

# Flexi-LoRA: Efficient LoRA Finetuning with Input-Adaptive Dynamic Ranks

Anonymous ACL submission

## Abstract

Parameter-efficient fine-tuning methods like Low-Rank Adaptation (LoRA) have become essential for deploying large language models, yet their static parameter allocation remains suboptimal for inputs of varying complexity. We present Flexi-LoRA, a novel framework that dynamically adjusts LoRA ranks based on input complexity during both training and inference. Through empirical analysis across question answering and mathematical reasoning tasks, we demonstrate that maintaining consistency between training and inference dynamics is important for effective adaptation, particularly for sequential reasoning tasks. Our findings reveal that input-dependent parameter allocation achieves superior performance with fewer parameters by optimally matching rank configurations to question complexity. Furthermore, task-specific sensitivity to rank dynamics varies, with mathematical reasoning tasks exhibiting higher sensitivity than QA tasks. Successful adaptation manifests not only in correctness but also in reasoning quality and instruction adherence. Flexi-LoRA consistently outperforms static LoRA while using fewer parameters, with performance gains more pronounced on tasks requiring strict reasoning chains. Our approach realizes key benefits of mixture-of-experts frameworks through a more streamlined implementation, reducing parameter redundancy while enhancing model capabilities. We provide comprehensive empirical studies across diverse tasks, establishing a foundation for future work in input-adaptive and efficient fine-tuning approaches.<sup>1</sup>

## 1 Introduction

As large language models grow in size, efficient fine-tuning methods like LoRA (Hu et al., 2022)

<sup>1</sup><https://github.com/Anonymous/Flexi-LoRA> (After the paper is published, this link will be de-anonymized. For related code, please see the supplementary materials.)

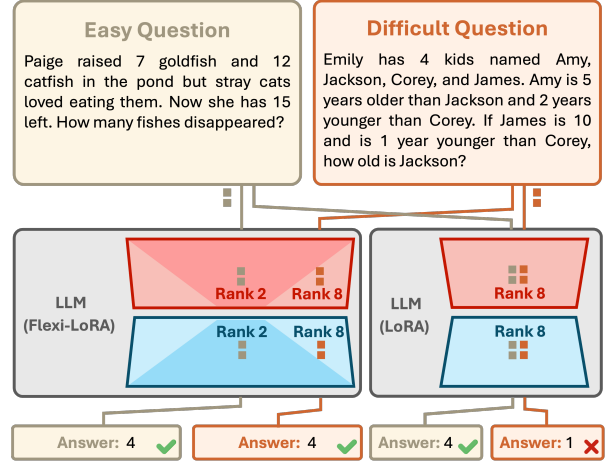


Figure 1: Flexi-LoRA’s input-adaptive rank allocation versus static normal LoRA. Flexi-LoRA (left) dynamically assigns rank 2 (dark trapezoid) for simple problems and rank 8 (light trapezoid) for complex ones, successfully solving both. LoRA (right) uses fixed rank 8 (light trapezoid) regardless of complexity, failing on difficult problems. This demonstrates the necessity of input-adaptive parameter allocation for handling varying question complexity.

have become essential for applications. However, their static parameter allocation remains suboptimal for questions of varying complexity, suggesting the need for input-adaptive approaches in parameter-efficient fine-tuning (Jiang et al., 2025).

Through empirical analysis, we observe two key **phenomena** in LoRA-based fine-tuning. **First**, there exists a notable performance gap when using static ranks during inference for models trained with dynamic ranks at fine-grained level, particularly in their ability to follow instructions precisely (DyLoRA vs DyLoRA+, Table 2 and 3, Figure 4). **Second**, while model performance generally saturates with increasing ranks, the optimal rank varies across different inputs: simple questions can be effectively handled with small ranks, while complex problems benefit substantially from larger ranks (Rank 4 vs 8, Table 2 and 3, Figure 3). These observations indicate that a one-size-fits-all approach to

rank selection is suboptimal, motivating the need for input-adaptive rank allocation.

Inspired by these observations, we propose **Flexi-LoRA**, a finetuning framework that dynamically adjusts LoRA ranks based on input complexity. Our approach not only achieves superior performance to high-rank LoRA while using fewer parameters, but also successfully solves some complex problems that static LoRA fails to handle even with equivalent rank, as shown in Figure 1.

Our **contributions** are threefold:

- **Novel Framework:** We introduce the first input-adaptive LoRA framework that maintains dynamic ranks during both training and inference, achieving higher performance with reduced parameter count compared to static LoRA.
- **Insights:** We demonstrate that (1) maintaining consistency between training and inference dynamics is important for LoRA adaptation, particularly for sequential reasoning tasks; (2) input-dependent parameter allocation achieves superior performance with fewer parameters by optimally matching rank configurations to question complexity; (3) task-specific sensitivity to rank dynamics varies, with mathematical reasoning tasks exhibiting higher sensitivity than QA tasks; (4) successful adaptation manifests not only in correctness but also in reasoning quality and instruction adherence; and (5) our approach realizes benefits of mixture-of-experts through a more streamlined implementation, reducing parameter redundancy while enhancing model capabilities.
- **Comprehensive Analysis:** We provide comprehensive empirical studies across diverse tasks, establishing a foundation for future work in input-adaptive and efficient finetuning approaches.

