The Self-Improvement Paradox: Can Language Models Bootstrap Reasoning Capabilities without External Scaffolding?

Anonymous ACL submission

Abstract

Self-improving large language models (LLMs) - i.e., to improve the performance of an LLM by fine-tuning it with synthetic data generated by itself - is a promising way to advance the capabilities of LLMs while avoiding extensive supervision. Existing approaches to selfimprovement often rely on external supervision signals in the form of seed data and/or assistance from third-party models. This paper presents CRESCENT - a simple yet effective framework for generating high-quality synthetic question-answer data in a fully autonomous manner. CRESCENT first elicits the LLM to generate raw questions via a bait prompt, then diversifies these questions leveraging a rejection sampling-based selfdeduplication, and finally feeds the questions to the LLM and collects the corresponding answers by means of majority voting. We show that CRESCENT sheds light on the potential of true self-improvement with zero external supervision signals for math reasoning; in particular, **CRESCENT-generated** question-answer pairs suffice to (i) improve the reasoning capabilities of an LLM while preserving its general performance (especially in the 0-shot setting); and (ii) distil LLM knowledge to weaker models more effectively than existing methods based on seed-dataset augmentation.

1 Introduction

001

017

037

041

In recent years, large language models (LLMs) such as GPT-40 (Hurst et al., 2024), Gemini (Anil et al., 2023), Llama (Touvron et al., 2023a), and DeepSeek-R1 (Guo et al., 2025) have demonstrated remarkable capabilities, revolutionizing natural language processing and various other tasks. The success of these models can be attributed to the scaling laws (Kaplan et al., 2020), which dictate the relationship between model parameters, computational resources, and training data size. For instance, the prominent performance of Llama-3.1 with 405B



Figure 1: Different schemes of self-improvement.

042

043

045

046

047

051

052

055

058

060

061

062

063

064

065

066

parameters (Dubey et al., 2024) roots in, amongst others, the massive, high-quality datasets for preand post-training. However, as models continue to scale, the available real-world (public) data quickly becomes exhausted; meanwhile, manually crafting high-quality data is time- and labor-intensive. Thus, data volume has become a key limiting factor for the effective scaling of new-generation models.

In response to this challenge, synthetic data generation and data augmentation have emerged as key methods to further improve the performance of LLMs while avoiding extensive supervision. These methods leverage the ability of LLMs to mirror real-world distributions and generate high-quality, pseudo-realistic data (Zhang et al., 2023). Following this line of research, the problem of *self*improvement naturally arises: Can we improve the performance of an LLM by fine-tuning it with synthetic data generated by itself? This problem has triggered a recent surge of research results (Wang et al., 2024). These methods, however, rely heavily on external seed datasets for augmentation (e.g., (Huang et al., 2023; Wang et al., 2023b)) and/or stronger third-party models as classifiers or reward agents (e.g., (Le et al., 2022; Xin et al.,

2024)); see Fig. 1. Such dependency on external 067 supervision signals limits their ability to achieve 068 true self-improvement. Orthogonally, the recently 069 proposed method Magpie (Xu et al., 2024) suffices to generate high-quality dialogue datasets (i.e., both responses and instructions) entirely through 072 the model itself. Nonetheless, the generated data is highly randomized and primarily dedicated to the alignment of base LLMs. Such data may improve instruction-following abilities but will degrade fundamental capabilities like math and rea-077 soning; see (Xu et al., 2024, Sect. 6). Recent discussions (Kambhampati et al., 2024; Shumailov et al., 2024) have explicitly questioned whether genuine self-improvement is feasible, suggesting that when trained solely on self-generated data, LLMs may fail. Can LLMs achieve true self-improvement? remains an open question in the literature.

This paper aims to provide the infrastructure to explore the self-improvement problem of LLMs: We present CRESCENT -a fully autonomous framework for generating high-quality synthetic question-answer (QA) data that suffice to improve the reasoning capabilities of an LLM while preserving its general performance. CRESCENT adopts a simple yet effective workflow: (i) It uses a bait prompt to guide the model to generate raw questions in a specific domain, such as math word problems; (ii) It applies a self-deduplication mechanism based on rejection sampling (Liu and Liu, 2001) to refine and diversify the question pool; and (iii) For each question, it performs majority voting (Wang et al., 2023a) to identify the most confident answer from the model (thus enhancing the consensus). The so-obtained QA pairs are then used to fine-tune the original LLM via, e.g., supervised fine-tuning (SFT), to improve its math-reasoning capability.

