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

Anonymous ACL submission

Abstract

With the widespread application of large language models (LLMs), the issue of generating non-existing facts, known as hallucination, has garnered increasing attention. Previous 004 research in enhancing LLM confidence estimation mainly focuses on the single problem setting. However, LLM awareness of its internal parameterized knowledge boundary under the more challenging multi-problem setting, which requires answering multiple problems accurately simultaneously, remains underexplored. To bridge this gap, we introduce a novel method, Multiple Answers and Confidence Stepwise Tuning (MAC-Tuning), that separates the learning of answer prediction and con-016 fidence estimation during fine-tuning. Extensive experiments across various base models 017 and different model sizes demonstrate that our method proposed outperforms baselines by up to 25% in average precision.¹

1 Introduction

021

034

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). 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). Notably, these work mainly focus on the single-problem setting, where users repeatedly input questions and context for models to answer one by one.



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). 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

064

039

¹Our code and resource will be released upon publication.

<sup>041
042
043
044
045
046
047
048
049
050
051
052
053
056
057
058
059
061
062
063</sup>

advocates for encouraging the LLM to express 065 confidence to reduce hallucinations, we introduce 066 Multiple Answers and Confidence Stepwise Tuning (MAC-Tuning) under multi-problem setting. Our approach involves several key steps. First, we identify the knowledge boundary between parametric knowledge and the multi-problem dataset to 071 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 multiproblem 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

087

880

100

101

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 origi-102 nal datasets to construct our initial Multi-Problem dataset. We utilize this to compare LLMs' out-104 puts with ground-truth answers, for distinguishing 105 the knowledge boundary between LLM parameters 106 and instruction data. Specifically, for each individ-107 108 ual problem in the multi-problem pair, we assign: "I am sure" if the output aligns with ground-truth 109 answer; "I am unsure" elsewise (e.g., Step 2 in 110 Figure 2). With the assigned confidence labels, we 111 construct Multi-Problem Tuning data as follows: 112 113 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

$$D_{MultQA} = [(q_1, a_1)...(q_i, a_i)...(q_n, a_n)]$$

Step1: Multi-Problem Dataset Step3: MAC-Tuning 🈂

Context a Solve seve Prompt parts did he character did	bout an actor De ral questions he e usually get? id he act?	ennis Farina> re. 1: What 3: Which	Multiple Question-Answer Pair Instruction <context about="" actor="" an="" dennis<br="">Farina> Question: Solve several</context>
Ground-Truth Ans	wer		questions here. 1: What parts did
1: Cops or gangs	ters 3: Joe Fo	ontana	character did he act?
(Status	Label
Output 3: Detect	gangsters		Answer: 1: Cops or gangsters 3: Joe Fontana
Step2: Assign	Confidence	Label	
Question	Answer	Confidence	Multiple QA-Confidence Pair
What parts did	Answer		Instruction
he usually get?	gangsters	I am sure	<multiple pair="" question-answer=""> Are you sure you accurately</multiple>
What happened	He joined	I am sure	answered the questions based on your internal knowledge?
111 2004 !	Law & Older		Label
Which character did he act?	Joe Fontana	I am unsure	1: I am sure 3: I am unsure
			\star

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 form².

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:

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

where Q, A, and C represent the sets of multiple questions, multiple answers, and multiple confidence levels, respectively. Θ_0 and Θ_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** (Hendrycks et al., 2021), and **ParaRel** (Elazar et al., 117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

²The template is in Appendix A.3

				Indep	endent	ndent				Sequential			
Model	Co	QA	QA ParaRel		GSM		MMLU		MTI-Bench		SQA		
	AP	ECE	AP	ECE	AP	ECE	AP	ECE	AP	ECE	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	69.8	7.33	76.1	3.61	79.9	3.16	63.1	12.5	64.0	13.4	65.0	14.6	

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.

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.

		Indepe	ndent		Sequentia	ıl
	CoQA	ParaRel	GSM	MMLU	MTI-Bench	SQA
Frain	5006	7500	7468	2448	2400	3985
Test	5011	5584	1319	2439	600	925
Туре	QA	QA	QA	MC	QA	QA

Table 2: Statistics of the datasets.