## 2 Related Work

**LoRA with dynamic ranks.** Recent works have explored dynamic rank adaptation in LoRA, with differences shown in Table 1. AdaLoRA (Zhang et al., 2023) performs importance-based parameter pruning at training checkpoints to gradually reduce ranks to a fixed target. DyLoRA (Valipour et al., 2023) randomly samples ranks from a predefined range for each training batch, with all samples in the batch sharing the same rank. Both approaches, while improving rank flexibility, are limited by either steps-level pruning or random batch-level assignment, and neither supports dynamic rank selection at inference. On the other hand, Flexi-LoRA

Method	Train	Level	Inference
LoRA	Fixed	All	Fixed
AdaLoRA	Selective	Steps	Fixed
DyLoRA	Random	Batch	Fixed
<b>DyLoRA+ (O)</b>	Random	Batch	<b>Random</b>
<b>Flexi-LoRA (O)</b>	<b>Router</b>	<b>Sample</b>	<b>Router</b>

Table 1: Comparison of rank adaptation strategies across different LoRA variants. "Train" indicates how ranks are determined during training, "Level" shows the level of rank assignment, and "Inference" specifies the rank selection mechanism at test time. Only Flexi-LoRA maintains consistent router-based sample-level dynamic rank allocation across both training and inference stages, while existing methods use fixed ranks during inference regardless of training dynamics. "O" is our method.

enables true sample-level rank selection by learning to map input complexity to appropriate ranks, maintaining this adaptive behavior during both training and inference.

**Input-Adaptive Methods.** Due to space constraints in the main text, a review of input-adaptive methods is provided in Appendix A.1.

## 3 Methods

Building upon previous work, we first introduce **DyLoRA+**, an enhanced variant of DyLoRA that maintains consistent rank dynamics by employing random batch-level rank selection during both training and inference stages. While DyLoRA+ demonstrates improved performance over the original DyLoRA, its random rank allocation remains suboptimal as it fails to account for input-specific complexity differences. We therefore propose **Flexi-LoRA**, a framework that automatically adjusts the rank based on input complexity. Our method consists of two key components: a difficulty-aware router that maps inputs to appropriate rank assignments and a flexible-rank LoRA training framework that maintains consistent dynamic rank allocation during both training and inference, as shown in Figure 2.

**Router** focuses on learning an optimal mapping  $R(h) : \mathbb{R}^d \rightarrow r_i$  from input embeddings to rank assignments. Given an input sequence  $x$  with mask  $m$ , we first compute its token embeddings  $H \in \mathbb{R}^{n \times d}$  and obtain a pooled embedding  $h = \sum_i (m_i H_i) / \sum_i m_i$ , where  $m_i$  masks padding tokens. We categorize training samples into difficulty classes based on task-specific metrics: F1 scores for MRQA datasets and accuracy for mathematical reasoning tasks. The router is then optimized using a noise-added cross-entropy

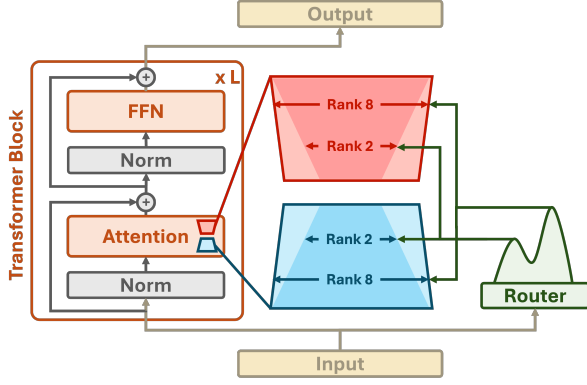


Figure 2: Flexi-LoRA framework with input-adaptive rank selection. Router analyzes input embeddings and outputs rank assignments (green arrows) for transformer layers. Red and blue trapezoids are LoRA’s A and B matrices, with color darkness indicating rank magnitude (darker = rank 2, lighter = rank 8). The router enables dynamic rank allocation based on input complexity while maintaining efficient gradient flow through residual connections.

objective:  $\mathcal{L}(\theta) = -\sum_i y_i \log(R(h_i + \epsilon))$ , where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  is Gaussian noise and  $y_i$  denotes the ground-truth difficulty label. The training data is balanced between easy and hard samples to ensure uniform difficulty evaluation.

**Input-adaptive LoRA** freezes the base model parameters and optimize only the LoRA matrices. For input  $x$ , we first obtain its token embeddings  $H^0$  and pooled embedding  $h$  following the same procedure as router training. The router then predicts rank  $r = R(h)$ , which is applied consistently across all transformer layers. Within each batch, different samples can be assigned different ranks based on their predicted difficulty, enabling dynamic resource allocation. For each transformer layer  $l$ , we compute the LoRA update as  $\Delta W_l = B_{l,r} A_{l,r}$ , where  $A_{l,r} \in \mathbb{R}^{r \times d}$  and  $B_{l,r} \in \mathbb{R}^{d \times r}$  are dynamically reduced to the first  $r$  rows/columns. The layer output is computed as  $H^l = W_l H^{l-1} + \alpha_r (B_{l,r} A_{l,r} H^{l-1})$ , where  $\alpha_r$  is a rank-specific scaling variable and  $H^{l-1}$  is the output from the previous layer. The model is trained to minimize the task-specific loss  $\mathcal{L}_{\text{task}} = -\sum_i \log p(y_i | x_i)$ , where  $x_i$  is the input sequence and  $y_i$  is the corresponding ground truth outputs. This design enables efficient batch processing while allowing for flexible input-dependent rank adaptation.

## 4 Experimental Design

**Datasets.** We evaluate Flexi-LoRA on both QA and mathematical reasoning tasks. For QA tasks, we conduct training on datasets from the MRQA (Fisch et al., 2019) training set, which unifies QA samples from SQuAD (Rajpurkar et al., 2016), TriviaQA (Joshi et al., 2017), NewsQA (Trischler et al., 2017), SearchQA (Dunn et al., 2017), HotpotQA (Yang et al., 2018), and NaturalQuestions (Kwiatkowski et al., 2019). Evaluation is performed on the MRQA test set consisting of BioASQ (Partalas et al., 2013), DROP (Dua et al., 2019), DuoRC (Saha et al., 2018), RACE (Lai et al., 2017), RelationExtraction (Levy et al., 2017), and TextbookQA (Kembhavi et al., 2017). For mathematical reasoning, we train on the GSM8K (Cobbe et al., 2021) subset of the MetaMathQA (Yu et al., 2024) dataset and evaluate on a diverse set of math benchmarks including GSM8K, SVAMP (Patel et al., 2021), MultiArith (Roy and Roth, 2015), and MAWPS (Koncel-Kedziorski et al., 2016). This design allows us to evaluate both in-distribution and out-of-distribution generalization capabilities of our method.

