## Towards A Unified View of Answer Calibration for Multi-Step Reasoning

Anonymous ACL submission

#### Abstract

Large Language Models (LLMs) employing Chain-of-Thought (CoT) prompting have broadened the scope for improving multi-step reasoning capabilities. We generally divide multi-step reasoning into two phases: path generation to generate the reasoning path(s); and answer calibration post-processing the reasoning path(s) to obtain a final answer. However, the existing literature lacks systematic analysis on different answer calibration approaches. In this paper, we summarize the taxonomy of recent answer calibration techniques and break them down into step-level and path-level strategies. We then conduct a thorough evaluation on these strategies from a unified view, systematically scrutinizing steplevel and path-level answer calibration across multiple paths. Experimental results reveal that integrating the dominance of both strategies tends to derive optimal outcomes. Our study holds the potential to illuminate key insights for optimizing multi-step reasoning with answer calibration.

#### 1 Introduction

003

007

014

017

024

034

Chain-of-Thought (CoT) prompting (Wei et al., 2022) has significantly improved multi-step reasoning capabilities of Large Language Models (LLMs) (Zhao et al., 2023b; Qiao et al., 2023). As seen from Figure 1, the process of multi-step reasoning generally contains two primary modules: *reasoning path generation* which generates one or multiple reasoning paths (Fu et al., 2023; Yao et al., 2023b); and *answer calibration* which post-processes the reasoning path(s) to calibrate the initial output (Wang et al., 2023i; Zhao et al., 2023a).

In practice, answer calibration is pluggable and can be integrated into path generation models. The answer calibration framework can be divided into step and path levels, applicable to single or multiple paths, as illustrated in Figure 1. For *steplevel* answer calibration on <u>a single path</u>, the model



For a <u>Single Path</u>: Self-revise the entire path to answer correctly For <u>Multiple Paths</u>: Select the answer that obtains the maximum **\*** Figure 1: Illustration of answer calibration for multistep reasoning with LLM. The methods of step/pathlevel answer calibration for *multiple paths* can employ answer calibration on *a single path* first.

rectifies errors in intermediate-step answers of a generated path (Zhao et al., 2023a). For step-level answer calibration on multiple paths, the model verifies each intermediate-step answer (Weng et al., 2023) or aggregates the correct step answers (Cao, 2023) from multiple paths. For *path-level* answer calibration on a single path, the model revises the entire rationale to obtain the correct answer (Baek et al., 2023). For path-level answer calibration on multiple paths, the model produces a result indicating the consensus of all candidate paths (Wang et al., 2023i; Yoran et al., 2023). As answer calibration can identify and rectify errors in the reasoning path, or even holistically utilize multiple candidate paths, it plays a vital role in multi-step reasoning to ensure a precise, consistent and reliable reasoning process (Pan et al., 2023).

However, we argue that the crucial factors driving the success of answer calibration strategies remain obscure, with a comprehensive systematic

060

analysis still underexplored. To bridge the gap, our study investigates: (1) The specific conditions where answer calibration notably boosts multistep reasoning performance; (2) The strengths and weaknesses of step-level versus path-level answer calibration, and the pathway to attaining optimal performance; (3) The robustness and generalizability of answer calibration strategies.

062

063

064

067

071

072

076

078

100

101

102

103

104

105

107

108

109

110

111

To address these questions, we dissect cuttingedge answer calibration techniques for multi-step reasoning with LLMs, and introduce a unified framework that elucidates step-level and path-level strategies. We define two thresholds to respectively signify the step-level and path-level dominance in the unified framework. We then undertake a comprehensive evaluation of answer calibration strategies, w.r.t. accuracy, faithfulness, informativeness, consistency, and perplexity over steps or paths. Through rigorous experiments on five representative multi-step reasoning tasks involving arithmetic (Ahn et al., 2024) and commonsense, we find that: (1) employing answer calibration can enhance accuracy, with the improvement being more noticeable in zero-shot scenarios ( $\S4.2$ ) and less significant on stronger backbone models ( $\S4.4$ ); (2) The optimal performance of the unified answer calibration strategy typically achieved by synthesizing step-level and path level dominance  $(\S4.3)$ ; (3) path-level answer calibration is more beneficial in improving accuracy, and step-level answer calibration is more effective for mitigating low-quality prompting  $(\S4.5)$ ; (4) answer calibration can improve consistency on arithmetic tasks but weakens faithfulness, informativeness and perplexity on both arithmetic and commonsense tasks ( $\S4.6$ ).

#### 2 Related Work

**Reasoning Path Generation.** Previous methods for reasoning path generation mostly focus on two aspects to improve reasoning process, including refining input query or prompts (*input refinement*) and polishing the reasoning path (*rationale polish*).

As for *input refinement*, Zero-shot CoT (Kojima et al., 2022) and Few-shot CoT (Wei et al., 2022) are classic methods to elicit multi-step reasoning ability of LLMs, with "Let's think step by step" prompts. To decouple planning and execution, Wang et al. (2023g); Sun et al. (2023) devise a plan by prompting and divide and conquer multistep tasks. To enrich prompts, Wang et al. (2023b) leverage structure triples as evidence, Kong et al. (2023) design role-play prompting, and Xu et al. (2023) employ re-reading instructions. Besides, LLM performance can also be affected by prompt complexity (Fu et al., 2023) and formats, such as program (Gao et al., 2023; Chen et al., 2023b; Sel et al., 2023; Jie et al., 2023; Lei and Deng, 2023; Bi et al., 2024) and table (Jin and Lu, 2023). Further, Wang et al. (2023c); Shi et al. (2023); Liang et al. (2023) propose to adaptively utilize prompts. Apart from refining prompts, Xi et al. (2023b) progressively refine the given questions, Wang et al. (2023j) convert semantically-wrapped questions to meta-questions, and Jie and Lu (2023) augment training data with program annotations.

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

In terms of rationale polish, recent work mainly focus on step-aware training (Wang et al., 2023k) and path-level optimization. For step-aware training, Zhang et al. (2023) introduce step-by-step planning and Lee and Kim (2023) recursively tackle intermediate steps; Jiang et al. (2023a) reconstruct the reasoning rationale within prompts by residual connections; Paul et al. (2023) iteratively provide feedback on step answers; Lanchantin et al. (2023) leverage self-notes as intermediate steps and working memory; Li et al. (2023b); Ling et al. (2023); Lightman et al. (2023) propose to verify on intermediate step answers; Li et al. (2023a); Wang et al. (2023e) process step-aware verification by knowledge base retrieval. For path-level optimization, Li and Qiu (2023) enable LLMs to self-improve via pre-thinking and recalling relevant reasoning paths as memory; Wang et al. (2023d); Yue et al. (2023) leverage hybrid rationales in formats of natural language and program. Some work also generate deliberate rationales beyond CoT, such as Tree-of-Thought (Yao et al., 2023b; Long, 2023), Graph-of-Thought (Yao et al., 2023e; Besta et al., 2023) and Hypergraph-of-Thought (HoT) (Yao et al., 2023a).

**Answer Calibration.** Given generated reasoning path(s), answer calibration methods *post-process* the path(s) to calibrate the answer, involving stepor path-level calibration on one or multiple path(s).