3.2 Evaluation Metrics

We directly compare the LLM generation to the 155 ground-truth answer for the Question-Answer for-156 mat. For Multiple-Choice format, we check the choice (A, B, C, D) and the option in the LLM 158 159 generation. Across both types of answer generation tasks, we consider three evaluation metrics: (1) 160 Average Precision (AP): We use AP to measure 161 the precision in identifying and ranking relevant predictions. A higher AP score means the model 163 has high certainty about correct answers and high 164 uncertainty about wrong answers. (2) Expected 165 Calibrated Error (ECE): We use ECE to mea-166 sure how closely the predicted certainty reflects the true certainty of LLM (Chen et al., 2023). Low 168 ECE indicates better-calibrated predictions. (3) Accuracy: We compute accuracy as the fraction 170 of correct responses amongst questions in which 171 LLMs expressed certainty towards their answers. 172

3.3 Baselines 173

We compare MAC-Tuning with the base model and 174 its variants in the multi-problem settings. We use 175

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 levels³.

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

197

198

199

201

202

203

204

206

207

208

210

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.

Ablation on Different Component We further test three variants of the MAC-Tuning method in the multi-problem setting: QA-Only, which is

³Baseline examples are in Appendix A.8. Implementation details are in Appendix A.6.

MAC-Tuning without the confidence component; 211 Single-QA, where we evaluate MAC-Tuning with 212 single problem data; and Merge-AC, where we 213 evaluate MAC-Tuning without separating the learn-214 ing process of ground-truth answers and confidence. As seen from the results in Table 1, MAC-Tuning 216 has up to 25% and, on average, 11% AP improve-217 ment compared with Merge-AC, reflecting that 218 separating the learning process of ground-truth an-219 swers 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 222 base model but worse than **QA-Only** in most cases, showing that LLM can aware the knowledge bound-224 ary 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 *Se-quential* 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

237

239

240

241

242

243

244

245

247

253

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 multiproblem setting. However, for other datasets like MMLU, MAC-Tuning performs slightly worse 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. Further studies are detailed in A.10.



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

3.7 Analysis on Different Base Model

Table 4 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. Results for Llama-3.2-3B and Phi-3.5-mini-Instruct is shown in in Appendix A.9.

		Indep	endent	t		Sequ	iential	
Approach	Par	aRel	MN	1LU	MTI-	Bench	S	QA
	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	78.7	9.59	73.0	17.1	53.3	18.6	47.7	29.2

Table 4: Confidence calibration result (%) for Qwen2-7B-Instruct, with **bold** denoting the top performance.

4 Conclusion

4

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. 264

265

266 267

268

269

270

271

272

273

274

275

277

278

279

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 289 prompt-based LLM studies, we cannot ensure that slight changes in prompts would not significantly 290 alter the results. Second, due to constraints of cost, 291 time, and computational resources, we selected a subset of experiments that we believe to be informative and representative. However, additional 294 experiments across a wider range of datasets and 295 LLMs might provide further insights. Lastly, in this 297 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 300 potential factors affecting performance.