**Evaluation Metrics.** We evaluate QA performance using **F1** and **Exact Match (EM)** scores. F1 computes the balanced average of precision and recall between prediction and ground truth, while EM measures exact string match. For **mathematical reasoning** tasks, we use **accuracy** for evaluation.

**Gold Standard & Baselines.** We compare against two **gold** standards: **full** model fine-tuning and standard **LoRA** with fixed rank. For **baselines**, we include **AdaLoRA**, which adapts ranks through importance-based parameter pruning while maintaining fixed inference ranks, and **DyLoRA**, which randomly samples ranks from a predefined range for each training batch but uses fixed ranks during inference. Our Flexi-LoRA differs by enabling input-adaptive rank selection during both training and inference.

**Models.** We employ **LLaMA-3.2-1B-Instruct** (Grattafiori et al., 2024) as the base model for our main results and include **LLaMA-3.2-3B-Instruct** to analyze model size in ablation studies.

## 5 Results

### 5.1 Overview

Figure 3 illustrates the performance-efficiency trade-offs across different parameter-efficient fine-tuning methods on QA and mathematical reasoning

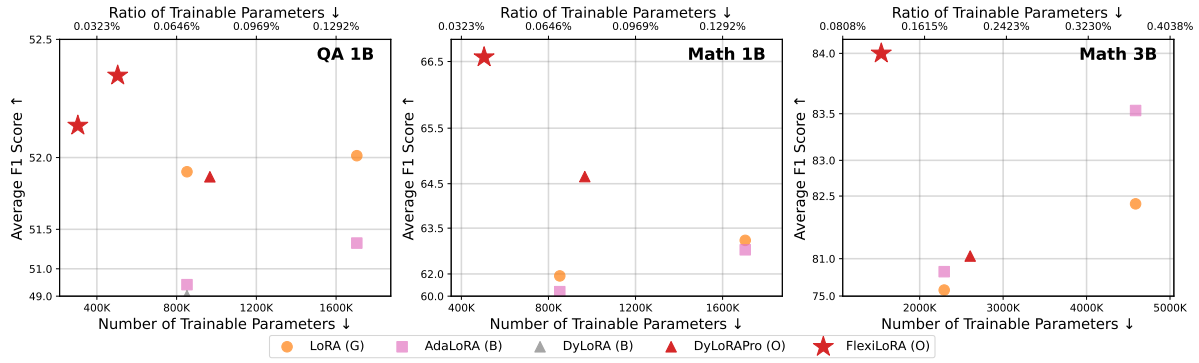


Figure 3: Performance-efficiency trade-off across different parameter-efficient fine-tuning methods. Our methods (O) achieve superior performance with fewer trainable parameters compared to both baseline methods (B) as well as gold standards (G). Results are shown for QA tasks and mathematical tasks using LLaMA-3.2-1B-Instruct and 3B.

Rank #	Gold Standard LoRA			Baselines				Ours		
	Full 1.2B	4 851K	8 1703K	4 851K	8 1703K	4 851K	8 1703K	DyLoRA+ 1-8 966K	Flexi-LoRA 1, 8 304K	2, 8 504K
<b>F1 Score</b>										
BioASQ	69.81	64.85	66.22	63.62	65.12	60.01	52.40	65.21	65.75	65.82
DROP	47.15	37.88	36.03	32.79	34.35	43.59	38.26	36.32	37.27	37.52
DuoRC	45.21	43.85	43.72	44.00	42.49	39.10	36.89	43.86	42.92	43.22
RACE	41.49	38.64	37.49	34.58	35.81	36.89	33.24	39.26	38.83	39.10
RE	84.11	74.97	76.41	74.51	75.73	81.02	78.70	75.47	76.47	76.83
TextbookQA	49.56	51.35	52.20	54.14	54.66	37.90	32.94	51.22	51.74	51.75
Average	56.22	51.92	52.01	50.61	51.36	49.75	45.40	51.89	52.16	52.37
<b>Exact Match</b>										
BioASQ	49.13	42.02	42.61	40.29	40.35	34.24	27.52	41.88	41.42	41.48
DROP	35.46	25.81	23.08	20.15	21.49	30.07	24.01	23.61	24.55	25.01
DuoRC	35.57	32.37	31.51	32.37	30.77	27.71	24.58	32.31	30.44	30.71
RACE	29.37	24.03	22.99	20.62	22.10	21.81	16.61	24.92	24.48	24.77
RE	72.01	57.73	60.27	57.25	58.44	69.13	66.31	59.15	59.90	60.44
TextbookQA	40.98	41.91	42.38	44.97	45.10	27.01	21.29	41.91	42.04	42.04
Average	43.75	37.31	37.14	35.94	36.38	34.99	30.05	37.30	37.14	37.41

Table 2: Performance comparison on out-of-domain QA tasks from the MRQA benchmark using LLaMA-3.2-1B-Instruct. F1 and Exact Match (EM) scores are reported, comparing our proposed methods (Flexi-LoRA and DyLoRA+) against gold standards and baselines. Flexi-LoRA (2,8) achieves the best average performance on both metrics while using fewer parameters than standard approaches. The "#" row indicates the number of trainable parameters. Green (teal) and red (maroon) cell coloring is higher and lower scores respectively, with deeper colors indicating larger performance differences.

tasks. Flexi-LoRA consistently achieves superior performance while requiring fewer parameters than competing approaches. Notably, DyLoRA results are not visible from mathematical reasoning figures due to substantially decreased performance, highlighting the importance of maintaining consistency between training and inference dynamics. AdaLoRA shows competitive results on specific tasks but fails to achieve a consistent advantage across different domains. From a Pareto optimality perspective, Flexi-LoRA dominates all baselines and gold standards by offering better performance at lower parameter counts, positioning itself closest to the full fine-tuning performance while maintaining parameter efficiency below 0.1% of total model parameters. The following sections analyze these results in detail across different task categories and provide ablation studies to analyze the contribu-

tions of individual components.

