

Reading Comprehension LLM

Anonymous ACL submission

Abstract

We propose Read-ComprehensionLLM, an intelligent system utilizing large language models (LLMs) to provide outstanding ability in Reading Comprehension. We adopt syllogism prompting strategies to construct supervised fine-tuning datasets in the reading comprehension domain named Read-Comprehension50k, including instruction samples of two categories(Reading comprehension Multiple choice questions and Reading comprehension short answer questions). Besides we use LoRA and QLoRA method to fine-tune LLMs . Evaluations conducted on multiple benchmarks demonstrate that our model performs better than baseline models. Further resources are available at <https://huggingface.co/datasets/KashiwaByte/DISC-Assignment>

1 Introduction

The recent NLP literature has witnessed a tremendous amount of activity in building models that can follow natural language instructions(Sanh et al., 2022); (Ouyang et al., 2022); (Chung et al., 2022) . These developments are powered by two key components: large pretrained language models (LM) and human-written instruction data (e.g., Prompt-Source and Super-NaturalInstructions , SuperNI for short).

Prior work on instruction finetune need high memory usage, while a new method named LoRA can reduces memory requirements by using a small set of trainable parameters, often termed adapters, while not updating the full model parameters which remain fixed, based on it, another efficient finetuning approach named QLoRA, can backpropagates gradients through a frozen, 4-bit quantized pretrained language model into Low Rank Adapters. These two approach would be utilized in out work.

Reading comprehension tasks require thorough reading of the given context and step-by-step analysis to arrive at the correct option. This poses a

significant challenge for large models, and enhancing reading comprehension abilities can effectively reduce the illusion of large models. Therefore, it becomes imperative to develop a domain-specific large model with reading comprehension capabilities.

To this end,we present our finetuned large language model tailored for effectively solving reading comprehension problems by analyzing and reasoning. We begin by adopting the reading comprehension syllogism prompting strategy to construct supervised fine-tuning datasets in the reading comprehension domain, named Read-Comprehension50k. These datasets are then employed to train Read-ComprehensionLLM with analyzing and reasoning on top of a general domain LLM with 2B parameters.

In order to evaluate the effectiveness Read-ComprehensionLLM, we utilize multiple evaluation benchmarks and experimental results show that Read-ComprehensionLLM outperforms significantly better than the base foundation model in all downstream tasks,In some domains, it even surpasses ChatGPT3.5, which demonstrates the advantage of our work.

2 Related Work

2.1 Traditional NLP models and Limitations

Traditional NLP models have made progress in various reading comprehension scenarios.

However ,The application of NLP models to the reading comprehension sector presents a unique set of challenges. Firstly, reading comprehension tasks require a thorough understanding of long texts. Secondly, reading comprehension tasks require a certain level of logical reasoning ability in order to derive answers from questions and context.Finally, many NLP models show poor adaptability, being designed for single-task performance and lacking cross-task generalization(Mishra et al.,

2022) These challenges underscore the need for future research to develop more robust and adaptable NLP models for the ever-evolving reading comprehension sector.

2.2 Large Language Models for Reading Comprehension

The proposal of LLM-based dialogue systems like ChatGPT OpenAI(OpenAI, 2023), GPT-4 OpenAI(OpenAI et al., 2024) have subverted previous dialogue systems (Zhang et al., 2020). These systems are famous for their zero-shot generalization ability(Zhao et al., 2023). One of the key technologies is instruction-tuning(Wei et al., 2022). Fine-tuning pre-trained LLM through diverse instruction data to obtain the desired behavior pattern has become a common way to domainize LLM like Disc-medllm(Bao et al., 2023); Disc-lawllm(Yue et al., 2023). However, General domain LLMs reveal serious problems of hallucination by generating irrelevant content to the specific case.

2.3 Methods to enhance performance of Large Language Models

Prompt engineering, which involves designing effective prompts for large language models to guide their responses and improve performance in specific tasks or domains. Domain-specific instruction dataset construction, focusing on creating datasets that are tailored to specific domains or tasks, including task design and data partitioning strategies.

