# ARIES: Autonomous Reasoning with Large Language Models on Interactive Thought Graph Environments

**Anonymous ACL Submission** 

#### Abstract

Recent research has shown that the performance of Large Language Models (LLMs) on reasoning tasks can be enhanced by scaling testtime compute. One promising approach, particularly with decomposable problems, involves arranging intermediate solutions as a graph on which transformations are performed to explore the solution space. However, prior works rely on pre-determined, task-specific transformation schedules defined by a set of searched hyperparameters. In this work, we view thought graph transformations as actions in a Markov decision process, and investigate whether LLMs can act as a policy agents on this graph-based environment. We introduce ARIES, a multiagent architecture for reasoning with LLMs in which LLM policy agents maintain visibility of the thought graph states and dynamically adapt the problem-solving strategy. We observe that using off-the-shelf LLMs as policy agents with ARIES can yield up to 29% higher accuracy on HumanEval relative to static transformation schedules despite bearing no search cost.

## 1 Introduction

001

011

017

019

021

024

025

027

034

039

042

Recent work showed that under a fixed compute budget for training and inference, LLM performance on reasoning tasks can be enhanced by allocating a higher proportion of compute to inference rather than training (Snell et al., 2024). This has been achieved using techniques such as Chain-of-Thoughts or Tree-of-Thought prompting (Wei et al., 2023; Yao et al., 2023a). Another promising approach involves arranging intermediate solutions (or "thoughts") in a graph, i.e. topological reasoning (Besta et al., 2024b). This proves particularly beneficial when problems can be decomposed into subproblems to be solved independently then aggregated through a sequence of graph transformations (Besta et al., 2024a), and intermediate solutions can be reliably scored using a Process Reward Model (Snell et al., 2024). Despite the benefits

of topological reasoning, prior works rely on predetermined traversal strategies parametrized by a discrete set of hyperparameters (Yao et al., 2023a; Besta et al., 2024a). This approach lacks generality, as these parameters must be tuned manually or through extensive Bayesian search to achieve high query efficiency, due to the varying characteristics of each task. To overcome this limitation, we propose viewing thought graphs as an interactive environment where graph transformations are seen as actions in a Markov Decision Process (MDP). As such, we turn to investigating action policies to effectively explore the solution space and yield a solution while learning from external feedback. 043

045

047

049

051

054

055

057

059

060

062

063

064

065

066

067

068

069

070

071

072

073

074

077

078

081

Motivated by recent improvements in LLMs, we aim to investigate whether their planning and reasoning abilities extend to selecting and executing a sequence of transformations (such as decomposition, solving, evaluation, refinement and aggregation) within the aforementioned thought graph environments. We implement ARIES, a multi-agent framework for solving reasoning problems formulated as thought graphs.

Figure 1 provides a summary of our approach - in each iteration, the policy agent monitors the thought graph state and samples from the action space to choose a graph transformation. The reasoning agent then performs these transformations and updates the thought graph state. Through a series of carefully controlled experiments against a number of benchmarks, we show that LLM-guided thought graph exploration can lead to up to 29% higher accuracy, as well as obviating search cost.

# 2 Topological Reasoning with Large Language Models

Given a reasoning problem defined as an ordered tuple of tokens  $p = (t_1, \ldots, t_m)$ , we define a thought  $\tau = (t_1, \ldots, t_j)$  as a sequence of tokens sampled autoregressively from an LLM parametrized by  $\theta$ ,



Figure 1: ARIES workflow in answering the HumanEval prompt: "Find the shortest palindrome that begins with a supplied string". First, the policy agent selects the split action, guiding the reasoning agent to decompose the problem by generating a skeleton implementation with yet-to-implement subfunctions. Since one of the solutions doesn't pass its testcases, the reasoning agent is instructed to refine it based on execution feedback.

i.e.  $t_i \sim P(t_i \mid t_1, \ldots, t_{i-1}; \theta)$ . This consists of a language representation of an intermediate step towards the solution. A thought sequence can hence be represented as an ordered tuple of thoughts  $S = (\tau^1, \tau^2, \dots, \tau^k)$  of length k, such that the final thought  $\tau^k$  represents a candidate solution to the problem p. A thought graph  $G_{\tau}$  can be represented as  $(\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is a set of thought nodes and  $\mathcal{E}$  is a set of edges connecting them. Additionally, each thought  $\tau^{ij}$  (*j*-th thought at depth *i*) has a value  $\lambda(\tau^{ij})$  such that nodes with higher values yield valid solutions to the problem with higher probability. Thought graph exploration can be regarded as a sequence of m graph transformations as follows, where each  $\phi_i: G^i_\tau \to G^{i+1}_\tau$  modifies the set of nodes and edges.