## 5.2 Question Answering

Table 2 presents the performance of different parameter-efficient fine-tuning methods on six out-of-domain QA datasets from the MRQA benchmark. We analyze these results from multiple perspectives: (1) **Overall Performance**: Flexi-LoRA (2,8) achieves the highest average F1 score (52.37%) and EM score (37.41%), outperforming both LoRA-8 with only 29.59% of LoRA-8’s parameters. The performance gap between Flexi-LoRA and full fine-tuning is considerably smaller than that of other parameter-efficient methods, demonstrating its effectiveness in approaching full fine-tuning. (2) **Stability Across Metrics**: Flexi-LoRA demonstrates consistency by achieving the best average performance on both



Rank	Gold Standard			Baselines				Ours	
	Full	LoRA		AdaLoRA		DyLoRA		DyLoRA+	Flexi-LoRA
	-	4	8	4	8	4	8	1-8	2, 8
LLaMA-3.2-1B-Instruct									
#	1.2B	851K	1703K	851K	1703K	851K	1703K	953K	533K
GSM8K	57.31	42.15	41.31	45.71	42.22	19.78	21.00	41.77	42.30
SVAMP	57.29	52.87	51.18	53.37	50.52	20.08	19.57	56.03	52.02
MultiArith	93.88	82.77	85.00	78.88	84.44	43.88	42.77	85.55	92.22
MAWPS	80.00	69.85	75.21	65.63	74.36	22.25	17.46	75.21	79.71
Average	72.12	61.91	63.17	60.90	62.89	26.50	25.20	64.64	66.56
LLaMA-3.2-3B-Instruct									
#	1.2B	2.29M	4.58M	2.29M	4.58M	2.29M	4.58M	2.60M	1.53M
GSM8K	75.20	65.95	69.37	71.49	72.32	60.04	60.87	70.43	69.90
SVAMP	78.30	66.29	74.47	75.91	77.80	67.48	64.87	74.71	77.09
MultiArith	100	97.77	99.44	90.55	95.55	75.00	82.22	92.77	100
MAWPS	89.85	84.50	86.19	83.66	88.45	69.57	65.63	86.47	89.01
Average	85.84	78.63	82.37	80.40	83.53	68.02	68.40	81.10	84.00

Table 3: Performance comparison on **mathematical tasks** using LLaMA-3.2-1B-Instruct and LLaMA-3.2-3B-Instruct models. Accuracy scores (%) are reported across four benchmark datasets (GSM8K as in-domain and the others as out-of-domain), comparing our proposed methods against gold standards and baselines. Flexi-LoRA achieves the best average performance on both model sizes while using fewer parameters than standard approaches. The "#" row indicates the number of trainable parameters. Green (teal) and red (maroon) cell coloring is higher and lower scores respectively, with deeper colors indicating larger performance differences.

F1 and EM metrics simultaneously, unlike other methods that typically excel in one metric. This dual-metric superiority indicates that Flexi-LoRA produces outputs that are both semantically close to ground truth (high F1) and syntactically precise (high EM). (3) **Dataset-Specific Analysis:** Flexi-LoRA performs particularly well on domain-specific and knowledge-based tasks, while showing moderate improvements on information extraction tasks. The adaptive rank allocation proves beneficial for datasets requiring diverse reasoning capabilities, suggesting Flexi-LoRA’s ability to assign appropriate computational resources based on input complexity. (4) **Comparison with Other Dynamic Approaches:** Flexi-LoRA consistently outperforms DyLoRA+, despite both using dynamic ranks during inference, demonstrating the importance of learned input-adaptive rank assignment versus the random assignment in DyLoRA+. DyLoRA shows high variance across datasets (from 81.02% F1 on RE to 32.94% F1 on TextbookQA), highlighting instability issues when training and inference dynamics are inconsistent. (5) **Rank Configuration Influence:** The comparison between Flexi-LoRA (1,8) and (2,8) configurations reveals that a slight increase in the minimum rank (from 1 to 2) provides modest but consistent performance improvements (from 52.16% to 52.37% F1), with a reasonable parameter increase (from 304K to 504K). This suggests that while low ranks can handle simpler questions, maintaining a slightly higher minimum rank improves robustness across diverse question types. (6) **Cross-Domain Generalization:** On out-of-domain test datasets, baseline methods exhibit inconsistent performance, ex-

celling on specific domains while decreasing on others. On the other hand, Flexi-LoRA maintains strong performance across all test datasets, suggesting its input-adaptive parameter allocation learns more generalized knowledge than static approaches that could overfit to training domain characteristics. (7) **Key Insights:** These results demonstrate that input-adaptive parameter allocation provides dual benefits in QA tasks: improved performance through parameter allocation and enhanced parameter efficiency through optimization of rank selection. The consistent improvement across diverse datasets suggests that question complexity varies substantially even within the same task category, validating our input-adaptive approach.

