EoT: Evolution of Thoughts for Complex Reasoning Tasks

Anonymous ACL submission

Abstract

Knowledge-based complex reasoning remains a significant challenge for large language models (LLMs) with in-context learning. To tackle this issue, previous studies focus on ensuring behavior fidelity, factuality, or reliability in generated reasoning processes that guide LLMs to produce solutions. However, these studies often neglect the simultaneous optimization on all these three aspects for each thought. The main challenges are the lack of comprehensive assessment mechanisms and the difficulty of efficient thought-level optimization. This paper introduces the Evolution of Thoughts (EoT) framework, which enhances the factuality, fidelity, and reliability of each thought in the reasoning process through a few LLM inferences. We propose a thought assessment method that is sensitive to knowledge and LLM behaviors, using three scorers to evaluate each thought by considering domain context, semantic alignment, and behavior impact. Additionally, we establish a self-reflective evolution mechanism to facilitate each reasoning process generation in a single-forward inference. Extensive experiments demonstrate that, for knowledge-based complex tasks, EoT improves the factuality and fidelity of reasoning processes by approximately 16.5% and 48.8%, respectively, while enhancing LLM reasoning capability by about 6.2%, outperforming advanced approaches.¹

1 Introduction

002

006

011

012

014

017

027

037

041

Nowadays, large language models (LLMs) such as GPTs (Achiam et al., 2024) and DeepSeeks (DeepSeek-AI, 2025) exhibit remarkable performance in various natural language processing (NLP) tasks. These models often employ in-context learning (ICL) schemes, enabling them to learn from contextual examples without updating billions of parameters (Brown et al., 2020). However, complex reasoning tasks requiring the comprehension of long-context domain knowledge and the generation of intricate solutions remain challenging for LLMs with ICL prompting (Chen et al., 2024). 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

Many studies have shown that guiding LLMs through step-by-step thought prompts significantly enhances their reasoning capabilities (Chen et al., 2024; Lyu et al., 2023). These prompts inspire LLMs to create reasoning processes that explain their behaviors and aid in task resolution. Each reasoning process, such as the chain of thoughts (CoT) and its variants (Wang et al., 2023), comprises a sequence of coherent text units known as thoughts, which serve as intermediate reasoning steps. In these mechanisms, the reasoning capability of LLMs is often reflected in the reliability of reasoning processes, which assesses the confidence or correctness of the solutions produced by LLMs under the guidance of reasoning processes (Zhang et al., 2024; Madaan et al., 2023). Therefore, many efforts have been dedicated to exploring logical reasoning processes with high reliability (Radhakrishnan et al., 2023; Besta et al., 2024).

Existing studies on optimizing reasoning processes can be categorized into three main areas. Firstly, some studies focus on enhancing behavior fidelity (Chuang et al., 2024; Lyu et al., 2023). Previous research (Liang et al., 2024) indicates that LLMs often provide reasoning processes that differ significantly from their actual reasoning behaviors, which compromises their reasoning capabilities. This discrepancy arises from the lack of awareness regarding the internal knowledge state in black-box LLMs. To address this issue, xLLM (Chuang et al., 2024) revises reasoning processes to revises reasoning processes through evolutionary iterations to improve their behavior fidelity. Secondly, some work address non-factual errors in reasoning processes to mitigate the hallucination phenomenon in LLMs. For example, Ye et al. (Ye and Durrett, 2022) extract the most factual reasoning process in each generation. Thirdly, many studies directly

¹The code of EoT and the case of experiment data can be found in our supplementary material files.

100

101

102

103

105

106

107

109

110

111

112

113

114

115

116

117

118

119

121

122

123

124

125

126

127

128

130

131

132

134

explore reasoning processes with high reliability (Madaan et al., 2023; Besta et al., 2024). For instance, BoT (Chen et al., 2024) uses an iterative method to explore many trees of thoughts, and accumulates trial-and-error experiences to derive a reasoning process yielding a reliable answer.

Nevertheless, enhancing the reasoning process for complex tasks presents three major challenges. multi-objective optimization regarding factuality, fidelity, and reliability has been a long-standing issue, as highlighted by previous research (Ye and Durrett, 2022; Liang et al., 2024). Secondly, the complexity of long-textual tasks exacerbates the difficulty in optimizing reasoning thoughts. These tasks require detailed knowledge across domains and expect complex solutions, it necessitates a reasoning process with multiple thought steps, involving intricate logic and precise extraction of relevant facts. This complexity hinders the assessment and improvement of the effectiveness of each thought. Thirdly, achieving thought-level optimization for reasoning processes while maintaining efficiency is challenging. Some studies (Chen et al., 2024; Radhakrishnan et al., 2023) improve thought quality by generating individual reasoning thoughts through extensive explorations, which incurs significant overhead. Conversely, studies like (Chuang et al., 2024; Ye and Durrett, 2022) strive to produce an improved reasoning process using a single LLM inference, it enhances time efficiency but compromises the fine-grained optimization for each thought step.

To address these challenges, we propose EoT, a framework that evolves reasoning processes to achieve multi-objective optimization in factuality, fidelity, and reliability. Firstly, EoT includes an assessment mechanism designed for complex reasoning tasks, which evaluates each thought in reasoning processes from all three perspectives. Secondly, we introduce a prompting mechanism to facilitate the creation of evolved reasoning processes with a single-forward LLM inference. This enables LLMs to comprehend thought assessment outcomes and ensure collaborative optimization for each thought within a few rounds of self-reflective evolution.

As highlighted in Table 1, EoT distinguishes itself from existing evolution frameworks with two key features. Firstly, EoT evaluates reasoning processes more comprehensively in three critical dimensions. Secondly, EoT effectively manages both time efficiency and thought-level optimization, producing a complete reasoning process in each iteration through single forward inference and ensuring

Table 1: Comparison of existing studies

	Reasonir	Reasoning Process		Optimization & Assessment			
Framework	Componente	Generation	Factors			Laval	
	components		Reliability	Factuality	Fidelity	LCVCI	
BoT (Chen et al., 2024)	ToT	ISTE, Node	~	X	X	Thought	
GoT (Besta et al., 2024)	GoT	ISTE, Node	✓	X	X	Thought	
Factor-dec (Radhakrishnan et al., 2023)	SQAT	ISTE, SQA	✓	X	X	Thought	
CoT-dec (Radhakrishnan et al., 2023)	SQAT	SFI	√	X	X	CRP	
Ye et al. (Ye and Durrett, 2022)	CoT	SFI	X	√	~	CRP	
xLLM (Chuang et al., 2024)	CoT	SFI	X	X	~	CRP	
EoT (ours)	CoT	SFI	✓	~	~	Thought	
⁻¹ CoT, ToT, GoT, SQAT and SQA stand for chain of thoughts, tree of thoughts, graph of thoughts, sub question-answer as thoughts, and sub							

question-answer respectively. 2 ISTE, SF1 and CRP stand for iterations of each step of thought exploration, single forward inference, and complete reasoning process respectively.

granular optimization at each reasoning step.

We conducted extensive experiments on two datasets, including one with 40 production operational maintenance tasks. Compared to five advanced frameworks, EoT improves reliability, factuality, and fidelity of reasoning processes by about 6.2%, 16.5%, and 48.8% respectively. The contributions of this work are summarized as follows:

- We consider reliability, factuality, and fidelity of reasoning processes to enhance LLMs' reasoning capabilities, and assess each thought in reasoning processes based on these factors.
- We propose an evolving mechanism that efficiently facilitates multi-objective optimization at the thought level via single-forward generation of the reasoning process in each iteration.
- We evaluate the effectiveness of EoT on a real production dataset and verify its generality on the LongBench dataset including questions obtained from diverse fields, such as MultifieldQA and HotpotQA. Each question involves a context with thousands of tokens.

2 Problem Setup

This section formulates reasoning processes, introduces three key factors that impact LLMs' reasoning capabilities, and defines the evolution problem.

2.1 Formalization of Reasoning Process

We first formalize LLM-generated reasoning processes used for solving complex tasks.

Let $Q = (q_1, q_2, \dots)$, $X = (x_1, x_2, \dots)$ and $A = (a_1, a_2, \dots)$ represent a question description, knowledge context and a reference answer for any complex task. Each component consists of sequential statements, denoted by q_i , x_i , and a_i respectively. When an LLM p_θ receives a task (X, Q), it generates a reasoning process R to explain its reasoning behaviors. Using R, the LLM guides itself to produce a solution \hat{Y} for the task, formalized as:

$$Y = p_{\theta}(X, Q|R) \tag{1}$$

135

136

138

139

140

141

142

143

144

157

158

154

155

156

159 160

161

162

163

164

165

166

167

168

169

170

171

172

223 224

225

226

227

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

 $Factual(T_i) = \frac{\sum_{t_{i,j} \in T_i} Ground(X, t_{i,j})}{J}$ $Ground(X, t_{i,j}) = \begin{cases} 1 & t_{i,j} \text{ is grounded in } X\\ 0 & Otherwise \end{cases}$ (3)

 $T_i = (t_{i,1}, t_{i,2}, \cdots, t_{i,J})$ represent a thought step

consisting of J statements in R for question Q with

context X. The factuality of T_i is expressed as:

Definition 3: Fidelity evaluates the faithfulness of thoughts in the reasoning process to explain the actual behaviors of the LLM when generating answers. Based on previous studies (Lopardo et al., 2023; Chuang et al., 2024), fidelity is defined by the extent to which an explained reasoning thought influences the LLM-generated answers, assuming ground-truth reasoning behaviors are typically available. A greater degree indicates higher fidelity. Specifically, for a reasoning process R, the fidelity of each thought $T_i \in R$ is assessed by comparing the answers generated by the LLM guided by R with and without T_i . This is expressed as:

$$Fidelity(T_i) = Diff(p_{\theta}((X,Q)|R), \\ p_{\theta}((X,Q)|(R \setminus T_i)))$$
(4)

where $(R \setminus T_i)$ represents the reasoning process without thought T_i , and $Diff(\cdot, \cdot)$ denotes the difference estimation. Appendix C.2 illustrates thoughts with high and low fidelity.

