

# 🍌 MAC-Tuning: LLM Multi-Compositional Problem Reasoning with Enhanced Knowledge Boundary Awareness

Junsheng Huang<sup>1,2</sup>, Zhitao He<sup>1</sup>, Yuchen Huang<sup>1</sup>, Sandeep Polisetty<sup>3</sup>  
Qingyun Wang<sup>4</sup> Yi R. (May) Fung<sup>1\*</sup>

<sup>1</sup>Hong Kong University of Science and Technology

<sup>2</sup>University of Illinois Urbana-Champaign  
jh103@illinois.edu

<sup>3</sup>UMass Amherst  
yrfung@ust.hk

<sup>4</sup>William & Mary

## Abstract

The hallucination of non-existent facts by LLMs is an important problem given its widespread adoption across various applications. Previous research addresses this problem by analyzing the internal parameterized knowledge boundaries to estimate confidence. However, these studies focus on the single-problem setting and have not explored the more challenging multi-problem setting, which requires accurately answering multiple questions simultaneously. We introduce a novel method for the multi-problem setting, **Multiple Answers and Confidence Stepwise Tuning (MAC-Tuning)**, that separates the learning of answer prediction and confidence estimation during fine-tuning on instruction data. Extensive experiments demonstrate that our method outperforms baselines by up to 25% in average precision.<sup>1</sup>

## 1 Introduction

Large language models (LLMs) are widely used in knowledge-intensive scenarios, such as question answering (Gu et al., 2023), information retrieval (Ren et al., 2023), and recommendation systems (Liu et al., 2023). Yet, they often produce non-existing facts when faced with questions outside their parametric knowledge, which undermines their reliability (Maynez et al., 2020; He et al., 2025c). Many efforts have been dedicated to mitigating LLM hallucination, such as leveraging knowledge boundaries to constrain the reasoning scope of LLMs to help them better distinguish between reliable and unreliable information (Chen et al., 2024; Liang et al., 2024a; Zhang et al., 2024; Jin et al., 2024). Notably, these work mainly focus on the **single-problem setting**, where users repeatedly input questions and context for models to answer one by one.

\*Corresponding author.

<sup>1</sup>We release our code and resource at [MAC-Tuning](#).

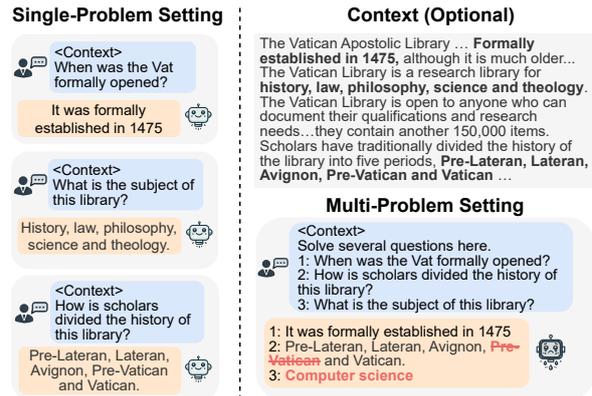


Figure 1: An illustration of the multi-problem setting. *Red* indicates that the LLM’s output is inaccurate.

LLM hallucination in the **multi-problem setting** — in which a single input contains multiple distinct sub-questions with optional context for the model to extract and address — remains relatively underexplored. As seen in Figure 1, this is a fundamentally challenging setting because the model must distinguish each sub-question, reason over different knowledge, and synthesize results cohesively. Undesirable overshadowing of context from one sub-question with another, and propagation of reasoning confusion, may compromise the reliability of LLMs in multi-problem answering (Cheng et al., 2023a, Wang et al., 2024, Son et al., 2024; Li et al., 2024; He et al., 2025b). As LLM-based multi-problem reasoning becomes increasingly widespread due to its efficiency benefits in scenarios involving extensive shared contexts (e.g., task instructions, exemplars), reduced model access, and lower API costs, enhancing model confidence estimation calibration for this emerging class of reasoning demands growing attention and effort as well.

In this paper, we investigate the hallucinations in LLMs within the multi-problem setting and propose leveraging the knowledge boundary to simultaneously handle the composition of multiple

problems. Inspired by Zhang et al. (2024), which advocates for encouraging the LLM to express confidence to reduce hallucinations, we introduce **Multiple Answers and Confidence Stepwise Tuning (MAC-Tuning)** under the multi-problem setting. Our approach involves several key steps. First, we identify the knowledge boundary between parametric knowledge and the multi-problem dataset to extract uncertain questions. Next, we automatically label the model’s confidence for both certain and uncertain data. These labeled data are then used to create multiple question-answer data and multiple QA-Confidence data so we can train the original model by separating the learning process of ground-truth answers and confidence, which enhances performance and reliability.

Our contributions can be summarized as follows:

- We are the first to explore LLM confidence estimation under the more challenging multi-problem setting, where LLMs must handle multiple problems simultaneously.
- We propose MAC-Tuning, which separates the learning process of answer and confidence predictions for enhancing knowledge boundary awareness and reducing hallucination.
- Through extensive experiments with different base models of varying sizes and various datasets, MAC-Tuning achieves an AP score gain of up to 25% over baselines in LLM multi-problem reasoning. Finally, we share our insights discovered to motivate future work.

## 2 Methodology

Figure 2 shows the data construction process for **Multiple Answers and Confidence Stepwise Tuning (MAC-Tuning)**.

### 2.1 Multi-Problem Tuning Data Construction

First, we combine  $n$  single problems from original datasets to construct our initial Multi-Problem dataset. We utilize this to compare LLMs’ outputs with ground-truth answers, for distinguishing the knowledge boundary between LLM parameters and instruction data. Specifically, for each individual problem in the multi-problem pair, we assign: “*I am sure*” if the output aligns with ground-truth answer; “*I am unsure*” otherwise (e.g., Step 2 in Figure 2). With the assigned confidence labels, we construct Multi-Problem Tuning data as follows: **Multiple QA pair**  $D_{MultQA}$ : We directly combine the questions and answers together, with *Question*

$q_i$  as input and *Answer*  $a_i$  as output label, to form  $D_{MultQA} = [(q_1, a_1) \dots (q_i, a_i) \dots (q_n, a_n)]$ .