084

099

100

102

103

104

105

106

108

110

$$G_{\tau}^* = \phi_m(\dots(\phi_1(\phi_0(G_{\tau}^0)))) \tag{1}$$

We consider a finite set of transformations;  $\phi_{dec}$ decomposes a reasoning problem into subproblems to be solved individually, creating new nodes in the thought graph.  $\phi_{sol}$  generates a candidate solution to a subproblem.  $\phi_{ref}$  considers an incorrect subproblem solution, utilizing further LLM queries to refine it.  $\phi_{red}$  removes nodes in the graph according to their values. Finally,  $\phi_{agg}$  performs node merging to aggregate subproblem solutions into a coherent solution to the original problem. Formal definitions of each transformation are shown in Table 2 in our Appendix.

111A standard (static) divide-and-conquer112strategy can be parametrized by the tuple113 $(R_{ed}, R_{ef}, S^m, A^m, R_{ef}^m)$ , where  $S^m, A^m, R_{ref}^m$ 114represents the multiplicity (i.e. number of attempts)115of  $\phi_{sol}$ ,  $\phi_{agg}$  and  $\phi_{ref}$ , respectively. First, the116 $\phi_{dec}$  transformation decomposes the starting

problem into B subproblems, which are then solved individually using  $\phi_{sol}$  then aggregated using  $\phi_{agg}$ .  $R_{ed}, R_{ef} \in \{0, 1\}$  indicate whether the  $\phi_{red}$  and  $\phi_{ref}$  transformations are applied after aggregation. If  $R_{ed} = 1$ , a single aggregation attempt is kept, while others are removed from the graph. If  $R_{ef} = 1$ , the remaining aggregation attempts are then refined with  $\phi_{ref}$ , and the highest-scoring attempt is kept as the final solution. See Appendix B for further details on the divide-and-conquer schedule. 117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

150

# **3** Thought Graph Exploration as a Markov Decision Process

Beyond static schedules, the transformation of a thought graph can be generalized as a Markov decision process  $(S, A, P_a)$  where the state  $s_t \in S$ represents an arrangement of nodes and edges in the thought graph, the action  $a \in \mathcal{A}$  indicates a transformation and target node subset (i.e.  $\mathcal{A} = \{(\mathcal{V}_s, \phi) \mid \mathcal{V}_s \subset \mathcal{V}, \phi \in \Omega\}, \text{ where } \Omega \text{ is the }$ set of transformations) and the transition probability  $\mathcal{P}_a(s, s')$  represents the probability that an action a applied at state s yields the expected new state s'. The optimal transformation sequence  $\Phi$  is then defined as the sequence of actions that maximize the conditional probability of reaching a solution state  $s^+$ , i.e.  $\Phi = (\phi_0, \ldots, \phi_n)$  that solves the following optimization problem. We bound the number of queries ( $|\Phi|$ ) by the constant  $\epsilon$ , as in the limit  $|\Phi| \to \infty$ ,  $P(s^+|s^0, \Phi) \to 1$ .

$$\underset{\Phi}{\arg\min} P(s^+ \mid s^0, \Phi) \quad \text{s.t.} \quad |\Phi| < \epsilon \quad (2)$$

In this work, we hypothesize that LLMs can approximate a solution to the stated optimization problem by acting as policy agents. We develop an



Figure 2: Multi-agent framework for reasoning over thought graphs. First, (1) the policy agent samples an action and subset of nodes given a prompt including (i-ii) general instructions and (iii-iv) an overview of the exploration state. The sample is then (2) passed to the reasoning agent, which finally (3) updates the thought graph state.

interactive framework consisting of a policy agent and a reasoning agent, as shown in Figure 2. In each iteration, (1) the policy agent selects an action from the action space, (i.e. the transformations in Table 2). The policy agent then (2) directs the reasoning agent to perform the selected action. Finally, (3) the reasoning agent updates the thought graph. The process is repeated until a solution is found or a maximum number of iterations is reached.