In summary, EoT guides LLMs to find a refined *n*-step reasoning process $R^* = (T_1^*, T_2^*, ..., T_n^*)$ via iterative self-reflecting evolution. In each iteration, the LLM p_{θ} generates an enhanced reasoning process to address multiple-objective optimization in reliability, factuality and fidelity for each reasoning thought, as defined in Eq. (5):

$$R^{*} \leftarrow \arg \max_{R} Reliability(R)$$

$$maxmize \ Factual(T_{i}^{*}) \ \text{for} \ \forall T_{i}^{*} \in R^{*}$$

$$maxmize \ Fidelity(T_{i}^{*}) \ \text{for} \ \forall T_{i}^{*} \in R^{*}$$

$$(5)$$

3 Evolution of Thoughts

3.1 Overview

This section introduces the Evolution of Thoughts (EoT) framework to enhance reasoning processes. EoT refines thoughts evolutionarily to simultaneously improve their reliability, factuality, and behavior fidelity, thereby boosting LLMs' reasoning

A generated reasoning process R comprises mul-174 tiple thoughts, denoted as $R = \{T_1, T_2, \cdots T_n\}$. 175 Each thought includes several statements, i.e., $T_i =$ 176 $(t_{i,1}, t_{i,2}, \cdots)$. This study explores an evolved rea-177 soning process R^* for each task (X, Q), to enhance the reasoning ability of p_{θ} . The guidance 179 of R to produce \hat{Y} , as depicted in Eq.(1), is real-180 ized through instructions to prompt LLM to reason 181 based on thoughts in R, as shown in Figure 4 in Ap-182 pendix B. Figures 5 and 7 in Appendix B presents 183 examples of question-solving with thoughts.

2.2 Reliability, Factuality and Fidelity of Thoughts

186

187

188

190

191

194

195

196

197

198

201

202

204

207

210

211

212

213

214

215

216

217 218

219

222

Next, we outline and define three key factors of reasoning processes that impact LLM's reasoning capabilities. To aid comprehension, we illustrate them with examples in Appendix C.

Previous studies (Chen et al., 2024; Zheng et al., 2023; Paul et al., 2023) have revealed that the reliability of reasoning processes is essential in evaluating LLMs' question-solving abilities. To elucidate, we define the reliability of a reasoning process R:

Definition 1: Reliability of R is measured by the similarity between the answer from LLMs guided by R and the reference answer. Greater similarity indicates higher reliability. Given a task, a generated answer \hat{Y} guided by R, and reference answer A, the reliability of R is defined as:

$$Reliability(R) = Similarity(\hat{Y}, A) \quad (2)$$

The similarity metric can be computed using various methods, such as token overlaps (Lin, 2004), learning-based distance (Sellam et al., 2020), and entailment extent in natural language inference (NLI) (Gao et al., 2023). Intuitively, Figure 6 and Figure 8 in Appendix C illustrate reasoning processes with different levels of reliability.

As outlined in Section 1, enhancing behavior fidelity and reducing hallucinations of reasoning processes can improve LLMs' reasoning capabilities. Inspired by prior studies (Ye and Durrett, 2022; Chuang et al., 2024), we define two key factors of reasoning processes, factuality and fidelity:

Definition 2: Factuality pertains to how well thoughts in reasoning processes are grounded in the relevant knowledge context. As shown in Figure 6 in Appendix C, a fully factual reasoning process excludes hallucinations that contradict the context. Conversely, a non-factual reasoning process with hallucinations can lead to erroneous solutions. Let

300 301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

330

331

332

333

334

335

336

337

338

339

340

341

343

344



Figure 1: The framework of EoT.

 \bigcirc

260

263

264

265

271

272

273

276

278

279

284

285

290

291

capabilities. Each iteration of evolution enables LLMs to comprehend the assessment of the current reasoning process and provides feedback to refine a group of individual thoughts in the reasoning process via a single-forward inference. In this way, EoT effectively and efficiently guides reasoning process optimization toward multiple objectives.

Figure 1 illustrates the framework of EoT, comprising three modules: an assessor, a prompter and a generator. These modules collaborate to achieve multi-objective optimization of reasoning processes over iterations of evolution. Initially, the assessor uses three scorers to evaluate the reliability, fidelity and factuality of thoughts in the reasoning process during each iteration. The prompter then guides LLMs to thoroughly comprehend the performance of the current reasoning process accross the three aspects, prompting self-reflection in LLMs to ensure thought-level optimization. Lastly, in each iteration, the generator served by task-solving LLMs, uses our crafted prompts to generate a refined reasoning process through a single-forward inference. After N iterations, EoT produces an improved reasoning process, denoted as R^* . Further details of the three modules are provided below.

3.2 Assessor

EoT focuses on optimizing the factuality, reliability and fidelity of reasoning processes. Therefore, we design three scorers in the assessor to quantify the performance of thoughts across these three aspects.

3.2.1 Factuality Scorer

As defined in Eq.(3), the factuality of reasoning process for complex tasks is expressed by how well the thoughts in it can be supported by relevant domainknowledge facts. For tasks with long-textual context, EoT leverages the impressive capability of LLMs in natural language inference (NLI) to tackle the scoring of factuality. We employ an NLI model ϕ to estimate $Ground(\cdot, \cdot)$ in the factuality definition in Eq. (3). Specifically, given a pair of premise and hypothesis statements, represented as x^{pre} and t^{hyp} respectively, ϕ infers their relationship through a triple-label classification task:

$$\phi(x^{pre}, t^{hyp}) = \begin{cases} 0 & t^{hyp} \text{ contradicts with } x^{pre} \\ 1 & x^{pre} \text{ entails } t^{hyp} \\ 2 & x^{pre} \text{ is neutral with } t^{hyp} \end{cases}$$
(6)

Definitions of "Entailment", "Contradiction" and "Neutral" are provided in Appendix E.1. Inspired by prior work (Gao et al., 2023), consider a reasoning process $R = (T_1, T_2, \cdots)$ for task (X, Q), where each thought $T_i = (t_{i,1}, t_{i,2}, \cdots, t_{i,J})$ contains J statements. For $\forall t_{i,j} \in T_i$, we define $Ground(X, t_{i,j}) = 1$ iff $\phi(X, t_{i,j}) = 1$; otherwise, it is 0. We use a GPT-4 model with few-shot learning for ϕ . The efficacy of this NLI scheme is presented in Appendix E.2. Then, the factuality score $S_{fac}(T_i)$ for $\forall T_i \in R$ is computed as:

$$S_{fac}(T_i) = Factual(T_i) \tag{7}$$

3.2.2 Reliability Scorer

According to Eq.(2), assessing the reliability of the reasoning process R largely relies on measuring the similarity between the answer generated under R and the reference answer. However, existing similarity metrics struggle to ensure both sufficient variation sensitivity and robust human-level alignment simultaneously. To address these issues, EoT proposes a novel similarity metric.

Recent studies propose learning-based methods, such as SimCSE (Gao et al., 2021) and BLEURT (Sellam et al., 2020), to improve sensitivity to semantic and syntactic variations beyond handcrafted metrics like ROUGE (Lin, 2004) and BLEU (Papineni et al., 2002). These methods utilize language models such as BERT (Devlin et al., 2019) to encode statement pairs, producing scalar similarity scores based on the encoded representations.

However, learning-based metrics struggle to robustly align with human judgment in domainspecific scenarios. These metrics typically rely on end-to-end predictions from models trained on synthetic or commonsense datasets (Zhang et al., 2019). Discrepancies between distributions of training data and domain knowledge can lead to misalignment due to domain drift (Honovich et al., 2022). Moreover, while syntactic alignment strategies are prevalent, achieving semantic alignment akin to human judgment remains challenging. To enhance robustness while maintaining high sensitivity, we introduce Hssim, a hybrid metric to score similarity between statements, which combines the learned metric BLEURT with the NLI judgement of model ϕ , as detailed in Eq. (6). Specifically, given a reasoning process R for task (X, Q), the LLM generates an answer of n_1 statements $\hat{Y} = (\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_{n_1})$ guided by R. Let $A = (a_1, a_2, \cdots a_{n_2})$ represent a reference answer n_2 statements. The similarity between \hat{Y} and A is evaluated using $Hssim(\hat{Y}, A)$ calculated as:

$$\begin{aligned} Hssim(\hat{Y}, A) &= \frac{\sum_{\hat{y}_i \in \hat{Y}} sim(\hat{y}_i, A)}{n_1} \\ sim(\hat{y}_i, A) &= \begin{cases} 0 & \phi(A, \hat{y}_i) = 0 \\ \mu + (1 - \mu)\beta_w(A, \hat{y}_i) & \phi(A, \hat{y}_i) = 1 \\ (1 - \mu)\beta_w(A, \hat{y}_i) & Otherwise \end{cases} \end{aligned}$$
(8)

where $\mu \in (0, 1)$ is the weight coefficient for measuring similarity between the NLI judgement and the learning-based scalar metric, and $\beta_w(\cdot, \cdot)$ denotes an approximation of the BLEURT score at the statement window level, calculated as:

$$\beta_w(A, \hat{y}_i) = \max(\{BLEURT(A[i_1 : i_2], \hat{y}_i) | \\ \forall 0 < i_1 < i_2 \le n_2\})$$
(9)

Leveraging the exceptional in-context learning capabilities of LLMs, their NLI judgment can be robustly aligned with human judgement even facing domain drift, as noted in previous work (Ye and Durrett, 2022). Furthermore, the window-based estimation of BLEURT scores ensures sensitivity and alignment at the syntactic level. This allows for a rational and effective assessment of the reliability of any reasoning process R for question (X, Q), using Hssim, which balances variation sensitivity and semantic alignment robustness:

$$S_{rel}(R) = Hssim(p_{\theta}(X, Q|R), A)$$
(10)

3.2.3 Fidelity Scorer

360

362

367

374

375

385

Based on the fidelity definition in Eq.(4), scoring the fidelity of thoughts involves two key aspects: 1) efficiently excluding a thought from LLMs' reasoning behaviors, and 2) accurately measuring the difference between answers produced with and without the thought. We use a thought masking scheme and metric Hssim to achieve these goals.

Given a reasoning process R, each thought $T_i \in R$ is removed from R to create a masked reasoning process $(R \setminus T_i)$, as shown in Figure 8 in Appendix



Figure 2: Prompt template to evolve reasoning process.

C. According to Eq.(1), a prompt directs the LLM to generate answers \hat{Y} and $\hat{Y}_{\backslash T_i}$ for task (X, Q), using instructions from R and $(R \backslash T_i)$ respectively. An example is provided in Figure 4 in Appendix B. Finally, with \hat{Y} as the reference, the fidelity score of T_i is assessed using the *Hssim* metric:

$$S_{fid}(T_i) = 1 - Hssim(Y_{\backslash T_i}, \hat{Y})$$
(11)

386

387

388

390

391

393

394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

Moreover, enhancing LLMs' reasoning capabilities requires fidelity optimization to focus on thoughts that faithfully represent the reasoning behaviors essential for reliable answers. Thus, we propose a weighted fidelity score for thought T_i , accounting for the reliability of the reasoning process:

$$S_{fid}^{rel}(T_i) = S_{fid}(T_i) \times S_{rel}(R)$$
(12)

Ultimately, we achieve a comprehensive assessment of the reasoning process and its thoughts:

$$\left(S_{rel}(R), \{S_{fac}(T_i), S_{fid}^{rel}(T_i) | \forall T_i \in R\}\right)$$
(13)

3.3 Prompter

The prompting mechanism of EoT is designed to enhance the reasoning process by enabling LLMs to comprehend the scoring outcomes from previous iterations. This prompt LLMs to self-reflect and refine reasoning processes, optimizing reliability, fidelity, and factuality at the thought level.