Zero-Shot Prompting is an important innovation in the field of Large Language Models (LLMs). Introduced by Radford(Radford et al., 2019), this technique allows us to guide the model to perform new tasks in the absence of large-scale specialized training data by using cleverly designed prompts.

Few-Shot Prompting was proposed by Brown(Brown et al., 2020) and, compared to Zero-Shot Prompting, it helps the model learn specific tasks by providing a small number of input-output examples. The paper describes that through carefully selected high-quality examples, the model’s performance in executing complex tasks can be significantly improved, especially in cases where no examples are available at all.

To overcome the limitations of Large Language Models (LLMs) in handling complex reasoning tasks, Wei (Wei et al., 2023) proposed an innovative approach called CoT. This technique introduces a special prompting strategy aimed at facilitating a more continuous and step-by-step thinking process

in the model. In comparison to traditional prompting methods, the primary contribution of CoT lies in its ability to more effectively prompt LLMs to generate structured and deeply considered answers.

SFT is an instruction fine-tuning dataset that provides explicit guidance for training language models. The primary characteristic of the SFT dataset is its collection in real-world environments, with a focus on instructions that align with human understanding and generation.

Some efficient fine-tuning strategies such as LoRA(Hu et al., 2021) (Layer-wise Adaptive Rate Scaling) and QLoRA(Dettmers et al., 2023) (Quantized Layer-wise Adaptive Rate Scaling), which aim to optimize the fine-tuning process for large language models, can reduce memory requirements by using a small set of trainable parameters

3 Method

3.1 Read-Comprehension50k Datasets

To train Read-ComprehensionLLM, we construct a high-quality supervised fine-tuning dataset, Read-Comprehension50k with two subsets, namely CosmosQA25K and TriviaQA25K. The former aims to enhance the capabilities of LLMs in multiple choice questions and standardized outputs, while the latter seeks to bolster the abilities of LLMs in responding to long-text short answer questions. The core of dataset construction involves creating <instruction, input, output> triplets based on prompt and the content of the original dataset. For specific format of datasets, please refer to Appendix A.

Dataset	Samples	Input Token
Multiple choice	25k	230
Short answer	25k	350
Total	50k	300

Table 1: Data statistics of the Read-Comprehension50k dataset.

3.1.1 CosmosQA25k

The CosmosQA(Huang et al., 2019) dataset comprises over 35,000 questions and more than 16,000 article paragraphs, sourced from diverse fields such as Wikipedia, news, fiction, and history. Each question is accompanied by one correct answer and also includes three incorrect answers as distractors. This design makes Cosmos QA a dataset suited for Multiple-Choice Question Answering (MCQA) tasks.