151

152

153

154

155

157

158

159

160

161

162

163

165

166

169

170

171

172

173

174

175

176

177

178

179

181

182

184

185

188

The policy agent is invoked using the prompt template shown in Figure 2. (i) The system prompt outlines the problem setting, input format and expected behaviour from the policy agent. (ii) A taskspecific list of actions, describing the preconditions and effects of each transformation, provides a semantic understanding of the action space. (iii) The current state of the graph is provided in a textual format, enumerating all nodes and edges. Finally, (iv) the action history in the current trial is included, promoting continuity in the strategies outlined in previous steps.

**In-Context Action Selection**: Prior work has shown that reasoning abilities of LLMs are enhanced when prompted to output a verbose sequence of steps before the solution (Wei et al., 2023; Wang et al., 2023). This mechanism can be seen as enabling in-context task learning from some extracted innate world knowledge. Hence, our policy agent is instructed to generate a detailed analysis on the state of the thought graph and exploration history before sampling the action space.

**Policy Agent Ensembles**: Given the stochastic nature of LLM next-token prediction, we observe high variability in the chosen action over several invocations of a policy agent from the same thought graph state. To enhance robustness, we democratize action selection over an ensemble of agents, meaning a parametrizable number of LLM queries are performed concurrently at every iteration. The selected action is takes as the most frequent proposal among the ensemble. See Appendix G for ablation studies on the impact of policy agent ensemble size on reasoning performance. 189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

## 4 Experiments

In this section, we outline the benchmarks, baselines and metrics used to evaluate ARIES against state-of-the-art methods for reasoning with LLMs. See also Appendix F for an evaluation of failure modes of our methodology.

**Experimental Setup:** We evaluate Llama-3.1-70B and Llama-3.1-405B as policy and reasoning agents, hosted with SGLang at a temperature of 1. We set the policy agent ensemble size to 5 in all experiments, as explained in Section G. Llama-3.1-70B was hosted with  $8 \times$  A6000 GPUs. Llama-3.1-405B was hosted using  $16 \times$  H100 GPUs distributed over 4 nodes. The total cost was approximately 3k GPU hours.

**Benchmarks:** We run our main evaluation on HumanEval (Chen et al., 2021). Additionally, we consider the set intersection problem at various levels of difficulty quantified by the size of the sets, leading to three benchmarks: setintersection32/64/128. Despite its simplicity, this has been shown to be a challenging benchmark for LLMs with direct prompting (Besta et al., 2024a).

**Baselines:** We use static transformation schedules as the baseline, following (Besta et al., 2024a). For each individual task, we carefully tune the hyperparameters using Bayesian optimization. We compare against baselines with several search compute budgets by considering three variants:  $GoT_{25\%}$ ,  $GoT_{50\%}$  and  $GoT_{100\%}$ , where the percentage corresponds to the number of trials spent until the hyperparameter search converges. See Ap-



Figure 3: Pareto frontiers in total query cost  $(C_{s+i})$  and task error  $(\mathcal{E})$  for set intersection tasks at various difficulty levels. The total cost is the number of queries expended at search and inference time. Llama-3.1-405B was used for the reasoning and policy agents. Our results (ARIES) have pushed the Pareto frontiers forward in each task.

Table 1: Task accuracy ( $\uparrow$ ), search and inference costs ( $C_s$  and  $C_i$ ,  $\downarrow$ ) on the HumanEval task.  $C_s$  and  $C_i$  are measured in number of queries. We use LLaMA-3.1-405B and GPT-40 as the underlying reasoning and policy agents. Direct IO means direct prompting, GoT follows Besta et al. (2024a), and ReACT follows Yao et al. (2023b).

Method		LLaMA-3.1-	405B	GPT-40			
	Acc (%)	Search Cost $(C_s)$	Inference Cost $(C_i)$	Acc (%)	Search Cost $(C_s)$	Inference Cost $(C_i)$	
Direct IO	77.4	0	1	84.1	0	1	
$GoT_{25\%}$	66.3	1160	34.8	75.6	1160	19.8	
$GoT_{50\%}$	67.5	2368	24.3	73.8	2368	18.5	
$GoT_{100\%}$	60.1	4742	8.17	69.5	4742	6.8	
ReACT	79.9	0	5.6	60.4	0	16.4	
ARIES	89.0	0	31.6	95.1	0	35.4	

pendix D for details on the full search methodology.