*Step-level answer calibration.* Xue et al. (2023); Cao (2023) propose to rectify factual inconsistency and reasoning logic between intermediate steps. Miao et al. (2023); Wu et al. (2024) check the correctness of each intermediate step. Zhao et al. (2023a) post-edit multi-step reasoning paths with external knowledge. Yao et al. (2023c); Hao et al. (2023); Shinn et al. (2023); Yao et al. (2023d); Chen et al. (2023a); Aksitov et al. (2023) draw up a plan and act step by step with LLMs as agents (Wang et al., 2023f; Xi et al., 2023a), encouraging interaction with the environment to provide feedback. Weng et al. (2023); Jiang et al. (2023b) unleash the self-verification ability of LLMs, by forward reasoning and backward verification on intermediate step answers. Zhou et al. (2023) propose code-based self-verification on reasoning steps.

163

164

165

166

168

169

170

171

172

173

174

175

176

177

178

179

181

182

184

185

188

189

190

193

194

195

196

197

198

199

200

204

Path-level answer calibration. Zelikman et al. (2022) present a self-taught reasoner to iteratively generate rationales. Zheng et al. (2023) progressively use the generated answers as hints to make double-check. Mountantonakis and Tzitzikas (2023) enrich generated reasoning paths with hundreds of RDF KGs for fact checking. Baek et al. (2023) iteratively rectify errors in knowledge retrieval and answer generation for knowledgeaugmented LMs. To cultivate the reasoning ability of smaller LMs, Ho et al. (2023); Wang et al. (2023h,l) propose to fine-tune CoT for knowledge distillation. Huang et al. (2022) demonstrate that LLMs can self-improve with high-confidence rationale-augmented answers. Yoran et al. (2023) prompt LLMs to meta-reason over multiple paths. Liu et al. (2023); Madaan et al. (2023) leverage feedback to improve model initial outputs. Wan et al. (2023) adaptively select in-context demonstrations from previous outputs to re-generate answers. Wang et al. (2023i) leverage self-consistency decoding strategy to majority vote on multiple path answers. Aggarwal and Yang (2023) propose adaptive-consistency to reduce sample budget.

#### 3 Comprehensive Analysis of Answer Calibration

#### 3.1 Formulation of Answer Calibration

Given a question denoted as Q and its associated prompt P, we leverage the LLM to generate the result  $\mathcal{R}$ .  $\mathcal{R}$  can either encompass a single reasoning path  $\mathcal{P}$  with an initial answer  $\mathcal{A}$  or multiple reasoning paths  $\mathbb{P} = {\mathcal{P}_i}_{i \in [1,N]}$  with a corresponding answer set  $\mathbb{A} = {\mathcal{A}_i}_{i \in [1,N]}$ . The total number of paths in  $\mathbb{P}$  is N. In this paper, we analyze under the assumption that each reasoning path comprises a maximum of M steps. Paths exceeding M steps are truncated, and those with fewer steps are padded. The intermediate step answers for each reasoning path  $\mathcal{P}_{(i)}$  are represented as  $\{a_j\}_{j\in[1,M]}^{(i)}$ .

210Step-Level Answer Calibration.Given a single211reasoning path  $\mathcal{P}$  with an initial final path answer

 $\mathcal{A}$  and intermediate step answers  $\{a_j\}_{j\in[1,M]}$ , the objective of step-level answer calibration is to rectify any erroneous  $a_j$ , so that deriving the correct  $\mathcal{A}$ . For multiple reasoning paths  $\mathbb{P}$ , step-level answer calibration seeks to either select the reasoning path with the maximum correct intermediate step answers or aggregate the verified correct steps to form the most accurate reasoning path, leading to a correct final path answer. *Self-verification* (Weng et al., 2023) is an effective approach for step-level answer calibration on multiple reasoning paths.

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

229

230

231

232

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

**Path-Level Answer Calibration.** Given a single reasoning path  $\mathcal{P}$  with an initial final path answer  $\mathcal{A}$ , the goal of path-level answer calibration is to revise the wrong  $\mathcal{A}$ . For multiple reasoning paths  $\mathbb{P} = {\mathcal{P}_i}_{i \in [1,N]}$  with corresponding answers  $\mathbb{A} = {\mathcal{A}_i}_{i \in [1,N]}$ , path-level answer calibration is designed to select the reasoning path from  $\mathbb{P}$  with the most consistent answer in  $\mathbb{A}$ . *Self-consistency* (Wang et al., 2023i) is a widely-used efficacious technique for path-level answer calibration on multiple reasoning paths.

#### 3.2 Unified View of Answer Calibration

Considering the advantages of both step-level and path-level answer calibration, we propose to integrate the two strategies on multiple paths. Given the multiple generated reasoning paths  $\mathbb{P} = \{\mathcal{P}_i\}_{i \in [1,N]}$ , we define a unified score  $\mathcal{D}_i$  for each  $\mathcal{P}_i$  (with the final path answer:  $\mathcal{A}_i$  and intermediate step answers:  $\{a_j\}_{i \in [1,M]}^{(i)}$ ):

$$\mathcal{D}_i = \underbrace{\alpha \frac{n_i}{N}}_{path-level} + \underbrace{(1-\alpha) \frac{m_i}{M}}_{step-level} \qquad (1)$$

where  $n_i \in [1, N]$  is the frequency of  $\mathcal{A}_i$  existing in  $\mathbb{A}, m_i \in [0, M]$  is the number of correct intermediate steps in  $\mathcal{P}_i$ , and  $\alpha$  is a hyper-parameter. The final answer is  $\mathcal{A}_{i^*}$  satisfying  $i^* = \underset{i \in [1,N]}{\arg \max(\mathcal{D}_i)}$ .

To better analyze the effects of varying  $\alpha$  in the unified framework, we then define particular choices for  $\alpha$  which we call *step and path level dominant answer calibration*.

**Definition 1.** Step-Level Dominant Answer Calibration: This choice refers to the level of  $\alpha$  at which the step-level score is used as the dominant criterion, with the path-level score given much smaller weight and only serving to break ties when necessary. Specifically, larger  $m_i$  always results

315

316

317

in larger  $\mathcal{D}_i$ , no matter how small  $n_i$  is. We denote this as:  $\forall n_j, n_k \in [1, N]$  and  $m_j, m_k \in [0, M]$ , where  $n_j < n_k$  and  $m_j > m_k$ , the scores  $D_j$  and  $D_k$  should satisfy

$$\alpha \frac{n_j}{N} + (1-\alpha)\frac{m_j}{M} > \alpha \frac{n_k}{N} + (1-\alpha)\frac{m_k}{M}$$

Thus we can obtain

/

254

259

261

262

265

266

 $\alpha < \frac{1}{\frac{M(n_k - n_j)}{N(m_j - m_k)} + 1}$ 

(2)

If Eq (2) is constant, we can infer that

$$\alpha < \min\left(\frac{1}{\frac{M(n_k - n_j)}{N(m_j - m_k)} + 1}\right) = \frac{1}{\frac{M\max(n_k - n_j)}{N\min(m_j - m_k)} + 1}$$
(3)

As  $1 \le n_j < n_k, n_j+n_k \le N$ , and  $0 \le m_k < m_j$ , we can deduce that  $\min(m_j - m_k) = 1, \max(n_k - n_j) = N - 2$ . From the above, we deduce:

$$\alpha < \frac{1}{\frac{M(N-2)}{N} + 1} \tag{4}$$

**Definition 2.** *Path-Level Dominant Answer Calibration:* For this choice,  $D_i$  gives priority to the path-level score, with the step-level score given much smaller weight and only serving to break ties when necessary. Concretely, larger  $n_i$  always conduces larger  $D_i$ , no matter how small  $m_i$  is. We denote this as:  $\forall n_j, n_k \in [1, N]$  and  $m_j, m_k \in [0, M]$ , where  $n_j > n_k$  and  $m_j < m_k$ , the scores  $D_j$  and  $D_k$  should satisfy

$$\alpha \frac{n_j}{N} + (1-\alpha)\frac{m_j}{M} > \alpha \frac{n_k}{N} + (1-\alpha)\frac{m_k}{M}$$

