

First Step Advantage: Importance of Starting Right in Multi-Step Reasoning

Anonymous ACL submission

Abstract

Large Language Models (LLMs) can solve complex reasoning tasks by generating rationales for their predictions. Distilling these capabilities into a smaller, compact model can facilitate the creation of specialized, cost-effective models tailored for specific tasks. However, smaller models often face challenges in complex reasoning tasks and often deviate from the correct reasoning path. We show that LLMs can guide smaller models and bring them back to the correct reasoning path only if they intervene at the right time. We show that smaller models fail to reason primarily due to their difficulty in initiating the process, and that guiding them in the right direction can lead to a performance gain of over 100%. We explore different model sizes and evaluate the benefits of providing guidance to improve reasoning in smaller models.

1 Introduction

Over the years, Large Language Models (LLMs) have improved their reasoning skills by explaining their intermediate thoughts (Wei et al., 2022). This allows LLMs to transfer intermediate knowledge to student models¹ to improve their reasoning skills; often referred to as knowledge distillation (Yuan et al., 2023; Magister et al., 2023; Shridhar et al., 2023b; Hsieh et al., 2023). While training student models can certainly improve their reasoning skills, there are instances where teacher intervention remains essential to guide the student model when it encounters uncertainty or confusion. This is similar to the situation in a classroom, where a student can acquire knowledge independently by learning from textbooks, but often benefits from the guidance of a teacher (Wood et al., 1976; Van de Pol et al., 2015).

¹We refer to smaller models as student models and larger models as teachers. The distinction is not based on the number of parameters, but rather on relative size, with smaller models often referred to as students.

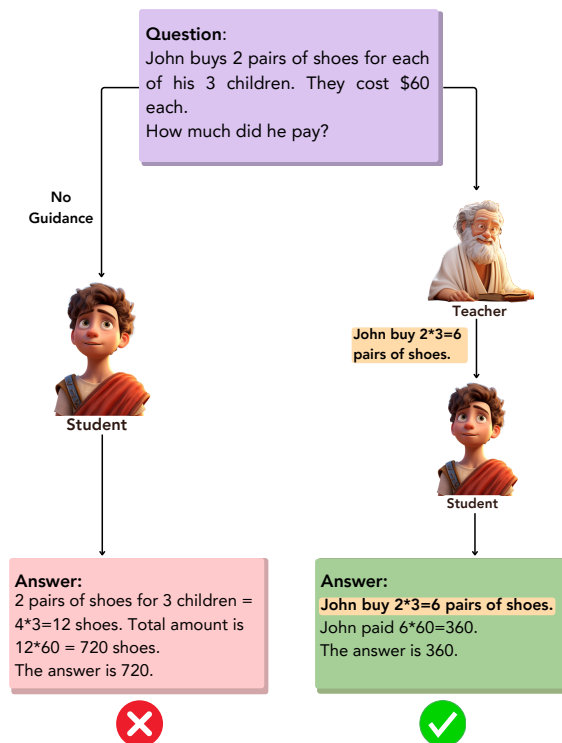


Figure 1: **The First Step Advantage:** Figure demonstrating the effect of first step guidance from a teacher on the student’s overall performance (right) versus no guidance (left).

While teacher intervention can provide valuable assistance to students, understanding *when* and *how* to provide guidance plays a critical role. In general, more guidance typically leads to a consistent improvement in student performance (Saha et al., 2023), but the question arises: *should guidance at different stages be given equal importance?*

Our observations, as shown in Figure 1, suggest that for multi-step reasoning tasks, intervening at the first step yields the most significant benefits, with the effects diminishing at subsequent steps. As expected, expert teachers (models with superior performance) tend to provide more effective guidance, resulting in better student performance.

051 On a mathematical dataset of multi-step word 101
052 problems, we demonstrate the effectiveness of *first* 102
053 *step guidance* on various combinations of teacher- 103
054 student pairs. A consistent improvement in student 104
055 performance was observed with first step guidance. 105
056 We show that a smaller student model (LLaMA 106
057 13B (Touvron et al., 2023)), when correctly guided 107
058 at its first step, can achieve the same performance 108
059 as a larger student model (LLaMA 70B) without 109
060 any guidance. Furthermore, our results show a con- 110
061 tinuous and upward improvement of the student 111
062 model’s performance with expert guidance, with 112
063 GPT-4 (OpenAI, 2023) as the teacher reaching a 113
064 performance level similar to that of a human in- 114
065 structor.