### 5.3 Math Problems

Table 3 presents the performance of different fine-tuning methods on mathematical reasoning tasks across two model sizes. We analyze these results for task-specific characteristics: (1) **Overall Performance:** Flexi-LoRA achieves the highest average accuracy on both LLaMA-3.2-1B-Instruct (66.56%) and LLaMA-3.2-3B-Instruct (84.00%), outperforming LoRA-8 (63.17% and 82.37%) while using only 31.29% and 33.40% of its parameters, respectively. The performance gap between Flexi-LoRA and full fine-tuning is narrower than other parameter-efficient methods, particularly on the larger model. (2) **Model Size Influence:** Increasing model size from 1B to 3B improves performance across all methods, with Flexi-LoRA maintaining its advantage. Notably, the absolute performance gap between Flexi-LoRA and full fine-tuning decreases from 5.56% to 1.84%, suggest-

ing that input-adaptive approaches become even more effective with larger models. (3) **In-Domain vs. Out-of-Domain:** Flexi-LoRA shows strong generalization from GSM8K (in-domain) to out-of-domain datasets. On the 1B model, Flexi-LoRA achieves an average accuracy of 74.65% on out-of-domain tasks (SVAMP, MultiArith, MAWPS), compared to 42.30% on in-domain GSM8K, demonstrating cross-domain robustness. This is consistent across both model sizes. (4) **Dataset Complexity:** Flexi-LoRA performs exceptionally well on both elementary arithmetic (MultiArith) and complex multi-step reasoning (GSM8K), indicating its ability to effectively handle varying levels of mathematical complexity. The exceptional performance on MultiArith (92.22%) approaches full fine-tuning (93.88%), showcasing the method’s ceiling capability on reasoning tasks. (5) **DyLoRA Performance Decrease:** DyLoRA exhibits performance decrease on mathematical tasks (average 26.50% on 1B model). The 40.06% performance gap between DyLoRA and Flexi-LoRA on math tasks highlights the influence of training-inference inconsistency on sequential reasoning tasks. This decrease is substantially more pronounced than in QA tasks, suggesting that mathematical reasoning is particularly sensitive to dynamic rank consistency. (6) **Key Insights:** These results demonstrate that maintaining consistent training-inference dynamics is important for mathematical reasoning tasks. The substantial performance improvements (Flexi-LoRA outperforms LoRA-8 by 3.39% on 1B) illustrate that input-adaptive parameter allocation provides greater benefits for problems with higher complexity variance and stricter evaluation criteria, compared to the modest gains observed on QA tasks.

#### 5.4 Cross-Task Analysis

Comparing Flexi-LoRA’s performance on QA and mathematical reasoning tasks reveals important insights about the relationship between task characteristics and parameter-efficient fine-tuning approaches: (1) **Task Nature Influence:** Task characteristics influence adaptation strategy efficacy. QA primarily involves information extraction, whereas mathematical reasoning demands coherent computational chains where early errors go through solutions. This sequential dependency in math tasks necessitates consistency between training and inference dynamics, unlike the more information retrieval in QA. (2) **Performance Gain Differ-**

**ence:** Flexi-LoRA’s improvement over LoRA-8 is substantially larger on mathematical tasks than on QA tasks. This difference suggests that input-adaptive parameter allocation yields greater benefits for sequential reasoning tasks where capacity requirements vary between simple and complex problems. (3) **Error Analysis:** QA and mathematical reasoning tasks demonstrate distinct error behaviors. QA permits partial correctness (evidenced by F1 scores exceeding EM scores), whereas math tasks produce binary outcomes. This evaluation in mathematical reasoning highlights the benefits of input-adaptive approaches that allocate capacity proportionally to problem complexity. (4) **Input Complexity Distribution:** Flexi-LoRA’s superior performance on mathematical tasks indicates more pronounced complexity variations compared to QA. Math problems encompass a broader difficulty distribution, from elementary arithmetic to multi-step reasoning, making them optimal candidates for adaptive parameter allocation. (5) **Cross-Task Generalization:** Flexi-LoRA’s consistent superior performance across both QA and mathematical reasoning tasks, despite their distinct nature and evaluation metrics, validates the robustness and generality of our input-adaptive approach. This multi-domain effectiveness indicates that the underlying complexity-aware parameter allocation mechanism learns fundamental aspects of language model adaptation rather than utilizing task-specific characteristics. These cross-task comparisons highlight that the benefits of input-adaptive parameter allocation are not uniform across tasks, but instead depend on the task’s inherent complexity distribution and error passing dynamics. The more sequential the reasoning required by a task, the greater the advantage offered by Flexi-LoRA’s dynamic rank allocation.

#### 5.5 Case Study

We present a qualitative analysis of model outputs to understand the influence of different fine-tuning approaches on mathematical reasoning tasks. Table 4 showcases two typical cases that show key differences between our proposed methods and existing baselines. (1) **Training-Inference Consistency Influence:** Table 4 reveals DyLoRA’s outputs mirror the Original model, indicating finetuning failure due to training-inference inconsistency (random rank sampling during training versus fixed rank at inference). On the other hand, DyLoRA+ successfully fine-tunes by maintaining consistent