The prompt template for the *K*-th iteration $(1 \le K \le N)$ is presented in Figure 2. Each evolution iteration includes a prompt with three components: (1) the question description *Q*, knowledge context *X*, and reference answer *A* for the task; (2) the evolution logic, encompassing: a) the significance of three score types and b) instructions for multi-objective optimization; and (3) the assessing results of reasoning processes produced in prior iterations, as defined in Eq.(13). An example of a complete prompt is shown in Figure 9 in Appendix D.

Table 2: The parameter settings in evaluations
--

Parameter	Value	Description
N	10	The number of iterations
μ	0.5	The alignment weight in Eq.(8)

3.4 Self-reflective Generator

To align the LLM's reasoning capability with the fidelity of its behavioral patterns, EoT employs the LLM used for question answering as the generator to iteratively refine reasoning processes via self-reflection. Each iteration uses prompts to guide the generator to produce an enhanced reasoning process, leading to multi-objective improvements recognized by LLMs. After N iterations, the reasoning process with the highest reliability score is chosen as the final evolved process R^* .

4 Experiments

4.1 Setup

Datasets To evaluate the effectiveness of EoT, we conduct experiments on two datasets: 1) OpsOA: an question-answering dataset with 40 complex operational maintenance tasks covering cloud computing, code management, application upgrades, deployment, and more. 2) LongBench (Bai et al., 2023b): a benchmark comprising problems from various open-source datasets such as HotpotQA, MultifieldQA, WikimQA. We select 60 representative questions from LongBench for evaluation. All tasks from OpsQA and LongBench involve complex reasoning requiring long-context knowledge. Each task includes a domain knowledge context with thousands or tens of thousands of tokens, a question description, and a reference answer, necessitating intricate solutions. Context length statistics are detailed in Table 7 in Appendix F.

Testbed and Parameter Settings We utilize two widely used LLMs, Qwen2-72B (Bai et al., 2023a) and GPT-4 Turbo (Openai, 2023), each severing as the generator of both reasoning processes and solutions. For the assessor, EoT uses GPT-4 Turbo to implement the NLI model ϕ , due to its advanced NLI accuracy as evaluated in Appendix E.2, and compute window-based BLEURT scores in Eq. (9) using V100 GPUs. Table 2 outlines the hyperparameter settings in our evaluations.

461MetricsTo assess LLMs' reasoning capabilities,462we compute four metrics by comparing answers463generated by the reasoning process \hat{Y} with the ref-464erence answer A: 1) reliability defined in Eq. (10),4652) BLEURT, 3) ROUGE-L (Lin, 2004), and 4) NLI

results from GPT-4 Turbo, denoted as entail(%), which measures the percentage of statements in \hat{Y} entailed by A. Additionally, following prior studies (Ye and Durrett, 2022; Lyu et al., 2023), we use S_{fac} and S_{fid} from Eq.(7) and Eq.(11) to assess performance in factuality and behavior fidelity, respectively. Scores of a reasoning process are averaged from those of thoughts. Higher values indicate better performance for each metric.

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

Baselines We compare EoT with five advanced frameworks for evolving reasoning processes: 1) BoT (Chen et al., 2024), CoT-dec (Radhakrishnan et al., 2023), and its variant Factor-dec, which emphasize reliability optimization; 2) xLLM (Chuang et al., 2024) designed for fidelity optimization; and 3) Calibrator (Ye and Durrett, 2022) targeting factuality and fidelity optimization. EoT, CoT-dec, and Calibrator generate reasoning via single-forward inference, while BoT and Factor-dec produce thought steps as nodes of ToT and subquestions during each LLM inference, respectively. For more details on baselines, see Appendix G.

4.2 Overall Performance

We conduct experiments on the OpsQA and Long-Bench datasets to evaluate the performance² of the six frameworks in terms of the reasoning capability, factuality and fidelity of reasoning processes and the time efficiency of evolution. Moreover, ablation study of EoT is detailed at Appendix I.

4.2.1 Reasoning Capability (Reliability)

Table 3 presents the average reasoning capability of two LLMs utilizing the six frameworks. EoT surpasses the five baselines in three areas. First, compared to the leading baseline, CoT-dec, EoT boosts reliability scores on the OpsQA and Long-Bench datasets by about 11.7% and 5.8% using Qwen2-72B, and by about 3.3% and 3.7% using GPT-4 Turbo. Second, in terms of sensitivity to token and semantic variation, EoT improves the ROUGE-L and BLEURT on LongBench dataset by about 7.2% and 16.2% using Qwne2-72B, and 2.6% and 6.8% using GPT-4 Turbo. Third, regarding semantic alignment robustness, on the OpsQA dataset, EoT improves entail(%) by around 14.4%and 6.7% using Qwen2-72B and GPT-4 Turbo, respectively. These findings indicate that EoT empowers diverse LLMs to enhance overall reasoning

428

429

430

431

421

435 436 437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

²Appendix J applies the T-test method to examine the significance of performance variations.

			-				-				
		Ops		sQA	sQA		LongBench				
LLMs	Models	DIFUDT	POLICE I	NLI results	Reliability	DIFUDT	POLICE I	NLI results	Reliability		
BLE	DLEUKI	KOUGE-L	entail(%)	Score	DLEUKI	KOUGE-L	entail(%)	Score			
	BoT	0.499	0.297	43.29	0.401	0.494	0.677	62.92	0.526		
	xLLM	0.554	0.355	65.83	0.576	0.572	0.739	87.50	0.721		
Owen?	Calibrator	0.559	0.365	70.27	0.597	0.575	0.744	80.88	0.675		
Qwell2	CoT-dec	0.579	0.361	73.35	0.632	0.586	0.769	92.63	0.757		
	Factor-dec	0.545	0.324	62.55	0.569	0.570	0.745	84.47	0.729		
	EoT (ours)	0.595	0.395	83.90	0.706	0.681	0.824	91.52	0.801		
	BoT	0.486	0.307	48.50	0.443	0.530	0.704	67.99	0.568		
	xLLM	0.575	0.408	72.16	0.609	0.629	0.801	85.09	0.729		
GPT-4	Calibrator	0.554	0.380	68.50	0.589	0.592	0.775	83.98	0.709		
011-4	CoT-dec	0.613	0.446	83.02	0.694	0.617	0.805	91.65	0.767		
	Factor-dec	0.541	0.358	74.48	0.632	0.587	0.767	89.44	0.739		
	EoT (ours)	0.602	0.427	88.58	0.717	0.659	0.826	92.47	0.795		

Table 3: Reliability performance for reasoning processes evolved by six frameworks with two LLMs on two datasets.

¹ Qwen2 and GPT-4 represent the LLMs of Qwen2-72B and GPT-4 Turbo, respectively



Figure 3: The CDF results of the reliability, factuality and fidelity scores of reasoning processes evolved by the six frameworks using GPT-4 Turbo on 40 tasks in the OpsQA dataset.

capability in complex tasks.

513

514

515

516

517

518

519

520

521

522

524

525

526

528

529

531

533

535

539

Figure 3a intuitively shows that, with GPT-4 Turbo, EoT achieves a reliability score > 0.75 for 52.5% of tasks in the OpsQA dataset, surpassing BoT, xLLM, Calibrator, CoT-dec, and Factor-dec by about 50%, 27.5%, 22.5%, 7.5% and 35%, respectively. Furthermore, for tasks in OpsQA and LongBench, the complete CDF results of reasoning capability using Qwen-72B and GPT-4 Turbo, equipped with the six frameworks, are detailed in Appendix H.1. These findings confirm EoT's significant generality in enhancing reasoning capability across diverse fields with various LLMs.

EoT's enhanced reasoning capability stems from two key aspects. Firstly, we introduce an effective similarity scoring metric *Hssim* to measure the reliability of reasoning processes. By accounting for sensitivity to semantic variations and robustness in semantic alignment, *Hssim* promote the reliability improvement of EoT, facilitating comprehensive optimization of LLMs' reasoning capabilities across various semantic aspects. For instance, on the OpsQA dataset, despite EoT's BLEURT performance being about 1.79% lower than CoT-dec using GPT-4 Turbo, EoT improves NLI results by about 6.7%, thereby raising the reliability score by 3.3%. Secondly, unlike baselines focusing on individual or partial factors, EoT efficiently achieves multi-objective optimization in reliability, factuality, and fidelity during reasoning process evolution, thereby enhancing LLMs' reasoning capabilities.

Table 4: The performance on factuality and fidelity for reasoning processes evolved by six frameworks.

LIM	Madala	Ops	QA	LongBench		
LLIVIS	Widdels	Factuality Fidelity		Factuality	Fidelity	
		S_{fac}	S_{fid}	S_{fac}	S_{fid}	
	BoT	0.752	0.234	0.879	0.208	
	xLLM	0.720	0.216	0.864	0.152	
Owan?	Calibrator	0.860	0.258	0.935	0.141	
Qwell2	CoT-dec	0.638	0.167	0.839	0.179	
	Factor-dec	0.706	0.208	0.901	0.197	
	EoT (ours)	0.823	0.267	0.943	0.243	
	BoT	0.775	0.203	0.896	0.237	
	xLLM	0.607	0.238	0.848	0.169	
CPT 4	Calibrator	0.878	0.257	0.942	0.180	
011-4	CoT-dec	0.685	0.199	0.934	0.171	
	Factor-dec	0.769	0.228	0.936	0.208	
	EoT (ours)	0.906	0.281	0.956	0.275	

4.2.2 Factuality & Fidelity Performance.

The central concept of EoT is to improve the factuality and behavior fidelity of thoughts, in a direction of enhancing LLMs' reasoning capabilities. To assess these efforts, Table 4 presents the average performance in fidelity and factuality of reasoning processes evolved by six frameworks, using Qwen2-72B and GPT-4 Turbo on the two datasets. These results highlight two key aspects. 544

545

546

547

548

549

550

551

552

Firstly, EoT significantly reduces non-factual 553 errors in thoughts while boosting the behavior fi-554 delity of LLMs. Compared to the leading baseline 555 Calibrator, using GPT-4 Turbo, EoT improves the factuality score by about 3.2% and 1.3% on the OpsQA and LongBench datasets respectively. With 558 Owen2-72B, EoT shows a factuality improvement 559 of around 4.0% on the LongBench dataset. Additionally, on the OpsQA dataset, EoT surpasses Calibrator by enhancing the fidelity score by about 562 9.3% using GPT-4 Turbo. Similarly, on the Long-563 Bench dataset, compared to BoT, EoT improves the 564 fidelity score by roughly 16.8% and 16.0% using 565 Qwen-72B and GPT-4 Turbo respectively. These results confirm EoT's capability to effectively en-567 hance both factuality and behavior fidelity.