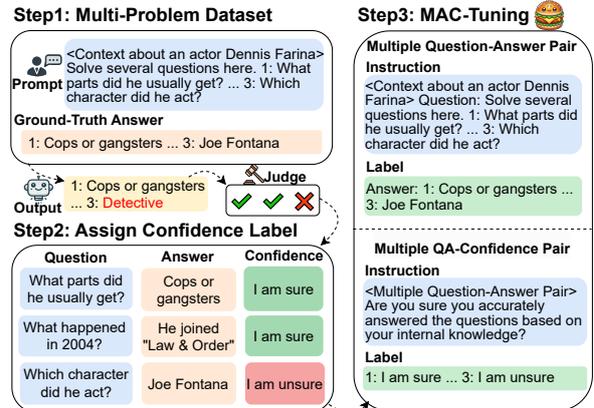


Figure 2: We first construct the Multi-Problem dataset, and then use it to generate Multi-Problem Tuning data.

**Multiple QA-Confidence pair**  $D_{MultQA,C}$ : The input consists of an instruction for the LLM to express its confidence (*i.e.*, certainty in correctness) for a given question-answer pair, while the output is the confidence level in linguistic form<sup>2</sup>.

### 2.2 Training and Inference

Using the Multi-Problem Tuning data, we conduct a two-step supervised fine-tuning process to train the model to answer questions and express confidence in a multi-problem setting. The objective for the first step, in answering question, is:

$$\max_{\Theta_0} \sum_{(Q,A) \in D_{MultQA}} \log P(A|Q; \Theta_0) \quad (1)$$

The objective for the second step, in expressing confidence, is:

$$\max_{\Theta_1} \sum_{(Q,A,C) \in D_{MultQA,C}} \log P(C|Q, A; \Theta_1) \quad (2)$$

where  $Q$ ,  $A$ , and  $C$  represent the sets of multiple questions, multiple answers, and multiple confidence levels, respectively.  $\Theta_0$  and  $\Theta_1$  represent the parameters of the base model and the model after the first step of fine-tuning, respectively.

## 3 Experiment

### 3.1 Dataset

We validate the effectiveness of our method across different problem settings and datasets: for the *Independent* setting, where the questions are not related to each other, we use the **CoQA** (Reddy et al., 2019), **GSM** (Cobbe et al., 2021), **MMLU**

<sup>2</sup>The template is in Appendix A.3

Model	Independent								Sequential			
	CoQA		ParaRel		GSM		MMLU		MTI-Bench		SQA	
	AP	ECE										
LLaMA3	54.6	22.6	45.1	40.8	79.3	52.8	50.3	43.8	37.4	17.7	44.9	35.4
QA-Only	66.3	15.1	53.7	12.6	75.3	36.1	58.5	17.9	45.0	16.9	56.6	21.0
Single-QA	65.5	28.9	73.5	10.7	56.6	44.5	58.3	25.7	N/A	N/A	N/A	N/A
Merge-AC	67.4	17.0	73.0	65.3	75.1	44.8	58.5	18.3	38.3	33.7	49.2	31.7
MAC-Tuning	<b>69.8</b>	<b>7.33</b>	<b>76.1</b>	<b>3.61</b>	<b>79.9</b>	<b>3.16</b>	<b>63.1</b>	<b>12.5</b>	<b>64.0</b>	<b>13.4</b>	<b>65.0</b>	<b>14.6</b>

Table 1: This is the confidence calibration result (%). We use one-shot CoT for GSM results. **Bold** font highlights the best performance for the dataset across different methods. We don’t apply Single-QA to the *Sequential* setting dataset, as doing so would disrupt the logical connections among the questions.

(Hendrycks et al., 2021), and **ParaRel** (Elazar et al., 2021) datasets; for the *Sequential* setting, where the questions are logically related to each other, we use the **MTI-Bench** (Son et al., 2024) and **SQA** (Iyyer et al., 2017) datasets. These datasets are either Question Answer (QA) or Multiple Choice (MC) formats. Table 2 shows the details of the dataset. Further information on the distribution of certain and uncertain data among the training set across different datasets is detailed in Appendix A.4.

	Independent				Sequential	
	CoQA	ParaRel	GSM	MMLU	MTI-Bench	SQA
<b>Train</b>	5006	7500	7468	2448	2400	3985
<b>Test</b>	5011	5584	1319	2439	600	925
<b>Type</b>	QA	QA	QA	MC	QA	QA

Table 2: Statistics of the datasets.

### 3.2 Evaluation Metrics

We directly compare the LLM generation to the ground-truth answer for the Question-Answer format. For Multiple-Choice format, we check the choice (A, B, C, D) and the option in the LLM generation. Across both types of answer generation tasks, we consider three evaluation metrics: (1) **Average Precision (AP)**: We use AP to measure the precision in identifying and ranking relevant predictions. A higher AP score means the model has high certainty about correct answers and high uncertainty about wrong answers. (2) **Expected Calibrated Error (ECE)**: We use ECE to measure how closely the predicted certainty reflects the true certainty of LLM (Chen et al., 2023). Low ECE indicates better-calibrated predictions. (3) **Accuracy**: We compute accuracy as the fraction of correct responses amongst questions in which LLMs expressed certainty towards their answers.

### 3.3 Baselines

We compare MAC-Tuning with the base model and its variants in the multi-problem settings. We use LLaMA3-8B-Instruct (**LLaMA3**) (Dubey et al., 2024) as the backbone. For baseline **QA-Only**, we fine-tune the base model directly using the Multiple Question-Answer pairs to evaluate the effectiveness of the traditional instruction tuning method under the multi-problem setting. For baseline **Single-QA**, we use single-problem data to fine-tune and directly apply it to the multi-problem setting. For baseline **Merge-AC**, instead of separating the learning process of ground-truth answers and confidence, we directly let the model learn multiple answers along with their corresponding confidence levels<sup>3</sup>.

### 3.4 Overall Performance

In Table 1, we report the results on multi-problem setting from three single questions combined together. MAC-Tuning achieves the best AP score across all datasets, showing up to a 15% improvement, along with a lower ECE. This suggests that after MAC-Tuning, the model becomes more adept at distinguishing between certain and uncertain questions, delivering more reliable results through improved confidence estimation in answer prediction. We also evaluate each model’s accuracy on every dataset. MAC-Tuning consistently outperforms the base model in accuracy by up to 45.8% and, on average, 23.7%. The reason is that we separate the tasks of learning correct answers and confidence within a multi-problem setting. After learning the ground-truth answer, the LLM can better understand confidence, while still retaining its ability to extract information, respond accurately, and address multiple problems simultaneously.