172	Utilizing ChatGPT, we designed a prompt and	of 1e-4, LoRA rank of 8, dropout parameters of 0.1	221
173	refined it according to the principles of Prompt	, 3 epochs training stage, maximum target length	222
174	Engineering to serve as the instruction. We then	of 512 tokens. The training process was carried	223
175	consolidated the context, question, answer0, an-	out on an 3090 GPU and the training cost is further	224
176	swer1, answer2, and answer3 into a single JSON	reduced with the help of deepspeed.	225
177	pair as the input. The label of the correct answer is		
178	used as the output.		
179	The crafted prompt is as follows: "As a reading	3.2.2 QLoRA Finetune for ChatGLM3-6B	226
180	comprehension expert, you will receive context,	We used the Xtuner(Contributors, 2023) framework	227
181	question, and four answer options. Please under-	for QLoRA fine-tuning, with the following hyper-	228
182	stand the given context first and then output the	parameter settings: global batch size of 1, bit quan-	229
183	label of the correct option as the answer to the	tization of 4, learning rate of 2e-4, LoRA rank of 64,	230
184	question based on the context."	dropout parameters of 0.1, 3 epochs training stage,	231
		maximum source length of 512 tokens, The train-	232
185	3.1.2 TriviaQA25K	ing process was carried out on an 3090 GPU and	233
186	TriviaQA(Joshi et al., 2017) is a challenging	the training utilized deepspeed Rasley et al.(2020).	234
187	reading comprehension dataset that contains over		
188	650,000 question-answer-evidence triplets. Triv-	4 Experiment	235
189	iaQA includes 95,000 question-answer pairs au-	4.1 Evaluation Setup	236
190	thored by trivia enthusiasts, accompanied by inde-	To evaluate the overall performance of the fine-	237
191	pendently collected evidence documents. However,	tuned model on reading comprehension tasks, we	238
192	the original TriviaQA dataset is not suitable for	primarily tested it on two types of tasks: Multiple	239
193	use as an instructional dataset. Therefore, we	Choice questions and Short answer questions. We	240
194	have reformatted its data structure to align with	employed various prompt strategies such as Zero-	241
195	the format of the Stanford Question Answering	shot, fewshot, and CoT. For specific prompt tem-	242
196	Dataset(Rajpurkar et al., 2016) (SQuAD). This ad-	plates, please refer to Appendix B.	243
197	justment involves structuring the data to better fa-		
198	cilitate the development and evaluation of models	4.1.1 Multiple Choice Questions task	244
199	on question answering tasks.	For Multiple-Choice questions, we conducted ex-	245
200	Leveraging ChatGPT and the principles of	periments using the official testing link and test	246
201	Prompt Engineering, we designed a prompt to serve	set provided by CosmosQA. We conducted mul-	247
202	as the instruction for processing the data. We then	multiple rounds of testing using the methods of CoT,	248
203	paired <context, question> as the input and desig-	ZeroShot, and FewShot.	249
204	nated the answer as the output.		
205	The crafted prompt is as follows: "As a reading	4.1.2 Short Answer Questions task	250
206	comprehension expert, you will receive context and	For Short answer questions, we partitioned a 7k	251
207	a question. Please understand the given context first	SQuAD-formatted TriviaQA dataset for testing, and	252
208	and then output the answer to the question based	modified the official validation code to align with	253
209	on the context."	our format while retaining core metrics (Exact and	254
		F1).	255
210	3.2 LLM Finetuning	4.2 Main Results	256
211	We utilized MiniCPM-2B(Hu et al., 2024) and	In this section, we present evaluation results of our	257
212	ChatGLM3-6B(Du et al., 2022) as the base models	model on above two tasks in the reading compre-	258
213	for fine-tuning, applying the LoRA method to the	hension domain.	259
214	former and the QLoRA method to the latter. For		
215	both fine-tuning processes, we used a 3090 GPU	4.2.1 Multiple Choice Questions task	260
216	and leveraged the DeepSpeed framework (Rasley	We conducted both ablation studies and compar-	261
217	et al., 2020). to accelerate training.	ative research on the fine-tuned models to assess	262
218	3.2.1 LoRA Finetune for MiniCPM-2B	their performance and identify the impact of differ-	263
219	The hyperparameters setting of this training process	ent modifications.	264
220	are as follows: global batch size of 4, learning rate		

Model	Score
Fewshot MiniCPM	0.3251
Fewshot LoRA MiniCPM	0.7773
Fewshot CoT LoRA MiniCPM	0.7790
CoT LoRA MiniCPM	0.8211
ZH LoRA MiniCPM	0.8215
LoRA MiniCPM	0.8291

Table 2: Ablation Studies

Unfinetuned Model Performance In ablation studies, we observed that the basic MiniCPM-2B model essentially lacks the capability for reading comprehension and selection. Under ZeroShot conditions, it is utterly incapable of completing tasks, and even with FewShot, its performance is only slightly better than pure randomness.

Multilingual Capability After LoRA fine-tuning, the MiniCPM-2B was able to achieve respectable results, ranking Top 36 on the evaluation leaderboard. When using Chinese prompts, the performance of the LoRA fine-tuned model showed only a minor decrease, reaching Top 42. This demonstrates that MiniCPM-2B possesses multilingual capabilities and can also be applied to Chinese reading comprehension tasks.

Ineffectiveness of Prompt It appears that smaller models have lower receptivity to prompts and FewShot learning. Experiments indicate that incorporating FewShot and CoT tends to degrade performance.