066 2 Related Work

067 Previous work has shown that it is possible to elicit 115
068 reasoning abilities from LLMs through in-context 116
069 learning (Wei et al., 2022; Zhou et al., 2022). A key 117
070 recipe in most methods is to spread the reasoning 118
071 process over multiple tokens, rather than expecting 119
072 them to provide an immediate response token. One 120
073 way is to provide the model with the intermediate 121
074 steps, or chain of thought (CoT), that leads to the 122
075 final answer (Wei et al., 2022; Kojima et al., 2023; 123
076 Yang et al., 2023; Wang et al., 2023). In parallel, in- 124
077 context learning has been used to teach the model 125
078 how to break a problem down into smaller, easier 126
079 sub-problems, and then solve those sub-problems 127
080 that eventually lead to the final answer (Shridhar 128
081 et al., 2022; Zhou et al., 2023).

082 However, if the problem is misinterpreted, it 129
083 can lead to a cascade of errors in subsequent steps. 130
084 To counter this, several techniques have been pro- 131
085 posed to intervene and correct intermediate steps 132
086 (Welleck et al., 2022) or to provide feedback on 133
087 their own generations, essentially “self-correcting” 134
088 their own generations (Madaan et al., 2023; Shrid- 135
089 har et al., 2023a). It is important to note that com- 136
090 plex reasoning and self-correcting capabilities only 137
091 emerge in very large language models. For smaller 138
092 models, CoT follows a rather flat scaling curve 139
093 (Wei et al., 2022). Our work presents an effec- 140
094 tive way to transfer such reasoning capabilities into 141
095 smaller models that does not require large scale pre- 142
096 training, making it more accessible for researchers 143
097 with limited compute. 144

098 While the LLM’s ability to revise its own gener- 145
099 ations may prove helpful in many cases, it some- 146
100 times leads to worse outcomes upon refinement, 147

101 requiring a “rolling-back” to the previous out- 102
102 put (Shridhar et al., 2023a). The “rolling-back” 103
103 dilemma can be avoided if we can know *when* to in- 104
104 tervene. (Saha et al., 2023) presented an approach 105
105 based on ToM (Kosinski, 2023; Kadavath et al., 106
106 2022), where a teacher model intervenes in a stu- 107
107 dent model only for harder questions by creating 108
108 an implicit mental model of the student’s under- 109
109 standing. In contrast, our work shows that it is not 110
110 necessary to intervene and help a student with the 111
111 entire solution. Rather, just starting correctly has 112
112 a significant impact on the student’s performance 113
113 and avoids the need for backtracking and correcting 114
114 mistakes, saving time and effort.

115 3 Experimental Design

116 **Dataset** We examine our intervention on multi- 116
117 step mathematical dataset : GSM8K (Cobbe et al., 117
118 2021). The dataset is a grade-school-level math 118
119 word problem dataset with a training set of 7473 119
120 samples and a test set of 1319 samples, each requir- 120
121 ing two to eight steps to solve.

122 **Setup** We used variants of LLaMA models 122
123 (LLaMA 7B, LLaMA 13B and LLaMA 70B) (Tou- 123
124 vron et al., 2023) and its variants (Mistral 7B 124
125 (Jiang et al., 2023) and MetaMath 7B and 13B (Yu 125
126 et al., 2023)) as student models. On the other hand, 126
127 we used LLaMA 13B and LLaMA 70B as teach- 127
128 ers alongside ChatGPT (gpt-3.5-turbo) (Brown 128
129 et al., 2020) and GPT-4 (gpt-4) (OpenAI, 2023). 129
130 For pre-trained models as students in few-shot set- 130
131 tings, 4-shot demonstrations were provided, which 131
132 were chosen randomly from the train set. Fine- 132
133 tuned students were trained on the training set 133
134 with no modifications for 3 epochs. All models 134
135 were evaluated in the greedy approach (temp=0, 135
136 top p=1). Fine-tuning was performed on 1 node of 136
137 8 A100 GPUs. We report the accuracy (maj@1) on 137
138 the test set. 138

139 4 Results and Discussion

140 **Early intervention is key** Figure 2 compares the 140
141 effect of intervention by humans as teachers for 141
142 LLaMA 7B student model on the GSM8K dataset. 142
143 Intervention at the first step (in blue) proves to be 143
144 most beneficial with maximum gains over baseline 144
145 without any interventions (dotted line). The gains 145
146 go down when the intervention is done at “step 2” 146
147 (in orange) and intervening at later stages (“step 147
148 3” and beyond) leads to diminishing returns. This 148
149 is because students may have already internalized 149