Analogously, we can obtain

 $\alpha > \frac{1}{\frac{M(n_j - n_k)}{N(m_k - m_j)} + 1}$   $\tag{5}$ 

If Eq (5) is constant, we can infer that

$$\alpha > \max\left(\frac{1}{\frac{M(n_j - n_k)}{N(m_k - m_j)} + 1}\right) = \frac{1}{\frac{M\min(n_j - n_k)}{N\max(m_k - m_j)} + 1}$$
(6)

As  $1 \le n_k < n_j$ , and  $0 \le m_j < m_k \le M$ , we deduce that  $\min(n_j - n_k) = 1$ ,  $\max(m_k - m_j) = M - 0 = M$ . From the above, we deduce:

$$\alpha > \frac{1}{\frac{1}{N} + 1} \tag{7}$$

267In general, considering step-level and path-level268answer calibration dominance, we can obtain two269thresholds:  $\frac{1}{\frac{M(N-2)}{N}+1}$  and  $\frac{1}{\frac{1}{N}+1}$ . Note that  $\alpha = 0$ 270and  $\alpha = 1$  are respectively equivalent to the271self-verification and self-consistency strategies.

#### 3.3 Evaluation of Answer Calibration

Calculation of ROSCOE Scores. In addition to the classical evaluation metric: Accuracy, Golovneva et al. (2023) have proposed **ROSCOE**, a suite of metrics for multi-step reasoning, under four perspectives: semantic alignment (ROSCOE-SA), semantic similarity (ROSCOE-SS), logical inference, and (ROSCOE-LI) and language coherence (ROSCOE-LC). Due to space limits, we select some representative scores from ROSCOE as evaluation metrics in the experiments.

Given source ground truth rationale (s) and generated rationale (h) with multiple steps ( $h_i$ ), we calculate five scores (All scores satisfy the principle that larger is better):

(1) Faithfulness<sub>step</sub>  $(h \rightarrow s)$ : To assess whether the model misconstrues the problem statement, or if the reasoning path is too nebulous, irrelevant, or improperly employs input information.

$$\sum_{i=1}^{N} r$$
-align $(h_i \to s)/N$  (8)

where N is the number of steps and r-align is used to measure how well  $h_i \in h$  can be aligned with any one of the steps in the ground truth path s.

(2) Informativeness<sub>path</sub>  $(h \rightarrow s)$ : To measure the level of concordance between the generated path and the source, and if the generated reasoning path is well-grounded with respect to the source.

$$[1 + \cos(h, s)]/2$$
 (9)

where  $\cos(\cdot, \cdot)$  is a function for cosine similarity.

(3) Consistency<sub>steps</sub>  $(h_i \leftrightarrow h_j)$ : To measure logical entailment errors *within* the reasoning steps.

 $1 - \max_{i=2..N} \max_{j < i} p_{\text{contr}}(h_i, h_j)$ (10)

where  $p_{\text{contr}}$  is used to assess the likelihood of step pairs contradicting each other.  $h_i \in h$  and  $h_j \in h$ . (4) Consistency<sub>path</sub> ( $h \leftrightarrow s$ ): To evaluate mistakes in logical entailment between the generated reasoning path h and source context s:

$$1 - \max_{i=1..N} \max_{j=1..T} p_{\text{contr}}(h_i, s_j)$$
 (11)

where  $p_{\text{contr}}$  is the likelihood of source and generated steps contradicting each other.  $s_i \in \mathbf{s}$ ;  $h_i \in \mathbf{h}$ .

(5) Perplexity<sub>path</sub> (h): As an indicator of language coherence, it calculates average perplexity of all tokens in the generated reasoning path steps.

$$1/\text{PPL}(\boldsymbol{h}) \tag{12}$$

where PPL denotes the perplexity.