We also consider an IO (Input-Output) baseline, i.e. direct LLM prompting, and *smolagents* (Roucher et al., 2025), which is an implementation of the ReACT algorithm (Yao et al., 2023b). In the latter, the reasoning agent iteratively generates a candidate solution then reasons over execution feedback to refine subsequent candidates.

**Reported metrics:** For HumanEval, we report the task accuracy, while for list sorting and set intersection we report error function value  $\mathcal{E}$ . Details on the definition for the error function for each task can be found in Appendix C. Additionally, we report both the search cost  $C_s$  and inference cost  $C_i$ , measured as the number of LLM queries.

## 4.1 Results

226

227

231

235

239

240

241

243

245

246

247

249

**HumanEval:** Our key findings are shown in Table 1. It can be seen that by formulating this code generation task as a Markov decision process with an off-the-shelf LLM policy agent, we achieve up to 28.9% higher accuracy than the most query-efficient static schedule baseline with Llama-405b. Across both models, the inference cost is comparable to GoT<sub>25%</sub>, although the latter is obtained after

1160 LLM queries while ARIES avoids any search time requirement. Relative to the ReACT baseline, we achieve 9.1% higher accuracy with Llama-405b. However, this baseline does not generalize well across models, leading to a 9.1% accuracy degradation relative to GoT<sub>100%</sub> with GPT-40. 250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

Set Intersection: In Figure 3, we plot a Pareto curve showing viable trade-off points in task error and query cost for the set intersection task. Our approach extends the existing Pareto frontier constructed by considering static schedule baselines and direct prompting. In the set-intersection32 task, we achieve a  $2.3 \times$  error reduction relative to GoT<sub>25</sub> while also achieving 116× lower overall cost.

# 5 Conclusion

We introduce ARIES, a multi-agent architecture for topological reasoning. By viewing thought graph transformations as actions in a Markov decision process, we show off-the-shelf LLMs can drive efficient action policies without task-specific tuning, leading to state-of-the-art accuracy at reduced inference costs.

### 6 Limitations

272

273

274

276

277

281

289

297

303

310

311

313

314

315

317

318

319

321

#### 6.1 Assumptions and Robustness

The ARIES framework introduces a novel approach to reasoning with large language models (LLMs) through interactive thought graph environments. However, several strong assumptions underlie our methodology. Firstly, we assume that thought graph transformations can be effectively modeled as a Markov decision process (MDP) with welldefined state transitions. While this formulation enables structured reasoning, it may not fully capture the complexities of more ambiguous or highly interconnected problems. Additionally, our approach assumes that off-the-shelf LLMs can act as reliable policy agents without additional fine-tuning. This assumption holds for certain problem domains but may degrade in tasks requiring domain-specific knowledge or long-horizon planning.

Our empirical evaluation is constrained to specific reasoning tasks, including HumanEval and set intersection. While these benchmarks serve as valuable test cases for structured reasoning, they do not necessarily generalize to all problem types, particularly those with weakly defined intermediate states or multi-modal reasoning requirements. Furthermore, our evaluation primarily focuses on LLaMA-3.1 and GPT-40 models, and results may not be directly transferable to other architectures.

### 6.2 Potential Risks

The ARIES framework introduces both opportunities and challenges in autonomous reasoning. One primary risk is the potential for incorrect or biased reasoning paths due to the stochastic nature of LLM-generated decisions. Although our policy agent ensembles mitigate some of this variability, they do not fully eliminate erroneous transformations, particularly in deeper decomposition settings. The framework's reliance on existing LLMs also means that any biases present in the underlying models could propagate into the reasoning process, potentially leading to unfair or misleading outcomes.

Another concern is the environmental impact associated with inference-heavy approaches. While ARIES improves query efficiency relative to static transformation schedules, it still necessitates a significant number of LLM queries to achieve high accuracy. As LLMs scale, the energy consumption required for these inference tasks could become a sustainability concern, particularly in highthroughput applications.