Model	Score	Model Size
ChatGPT3.5	0.7233	175B
LoRA MiniCPM	0.8291	2B
QLoRA Chatglm3	0.8416	6B

Table 3: Comparative Experiments

Effectiveness of SFT In the comparative experiments, we contrasted the fine-tuned MiniCPM-2B with fine-tuned ChatGLM3-6B and ChatGPT3.5. The results demonstrated that small-parameter models, after undergoing instruction-based fine-tuning, were able to surpass the performance of ChatGPT3.5 in specific tasks. This highlights the potential of smaller models to achieve competitive results in targeted applications when effectively fine-tuned.

4.2.2 Short Answer Questions task

In the Short Answer questions task, we repeatedly tested eight scenarios including FewShot and ZeroShot LoRA fine-tuned MiniCPM-2B, the original MiniCPM-2B model, QLoRA fine-tuned ChatGLM3-6B, and ChatGLM3-6B itself.

Limitations of small-sized Models Unfortunately, none of these configurations yielded practically usable results. This outcome suggests that despite the fine-tuning efforts, the models may still face challenges in handling the complexity or specific requirements of short answer question tasks, indicating a need for further model optimization or exploring alternative approaches.

Importance of LongText Finetuning We speculate that the reason for the training failures is that the MAX Token LENGTH setting during model fine-tuning was too small, preventing the models from effectively handling reading comprehension tasks with long contexts. This limitation likely restricted the models’ ability to process and understand the full scope of the provided texts, thereby impacting their performance on tasks requiring detailed comprehension of lengthy passages. Adjusting the MAX Token LENGTH parameter to accommodate longer contexts could potentially improve model performance in future experiments.

Performance of ChatGPT3.5 Additionally, we tested ChatGPT3.5 which is able to accommodate 4k tokens and the results indicated that it performed well. On a test set of 1,000 entries, it achieved an Exact Match score of 0.157 and an F1 score of 0.377, with only 73 instances of misunderstanding. This performance highlights the model’s effectiveness in grasping and responding to short answer questions, suggesting a robust comprehension capability compared to the earlier tested models.

5 Conclusion

In this paper, we constructed an instruction fine-tuning dataset specifically for the reading comprehension domain and developed two domain-fine-tuned models using LoRA and QLoRA methods. Our evaluation results demonstrate the effectiveness of our models on the Multiple Choice Questions task. Additionally, we identified significant limitations in small-parameter models when dealing with long-text reading comprehension challenges. Besides, we found that prompt strategies

346 such as CoT and Fewshot did not achieve the ex-
 347 pected results on small-sized models. These find-
 348 ings have provided us with two research directions:
 349 1. Identifying the scale boundary of models that can
 350 handle long texts. 2. Modifying prompt strategies
 351 to adapt to different model scale. For example, de-
 352 signing a strategy that can automatically optimize
 353 prompt based on the model scale.

354 References

355 Zhijie Bao, Wei Chen, Shengze Xiao, Kuang Ren, Jiaao
 356 Wu, Cheng Zhong, Jiajie Peng, Xuanjing Huang, and
 357 Zhongyu Wei. 2023. *Disc-medllm: Bridging gen-
 358 eral large language models and real-world medical
 359 consultation.*

360 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
 361 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
 362 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
 363 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
 364 Gretchen Krueger, Tom Henighan, Rewon Child,
 365 Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens
 366 Winter, Christopher Hesse, Mark Chen, Eric Sigler,
 367 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack
 368 Clark, Christopher Berner, Sam McCandlish, Alec
 369 Radford, Ilya Sutskever, and Dario Amodei. 2020.
 370 *Language models are few-shot learners.* *ArXiv*,
 371 abs/2005.14165.

372 Hyung Won Chung, Le Hou, Shayne Longpre, Barret
 373 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
 374 Wang, Mostafa Dehghani, Siddhartha Brahma, Al-
 375 bert Webson, Shixiang Shane Gu, Zhuyun Dai,
 376 Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-
 377 ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,
 378 Dasha Valter, Sharan Narang, Gaurav Mishra, Adams
 379 Yu, Vincent Zhao, Yanping Huang, Andrew Dai,
 380 Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-
 381 cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,
 382 and Jason Wei. 2022. *Scaling instruction-finetuned
 383 language models.*

384 XTuner Contributors. 2023. *Xtuner: A toolkit for
 385 efficiently fine-tuning llm.* [https://github.com/
 386 InternLM/xtuner](https://github.com/InternLM/xtuner).