References

302

303

307

309

310

311

312

313

314

315

317

318

319

320

321

322

324

325

326 327

328

332

- 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, 334 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, 335 Akhil Mathur, Alan Schelten, Amy Yang, Angela 336 Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, 337 Archi Mitra, Archie Sravankumar, Artem Korenev, 339 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien 340 Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie 341 Chern, Charlotte Caucheteux, Chaya Nayak, Chloe 342 Bi, Chris Marra, Chris McConnell, Christian Keller, 343 Christophe Touret, Chunyang Wu, Corinne Wong, 344 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-345 lonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, 347 Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, 348 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Geor-351 gia Lewis Anderson, Graeme Nail, Gregoire Mi-352 alon, Guan Pang, Guillem Cucurell, Hailey Nguyen, 353 Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan 354 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 355 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, 358 Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph 361 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, 362 Kalyan Vasuden Alwala, Kartikeya Upasani, Kate 363 Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, 364 Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuen-365 ley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Lau-366 rens van der Maaten, Lawrence Chen, Liang Tan, Liz 367 Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, 368 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, 369 Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, 370 Manohar Paluri, Marcin Kardas, Mathew Oldham, 371 Mathieu Rita, Maya Pavlova, Melanie Kambadur, 372 Mike Lewis, Min Si, Mitesh Kumar Singh, Mona 373 Hassan, Naman Goyal, Narjes Torabi, Nikolay Bash-374 lykov, Nikolay Bogoychev, Niladri Chatterji, Olivier 375 Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan 376 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, 378 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao 379 Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon 380 Calderer, Ricardo Silveira Cabral, Robert Stojnic, 381 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, 383 Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar 384 Hosseini, Sahana Chennabasappa, Sanjay Singh, 385 Sean Bell, Seohyun Sonia Kim, Sergey Edunov, 386 Shaoliang Nie, Sharan Narang, Sharath Raparthy, 387 Sheng Shen, Shengye Wan, Shruti Bhosale, Shun 388 Zhang, Simon Vandenhende, Soumya Batra, Spencer 389 Whitman, Sten Sootla, Stephane Collot, Suchin Gu-390 rurangan, Sydney Borodinsky, Tamar Herman, Tara 391 Fowler, Tarek Sheasha, Thomas Georgiou, Thomas 392 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong 393 Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor 394 Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent 395 Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-396 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-397

ney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine 400 Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue 402 Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng 403 Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, 404 Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam 405 Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesen-406 berg, Alex Vaughan, Alexei Baevski, Allie Feinstein, 407 Amanda Kallet, Amit Sangani, Anam Yunus, An-408 409 drei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew 410 Ryan, Ankit Ramchandani, Annie Franco, Apara-411 jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, 412 413 Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, 414 Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi 415 416 Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, 417 Brian Gamido, Britt Montalvo, Carl Parker, Carly 418 Burton, Catalina Mejia, Changhan Wang, Changkyu 419 Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, 420 Chris Cai, Chris Tindal, Christoph Feichtenhofer, Da-421 mon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, 422 423 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, 424 425 Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Mont-426 gomery, Eleonora Presani, Emily Hahn, Emily Wood, 427 428 Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan 429 Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat 430 Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Flo-431 432 rez, Gabriella Schwarz, Gada Badeer, Georgia Swee, 433 Gil Halpern, Govind Thattai, Grant Herman, Grigory 434 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, 435 Hamid Shojanazeri, Han Zou, Hannah Wang, Han-436 wen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim 437 Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena 438 Veliche, Itai Gat, Jake Weissman, James Geboski, 439 440 James Kohli, Japhet Asher, Jean-Baptiste Gaya, 441 Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, 442 Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, 443 Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, 444 445 Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou 446 U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik 447 Veeraraghavan, Kelly Michelena, Keqian Li, Kun 448 449 Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, 450 Lailin Chen, Lakshya Garg, Lavender A, Leandro 451 Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng 452 Yu, Liron Moshkovich, Luca Wehrstedt, Madian 453 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsim-454 poukelli, Martynas Mankus, Matan Hasson, Matthew 455 Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. 456 Seltzer, Michal Valko, Michelle Restrepo, Mihir 457 Patel, Mik Vyatskov, Mikayel Samvelyan, Mike 458 Clark, Mike Macey, Mike Wang, Miquel Jubert Her-459 460 moso, Mo Metanat, Mohammad Rastegari, Mun-461 ish Bansal, Nandhini Santhanam, Natascha Parks,

401

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, Vítor 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.

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha 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.
- 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).

619

620

621

622

623

624

625

626

627

628

577

578

- 521 522

- 532 533
- 534 535 536
- 538
- 539
- 540

541

543 545

- 550 551

562 563

- 564
- 565 566

567

571

573

574

575 576

- 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.
- 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.
 - 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^2GC: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 parameterefficient 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.

- 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, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2023. Principle-driven selfalignment 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.

Appendix Α

A.1 Full Case for Examples of Introduction

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

Context

The Vatican Apostolic Library, more commonly called the Vatican Library or simply the Vat, is the library of the Holy See, located in Vatican City. **Formally established in 1475**, although it is much older, it is one of the oldest libraries in the world and contains one of the most significant collections of historical texts. It has 75,000 codices from throughout history, as well as 1.1 million printed books, which include some 8,500 incunabula.