090

097

100 101

102

103

105

106

107

108

110

111

112

113 114

115

116

117

118

Experiments with CRESCENT demonstrate evident self-improvement of LLMs consistently for three benchmarks on mathematical word problems in both 0-shot and 5-shot settings, without trading off their general capabilities. The improvement is especially prominent for the 0-shot case, thus improving the generalization ability of the model to real-world tasks. Ablation studies further demonstrate the superiority of CRESCENT over Magpie (Xu et al., 2024) in the generation of themed data: the latter tends to generate math-related dialogues, e.g., "Could you tell me what type of mathematics you like?" – rather than proper mathematical problems. Moreover, our experiments show that CRESCENT can serve as a highly effective and efficient distillation method, surpassing the baselines using external data and stronger models.

Contributions. Our main contributions include: 121

119

120

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

- We present a simple yet effective framework CRESCENT – utilizing the techniques of bait prompting, diversification, and consensus enhancement – to investigate the selfimprovement problem of LLMs.
- We show that CRESCENT-generated QA pairs suffice to improve the reasoning capabilities of an LLM with zero supervision signals while preserving its general performance, thereby providing an affirmative answer to the selfimprovement problem in the domain of mathematical reasoning (math word problems).
- Experiments demonstrate significant improvements achieved by CRESCENT compared to multiple prompting methods. As a by-product, we show CRESCENT facilitates more effective LLM knowledge distillation than existing approaches based on seed-dataset augmentation.

2 The CRESCENT Approach

This section presents CRESCENT – a framework for <u>controlled QA self-generation via diversification</u> and <u>consensus enhancement</u>. CRESCENT suffices to generate high-quality domain-specific QA pairs leveraging only the model itself, with zero external data, nor assistance from third-party models.

Fig. 2 sketches the general workflow of CRES-CENT, which consists of three main steps: (I) Bait prompting: We use a bait prompt to instruct the original, aligned LLM to produce a set of raw questions within a specific domain; (II) Diversification: The raw questions may be semantically analogous to each other (as per some similarity metric), and thus we employ a rejection sampling mechanism to attain a diverse pool of representative questions through self-deduplication; (III) Consensus enhancement: We treat the generated questions as query prompts and feed them back to the LLM. Then, by majority vote, we obtain the final set of synthetic QA pairs. We show that such QA pairs are of high quality in the sense that they suffice to improve the domain-specific capabilities (mathematical reasoning, in our case) by fine-tuning the original LLM with these QA pairs while preserving its general capabilities.



Figure 2: The general workflow of CRESCENT in mathematical reasoning.

Below, we first present the technical details of Steps (I) to (III) and then provide the rationale behind the self-improvement achieved by these steps.

166

169

170

171

172

173

175

176

178

179

181

182

185

2.1 Question Generation (Steps (I) and (II))

We begin by utilizing a simple *bait prompt* to elicit the LLM to generate a bunch of domain-specific questions, such as math word problems illustrated in Fig. 2, denoted as *raw questions*. As some of them may be semantically analogous to each other, we optimize diversity of the questions in an iterative manner: Each generated question is vectorized and compared against the (embeddings of) other questions. If there exists a question that is deemed sufficiently similar (i.e., the similarity score is below a prescribed threshold), we apply the following *deduplication prompt* to modify it:

{question} is very similar to {question}, please modify the latter to make it different.

This iterative process ensures that the question pool remains diverse and representative across the specific domain through redundancy-aware selection.

Formally, the question-generation phase can be 187 described as follows: Let $Q = \{q_1, q_2, \dots, q_n\}$ be 188 the set of raw questions generated by the LLM per the bait prompt. For each question q_i , we embed 190 it as a real-valued vector v_i and compare it against 191 the vector representations $\{v_1, v_2, \ldots, v_{i-1}\}$ of the previously generated questions. The similarity be-194 tween the two questions is determined by the distance between their respective vector embeddings in the inner product space, e.g., the L^2 distance. If 196 the distance is below a given threshold θ , then q_i with (i > j) is considered as a *duplicate* and thus 198

needs to be modified via the deduplication prompt, i.e.,

If $d(v_i, v_j) < \theta$ then $q_i^* = \text{Deduplicate}(q_i)$. (†)

199

200

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

222

223

224

225

226

229

Such similarity-based deduplication incorporates the *maximal marginal relevance* (MMR) criterion (Carbonell and Goldstein, 1998) to minimize repetition while preserving content relevance. Moreover, the iterative refining process falls into the paradigm of *rejection sampling* (cf. e.g., (Liu and Liu, 2001)), which ultimately yields a diversified question pool featuring relevance and representativeness w.r.t. the target domain with negligible redundancy; see Section 2.3.