#### 6.3 Failure Modes

Our empirical findings regarding the failure modes 324 of our methodology (Appendix F) highlight two 325 major failure modes: (1) inadequate LLM param-326 eter sizes and (2) increasing decomposition depth. 327 Smaller models (e.g., LLaMA-3.1-70B) struggle to 328 act as policy agents effectively, demonstrating sub-329 par reasoning capabilities compared to larger coun-330 terparts. This suggests that autonomous policy-331 driven thought graph exploration may require mod-332 els beyond a certain scale threshold to function 333 reliably. Additionally, as the depth of problem decomposition increases, ARIES exhibits a decline in 335 performance, primarily due to errors in aggregating 336 intermediate solutions. This limitation indicates 337 that current LLMs may have difficulties managing extended reasoning chains, which presents a barrier 339 to scalability.

322

323

### 341 References

342

343 344

345

346

347

357

361

373

375

376

377

379

381

396

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024a. Graph of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690.

Maciej Besta, Florim Memedi, Zhenyu Zhang, Robert Gerstenberger, Guangyuan Piao, Nils Blach, Piotr Nyczyk, Marcin Copik, Grzegorz Kwaśniewski, Jürgen Müller, Lukas Gianinazzi, Ales Kubicek, Hubert Niewiadomski, Aidan O'Mahony, Onur Mutlu, and Torsten Hoefler. 2024b. Demystifying chains, trees, and graphs of thoughts. *Preprint*, arXiv:2401.14295.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. Evaluating large language models trained on code. *Preprint*, arXiv:2107.03374.

- Aymeric Roucher, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunismäki. 2025. 'smolagents': a smol library to build great agentic systems. https://github.com/ huggingface/smolagents.
  - Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *Preprint*, arXiv:2408.03314.
  - 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.
  - Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent abilities of large language models. *Preprint*, arXiv:2206.07682.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models. *Preprint*, arXiv:2305.10601.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b.

React: Synergizing reasoning and acting in language models. *Preprint*, arXiv:2210.03629.

397 398

Table 2: Thought graph transformations. Each transformation is defined as  $\phi(G_{\tau}, m, S) = (V \cup V^+ \setminus V^-, E \cup V^-)$  $E^+ \setminus E^-$ ), where  $G_\tau = (V, E)$  is a thought graph,  $S \subset V$  is a subset of nodes, m is the multiplicity (number of attempts), and  $\mathcal{E}$ ,  $\mathcal{R}$ ,  $\mathcal{A}$  represent arbitrary functions for node expansion, refinement and aggregation, respectively. The sets  $V^+, V^-, E^+, E^-$  are defined as follows.

Transformation	Symbol	$\mathbf{V}^+$	$\mathbf{V}^{-}$	$\mathbf{E}^+$	$\mathbf{E}^{-}$
Decompose	$\phi_{dec}$	$\{\mathcal{E}(v) v\in S\}$	Ø	$\{(u,v) u\in S, v\in V^+\}$	Ø
Solve	$\phi_{sol}$	$\{\mathcal{S}(v) v\in S\}$	Ø	$\{(u,v) u\in S, v\in V^+\}$	Ø
Refine	$\phi_{ref}$	$\{\mathcal{R}(t) t\in S\}$	Ø	$\{(u,v) u\in S, v\in V^+\}$	Ø
Reduce	$\phi_{red}$	Ø	S	Ø	$\{(u,v) u\in S \lor v\in S\}$
Aggregate	$\phi_{agg}$	$\mathcal{A}(S)$	Ø	$\{(u,v) u\in S, v\in V^+\}$	Ø

#### **Thought Graph Transformations** Α

The full set of considered transformations is shown in Table 2.