The Vatican Library is a research library for **history**, **law**, **philosophy**, **science and theology**. The Vatican Library is open to anyone who can document their qualifications and research needs. Photocopies for private study of pages from books published between 1801 and 1990 can be requested in person or by mail.

In March 2014, the Vatican Library began an initial four-year project of digitizing its collection of manuscripts, to be made available online. The Vatican Secret Archives were separated from the library at the beginning of the 17th century; they contain another 150,000 items. Scholars have traditionally divided the history of the library into five periods, **Pre-Lateran, Lateran, Avignon, Pre-Vatican and Vatican**. The Pre-Lateran period, comprising the initial days of the library, dated from the earliest days of the Church. Only a handful of volumes survive from this period, though some are very significant.



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.

A.2 Related Work

631

643

Hallucination: 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). 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 (Du et al., 2023,Sun et al., 2023), and model confidence calibration (Zhang et al., 2024, Hu et al., 2023).
Knowledge Boundary: There are many different

ways to utilize knowledge boundary to reduce LLM 644 hallucination. (Liang et al., 2024b)'s work uses 645 merged knowledge probing and consistency check-646 ing methods to help LLM express their internal 647 knowledge. (Chen et al., 2024)'s work leverages 648 LLM internal signals to let LLM know their un-649 knowns. (Zhang et al., 2024) utilize knowledge 650 boundary to instruct LLM say "I don't know". It is 651 a popular way to use confidence to express knowl-652 edge boundary of LLMs and we also follow this. 653

Multiple Problem Setting: Current LLM research 654 has predominantly focused on single problem set-655 ting. There are only a few works focusing on this 656 new setting. (Cheng et al., 2023a) propose batch 657 prompting that prompts LLMs with single indepen-658 dent problems batched together following few-shot 659 exemplars together. Son et al. (2024) goes further 660 by researching sequential datasets and develops the 661 first multi-task benchmark (MTI-Bench). Wang 662 et al. (2024) pays attention to zero-shot cases of 663 multi-problem setting and design a new benchmark 664 ZeMPEB. Li et al. (2024) analyze different strategy 665 under independent setting, where single questions 666 are combined into various constraint formats with-667 out sharing context between them. Despite these 668 efforts, the multi-problem setting presents signifi-669 cant challenges. For example, Wang et al. (2024) 670 shows that in zero-shot setting, LLMs consistently 671 perform worse when selecting indices of texts for 672 a given class label with multiple mixed-source rea-673 soning problems. Similarly, for few-shot setting, 674 Cheng et al. (2023b) and Lin et al. (2024) have 675 found that the overall accuracy decreases with the 676 increase in batch size. Notably, this setting is also 677 meaningful in real-world applications: for inde-678 pendent scenario, batching unrelated queries can 679 reduce model calls and API costs; for sequential 680 senario, where questions share context-such as in 681 math problem solving, data processing, or software 682 debugging-the correctness of each intermediate 683 reasoning step is critical. Overall, hallucination and 684 performance instability under the multi-problem 685 setting are still under-explored and present signifi-686 cant challenges for current LLMs. 687

688

690

691

698

701

702

703

707

709

710

711

712

714

715

716

717

718

719

721

723

724

727

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>* 3:

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 breath. 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:

$$AP = \sum_{k=1}^{n} \left(R_k - R_{k-1} \right) \times P_k$$
73

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:

$$\text{ECE} = \sum_{m=1}^{10} \frac{|B_m|}{n} |\overline{p}_m - \overline{y}_m|$$
 75:

where *m* is the bin number with corresponding average predicted probability \overline{p}_m and actual proportion of positive samples \overline{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 730 731 732

728

729

733 734

735

736

737

739

740

741

742

743

744

745

746

747

748

749

750

751

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

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.

772

773

774

775

776

778

779

781

786

787

788

790

791

792

797

798

801

A.8 Detailed Information for Variant Methods

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

A.9 Impact of Model Size on Performance

The confidence calibration results for Llama-3.2-3B and Phi-3.5-mini-Instruct is shown in table 5 and table 6 respectively. Despite using different base models of varying sizes, the results follow the same trend as reported in the main paper. For smaller models, the accuracy improvement after **MAC-Tuning** is more evident, indicating enhanced ability to differentiate between certain and uncertain questions.