2.2 Answer Generation (Step (III))

Let $Q^* = \{q_1^*, q_2^*, \dots, q_n^*\}$ be the deduplicated set of questions generated through the previous step. The phase of answer generation aims to synthesize the corresponding high-quality answers w.r.t. each $q_i^* \in Q^*$. We achieve this by means of *consensus enhancement*, namely, we feed each question q_i^* back to the LLM and collect *m independently* produced answers, denoted by the set $A_i = \{a_1, a_2, \dots, a_m\}$, where each a_j contains integrated chain-of-thought (CoT) processes (Wei et al., 2022) generated for question q_i^* . We then select the final answer a_i^* for question q_i^* using *majority voting* (Wang et al., 2023a). That is, we first identify the set \overline{A}_i of *most frequent answers*:

$$\bar{A}_{i} \triangleq \left\{ a_{j} \in A_{i} \mid f\left(a_{j}\right) = \max_{a_{k} \in A_{i}} f\left(a_{k}\right) \right\} , \qquad 227$$

where $f(a_j)$ denotes the *frequency* (i.e., the number of occurrences) of answer a_j in A_i . Then, we



Figure 3: The intuition of CRESCENT. Let the black dots be question embeddings and distribution curve be conditional answer distribution. (1) Our diversification step modifies question samples violating the minimal distance criterion per (†) (the middle plot). (2) the consensus enhancement step selects the majority mode answer. (the green X in the left and right plots.)

uniformly sample an answer from A_i as the final answer a_i^* paired with question q_i^* . By repeating the majority voting procedure for every question, we obtain the final set of synthetic QA pairs:

$$(Q^*, A^*) = \{(q_1^*, a_1^*), (q_2^*, a_2^*), \dots, (q_n^*, a_n^*)\}.$$

2.3 Rationale for Self-Improvement

230

231

236

240

241

242

243

244

245

247

250

251

Next, we provide the intuition on why selfgenerated QA pairs using the CRESCENT framework can be used to improve the capabilities of the underlying LLM. This observation will be further justified by extensive experiments in Section 3. The intuition is three-fold (see Fig. 3):

- (i) *Relevance by bait prompting*: The initial bait prompt restricts the considered space of questions and answers to a specific domain and hence all the generated QA pairs within the CRESCENT scope are pertinent to this domain.
- (ii) Diversity by rejection sampling-based deduplication: Our diversification step explores the question space while maintaining a minimal pair-wise distance to alleviate redundancy. This is achieved by a rejection sampling loop where question samples violating the distance criterion per (†) are modified and, therefore, the generated questions exhibit a scattered distribution stretching over the space.
- (iii) Accuracy by majority voting: Based on the
 observation that a complex reasoning problem typically admits multiple distinct ways
 of thinking yielding its unique correct answer (Wang et al., 2023a), our consensus enhancement step selects, for each question, the
 most frequent answer that may coincide with
 the correct one with high likelihood.

As a consequence, fine-tuning the original LLM with the so-obtained QA pairs will strengthen its domain-specific capabilities by *enforcing a reduction in the variance of answer generation for a diverse set of domain-relevant questions.* 264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

282

283

284

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

3 Experiments

3.1 Experimental Setups

Benchmarks. We adopt three benchmarks on math word problems (MWPs): (i) **GSM8K** (Cobbe et al., 2021): 8.5K grade school math problems with step-by-step solutions; (ii) **ASDiv** (Miao et al., 2020): 2,305 diverse MWPs covering multiple difficulty levels; and (iii) **GSM-Plus** (Li et al., 2024): an enhanced version of GSM8K with 12K problems incorporating robustness checks. In order to accelerate the evaluation, we use **GSM-Plus-mini** – a subset of GSM-Plus containing 2,400 questions. It should be noted that the GSM-Plus-mini and GSM8K datasets do not overlap.

Baseline Models. We conduct self-improvement experiments with two different LLM models: (i) Llama3-8B-Instruct: the instruction-tuned version of Llama3-8B (Dubey et al., 2024); and (ii) Llama2-7B-Chat (Touvron et al., 2023b): a instruction-tuned version of Llama2-7B.

Generation Configurations. For each model, we generate MWP QA pairs following these settings:

Question Generation: Bait prompt: "Generate a diverse math word problem requiring multistep reasoning". We generate 50K candidate questions for Llama2-7B-Chat and 75k for Llama3-8B-Instruct, both with temperature T = 0.95. Diversification: We use sentence embeddings generated by the all-MiniLM-L6-v2 model from the Sentence-BERT (Reimers and Gurevych, 2019) family; we eliminate semantically similar questions using the L^2 distance with threshold $\theta = 0.25$. We employ FAISS (Douze et al., 2024) to accelerate vector computation and comparisons.