> Secondly, EoT's optimization on factuality and fidelity effectively boosts reasoning capabilities of LLMs. On the OpsQA dataset, EoT surpasses CoTdec in factuality and fidelity by about 28.9% and 59.9% respectively, resulting in a reliability improvement of around 11.7%. This enhanced reasoning capability is largely due to substantial advancements in factuality and fidelity of reasoning processes. Moreover, although Calibrator improves factuality by about 4.4% over EoT on the OpsQA dataset, EoT improves fidelity and reliability by about 3.5% and 18.3% respectively.

571

573

576

579

581

583

584

587

588

589

591

592

593

594

596

601

To illustrate EoT's generality in optimizing factuality and fidelity, Figures 3b and 3c show the CDF of factuality and fidelity scores for reasoning processes evolved by six frameworks using GPT-4 Turbo on the OpsQA dataset, respectively. EoT achieves a factuality score of 1.0 for 50% of OpsQA questions, surpassing BoT, xLLM, Calibrator, CoT-dec, and Factor-dec by around 15%, 47.5%, 20%, 35%, and 42.5%, respectively. Moreover, EoT exceeds a fidelity score of 0.25 for 60% of questions, surpassing BoT, xLLM, Calibrator, CoTdec, and Factor-dec by about 32.5%, 20%, 25%, 50%, and 37.5%, respectively. Further CDF results for fidelity and factuality on OpsQA and Long-Bench datasets are provided in Appendix H.2 and H.3. These results confirm that EoT's improvements in fidelity and factuality are generalized to tasks across various domains, including operational maintenance and academic literature.

These improvements stem from two aspects. Firstly, EoT effectively assesses factuality and fidelity at a granular thought level, providing a foundation for effective optimization. Secondly, our prompter fully leverages ICL capabilities of LLMs, Table 5: The average performance on time efficiency of the six frameworks on the OpsQA dataset.

Models	Overhead(s)	Iterations
BoT	945.85	4.48
Factor-dec	114.77	5.21
xLLM	70.89	5.25
Calibrator	67.75	4.92
CoT-dec	70.22	5.34
EoT (ours)	75.46	4.62

enabling them to rationally comprehend performance experiences and simultaneously optimize across three objectives, as defined in Eq.(5).

605

606

607

608

609

610

611

612

613

614

615

616

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

4.2.3 Time Efficiency

We assess the time efficiency of six frameworks using two criteria: 1) the average overhead per iteration, where lower values denote higher efficiency; and 2) the convergence rate, measured as the average number of iterations needed to produce the final enhanced reasoning process for various questions. Lower values signify quicker convergence.

Table 5 showcases the time efficiency of the six frameworks on the OpsQA dataset. EoT reduces iteration overhead by about 92.0% compared to BoT and 33.4% versus Factor-dec, both of which involve multiple thought explorations per iteration. In addition, compared to three other baselines employing single-forward inference, EoT uses fewer iterations to generate refined reasoning processes with similar iteration time. These results validate that EoT achieves comprehensive, fine-grained optimization of each thought with high time efficiency.

5 Conclusion

Through in-depth analysis, we outline three key factors in reasoning processes to enhance LLMs' capabilities, namely factuality, fidelity and reliability. We propose EoT to evolve LLM-generated reasoning processes across these dimensions. An assessor is designed to quantify performance on these aspects for each thought in reasoning processes. Additionally, we propose a prompting mechanism that efficiently guides LLMs to comprehend assessments and trigger self-reflection to achieve thoughtlevel and multi-objective optimization of reasoning processes via single-forward inference. EoT offers two key advantages. First, it considers comprehensive factors to improve reasoning capability. Second, it ensures thought-level evolution with high time efficiency. In the future, we will address unsupervised evolution of reasoning processes to further enhance reasoning capabilities of LLMs for out-ofdomain tasks without reference solutions.

References

647

651

652

653

654

655

657

661

664 665

670

674

675

678

679

681

684

687

701

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, and et al. 2023a. Qwen technical report. *Preprint*, arXiv:2309.16609.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023b. Longbench: A bilingual, multitask benchmark for long context understanding. *Preprint*, arXiv:2308.14508.
 - Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and 1 others. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and et al. 2020. Language models are fewshot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.
- Sijia Chen, Baochun Li, and Di Niu. 2024. Boosting of thoughts: Trial-and-error problem solving with large language models. *arXiv preprint arXiv:2402.11140*.
 - Yu-Neng Chuang, Guanchu Wang, Chia-Yuan Chang, Ruixiang Tang, Fan Yang, Mengnan Du, Xuanting Cai, and Xia Hu. 2024. Large language models as faithful explainers. *arXiv preprint arXiv:2402.04678.*
- DeepSeek-AI. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *Preprint*, arXiv:2501.12948.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling large language models to generate

text with citations. In *Empirical Methods in Natural* Language Processing (EMNLP). 702

703

704

705

706

708

709

710

711

712

713

715

717

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

- Or Honovich, Roee Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022. TRUE: Re-evaluating factual consistency evaluation. In *Proceedings of the 2nd DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering*, pages 161– 175, Dublin, Ireland. Association for Computational Linguistics.
- Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiaxing Zhang. 2024. Learning to trust your feelings: Leveraging self-awareness in llms for hallucination mitigation. *arXiv preprint arXiv:2401.15449*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Gianluigi Lopardo, Frederic Precioso, and Damien Garreau. 2023. Faithful and robust local interpretability for textual predictions. *arXiv preprint arXiv:2311.01605*.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-ofthought reasoning. *Preprint*, arXiv:2301.13379.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, and et al. 2023. Self-refine: Iterative refinement with self-feedback. In *Advances in Neural Information Processing Systems*, volume 36, pages 46534–46594.

Openai. 2023. Gpt-4 turbo in the openai api.

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting on Association for Computational Linguistics, ACL '02, page 311–318, USA. Association for Computational Linguistics.
- Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. 2023. Refiner: Reasoning feedback on intermediate representations. *arXiv preprint arXiv:2304.01904*.
- Ansh Radhakrishnan, Karina Nguyen, Anna Chen, Carol Chen, Carson Denison, Danny Hernandez, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilė Lukošiūtė, and 1 others. 2023. Question decomposition improves the faithfulness of model-generated reasoning. *arXiv preprint arXiv:2307.11768*.
- Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. *Preprint*, arXiv:2203.11171.

755

756

758

759

760

761

769

770

773

774

775

776

777

778

779 780

781

- Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. *Preprint*, arXiv:1704.05426.
- Xi Ye and Greg Durrett. 2022. The unreliability of explanations in few-shot prompting for textual reasoning. In *Advances in Neural Information Processing Systems*, volume 35, pages 30378–30392.
- Chenrui Zhang, Lin Liu, Chuyuan Wang, Xiao Sun, Hongyu Wang, Jinpeng Wang, and Mingchen Cai.
 2024. Prefer: Prompt ensemble learning via feedback-reflect-refine. Proceedings of the AAAI Conference on Artificial Intelligence, 38(17):19525– 19532.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo Li, and Yu Li. 2023. Progressive-hint prompting improves reasoning in large language models. *arXiv preprint arXiv:2304.09797*.

A Limitations

The EoT framework exhibits two principal limitations that warrant discussion. First, its assessing mechanism inherently relies on reference solutions to assess reasoning reliability and subsequently refines the reasoning processes based on these supervised signals. This paradigm creates a dependency on domain-specific answer templates, potentially constraining the framework's capacity to generalize and enhance LLM reasoning performance for out-of-domain tasks where authoritative references are unavailable or undefined.

Second, the *Hssim* metric implementation introduces an intrinsic dependency on the LLM's natural language inference (NLI) capabilities. The effectiveness of *Hssim* measurements becomes contingent upon the model's cross-domain NLI accuracy, introducing potential error propagation across different knowledge domains. To address this limitation, our future work will focus on integrating state-of-the-art LLMs with enhanced NLI datasets while developing domain-agnostic verification mechanisms to strengthen metric robustness.

B Examples of Reasoning Process Formalization

As formulated in Section Problem Setup, for each complex task (X, Q), an LLM expresses its reasoning behaviors in the form of a reasoning process *R*. This process consists of multiple thought steps and helps guide the LLM in producing the task's solution, denoted as \hat{Y} in Eq. (1) of our submitted manuscript. Specifically, LLMs that respond to questions by following a reasoning process operate by using specific prompting instructions. The prompt template used in these invocations is presented in Figure 4. This prompting strategy encourages LLMs to execute their genuine reasoning behaviors in a way that closely reflects selfexplanatory thoughts, with the goal of enhancing the influence of these thoughts on the LLMs' solution generation.

On this basis, to intuitively illustrate the selfexplained reasoning process, Figure 5 and 7 present two examples of reasoning processes generated by the widely used LLM Qwen2-72B. These examples are specifically aimed at guiding the LLM in addressing two distinct tasks in LongBench. 783

784

785

786

787

793

794

795

796

797

798 799 800

801 802 803

804

805

806

807 808

809

810

811 812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

/**Instruction to emphasize LLMs to sovie questions following self-explained thoughts**/
Hello, you are the intelligent robot of the smart Q&A project. Now, your task is to observe the
currently explained reasoning process, Context, Question, guide yourself to refer to the valid content
in Context based on the thinking and answering methodologies provided by the provided reasoning
process, and achieve an answer to the Question (i.e., generating the Answer) subsequently.
Please pay attention, you need to ensure the following requirements:
(1) In the process of generating the answer, please ensure as much as possible to think and deduce
answers step-by-step according to the provided reasoning process's STEP-BY-STEP Thought;
(2) Please ensure that you do not introduce new reasoning process or step of thoughts to the process
of reasoning the answer, and ensure that you do not revise the provided reasoning process. On this
basis, please ensure as much as possible that the thinking method or STEP-BY-STEP Thought used
in reasoning is consistent with the provided reasoning process. Specifically, please ensure that the
way you refer to/quote Context and the way of question answering is consistent with the provided
reasoning process.
(3) Please note that your task is to refer to or quote the text according to the thoughts of the reasoning
process to achieve question answering, and ultimately your output is the generated answer namely
Answer. Therefore, please note, whether the generated Answer is right or wrong, you can not offer
additional modification suggestions for the answer generated under the guideline of the provided
reasoning process, nor modify the generated answer."
/**Description of Complex Task**/
Given Context: < Domain Knowledge Context (X)>
Given Question: <question (q)="" description=""></question>
/** Given a Reasoning Process to guide LLMs**/
Your explained Reasoning Process: < Explained Reasoning Process (R)>
/**Instruction to Trigger the Question Answering**/
Following the thoughts and answering method of this reasoning process to answer the above
Question. Your Answer to the Question is:

Figure 4: An example prompt template that guides LLMs to generate task solutions under the guideline of specified self-explanatory reasoning processes.

