate between high-quality and low-quality retrieved
631 content by being trained on a carefully curated mix
632 of passages with different relevance levels.

633 **RAAT** A retrieval-augmented model that intro-
634 duces a novel noise-aware training strategy. It
635 specifically targets the challenge of distinguishing
636 between helpful and misleading retrieved informa-
637 tion by incorporating an adaptive training mecha-
638 nism that exposes the model to varying types of
639 retrieval noise.

640 **Self-RAG** utilizes instruction fine-tuning to adap-
641 tively retrieve passages based on the question and
642 determine if the passage contains useful informa-
643 tion for answering the question.

644 **C Hyperparameter Analysis**

645 **Learning Rate Analysis** We investigate the sen-
646 sitivity of our method to different learning rates
647 during test-time adaptation with number of adapta-
648 tion pairs of 3. As shown in Figure 1, we evaluate
649 learning rates ranging from 1e-6 to 1e-4 across
650 all three model architectures. Llama-3.1-8B-it
651 achieves optimal performance at 1e-5 (31.9% ac-
652 curacy), with performance gradually declining at
653 higher learning rates. ChatGLM-6B shows more ro-
654 bust behavior across different learning rates, reach-
655 ing peak performance at 5e-6 to 1e-5 (25.8% ac-
656 curacy). Llama-2-7B-chat demonstrates the most
657 stable performance curve, with accuracy varying
658 only slightly (30.4-30.8%) across all tested learn-
659 ing rates, peaking at 1e-6 (30.8% accuracy). These
660 results suggest that smaller learning rates (1e-6 to
661 1e-5) generally provide better and more stable adap-
662 tation, likely because they prevent over-aggressive
663 parameter updates that could disrupt the model’s
664 pre-trained knowledge. All models show consist-
665 ent improvement over their original performance
666 (indicated by dashed lines) across most learning
667 rates, validating the robustness of our approach.

668 **Number of Adaptation Passages** We examine
669 how the number of retrieved passages used for
670 adaptation affects performance. This study helps
671 determine the optimal amount of context needed
672 for effective adaptation while considering computa-
673 tional efficiency. As shown in Figure 2, we observe
674 different optimal points across model architectures.
675 Llama-3.1-8B-it achieves peak performance with
676 3 adaptation pairs (31.7% accuracy), while Llama-
677 2-7B-chat shows optimal results at 4 pairs (31.7%
678 accuracy). ChatGLM-6B maintains relatively sta-
679 ble performance between 2-5 pairs, peaking at 5
680 pairs (25.8% accuracy). Notably, all models show
681 performance degradation when using 10 pairs. This
682 degradation likely stems from over-aggressive pa-
683 rameter updates that disrupt the model’s pre-trained
684 knowledge. Too many adaptation pairs may cause
685 excessive deviation from the original parameters,
686 compromising the valuable knowledge acquired
687 during pre-training. These results indicate that a
688 moderate number of adaptation pairs (3-5) gen-
689 erally provides the best balance between adapta-
690 tion effectiveness and preserving the model’s pre-
691 trained knowledge.