Answer Generation: For each question, sample 5 answers with temperature T = 0.95, then select the most frequent answer as the final answer. We use the same answer generation settings for both models. We use the vLLM (Kwon et al., 2023) inference framework for both generation stages.

GPU hours: It took 30.0 GPU hours to generate 75k QA pairs with Llama3-8B-Instruct and 42.9 GPU hours for the 50k pairs with Llama2-7B-Chat.

Table 1: Main results comparing original models vs. CRESCENT versions. Best results in **bold** (accuracy %).

Model	Training		0-shot		5-shot		
		GSM8K	ASDiv	GSM+	GSM8K	ASDiv	GSM+
Llama2-7B-Chat	Original	18.8	41.7	11.3	23.0	45.9	13.5
	CRESCENT	23.2	46.0	13.0	25.1	45.2	14.8
Llama3-8B-Inst.	Original	34.5	43.6	23.1	75.8	62.3	51.2
	CRESCENT	63.3	65.9	48.6	77.6	63.8	52.8

SFT Implementation. Our SFT procedure uses 312 single-epoch training with max sequence length of 2,048 tokens. Optimization is performed using AdamW (Loshchilov and Hutter, 2019) ($\beta_1 =$ $0.9, \beta_2 = 0.95$) under a linear learning rate schedule (initial LR = 1e-5, 3% warm-up), and the batch 318 size is set to 128 through 8-way parallelization on NVIDIA A100-80GB GPUs with 16-step gradient accumulation. We use DeepSpeed Stage3 (Rasley et al., 2020) and bfloat16 for mitigating memory constraints, and FlashAttention-2 (Dao, 2024) for efficient attention computation.

> **Evaluation Protocol.** We use LM-Evaluation-Harness (Gao et al., 2024) library; all datasets are evaluated under **0-shot** and **5-shot** settings. Fewshot examples are randomly selected from training sets, excluding test samples. We use two answer extractors: one identifies the number appearing after "####" and the other extracts the last number in the output. An answer is considered correct if either of the extractors retrieves the correct answer.

3.2 **Main Results**

313

314

315

317

319

321

322

323

324

326

328

330

331

334

336

337

338

340

341

342

352

The experimental results shown in Table 1 validate our core hypothesis: self-generated reasoning QA pairs – boosted through diversification and consensus enhancement – enable model improvement without external supervision signals. For GSM8K, Llama2-7B-Chat shows improvements of +4.4% (0-shot) and +2.1% (5-shot), while Llama3-8B-Instruct achieves noticeable gains of +28.8% (0shot) and +1.8%[†] (5-shot). Similar observations apply consistently to ASDiv and GSM-Plus-mini featuring different QA distributions.

It is noteworthy that CRESCENT leads to substantial improvements in the 0-shot setting across all three datasets, with performance on certain datasets surpassing even the 5-shot counterparts for the original models. This observation highlights the potential of 0-shot learning in reducing dependency on task-specific examples, thus indicating better generalization to real-world unseen problem types.



Figure 4: Accuracies w.r.t. the ablation study.

354

357

358

359

360

361

362

363

364

365

366

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

384

387

388

389

390

392

3.3 Ablation Study

To justify the pivotality of CRESCENT's core components, we conduct comprehensive ablation experiments over Llama3-8B-Instruct under 5-shot GSM8K evaluation. As depicted in Fig. 4, (i) full method of CRESCENT achieves accuracy of 77.6%, outperforming all ablated variants and the baseline; (ii) removing consensus enhancement (w/o CE) reduces performance to 73.0% (-4.6%); (iii) excluding diversification (w/o DV) yields a more severe drop to 71.1% (-6.53%); (iv) using only bait prompting (BP only) results in 70.6% (-7.0%). The results demonstrate the significance of both diversification and consensus enhancement.

Notably, CRESCENT surpasses the Magpie variants by substantial margins: (i) +5.6% over Magpie-Common (Magpie-C) (72.0%); (ii) +11.0% over Magpie-Math (Magpie-M) (66.6%).