C Illustration of Factuality, Fidelity and Reliability

In Section **Problem Setup** of our submitted manuscript, we define three critical factors that influence the reasoning capability of LLMs: factuality, fidelity and reliability. To aid in the intuitive understanding of the measurement of these three factors, this appendix presents both positive and negative examples of reasoning thoughts for each of them. For each factor, the positive examples demonstrate strong performance, while the negative ones illustrate poor performance.

C.1 Factuality

829

830

831

834

835

838

842

847

852

853

Figure 6 presents a completely factual reasoning process R_{fac} , alongside a reasoning process that contains non-factual errors R_{nonfac} . Specifically, all four thought steps in R_{fac} are factually grounded in the knowledge context. In contrast, T_2 and T_3 in R_{nonfac} contain errors that contradict established domain knowledge, as highlighted in the underlined part, which demonstrates that T_2 and T_3 exhibit poor factual performance. On this basis, Qwen2-72B generates two answers to the question "Which film came out first", denoted as \hat{Y}_{fac} and \hat{Y}_{nonfac} , under the guideline of R_{fac} and R_{nonfac} , respectively.

Then, we observe that, when Qwen2-72B responds to the question following the guideline of reasoning process R_{nonfac} which contains nonfactual date errors, non-factual hallucination may emerge in its generated solution \hat{Y}_{nonfac} . In contrast, a correct answer \hat{Y}_{fac} can be produced by LLMs under the guideline of R_{fac} , where the factuality of the explanatory reasoning thoughts is ensured. These clear examples highlight the importance of ensuring the factuality of reasoning processes during optimization, as this is crucial to prevent the reasoning capability of LLMs being compromised by non-factual hallucinations. 854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

C.2 Fidelity

We review the reasoning process R that guides Qwen2-72B in solving the question, "What are some fields in which the inverse problem is encountered?", as presented in Figure 7. This reasoning process consists of four steps of thought. To investigate whether each step faithfully reflects the actual reasoning behaviors of LLMs, we demonstrate the measurement of behavior fidelity for each reasoning step as defined in Eq.(4). Specifically, we mask thoughts T_1 and T_3 in R to obtain two masked rea-

Question (Q): Which film came out first, "Donne con le gonne" or "The Bag Man"?
Knowledge Context (X): Passage 1: "Bagman" (video game), a 1982 French platform arcade game "Begman" (film), a 2010 film about Jack Abramoff
Passage 3.
 Donne con le gonne Women in Skirts (Italian: Donne con le gonne) is <mark>a 1991 Italian romantic comedy</mark> film directed by Francesco Nuti
 Passage 7: The Bag Man (also known as Motel or The Carrier) is a 2014 neo-noir crime thriller film directed by
Reference Answer (A): Donne con le gonne.
Factual Reasoning Process (R): Thought T_1 : Task Identification: The task necessitates establishing which movie, "Donne con le gonne" or "The Bag Man", was released first by examining their respective release dates. Thought T_2 : Date Extraction: From Passage 3, I extract that "Donne con le gonne" was released in 1991. Meanwhile, Passage 7 indicates "The Bag Man" as a 2014 film. Thought T_3 : Temporal Sequence: Comparing the two years, it's evident that 1991 comes before 2014 confirming "Donne con le gonne" was released prior to "The Bag Man." Thought T_4 : Conclusion Formation: Based on the chronological order of their release dates, "Donne con le gonne" was released first.
Answer Guided by $R(\hat{Y} \text{ in Eq.}(1))$: "Donne con le gonne" was released first.

Figure 5: A factual reasoning process produced by Qwen2-72B for problem "Which film came out first?". The highlighted part represents the critical information and claims in the context and reasoning behaviors.

Question (Q): Which film came out first, "Donne con le gonne" or	Question (Q): Which film came out first, "Donne con le gonne" or "The Bag Man"?				
Knowledge Context (X): Passage 1: <u>"Bagman"</u> (video game), a <u>1982</u> French platform arcade game"Begman" <u>(film)</u> , a 2010 film about Jack Abramoff 					
Passage 3.					
 Donne con le gonne: Women in Skirts (Italian: Donne con le gonne) is <mark>a 1991 Italian romantic comedy film</mark> directed by Francesco Nuti					
 Passage 7: The Bag Man (also known as Motel or The Carrier) is <mark>a 2014 neo-noir crime thriller film</mark> 					
Factual Reasoning Process (R_{fac}) :	Non-Factual Reasoning Process (Rnonfac):				
Thought T_1 : Task Identification: The task necessitates	Thought T_1 : Task Identification: The task necessitates				
establishing which movie, "Donne con le gonne" or "The Bag	establishing which movie, "Donne con le gonne" or "The Bag				
Man", was released first by examining their respective release	Man", was released first by examining their respective release				
dates within the given texts.	dates within the given texts.				
"Donne con la conne" was released in 1001 Meanwhile Bassage 7	"Denne con la conne" was released in 1001. Meanwhile Bassage 1				
indicates "The Bag Man" as a 2014 film	indicates "The Bag Man" as a 1982 film				
Thought T_2 : Temporal Sequence: Comparing the two years.	Thought T_2 : Temporal Sequence: Comparing the two years.				
it's evident that 1991 comes before 2014, confirming "Donne con le	it's evident that 1982 comes before 1991, confirming "The Bag				
gonne' was released prior to "The Bag Man". Man' was released prior to "Donne con le gonne".					
Thought T_4 : Conclusion Formation: Based on the	Thought T_4 : Conclusion Formation: Based on the				
chronological order of their release dates, "Donne con le gonne"	chronological order of their release dates, "The Bag Man" was				
was released first.	released first.				
Correct Answer: "Donne con le gonne" was released first.	Hallucination: "The Bag Man" was released first.				

Figure 6: Examples of total factual and non-factual reasoning processes produced by Qwen2-72B for problem "Which film came out first?". The highlighted part in factual reasoning process represents the key factual claims, while the underlined part is the claims with non-factual errors which cause the hallucination in answers.



physics, biology, social sciences, finance, and neuroscience, demonstrating its wide applicability in understanding intricate systems.

Figure 7: A reasoning process of Qwen2-72B for the problem "What are some fields in which the inverse problem is encountered?". The highlighted part represents the critical information and claims in the context and reasoning behaviors.

soning processes $R \setminus T_1$ and $R \setminus T_3$ respectively, as discussed in subsection **Fidelity Scorer**. Figure 8 shows that Qwen2-72B generates two answers to the question under the guidelines of $RP \setminus T_1$ and $RP \setminus T_3$, denoted as $\hat{Y}_{\setminus T_1}$ and $\hat{Y}_{\setminus T_3}$, respectively. A comparison reveals that, relative to the answer guided by R, represented as \hat{Y} , a significant change is observed in $\hat{Y}_{\setminus T_1}$, while only a minor modification is noted in $\hat{Y}_{\setminus T_3}$.

In light of this, we can derive two insights. Firstly, the reasoning behavior exhibited in T_1 , namely "Extract Key Concepts", faithfully reflects a critical step in the actual reasoning actions of Qwen2-72B. This indicates that LLMs delve into the significance of each key entity, including "inverse problem" and its related fields, at the beginning of inference. In other words, T_1 significantly affects the reasoning results of LLMs, which indicates that T_1 exhibits strong behavior fidelity for Qwen2-72B. Secondly, the reasoning claim in T_3 , termed "Contextual Understanding", inadequately uncovers the actual behaviors that Qwen2-72B executes, as omitting these claims has only a minor impact on the LLM's reasoning results. Thus, reasoning thoughts such as T_3 can be regarded as unfaithful thoughts in the reasoning processes explained by LLMs. In other words, thoughts like T_3 demonstrate poor behavior fidelity for the LLM. 895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

C.3 Reliability

As demonstrated in Section **Problem Setup**, the reliability of a reasoning process for LLMs in addressing any task can be measured by the simi-



Figure 8: Examples of masked reasoning processes produced by Qwen2-72B for problem "What are some fields in which the inverse problem is encountered?". The highlighted part represents the critical information and claims in the context and reasoning behaviors. We exclude the claims of thoughts T_1 and T_3 respectively, and find that T_1 has a higher fidelity performance compared with T_3 .

larity between the generated solution and the reference answer. Since solutions often consist of long-textual statements, the assessment of similarity should occur at both syntactic and semantic levels. In addition, all measurements should aim to align with the understanding of experts in the relevant domain.

911

912

913 914

915

917

Intuitively, the reasoning processes presented in 918 Figure 5 and Figure 7 perform satisfied reliability 919 as they effectively guide Qwen2-72B in producing solutions that adequately address the respective 921 questions. Consequently, we should evaluate the high similarities between the produced and refer-923 ence solutions for these two reasoning processes 925 accordingly. Conversely, the assessed similarity should be significantly low for unreliable reasoning processes. For example, as shown in Figure 6, the non-factual reasoning process R_{nonfac} results in an incorrect answer to the question. Therefore, we 929

should accurately measure a low similarity between these hallucinations and the correct reference. In addition, as presented in Figure 8, a reasoning process excluding faithful thoughts, denoted as $R \setminus T_1$, will diminish the reasoning capabilities of LLMs, and results in the emergence of redundant statements in generated answers, which are unrelated to problem-solving. Thus, a low similarity needs to be measured to represent the decline of reliability in terms of the precision of extracting useful knowledge. 930

931

932

933

934

935

936

937

938

939

940

941

D Prompting Mechanism of EoT

To illustrate the design of prompting scheme in EoT942that drives the evolution of the reasoning process,943Figure 9 provides an example prompt in an evo-944lution iteration. This prompt is intended to guide945a specific optimization of the reasoning process946explained by Qwen2-72B for the question "What947