Student		Teacher					
Model	Type	No	LLaMA 13B	LLaMA 70B	ChatGPT	GPT-4	Human
LLaMA 7B	Pre-Trained	10.53	14.86	19.48	21.00	23.27	22.74
LLaMA 7B	Fine-Tuned	34.19	38.26	45.90	45.94	47.61	47.15
LLaMA 13B	Pre-Trained	24.70	-	26.39	33.24	35.75	33.75
LLaMA 13B	Fine-Tuned	46.24	-	55.34	59.28	60.50	61.86
LLaMA 70B	Pre-Trained	58.90	-	-	63.53	67.40	67.55
LLaMA 70B	Fine-Tuned	63.30	-	-	70.05	72.32	74.14
Mistral 7B	Pre-Trained	40.25	-	46.17	48.82	49.50	50.34
MetaMath 7B	Pre-Trained	62.69	-	61.48	66.94	69.52	65.95
MetaMath 13B	Pre-Trained	67.85	-	66.79	71.41	75.51	73.00

Table 1: Accuracy comparison for different configurations of the student and the teacher models. No refers to *no intervention* by any teacher. Best results are presented in **bold**.

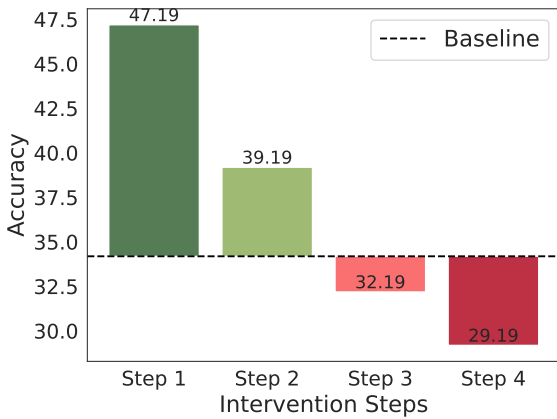


Figure 2: Accuracy of LLaMA 7B fine-tuned student model on GSM8K dataset with correct intervention at different steps by humans. The baseline is the represented by the dotted line with an accuracy of 34.19.

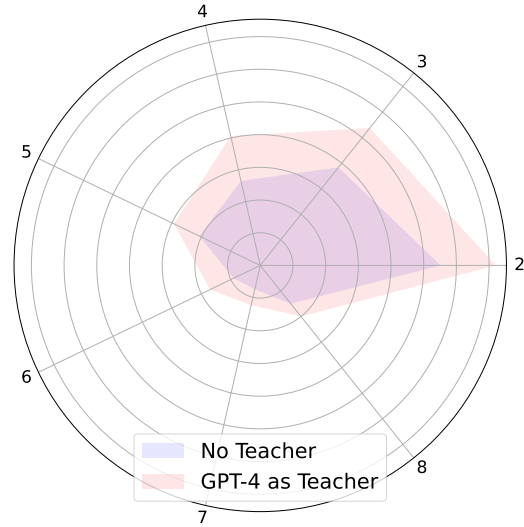


Figure 3: Accuracy comparison for No teacher vs GPT-4 as teacher for LLaMA 7B fine-tuned model across different steps for GSM8K dataset. 2-8 represents the number of steps needed to solve the problem.

incorrect concepts or approaches, making it harder to correct and replace them with the correct ones. This underscores the necessity of early and accurate guidance.

Starting right: first step to successful reasoning

Table 1 demonstrates the usefulness of a teacher guiding a student with first-step guidance. LLaMA 7B pre-trained model’s performance improves by 40% (10.5 \rightarrow 14.8) when a slightly larger model (LLaMA 13B) acts as a teacher to provide first-step guidance. The performance increases monotonically with larger and expert models (10.5 \rightarrow 19.5 with LLaMA 70B as teacher and to 21.0 with ChatGPT as teacher) and goes beyond 100% improvement with GPT-4 (10.5 \rightarrow 23.2) or Human teachers (10.5 \rightarrow 22.7). This demonstrates the importance of first-step guidance in reaching to the correct answer. Similar results are observed for LLaMA 13B as a student where the performance goes from

24.70 with no guidance to 35.75 with GPT-4 first step guidance. The trend follows similarly for even a relatively larger student model (LLaMA 70B) where the performance goes up from 58.9 without guidance to 67.4 with GPT-4 guidance.