Task	Method	Accuracy ↑	Faithfulness ↑ (Over Steps)	Informativeness ↑ (Over Path)	Consistency ↑ (Within Steps)	Consistency ↑ (Within I/O)	Perplexity ↑ (Over Path)
GSM8K	CoT CoT + SV CoT + SC	80.21 82.34 <sub>(+2.13)</sub> 87.11 <sub>(+6.90)</sub>	88.73 86.22 <sub>(-2.51)</sub> <b>88.83<sub>(+0.10)</sub></b>	96.38 95.19 <sub>(-1.19)</sub> <b>96.40<sub>(+0.02)~</sub></b>	<b>97.94</b> 96.78 <sub>(-1.16)</sub> 97.90 <sub>(-0.04)~</sub>	$96.94 \\93.46_{(-3.48)} \\97.44_{(+0.50)}$	9.14 <b>14.90</b> (+5.76) 8.90 <sub>(-0.24)</sub>
	ZS CoT ZS CoT + SV ZS CoT + SC	$\begin{array}{c} 62.85\\ 67.70_{(+4.85)}\\ \underline{71.42}_{(+8.57)}\end{array}$	$\frac{86.58}{86.24_{(-0.34)}}$ $\underline{86.70_{(+0.12)}}$	95.61 95.19 <sub>(−0.42)</sub> <u>95.67<sub>(+0.06)</sub>~</u>	$\frac{97.30}{96.78}_{(-0.52)}$ $97.19_{(-0.11)}$	93.07 93.44 <sub>(+0.37)</sub> $94.57_{(+1.50)}$	$\frac{15.67}{14.90}_{(-0.77)}$ $14.95_{(-0.72)}$
	CoT CoT + SV CoT + SC	78.20 <b>85.80</b> (+7.60) 84.40(+6.20)	$\begin{array}{c} \textbf{87.73} \\ \textbf{87.26}_{(-0.47)} \\ \textbf{87.60}_{(-0.13)} \end{array}$	$95.74 \\95.00_{(-0.74)} \\95.71_{(-0.03)}$	30.57 33.39 <sub>(+2.82)</sub> <b>33.51</b> <sub>(+2.94)</sub>	9.82 <b>10.41</b> <sub>(+0.59)</sub> 9.92 <sub>(+0.10)</sub>	$\begin{array}{c} \textbf{6.65} \\ \textbf{6.23}_{(-0.42)} \\ \textbf{6.22}_{(-0.43)} \end{array}$
SVAMP	ZS CoT ZS CoT + SV ZS CoT + SC	$\begin{array}{c c} 72.80\\ 81.20_{(+8.40)}\\ \underline{82.00}_{(+9.20)} \end{array}$	$\frac{\underline{87.46}}{86.92}_{(-0.54)}$ $87.40_{(-0.06)}$	95.77 95.05(−0.72) <u>95.81(+0.04)</u> ~	$\frac{31.71}{\underline{35.27}_{(+3.56)}}$ $\overline{34.73}_{(+3.02)}$	$\frac{18.39}{19.67_{(+1.28)}}$	$\frac{\underline{11.93}}{11.44_{(-0.49)}}\\11.68_{(-0.25)}$
	CoT CoT + SV CoT + SC	97.67 98.33 <sub>(+0.66)</sub> 98.17 <sub>(+0.50)</sub>	$\begin{array}{c} \textbf{88.53} \\ \textbf{88.36}_{(-0.17)} \\ \textbf{88.42}_{(-0.11)} \end{array}$	$94.91 \\ 94.38_{(-0.53)} \\ 94.82_{(-0.09)}$	7.77 <b>46.59</b> <sub>(+38.82)</sub> 10.22 <sub>(+2.45)</sub>	7.47 <b>24.56</b> <sub>(+17.09)</sub> 9.29 <sub>(+1.82)</sub>	5.51 <b>10.54</b> (+5.03) 5.33 <sub>(-0.18)</sub>
MultiArith	ZS CoT ZS CoT + SV ZS CoT + SC	$\begin{array}{c c} 87.00 \\ \underline{97.00}_{(+10.00)} \\ \underline{97.00}_{(+10.00)} \end{array}$	$\frac{89.32}{88.35}_{(-0.97)}$ $89.18_{(-0.14)}$	95.30 94.38 <sub>(−0.92)</sub> <u>95.32<sub>(+0.02)</sub></u> ~	$\frac{47.54}{46.26}_{(-1.28)}$ $47.42_{(-0.12)}$	$24.39 \\ \underline{24.58}_{(+0.19)} \\ 23.83_{(-0.56)}$	$\frac{10.75}{10.54}_{(-0.21)}$ $10.63_{(-0.12)}$
	CoT CoT + SV CoT + SC	52.83 <b>54.74</b> (+1.91) 54.47(+1.64)	<b>85.99</b> 85.93 <sub>(-0.06)</sub> 85.93 <sub>(-0.06)</sub>	$\begin{array}{c} \textbf{95.31} \\ \textbf{95.24}_{(-0.07)} \\ \textbf{95.20}_{(-0.11)} \end{array}$	49.57 51.39 <sub>(+1.82)</sub> <b>51.73</b> <sub>(+2.16)</sub>	23.78 24.61 <sub>(+0.83)</sub> <b>25.03</b> <sub>(+1.25)</sub>	$7.64 \\ 7.18_{(-0.46)} \\ 7.15_{(-0.49)}$
MathQA	ZS CoT ZS CoT + SV ZS CoT + SC	$\begin{vmatrix} 49.45 \\ \underline{52.86}_{(+3.41)} \\ 49.51_{(+0.06)} \end{vmatrix}$	$\frac{85.20}{85.93(+0.73)}$ $\frac{85.22}{(+0.02)} \sim$	$\frac{96.08}{95.24}_{(-0.84)}$ 96.08 $_{(-0.00)}$ ~	$23.50 \\ \underline{51.40}_{(+27.90)} \\ 23.66_{(+0.16)}$	$\frac{13.76}{\underline{24.63}_{(+10.87)}}$ $\overline{13.79}_{(+0.03)}$	13.44 7.19 <sub>(-6.25)</sub> <u>13.48<sub>(+0.04)</sub>~</u>
CSQA	CoT CoT + SV CoT + SC	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	81.40 80.89 <sub>(-0.51)</sub> <b>81.50<sub>(+0.10)</sub></b>	92.57 92.10 $_{(-0.47)}$ 92.71 $_{(+0.14)}$	<b>95.57</b> 92.77 <sub>(-2.80)</sub> 95.04 <sub>(-0.53)</sub>	<b>57.54</b> 56.05 <sub>(-1.49)</sub> 56.97 <sub>(-0.57)</sub>	2.46 2.47 <sub>(+0.01)</sub> ~ 2.43 <sub>(-0.03)</sub>
	ZS CoT ZS CoT + SV ZS CoT + SC	$\begin{array}{c c} 67.57\\ 66.42_{(-1.15)}\\ \underline{71.58}_{(+4.01)}\end{array}$	$\frac{\underline{79.77}}{\overline{79.06}_{(-0.71)}}\\ \underline{79.51}_{(-0.26)}$	$\frac{95.26}{94.65}_{(-0.61)}$ 95.21 $_{(-0.05)\sim}$	$\frac{\underline{25.81}}{\underline{25.36}_{(-0.45)}}$ $\underline{25.08}_{(-0.73)}$	$29.17 \\ 28.56_{(-0.61)} \\ \underline{29.69}_{(+0.52)}$	$\frac{\underline{9.90}}{9.06}_{(-0.84)}$ $8.96_{(-0.94)}$

Table 1: Comprehensive performance (%) with different strategies on GPT-3.5 (gpt-3.5-turbo). CoT: Fewshot CoT (Wei et al., 2022) with complex-prompting (Fu et al., 2023); **ZS-CoT**: Zero-Shot CoT (Kojima et al., 2022); **SV**: Self-Verification (Weng et al., 2023); **SC**: Self-Consistency (Wang et al., 2023i). **Best few-shot results** are marked in **bold**; best *zero-shot* results are <u>underlined</u>. I/O: input/output.  $\uparrow$ : larger is better.  $\sim$ ,  $\sim$ : comparable.

#### 4 Experiments

#### 4.1 Setup

**Evaluation Metrics.** In this paper, we aim to conduct comprehensive evaluation on multi-step reasoning, thus we select some scores from ROSCOE (Golovneva et al., 2023) as introduced in §3.3, which contains a suite of metrics allowing us to evaluate the quality of reasoning rationales, not limited to the correctness of final answers.

**Datasets.** We evaluate on five benchmark datasets involving arithmetic and commonsense multi-step reasoning: **GSM8K** (Cobbe et al., 2021), **SVAMP** (Patel et al., 2021), **MultiArith** (Roy and Roth, 2015), **MathQA** (Amini et al., 2019) and **CSQA** (Talmor et al., 2019).

Models. For reasoning *path generation*, we leverage Zero-shot CoT (ZS CoT) (Kojima et al., 2022) and Few-shot CoT (CoT) (Wei et al., 2022) with complexity-based prompting (Fu et al., 2023). For *answer calibration*, we employ Self-Verification (SV) (Weng et al., 2023) and Self-Consistency (SC) (Wang et al., 2023i) on multiple

paths. SV is a step-level strategy, which verifies intermediate-step answers and returns the path containing the maximum number of correct step answers. SC is a path-level strategy, which conducts majority voting on final answers of all generated paths and selects the most consistent result. **Implementation.** We release the codes and generated results anonymously<sup>1</sup>. In this paper, the number of reasoning paths N defined in Eq (1) is 10, and number of intermediate steps M is 3 on all datasets except for CSQA where M is 10. We utilize GPT-3.5 with gpt-3.5-turbo engine as the backbone LLM to generate reasoning paths (the model choice justification is elaborated in Appendix B), and the temperature is set to 0.7. We also leverage GPT-4 (OpenAI, 2023) with gpt-4 engine to generate ground-truth rationales given the ground-truth answers for all datasets excluding GSM8K (which already contains them). For evaluation referring to ROSCOE (Golovneva et al., 2023), we respectively lever-

<sup>&</sup>lt;sup>1</sup>https://anonymous.4open.science/r/Eval\_Multi-Step\_ Reasoning-4E60.



Figure 2: Accuracy under different integrated *step-level* and *path-level* answer calibration strategies, varying with the values of  $\alpha$  defined in Eq (1). Performance with two thresholds of  $\frac{1}{\frac{M(N-2)}{N-2}+1}$  and  $\frac{1}{\frac{1}{N}+1}$  are marked as  $\bigstar$ .

age all-MiniLM-L6-v2/SentenceTransformer, and pretrained gpt2-large (Radford et al., 2019) to obtain token/sentence embedding and calculate perplexity defined in Eq (12). All the reasoning paths for CoT and ZS CoT were generated during 8th to 23rd June 2023, and answer calibration on the generated reasoning paths was conducted during 12th October to 8th November 2023.