<sup>3</sup>Baseline examples are in Appendix A.8. Implementation details are in Appendix A.6.

**Ablation on Different Component** We further test three variants of the MAC-Tuning method in the multi-problem setting: **QA-Only**, which is MAC-Tuning without the confidence component; **Single-QA**, where we evaluate MAC-Tuning with single problem data; and **Merge-AC**, where we evaluate MAC-Tuning without separating the learning process of ground-truth answers and confidence. As seen from the results in Table 1, **MAC-Tuning** has up to 25% and, on average, 11% AP improvement compared with **Merge-AC**, reflecting that separating the learning process of ground-truth answers and confidence is crucial in multi-problem setting, as LLM cannot learn both in one time. The performance of **Single-QA** is better than the base model but worse than **QA-Only** in most cases, showing that LLM can aware the knowledge boundary under single-problem setting and transfer it to multi-problem setting, but it is not sufficient for LLM to answer multiple problems simultaneously.

### 3.5 Investigation on Out of Domain Settings

We perform MAC-Tuning on base model with *Sequential* setting dataset SQA and test it on other datasets, with the results as presented in Table 3. Even on out-of-domain datasets, MAC-Tuning still outperforms the base model, showing that it can effectively learn the multi-problem setting and generalize across different domains.

Metric	CoQA	Pararel	MMLU	MTI-Bench
Accuracy	59.3	70.3	52.6	57.8
AP score	62.2	58.7	53.8	81.7
ECE	10.4	9.64	8.95	16.1

Table 3: The result (%) for MAC-Tuning on SQA dataset and test on other datasets.

### 3.6 Analysis on Various Number of Questions

We explore different numbers of questions in the multi-problem setting to investigate how this varies the accuracy. We only do this for three *Independent* setting datasets, and the results are reported in Figure 3. MAC-Tuning consistently outperforms the base model in accuracy by at least 10.0% and, on average, 26.8%. For easy tasks like ParaRel, the ability of the base model to handle multiple problems simultaneously is even higher when compared with the traditional single-problem setting, indicating that LLM could leverage in-context learning and focus on relevant knowledge better under the multi-problem setting. However, for other datasets like MMLU, MAC-Tuning performs slightly worse

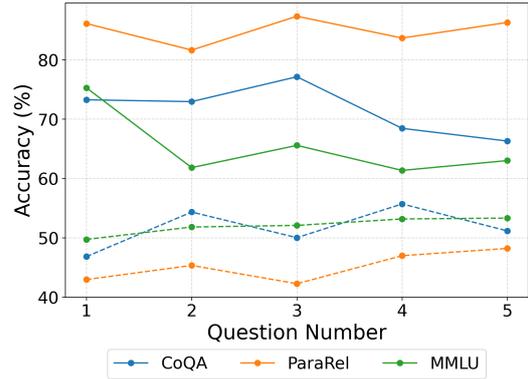


Figure 3: Accuracy for combining different number ( $n$ ) of single problem together. Solid lines represent MAC-Tuning, while dashed lines represent LLaMA3.

as the question number increases. A reasonable explanation is that it is out of the base model’s ability to learn too many hard tasks together but within effective scope to learn several easy tasks at the same time. The result for extremely large numbers of questions is in Appendix A.9.

For future work, we believe it is meaningful to further explore whether mixing questions of varying difficulty and diversity in multi-problem settings leads to better scaling behavior (Qin et al., 2025). This direction may help uncover strategies for enhanced model generalization.

### 3.7 Cross Task Transfer Study

We fine-tune the model with question number  $n = 3$ , and subsequently evaluate its performance on both single problem inputs ( $n = 1$ ) and more complex instances involving a higher number of questions ( $n = 5$ ). The evaluation on  $n = 1$  aims to examine whether the model retains its ability to accurately solve individual problems after being exposed to multi-problem setting training. Conversely, the evaluation on  $n = 5$  serves to assess the generation capability of MAC-Tuning when scaling to larger compositions beyond the training scope. This allows us to understand both the robustness and scalability of our proposed method across different levels of compositional complexity. The results are reported in Table 4.

From the result of single question inputs, we observe that accuracy increases on easier dataset (e.g. CoQA) but decreases on more challenging dataset like GSM, when compared to models fine-tuned specifically under that setting. This pattern indicates that LLM acquires underlying knowledge during MAC-Tuning rather than merely memorizing

Question Number	CoQA	ParaRel	GSM	MMLU
n = 1	78.8	84.2	71.1	54.6
n = 5	79.1	86.2	67.7	63.7

Table 4: Accuracy (%) for MAC-Tuning with question number  $n = 3$  transferring to question number  $n = 1$  and  $n = 5$ . We use one-shot CoT for GSM results.

the patterns for multi-problem setting. In contrast, when fine-tuned on five-question inputs ( $n = 5$ ), the model’s performance is comparable or even exceeds that of the baseline fine-tuned directly on  $n = 5$ . These findings strengthen the statement we make in Section 3.6: while LLMs can efficiently learn multiple easy tasks, they exhibit difficulty when faced with several difficult tasks simultaneously.

### 3.8 Analysis on Different Base Model

Table 5 shows the result from changing the base model to Qwen2-7B-Instruct (Yang et al., 2024). We observe that the performance trends remain consistent even with a different base model. MAC-Tuning continues to demonstrate an average precision (AP) gain of up to near 24% with a lower ECE, showcasing the effectiveness of learning ground-truth answers and confidence separately.

Approach	Independent				Sequential			
	ParaRel		MMLU		MTI-Bench		SQA	
	AP	ECE	AP	ECE	AP	ECE	AP	ECE
Vanilla	54.3	37.8	68.1	25.3	48.8	31.3	30.3	54.6
MAC-Tuning	<b>78.7</b>	<b>9.59</b>	<b>73.0</b>	<b>17.1</b>	<b>53.3</b>	<b>18.6</b>	<b>47.7</b>	<b>29.2</b>

Table 5: Confidence calibration result (%) for Qwen2-7B-Instruct, with **bold** denoting the top performance.