To investigate the discrepancy between CRES-CENT and Magpie-Math, we conduct a sampling analysis on the mathematical questions generated by CRESCENT, CRESCENT w/o DV, and Magpie-Math: For each method, we randomly sample 1,500 questions; Each question is then classified by difficulty using GPT-40 (Hurst et al., 2024), vectorized with the all-MiniLM-L6-v2 embedding model, and projected into a two-dimensional plane using t-SNE (Van der Maaten and Hinton, 2008). The visualization in Fig. 5 suggests that, even without diversification, CRESCENT can still generate highquality mathematical questions, albeit with reduced diversity and difficulty (Fig. 5b). In contrast, the vectors for Magpie-Math problems (Section 3.3) feature (i) a more agglomerate form exhibiting significantly low coverage than CRESCENT; and (ii) numerous gray points signifying non-mathematical problems; they are merely instructions related to the mathematics topic, e.g., "Could you tell me what type of mathematics you like?". The latter aligns with the observation in (Xu et al., 2024,



Figure 5: T-SNE visualization of synthetic math questions. Points colored from 1 to 9 represent mathematical questions with increasing difficulty; Gray marks math-related questions (rather than actual mathematical problems).

Table 2: General capability before/after CRESCENT (%).

Benchmark	#shots	before	after	Δ
ARC-C	0	52.9	52.3	0.6↓
MMLU	5	65.6	65.9	0.3↑
IFEval	-	50.9	52.5	1.6↑
HellaSwag	5	77.9	77.2	0.7↓
GPQA	0	31.2	31.5	0.3↑

Sect. 6) stating that Magpie-generated dialogues may degrade math and reasoning capabilities.

4 Detailed Analysis of CRESCENT

4.1 General-Capability Preservation

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

Will CRESCENT incur catastrophic forgetting of general capabilities? We address this problem by evaluating Llama3-8B-Instruct before and after CRESCENT on five non-mathematical benchmarks covering commonsense reasoning (ARC-C (Clark et al., 2018), HellaSwag (Zellers et al., 2019)), general knowledge preserving (MMLU (Hendrycks et al., 2021)), instruction following (IFEval (Zhou et al., 2023)), and graduate-level question answering (GPQA (Rein et al., 2023)). We use the CRES-CENT checkpoint directly from Section 3.2.

Table 2 shows that the CRESCENT-enhanced model exhibits performance comparable to that of the original model in all five tasks. This observation reveals that domain-specific self-enhancement through CRESCENT does not compromise general capabilities, a critical advantage over fine-tuning approaches using external data, which often exhibit significant capability trade-offs (Luo et al., 2023).

4.2 Analysis of Corrected Questions

417 Our results show significant improvements in the
418 O-shot setting. However, does this improvement re419 flect better generalization, or is it due to the lack of
420 formatting constraints in GSM8K's 0-shot evalua421 tion, which can lead to incorrect answer extraction?
422 To investigate, we analyze Llama3-8B-Instruct's 0-



Figure 6: Breakdown of the corrected questions after applying CRESCENT in the 0-shot setting.

shot results before and after applying CRESCENT, focusing on questions that were incorrect before but correct after (**corrected questions**). We use GPT-40 to classify and analyze these errors.

423

424

425

426

497

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

Fig. 6 shows the total number of corrected questions is 453. 390 (86%) of them are due to genuine improvement in mathematical reasoning ability. These corrected questions can be further broken down into the following: (i) Stepwise reasoning: 199 questions (44%) had errors in stepwise reasoning due to variable tracking (113), step sequence issues (47), and missing steps (39); (ii) Mathematical concept: 115 questions (25%) involved fundamental math errors, with 98 attributed to calculation mistakes and 17 to unit conversion failures; (iii) Redundant information: 37 questions (8%) were impacted by irrelevant information in the problem statement; (iv) Logical structure: 19 questions (4%) involved errors in logical reasoning, such as issues with propositions or set operations; (v) Other errors: 20 questions (4%) were due to other miscellaneous error types.

Meanwhile, there are 63 (14%) corrected questions due to a better output format. After finetuning with CRESCENT-generated QA pairs, these

6

Table 3: Comparison with prompting methods (%).

Method	0-shot	5-shot
Standard prompt	34.5	75.8
Standard prompt + SC	37.8	75.6
Random rephrased	36.9	75.8
CoT prompt	43.6	76.0
Optimized prompt	45.1	75.7
CRESCENT + standard CRESCENT + optimized	63.3 69.8	77.6 77.1

questions are correctly answered without generating redundant content, indicating that CRESCENT's high-quality QA data also improves the model's instruction-following capability.