361

374

375

391

#### 4.2 Analysis on Step-Level and Path-Level Answer Calibration Strategies

We respectively incorporate the effective step-level and path-level answer calibration strategies, Self-Verification (SV) and Self-Consistency (SC), into CoT-based models operating on multiple paths. We evaluate their performance using six evaluation metrics, with the results presented in Table 1.

Generally, in terms of accuracy, employing answer calibration is effective. Seen from Table 1, we find that models equipped with SV and SC obviously outperform vanilla methods, as both fewshot and zero-shot CoT employing SV/SC achieve significant accuracy improvements on almost all tasks. Notably, zero-shot CoT with SV and SC achieves much more significant outperformance of accuracy than few-shot settings on almost all tasks, demonstrating that answer calibration is more effective in zero-shot settings. As zero-shot CoT is relatively challenging due to the absence of task-specific in-context learning, answer calibration strategies essentially creating a feedback loop where the model assesses its own performance and adjusts accordingly, could help to mitigate biases and overfitting to specific patterns during inference, allowing the model to better generalize to

new types of problems and datasets.

395

396

397

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

Furthermore, in terms of other metrics, answer calibration can improve consistency on arithmetic tasks but weakens faithfulness, informativeness and perplexity on both arithmetic and commonsense tasks. Observed from Table 1, we find that SV and SC weaken the *perplexity* score (16 out of 20 cases), suggesting that the rationale generated from multiple paths is more complex than that from a single path with CoT models. However, these two strategies improve *consistency* scores on arithmetic tasks (10 out of 16 cases; 14 out of 16 cases), intuitively benefiting from multiple paths. As SV verifies answers for intermediate steps and SC considers answers for all paths, they naturally enhance consistency within steps and between input/output (I/O). Additionally, SV and SC worsen faithfulness and informativeness on almost all tasks (15 out of 20 cases for both). The possible reason is that answer calibration on multiple paths focuses more on answer accuracy, while its increased complexity of its rationales tends to result in lower alignment and concordance between the source content and the output path. Generally, despite the benefits of employing SV and SC to CoT-based methods, the improvements are taskdependent and vary across different metrics.

# 4.3 Analysis on Unified Answer Calibration Strategies

We then integrate step-level and path-level answer calibration strategies, varying  $\alpha$  as defined in Eq (1). We present the accuracy of the unified strategies in Figure 2. As observed, accuracy peaks at a specific value of  $\alpha$  between the two thresholds

4	3	8
4	3	9
4	4	0
4	4	1
4	4	2
4	4	3
4	4	4
4	4	5
4	4	6
4	4	7
4	4	8
4	4	9
4	5	0
4	5	1
л		0
4	0	~
4	5	3
4	5	4
4	5	5
4	5	6
4	5	7
4	5	8
4	5	9
4	6	0

463

464

429

430

431

432

433

434

435

436

437

Engine	Strategy	GSM8K	SVAMP	MultiArith	CSQA
GPT-3(175B) code-davinci-001	CoT CoT + SV CoT + SC CoT + SC + SV	13.84 13.92↑ 23.40↑ 23.59↑	38.42 38.96↑ 54.58↑ 54.68↑	45.85 46.19↑ 79.82↑ 80.01↑	46.75 47.68↑ 54.92↑ 55.09↑
Instruct-GPT (175B) code-davinci-002	CoT CoT + SV CoT + SC CoT + SC + SV	60.81 65.14↑ 78.00↑ 78.32↑	75.87 76.99↑ 86.77 <b>↑</b> 86.94 <b>↑</b>	96.13 99.15↑↑ 100.00↑↑ 100.00↑	77.42 77.83↑ 81.43↑ 81.53↑
GPT-3.5 gpt-3.5-turbo	CoT CoT + SV CoT + SC CoT + SC + SV	80.21 82.34↑ 87.11↑ 88.25↑	78.20 85.80↑ 84.40↑ 86.80↑	97.67 98.33↑ 98.17↑ 99.00↑	74.77 74.04↓ 75.27↑ 75.18↑

Table 2: Accuracy (%) with different backbone engines.  $\uparrow/\uparrow$ : slightly/significantly better;  $\downarrow$ : slightly worse than the baseline few-shot CoT. We refer to Weng et al. (2023) for results with GPT-3 and Instruct-GPT engines. As Weng et al. (2023) didn't test on MathQA dataset, we also exclude the results of MathQA here for fair comparisons.

defined in Eq (4) and (7) in almost all scenarios across all tasks (i.e., 8 out of 10 cases), demonstrating that optimal model performance should balance both step-level and path-level answer calibration dominance. Besides, we notice that for "CoT on SVAMP task" in Figure 2(b) and "zeroshot CoT on MathOA task" Figure 2(d), employing integrated answer calibration strategies reaches a peak with  $\alpha$  not between the two thresholds, and the overall performance remains stably lower than the initial best accuracy with  $\alpha = 0$  (*i.e.*, SV). The possible reason may related to *employing SV* (i.e.,  $\alpha = 0$ ) presenting more significant advantages than SC (i.e.,  $\alpha = 1$ ) in the two scenarios. Specifically, CoT on SVAMP respectively achieves accuracy of 85.80% and 84.40% when  $\alpha$  values 0 (SV) and 1 (SC), with the difference larger than 1%; Zero-shot CoT on MathQA employing SV and SC achieves accuracy of 52.86% v.s. 49.51%, where the difference is larger than 3%. Except for these two distinctive scenarios, others in Figure 2 obtain the optimal results by synthesizing step-level and path level answer calibration dominance.

In conclusion, the value of  $\alpha$  plays a significant role in the performance of both few-shot and zeroshot CoT. Optimal ranges of  $\alpha$  for each task are mostly between the two thresholds of step-level and path-level answer calibration dominance. The marked two thresholds represent boundaries for optimizing performance, which could guide further fine-tuning. Besides, the performance variance across datasets implies that the characteristics of each task, such as complexity, size, or the nature of the tasks. Models equipped with answer calibration strategies may require task-specific tuning to achieve the best performance.

#### 4.4 Effects of Backbone Models

We compare accuracy on CoT-based answer calibration strategies with different LLM backbone engines, and present results in Table 2. 465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

As observed from the results, for GPT-3 and Instruct-GPT, both self-verification (SV) and selfconsistency (SC) provide consistent improvements; while on the larger GPT-3.5 model, their improvements are observably weaker, particularly for SV, with which accuracy even slightly drops on the CSQA task. The possible reason is that GPT-3.5 is more prone to making mistakes when verifying on intermediate-step answers for multiple paths. Further, for integrated answer calibration strategies (SV+SC), the model's performance is close to the better one between SV and SC. Generally, path-level answer calibration is more advantageous than step-level one, with relatively higher accuracy and lower computation cost. Based on these observations, we can infer that **answer calibration** strategies, especially path-level self-consistency, provide benefits in many cases, particularly on less powerful LLMs.

We further speculate, if the path generation for CoT with strong backbone LLM is sophisticated enough, the answer calibration may be simplified. We can directly conduct *path-level* answer calibration for multiple paths. But these findings cannot indicate that step-level answer calibration is meaningless for stronger backbone LLMs. As seen from Table 1, LLM equipped with step-level answer calibration is relatively beneficial to improve consistency scores. Besides, as mentioned in Weng et al. (2023), step-level answer calibration can provide explainable answers by verifying on intermediatestep answers, making results more reliable.