#### B **Static Thought Graph Transformation** Schedule

Algorithm 1 represents a standard divide-andconquer strategy. The function  $\Delta(G^a_{\tau}, G^b_{\tau})$  outputs all nodes present in the first graph  $G_{\tau}^{a} = (\mathcal{V}_{a}, \mathcal{E}_{a})$ but not in the second  $G^b_{\tau} = (\mathcal{V}_b, \mathcal{E}_b)$ , defined formally as follows.

$$\Delta(G^a_\tau, G^b_\tau) = \{ v | v \in \mathcal{V}_a \& v \notin \mathcal{V}_b \}$$
(3)

Algorithm 1 Static Thought Graph Transformation Schedule

**Require:** Starting graph  $G_{\tau}^0$ , allow reduce  $R_{ed}$ , allow refine  $R_{ef}$ 

**Require:** Solve multiplicity  $S^m$ , aggregate multiplicity  $A^m$ , and refine multiplicity  $R^m_{\alpha f}$ 

$$\begin{split} & G_{\tau}^{dec} \leftarrow \phi_{dec}(G_{\tau}^{0}, 1, \{0\})) \\ & G_{\tau}^{sol} \leftarrow \phi_{sol}(G_{\tau}^{dec}, S^{m}, \Delta(G_{\tau}^{dec}, G_{\tau}^{0})) \\ & G_{\tau}^{agg} \leftarrow \phi_{agg}(G_{\tau}^{sol}, A^{m}, \Delta(G_{\tau}^{sol}, G_{\tau}^{dec})) \\ & \text{if } R_{ed} \text{ then} \\ & G_{\tau}^{red} \leftarrow \phi_{red}(G_{\tau}^{agg}, 1, \Delta(G_{\tau}^{agg}, G_{\tau}^{sol})) \\ & \text{else} \\ & G_{\tau}^{red} \leftarrow G_{\tau}^{agg} \\ & \text{end if} \\ & \text{if } R_{ef} \text{ then} \\ & G_{\tau}^{ref} \leftarrow \phi_{ref}(G_{\tau}^{red}, R_{ef}^{m}, \Delta(G_{\tau}^{red}, G_{\tau}^{agg})) \\ & G_{\tau}^{*} \leftarrow \phi_{red}(G_{\tau}^{ref}, 1, \Delta(G_{\tau}^{ref}, G_{\tau}^{red})) \\ & \text{else} \\ & G_{\tau}^{*} \leftarrow G_{\tau}^{red} \\ & \text{end if} \\ & \text{Return: } G_{\tau}^{*} \end{split}$$

#### С **Benchmarks**

We consider two popular tasks for topological reasoning with LLMs, which are amenable to a divideand-conquer strategy (i.e. decomposition, solving subproblems and merging): list sorting and set intersection. Despite their simplicity, prior works have shown that these tasks are extremely challenging for LLMs with direct prompting (Besta et al., 2024a).

As mentioned in Section 4, the core results are reported using HumanEval and set-intersection. The list-sorting task was used for illustration purposes in estimating transition probabilities, performing ablation studies and evaluating failure modes.

**Sorting**: involves sorting a list of numbers between 0 and 9 in ascending order. The error function  $\mathcal{E} = X + Y$  has its subterms defined in Equation 4, where a is the input list and b is a candidate solution. X corresponds to the number of incorrectly sorted pairs, while Y corresponds to the frequency difference between a and b for each digit.

$$X = \sum_{i=1}^{m-1} \operatorname{sign}(\max(b_i - b_{i+1}, 0))$$

$$Y = \sum_{i=0}^{9} ||\{b_p : b_p = i\}| - |\{a_q : a_q = i\}||$$
(4)

Set Intersection: involves finding the intersection of sets A and B. The error function is defined in Equation 5, where C is the candidate solution. The first and second terms correspond to missing and extra elements, respectively.

$$\mathcal{E} = |(A \cap B) \setminus C| + |C \setminus (A \cap B)|$$
 (5)

#### D Static Schedule Parameter Search

As described in Section 2, a static transformation can be characterized using a set of discrete param-

399

400

401

402

403

404

405

406

407

408

432

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

- 433
- 434 435

436 437

438

439

440

441

442 443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473 474

475

476

477

478

eters. We ran bayesian search using using Treestructured Parzen Estimator (TPE) sampling to determine each parameter, establishing strong baselines for each task.

Table 3: Search space for each parameter characterizing a static transformation.