387 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and
 388 Luke Zettlemoyer. 2023. *Qlora: Efficient finetuning
 389 of quantized llms.*

390 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding,
 391 Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. *Glm:
 392 General language model pretraining with autoregres-
 393 sive blank infilling.* In *Proceedings of the 60th An-
 394 nual Meeting of the Association for Computational
 395 Linguistics (Volume 1: Long Papers)*, pages 320–335.

396 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
 397 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
 398 Weizhu Chen. 2021. *Lora: Low-rank adaptation of
 399 large language models.*

Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu
 Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxi-
 ang Huang, Weilin Zhao, et al. 2024. *Minicpm:
 Unveiling the potential of small language models
 with scalable training strategies.* *arXiv preprint
 arXiv:2404.06395*.

Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and
 Yejin Choi. 2019. *Cosmos qa: Machine reading com-
 prehension with contextual commonsense reasoning.*

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke
 Zettlemoyer. 2017. *TriviaQA: A large scale distantly
 supervised challenge dataset for reading comprehen-
 sion.* In *Proceedings of the 55th Annual Meeting of
 the Association for Computational Linguistics (Vol-
 ume 1: Long Papers)*, pages 1601–1611, Vancouver,
 Canada. Association for Computational Linguistics.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, and
 Hannaneh Hajishirzi. 2022. *Cross-task generaliza-
 tion via natural language crowdsourcing instructions.*

OpenAI. 2023. *Chatgpt.* [https://openai.com/blog/
 chatgpt](https://openai.com/blog/chatgpt).

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,
 Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
 man, Diogo Almeida, Janko Altmenschmidt, Sam Alt-
 man, Shyamal Anadkat, Red Avila, Igor Babuschkin,
 Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-
 ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-
 wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,
 Christopher Berner, Lenny Bogdonoff, Oleg Boiko,
 Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-
 man, Tim Brooks, Miles Brundage, Kevin Button,
 Trevor Cai, Rosie Campbell, Andrew Cann, Brittany
 Carey, Chelsea Carlson, Rory Carmichael, Brooke
 Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully
 Chen, Ruby Chen, Jason Chen, Mark Chen, Ben
 Chess, Chester Cho, Casey Chu, Hyung Won Chung,
 Dave Cummings, Jeremiah Currier, Yunxing Dai,
 Cory Decareaux, Thomas Degry, Noah Deutsch,
 Damien Deville, Arka Dhar, David Dohan, Steve
 Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,
 Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
 Simón Posada Fishman, Juston Forte, Isabella Ful-
 ford, Leo Gao, Elie Georges, Christian Gibson, Vik
 Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-
 Lopes, Jonathan Gordon, Morgan Grafstein, Scott
 Gray, Ryan Greene, Joshua Gross, Shixiang Shane
 Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,
 Yuchen He, Mike Heaton, Johannes Heidecke, Chris
 Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,
 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin
 Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,
 Joanne Jang, Angela Jiang, Roger Jiang, Haozhun
 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-
 woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-
 mali, Ingmar Kanitscheider, Nitish Shirish Keskar,
 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,
 Christina Kim, Yongjik Kim, Jan Hendrik Kirchner,
 Jamie Kiros, Matt Knight, Daniel Kokotajlo,
 Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-
 stantinidis, Kyle Kopic, Gretchen Krueger, Vishal