4.3 Comparison with Prompt Engineering

448

449

450

451

452

453

454

455

456

457

458

459

460

461

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

Can prompt techniques achieve a similar performance with CRESCENT? We address this question by comparing CRESCENT-trained LLaMA3-8B-Instruct against five prompting methods: (i) Stan**dard prompt** from Llama3 official repository;¹ (ii) Standard prompt with self-consistency (SC, aka majority voting) following the settings in (Wang et al., 2023a); (iii) Random rephrased utilizes GPT-40 to randomly rephrase the standard prompt five times (where we select the best evaluation result). Considering the answer-extractor failures discussed in Section 4.2, we carefully craft each instruction to control the output format, such as requesting the answer to be placed after "####" or at the end of the output, ensuring that the prompt includes relevant formatting information compatible with our answer extractor when rephrased by GPT-40; (iv) CoT prompt following the settings in (Wei et al., 2022); (v) Optimized prompt by integrating CoT, the best candidate from random rephrased, and the SC process.

The comparison results are reported in Table 3. Overall, 0-shot outcomes demonstrate higher sensitivity to prompt variations compared to 5-shot configurations. For the original model, the optimized prompt achieves optimal performance, improving 0-shot accuracy by 10.6% over standard prompts while exhibiting comparable 5-shot results. However, this result remains *substantially inferior* (-18.2%) to CRESCENT using only standard prompts. Notably, when employing the same optimized prompts, the CRESCENT-enhanced model further improves 0-shot performance by 6.5%.

> The observed performance gap substantiates that the improvements achieved by CRESCENT *can*-

¹https://github.com/meta-llama/llama-cookbook

Table 4: 0-shot robustness w.r.t. rephrased prompts (%).

Method	R	Random rephrased trials					Std a
Method	T1	T2	Т3	'3 T4 T5		Wiedh	510 0
Original	29.9	19.9	28.6	36.9	24.4	27.9	5.69
CRESCENT	64.9	63.3	64.6	67.8	66.1	65.3	1.52
0.8							.
0.7							
0.0 A		/		-			
9 0.5				-	0-shot 0-shot	Baseline	
0.4					5-shot 5-shot	Baseline	
0.3 25	ōk 5	i0k T	75k `raining	100 Data S	< ize	1	50k

Figure 7: Accuracy in terms of synthetic data volume.

not be replicated through prompting techniques. Moreover, in random rephrased experiments (cf. Table 4), CRESCENT demonstrates *superior robustness* across five different prompts, exhibiting consistent performance with 37.4% higher accuracy and much lower standard deviation. This result indicates that CRESCENT not only enhances *domainspecific proficiency*, but also establishes *promptagnostic generalization* in 0-shot scenarios. 488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

4.4 Data Efficiency and Training Dynamics

Next, we investigate the effect of self-improvement in terms of the volume of synthetic data and the number of training epochs.

Data Volume: We perform one epoch of SFT using Llama3-8B-Instruct on CRESCENT data with data volumes of 25k, 50k, 75k, 100k, and 150k; we use the standard prompt for evaluation. As shown in Fig. 7, the model's performance improves consistently from 25k to 75k, but stabilizes between 75k and 150k, suggesting an upper limit to the improvement gained from increasing data volume.

Training Epochs: We perform SFT with Llama3-8B-Instruct on 50k CRESCENT data for 4 epochs. The evaluation is conducted using the standard prompt. Table 5 shows that, in both settings of 0-shot and 5-shot, the model exhibits a steady performance as the number of epochs increases.

4.5 **CRESCENT for Model Distillation**

Next, we explore the potential of using the CRES-CENT-generated data to distil the knowledge of an

Table 5: Accuracy in terms of number of epochs (%).

#epochs	1	2	3	4
0-shot	50.8	60.4	61.1	62.6
5-shot	74.3	75.7	75.3	75.9

Table 6: Comparison of distillation approaches (%).

Method	Teacher data	#Data	Teacher model	Acc (5-shot)	Acc (0-shot)
-	GSM8K	7k	-	38.4	38.4
MetaMath	GSM8K	50k	Llama3-8B-I.	41.7	22.0
ScaleQuest	GSM8K&MATH	50k	Mix	38.9	22.8
MMIQC	Mix	50k	GPT-4	33.7	28.3
CRESCENT	-	50k	Llama3-8B-I.	44.8	30.8