		Search
	Parameter	Space
$R_{ed}$	Allow reduction	$\{0,1\}$
$R_{ef}$	Allow refinement	$\{0,1\}$
$S^{m}$	Solve multiplicity	$\{1, 5, 10, 15, 20\}$
$A^m$	Aggregate multiplicity	$\{1, 5, 10, 15, 20\}$
$R^m_{ef}$	Refine multiplicity	$\{1, 5, 10, 15, 20\}$

The search space is shown in Table 3. We run multi-objective search to concurrently minimize the task-specific error function  $\mathcal{E}$  (Section C) and associated cost, measured as  $|\Phi(\omega)|$  where  $\Phi(\omega) =$  $(\phi_0, \ldots, \phi_m)$  is a tuple enumerating thought graph transformations, as a function of the schedule parameters  $\omega \in \Omega$ , where  $\Omega$  is the search space. Note that  $|\Phi(\omega)|$  correlates with the number of LLM queries, meaning this formulation aims to minimize exploration cost.

In selecting parameter configurations, we use the cost function in Equation 6, such that the objectives of cost and error minimization are balanced through the scalar constant  $\alpha \in (0, 1)$ . We aim to assign equal importance to the cost and error objectives by tuning  $\alpha$  independently for each task such that the mean value of the first term matches the second term, i.e.  $\alpha E[\mathcal{E}] = (1 - \alpha)E[|\Phi(\omega)|)]$ , or equivalently  $\alpha = \frac{E[|\Phi(\omega)|]}{E[\mathcal{E}+|\Phi(\omega)|]}$  where *E* denotes the expected value. The expectations are obtained with random sampling.

$$\min_{\omega} \left[ \alpha \mathcal{E} + (1 - \alpha) |\Phi(\omega)| \right] \tag{6}$$

Search was conducted separately on Llama-3.1-70B and Llama-3.1-405B. For sorting and set intersection tasks, search is conducted separately for each difficulty level, ensuring the chosen parameters are adapted to the task. Note that we present three search checkpoints  $GoT_n$  for  $n \in$ {25, 50, 100}, where *n* corresponds to the percentage of trials until convergence. We define the convergeance point as the first iteration where a rolling window *J* of size 20 matches the condition  $J^k = J^{k-1}$ . This enables comparing our proposed

Table 4: Results from GoT static schedule parametersearch on Llama-3.1-405B.

Task	Alpha ( $\alpha$ )	GoT <sub>25</sub>	GoT <sub>50</sub>	GoT <sub>100</sub>
sorting32	0.99	0.38	0.38	0.37
sorting64	0.96	4.85	4.49	3.84
sorting128	0.84	28.76	25.76	24.36
set32	0.99	0.16	0.16	0.12
set64	0.99	0.71	0.51	0.31
set128	0.98	3.51	3.51	2.99

LLM-guided approach to optimized search schedules at various search budgets.

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

513

514

515

516

The complete search results for Llama-3.1-405B are shown in Table 4. It can be seen that tasks with higher decomposition depth incur lower values of  $\alpha$  due to the higher magnitude of the error function. sorting64, sorting128 and set-intersection64 show a smooth decline in the cost function, while the remaining tasks remain at local minima until close to the end of the search. The non-convexity of the search space highlights the cost associated to optimize the parameter set associated with static transformations.

# **E** Transition Probability Profiling

In this section, we estimate the transition probabilities for each thought graph transformation across a number of tasks to gain insight into factors impacting a thought graph formulation of each reasoning problem. For  $\phi_{ref}$ , we define a successful transition when  $\mathcal{E} = 0$  for the resulting node, considering only cases when the transformation is executed on nodes previously containing errors. In transformations requiring LLM calls, the transition probability between two states is a random process governed by the token distribution parametrized by the LLM. When LLM calls are not required, i.e. the transformation is implemented through simple node manipulation, the transition probability is 1.

The results are summarized in Table 7. We observe the refinement transformation has notably low success probability, particularly in coding and sorting tasks. Additionally, sorting is the only task with non-deterministic aggregation, which is a potential error source. We note that the performance of a thought graph formulation depends on the ability of the policy agent to capture the success profile of various transformations for a task, and adapt the exploration strategy accordingly.

Table 5: Failure mode 1 results. Mean value of the error  $\mathcal{E}(\downarrow)$  for benchmarks with low decomposition depth. Llama-3.1-70B was used for the reasoning and policy agents.