<b>Question:</b> Cynthia eats one serving of ice cream every night. She buys cartons of ice cream with 15 servings of ice cream per carton at a cost of \$4.00 per carton. After 60 days, how much will she spend on ice cream? <b>Standard Answer:</b> 16		
<b>Original:</b> To find the total cost of ice cream, multiply the number of servings per carton by the number of cartons Cynthia buys. 15 servings per carton * 60 cartons = 900 servings. 900 servings * \$4.00 per serving = \$3,600. The answer is: \$3,600.	<b>DyLoRA (Baseline):</b> To find the total cost of ice cream, multiply the number of servings per carton by the number of cartons Cynthia buys. 15 servings per carton * 60 cartons = 900 servings. 900 servings * \$4.00 per serving = \$3600. The answer is: \$3600.	<b>DyLoRA+ (Ours):</b> Cynthia eats one serving of ice cream every night, so in 60 days, she will eat 60 servings of ice cream. Each carton of ice cream has 15 servings, so she will buy $60/15 = 4$ cartons of ice cream. Each carton costs \$4.00, so she will spend $4 * \$4.00 = \$16.00$ on ice cream. #### 16 The answer is: 16
<b>Question:</b> Meredith is a freelance blogger who writes about health topics and submits to clients each day as her permanent job. A blog article takes an average of 4 hours to research and write about. Last week, she wrote 5 articles on Monday and 2/5 times more articles on Tuesday than on Monday. On Wednesday, she wrote twice the number of articles she wrote on Tuesday. Calculate the total number of hours she spent writing articles in the three days. <b>Standard Answer:</b> 104		
<b>Original:</b> Meredith wrote 5 articles on Monday, $5 + 2(5) = 15$ articles on Tuesday, and $2(15) = 30$ articles on Wednesday. The total number of hours she spent writing articles is $4(5) + 4(15) + 4(30) = 20 + 60 + 120 = 200$ hours.	<b>LoRA (Gold Standard):</b> Meredith wrote 5 articles on Monday, which took 4 hours to research and write, so she spent $5 * 4 = 20$ hours writing on Monday. On Tuesday, she wrote 2/5 times more articles than on Monday, so she wrote $5 + (2/5 * 5) = 5 + 2 = 7$ articles. On Wednesday, she wrote twice the number of articles she wrote on Tuesday, so she wrote $7 * 2 = 14$ articles. In total, Meredith spent $20 + 7 + 14 = 41$ hours writing articles in the three days. #### 41 The answer is: 41	<b>Flexi-LoRA (Ours):</b> Meredith wrote 5 articles on Monday, and 2/5 times more articles on Tuesday, so she wrote $5 + (2/5)*5 = 5 + 2 = 7$ articles on Tuesday. On Wednesday, she wrote twice the number of articles she wrote on Tuesday, so she wrote $2*7 = 14$ articles on Wednesday. In total, Meredith wrote $5 + 7 + 14 = 26$ articles in the three days. Each article takes an average of 4 hours to research and write about, so she spent $26*4 = 104$ hours writing articles in the three days. #### 104 The answer is: 104

Table 4: **Case study** comparing reasoning performance across different LoRA variants on mathematical problems for LLaMA-3.2-1B-Instruct. Examples demonstrate how training-inference LoRA dynamics consistency influences problem solving: DyLoRA+ and Flexi-LoRA correctly solve problems that confound baseline methods, with errors highlighted in red and correct steps in green. Flexi-LoRA’s input-adaptive parameter allocation enables more accurate calculations and multi-step reasoning compared to static LoRA.

random rank selection across both stages. This difference demonstrates that dynamic consistency between training and inference is important for effective parameter-efficient finetuning. (2) **Fine-tuning Quality Indicators:** Beyond correctness, fine-tuned models show qualitative improvements in reasoning. DyLoRA+ and Flexi-LoRA produce detailed step-by-step solutions and follow output formatting ("The answer is:"), while Original and DyLoRA generate abbreviated responses that ignore instruction requirements. These differences, detailed reasoning and instruction adherence, serve as reliable indicators of successful finetuning. (3) **Error Analysis:** The cases reveal three error types in mathematical reasoning: (a) conceptual errors (Case 1: treating \$4/carton as \$4/serving, leading to \$3,600 vs \$16), (b) arithmetic misunderstanding (Case 2: "2/5 times more" computed as "2 times more"), and (c) process errors (correct intermediate steps but incomplete final calculations). These errors go through multi-step reasoning, showing ini-

tial mistakes. (4) **Reasoning Quality Comparison:** Successful methods maintain accuracy throughout reasoning chains. Case 2 illustrates this: LoRA correctly understands proportional relationships but fails in final answer, while Flexi-LoRA completes all computational steps correctly. Early-stage errors in baseline approaches go through the entire solution.

## 5.6 Ablation Studies

Figure 4 presents a comprehensive analysis analyzing key variables that contribute to Flexi-LoRA’s effectiveness. We examine four aspects of our approach: (1) **Training-Inference Rank Consistency** is essential for performance, as evidenced by DyLoRA’s decrease when using fixed ranks at inference despite dynamic training. This confirms our hypothesis that maintaining consistent rank dynamics is important, especially for sequential reasoning tasks. (2) **Input-Adaptive vs. Random Selection** comparison between Flexi-LoRA and DyLoRA+

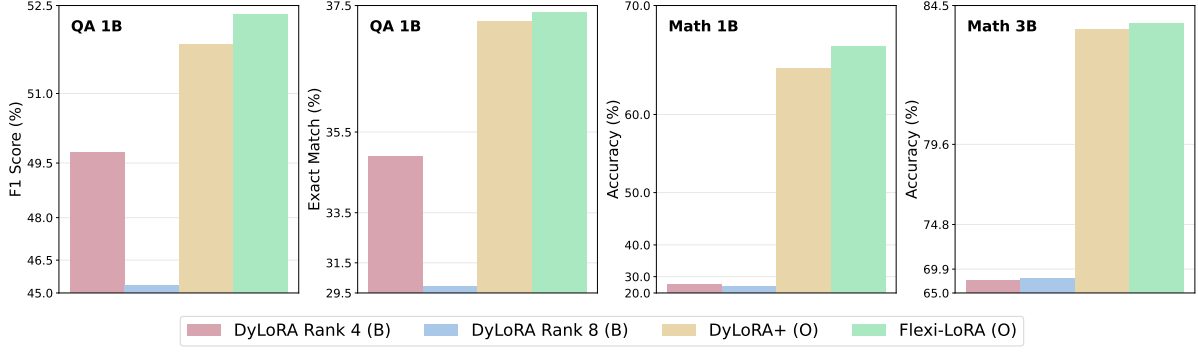


Figure 4: **Ablation study** comparing performance influence of rank selection strategies across model sizes. Charts show performance on QA tasks and mathematical reasoning tasks for LLaMA-3.2-1B-Instruct and 3B models. The methods are labeled as (B) for baseline methods and (O) for our proposed methods. All methods use dynamic ranks during training but differ in inference approach: DyLoRA uses fixed ranks (4 or 8), DyLoRA+ employs random rank selection, and Flexi-LoRA utilizes input-adaptive rank allocation. Results demonstrate that Flexi-LoRA maintains its advantage regardless of model size and task type, confirming the benefits of input-adaptive rank allocation and the consistency between training and inference rank dynamics.