 (**Logic of Reasoning Process Evolution**) Hello, you are the intelligent robot of the smart Q&A project. Now, your task is to observe the given Context, Question, and Reference Answer, and refer to the existing reasoning process in the provided History and the quantitative scores of these reasoning processes to generate an optimized reasoning process that you believe can achieve higher scores. Specifically, each generated reasoning process is what you think you need to use or follow in terms of thoughts or logical reasoning process (i.e., Reasoning Process) when you approach the reliable Answer better as you answer the Question based on the Context. Please note the following requirements when generating the optimized reasoning process: Please note the following requirements when generating the optimized reasoning process: Please note that the reasoning process you output should be an optimized reasoning process provided in History, and should not repeat the reasoning process you output should be an optimized reasing reasoning process provided in History in three dimensions: Please note that the reasoning process the provided Contexts and the provided Reference Answer, the higher the better. The higher the score, the more reliably the corresponding reasoning process and uide the generated answer guided by reasoning process is at the target Answer. Weighted Fidelity score: ranging from 0 to 1, a fidelity score is given for each step of thoughts in the reasoning process, the higher the better. This score for each step of thoughts represents the similar by generate an answer similar to the Reference Answer. The higher the score, the generated answer guided by reasoning process is at the target Answer. Weighted Fidelity score: ranging from 0 to 1, a fidelity score is given for each step of thoughts in the reasoning process, the higher the score, the generated answer guided by reasoning process is at the target Answer. Weighted Fidelity
represents the degree of factuality of the content expressed in this thought. Specifically, the higher the score, the more the reasoning basis of the thought comes from the reference or understanding of the provided Context or Question, accordingly, the higher the score, the lower the probability that the reasoning basis of the thought comes from your own guess and the lower the probability of the reasoning basis deviating from the provided Context and Question. 3. Please ensure that the optimized reasoning process as a whole can achieve a higher overall reliability score. Under this premise, please try to improve the weighted fidelity score and factuality score of each step of thought in your optimized reasoning process. Also, please note that the number and titles of thoughts you generate do not need to match those in the existing reasoning process in History. 4. Please note that the purpose of generating a reasoning process is to guide the large model to provide an answer that reliability solve the Question and ensure the factuality and fidelity to answer the Ouestion, nor to provide any optimization suggestions for the target Answer.
6. When generating the optimized reasoning process, please note that you need to output both the optimized reasoning process. The specific text content of each reasoning process is a sequence of multiple steps of thoughts, and each step of thoughts needs to provide a title
/**December Tack**/
(Omitting. Please Reiew Figure 6)
/** Textual Description of Assessment of Existing Reasoning Processes**/ History: Existing Reasoning Process R_1 : Thought T_1 : Extract Key Concepts: We learn about the inverse problem, which involves statistical inference for systems like the XY model Factuality Score of T_1 : 1.0; Weighted Fidelity Score of T_1 : 0.5265 Thought T_2 : Relate Inverse Problem to Fields: In physics, inverse problems are crucial for understanding complex systems. Biology utilizes these Factuality Score of T_2 : 1.0; Weighted Fidelity Score of T_2 : 0.1659 Thought T_3 : Contextual Understanding: The inverse problem's relevance across these fields underscores its importance in extracting meaningful information Factuality Score of T_3 : 1.0; Weighted Fidelity Score of T_3 : 0.1427 Thought T_4 : Formulate Answer: Based on the the text, the inverse problem is encountered in diverse fields such as physics, biology, social sciences, finance, Factuality Score of T_4 : 1.0; Weighted Fidelity Score of T_4 : 0.1724 Reliability Score of R_1 : 0.8843 Though Taise the Difference of Taise that is the state the text of the problem is encountered in the text is the text of the state the text of the state the text of the text of the state text of text

Figure 9: Instance of a specific complete prompt to evolve the reasoning process in Figure 7, involving the instructions of evolution and the textual description of assessment results of reasoning process shown in Figure 7

are some fields in which the inverse problem is encountered?". Specifically, this prompt follows the template illustrated in Figure 2 of our submitted manuscript, and includes a textual description of the assessment results from an existing reasoning process, as shown in Figure 7.

948

950

951

952

954

955

957

960

961

962

963

964

965

968

E NLI mechanism in Assessment of Reasoning process

In EoT, we prompt LLMs to perform natural language inference (NLI) judgements to evaluate the effectiveness of reasoning processes. Overall, leveraging the NLI of LLMs is intended to improve semantic alignment with the human-level perception, particularly when estimating the extent of estimation of knowledge entailment and assessing similarities. This approach further enhances the robustness of the scoring scheme in EoT. In this appendix, we first define the triple labels of the statement-pair relationships. Then, we evaluate the effectiveness of NLI mechanism when using various types of LLMs and discuss the threats to the EoT framework given the use of the GPT-4 model.

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

E.1 Definition of Triple Labels in NLI Mechanism

In Section 3.2 in our submitted manuscript, EoT regards the NLI judgements as triple-label classification tasks as shown in Eq.(6). Specifically, let x^{pre} and t^{hyp} denote a premise and hypothesis statement, respectively. In a canonical NLI mechanism, there are three category labels to represent the relationship between any statement pair (x^{pre}, t^{hyp}) , which is defined as follows.

Definition 4 Entailment: If hypothesis statement t^{hyp} is necessarily true or appropriate whenever the premise statement x^{pre} is true, we label the relationship between the statement pair (x^{pre}, t^{hyp}) as "Entailment". In other words, x^{pre} entails t^{hyp} .

Definition 5 Contradiction: If hypothesis statement t^{hyp} is necessarily false or inappropriate whenever the premise statement x^{pre} is true, we label the relationship between the statement pair (x^{pre}, t^{hyp}) as "Contradiction". In other words,

Table 6	5: [The .	Accuracy	v of	NLI	using	LLMs

LLM Model	Accuracy(%)
Qwen2-72B	76.14
GPT-3.5 Turbo	67.64
GPT-4 Turbo	85.80

 t^{hyp} contradicts with x^{pre} .

993

995

996

997

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1016

1017

1018

1020

1021

1022

1024

1025

1026

1027 1028

1029

1030

1032

Definition 6 Neutral: When neither "Entailment" nor "Contradiction" applies to relationship between the statement pair (x^{pre}, t^{hyp}) , we label this relationship between (x^{pre}, t^{hyp}) as "Neutral". In other words, x^{pre} is neutral with t^{hyp} .

E.2 Performance of NLI using Various Types of LLMs

To evaluate the effectiveness of NLI in LLMs using few-shot prompts, we conduct experiments on the MNLI dataset (Williams et al., 2018). This dataset is well-known and comprises 40, 000+ samples of NLI tasks collected from dozens of different domains, including transcribed speech, fiction, and reports. For our evaluation, we select 10, 000 representative samples from the MNLI dataset, considering the scale and overhead of experiments. On this basis, we evaluate the NLI performance of three widely used LLMs: Qwen2-72B, GPT-3.5 Turbo and GPT-4 Turbo, using the selected samples.

Table 6 presents the accuracy performance of NLI achieved by the three LLMs. It is evident that GPT-4 Turbo outperforms the others in NLI tasks that require context from various domains. Moreover, based on previous studies (Gao et al., 2023) and our experiment results, we observe that GPT-4 Turbo achieves state-of-the-art accuracy in NLI. This suggests that using GPT-4 Turbo for NLI can lead to superior alignment between the judgements of LLMs and human assessments. Consequently, EoT employs GPT-4 Turbo with tailored few-shot learning to perform NLI by default, aiming to better align with human judgement.

Finally, we discuss the limitations of the NLI mechanism based on the GPT-4 Turbo. As we can see, the classification accuracy of NLI in EoT still needs to be improved, which could result in threats to the precise alignment of the assessment of EoT with that of humans. To further enhance EoT's capabilities, we intend to improve NLI performance in EoT by integrating novel LLMs with advanced NLI capabilities or optimizing the few-shot prompting mechanism.

	Table 7:	Statistics	of Contex	xt Length	of Selected	Tasks
--	----------	-------------------	-----------	-----------	-------------	-------

Dataset	Context (token)				
Dataset	Mean	P95	Max		
OpsQA	6,052	11,640	21,045		
LongBench (selected)	5,700	11,782	14,640		

F Description of Context Complexity

In our evaluation, we select 40 questions from the OpsQA dataset and 60 questions from the Long-Bench dataset respectively. Resolving each question is a representative reasoning task, which needs to comprehend a series of long-textual domain knowledge contexts. In this appendix, we present a set of statistics related to the context length of selected questions, including the mean, the 95th percentile (P95) and the maximum (Max) value, to illustrate the complexity of tasks in our evaluation.

As presented in Table 7, it can be found that, the average length of knowledge context can achieve 6, 052 and 5, 700 tokens for selected tasks in OpsQA and LongBench datasets, respectively. Moreover, the maximum context length can even be up to 21, 045 and 14, 640 tokens for the selected task in OpsQA and LongBench datasets respectively. This significant complexity of the domain knowledge context for tasks could validate that, our experiments could appropriately evaluate the effectiveness of EoT for addressing the complex reasoning tasks that require the comprehension of long-context domain knowledge and the generation of intricate solutions.

G Description of Baselines

We compare EoT with five state-of-the-art frame-1059 works for reasoning process evolution. Accord-1060 ing to the schema of reasoning process generation, 1061 these frameworks can be categorized into three 1062 groups. 1) Exploration of structured thoughts: BoT 1063 (Chen et al., 2024) guides LLMs to explore ensem-1064 ble of trees of thoughts (ToTs) and acquire trial-1065 and-error reasoning experiences for each explored 1066 ToT, with the aim of improving the reliability of 1067 reasoning processes. In BoT, a reasoning thought 1068 is generated by the LLM as a node of ToT in each 1069 inference, and the complete reasoning process is 1070 formed by a series of such LLM inferences. 2) 1071 Revision of the complete reasoning process: a) 1072 xLLM (Chuang et al., 2024) provides feedback 1073 on the fidelity assessment for the complete rea-1074 soning process and guides LLMs to optimize rea-1075 soning processes through iterative improvements, 1076

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056



Figure 10: The CDF results of the reliability scores of evolved reasoning processes for various complex questions

aiming to enhance fidelity performance; b) Calibra-1077 tor (Ye and Durrett, 2022) focuses on revising rea-1078 soning processes to improve factuality and further 1079 enhance the reasoning capability of LLMs through iterations. In these studies, LLMs generate a complete reasoning process in textual paragraphs via 1082 1083 a single-forward inference. 3) Question decomposition: a) CoT-dec (Radhakrishnan et al., 2023) 1084 prompts LLMs to decompose the reasoning pro-1085 cess to a sequence of subquestion-subanswer pairs 1086 produced in single-forward inference in each itera-1087 1088 tion. It seeks decomposed results that achieve the best reliability performance as the final reasoning 1089 process after multiple iterations; b) Factor-dec is 1090 a variant of CoT-dec, aiming to enhance the reli-1091 ability of reasoning processes. In each evolution iteration, Factor-dec guides LLMs to decompose a reasoning process into a subquistion sequence 1094 and prompts LLMs to answer these subquestions 1095 through step-by-step inferences. In other words, 1096 each evolution iteration of Factor-dec involves mul-1097 tiple steps of question-answering in generating a 1098 reasoning process. 1099

H Experiment Results of Diverse Questions

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

In this appendix, we intuitively present the performance of the six evolution frameworks applied to various tasks within the OpsQA and LongBench datasets, utilizing two LLMs Qwen2-72B and GPT-4 Turbo. The results demonstrate that EoT significantly enhances the reasoning capabilities of diverse LLMs, as well as behavior fidelity and factuality of the explained reasoning processes across a broad spectrum of tasks.