Method	Direct Prompting	$GoT_{25\%}$	${ m GoT}_{50\%}$	$GoT_{100\%}$	ARIES	
sorting32	2.2	0.82	0.95	0.73	1.29	
set-intersection32	1.05	0.41	0.0	0.37	1.22	

Table 6: Failure mode 2 results. Mean value of the error  $\mathcal{E}(\downarrow)$  and search cost *C* in terms of number of queries ( $\downarrow$ ). Both the reasoning and policy agents are LLaMA-405B.

Method	Direct F	rompting	GoT	25%	GoT	50%	GoT	100%	ARI	ES
Metrics	E	C	${\mathcal E}$	C	E	C	${\mathcal E}$	C	E	C
sorting32	0.6	1	0.74	825	0.82	1650	0.28	3300	0.22	20
sorting64	5.07	1	2.22	1671	2.74	3343	3.46	6687	9.15	48
sorting128	12.75	1	13.96	2444	12.65	4888	18.65	9776	32.74	48

Table 7: Esimated transition probabilities for each thought graph transformation, taken as the number of successful state transitions in a static schedule.

	$\phi_{\mathbf{sol}}$	$\phi_{\mathbf{ref}}$	$\phi_{\mathbf{red}}$	$\phi_{\mathbf{agg}}$
HumanEval	0.77	0.29	1	1
sorting32	0.57	0.12	1	0.60
set-intersection32	0.75	0.71	1	1

## F Failure Modes

517

518

519

521

523

537

538

540

In this section, we perform a number of empirical studies aiming to understand the main limiting factors impacting the performance of LLM policy agents on interactive thought graphs. We find there are two major failure modes, described as follows. *Failure mode 1: LLM Parameter Count* 

We find that LLMs with insufficiently large parame-524 ter sizes exhibit limited performance when utilized 525 as policy agents on thought graph environments. 526 We deploy Llama-3.1-70B as policy and reasoning agents in sorting and set intersection tasks, against 528 which the larger LLM (Llama-405B) was shown to perform well as a policy agent. As shown in Ta-530 ble 5, LLM-guided graph exploration (ARIES) did not outperform static schedule baselines in this scenario. These findings are consistent with (Wei et al., 2022), which demonstrated that zero-shot chainof-thought reasoning abilities emerges in models 535 beyond 175B parameters.

✤ Failure mode 2: Decomposition Depth

We examine the impact of decomposition depth by analyzing the results in the sorting task, shown in Table 6. We observe LLM policy agents lead to a 21% performance improvement relative to the most optimized static baseline in sorting32, which has a decomposition depth of 2. However, as discussed in Appendix E, the sorting task presents a particular challenge due to the lower success probability of the aggregation transformation. As the complexity and decomposition depth of a task increases, the policy agent is required to apply a higher number of aggregation transformations. Therefore, we observe up to  $4.12 \times$  and  $2.6 \times$  performance deterioration in sorting64 and sorting128, respectively. Through empirical analysis, we observe that in the latter tasks, the  $\phi_{agg}$  transformation constitutes 86% and 68% of all policy agent errors, respectively. As such, we conclude that high decomposition depths present a significant failure mode for LLM-guided thought graph exploration, particularly in tasks with low success transition probabilities for the aggregation transformation.

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

563

565

566

567

568

569

570

571

572

573

## **G** Ablation Studies

As discussed in Section 3, two factors that impact the performance of LLMs as policy agents in interactive thought graph environments are the size of the ensemble and the use of chain of thought reasoning to enhance the planning abilities of the policy agent. In this section, we aim to understand the impact of each factor by evaluating sorting tasks over a range of ensemble sizes from 1 to 15, with and without CoT prompting in the policy agent.

As shown in Figure 4, as the ensemble size increases to 5, CoT prompting leads to large performance improvements, though the benefits start diminishing beyond this point. Without CoT prompt-



Figure 4: Mean error (y-axis) obtained in the sorting32 task over a sweep of ensemble sizes (x-axis). Llama-3.1-70B was used as the policy agent.

574	ing, the trend is less consistent, and larger ensemble
575	sizes sometimes yield worse performance. Addi-
576	tionally, errors without CoT are higher for both
577	tasks at any ensemble size. This highlights the ne-
578	cessity of CoT prompting in enhancing the LLM
579	policy agent's ability to adapt from feedback and
580	drive thought graph transformations.