		Indep	endent			Sequ	iential	
Approach	Par	arel	Co	QA	MTI-	Bench	S	QA
	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	55.5	33.5	62.4	33.5	35.5	33.5	44.3	33.5

Table 5: Confidence calibration result (%) for Llama-3.2-3B, with **bold** denoting the top performance across different methods.

		Indep	endent			Sequ	ential	
Approach	Par	arel	Co	QA	MTI-	Bench	S	QA
	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	70.2	14.2	68.2	29.0	68.7	22.6	52.3	23.9

Table 6: Confidence calibration result (%) for Phi-3.5mini-Instruct, with **bold** denoting the top performance across different methods.

A.10 Cross Task Transfer Study

We fine-tune the model with question number n = 3 and let it response to question number n = 1 (which is Single Problem setting) and question number n = 5. For the former one, we want to test if the model understands the single problem. For the latter one, we want to test if MAC-Tuning can generalize across different question number n. The results are reported in Table 7.

Question Number	CoQA	ParaRel	GSM	MMLU
n = 1	78.78	84.24	71.12	54.60
n = 5	79.12	86.16	67.67	63.69

Table 7: Accuracy (%) for MAC-Tuning with question number n = 3 transferring to question number n = 1(which is single-problem setting) and question number n = 5. We use one-shot CoT for GSM results.

From the result of n = 1, we observe that accuracy improves in easy dataset like CoQA but decreases in difficult dataset like GSM, comparing with using question number n = 1 to fine-tune. This indicates that LLM learn knowledge during fine-tuning rather than just memorizing the patterns for multi-problem setting. With n = 5, the model performs about the same or even better than using question number n = 5 to fine-tune. This strengthen the statement we make in Section 3.6 that LLMs are more efficient at learning multiple easy tasks but struggle with several difficult tasks simultaneously. 802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

A.11 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:

		Co	ntext						
Charac	ter	First Appeared	Home World	Powers	Prompt Solve	several questions	here:		
Night G	Sirl	2007	Kathoon	Super strength	1: A ne downlo	ew program had 6 bads in the secon	60 downloads in the f d month was three ti	irst month. mes as mar	The r ny as
Dragonw	ving	2010	Earth	Fire breath	How m	nany downloads c	lid the program have	total over the total cost \$90	he thi
Gates	s	2009	Vyrga	Teleporting	year it	will grow 7 lemor	ns, which he can sell the tree. How many y	for \$1.5 eac /ears will it t	ch. It ake b
XS		2009	Aarok	Super speed	starts e 3: Mela	earning money or anie is a door-to-c	n the lemon tree? door saleswoman. Sl	he sold a thi	rd of
Harmor	nia	2011	Earth	Elemental	vas lei	rs at the green ho ft at the orange h did she start with	ouse, 2 more to the r ouse. If Melanie has	ed house, a 5 vacuum c	nd ha leane
Context> Solve several questions here: 1: Who are all of the super heroes?			Let's th followin	nink step by step ng format:1: answ	and give me an ansv ver 2: answer 3: ansv	ver for each wer	n ques		
 Yho are all of the super heroes? Who are all of the super heroes? Which of them come from Earth? Of those, who appeared most recently? Directly Give me an answer without explanation for each question in following format: 1: answer 2: answer 		rth? t recently?	Ground	d-Truth Answer	1: 366 2: 13 3: 18		(P)		
		1: The number had 60	program had 60*3 r of downloads in +180+126=366 d	3=180 downloads in the third month was ownloads over the th	the second 180*0.7=12 ree months	monti 26. Th			
Ground	d-Trut	th Answer		Output	2: The 1	tree will grow 7 le o the tree will ear	emons per year, and n \$1.5*7 = \$10.5 per	each lemon year. The tr	can b ree wi
1: Night Gates, X 2: Drago 3: Harmo	Girl, E S, Ha nwing onia	Dragonwing, armonia g, Harmonia	1: Night Gir Gates, XS, 2: Dragonw 3: Dragonw	l, Dragonwing, Harmonia ing, Harmonia <mark>ing</mark>	per yea It will ta 3: Mela sold 2 r	ar to water and fea ake 90/7.5 = 12 ye anie sold 1/3*5 = 1 more at the red he She sold 1/2*5 =	ed, so the tree will ea ears for the tree to st 1.67 vacuum cleaner ouse, so she sold 1.6 = 2.5 vacuum cleaner	arn \$10.5 - \$ art earning i s at the gree 67+2 = 3.67	\$3 = \$ mone en ho vacui
_	<qu Are</qu 	uestions>. <ai you sure you</ai 	nswers> u accurately an	swered the	sold 5-3 with 3.6	3.67-2.5 = 1 vacu 67+1 = 4.67 vacu	um cleaner at the re um cleaners. #### 4	maining hou .67	ises.
Prompt	que Ans sure	stion based o swer in follow e/unsure 2: I	on your interna ing format: 1: I am sure/unsur	am e 3: I am		<questions>.<!--</td--><td>Answers> ou accurately answer</td><td>ed</td><td>٩</td></questions>	Answers> ou accurately answer	ed	٩
	sure 1: I a	e/unsure. <mark>m sure 2: I ar</mark>	m sure 3: I am	unsure Outpu	Prompt	the question ba knowledge? An 1: I am sure/un: 3: I am sure/un:	ised on your internal swer in following forr sure 2: I am sure/uns sure.	mat: sure	1: ar 2: ar 3: ar