Figure 3: Performance (%) of "Accuracy, Faithfulness (Over Steps) and Informativeness (Over Path)" on SVAMP and MultiArith with different prompting on CoT models. We didn't show full results of other tasks for space limits.

#### 4.5 Effects of Prompting

We further demonstrate the effects of prompting with few-shot demonstrations on answer calibration, evaluated on CoT models.

We respectively input prompts of *no coherence* and no relevance for few-shot CoT referring to Wang et al. (2023a) (examples are listed in Appendix C), and present performance on SVAMP and MultiArith in Figure 3. As seen, the deficiency of coherence and relevance in the prompting observably weaken the performance of all models, with no relevance having a more profound impact than no coherence. In addition, CoT+SV achieves comparable performance with CoT+SC when prompting is standard or not coherent. Further, CoT+SV tends to perform observably better than CoT+SC, when prompting with no relevance, indicating that step-level answer calibration strategy SV, is beneficial to maintain performance under adverse conditions. This observation suggests the robustness of step-level answer calibration. It also highlights the potential benefits of step-level answer calibration strategies to mitigate performance degeneration caused by poor prompting. The possible reason is that step-level answer calibration strategies break down the task into subtasks, and these subtasks are simple enough so that less likely to be influenced by the low-quality prompts.

#### 4.6 Analysis on Tasks

As seen from Table 1,2, and Figure 2, generally, SV and SC present more significant outperformance on arithmetic tasks than on the commonsense task (CSQA). Further, for CSQA, employing answer calibration tends to worsen the consistency scores, which is contrary to the trend observed in arithmetic tasks. The possible explanation lies in the characteristics of each task, such as complexity, size, or the nature of the tasks. In the CSQA task, correct intermediate steps may not always contribute to a coherent reasoning path due to potential irrelevance and redundancy. Specifically, even if we calibrate both intermediate step and path answers, there can be some correct commonsense statements while irrelevant to the question, resulting in worse consistency and perplexity. Conversely, in arithmetic tasks, correct intermediate answers almost guarantee a consistent reasoning path, as all intermediate answers are necessary and will contribute to a correct final answer.

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

554

555

556

557

558

559

560

561

562

563

564

566

567

568

570

#### 5 Conclusion and Future Work

In this paper, we dissect multi-step reasoning into path generation and answer calibration, and provide a unified view of answer calibration strategies through a comprehensive evaluation. We find that path-level answer calibration is particularly potent in improving accuracy, while step-level answer calibration is more suitable for addressing issues related to low-quality prompting. The improvement is more pronounced in zero-shot scenarios and less significant on stronger backbone models. We also define step-level and path-level answer calibration dominance with two thresholds, and propose to integrate of the two types of strategies, which is promising to achieve optimal performance. Our findings suggest that answer calibration is a versatile strategy that can be integrated into various models to bolster multi-step reasoning capabilities of LLMs. In the future, we aim to develop more sophisticated multi-step reasoning models, drawing on the insights and conclusions from this study.

588

589

591

593

594

595

596

598

599

601

610

611

612

613

615

616

## Technical Novelty Emphasis

We have conducted an empirical study of answer 572 calibration and proposed a unified method to ad-573 dress that Step-Level / Path-Level Answer Calibra-574 tion for a Single or Multiple Paths can be integrated 575 together, with the two thresholds of Step-/Path-Level Dominant Answer Calibration and a hyper-577 parameter  $\alpha$ . Our analysis has the potential to inspire further research and practical implications 579 on unified answer calibration, such as "how the 580 hyper-parameter  $\alpha$  can be optimally chosen across different tasks, like iterative tuning". Our paper is 582 based on an empirical study, and its main contributions are to unify multiple seemingly disparate types of approaches into a common framework, allowing us to investigate empirical questions to 586 obtain more insights, such as:

- Employing answer calibration can enhance accuracy, with the improvement being more noticeable in zero-shot scenarios and less significant on stronger backbone models;
- (2) The optimal performance of the unified answer calibration strategy typically achieved by synthesizing step-level and path level dominance;
- (3) Path-level answer calibration is more beneficial in improving accuracy, and step-level answer calibration is more effective for mitigating lowquality prompting;
- (4) Answer calibration can improve consistency on arithmetic tasks but weakens faithfulness, informativeness and perplexity on both arithmetic and commonsense tasks.

#### Limitations

The main limitation for this paper is that we didn't analyze more answer calibration strategies, such as step-/path-level methods on the single path, and varying the numbers of steps and paths in the unified answer calibration strategies. Besides, we can also employ answer calibration strategies to other path generation models, not limited to CoT-based methods. Further, we should also evaluate answer calibration strategies on more tasks to make the results more sufficient.

### 614 References

Aman Madaan Pranjal Aggarwal and Mausam Yiming Yang. 2023. Let's sample step by step: Adaptiveconsistency for efficient reasoning and coding with llms. In *EMNLP*, pages 12375–12396. Association for Computational Linguistics.

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

- Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. Large language models for mathematical reasoning: Progresses and challenges. *CoRR*, abs/2402.00157.
- Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang Li, Sheila Babayan, Kavya Kopparapu, Zachary Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Srinivasan, Manzil Zaheer, Felix Yu, and Sanjiv Kumar. 2023. Rest meets react: Selfimprovement for multi-step reasoning llm agent. *CoRR*, abs/2312.10003.
- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. In *NAACL-HLT (1)*, pages 2357–2367. Association for Computational Linguistics.
- Jinheon Baek, Soyeong Jeong, Minki Kang, Jong Park, and Sung Ju Hwang. 2023. Knowledge-augmented language model verification. In *EMNLP*, pages 1720–1736. Association for Computational Linguistics.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2023. Graph of thoughts: Solving elaborate problems with large language models. *CoRR*, abs/2308.09687.
- Zhen Bi, Ningyu Zhang, Yinuo Jiang, Shumin Deng, Guozhou Zheng, and Huajun Chen. 2024. When do program-of-thoughts work for reasoning? In *AAAI*. AAAI Press.
- Lang Cao. 2023. Enhancing reasoning capabilities of large language models: A graph-based verification approach. *CoRR*, abs/2308.09267.
- Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. 2023a. Fireact: Toward language agent fine-tuning. *CoRR*, abs/2310.05915.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023b. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168.
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. 2023. Complexity-based prompting for multi-step reasoning. In *ICLR*. OpenReview.net.