518

519

520

522

523

530

533

535

540

541

542

545

549

551

553

LLM into a weaker model. Specifically, we use 50k data generated by Llama3-8B-Instruct through CRESCENT to perform SFT on Llama2-7B-Chat, with settings inherited from Section 3.2. We compare this approach with the following distillation methods: (i) Directly using the GSM8K training set without external model enhancement, which contains only 7k samples; (ii) MetaMath (Yu et al., 2024): a method bootstraps existing math datasets by rewriting questions from multiple perspectives, generating a new dataset called MetaMathQA. For comparability, we use Llama3-8B-Instruct to generate 50k new QA pairs from GSM8K training set; (iii) ScaleQuest (Ding et al., 2024): a hybrid method combining multiple models, including Qwen2-Math-7B (Yang et al., 2024), DeepSeek-Math7B-RL (Shao et al., 2024), GPT-4o, and InternLM2-7B-Reward (Cai et al., 2024), along with datasets from GSM8K and MATH. We randomly sample 50k QA pairs from their open-source dataset;² (iv) **MMIQC** (Liu et al., 2024): a method leverages GPT-40 to enhance existing GSM8K, MATH and MetaMathQA datasets. We similarly sample 50k QA pairs from their open-source data³.

The results shown in Table 6 demonstrate that CRESCENT outperforms all other approaches that rely on external data or stronger models. This highlights that CRESCENT is an efficient and effective distillation approach, requiring no external datasets, let alone complex interactions with them. Furthermore, this result also suggests that excessive reliance on external data during distillation may limit the quality of the distilled data, in other words, the model inherently features the ability to produce data of higher quality than the seed dataset, but is constrained to merely modifying or enhancing the seed data; CRESCENT, in contrast, unleashes such554ability to achieve self-improvement.555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

600

5 Related Work

Synthetic Data from Scratch: Recent efforts to reduce reliance on external seed data have led to the exploration of generating data from scratch for fine-tuning LLMs. UltraChat (Ding et al., 2023) shows how to generate diverse, high-quality multi-turn conversations without human queries. Magpie (Xu et al., 2024) introduces a self-synthesis method to generate large-scale alignment data by utilizing only pre-defined chat templates. GenQA (Chen et al., 2024a) aims to generate large instruction datasets with minimal human oversight by prompting LLMs to create diverse instruction examples. Note note that these methods primarily focus on *creating alignment data to train the instruction-following capabilities of base models.*

LLM Self-Improvement: Recent methods exploring self-improvement demonstrate the potential of enhancing LLMs' capabilities through selfgenerated feedback. (Huang et al., 2023) demonstrates that LLMs can improve by sampling highconfidence answers from existing high-quality question sets. Similarly, CodeRL (Le et al., 2022) introduces reinforcement learning to program synthesis, where the model receives feedback from unit tests and critic scores from other models, aiming to optimize performance on unseen coding tasks. StaR (Zelikman et al., 2022) leverages small amounts of rationale examples and iteratively refines the reasoning ability through self-generated rationales. SPIN (Chen et al., 2024b) proposes a self-play fine-tuning method, where a model generates its training data from previous iterations.

6 Conclusion

We presented CRESCENT as a simple yet effective framework – leveraging techniques of bait prompting, diversification, and consensus enhancement – for exploring the self-improvement problem of LLMs. We show that CRESCENT suffices to improve the mathematical reasoning capabilities of an LLM with zero supervision signals while preserving its general performance. Moreover, it facilitates more effective and efficient LLM knowledge distillation than existing approaches based on seeddataset augmentation.

²https://huggingface.co/datasets/dyyyyyyy/ ScaleQuest-Math

³https://huggingface.co/datasets/Vivacem/MMIQC

601 Limitations

602 We observe the following limitations of this work:

603**Domain scalability.** Although CRESCENT can604generate a variety of domain-specific datasets, the605experiments in this paper are confined to evaluat-606ing its effectiveness in improving math reasoning607capabilities. Further extensions to other domains608are subject to future work.

Aligned model restriction. CRESCENT is designed for aligned chat models. In this paper, we did not investigate whether the same approach can be used to generate high-quality, domain-specific data for base models without instruction tuning.

614 References

615

616

617

618

619

620

621

622

623

625

626

627

631

632

633

637

638

639

641

643

647

650

651

653

654

- Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy P. Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul Ronald Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. 2023. Gemini: A family of highly capable multimodal models. CoRR, abs/2312.11805.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaxing Li, Jingwen Li, Linyang Li, Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong, Xiaomeng Zhao, and et al. 2024. InternIm2 technical report. CoRR, abs/2403.17297.
 - Jaime G. Carbonell and Jade Goldstein. 1998. The use of mmr, diversity-based reranking for reordering doc-

uments and producing summaries. In *SIGIR*, pages 335–336. ACM.