:

Figure 5: Specific cases for MAC-Tuning under the multiple problem setting with question number n = 3. Redhighlighted 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.

QA-Only

Input Instruction

<Shared Context> Question: Solve several questions here.

 Did Jane think she could control Lassiter? 2: Who is Bells?
 How did he get his name?

Directly Give me an answer without explanation for each question in following format: 1: answer 2: answer 3: answer

Label

Answer:1: No 2: One of her safe racers 3: because of the way he struck his iron shoes on the stones.

Merge-AC

Input Instruction

<Shared Context> Question: Solve several questions here.

Did Jane think she could control Lassiter? 2: Who is Bells?
How did he get his name?

Directly Give me an answer without explanation for each question in following format: 1: answer 2: answer 3: answer

Label

Answer: 1: No 2: One of her safe racers 3: because of the way he struck his iron shoes on the stones.

Are you accurately answered the question based on your internal knowledge? 1: I am sure 2: I am unsure 3: I am sure

Single-QA

Input Instruction

<Shared Context> Question: Solve several questions here. 1: Did Jane think she could control Lassiter? Directly Give me an answer without explanation for each question in following format: 1: answer

Label

Answer: 1: No Are you accurately answered the question based on your internal knowledge? 1: I am sure

Shared Context

CHAPTER VII. THE DAUGHTER OF WITHERSTEEN

"Lassiter, will you be my rider?" Jane had asked him. "I reckon so," he had replied. Few as the words were, Jane knew how infinitely much they implied. She wanted him to take charge of her cattle and horse and ranges, and save them if that were possible. Yet, though she could not have spoken aloud all she meant, she was perfectly honest with herself. Whatever the price to be paid, she must keep Lassiter close to her: she must shield from him the man who had led Milly Erne to Cottonwoods. In her fear she so controlled her mind that she did not whisper this Mormon's name to her own soul, she did not even think it. Besides, beyond this thing she regarded as a sacred obligation thrust upon her, was the need of a helper, of a friend, of a champion in this critical time. If she could rule this gun-man, as Venters had called him, if she could even keep him from shedding blood, what strategy to play his flame and his presence against the game of oppression her churchmen were waging against her? Never would she forget the effect on Tull and his men when Venters shouted Lassiter's name. If she could not wholly control Lassiter, then what she could do might put off the fatal day.

One of her safe racers was a dark bay, and she called him Bells because of the way he struck his iron shoes on the stones. When Jerd led out this slender, beautifully built horse Lassiter suddenly became all eyes. A rider's love of a thoroughbred shone in them. Round and round Bells he walked, plainly weakening all the time in his determination not to take one of Jane's favorite racers.

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