### 3.9 Analysis on Different Model Size

We compare base models of different sizes to study how model size affects performance and the confidence calibration results for Llama-3.2-3B (Dubey et al., 2024) and Phi-3.5-mini-Instruct (Abdin et al., 2024) is shown in Table 6 and Table 7 respectively. Despite using different base models of varying sizes, the results indicate consistent performance patterns. For smaller models, the accuracy improvement after **MAC-Tuning** is more evident, indicating enhanced ability to differentiate between certain and uncertain questions.

### 3.10 Human Evaluation

We randomly selected and evaluated 100 examples from the ParaRel dataset. The human annotator was shown only the query and the ground-truth answer, and asked to assess the factual correctness of each

Approach	Independent				Sequential			
	ParaRel		CoQA		MTI-Bench		SQA	
	AP	ECE	AP	ECE	AP	ECE	AP	ECE
Vanilla	30.7	70.3	46.0	45.0	34.3	70.3	35.6	45.0
MAC-Tuning	<b>55.5</b>	<b>33.5</b>	<b>62.4</b>	<b>33.5</b>	<b>35.5</b>	<b>33.5</b>	<b>44.3</b>	<b>33.5</b>

Table 6: Confidence calibration result (%) for Llama-3.2-3B, with **bold** denoting the top performance across different methods.

Approach	Independent				Sequential			
	ParaRel		CoQA		MTI-Bench		SQA	
	AP	ECE	AP	ECE	AP	ECE	AP	ECE
Vanilla	58.0	22.8	56.0	32.9	21.4	29.1	96.6	33.7
MAC-Tuning	<b>70.2</b>	<b>14.2</b>	<b>68.2</b>	<b>29.0</b>	<b>68.7</b>	<b>22.6</b>	<b>52.3</b>	<b>23.9</b>

Table 7: Confidence calibration result (%) for Phi-3.5-mini-Instruct, with **bold** denoting the top performance across different methods.

model’s output without knowing the confidence label (“I am sure” or “I am unsure”). We then compared the accuracy rates between responses that the model labeled as “I am sure” versus those labeled as “I am unsure.” For answers the model labeled as “I am sure”: Human evaluation confirmed a factual accuracy of **89.2%**. For answers the model labeled as “I am unsure”: The human-verified factual accuracy was only **41.2%**.

These results provide strong empirical evidence for our central claim. The substantial 48-percentage-point gap demonstrates that the confidence learned by MAC-Tuning is not merely an artifact of automatic metrics. Instead, it reflects a genuine, human-perceptible distinction in answer quality. This strong alignment with human judgment validates the reliability and real-world applicability of our method.

## 4 Conclusion

In this paper, we introduce a novel method, MAC-Tuning, to enhance large language model (LLM) confidence calibration and reasoning robustness in the challenging yet underexplored multi-problem scenario. Our proposed approach automatically constructs multi-problem setting question-answer pairs with confidence annotations for identifying the intrinsic knowledge gap between parametric knowledge and instructional data. With this data constructed, we guide the LLM to better reason on answer prediction and confidence estimation separately, in multi-problem setting. Extensive experiments across different datasets show that our method significantly improves performance in areas where the original LLM struggles.

## Limitation

While our work provides valuable insight on the new Multiple Question setting and introduces an innovative fine-tuning method, there are several limitations to acknowledge. First, although we experimented with various prompts, as is typical in prompt-based LLM studies, we cannot ensure that slight changes in prompts would not significantly alter the results. Second, due to constraints of cost, time, and computational resources, we selected a subset of experiments that we believe to be informative and representative. However, additional experiments across a wider range of datasets and LLMs might provide further insights. Lastly, in this new setting, there may be other underlying reasons for the experimental results. Future work will aim to address these limitations by expanding datasets and conducting new experiments to explore other potential factors affecting performance.

## Acknowledgment

This research was supported in part by WeBank (Grant WEB24EG01-L).

## References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone.](#)
- Lida Chen, Zujie Liang, Xintao Wang, Jiaqing Liang, Yanghua Xiao, Feng Wei, Jinglei Chen, Zhenghong Hao, Bing Han, and Wei Wang. 2024. [Teaching large language models to express knowledge boundary from their own signals.](#)
- Yangyi Chen, Lifan Yuan, Ganqu Cui, Zhiyuan Liu, and Heng Ji. 2023. [A close look into the calibration of pre-trained language models.](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1343–1367, Toronto, Canada. Association for Computational Linguistics.
- Zhoujun Cheng, Jungo Kasai, and Tao Yu. 2023a. [Batch prompting: Efficient inference with large language model APIs.](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 792–810, Singapore. Association for Computational Linguistics.
- Zhoujun Cheng, Jungo Kasai, and Tao Yu. 2023b. [Batch prompting: Efficient inference with large language model apis.](#)
- 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. *arXiv preprint arXiv:2110.14168*.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. 2023. [Improving factuality and reasoning in language models through multiagent debate.](#)
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic,

Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Milalon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papanikos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baeovski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi,

Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, DingKang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damraj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khadelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang

- Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. [The llama 3 herd of models](#).
- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhishava Ravichander, Ed Hovy, Hinrich Schutze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. *ArXiv*, abs/2102.01017.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. [Retrieval-augmented generation for large language models: A survey](#).
- Yu Gu, Xiang Deng, and Yu Su. 2023. [Don't generate, discriminate: A proposal for grounding language models to real-world environments](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4928–4949, Toronto, Canada. Association for Computational Linguistics.
- Zhitao He, Pengfei Cao, Yubo Chen, Kang Liu, Ruopeng Li, Mengshu Sun, and Jun Zhao. 2023. [LEGO: A multi-agent collaborative framework with role-playing and iterative feedback for causality explanation generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9142–9163, Singapore. Association for Computational Linguistics.
- Zhitao He, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. 2024. [Agentscourt: Building judicial decision-making agents with court debate simulation and legal knowledge augmentation](#).
- Zhitao He, Zijun Liu, Peng Li, Yi R Fung, Ming Yan, Ji Zhang, Fei Huang, and Yang Liu. 2025a. [Advancing language multi-agent learning with credit re-assignment for interactive environment generalization](#).
- Zhitao He, Zongwei Lyu, Dazhong Chen, Dadi Guo, and Yi R. Fung. 2025b. [Matp-bench: Can mllm be a good automated theorem prover for multimodal problems?](#)
- Zhitao He, Sandeep Polisetty, Zhiyuan Fan, Yuchen Huang, Shujin Wu, and Yi R. Fung. 2025c. [Mm-boundary: Advancing mllm knowledge boundary awareness through reasoning step confidence calibration](#).
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#).
- Shengding Hu, Yifan Luo, Huadong Wang, Xingyi Cheng, Zhiyuan Liu, and Maosong Sun. 2023. [Won't get fooled again: Answering questions with false premises](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5626–5643, Toronto, Canada. Association for Computational Linguistics.
- Yuchen Huang, Zhiyuan Fan, Zhitao He, Sandeep Polisetty, Wenyan Li, and Yi R. Fung. 2025. [Cultureclip: Empowering clip with cultural awareness through synthetic images and contextualized captions](#).
- Mohit Iyyer, Wen-tau Yih, and Ming-Wei Chang. 2017. [Search-based neural structured learning for sequential question answering](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1821–1831, Vancouver, Canada. Association for Computational Linguistics.
- Zhuoran Jin, Pengfei Cao, Chenhao Wang, Zhitao He, Hongbang Yuan, Jiachun Li, Yubo Chen, Kang Liu, and Jun Zhao. 2024. [Rwku: Benchmarking real-world knowledge unlearning for large language models](#).
- Ming Li, Pei Chen, Chenguang Wang, Hongyu Zhao, Yijun Liang, Yupeng Hou, Fuxiao Liu, and Tianyi Zhou. 2024. [Mosaic-it: Free compositional data augmentation improves instruction tuning](#).
- Qiuyu Liang, Weihua Wang, Feilong Bao, and Guanglai Gao. 2024a. [L<sup>2</sup>GC:lorentzian linear graph convolutional networks for node classification](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9988–9998, Torino, Italia. ELRA and ICCL.
- Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiaxing Zhang. 2024b. [Learning to trust your feelings: Leveraging self-awareness in LLMs for hallucination mitigation](#). In *Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP*, pages 44–58, Bangkok, Thailand. Association for Computational Linguistics.

- Jianzhe Lin, Maurice Diesendruck, Liang Du, and Robin Abraham. 2024. [Batchprompt: Accomplish more with less](#).
- Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. 2023. [A first look at llm-powered generative news recommendation](#). *ArXiv*, abs/2305.06566.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. [Peft: State-of-the-art parameter-efficient fine-tuning methods](#). <https://github.com/huggingface/peft>.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. [On faithfulness and factuality in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. 2023. [Check your facts and try again: Improving large language models with external knowledge and automated feedback](#). *arXiv preprint arXiv:2302.12813*.
- Zeyu Qin, Qingxiu Dong, Xingxing Zhang, Li Dong, Xiaolong Huang, Ziyi Yang, Mahmoud Khademi, Dongdong Zhang, Hany Hassan Awadalla, Yi R. Fung, Weizhu Chen, Minhao Cheng, and Furu Wei. 2025. [Scaling laws of synthetic data for language models](#).
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. [CoQA: A conversational question answering challenge](#). *Transactions of the Association for Computational Linguistics*, 7:249–266.
- Ruiyang Ren, Wayne Xin Zhao, Jing Liu, Hua Wu, Ji-Rong Wen, and Haifeng Wang. 2023. [Tome: A two-stage approach for model-based retrieval](#).
- Guijin Son, Sangwon Baek, Sangdae Nam, Ilgyun Jeong, and Seungone Kim. 2024. [Multi-task inference: Can large language models follow multiple instructions at once?](#)
- Zhiqing Sun, Yikang Shen, Qinzhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2023. [Principle-driven self-alignment of language models from scratch with minimal human supervision](#).
- Zhengxiang Wang, Jordan Kodner, and Owen Rambow. 2024. [Exploring the zero-shot capabilities of llms handling multiple problems at once](#).
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. [Qwen2 technical report](#).
- Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024. [R-tuning: Instructing large language models to say ‘I don’t know’](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7113–7139, Mexico City, Mexico. Association for Computational Linguistics.

## A Appendix

### A.1 Full Case for Examples of Introduction

Full case for the examples in introduction can be found in Figure 4.

### A.2 Related Work

**Hallucination:** Large language models (LLMs) are widely used in knowledge-intensive scenarios, such as question answering (Gu et al., 2023; He et al., 2025c), information retrieval (Ren et al., 2023; Huang et al., 2025) and recommendation systems (Liu et al., 2023). However, LLMs have tendency to generate non-existing facts when faced with questions that are out of their parametric knowledge (Maynez et al., 2020). Many efforts are dedicated to mitigating hallucinations in LLMs, such as retrieval-augmented generation (Gao et al., 2024; Peng et al., 2023), multi-agent debate (He et al., 2023; Du et al., 2023; He et al., 2024; Sun et al., 2023; He et al., 2025a), and model confidence calibration (Zhang et al., 2024; Hu et al., 2023; He et al., 2025c).

**Knowledge Boundary:** There are many different ways to utilize knowledge boundary to reduce LLM hallucination. Liang et al. (2024b)’s work uses merged knowledge probing and consistency checking methods to help LLM express their internal knowledge. Chen et al. (2024)’s work leverages LLM internal signals to let LLM know their unknowns. Zhang et al. (2024) utilize knowledge boundary to instruct LLM say "I don’t know". It is a popular way to use confidence to express knowledge boundary of LLMs and we also follow this.