- 672 673 674
- 6 6
- 6 6 6
- 68 68 68 68
- 68 68
- 69
- 6
- 6
- 6 6 7
- 7
- 704 705 706
- 7 7
- 7
- 712 713
- 714
- 716 717
- 7

- 72
- 72

- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. PAL: program-aided language models. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pages 10764–10799. PMLR.
- Olga Golovneva, Moya Peng Chen, S pencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2023. ROSCOE: A suite of metrics for scoring step-by-step reasoning. In *The Eleventh International Conference on Learning Representations*.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. On calibration of modern neural networks. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 1321–1330. PMLR.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023.
  Reasoning with language model is planning with world model. In *EMNLP*, pages 8154–8173. Association for Computational Linguistics.
- Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023. Large language models are reasoning teachers. In *ACL (1)*, pages 14852–14882. Association for Computational Linguistics.
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *CoRR*, abs/2210.11610.
- Song Jiang, Zahra Shakeri, Aaron Chan, Maziar Sanjabi, Hamed Firooz, Yinglong Xia, Bugra Akyildiz, Yizhou Sun, Jinchao Li, Qifan Wang, and Asli Celikyilmaz. 2023a. Resprompt: Residual connection prompting advances multi-step reasoning in large language models. *CoRR*, abs/2310.04743.
- Weisen Jiang, Han Shi, Longhui Yu, Zhengying Liu, Yu Zhang, Zhenguo Li, and James T. Kwok. 2023b. Forward-backward reasoning in large language models for verification. *CoRR*, abs/2308.07758.
- Zhanming Jie and Wei Lu. 2023. Leveraging training data in few-shot prompting for numerical reasoning. In ACL (Findings), pages 10518–10526. Association for Computational Linguistics.
- Zhanming Jie, Trung Quoc Luong, Xinbo Zhang, Xiaoran Jin, and Hang Li. 2023. Design of chain-of-thought in math problem solving. *CoRR*, abs/2309.11054.
- Ziqi Jin and Wei Lu. 2023. Tab-cot: Zero-shot tabular chain of thought. In *ACL (Findings)*, pages 10259– 10277. Association for Computational Linguistics.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *NeurIPS*.

Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong Qin, Ruiqi Sun, and Xin Zhou. 2023. Better zero-shot reasoning with role-play prompting. *CoRR*, abs/2308.07702. 725

726

727

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

753

754

756

757

759

760

761

762

763

764

765

766

767

768

769

770

771

773

- Jack Lanchantin, Shubham Toshniwal, Jason Weston, Arthur Szlam, and Sainbayar Sukhbaatar. 2023. Learning to reason and memorize with self-notes. *CoRR*, abs/2305.00833.
- Soochan Lee and Gunhee Kim. 2023. Recursion of thought: A divide-and-conquer approach to multicontext reasoning with language models. In *ACL* (*Findings*), pages 623–658. Association for Computational Linguistics.
- IokTong Lei and Zhidong Deng. 2023. Selfzcot: a self-prompt zero-shot cot from semantic-level to code-level for a better utilization of llms. *CoRR*, abs/2305.11461.
- Xiaonan Li and Xipeng Qiu. 2023. Mot: Memory-ofthought enables chatgpt to self-improve. In *EMNLP*, pages 6354–6374. Association for Computational Linguistics.
- Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Lidong Bing, Shafiq R. Joty, and Soujanya Poria. 2023a. Chain of knowledge: A framework for grounding large language models with structured knowledge bases. *CoRR*, abs/2305.13269.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023b. Making language models better reasoners with step-aware verifier. In *ACL* (1), pages 5315–5333. Association for Computational Linguistics.
- Zhenwen Liang, Dian Yu, Xiaoman Pan, Wenlin Yao, Qingkai Zeng, Xiangliang Zhang, and Dong Yu. 2023. Mint: Boosting generalization in mathematical reasoning via multi-view fine-tuning. *CoRR*, abs/2307.07951.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let's verify step by step. *CoRR*, abs/2305.20050.
- Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su. 2023. Deductive verification of chain-of-thought reasoning. In *NeurIPS*.
- Hao Liu, Carmelo Sferrazza, and Pieter Abbeel. 2023. Chain of hindsight aligns language models with feedback. *CoRR*, abs/2302.02676.
- Jieyi Long. 2023. Large language model guided treeof-thought. *CoRR*, abs/2305.08291.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Sean Welleck, Bodhisattwa Prasad Majumder,

884

885

886

- 836 837 838 839 840 841 842 843

- Shashank Gupta, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. CoRR, abs/2303.17651.
  - Ning Miao, Yee Whye Teh, and Tom Rainforth. 2023. Selfcheck: Using llms to zero-shot check their own step-by-step reasoning. CoRR, abs/2308.00436.
  - Michalis Mountantonakis and Yannis Tzitzikas. 2023. Using multiple RDF knowledge graphs for enriching chatgpt responses. In ECML/PKDD (7), volume 14175 of Lecture Notes in Computer Science, pages 324-329. Springer.
- Alexandru Niculescu-Mizil and Rich Caruana. 2005. Predicting good probabilities with supervised learning. In ICML, volume 119 of ACM International Conference Proceeding Series, pages 625-632. ACM.
- OpenAI. 2023. GPT-4 technical report. OpenAI.

790

795

796

810

811

812

814

815

816

818

- Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2023. Automatically correcting large language models: Surveying the landscape of diverse selfcorrection strategies. CoRR, abs/2308.03188.
  - Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In NAACL-HLT, pages 2080-2094. Association for Computational Linguistics.
  - Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, **REFINER:** reasoning and Boi Faltings. 2023. feedback on intermediate representations. CoRR, abs/2304.01904.
- Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. 2023. Reasoning with language model prompting: A survey. In ACL (1), pages 5368–5393. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. OpenAI.
- Subhro Roy and Dan Roth. 2015. Solving general arithmetic word problems. In EMNLP, pages 1743–1752. The Association for Computational Linguistics.
- Bilgehan Sel, Ahmad Al-Tawaha, Vanshaj Khattar, Lu Wang, Ruoxi Jia, and Ming Jin. 2023. Algorithm of thoughts: Enhancing exploration of ideas in large language models. CoRR, abs/2308.10379.
- Fobo Shi, Peijun Qing, Dong Yang, Nan Wang, Youbo Lei, Haonan Lu, and Xiaodong Lin. 2023. Prompt space optimizing few-shot reasoning success with large language models. CoRR, abs/2306.03799.

- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. In NeurIPS.
- Simeng Sun, Yang Liu, Shuohang Wang, Chenguang Zhu, and Mohit Iyyer. 2023. PEARL: prompting large language models to plan and execute actions over long documents. CoRR, abs/2305.14564.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In NAACL-HLT (1), pages 4149-4158. Association for Computational Linguistics.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D. Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In EMNLP, pages 5433-5442. Association for Computational Linguistics.
- Xingchen Wan, Ruoxi Sun, Hanjun Dai, Sercan Ö. Arik, and Tomas Pfister. 2023. Better zero-shot reasoning with self-adaptive prompting. In ACL (Findings), pages 3493-3514. Association for Computational Linguistics.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023a. Towards understanding chain-of-thought prompting: An empirical study of what matters. In ACL(1), pages 2717–2739. Association for Computational Linguistics.
- Jianing Wang, Qiushi Sun, Nuo Chen, Xiang Li, and Ming Gao. 2023b. Boosting language models reasoning with chain-of-knowledge prompting. CoRR, abs/2306.06427.
- Jinyuan Wang, Junlong Li, and Hai Zhao. 2023c. Self-prompted chain-of-thought on large language models for open-domain multi-hop reasoning. In EMNLP (Findings), pages 2717–2731. Association for Computational Linguistics.
- Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Lingi Song, Mingjie Zhan, and Hongsheng Li. 2023d. Mathcoder: Seamless code integration in llms for enhanced mathematical reasoning. CoRR. abs/2310.03731.
- Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. 2023e. Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering. CoRR, abs/2308.13259.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. 2023f. A survey on large

943

944

945

language model based autonomous agents. *CoRR*, abs/2308.11432.