Next, we fine-tune all the student models on the training data to see if first step guidance can be learned during training. Although fine-tuning improves the baseline performance significantly, the guidance still proved useful in boosting the performance for all models. LLaMA 7B fine-tuned model’s performance improves by 40% (34.1 \rightarrow 47.6) with expert guidance from GPT-4 as the teacher. A qualitative example is presented in Figure 4 where if LLaMA 7B model is guided to start correctly, it will to the correct answer. The trend is

Question Gene is sewing a quilt out of old souvenir t-shirts. He has one shirt from each vacation he has been on. Every shirt is its own quilt block. Each row is made of blocks from a different year of vacations. He goes on four vacations a year and has been vacationing since he was 23 years old. He is now 34. How many quilt blocks does he have in total?

No Guidance He has been on 4 vacations a year for 34 years so he has been on $4 \times 34 = \ll 4 \times 34 = 136 \gg$ 136 vacations. Each vacation is a quilt block so he has 136 quilt blocks. Each quilt block is one shirt so he has 136 shirts. The answer is 136. ✗

GPT4 Guidance Gene has been on $34 - 23 = \ll 34 - 23 = 11 \gg$ 11 years of vacations, so he has 11 rows of blocks. Each row is made of 4 blocks, so he has $11 \times 4 = \ll 11 \times 4 = 44 \gg$ 44 blocks. The answer is 44. ✓

Figure 4: Example demonstrating the usefulness of first step guidance.

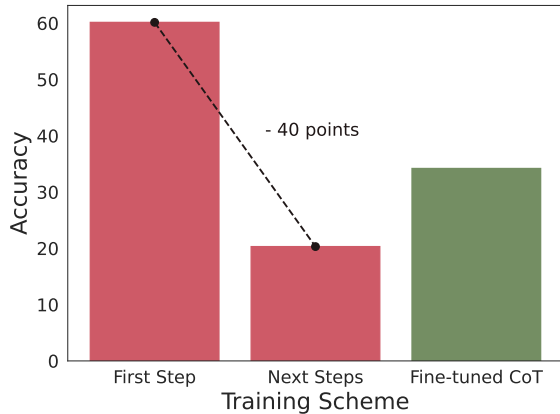


Figure 5: Comparison of curriculum style training for a 7B student model (first step training followed by next steps training) vs fine-tuning (Fine-tuned CoT) on GSM8K dataset.

similar for LLaMA 13B fine-tuned version with a gain of 30% with GPT-4 guidance (46.2 → 60.5) and LLaMA 70B goes up by 15% (63.3 → 72.3).

Finally, we test the applicability of first-step guidance across better student models: Mistral 7B and MetaMath 7B and 13B. Mistral 7B achieves a performance boost of 25% with GPT-4 as teacher (40.25 → 49.50) while MetaMath 7B and 13B gain more than 10% each (62.69 → 69.52 and 67.85 → 75.51 respectively). In all the cases above, it is worth noting that the guidance of GPT-4 is very close to human’s guidance and in many cases surpasses it. This demonstrates the capabilities of GPT-4 as an alternative to teachers in educational domains.

Starting right helps even for longer reasoning chains Figure 3 demonstrates the performance of LLaMA 7B model with and without first step guidance from a teacher for different steps in GSM8K dataset. Across all steps (2 to 8), guidance improves the performance and suggests that starting with a solid foundation can help over longer context. However, the improvement is higher for prob-

lems with 2 to 5 steps compared to 6 to 8, suggesting that over a longer reasoning chain, the chances of making mistakes increase with the increase in the number of steps required to solve it. Nonetheless, starting right has a positive impact on longer reasoning chains too.

Can smaller models be aligned to start better?

Each problem can be broken down into a first step and the next consecutive steps, where the first step can serve as a guidance for the consecutive steps. We train the student model to first learn the initial step and then fine-tune it further to learn the next steps required to solve the problem. This two-step training mechanism has similarity with curriculum learning (Platanios et al., 2019; Xu et al., 2020) where the simpler first step is learnt first, followed by the subsequent more difficult steps. Figure 5 shows a drop of 40 points once the next steps are learned and overall performance gets worse than learning all steps at once in a fine-tuning style. Since the first-step accuracy is close to 60%, only 3/5 samples get the correct guidance during the next steps training and we suspect this might be the reason for a lower overall performance.