**Multiple Problem Setting:** Current LLM research has predominantly focused on single problem set-

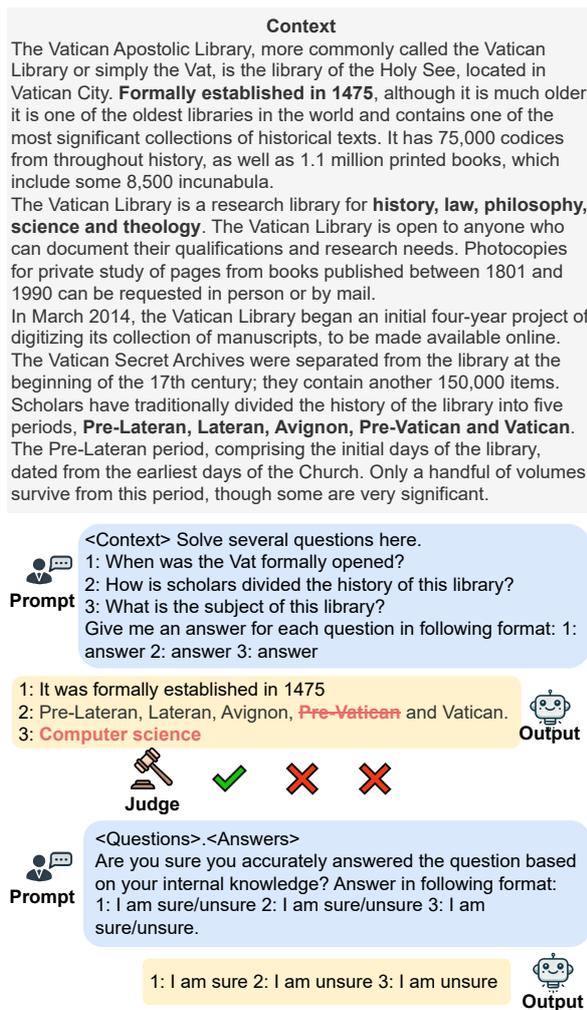


Figure 4: The full case of examples in introduction in Multiple Problem setting. Red context indicates that LLM’s output is inaccurate. The second answer lacks the information of "Pre-Vatican" and the third answer contains a completely factual error. After MAC-Tuning, LLM show uncertainty towards answering this two previously incorrect questions.

ting. There are only a few works focusing on this new setting. Cheng et al. (2023a) propose batch prompting that prompts LLMs with single independent problems batched together following few-shot exemplars together. Son et al. (2024) goes further by researching sequential datasets and develops the first multi-task benchmark (MTI-Bench). Wang et al. (2024) pays attention to zero-shot cases of multi-problem setting and design a new benchmark ZEMPEB. Li et al. (2024) analyze different strategy under independent setting, where single questions are combined into various constraint formats without sharing context between them. Despite these efforts, the multi-problem setting presents significant challenges. For example, Wang et al.

(2024) shows that in zero-shot setting, LLMs consistently perform worse when selecting indices of texts for a given class label with multiple mixed-source reasoning problems. Similarly, for few-shot setting, Cheng et al. (2023b) and Lin et al. (2024) have found that the overall accuracy decreases with the increase in batch size. Notably, this setting is also meaningful in real-world applications: for independent scenario, batching unrelated queries can reduce model calls and API costs; for sequential scenario, where questions share context—such as in math problem solving, data processing, or software debugging—the correctness of each intermediate reasoning step is critical. Overall, hallucination and performance instability under the multi-problem setting are still under-explored and present significant challenges for current LLMs.

### A.3 Template for QA-Confidence pair

Question: <Question>. Answer: <Answer>. Are you sure you accurately answered the question based on your internal knowledge?  
 1: <Confidence> 2: <Confidence> 3: <Confidence>

### A.4 Dataset Details

We carry out our experiments across six datasets, described as follows.

- **GSM** (Cobbe et al., 2021): a dataset containing high-quality grade school math problems created by the OpenAI group. These problems require between 2 and 8 steps to solve, primarily involving a sequence of elementary calculations with basic arithmetic operations such as addition, subtraction, multiplication, and division to arrive at the final answer. We directly use 7.5k training data and 1k testing data in our Question Answer setting.
- **Pararel** (Elazar et al., 2021): a dataset containing factual knowledge with a variety of prompts and relationships, originally created for mask prediction. In Question Answer setting, we employ the modified dataset from Zhang et al. (2024).
- **MMLU** (Hendrycks et al., 2021): a dataset covering different subjects and difficulty. It tests both world knowledge and problem solving ability, which has good granularity and breadth. We directly use the modified dataset from Zhang et al. (2024) in our Multiple Choice setting.

- **CoQA** (Reddy et al., 2019): a dataset designed to evaluate the ability of models to understand and generate answers in a conversational setting. We randomly pick 5k training dataset from theirs. In Question Answer setting, we combine multiple questions together under the same "story" category in the dataset.
- **MTI Bench** (Son et al., 2024): a comprehensive evaluation benchmark encompassing 5,000 instances across 25 tasks. We pick the sequential part of this benchmark and divide it into 800 training data and 200 test data.
- **SQA** (Iyyer et al., 2017): a dataset designed to explore the task of answering sequences of inter-related questions on HTML tables. We pick 5 sequential questions for each HTML table and have 3985 training data.

### A.5 Formula and Calculation Details

**Average Precision (AP) Score** measures the performance of a binary classifier’s confidence rankings. It corresponds to the area under the Precision-Recall curve. It is calculated as follows:

$$AP = \sum_{k=1}^n (R_k - R_{k-1}) \times P_k$$

where  $k$  is the number of data at current thread with precision  $P_k$  and recall  $R_k$ .  $n$  is the total data number. The confidence is the weighted average of certain prediction probability and uncertain prediction probability.

**Expected Calibrated Error (ECE)** indicates how well a model’s predicted probabilities match the true likelihood of an event. We split the predictions into 10 bins based on the certain prediction probability, then compare the average predicted probability to the actual proportion of positive samples (correct cases) in each bin. It is calculated as follows:

$$ECE = \sum_{m=1}^{10} \frac{|B_m|}{n} |\bar{p}_m - \bar{y}_m|$$

where  $m$  is the bin number with corresponding average predicted probability  $\bar{p}_m$  and actual proportion of positive samples  $\bar{y}_m$ .

### A.6 Implementation

We use HuggingFace PEFT (Mangrulkar et al., 2022) to conduct LoRA fine-tuning (Hu et al.,

2021). We set the training epoch to 3, learning rate to  $1e^{-5}$ , LoRa rank to 8, and LoRa scaling factor to 32. The batch size is 1 and the temperature is 0. All experiments are implemented on Nvidia A100-40GB GPUs.

### A.7 Case Study

We show two specific cases for MAC-Tuning under the multiple problem setting with question number  $n = 3$  in Figure 5. The example on the left is from the SQA dataset, in which a table context is given and the LLM need to answer sequential questions based on the table. LLM answers correctly and shows certainty to first two questions, so these two questions will be counted into accuracy calculation. It answers wrong and shows uncertainty to the third question, which achieves the refusal behavior that we aim to see. The example on the right is from the GSM dataset. The LLM gives wrong answers to the second question but indicates certainty, which means this is a failure case.