H.1 Reasoning Capability (Reliability)

1112Figure 10a ~ Figure 10d present the CDF results1113of reliability scores of reasoning processes evolved1114by the six frameworks on the two LLMs for all1115questions in the two datasets. Notably, EoT demon-1116strates superior reasoning capabilities compared

to the five baselines for most tasks across diverse 1117 LLMs. For instance, in the case of the 40 questions 1118 in the OpsQA dataset, Figure 10a shows that EoT 1119 achieve a reliability score exceeding 0.6 for 85.0%1120 of the questions. This performance surpasses that 1121 of BoT, xLLM, Calibrator, CoT-dec and Factor-dec 1122 by about 75.0%, 35.0%, 22.5%, 15.0% and 40.0%, 1123 respectively, when using Qwen2-72B. As for the 1124 60 questions in the LongBench dataset, when us-1125 ing Qwen2-72B, Figure 10c illustrates that EoT 1126 achieves a reliability score exceeding 0.8 for about 1127 61.67% of the questions. This performance sur-1128 passes that of BoT, xLLM, Calibrator, CoT-dec 1129 and Factor-dec by about 43.43%, 33.34%, 43.34%, 1130 21.67% and 38.34%, respectively. In addition, 1131 when employing GPT-4 Turbo, Figure 10d indi-1132 cates that EoT attains a reliability score higher than 1133 0.8 for roughly 63.3% of these 60 questions, which 1134 exceeds the performance of BoT, xLLM, Calibrator, 1135 CoT-dec and Factor-dec by roughly 45.0%, 20.0%, 1136 34.67%, 16.67% and 31.67%, respectively. 1137

In summary, these results confirm that EoT achieves exceptional reliability across a wide range of tasks. This suggests that the improvement of reasoning capabilities in LLMs, facilitated by EoT, is well-suited for problem-solving in various domains with leading generality.

1138

1139

1140

1141

1142

1143

1144

H.2 Factuality

Figure $11a \sim$ Figure 11d present the CDF results 1145 of factuality scores of reasoning processes evolved 1146 by the six frameworks on the two LLMs for all 1147 the questions in the two datasets. We find that, 1148 compared with the five baselines, EoT remark-1149 ably optimizes the factuality of reasoning processes 1150 for a wider range of questions on diverse LLMs. 1151 Specifically, as for the 40 questions in OpsQA, 1152 Figure 11b illustrates that EoT achieves the upper 1153 limit of factuality score 1.0 for about 50.0% of 1154 the questions, which exceeds the performance of 1155 BoT, xLLM, Calibrator, CoT-dec and Factor-dec 1156 by about 15.0%, 47.5%, 20.0%, 35.0% and 42.5%, 1157



(a) OpsQA + Qwen2-72B (b) OpsQA + GPT-4 Turbo (c) LongBench + Qwen2-72B (d) LongBench + GPT-4 Turbo Figure 11: The CDF results of the factuality scores of evolved reasoning processes for various complex questions



Figure 12: The CDF results of the fidelity scores of evolved reasoning processes for various complex questions

respectively when using GPT-4 Turbo. As for 1158 the 60 questions in the LongBench dataset, when 1159 using Qwen-72B, Figure 11c illustrates that EoT 1160 achieves factuality scores of 1.0 for about 73.3%1161 of the questions, which exceeds the performance 1162 1163 of BoT, xLLM, Calibrator, CoT-dec and Factor-dec by about 25.0%, 20.0%, 3.3%, 13.3% and 26.7%, 1164 respectively. Moreover, when using GPT-4 Turbo, 1165 Figure 11d shows that EoT obtains factuality scores 1166 of 1.0 for 75% of these 60 questions, which sur-1167 passes the performance of BoT, xLLM, Calibrator, 1168 CoT-dec and Factor-dec by about 26.7%, 36.7%, 1169 3.3%, 5.0% and 25.0%, respectively. 1170

These results suggest that EoT demonstrates remarkable generality in ensuring the factuality of reasoning processes, proving to be highly robust across various knowledge domains. This assurance stems from the detailed evaluation of thoughts and the effective evolving framework that enables LLMs to eliminate non-factual errors during the explanation of these thoughts.

H.3 Fidelity

1171

1172

1173

1174

1175

1176

1177

1178

1179

Figure 12a \sim Figure 12d illustrate the behavior fi-1180 delity performance of reasoning processes evolved 1181 by the six frameworks on the two LLMs for all 1182 the questions in the two datasets. Compared to 1183 the five baselines, EoT implements an evolution 1184 1185 mechanism that guides diverse LLMs in generating reasoning processes, resulting in state-of-the-art 1186 performance in behavior fidelity across most tasks 1187 from various knowledge domains. These enhance-1188 ments demonstrate the exceptional generality of 1189

EoT in improving behavior fidelity. Specifically, as 1190 for the 40 questions in the OpsQA dataset, when 1191 using Qwen2-72B, Figure 12a illustrates that EoT 1192 obtains fidelity scores higher than 0.225 for 71.7%1193 of the questions, which exceeds the performance 1194 of BoT, xLLM, Calibrator, CoT-dec and Factor-1195 dec by about 36.7%, 53.3%, 5.0%, 58.3% and 1196 43.3%, respectively. Moreover, when using GPT-4 1197 Turbo, Figure 12b presents that EoT achieves fi-1198 delity scores higher than 0.25 for 60.0% of these 1199 40 questions, which surpasses the performance 1200 of BoT, xLLM, Calibrator, CoT-dec and Factor-1201 dec by roughly 32.5%, 20.0%, 25.0%, 50.0% and 1202 37.5%, respectively. As for the 60 questions in the 1203 LongBench dataset, when using Qwen2-72B, Fig-1204 ure 12c presents that EoT achieves fidelity scores 1205 higher than 0.2 for 40% of the questions, which 1206 surpasses the performance of BoT, xLLM, Cali-1207 brator, CoT-dec and Factor-dec by roughly 10%, 1208 26.7%, 26.7%%, 21.7% and 6.7%, respectively. In 1209 addition, when using GPT-4 Turbo, Figure 12d il-1210 lustrates that EoT obtains fidelity scores higher 1211 than 0.2 for 75.0% of the questions, which exceeds 1212 the performance of BoT, xLLM, Calibrator, CoT-1213 dec and Factor-dec by about 16.7%, 41.7%, 36.7%, 1214 51.7% and 26.7%, respectively. 1215

These results demonstrate that EoT significantly1216improves the fidelity of reasoning processes as explained by various LLMs, and this improvement is1217plained by various LLMs, and this improvement is1218broadly applicable across diverse tasks rooted in1219different domains of knowledge. In summary, this1220enhancement in behavior fidelity can be attributed1221

to two key factors. Firstly, building on prior studies, EoT implements an effective mechanism for assessing the behavior fidelity of explained reasoning thoughts in a detailed manner. Secondly, the evolution mechanism of EoT effectively encourages LLMs to clarify their reasoning with enhanced faithfulness.

1222

1223

1224

1225

1227

1228

1229

1230

1231

1232

1233

1234

1235

The optimization of reasoning processes regarding factuality and fidelity, as shown in Appendices H.2 and H.3, is beneficial for further enhancing the reasoning capabilities of LLMs, as illustrated in Appendix H.1.

Table 8: The performance on reliability for reasoning processes evolved by EoT and its variants using two LLMs on the OpsQA dataset.

		OpsQA					
LLMs	Models	BLEURT	ROUGE-L	NLI results	Reliability		
				entail(%)	Score		
Qwen2	EoT_w/o_fact	0.593	0.362	70.16	0.607		
	EoT_w/o_fide	0.602	0.377	76.32	0.651		
	EoT_w/o_reli	0.609	0.386	77.33	0.655		
	EoT	0.595	0.395	83.90	0.706		
GPT-4	EoT_w/o_fact	0.573	0.403	78.75	0.658		
	EoT_w/o_fide	0.597	0.417	80.03	0.675		
	EoT_w/o_reli	0.593	0.402	75.84	0.634		
	FoT	0.602	0.427	88 58	0 717		

¹ Qwen2 and GPT-4 represent the LLMs of Qwen2-72B and GPT-4 Turbo, respectively.

Table 9: The performance on factuality and fidelity for reasoning processes evolved by EoT and its variants using two LLMs on the OpsQA dataset.

		OpsQA			
LLMs	Models	Factuality	Fidelity		
		S_{fac}	S_{fid}		
Qwen2	EoT_w/o_fact	0.737	0.231		
	EoT_w/o_fide	0.840	0.226		
	EoT_w/o_reli	0.723	0.272		
	EoT	0.823	0.267		
GPT-4	EoT_w/o_fact	0.722	0.277		
	EoT_w/o_fide	0.869	0.262		
	EoT_w/o_reli	0.882	0.296		
	EoT	0.906	0.281		

I Ablation Study

I.1 Setup

We conduct ablation studies for EoT on the OP-1236 SQA dataset. To evaluate the effectiveness of the 1237 optimization achieved by EoT concerning the three 1238 factors, namely reliability, factuality and behavior 1239 fidelity of reasoning processes, we design three 1240 variants of EoT: EoT_w/o_fact, EoT_w/o_fide, and 1241 1242 EoT_w/o_reli. Each variant removes the optimization related to factuality, behavior fidelity, and reli-1243 ability respectively from the evolution of thoughts. 1244 We then compare the performance of the canonical 1245 EoT with that of these three variants. 1246

Table 8 presents the reasoning capability of Qwen2-72B and GPT-4 Turbo equipped with the canonical EoT and its three variants. Additionally, Table 9 shows the fidelity and factuality performance of the reasoning processes evolved by the EoT and its variants.