demonstrates that learned complexity-aware allocation consistently outperforms random selection, validating the effectiveness of our router-based approach. (3) **Model Size Influences** show that while scaling from 1B to 3B parameters improves overall performance across all methods, Flexi-LoRA maintains its relative advantage, indicating that input-adaptive allocation remains beneficial regardless of model size. (4) **Task-Dependent Sensitivity** analysis reveals that mathematical reasoning tasks exhibit higher sensitivity to rank dynamics than QA tasks, with up to 39.44% performance gap between consistent and inconsistent methods on math compared to 6.49% on QA, illustrating how tasks with stricter evaluation criteria and error passing benefit more from adaptive parameter allocation.

## 5.7 Efficiency Analysis

Our analysis demonstrates Flexi-LoRA’s superior parameter efficiency across multiple aspects. Figure 3 and Tables 2 and 3 quantify this advantage: Flexi-LoRA (2,8) achieves the highest QA performance with only 504K trainable parameters, compared to 1703K for LoRA-8, a 70.40% parameter reduction while improving performance. Flexi-LoRA (1,8) further reduces parameter count to 304K (17.85% of LoRA-8) while maintaining competitive performance. This efficiency improvement extends to mathematical reasoning tasks, where Flexi-LoRA outperforms LoRA-8 on 1B models (66.56% vs. 63.17% accuracy) using only 533K parameters (31.29% of LoRA-8’s 1703K). Regarding computational overhead, Flexi-LoRA’s router

consists of only two layers that process the pooled input embedding once per sequence, introducing negligible additional computation compared to the base model computation. This minimal overhead is substantially balanced by the parameter efficiency gains, resulting in an overall more efficient finetuning framework that demonstrates consistent advantages across model sizes and tasks.

## 6 Conclusions

This paper introduces Flexi-LoRA, an input-adaptive framework that dynamically adjusts LoRA ranks based on question complexity. We demonstrate that maintaining consistent rank dynamics between training and inference is important for finetuning models, particularly for sequential reasoning tasks. Flexi-LoRA outperforms static LoRA while using fewer parameters (29.6% for QA and 31.3% for math reasoning), with performance gains more pronounced on mathematical tasks requiring reasoning chains. These results confirm that input-dependent parameter allocation enables efficient capacity allocation while reducing parameter redundancy, achieving benefits similar to mixture-of-experts frameworks through a more streamlined method.

Future work could have several directions: (1) layer-specific dynamic ranks to optimize parameter utilization at finer level; (2) router frameworks learning hierarchical aspects of input complexity; (3) integration with other parameter-efficient techniques such as sparse fine-tuning.



## Limitations

Although Flexi-LoRA exhibits encouraging results, several important considerations warrant attention. As with other fine-tuning methodologies, users should exercise caution regarding training data licensing and usage rights. Additionally, despite our documented strong performance across benchmark evaluations, developing deeper insights for varied applications would greatly benefit from expanded community participation and open-source collaborative initiatives.

## Ethics Statement

No ethical approval was required for this study. No ethical concerns are present.

## Availability Statement

The codes and models related to this paper are uploaded to the open-source community at <https://github.com/Anonymous/Flexi-LoRA> (After the paper is published, this link will be de-anonymized. For related code, please see the supplementary materials.).

## References

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168.
- Mehul Damani, Idan Shenfeld, Andi Peng, Andreea Bobu, and Jacob Andreas. 2025. [Learning how hard to think: Input-adaptive allocation of LM computation](#). In *The Thirteenth International Conference on Learning Representations*.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.
- Matthew Dunn, Levent Sagun, Mike Higgins, V. Ugur Guney, Volkan Cirik, and Kyunghyun Cho. 2017. [Searchqa: A new q&a dataset augmented with context from a search engine](#). *Preprint*, arXiv:1704.05179.
- Adam Fisch, Alon Talmor, Robin Jia, Minjoon Seo, Eunsol Choi, and Danqi Chen. 2019. MRQA 2019 shared task: Evaluating generalization in reading comprehension. In *Proceedings of 2nd Machine Reading for Reading Comprehension (MRQA) Workshop at EMNLP*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, ..., and Zhiyu Ma. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*.
- Shuyang Jiang, Yusheng Liao, Ya Zhang, Yanfeng Wang, and Yu Wang. 2025. [Fine-tuning with reserved majority for noise reduction](#). In *The Thirteenth International Conference on Learning Representations*.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi. 2017. [Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5376–5384.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. [MAWPS: A math word problem repository](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1152–1157, San Diego, California. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [RACE: Large-scale ReAding comprehension dataset from examinations](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.