### A.8 Detailed Information for Variant Methods

The detailed example for different baseline methods is shown in Figure 6.

### A.9 Extremely Large Multi-problem Scenarios

We also conduct experiments to include larger values of question numbers ( $n = 10$  and  $15$ ). These results are shown in table 8.

Pararel	n=1	n=2	n=3	n=4	n=5	n=10	n=15
LLaMA3	42.9	45.3	42.2	47.0	48.2	48.2	49.6
MAC-Tuning	86.1	81.6	88.0	83.7	86.3	82.6	77.1

Table 8: Accuracy (%) comparison between LLaMA3 and MAC-Tuning on ParaRel dataset with extremely large question number of combined questions  $n$ .

The extended results show that MAC-Tuning consistently and significantly outperforms the Llama3 baseline across all tested values of  $n$ . Notably, even at  $n = 15$ , our method achieves 77.1% accuracy, maintaining a substantial margin of over 27.5 percentage points against the baseline (49.6%).

One notable point is that for different values of  $n$ , we still use 3 training epochs for fairness. However, we only have around 5000 training questions for the Pararel dataset. Thus, for a large value of  $n$  (taking  $n = 10$  as an example), each epoch only has

Context			
Character	First Appeared	Home World	Powers
Night Girl	2007	Kathoon	Super strength
Dragonwing	2010	Earth	Fire breath
Gates	2009	Vyrge	Teleporting
XS	2009	Aarok	Super speed
Harmonia	2011	Earth	Elemental

**Prompt**

<Context> Solve several questions here:  
 1: Who are all of the super heroes?  
 2: Which of them come from Earth?  
 3: Of those, who appeared most recently?  
 Directly Give me an answer without explanation for each question in following format: 1: answer 2: answer 3: answer

**Ground-Truth Answer**

1: Night Girl, Dragonwing, Gates, XS, Harmonia  
 2: Dragonwing, Harmonia  
 3: Harmonia

**Output**

1: Night Girl, Dragonwing, Gates, XS, Harmonia  
 2: Dragonwing, Harmonia  
 3: **Dragonwing**

**Prompt**

<Questions>.<Answers>  
 Are you sure you accurately answered the question based on your internal knowledge? Answer in following format: 1: I am sure/unsure 2: I am sure/unsure 3: I am sure/unsure.

**Output**

1: I am sure 2: I am sure 3: I am unsure

**Prompt**

Solve several questions here:  
 1: A new program had 60 downloads in the first month. The number of downloads in the second month was three times as many as the downloads in the first month, but then reduced by 30% in the third month. How many downloads did the program have total over the three months?  
 2: Carlos is planting a lemon tree. The tree will cost \$90 to plant. Each year it will grow 7 lemons, which he can sell for \$1.5 each. It costs \$3 a year to water and feed the tree. How many years will it take before he starts earning money on the lemon tree?  
 3: Melanie is a door-to-door saleswoman. She sold a third of her vacuum cleaners at the green house, 2 more to the red house, and half of what was left at the orange house. If Melanie has 5 vacuum cleaners left, how many did she start with?  
 Let's think step by step and give me an answer for each question in following format: 1: answer 2: answer 3: answer

**Ground-Truth Answer** 1: 366 2: 13 3: 18

**Output**

1: The program had  $60 \times 3 = 180$  downloads in the second month. The number of downloads in the third month was  $180 \times 0.7 = 126$ . The program had  $60 + 180 + 126 = 366$  downloads over the three months. ##### 366  
 2: The tree will grow 7 lemons per year, and each lemon can be sold for \$1.5, so the tree will earn  $\$1.5 \times 7 = \$10.5$  per year. The tree will cost \$3 per year to water and feed, so the tree will earn  $\$10.5 - \$3 = \$7.5$  per year. **It will take  $90 / 7.5 = 12$  years for the tree to start earning money. ##### 12**  
 3: **Melanie sold  $1/3 \times 5 = 1.67$  vacuum cleaners at the green house. She sold 2 more at the red house, so she sold  $1.67 + 2 = 3.67$  vacuum cleaners in total. She sold  $1/2 \times 5 = 2.5$  vacuum cleaners at the orange house. She sold  $5 - 3.67 - 2.5 = 1$  vacuum cleaner at the remaining houses. She started with  $3.67 + 1 = 4.67$  vacuum cleaners. ##### 4.67**

**Prompt**

<Questions>.<Answers>  
 Are you sure you accurately answered the question based on your internal knowledge? Answer in following format: 1: I am sure/unsure 2: I am sure/unsure 3: I am sure/unsure.

**Output**

1: I am sure  
 2: I am sure  
 3: I am unsure

Figure 5: Specific cases for MAC-Tuning under the multiple problem setting with question number  $n = 3$ . Red-highlighted context indicates inaccuracies in the LLM’s output. The left example is drawn from the *Sequential* setting dataset (SQA), while the right example is from the *Independent* setting dataset (GSM), with one-shot context omitted for conciseness.

500 training data and the model might be under training because of this.  
 Overall, it provides a clearer picture of our method’s performance scalability and demonstrates its robustness in handling more complex, large-scale multi-problem inputs.

**A.10 Certainty Distribution of the Training Dataset**

We demonstrate the certainty distribution of the training dataset under Multiple Problem setting with question number  $n = 3$  in Figure 7:

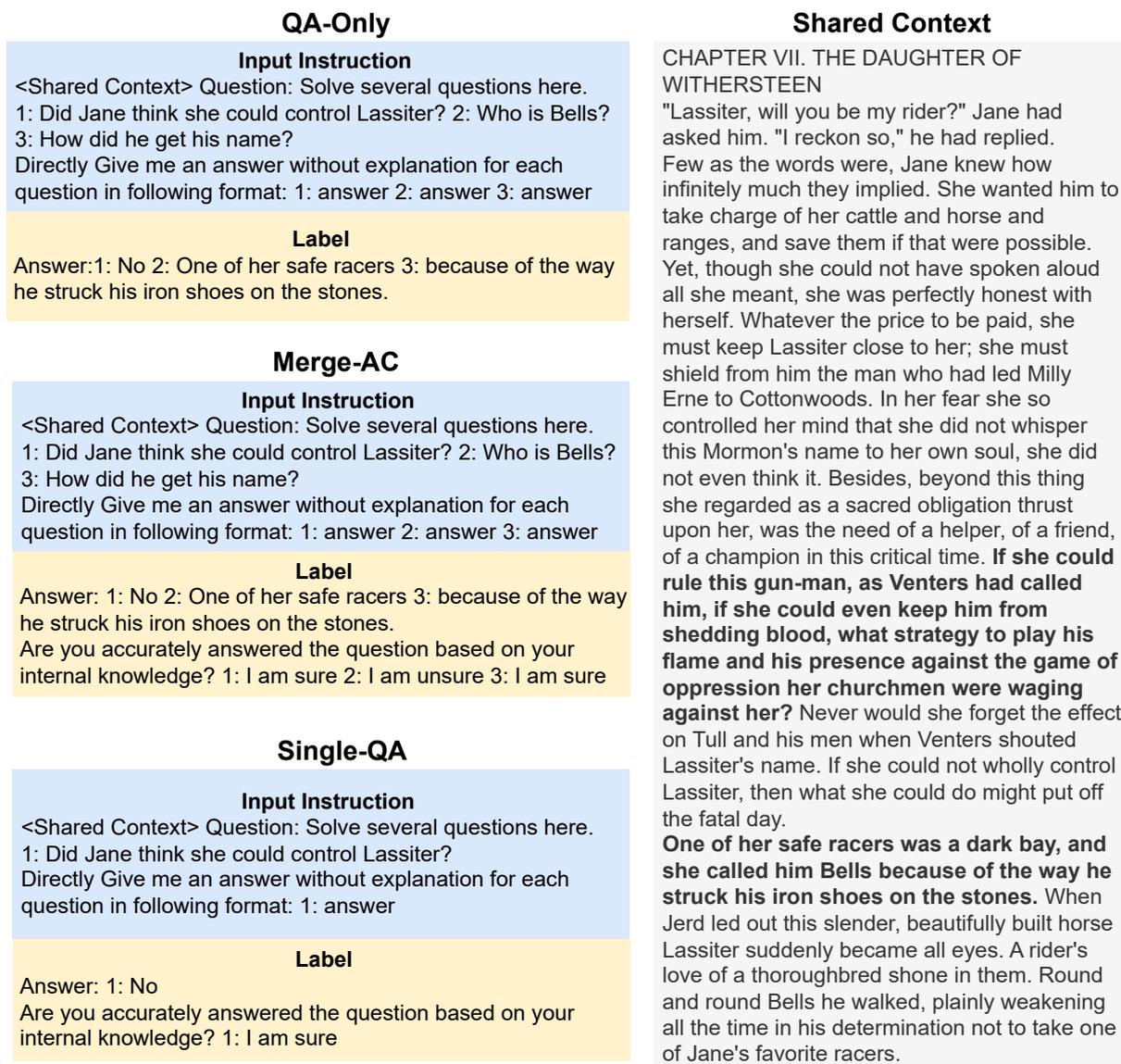


Figure 6: A specific case to show how baseline methods are doing the fine-tuning. The answers are derived from the highlighted portions of the context. In QA-Only, the input is the Question instruction, and the output is the Answer. In Merge-AC, the output includes both the Answer and its Confidence. Single-QA is the single-problem variant of Merge-AC.

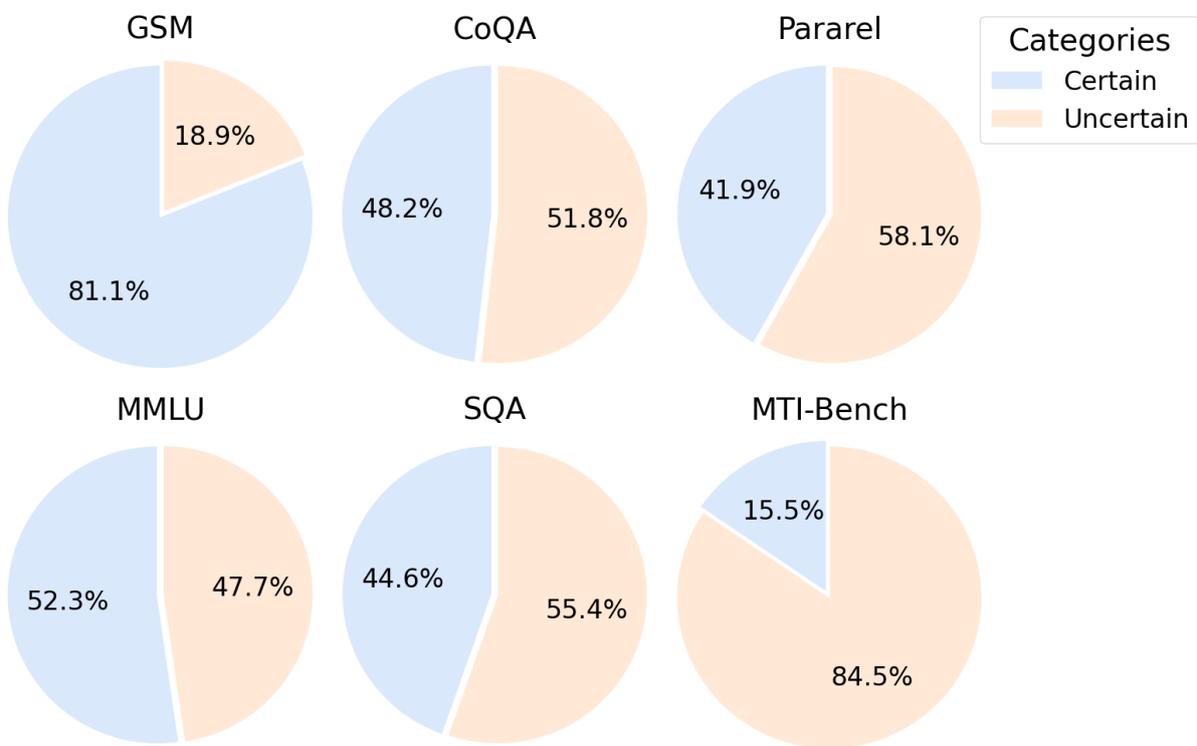


Figure 7: Certainty distribution of the training set under multi-problem setting with  $n = 3$