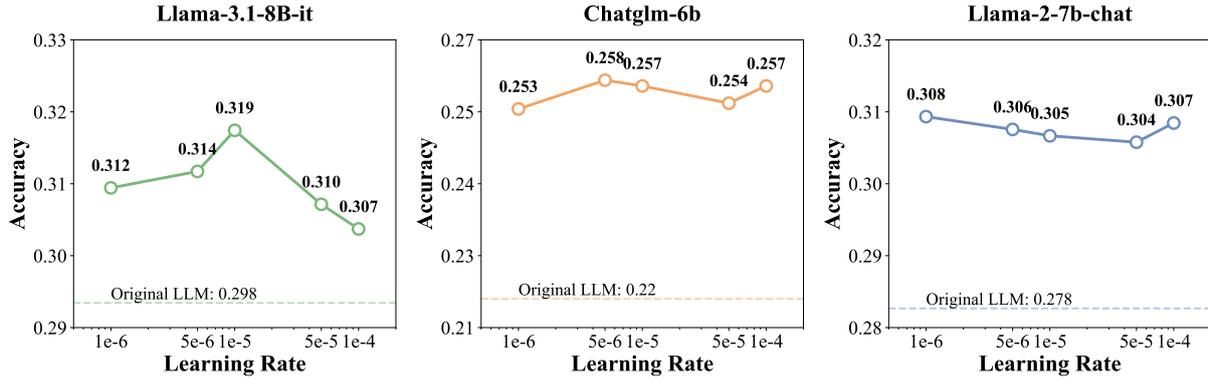


Figure 1: Accuracy vs. Learning Rate

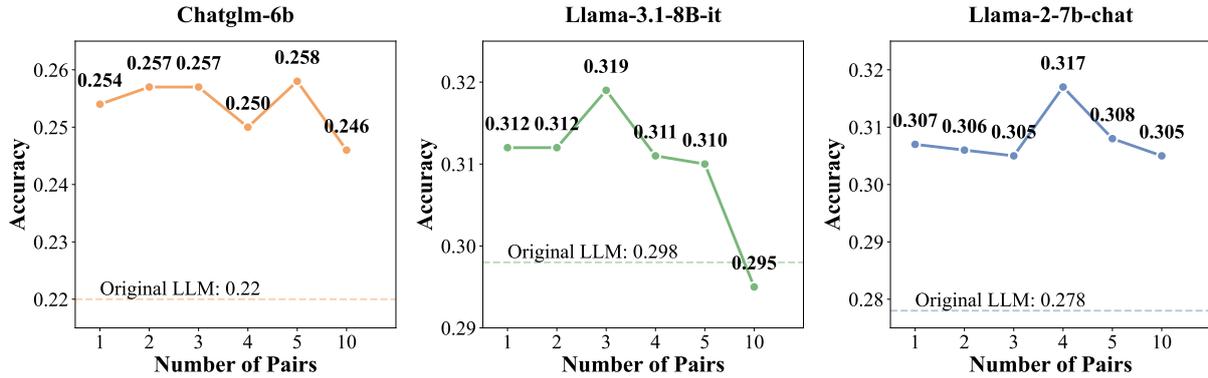


Figure 2: Accuracy vs. Number of Adaptation Pairs

D Implementation Details

We use Llama-3.1-8b-instruct, ChatGLM-3-6b, Llama-2-7b-chat as our backbone models. Here we detail the hyperparameters and configuration settings used in our implementation. For the optimization process, we employ a learning rate of $1e-5$, which provides a balance between adaptation effectiveness and stability. To improve training efficiency while managing memory constraints, we implement gradient accumulation with 2 steps. Gradient clipping is set at 0.1 to prevent gradient explosions, particularly important during rapid adaptation to new contexts. We use the AdamW optimizer with weight decay of 0.01 and epsilon of $1e-8$, which helps prevent overfitting while maintaining numerical stability. Additional controls include filtering out sentences shorter than 6 tokens and limiting adaptation to 3 pairs per step. These parameters were determined through extensive experimentation across various domains, optimizing for both adaptation performance and computational efficiency. All experiments were conducted three times and the average results are reported.

All experiments are conducted on NVIDIA

A100 GPUs with 80GB of memory. We utilize a fixed random seed of 42, and the experimental results are reported within a single run. For implementation, we use the following library versions: transformers 4.30.2, torch 2.1.0.

Table 5: Number of samples in each domain of CRAG dataset.

Domain	Finance	Sports	Music	Movie	Open
#Samples	661	519	373	611	542

E Dataset Statistics

The statistics of the CRAG dataset are shown in Table 5.

F Licensing

The CRAG, BioASQ and PubMedQA datasets are released for academic usage. These datasets are designed for evaluating RAG systems. Thus, our use of these datasets is consistent with their intended use.

The language models used in our experiments are released under the following licenses: **Llama-**

732 **2-7b-chat** (Touvron et al., 2023) is released un-
733 der the Meta Llama 2 Community License Agree-
734 ment. It is a variant of the Llama 2 family released
735 in July 2023, featuring 7 billion parameters and
736 specifically optimized for dialogue applications.
737 **Llama-3.1-8b-it** (Meta-AI, 2024) is released un-
738 der the Llama 3 License. Released in April 2024,
739 it features 8 billion parameters and is specifically
740 designed for instruction-following tasks, represent-
741 ing one of the most advanced open-source LLMs.
742 **ChatGLM3-6b** (GLM et al., 2024) is released un-
743 der the Apache 2.0 License. It is a bilingual conver-
744 sational language model featuring 6 billion param-
745 eters, demonstrating strong performance in both
746 English and Chinese tasks. All these models are
747 open for academic usage.