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

- Jiuhai Chen, Rifaa Qadri, Yuxin Wen, Neel Jain, John Kirchenbauer, Tianyi Zhou, and Tom Goldstein. 2024a. Genqa: Generating millions of instructions from a handful of prompts. *CoRR*, abs/2406.10323.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024b. Self-play fine-tuning converts weak language models to strong language models. In *ICML*. OpenReview.net.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? Try ARC, the AI2 reasoning challenge. *CoRR*, abs/1803.05457.
- 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*.
- Tri Dao. 2024. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations* (*ICLR*).
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. In *EMNLP*, pages 3029–3051. Association for Computational Linguistics.
- Yuyang Ding, Xinyu Shi, Xiaobo Liang, Juntao Li, Qiaoming Zhu, and Min Zhang. 2024. Unleashing reasoning capability of LLMs via scalable question synthesis from scratch. *CoRR*, abs/2410.18693.
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic,

820

821

Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The Llama 3 herd of models. CoRR, abs/2407.21783.

711

713

714

718

719

721

722

726

727

728

729

730

731

732

733

736

737

739 740

741

742

743

744

745

746

747

748

750

751

752

755

759

760

761

- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. A framework for few-shot language model evaluation.
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
 2021. Measuring massive multitask language understanding. In *ICLR*. OpenReview.net.
 - Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023. Large language models can self-improve. In *EMNLP*, pages 1051–1068. Association for Computational Linguistics.
 - Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, et al. 2024. Gpt-4o system card. *CoRR*, abs/2410.21276.
 - Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Saldyt, and Anil Murthy. 2024. Position: LLMs can't plan, but can help planning in LLMmodulo frameworks. In *ICML*. OpenReview.net.
 - Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *CoRR*, abs/2001.08361.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
 Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model

serving with pagedattention. In *Proceedings of the* ACM SIGOPS 29th Symposium on Operating Systems Principles.

- Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu-Hong Hoi. 2022. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. In *NeurIPS*.
- Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. 2024. Gsm-plus: A comprehensive benchmark for evaluating the robustness of LLMs as mathematical problem solvers. In *ACL* (1), pages 2961–2984. Association for Computational Linguistics.
- Haoxiong Liu, Yifan Zhang, Yifan Luo, and Andrew Chi-Chih Yao. 2024. Augmenting math word problems via iterative question composing. *CoRR*, abs/2401.09003.
- Jun S Liu and Jun S Liu. 2001. *Monte Carlo strategies in scientific computing*, volume 10. Springer.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *ICLR (Poster)*. Open-Review.net.
- Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2023. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *CoRR*, abs/2308.08747.
- Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2020. A diverse corpus for evaluating and developing english math word problem solvers. In *ACL*, pages 975–984. Association for Computational Linguistics.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *KDD*, pages 3505– 3506. ACM.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *EMNLP/IJCNLP (1)*, pages 3980–3990. Association for Computational Linguistics.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. GPQA: A graduate-level google-proof q&a benchmark. *CoRR*, abs/2311.12022.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300.
- Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross J. Anderson, and Yarin Gal. 2024. AI models collapse when trained on recursively generated data. *Nat.*, 631(8022):755–759.

822

- 831 832 834 839
- 841
- 845

- 851 852

855

858 859 860

857

- 865
- 869
- 873

874 875

- 876 877
- 878
- 879

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. LLaMA: Open and efficient foundation language models. CoRR, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumva Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem 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, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. CoRR, abs/2307.09288.
 - Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. Journal of machine learning research, 9(11).
 - Ke Wang, Jiahui Zhu, Minjie Ren, Zeming Liu, Shiwei Li, Zongye Zhang, Chenkai Zhang, Xiaoyu Wu, Qiqi Zhan, Qingjie Liu, and Yunhong Wang. 2024. A survey on data synthesis and augmentation for large language models. CoRR, abs/2410.12896.
 - Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023a. Self-consistency improves chain of thought reasoning in language models. In ICLR. OpenReview.net.
 - Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. Self-instruct: Aligning language models with self-generated instructions. In ACL(1), pages 13484–13508. Association for Computational Linguistics.
 - 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 NeurIPS.
 - Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. 2024. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. CoRR, abs/2405.14333.

Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024. MAGPIE: Alignment data synthesis from scratch by prompting aligned LLMs with nothing. CoRR, abs/2406.08464.

881

882

883

884

885

886

887

888

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. CoRR, abs/2409.12122.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. Metamath: Bootstrap your own mathematical guestions for large language models. In ICLR. OpenReview.net.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. 2022. Star: Bootstrapping reasoning with reasoning. In NeurIPS.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In ACL (1), pages 4791-4800. Association for Computational Linguistics.
- Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu Ding, Joshua B. Tenenbaum, and Chuang Gan. 2023. Planning with large language models for code generation. In ICLR. OpenReview.net.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. CoRR, abs/2311.07911.