638	Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. <a href="#">Zero-shot relation extraction via reading comprehension</a> . In <i>Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)</i> , pages 333–342, Vancouver, Canada. Association for Computational Linguistics.	694
639		695
640		696
641		
642		697
643		698
644	Rohin Manvi, Anikait Singh, and Stefano Ermon. 2024. <a href="#">Adaptive inference-time compute: Lms can predict if they can do better, even mid-generation</a> . <i>Preprint</i> , arXiv:2410.02725.	699
645		700
646		701
647		
648	Ioannis Partalas, Eric Gaussier, Axel-Cyrille Ngonga Ngomo, et al. 2013. Results of the first bioasq workshop. In <i>BioASQ@ CLEF</i> , pages 1–8.	702
649		703
650		704
651	Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. <a href="#">Are NLP models really able to solve simple math word problems?</a> In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2080–2094, Online. Association for Computational Linguistics.	705
652		706
653		707
654		708
655		709
656		
657		
658	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <a href="#">SQuAD: 100,000+ questions for machine comprehension of text</a> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	710
659		711
660		712
661		713
662		714
663		715
664	Subhro Roy and Dan Roth. 2015. <a href="#">Solving general arithmetic word problems</a> . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 1743–1752, Lisbon, Portugal. Association for Computational Linguistics.	716
665		717
666		718
667		719
668		720
669	Amrita Saha, Rahul Aralikatte, Mitesh M. Khapra, and Karthik Sankaranarayanan. 2018. <a href="#">DuoRC: Towards complex language understanding with paraphrased reading comprehension</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1683–1693, Melbourne, Australia. Association for Computational Linguistics.	721
670		722
671		723
672		724
673		
674		
675		
676		
677	Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2025. <a href="#">Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning</a> . In <i>The Thirteenth International Conference on Learning Representations</i> .	725
678		726
679		727
680		728
681		729
682	Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordani, Philip Bachman, and Kaheer Suleman. 2017. <a href="#">NewsQA: A machine comprehension dataset</a> . In <i>Proceedings of the 2nd Workshop on Representation Learning for NLP</i> , pages 191–200, Vancouver, Canada. Association for Computational Linguistics.	
683		
684		
685		
686		
687		
688		
689	Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobzyev, and Ali Ghodsi. 2023. <a href="#">DyLoRA: Parameter-efficient tuning of pre-trained models using dynamic search-free low-rank adaptation</a> . In <i>Proceedings of the 17th Conference of the European</i>	
690		
691		
692		
693		

## A Appendix

### A.1 Related Work for Input-adaptive Methods

**Input-adaptive methods.** Recent work has explored various approaches to dynamically allocate computational resources based on input complexity. Wu et al. (2024) and Snell et al. (2025) investigate compute-optimal scaling strategies that adapt test-time computation allocation per prompt, showing that different inference strategies yield optimal performance depending on problem difficulty. Manvi et al. (2024) propose using self-evaluations to adaptively determine necessary inference computation, while Damani et al. (2025) explore difficulty estimation for optimal compute allocation. Zhang et al. (2025) introduce OSCA, an algorithm that optimizes sample compute allocation across different inference configurations. These approaches primarily focus on inference-time adaptation rather than model framework adaptation. Different from these methods which modify inference strategies, our work centers on input-adaptive framework modifications during both training and inference by dynamically adjusting LoRA ranks based on input complexity. This provides a more streamlined approach to input-adaptive computation than methods requiring test-time search or complex verifier frameworks.

### A.2 Implementation Details

Implementation details, including inference parameters in Table 5, training hyperparameters in Table 6, and method-specific configurations in Tables 7, are provided.

	QA	Math
do_sample	False	False
early_stopping	True	True
length_penalty	1.0	1.0
max_new_tokens	20	768
num_beams	1	1
pad_token_id	pad_token_id	pad_token_id
temperature	1.0	1.0
top_p	1.0	1.0

Table 5: **Inference parameters** for QA and mathematical tasks.

Prompt for Math datasets:

```
<|start_header_id|>user<|end_header_id|>\n\nSolve the question and your response should end with \"The answer is: [answer]\".\n\nQuestion: {}<|eot_id|><|start_header_id|>assistant<|end_header_id|>\n\nAnswer:
```

### A.3 Environment

```
datasets==2.18.0
deepspeed==0.15.3
huggingface_hub==0.24.2
numpy==1.23.5
python==3.11.5
torch==2.3.1+cu118
tqdm==4.66.4
transformers==4.46.0
wandb==0.14.2
```

This paper was refined with ChatGPT and Claude.

Prompt for MRQA:

```
<|start_header_id|>user<|end_header_id|>\n\nExtract the exact text span from the given context that directly answers the question, without modifying or combining multiple parts of the text.\n\nContext: {}\n\nQuestion: {}<|eot_id|><|start_header_id|>assistant<|end_header_id|>\n\nAnswer:
```

	QA	Math
Train steps	{1000} for LLaMA-3.2-1B-Instruct	{200, 400, 600, 800, 1000} for LLaMA-3.2-1B-Instruct {200, 800} for LLaMA-3.2-3B-Instruct
Optimizer	AdamW	AdamW
Max length	512	768
warmup_steps	100	100
learning_rate	2e-5	2e-5
per_device_train_batch_size	16	16
lr_scheduler_type	constant_with_warmup	constant_with_warmup
gradient_accumulation_steps	2	2

Table 6: **Training hyperparameters** for QA and mathematical tasks. { } shows parameter search.

Method	Details
LoRA	r=4/8; lora_alpha=16; target_modules=["q_proj", "k_proj", "v_proj", "o_proj"]; lora_dropout=0
AdaLoRA	r=16/8; target_r=8/4; lora_alpha=16; target_modules=["q_proj", "k_proj", "v_proj", "o_proj"]; lora_dropout=0; stage=1/2
DyLoRA	rank_choices=[1, 2, 3, 4, 5, 6, 7, 8]; lora_alpha=16; target_modules=["q_proj", "k_proj", "v_proj", "o_proj"]; lora_dropout=0
Flexi-LoRA	max_r=8; min_r=1; rank_mapping=0: 2, 1: 8; lora_alpha=16; target_modules=["q_proj", "k_proj", "v_proj", "o_proj"]; lora_dropout=0

Table 7: **Method configurations** for various PEFT methods.