887

891

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

928

929

930

931

933

935

936

937

938

941

- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023g. Plan-and-solve prompting: Improving zeroshot chain-of-thought reasoning by large language models. In ACL (1), pages 2609–2634. Association for Computational Linguistics.
- Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023h. SCOTT: selfconsistent chain-of-thought distillation. In *ACL* (1), pages 5546–5558. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023i. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Yiming Wang, Zhuosheng Zhang, and Rui Wang. 2023j. Meta-reasoning: Semantics-symbol deconstruction for large language models. *CoRR*, abs/2306.17820.
- Zekun Wang, Ge Zhang, Kexin Yang, Ning Shi, Wangchunshu Zhou, Shaochun Hao, Guangzheng Xiong, Yizhi Li, Mong Yuan Sim, Xiuying Chen, Qingqing Zhu, Zhenzhu Yang, Adam Nik, Qi Liu, Chenghua Lin, Shi Wang, Ruibo Liu, Wenhu Chen, Ke Xu, Dayiheng Liu, Yike Guo, and Jie Fu. 2023k. Interactive natural language processing. *CoRR*, abs/2305.13246.
- Zhaoyang Wang, Shaohan Huang, Yuxuan Liu, Jiahai Wang, Minghui Song, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, and Qi Zhang. 20231. Democratizing reasoning ability: Tailored learning from large language model. In *EMNLP*, pages 1948–1966. 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*.
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao.
  2023. Large language models are better reasoners with self-verification. In *EMNLP (Findings)*, pages 2550–2575. Association for Computational Linguistics.
- Zhenyu Wu, Meng Jiang, and Chao Shen. 2024. Get an a in math: Progressive rectification prompting. In *AAAI*. AAAI Press.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou,

Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huan, and Tao Gui. 2023a. The rise and potential of large language model based agents: A survey. *CoRR*, abs/2309.07864.

- Zhiheng Xi, Senjie Jin, Yuhao Zhou, Rui Zheng, Songyang Gao, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023b. Self-polish: Enhance reasoning in large language models via problem refinement. *CoRR*, abs/2305.14497.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. In *ICLR*. OpenReview.net.
- Xiaohan Xu, Chongyang Tao, Tao Shen, Can Xu, Hongbo Xu, Guodong Long, and Jianguang Lou. 2023. Re-reading improves reasoning in language models. *CoRR*, abs/2309.06275.
- Tianci Xue, Ziqi Wang, Zhenhailong Wang, Chi Han, Pengfei Yu, and Heng Ji. 2023. RCOT: detecting and rectifying factual inconsistency in reasoning by reversing chain-of-thought. *CoRR*, abs/2305.11499.
- Fanglong Yao, Changyuan Tian, Jintao Liu, Zequn Zhang, Qing Liu, Li Jin, Shuchao Li, Xiaoyu Li, and Xian Sun. 2023a. Thinking like an expert: Multimodal hypergraph-of-thought (hot) reasoning to boost foundation modals. *CoRR*, abs/2308.06207.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023b. Tree of thoughts: Deliberate problem solving with large language models. In *NeurIPS*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023c. React: Synergizing reasoning and acting in language models. In *ICLR*. OpenReview.net.
- Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, Ran Xu, Phil Mui, Huan Wang, Caiming Xiong, and Silvio Savarese. 2023d. Retroformer: Retrospective large language agents with policy gradient optimization. *CoRR*, abs/2308.02151.
- Yao Yao, Zuchao Li, and Hai Zhao. 2023e. Beyond chain-of-thought, effective graph-of-thought reasoning in large language models. *CoRR*, abs/2305.16582.
- Ori Yoran, Tomer Wolfson, Ben Bogin, Uri Katz, Daniel Deutch, and Jonathan Berant. 2023. Answering questions by meta-reasoning over multiple chains of thought. In *EMNLP*, pages 5942–5966. Association for Computational Linguistics.

- 999
- 1002 1003
- 1007 1008 1009 1010 1011
- 1012 1013 1014 1015
- 1016 1017
- 1018 1019 1020
- 1021 1022 1023
- 1024
- 1025

- 1029 1030
- 1031 1032

1037

1039

1040

1042 1044

1046

1047

1048

1049

1050

1026

1033

1035

1038

1043

Mengxue Zhang, Zichao Wang, Zhichao Yang, Weiqi Feng, and Andrew S. Lan. 2023. Interpretable math word problem solution generation via step-by-step planning. In ACL (1), pages 6858-6877. Association for Computational Linguistics.

Ruochen Zhao, Xingxuan Li, Shafiq R. Joty, Chengwei Qin, and Lidong Bing. 2023a. Verify-and-edit: A knowledge-enhanced chain-of-thought framework. In ACL (1), pages 5823–5840. Association for Computational Linguistics.

Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wen-

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Good-

man. 2022. STar: Bootstrapping reasoning with rea-

els through hybrid instruction tuning.

abs/2309.05653.

soning. In NeurIPS.

hao Huang, Huan Sun, Yu Su, and Wenhu Chen.

2023. Mammoth: Building math generalist mod-

CoRR,

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023b. A survey of large language models. CoRR, abs/2303.18223.

Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo Li, and Yu Li. 2023. Progressive-hint prompting improves reasoning in large language models. CoRR, abs/2304.09797.

Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023. Solving challenging math word problems using GPT-4 code interpreter with code-based self-verification. CoRR, abs/2308.07921.

#### Appendices

#### Α **Terminology Clarification of Answer Calibration and Model Calibration**

To avoid the confusion caused by the usage of the already-existing concept "calibration", we provide a terminology clarification. We emphasize that "answer calibration" defined in our paper differs from "model calibration" (Niculescu-Mizil and Caruana, 2005; Guo et al., 2017; Tian et al., 2023; Xiong et al., 2024). "Answer Calibration" refers to the post-processing methods applied to one or more reasoning path(s), to obtain a final answer. We categorize answer calibration methods as 'steplevel' if they break down the reasoning path(s) into their individual steps, and 'path-level' otherwise. In most cases, "Answer Calibration" is more akin to "Answer Correction" (Pan et al., 2023), involves

correcting mistakes in the initial output. We did 1051 give a definition like this in the Abstract, Introduc-1052 tion, and we have already provided clear definitions 1053 of "Answer Calibration" in §3. 1054

1055

1067

1068

#### **Model Choice Justification** B

The choice of GPT-3.5 was driven by its relevance and accessibility for our research objectives. Our research includes an empirical study of answer cal-1058 ibration and a proposal of a unified method, where 1059 the backbone LLM is pluggable. To facilitate reproducibility, we have already released the code 1061 and LLM-generated data anonymously<sup>1</sup> (provided 1062 at the bottom of Page 5 in §4), aiming to enhance 1063 transparency to some extent and facilitate further 1064 research in this area. We remain committed to ex-1065 ploring more transparent models in future work. 1066

#### С **Cases of Low-Quality Prompts**

We list some examples of prompts in Table 3.

Prompt Setting	<b>Example Query (Arithmetic Reasoning)</b> Leah had 32 chocolates and her sister had 42. If they ate 35, how many pieces do they have left in total?	
Standard CoT	Originally, Leah had 32 chocolates and her sister had 42. So in total they had $32 + 42 = 74$ . After eating 35, they had 74 - 35 = 39 pieces left in total. The answer is 39.	
No Coherence	After eating $32 + 42 = 74$ , they had 32 pieces left in total. Originally, Leah had $74 - 35 = 39$ chocolates and her sister had 35. So in total they had 42. The answer is 39.	
No Relevance	Patricia needs to donate 19 inches, and wants her hair to be 31 inches long after the donation. Her hair is 29 inches long currently. Her hair needs to be $19 + 31 = 50$ inc long when she cuts it. So she needs to grow $50 - 29 = 21$ more inches. The answer is 21.	

Table 3: Examples of prompts (standard, no coherence and no relevance) in our experiments.