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1288

1289

1290

1291

1292

1293

1294

1296

I.2 Effectiveness of Factuality Optimization

We compare the performance of EoT with that of EoT_w/o_fact to assess the effectiveness of evolution in terms of factuality. First of all, it can be observed that EoT significantly enhances the factuality of reasoning processes through the finegrained evolution explicitly optimizing factuality. For instance, Table 9 shows that, when using GPT-4 Turbo, EoT increases the factuality score by about 25.5% compared to EoT_w/o_fact. Furthermore, the experiment results validate that the optimization on the factuality of reasoning processes in EoT further enhances the reasoning capabilities of diverse LLMs. For instance, Table 8 demonstrates that compared to EoT_w/o_fact, EoT improves BLEURT, ROUGE-L and NLI result entail(%)by about 5.1%, 6.0%, and 12.5% respectively, and finally obtains the improvement of reliability score by about 9.0% when using GPT-4 Turbo.

I.3 Effectiveness of Fidelity Optimization

We compare the performance of EoT with that of EoT_w/o_fide, and then there are two insights can be observed, which indicates the effectiveness of behavior fidelity evolution in EoT. Firstly, compared to EoT_w/o_fide, EoT produces reasoning processes that achieve a significantly improved behavior fidelity. Specifically, as shown in Table 9, when using Qwen2-72B, EoT increases the fidelity score by about 18.1%. This enhancement is attributed to the fine-grained assessment of fidelity of reasoning thoughts and the effective prompting scheme aimed at faithfully explaining the behaviors of LLMs.

Secondly, the mechanism that enhances behavior fidelity in EoT facilitates the positive emergence of reasoning capabilities in LLMs. For instance, compared to EoT_w/o_fide, EoT increases the reliability score by about 8.5% when using Qwen2-72B. In detail, this improvement is mainly attributed to that the fidelity enhancement achieved by EoT simultaneously enhances the reasoning capabilities of LLMs from aspects of token overlapping and semantic alignment with human judgement. Specifically, EoT improves ROUGE-L and NLI result

LLMs		OpsQA LongB				Bench			
	Models	DIEUDT	POLICE I	NLI results	Reliability	Reliability Score BLEURT	ROUGE-L	NLI results	Reliability
		BLEUKI	KOUGE-L	entail(%)	Score			entail(%)	Score
Qwen2	BoT	8.602×10^{-5}	3.720×10^{-6}	6.449×10^{-10}	6.463×10^{-10}	6.007×10^{-9}	2.843×10^{-7}	3.143×10^{-6}	2.308×10^{-9}
	xLLM	0.011	0.003	1.362×10^{-5}	3.263×10^{-7}	1.561×10^{-9}	1.598×10^{-7}	0.026	0.001
	Calibrator	0.028	0.036	0.006	9.268×10^{-6}	1.032×10^{-10}	8.322×10^{-8}	0.003	3.122×10^{-7}
	CoT-dec	0.023	0.043	0.008	0.043	1.690×10^{-7}	6.397×10^{-5}	0.221	0.038
	Factor-dec	0.014	1.343×10^{-5}	1.657×10^{-6}	2.380×10^{-7}	5.173×10^{-9}	2.270×10^{-6}	0.030	0.003
GPT-4	BoT	9.132×10^{-7}	1.546×10^{-7}	2.110×10^{-9}	5.074×10^{-11}	5.971×10^{-9}	2.993×10^{-9}	1.892×10^{-6}	3.809×10^{-9}
	xLLM	0.036	0.045	5.120×10^{-8}	3.042×10^{-8}	0.047	0.034	0.013	9.875×10^{-4}
	Calibrator	0.034	0.037	3.330×10^{-4}	2.664×10^{-4}	4.896×10^{-5}	1.813×10^{-4}	7.135×10^{-4}	6.534×10^{-6}
	CoT-dec	0.014	0.446	0.039	0.040	0.005	0.041	0.235	0.043
	Factor-dec	0.001	1.073×10^{-8}	0.001	7.063×10^{-4}	5.572×10^{-6}	3.859×10^{-5}	$3.864 imes 10^{-4}$	0.001

Table 10: The p-value of performance improvement or decrease in terms of reasoning capability of LLMs between the five baselines and EoT in T-test.

entail(%) by about 4.8% and 9.9% respectively when using Qwen2-72B.

I.4 Effectiveness of Reliability Optimization

1297

1298

1299

1300

1302

1303

1304

1305

1306

1307

1308

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1322

1323

1324

1325

1327

1328

1331

1332

1333

1334

1335

We assess the performance of EoT in comparison to EoT_w/o_reli to determine the effectiveness of the evolution incorporated in the canonical EoT, which aims to directly improve the overall reliability of reasoning processes. As shown in Table 8, EoT increases the reliability score by about 13.1% when using GPT-4 Turbo, compared to EoT_w/o_reli. For further details, EoT improves the BLEURT, ROUGE-L and entail(%) by about 1.5%, 6.2% and 16.8% respectively. These findings suggest that explicitly addressing the reliability factor can lead to a significant improvement in the overall reasoning capabilities of LLMs, in contrast to the unsupervised evolution that does not consider the solving reliability under the guideline of previously produced reasoning processes.

Effectiveness of Multi-objective I.5 **Optimization on Three Factors**

According to the experiment results of ablation studies as mentioned above, it can be observed that EoT achieves the best collaborative optimization among the three factors compared with the three variants. In other words, when considering the three factors simultaneously, EoT produces the evolved reasoning processes that achieves the best trade-off among the three factors under the promise of reasoning capability improvement on LLMs.

Firstly, excluding a factor from EoT can make evolved reasoning processes achieve a better performance on some other factors. However, this improvement of performance on other factors is always along with the degradation of the performance on the excluded factor, and would further compromise reasoning capabilities of LLMs. For instance, as shown in Table 8 and Table 9, when using Qwen2-72B, EoT_w/o_fide improves the factuality performance by about 2.1% compared to EoT. However, EoT attains the improvement on fidelity performance and the overall reasoning capabilities by about 18.1% and 8.5% respectively.

1336

1337

1338

1339

Secondly, evolving reasoning process in terms 1340 of fidelity and factuality in a unsupervised way, 1341 can effectively optimize the behavior fidelity of 1342 reasoning processes produced by diverse LLMs. 1343 Nevertheless, EoT which evolves reasoning pro-1344 cess with supervised awareness of reliability perfor-1345 mance further enhance the factuality of produced 1346 reasoning processes and the reasoning capability of 1347 LLMs simultaneously. Specifically, compared to 1348 EoT, EoT w/o reli improves the fidelity score by 1349 about 1.9% and 5.3% when using Qwen2-72B and 1350 GPT-4 Turbo, respectively. In contrast to that, EoT 1351 improves the factuality score and reliability score 1352 by about 13.8% and 7.8% when using Qwen2-72B, 1353 and by about 2.7% and 13.1% when using GPT-4 1354 Turbo, respectively. 1355

lines and EoT in T-test. LongBench OpsQA LLMs Models Fidelity Factuality Factuality Fidelity S_{fac} Sfid Sfid Bol 0.039 0.041 0.017 0.045 xLLM 0.003 1.215×10 1.590×10^{-1} 047×10 4.258×10^{-1} Calibrator Owen2 0.046 0.07 2.550×10^{-4} 3.172×10^{-1} CoT-dec 1.142×10^{-9} 0.007 0.002 Factor-dec 8.081×10^{-5} 1.488×10^{-1} 0.0450.032 0.001 2.604×10^{-6} 0.009 xLLM 3.667×10^{-5} 1.010×10^{-1} 2.988×10 1.470×10^{-1} GPT-4 Calibrator 0.039 0.041 0.087 6.041×10^{-10} 1.644×10 8.337×10 1.372×10

 4.683×10^{-10}

CoT-dec

Factor-dec

 8.267×10^{-5}

Table 11: The p-value of performance improvement or

decrease in terms of factuality and fidelity for reasoning processes of reasoning processes between the five base-

J **Significance of Performance** 1356 **Improvement or Decrease** 1357

0.062

 8.401×10^{-6}

In this appendix, we use the T-test, an appropriate 1358 method of statistical test, to evaluate the signif-1359 icance of performance improvement or decrease presented in Table 3 and Table 4 in our submitted 1361

manuscript. In general, for each of the five base-1362 lines, we compute the p-value in T-test between 1363 the distribution of each performance metric ob-1364 tained by the baseline and EoT respectively when 1365 using each type of LLMs on each dataset. Specif-1366 ically, the performance metric includes BLEURT, 1367 ROUGE-L, NLI score and reliability score for 1368 reasoning capability evaluation and the factuality 1369 1370 score and fidelity score for the evaluation of factuality and behavior fidelity of reasoning processes. 1371 When computing the p-value of any metric for each 1372 baseline on a specific dataset, the performance dis-1373 tribution is estimated among the metric value for 1374 each question in the dataset. A lower p-value rep-1375 resents a higher significance for performance im-1376 provement and decrease, and p-value < 0.05 often 1377 indicates that the evaluation results achieve the sufficient significance in statistical tests. 1379

1380

1382

1384

1385

1386

1388

1390

1391

1392

1393

1394

1395 1396

1397

1398

1400

1401

1402

1403

1404

1405

1406

1407

1408 1409

1410

1411

1412

1413

As for the performance metrics reflecting the reasoning capability of LLMs, namely BLEURT, ROUGE-L, NLI results and reliability score, As presented in Table 10, the p-value between performance achieved by the baselines and that obtained by EoT can be limited below 0.05 in most scenarios. These results indicate that the performance improvement or decrease achieved by EoT in terms of reasoning capability of LLMs which are presented in our submitted manuscript, have the outstanding significance. Thus, our conducted evaluations and the conclusion that EoT effectively achieves the enhancement of reasoning capability have a reasonable robustness. Nevertheless, we find that p-value significantly exceeding 0.05 only occurs between the value of NLI results achieved by the CoT-dec and EoT respectively on LongBench dataset. This is mainly attributed to that CoT-dec and EoT obtains the similar performance of reasoning capability on the aspects of semantic alignment robustness on LongBench dataset. Since EoT achieves adequate significance of performance improvement for all the remaining metrics i.e., BLEURT, ROUGE-L and reliability score, it still can be regarded that EoT enhances the reasoning capability of LLMs compared with CoT-dec in a robust way on Long-Bench dataset. In future, we will evaluate the performance difference between NLI results of CoTdec and EoT on a larger scale of datasets.

As for the performance metrics mirroring factuality and behavior fidelity of evolved reasoning processes, It can be observed in Table 11 that, on OpsQA dataset, p-value ≤ 0.05 occurs for each metric achieved by the five baselines in our evaluation. Nevertheless, on LonBench dataset, we find 1414 that p-value for factuality score achieved by Cal-1415 ibrator, CoT-dec and Factor-dec slightly exceeds 1416 0.05 when using GPT-4 Turbo. This is because 1417 that these three baselines and EoT all obtain the 1418 outstanding but close performance on factuality of 1419 reasoning processes. We will use more open-source 1420 dataset to further evaluate the performance differ-1421 ence on factuality among these four frameworks in 1422 the future. 1423