5 Conclusion

Distilling reasoning capabilities in smaller models is a challenging task due to their limited abilities to learn complex reasoning strategies. To make these skills more accessible to smaller models, we present an effective way of first-step guidance, where LLMs can guide smaller models in the right direction to solve a reasoning task. On a multi-step reasoning dataset, we show the importance of starting right with a performance improvement of over 100%. Finally, our experiments reveal the quality of guidance which the monotonically increases with the size of the expertise of the teacher model.

6 Limitations

Our work has been tested on one multi-step mathematical reasoning dataset, and while the method can be extended to other reasoning datasets, we have not explicitly tested this in this work. LLMs in general are vulnerable to adversarial attacks and are often very sensitive to hyperparameter changes. We do not see any real-world application of our work.

References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

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](#).

Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander J. Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. [Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes](#). *ArXiv*, abs/2305.02301.

Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L’elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *ArXiv*, abs/2310.06825.

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared

Kaplan. 2022. [Language models \(mostly\) know what they know](#). 300
301

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. [Large language models are zero-shot reasoners](#). 302
303
304

Michal Kosinski. 2023. [Theory of mind might have spontaneously emerged in large language models](#). 305
306

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. [Self-refine: Iterative refinement with self-feedback](#). 307
308
309
310
311
312
313

Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adamek, Eric Malmi, and Aliaksei Severyn. 2023. [Teaching small language models to reason](#). 314
315
316

OpenAI. 2023. [Gpt-4 technical report](#). 317

Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabas Poczos, and Tom Mitchell. 2019. [Competence-based curriculum learning for neural machine translation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1162–1172, Minneapolis, Minnesota. Association for Computational Linguistics. 318
319
320
321
322
323
324
325
326

Swarnadeep Saha, Peter Hase, and Mohit Bansal. 2023. [Can language models teach? teacher explanations improve student performance via personalization](#). In *Thirty-seventh Conference on Neural Information Processing Systems*. 327
328
329
330
331

Kumar Shridhar, Harsh Jhamtani, Hao Fang, Benjamin Van Durme, Jason Eisner, and Patrick Xia. 2023a. [Screws: A modular framework for reasoning with revisions](#). 332
333
334
335

Kumar Shridhar, Jakub Macina, Mennatallah El-Assady, Tanmay Sinha, Manu Kapur, and Mrinmaya Sachan. 2022. [Automatic generation of socratic subquestions for teaching math word problems](#). 336
337
338
339

Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. 2023b. [Distilling reasoning capabilities into smaller language models](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7059–7073, Toronto, Canada. Association for Computational Linguistics. 340
341
342
343
344
345

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor

355	Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>ArXiv</i> , abs/2307.09288.	410
356		411
357		412
358		413
359		414
360		
361		415
362		416
363		417
364		418
365		
366		
367		
368		
369	Janneke Van de Pol, Monique Volman, Frans Oort, and Jos Beishuizen. 2015. The effects of scaffolding in the classroom: support contingency and student independent working time in relation to student achievement, task effort and appreciation of support. <i>Instructional Science</i> , 43:615–641.	
370		
371		
372		
373		
374		
375	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models .	
376		
377		
378		
379	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models . In <i>Advances in Neural Information Processing Systems</i> .	
380		
381		
382		
383		
384	Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khashabi, and Yejin Choi. 2022. Generating sequences by learning to self-correct .	
385		
386		
387		
388	David Wood, Jerome S Bruner, and Gail Ross. 1976. The role of tutoring in problem solving. <i>Child Psychology & Psychiatry & Allied Disciplines</i> .	
389		
390		
391	Benfeng Xu, Licheng Zhang, Zhendong Mao, Quan Wang, Hongtao Xie, and Yongdong Zhang. 2020. Curriculum learning for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6095–6104, Online. Association for Computational Linguistics.	
392		
393		
394		
395		
396		
397		
398	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2023. Large language models as optimizers .	
399		
400		
401	Long Long Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zheng Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models . <i>ArXiv</i> , abs/2309.12284.	
402		
403		
404		
405		
406	Zheng Yuan, Hongyi Yuan, Cheng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. 2023. Scaling relationship on learning mathematical reasoning with large language models . <i>ArXiv</i> , abs/2308.01825.	
407		
408		
409		
	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and Ed Chi. 2023. Least-to-most prompting enables complex reasoning in large language models .	
	Hattie Zhou, Azade Nova, Hugo Larochelle, Aaron Courville, Behnam Neyshabur, and Hanie Sedghi. 2022. Teaching algorithmic reasoning via in-context learning .	