460	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	Yuxiong He. 2020. Deepspeed: System optimiza-	521
461	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	tions enable training deep learning models with over	522
462	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	100 billion parameters. In <i>Proceedings of the 26th</i>	523
463	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	<i>ACM SIGKDD International Conference on Knowl-</i>	524
464	Anna Makanju, Kim Malfacini, Sam Manning, Todor	<i>edge Discovery & Data Mining (KDD '20, Tutorial).</i>	525
465	Markov, Yaniv Markovski, Bianca Martin, Katie		
466	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	Victor Sanh, Albert Webson, Colin Raffel, Stephen	526
467	McKinney, Christine McLeavey, Paul McMillan,	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	527
468	Jake McNeil, David Medina, Aalok Mehta, Jacob	Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey,	528
469	Menick, Luke Metz, Andrey Mishchenko, Pamela	M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya	529
470	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	Sharma, Eliza Szczechla, Taewoon Kim, Gunjan	530
471	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan	531
472	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	Chang, Mike Tian-Jian Jiang, Han Wang, Matteo	532
473	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	Manica, Sheng Shen, Zheng Xin Yong, Harshit	533
474	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	Pandey, Rachel Bawden, Thomas Wang, Trishala	534
475	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	Neeraj, Jos Rozen, Abheesht Sharma, Andrea San-	535
476	tista Parascandolo, Joel Parish, Emy Parparita, Alex	tilli, Thibault Fevry, Jason Alan Fries, Ryan Tee-	536
477	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	han, Teven Le Scao, Stella Biderman, Leo Gao,	537
478	man, Filipe de Avila Belbute Peres, Michael Petrov,	Thomas Wolf, and Alexander M Rush. 2022. Multi-	538
479	Henrique Ponde de Oliveira Pinto, Michael, Poko-	task prompted training enables zero-shot task gener-	539
480	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	alization. In <i>International Conference on Learning</i>	540
481	ell, Alethea Power, Boris Power, Elizabeth Proehl,	<i>Representations (ICLR).</i>	541
482	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,		
483	Cameron Raymond, Francis Real, Kendra Rimbach,	Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin	542
484	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	Guu, Adams Wei Yu, Brian Lester, Nan Du, An-	543
485	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	drew M. Dai, and Quoc V. Le. 2022. <i>Finetuned</i>	544
486	Girish Sastry, Heather Schmidt, David Schnurr, John	<i>language models are zero-shot learners.</i>	545
487	Schulman, Daniel Selsam, Kyla Sheppard, Toki		
488	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	546
489	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	547
490	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	Denny Zhou. 2023. <i>Chain-of-thought prompting elic-</i>	548
491	Sokolowsky, Yang Song, Natalie Staudacher, Fel-	<i>its reasoning in large language models.</i>	549
492	ipe Petroski Such, Natalie Summers, Ilya Sutskever,		
493	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	Shengbin Yue, Wei Chen, Siyuan Wang, Bingxuan Li,	550
494	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	Chenchen Shen, Shujun Liu, Yuxuan Zhou, Yao Xiao,	551
495	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fel-	Song Yun, Xuanjing Huang, and Zhongyu Wei. 2023.	552
496	ipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	<i>Disc-lawllm: Fine-tuning large language models for</i>	553
497	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	<i>intelligent legal services.</i>	554
498	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,		
499	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,	555
500	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing	556
501	Clemens Winter, Samuel Wolrich, Hannah Wong,	Liu, and Bill Dolan. 2020. <i>Dialogpt: Large-scale</i>	557
502	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	<i>generative pre-training for conversational response</i>	558
503	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	<i>generation.</i>	559
504	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong		
505	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	560
506	Zheng, Juntang Zhuang, William Zhuk, and Barret	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	561
507	Zoph. 2024. <i>Gpt-4 technical report.</i>	Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen	562
		Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang,	563
508	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu,	564
509	roll L Wainwright, Pamela Mishkin, Chong Zhang,	Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. <i>A</i>	565
510	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	<i>survey of large language models.</i>	566
511	2022. Training language models to follow instruc-		
512	tions with human feedback. In <i>Advances in Neural</i>	A Read-Comprehension50k Datasets	567
513	<i>Information Processing Systems (NeurIPS).</i>	Format	568
514	Alec Radford, Jeff Wu, Rewon Child, David Luan,	In this section, we will demonstrate the format of	569
515	Dario Amodei, and Ilya Sutskever. 2019. <i>Language</i>	fine-tuning instructions for two subcategories of	570
516	<i>models are unsupervised multitask learners.</i>	datasets for model fine-tuning.	571
517	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and		
518	Percy Liang. 2016. <i>Squad: 100,000+ questions for</i>		
519	<i>machine comprehension of text.</i>		
520	Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and		

SFT CosmosQA25K Dataset Format

"instruction": "As a reading comprehension expert, you will receive context, question and four answer options. Please understand the given Context first and then output the label of the correct option as the answer to the question based on the Context",
 "input": (

```
'context':context,
  'question':question,
  "answer0":answer0,
  "answer1":answer1,
  "answer2":answer2,
```

"output":label

SFT TriviaQA25K Dataset Format

"instruction": "As a reading comprehension expert, you will receive context and question. Please understand the given Context first and then output the answer of the question based on the Context",

```
"input": str(
  'context':context,
  'question':question
),
```

"output":answer

ZeroShot Prompt for TriviaQA

As a reading comprehension expert, you will receive context and question. Please understand the given Context first and then output the answer of the question based on the Context ,

Context: (input Context)

Answer:

CoT Prompt

As a reading comprehension expert, you will receive context, question and four answer options. Please understand the given Context first and then output the label of the correct option as the answer to the question based on the Context,

Context: Good Old War and person L : I saw both of these bands Wednesday night , and they both blew me away . seriously . Good Old War is acoustic and makes me smile . I really can not help but be happy when I listen to them ; I think it 's the fact that they seemed so happy themselves when they played .
 , 'question': 'In the future , will this person go to see other bands play ?', 'answer0': 'None of the above choices .', 'answer1': 'This person likes music and likes to see the show , they will see other bands play .', 'answer2': 'This person only likes Good Old War and Person L , no other bands .', 'answer3': 'Other Bands is not on tour and this person can not see them .',

Answer: "The person enjoys music and attending live shows, as indicated by their positive experience watching Good Old War and Person L. Given their enjoyment and the positive impact of the performances, it is likely that they will go see other bands play in the future. Therefore, the correct answer is label 1",

Context: (input Context)

Answer:

B Prompt Engineering

In this section, we will sequentially showcase the ZeroShot Prompt for CosmosQA, ZeroShot Prompt for TriviaQA, CoT Prompt, and Fewshot Prompt.

ZeroShot Prompt for CosmosQA

As a reading comprehension expert, you will receive context, question and four answer options. Please understand the given Context first and then output the label of the correct option as the answer to the question based on the Context,

Context: (input Context)

Answer:

572

573

574

575

576

577

578

579

580

Fewshot Prompt

As a reading comprehension expert, you will receive context, question and four answer options. Please understand the given Context first and then output the label of the correct option as the answer to the question based on the Context,

Context: Good Old War and person L : I saw both of these bands Wednesday night , and they both blew me away . seriously . Good Old War is acoustic and makes me smile . I really can not help but be happy when I listen to them ; I think it 's the fact that they seemed so happy themselves when they played .

, 'question': 'In the future , will this person go to see other bands play ?', 'answer0': 'None of the above choices .', 'answer1': 'This person likes music and likes to see the show , they will see other bands play .', 'answer2': 'This person only likes Good Old War and Person L , no other bands .', 'answer3': 'Other Bands is not on tour and this person can not see them .',

Answer: "1",

Context: ""A hot girl appears who flirts with me to convince me to help her dig the deeper well (Whoa ... even my dreams reference Buffy . Whedon is my God) . I ' m in a narrow tunnel helping to pull a horse attached to a plow of some sort , while the brother of the hot chick is riding on the back of the plow . By the time we 're done I can tell it 's becoming night , though oddly when I get out of the tunnel it 's still light outside ."" , 'question': 'Why do even my dreams reference Buffy ?', 'answer0': ""I like Joss Whedon 's work a little bit ."" , 'answer1': 'I think Buffy the Vampire Slayer is an alright TV show .', 'answer2': 'I love the TV show Buffy the Vampire Slayer .', 'answer3': 'None of the above choices .',

Answer: "2",

Context: (input Context)

Answer: