EFFICIENT LLM COLLABORATION VIA PLANNING

Anonymous authors

000

001 002 003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

025

031

032

033

034

036

038

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Recently, large language models (LLMs) have demonstrated strong performance, ranging from simple to complex tasks. However, while large proprietary models (e.g., models with over 100B parameters) achieve remarkable results across diverse tasks, they are often accessible through costly APIs, making frequent use too costly for many applications. In contrast, small open-source models (e.g., models with fewer than 3B parameters) are freely available and easy to deploy locally, but their performance on complex tasks remains limited. This trade-off raises a natural question: how can small and large models efficiently collaborate to combine their complementary strengths? To bridge this trade-off, we propose COPE, a test-time collaboration framework. A planner model first generates a plan, a high-level abstraction of the task, and this plan serves as a lightweight intermediate that guides a downstream executor model. Small and large models take turns acting as planner and executor, exchanging plans in a multi-stage cascade to collaboratively solve tasks. Through comprehensive experiments on benchmarks spanning mathematical reasoning, code generation, open-ended tasks, and agent tasks, we demonstrate that COPE achieves performance comparable to large proprietary models, while drastically reducing the inference API cost. These results highlight planning as an effective prior for cost-efficient inference.

1 Introduction

Large language models (LLMs) have achieved impressive performance on complex tasks (Rae et al., 2021; Lewkowycz et al., 2022). However, deploying them at scale presents a critical trade-off between model capability and computational cost. As larger models continue to improve in capability, they also grow increasingly costly to run (Achiam et al., 2023; Schick et al., 2023). Smaller models, in contrast, are far more deployable, but their limited capacity hinders them on complex tasks (Kojima et al., 2022; Wei et al., 2022).

This trade-off is especially clear in realistic deployment scenarios. Users typically interact with language models via mobile or edge devices, which cannot host large models (Leviathan et al., 2023; Chowdhery et al., 2023). Instead, large models often reside in the cloud and incur substantial inference costs. In this setting, relying solely on large models for every request is both economically and practically infeasible. These constraints call for inference strategies that use small models whenever possible and selectively leverage large models only when necessary. For clarity, we use the terms free vs. costly, small vs. large, and weak vs. strong interchangeably to refer to models with different capacities and cost characteristics, depending on context.

There are few methods that address cost-aware small-large model collaboration, where large models are used sparingly at test time (Yue et al., 2023; Kolawole et al., 2024; Chen et al., 2023). These methods essentially operate as cascades: a small model first attempts the task, and if its output fails to satisfy a confidence-based criterion (e.g., majority voting threshold), the task escalates to the large model in the next stage. However, they remain a form of independent delegation, as only one model operates in each stage, without collaborating within a stage to produce outputs together. This limits their ability to jointly perform complex tasks in a structured and interactive manner.

In contrast, we propose to structure this collaboration through *planning*. Instead of directly performing a task, a model first produces a high-level abstraction of the task (such as a goal or a guideline), which we call a plan. A plan guides downstream execution by another model. This planning step

serves as a lightweight, transferable intermediate that allows models to scaffold each other's thinking. Thus, we explore the following research question:

Research Question

Can planning function as a mediator of collaboration between large and small models for efficient inference?

To answer this question, we introduce COPE (Collaborative Planning and Execution), a test-time inference framework in which free and costly models share the roles of planner and executor. The process begins with a free model attempting both planning and execution. If the model lacks confidence in its output, planning and then execution are gradually escalated to a costly model. This structure allows models to collaborate adaptively: easy tasks are handled by free models, and harder tasks are handled by costly models.

COPE significantly reduces inference cost while matching or even surpassing the performance of the larger model. For example, on the MATH-500 dataset, COPE achieves 75.8% accuracy (higher than GPT-4o's 75.2%) while reducing cost by nearly 45%. On the MBPP code generation benchmark, COPE improves accuracy to 66.4% compared to GPT-4o's 64.0%, while cutting inference cost by nearly 75%. Similar trends hold across diverse tasks including open-ended tasks and agent tasks. These results highlight COPE as a broadly effective and scalable framework for efficient LLM inference.

2 RELATED WORK

Planning for LLM inference. Planning has been primarily studied in robotics and vision-language-action models as a core mechanism for structured decision-making (Kim et al., 2024; Black et al., 2023). A high-level planning module decomposes a task into subgoals, which are then carried out through low-level execution by downstream components. Planning can be seen as a feature of System-2 thinking (Ji et al., 2025).

In contrast, planning remains relatively underexplored in the context of LLM inference. Wang et al. propose the use of planning tokens prepended to each chain-of-thought step. However, this method requires supervised training to learn the tokens and focuses on local step-wise annotations rather than providing global guidance. Saha et al. (2025) introduce EvalPlanner, a framework for LLM-as-a-Judge (Zheng et al., 2023) evaluation. Given a pair of responses, the model generates an evaluation plan to structure its judgment. However, the plan is generated only after seeing the responses and serves solely to structure the judgment process, not to guide generation. The method further requires training on both planning and reasoning.

In our method, planning precedes execution in LLM inference. The plan is a purpose-driven representation of the task (such as a goal or guideline). Here, planning is not learned and decoupled from execution, so it serves as a lightweight interface between models of different capacities.

Test-time compute scaling. Test-time compute scaling (Snell et al., 2024) refers to how computational resources are allocated and utilized at inference time to improve accuracy. These efforts typically fall into two directions: parallel sampling, which generates multiple outputs and selects the best among them (Wei et al., 2022; Wang et al., 2022), and sequential sampling (or iterative refinement), where a model revises its outputs over multiple passes (Madaan et al., 2023; Zelikman et al., 2022; Lee et al., 2025; Kumar et al., 2024; Muennighoff et al., 2025).

These techniques typically aim to improve performance by allocating more compute to a single model. However, they remain limited by the capacity of the model itself. Once a model's abilities are exceeded, performance saturates (Snell et al., 2024). This motivates a different approach: rather than pushing a single model further, we explore collaboration across models of different capacities. This form of cross-model scaling can be seen as a new axis of test-time compute scaling.

Cost-aware LLM inference. Some recent approaches aim to reduce the cost of LLM inference by using costly models only when necessary. Given a problem, LLM Cascade (Yue et al., 2023) samples multiple answers from a free model and estimates consensus ratio, which is the degree

Table 1: Accuracy (%) of model combinations on the MATH-500 dataset. 'Base' shows accuracy without a planner; other columns use GPT-mini, Llama-3B, or Llama-1B as planners.

Table 2: Accuracy of Llama-1B with different plan types on the MATH-500 dataset.

Executor	Base	GPT-mini	Llama-3B	Llama-1B
GPT-mini	73.8	76.2	70.6	69.6
Llama-3B	42.8	53.0	37.6	32.8
Llama-1B	25.2	36.4	26.0	23.2

Plan Type	Accuracy (%)
None	25.2
Guideline	23.2
Goal	30.2

of agreement among the answers. If the consensus ratio fails to meet a criterion, the problem is escalated to a costly model. It leverages few-shot examples with a mixture of reasoning formats (Chain of Thought (CoT) and Program of Thought (PoT)). ABC (Kolawole et al., 2024) routes among multiple models without relying on curated few-shot examples. FrugalGPT (Chen et al., 2023) trains a lightweight scoring function for a routing policy. However, the limited capacity of the scoring function may hinder performance on complex tasks. In contrast, our method requires neither few-shot examples nor training, and it extends beyond reasoning to open-ended and agent tasks.

3 MOTIVATION

In this section, we hypothesize that generating a simple plan can guide downstream execution (e.g., reasoning), and introduce motivational examples and the insights drawn from them. In this work, a *plan* refers to a high-level abstraction of the task. More precisely, it is a purpose-oriented representation that serves as a guide for execution (refer to Figure 2).

Given a task, we first prompt the planner model to generate a brief guideline, typically one or two sentences, that outlines how to approach the task. This guideline is then provided, along with the original task, to the executor model. The following observations illustrate how this planning stage enables collaboration between models.

We use Llama-3.2-1B-Instruct, Llama-3.2-3B-Instruct, and GPT-4o-mini, hereafter referred to as Llama-1B, Llama-3B, and GPT-mini, respectively. Llama-1B and Llama-3B are open-source and freely accessible, while GPT-mini is accessed via a paid API.

Observation 1. Larger planners help smaller executors.

We begin by examining whether larger models can improve the performance of smaller ones by providing plans. As shown in Table 1, when a larger model acts as a planner for a smaller executor, such as GPT-mini planning for Llama-3B or Llama-3B planning for Llama-1B, the executor's accuracy improves over the no-plan baseline. For example, Llama-3B achieves 42.8% accuracy without planning, which rises to 53.0% when guided by GPT-mini. Likewise, Llama-1B improves from 25.2% to 36.4% when using plans from Llama-3B. This suggests that larger models are effective at guiding smaller models through planning.

Observation 2. Smaller planners degrade larger executors.

We next examine whether smaller models can effectively serve as planners for larger executors. As shown in Table 1, when smaller models such as Llama-1B or Llama-3B act as planners for larger models like GPT-mini or Llama-3B, the resulting performance is often lower than the no-plan baseline. For instance, GPT-mini achieves 73.8% without a plan, but drops to 70.6% and 69.6% when using plans from Llama-3B and Llama-1B, respectively. This suggests that low-quality plans generated by smaller models can hinder the execution ability of larger models.

Observation 3. A model benefits from plans aligned with its capacity.

We investigate whether a model can effectively act as their own planners, i.e., given a task, the model first generates a plan and then performs the task based on the plan. For large models such as GPT-mini, self-generated plans improve performance: GPT-mini achieves 73.8% accuracy without planning, which increases to 76.2% when using its own plans (Table 1). This suggests that large models can successfully scaffold their own execution through planning. In contrast, small models

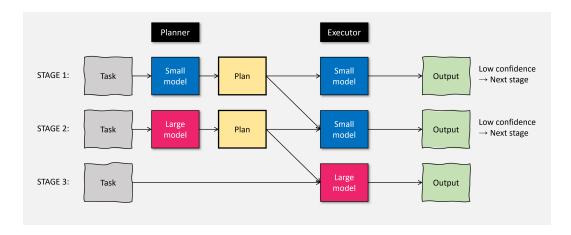


Figure 1: Overall framework of COPE. The system proceeds in up to three stages of inference, where small and large models alternate roles as planner and executor. In each stage, given a task, a plan is generated by the planner, and the executor produces candidate outputs. If a task-specific confidence falls below the stage threshold, the task escalates to the next stage. Plans generated in earlier stages are retained and reused in later stages.

show mixed results when acting as their own planners. For example, Llama-1B performs worse with self-generated plans (23.2%) than with no plan at all (25.2%).

To explore whether simpler planning might help, we prompt the model to generate a goal (what to achieve) rather than a guideline (how to solve the task) as before. In this setting, accuracy improved to 30.2% (Table 2). This suggests that planning must be tailored to the model's capacity: small models benefit more from simpler, purpose-oriented abstractions than from complex instructions.

These findings motivate a broader framework where planning functions as a modular component in test-time collaboration.

4 Method

We introduce COPE, a collaborative inference framework in which small and large models alternate roles as planners and executors across multiple stages. The core idea is to structure model collaboration by first generating an intermediate *plan*, which then guides downstream execution. Figure 1 illustrates the overall structure of COPE. Figure 2 shows an example of plan generation and execution based on that plan.

COPE operates in three stages. In Stage 1, a small model attempts to perform the task on its own by generating both a plan and an output. In Stage 2, a large model provides a new plan, and the small model reattempts execution with this guidance. In Stage 3, the large model assumes both roles to directly perform the task. Each stage is triggered only when the previous stage fails to reach sufficient confidence, enabling adaptive allocation of compute. Here, confidence is a general notion that can be instantiated in different ways depending on the task (e.g., majority voting for reasoning, test case pass rate for coding, perplexity for open-ended generation). Below we describe each stage in detail, focusing on majority voting as the confidence measure for reasoning tasks.

Stage 1: small model planning and execution In Stage 1, the small model acts as both planner and executor. Given a problem x, it first samples n goal-type plans $g_1, \ldots, g_n \sim P_{\text{small}}(g \mid x)$. For each plan g_i , it then generates a solution $y_i \sim P_{\text{small}}(y \mid x, g_i)$, and extracts the corresponding answer $a_i = f(y_i)$, where f denotes an answer extraction function.

The final answer is determined by majority vote over the answer set $\mathcal{A} := \{a_1, \dots, a_n\}$, treating agreement among solutions as an indicator of correctness (Wang et al., 2022; Du et al., 2023). The *consensus ratio*, defined as the fraction of samples agreeing on the most frequent answer, serves as the confidence score:

Vanilla

Problem: You had 7 bags of gold coins, each containing the same number of coins. You found another bag that contained 53 coins, and you decided to redistribute all your coins evenly across the 8 bags. You successfully did this, and you also know that the total number of coins you have is more than 200. What is the smallest number of coins you could have had before finding the extra bag? **Solution (executor):** Step 1: Determine the Target Number of Coins per Bag. Since you end up with eight bags and all bags have the same number of coins, let x be the number of coins in each of the original seven bags. After adding the 53 coins, the total number of coins becomes 8x. Thus, we have: 7x + 53 = 8x [...]

COPE

Problem: You have seven bags of gold coins. [...] Guideline (planner): You want to redistribute these coins into 8 bags, meaning 7x+53 must be divisible by 8. Then, set up the equation $7x+53 \geq 201$ to ensure you have more than 200 coins. Solve for x to find the smallest integer that satisfies both conditions. Solution (executor): Step 1: Set up the divisibility condition. The total coins 7x+53 must be divisible by 8. This can be expressed as: $7x+53\equiv 0 \pmod 8$. Simplifying modulo 8: $7x+53\equiv 7x+5\equiv 0 \pmod 8$, $7x\equiv -5 \pmod 8$. Since $-5\equiv 3 \pmod 8$: $7x\equiv 3 \pmod 8$. [...] Step 3: Ensure 7x+53>200. Substitute x=8k+5 into the inequality: [...] The smallest integer k that satisfies this inequality is k=3. Therefore, the total number of coins before finding the bag of 53 coins is: $7x=7\times 29=203$. Thus, the smallest number of coins you could have had before finding the bag of 53 coins is: 203

Figure 2: Comparison between vanilla inference and planning-guided inference with COPE. The vanilla model fails to account for the divisibility and inequality constraints, leading to incorrect reasoning (red). In contrast, COPE's planner highlights these constraints explicitly (green), allowing the executor to follow a structured solution path.

Consensus(
$$\mathcal{A}$$
) := $\frac{\max_{a' \in \mathcal{A}} |\{i : a_i = a'\}|}{n}$. (1)

If the consensus exceeds a threshold τ_1 , the answer is accepted. Otherwise, the problem is escalated to the next stage. In that case, one of the plans g_i that produced the majority-voted answer is randomly selected and designated as g^S , which is passed to Stage 2.

Stage 2: large model planning, small model execution If Stage 1 fails to produce a confident answer, Stage 2 begins with the large model generating a new guideline-type plan $g^L \sim P_{\text{large}}(g \mid x)$. This plan is passed to the small model, which then attempts to solve the problem again by leveraging both the original plan g^S from Stage 1 and the new plan g^L . These two plans are provided jointly as input to the small model in the prompt.

Then, the small model generates a solution $y_i \sim P_{\text{small}}(y \mid x, g^S, g^L)$ and corresponding answer $a_i = f(y_i)$. Majority voting is applied over the answer set $\{a_1, \ldots, a_n\}$ as in Stage 1. If the consensus ratio exceeds a second threshold τ_2 , the answer is accepted. Otherwise, the problem is escalated to Stage 3.

Stage 3: large model planning and execution If the small model still fails in Stage 2, the problem is escalated to Stage 3. Here, the large model takes full control of both planning and execution.

Table 3: Performance on the MATH-500 dataset for mathematical reasoning.

Small	Large	Method	Shot	Train	To	otal
Siliuii	6		Shot	Truin	Acc↑	Cost↓
	_	Small (single)	Zero	×	52.8	_
	- Small (majority)		Zero	×	65.2	_
		Large	Zero	×	75.2	5025
		Cascade (vote)	Few	×	72.4	3017
	GPT-40	Cascade (verify)	Few	×	73.0	3021
EVAQNE 2.5	GP 1-40	ABC	Zero	×	74.4	4297
EXAONE-3.5		FrugalGPT	Few	\checkmark	68.0	4163
-2.4B		COPE (ours)	Zero	×	75.8	2766
		Large	Zero	×	73.8	355
	GPT-40 -mini	Cascade (vote)	Few	×	70.6	224
		Cascade (verify)	Few	×	72.2	250
		ABC	Zero	×	72.4	304
		FrugalGPT	Few	\checkmark	65.6	298
		COPE (ours)	Zero	×	74.4	212
	_	Small (single)	Zero	×	49.0	_
	_	Small (majority)	Zero	×	58.0	_
		Large	Zero	×	75.2	5025
		Cascade (vote)	Few	×	71.4	5150
	GPT-40	Cascade (verify)	Few	×	72.6	5046
Ovven 2.5	GF 1-40	ABC	Zero	×	72.6	4723
Qwen-2.5 -1.5B		FrugalGPT	Few	\checkmark	68.2	5472
-1.58		COPE (ours)	Zero	×	75.2	3725
		Large	Zero	×	73.8	355
		Cascade (vote)	Few	×	71.6	327
	GPT-4o	Cascade (verify)	Few	×	72.8	319
	-mini	ABC	Zero	×	72.0	298
		FrugalGPT	Few	\checkmark	69.4	349
		COPE (ours)	Zero	×	74.0	285

Specifically, it reuses the plan g^{L} that it generated in Stage 2, and samples $y \sim P_{\text{large}}(y \mid x, g^{L})$. The final answer is extracted as a = f(y). Since this is the final stage, no further escalation occurs, and the answer is returned as the system's output.

EXPERIMENTS

We evaluate COPE on multiple tasks spanning mathematical reasoning, code generation, openended tasks, and agent tasks. We find $n=8, \tau_1=0.75, \text{ and } \tau_2=0.5$ are effective in many of experiments. Please refer to the empirical study on the effect of these hyperparameters in Appendix A.

Baselines. We compare COPE against several recent methods for cost-aware inference. We implement the method proposed in Yue et al. (2023), referred to as Cascade. Cascade includes two decision strategies: vote and verify. In the vote setting, multiple answers are sampled from two different few-shot prompts, and the final answer is selected by majority voting across all the answers. If the consensus ratio exceeds a threshold, the answer is accepted. Otherwise, the problem is escalated to a stronger LLM. In the verify setting, majority voting is applied separately to answers from each few-shot prompt, and the final decision is accepted only if the two majority-voted answers agree. Otherwise, the problem is escalated to the stronger LLM.

324 325

Table 4: Performance on the AIME-2024 dataset for mathematical reasoning.

334 335 336

337 338

339

340 341 342

343 344 345

347 348 349

346

350351352353

354

355

361362363

364

360

366367368369

370

371372373374375

376

377

Method Shot Acc ↑ Train Cost ↓ Small (single) Zero 3.33 X Small (majority) Zero 6.67 X Large Zero 33.3 502 X 23.3 568 Cascade (vote) Few X Cascade (verify) 26.6 570 Few X ABC Zero X 33.3 485 FrugalGPT Few 33.3 513 40.0 478 COPE (ours) Zero X

ABC (Kolawole et al., 2024) is another agreement-based method. Unlike Cascade, it does not rely on curated few-shot exemplars. Instead, it proceeds through multiple stages with progressively stronger models as the stages advance. We report results using three models, following the original setup.

FrugalGPT (Chen et al., 2023) trains a lightweight scoring function which is implemented with DistilBERT (Sanh et al., 2019). Specifically, the scoring function is trained to estimate whether an answer is correct for a given problem. A routing policy based on the score decides whether to accept the answer or escalate to a stronger model.

We also consider three simpler baselines. The first uses only the small model in a single-pass setting without any sampling or aggregation (*Small (single)*). The second baseline uses the small model to generate multiple outputs and applies majority voting to select the most frequent answer (*Small (majority)*). The third baseline runs the large model once without planning or verification (*Large*). For each task, we evaluate against the baselines that are applicable to that task.

5.1 MATHEMATICAL REASONING

Datasets and models. We use the MATH-500 dataset (Lightman et al., 2023) and the AIME-2024 dataset (Art of Problem Solving, 2024). MATH-500 is a filtered subset of the MATH dataset (Hendrycks et al., 2021) consisting of 500 math problems. The problems cover a broad range of topics and difficulty levels. Compared to MATH-500, AIME-2024 is a more challenging benchmark. AIME-2024 consists of 30 curated olympiad-level problems.

We use EXAONE-3.5-2.4B-Instruct (Research, 2024) and Qwen-2.5-1.5B-Instruct (Team, 2024) as the small models, and GPT-4o and GPT-4o-mini as the large models. For the more challenging AIME-2024 dataset, considering the increased difficulty of the problems, we use Qwen-2.5-3B-Instruct as the small model, and Gemini-2.0-Flash as the large model. For ABC, which requires three models, we additionally include Qwen-2.5-0.5B-Instruct alongside the above models.

Evaluation metrics. We evaluate accuracy (%) using exact match. For mathematical answers, we use symbolic comparison to accept equivalent expressions (e.g., 1/2 and 0.5). For non-mathematical outputs, exact string match is used. For cost, we use the actual per-problem API price based on real-world model pricing (OpenAI, 2025; Google DeepMind, 2024). We scale the values by one million for readability, reporting cost in USD per million problems.

Results. Table 3 shows that COPE outperforms the large model baseline in overall accuracy, with 75.8% vs. 75.2% for GPT-40 and 74.4% vs. 73.8% for GPT-40-mini. This suggests that collaborative planning enables small models to solve certain problems that large models alone cannot. COPE also reduces inference cost significantly, with 2766 compared to 5025 for GPT-40, and 212 compared to 355 for GPT-40-mini. This corresponds to nearly a 50% reduction in cost while maintaining or improving accuracy. Moreover, COPE consistently outperforms few-shot, multi-model, and training-based baselines in both accuracy and cost.

In Table 4, COPE achieves the highest accuracy while incurring the lowest cost, outperforming the baselines. This suggests that COPE excels even on significantly more difficult reasoning tasks.

Table 5: Performance on the MBPP dataset for code generation tasks.

Small	Large	Method	Shot	Train	To	otal	
Sillari	Luige	111041104			Acc↑	Cost ↓	
	_	Small (single)	Zero	×	42.2		
	_	Small (multi)	Zero	×	57.2	_	
		Large	Zero	×	64.0	4889	
		Cascade	Few	×	64.6	2706	
EXAONE-3.5	GPT-40	ABC	Zero	×	63.2	1267	
-2.4B		FrugalGPT	Few	\checkmark	63.6	4012	
-2.4D		COPE (ours)	Zero	×	66.4	1279	
		Large	Zero	×	57.6	276	
	GPT-40	Cascade	Few	×	59.2	159	
	-mini	ABC	Zero	×	59.4	42.2 - 57.2 - 64.0 4889 64.6 2706 63.2 1267 63.6 4012 66.4 1279 57.6 276 59.2 159	
	-1111111	FrugalGPT	Few	\checkmark	240		
		COPE (ours)	Zero	×	62.2	77	
	_	Small (single)	Zero	×	38.6	_	
	_	Small (multi)	Zero	×	54.2	_	
		Large	Zero	×	64.0	4889	
		Cascade	Few	×	62.2	2263	
Owen-2.5	GPT-4o	ABC	Zero	×	62.6	1542	
C		FrugalGPT	Few	\checkmark	59.4	5029	
-1.5B		COPE (ours)	Zero	×	66.2	1469	
		Large	Zero	×	57.6	276	
	GPT-40	Cascade	Few	×	57.4	147	
	-mini	ABC	Zero	×	59.2	3.4 240 22.2 77	
	1111111	FrugalGPT	Few	\checkmark	57.2	493	
		COPE (ours)	Zero	×	61.4	94	

Compared to Table 3, the performance gap between the Small baselines and COPE is larger on these harder problems, which is consistent with the difficulty-based analysis in Appendix A.2.

5.2 Code Generation

For code generation, we adopt a pass-based evaluation instead of majority voting, since multiple correct outputs may differ syntactically. Each generated output is executed against the task's test cases. If at least one output passes all test cases, a passing output is randomly selected as the final output. Otherwise, the task is escalated to the next stage.

Datasets and models. We use the 500-problem test set of the MBPP dataset (Austin et al., 2021), which consists of crowd-sourced Python code generation tasks. Each task includes a natural language task description and test cases. We use EXAONE-3.5-2.4B-Instruct and Qwen-2.5-1.5B-Instruct as the small models, and GPT-40 and GPT-40-mini as the large models.

Evaluation metrics. We evaluate performance using the separate hidden test cases from the MBPP-ET dataset (Dong et al., 2025). A generated program is considered correct if it passes all hidden test cases.

Results. Table 5 reports the results on code generation. We find that, due to its limited capacity, the scoring function of FrugalGPT is difficult to train effectively on complex tasks. COPE achieves both the highest accuracy and the lowest cost. This shows that the benefits of collaborative planning extend beyond mathematical reasoning.

Table 6: Performance on the MT-Bench dataset for open-ended tasks.

Method	Win rate (%) ↑	Cost↓
Large	35.0	225
COPE (ours)	36.9	160

Table 7: Performance on the ALFWorld benchmark for agent tasks.

Method	Success rate ↑	Cost↓
Large	42.5	11.06
COPE (ours)	44.8	8.64

5.3 OPEN-ENDED TASKS

To examine whether COPE extends beyond reasoning domains, we further evaluate it on open-ended tasks. In each stage, the small model generates multiple candidate outputs. The one with the lowest perplexity is selected as the final answer. If this perplexity is above a threshold, indicating low confidence, the task is escalated to the next stage.

Datasets and models. We use the MT-Bench dataset (Zheng et al., 2023), which consists of multi-turn open-ended prompts spanning domains such as writing and roleplay. We use EXAONE-3.5-2.4B-Instruct as the small model, GPT-40-mini as the large model, and GPT-40 as the judge model.

Evaluation metrics. Since open-ended tasks lack unique ground-truth answers, we adopt the LLM-as-a-Judge protocol from MT-Bench (Zheng et al., 2023). A judge model compares two responses and assigns a win, loss, or tie. We report the win rate between COPE and the large model. We alternate the positions of outputs in pairwise comparisons to mitigate positional bias.

Results. Table 6 reports the win rate and inference cost on MT-Bench. COPE achieves a higher win rate than the large model baseline with much lower inference cost. Specifically, COPE obtains 36.9% wins, 28.1% ties, and 35.0% losses against the large model. This highlights that perplexity can serve as an effective confidence signal, and that COPE is applicable even to open-ended tasks.

5.4 AGENT TASKS

We further demonstrate the effectiveness of COPE on agent tasks, which require multi-step decision-making. At every action step, the planner and executor collaborate to decide on next action, where the consensus status is determined by majority voting.

Datasets and models. We conduct experiments on the evaluation split of the ALFWorld benchmark (Shridhar et al., 2020), which comprises 134 household tasks. In this benchmark, an LLM is employed as an agent, and it must predict actions over multiple steps by interacting with a text-based simulator. If the LLM-based agent successfully achieves the task, it receives the reward of 1.0; otherwise, 0.0. We use Qwen3-8B as the small model and GPT-4.1 as the large model.

Evaluation metrics. We report the success rate, which measures the percentage of tasks successfully resolved by the agent. In addition, we report the cost in USD, representing the total expense of running the evaluation on all 134 tasks.

Results. Table 7 shows that COPE improves success rate while substantially reducing cost. This demonstrates the effectiveness of COPE in agent tasks, where cost-efficiency is critical due to the long sequence of actions.

6 Conclusions

We introduce COPE, a test-time inference framework where small and large language models collaborate via planning. COPE improves accuracy and reduces cost by exchanging lightweight plans across planner-executor roles, without relying on few-shot exemplars or additional training. Across diverse benchmarks, COPE consistently outperforms prior baselines and large models alone, while significantly cutting inference cost. These findings highlight planning as a key mediator of efficient cross-model collaboration, even when human priors or training data are unavailable.

7 ETHICS STATEMENT

This work complies with the ICLR Code of Ethics. All datasets used are publicly available, and no human subjects or private data were involved. All models were used under their respective licenses (see Appendix D).

8 REPRODUCIBILITY STATEMENT

We provide details of datasets, model versions, prompts, and hyperparameters in the main text and appendices (see Appendix C and Appendix G). Code is released as the supplementary material to enable replication of results.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Art of Problem Solving. AIME Problems and Solutions. https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions, 2024. Accessed: 2025-05-19.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. arXiv preprint arXiv:2108.07732, 2021.
- Edward Beeching, Lewis Tunstall, and Sasha Rush. Scaling test-time compute with open models. URL https://huggingface.co/spaces/HuggingFaceH4/blogpost-scaling-test-time-compute.
- Kevin Black, Mitsuhiko Nakamoto, Pranav Atreya, Homer Walke, Chelsea Finn, Aviral Kumar, and Sergey Levine. Zero-shot robotic manipulation with pretrained image-editing diffusion models. *arXiv preprint arXiv:2310.10639*, 2023.
- Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on mixture of experts. *arXiv preprint arXiv:2407.06204*, 2024.
- Lingjiao Chen, Matei Zaharia, and James Zou. Frugalgpt: How to use large language models while reducing cost and improving performance. *arXiv preprint arXiv:2305.05176*, 2023.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- Yihong Dong, Jiazheng Ding, Xue Jiang, Ge Li, Zhuo Li, and Zhi Jin. Codescore: Evaluating code generation by learning code execution. *ACM Transactions on Software Engineering and Methodology*, 34(3):1–22, 2025.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2023.
- Google DeepMind. Gemini api pricing, 2024. URL https://ai.google.dev/gemini-api/docs/pricing. Accessed: 2025-05-19.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv* preprint arXiv:2103.03874, 2021.
- Yixin Ji, Juntao Li, Hai Ye, Kaixin Wu, Jia Xu, Linjian Mo, and Min Zhang. Test-time computing: from system-1 thinking to system-2 thinking. *arXiv preprint arXiv:2501.02497*, 2025.

- Michael I Jordan and Robert A Jacobs. Hierarchical mixtures of experts and the em algorithm.

 Neural computation, 6(2):181–214, 1994.
 - Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
 - Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
 - Steven Kolawole, Don Dennis, Ameet Talwalkar, and Virginia Smith. Agreement-based cascading for efficient inference. *arXiv preprint arXiv:2407.02348*, 2024.
 - Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, et al. Training language models to self-correct via reinforcement learning. *arXiv preprint arXiv:2409.12917*, 2024.
 - Hyunseok Lee, Seunghyuk Oh, Jaehyung Kim, Jinwoo Shin, and Jihoon Tack. Revise: Learning to refine at test-time via intrinsic self-verification. *arXiv preprint arXiv:2502.14565*, 2025.
 - Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pp. 19274–19286. PMLR, 2023.
 - Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.
 - Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
 - Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.
 - Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.
 - OpenAI. Openai api pricing. https://platform.openai.com/docs/pricing, 2025. Accessed: 2025-05-13.
 - Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
 - LG AI Research. Exaone 3.5: Series of large language models for real-world use cases. *arXiv* preprint arXiv:https://arxiv.org/abs/2412.04862, 2024.
 - Swarnadeep Saha, Xian Li, Marjan Ghazvininejad, Jason Weston, and Tianlu Wang. Learning to plan & reason for evaluation with thinking-llm-as-a-judge. *arXiv preprint arXiv:2501.18099*, 2025.
 - Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv* preprint arXiv:1910.01108, 2019.
 - Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.

- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv* preprint arXiv:2010.03768, 2020.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv* preprint arXiv:2408.03314, 2024.
- Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/qwen2.5/.
- Xinyi Wang, Lucas Caccia, Oleksiy Ostapenko, Xingdi Yuan, and Alessandro Sordoni. Guiding language model reasoning with planning tokens, december 2023b. *URL http://arxiv.org/abs/2310.05707*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.
- Murong Yue, Jie Zhao, Min Zhang, Liang Du, and Ziyu Yao. Large language model cascades with mixture of thoughts representations for cost-efficient reasoning. *arXiv* preprint arXiv:2310.03094, 2023.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.

Table 8: Truncated COPE. Using COPE with only Stage 1 and Stage 2 yields a dramatic reduction in inference cost.

Setting	Accuracy (%)	Cost
Large	73.8	355
COPE (S1+S2+S3)	74.4	212
COPE (S1+S2)	71.8	50

Table 9: Performance by difficulty. Improvements are larger for more difficult problems.

Method	Leve	1 1–2	Leve	1 3–4	Lev	el 5	To	otal
1/10/11/00	Acc	Cost	Acc	Cost	Acc	Cost	Acc ↑	Cost↓
Small (single)	78.9	_	54.9	_	23.1	_	52.8	_
Small (majority)	86.5	_	70.8	_	34.3	_	65.2	_
COPE (ours)	91.0	70	79.8	190	49.3	392	74.4	212

A EMPIRICAL STUDY

All experiments in this section are conducted on the MATH-500 dataset using EXAONE-3.5-2.4B-Instruct as the small model and GPT-4o-mini as the large model.

A.1 TRUNCATING AT STAGE 2: BALANCING EFFICIENCY AND ACCURACY

In Table 8, we compare three settings: a large model alone, full 3-stage COPE, and a truncated version that stops at Stage 2.

The 2-stage variant yields a dramatic reduction in inference cost (202 vs 50) while maintaining competitive accuracy (74.4% vs. 71.8%). This shows that, in many cases, the small model, once guided by a plan from the large model in Stage 2, produces reliable answers without further escalation.

This makes the truncated version a practical alternative when minimizing inference cost is a priority, offering a strong balance between performance and efficiency.

A.2 Performance by difficulty

Table 9 shows the performance breakdown by difficulty level. Compared to the Small (majority) baseline, COPE achieves consistent improvements across all segments. On the easiest Level 1–2 problems, accuracy increases from 86.5% to 91.0% (+4.5). For intermediate Level 3–4 problems, the gain is larger, from 70.8% to 79.8% (+9.0). The improvement is most pronounced on the most challenging Level 5 problems, where accuracy rises from 34.3% to 49.3% (+15.0). These results indicate that the benefits of COPE become more significant as problem difficulty increases.

A.3 EFFECT OF SAMPLE COUNT ON VOTING STABILITY

In Table 10, we study how the number of sampled solutions n influences performance in our framework. Since the final answer is selected via majority voting under a consensus threshold, increasing n does not always lead to higher overall accuracy. However, it introduces greater diversity in the reasoning paths explored, allowing the model to consider a broader range of plausible solutions.

Moreover, increasing n improves the stability and reliability of voting outcomes. With more samples, the consensus ratio becomes less sensitive to sampling noise, and the decisions become more consistent. We observe that the conditional accuracy under consensus also increases with n, suggesting that when agreement is reached, it is more likely to reflect a correct answer. Overall, larger sample counts enhance both the breadth of exploration and the trustworthiness of decisions.

Table 10: Effect of number of samples n on accuracy and conditional accuracy under consensus.

\overline{n}	Accuracy (%)	Accuracy @ Consensus (%)
4	71.4	78.5
8	74.4	84.1
32	74.4	85.7

Table 11: Effect of consensus thresholds. Each row shows the number of problems resolved in each stage under different threshold settings for Stage 1 and Stage 2.

Thre	shold	Resolved Problems		Results		
Stage 1	Stage 2	Stage 1	Stage 2	Stage 3	Accuracy (%)	Cost
0.5	0.5	312	72	116	70.0	173
0.5	0.75	326	37	137	70.8	175
0.75	0.5	204	150	146	74.4	212
0.75	0.75	205	96	199	74.2	258

A.4 EFFECT OF CONSENSUS THRESHOLD

Table 11 shows how varying the consensus thresholds in Stage 1 and Stage 2 affects both accuracy and cost. When both thresholds are low (e.g., 0.5/0.5), the small model tends to accept answers even when the agreement among solutions is weak. This leads to lower cost, since fewer problems are escalated, but often results in reduced accuracy due to incorrect answers being accepted too early.

On the other hand, setting both thresholds high (e.g., 0.75/0.75) makes the system more cautious, escalating more problems to the large model. While this improves accuracy, it also increases inference cost. Among the tested combinations, the 0.75/0.5 setting offers a more favorable balance. This threshold configuration effectively reduces unnecessary use of the large model, without sacrificing much in terms of final accuracy.

A.5 ERROR BARS

We report error bars in Table 12 as standard deviations computed over five independent runs with different random seeds.

B THEORETICAL INTERPRETATION

Each inference stage (Figure 3) follows three steps: (i) sample a $plan\ g$, (ii) generate a solution y conditioned on g, and (iii) aggregate the extracted answers by majority vote. We formalize the first stage (the template for all stages) with a concise latent-variable model.

Latent variable model We formalize prediction as computing the most likely answer a = f(y), where y is a solution and x is the input problem:

$$\hat{a} = f\left(\arg\max_{y} P(y \mid x)\right),\tag{2}$$

where $f(\cdot)$ denotes an answer extraction function that maps a complete solution y (e.g., a reasoning trace) to its final answer a (e.g., a scalar result).

Although one could sample solutions directly from $P(y \mid x)$, we instead introduce a latent plan g to structure the generation process. This yields the decomposition:

$$P(y \mid x) = \sum_{g} P(y \mid g, x) \cdot P(g \mid x). \tag{3}$$

Table 12: Mean accuracy and cost (with standard deviation).

	Method	Acc ↑	Cost↓	•
	Large COPE (ours)	73.8 ± 0.93 74.4 ± 0.88	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	_
Problem x	$\Rightarrow \begin{array}{c} \text{Plan } g \\ P(g x) \end{array}$	$\Rightarrow \begin{cases} Sc \\ P(y) \end{cases}$	$ g,x\rangle \longrightarrow$	Answer $a = f(y)$

Figure 3: Common inference module in COPE.

This decomposition enables exploration of diverse reasoning strategies, each corresponding to a distinct plan. We sample from the resulting process by drawing $g_i \sim P(g \mid x)$ and then $y_i \sim P(y \mid g_i, x)$.

This process yields a set of solutions y_1, \ldots, y_n sampled from $P(y \mid x)$. From each, we extract an answer $a_i = f(y_i)$. We then approximate the objective in Eq. equation 2 via majority voting:

$$f\left(\arg\max_{y} P(y\mid x)\right) \approx \operatorname{mode}(f(y_1),\dots,f(y_n)).$$
 (4)

This procedure offers a simple approximation to the most likely answer under $P(y \mid x)$. Since higher-probability y_i are sampled more often, their corresponding answers $f(y_i)$ dominate the vote. The latent plan g encourages diversity by inducing varied reasoning strategies. This improves robustness, helping recover correct answers that might be missed by direct decoding.

Bayesian inference From a Bayesian perspective, we are interested in the posterior predictive distribution $P(a \mid x)$ over answers:

$$\sum_{g} \int_{\mathcal{Y}} \delta(f(y) = a) P(y \mid g, x) P(g \mid x) dy, \tag{5}$$

where $P(g \mid x)$ serves as a prior over plans and $P(y \mid g, x)$ is the likelihood of a solution path.

We approximate this distribution via Monte Carlo sampling: $g_i \sim P(g \mid x)$, $y_i \sim P(y \mid g_i, x)$, and set $a_i = f(y_i)$. Majority voting over $\{a_i\}$ yields an estimate $\hat{a} \approx \arg \max_a P(a \mid x)$.

We treat the consensus ratio as an empirical estimate of $P(a = \hat{a} \mid x)$, and use it to decide whether to accept the prediction or escalate.

Mixture-of-experts (MoE) Our approach shares several structural similarities with MoE models (Jordan & Jacobs, 1994; Cai et al., 2024):

First, sampling a plan $g \sim P(g \mid x)$ corresponds to selecting an expert based on the input, akin to input-dependent gating. Second, each plan defines a distinct strategy for solving the problem, much like how individual experts specialize in different regions of the input space. Third, given a selected plan, the model generates a solution via $P(y \mid g, x)$, which parallels how an expert produces an output conditioned on its assignment.

C FURTHER IMPLEMENTATION DETAILS

We conduct all experiments using 8 NVIDIA RTX-3090 GPUs. Our implementation builds upon and modifies the codebase from the search-and-learn repository by Hugging Face (Beeching et al.). We provide our code in the supplementary material.

Model version We use the default API versions of each large model and ensure consistency across all experiments. Specifically, GPT-40-2024-08-06 is used for GPT-40, GPT-40-mini-2024-07-16 for GPT-40-mini, and Gemini-2.0-Flash-001 for Gemini-2.0-Flash.

Hyperparameters For code generation, we use a batch size of 32, and for mathematical reasoning, a batch size of 160. In COPE, decoding temperatures are set as follows: in Stage 1, the planner uses a temperature of 0.8 and the executor uses 0.4; in Stage 2, both the planner and executor use 0.8; and in Stage 3, the executor uses greedy decoding (temperature = 0).

Prompts For mathematical reasoning, we adopt the system prompt from Beeching et al. for the executor. All other prompts, including those for planners are designed by us. The system prompts are provided in Appendix G.

D LICENSE

 We use GPT-40 and GPT-40-mini via their official APIs, which are subject to OpenAI's commercial license. LLaMA, EXAONE, Qwen models are used under Llama 3.2 Community License Agreement, EXAONE AI Model License Agreement, and Apache license 2.0, respectively. All datasets used in this work are publicly available and used under their respective licenses.

E BROADER IMPACT

As AI systems increasingly span both edge devices and the cloud, collaborative inference frameworks like COPE may help shape efficient and scalable deployment strategies. By reducing unnecessary compute, this approach can also contribute to more sustainable AI usage.

In addition, COPE reduces reliance on expensive proprietary models, making high-quality language model inference more accessible to users with limited resources. This may support broader adoption of AI tools across diverse settings, including education, research, and local deployment.

F QUALITATIVE EXAMPLES

Vanilla

Problem: Six cars pull up to a red light, one at a time. At the light, there are three lanes, one left-turn lane, one straight-going lane, and one right-turn lane. How many ways can the cars stack up so that all three lanes are occupied? Note that if the first car turns left and the second goes straight, this is considered different from the first car going straight and the second car turning left. In other words, the cars are distinguishable, but pull up to the intersection in a fixed order.

Solution (executor): To solve this problem, we need to determine the number of ways six distinguishable cars can be arranged such that each lane (left-turn, straight-going, right-turn) is occupied by exactly one car. Step-by-Step Solution: Step 1: Assigning Cars to Lanes Since there are three lanes and six cars, one lane will remain empty. We need to choose which lane will be empty. There are 3 ways to choose which lane will be empty (left-turn, straight-going, or right-turn). Step 2: Arranging Cars in Chosen Lanes Once we have chosen which lane will be empty, we need to arrange the remaining five cars in the three lanes. The number of ways to arrange n distinguishable objects in k distinguishable bins is given by k^n . Here, we have 5 cars to arrange in 3 lanes: $3^5 = 243$ Step 3: Combining Both Steps Since there are 3 ways to choose which lane will be empty and for each choice, there are 3^5 ways to arrange the cars in the remaining lanes: $3 \times 243 = 729$ Therefore, the total number of ways the cars can stack up so that all three lanes are occupied is: $\boxed{729}$

COPE

Problem: Six cars pull up to a red light, one at a time. [...]

Goal (planner): The goal of the problem is to determine the number of distinct permutations of six distinguishable cars parking in three lanes (left turn, straight, right turn) such that each lane is occupied by at least one car.

Guideline (planner): To solve this problem, first determine the total number of ways to assign the cars to the three lanes such that each lane has at least one car. This can be achieved using the principle of inclusion-exclusion: calculate the total arrangements of cars in the lanes, subtract the arrangements where at least one lane is empty, and then add back arrangements where two lanes are empty. Finally, ensure that each arrangement respects the order of the cars.

Solution (executor): To solve this problem efficiently using combinatorial principles, we need to ensure that each lane (left-turn, straight-going, right-turn) has at least one car while accommodating six distinguishable cars in three lanes sequentially. Step-by-Step Solution: Step 1: Total Arrangements Without Restrictions. First, calculate the total number of ways to assign six cars into three lanes without any restrictions. Each car has 3 choices (left, straight, right): $3^6 = 729$. Step 2: Subtract Arrangements Where At Least One Lane is Empty. Next, subtract the cases where at least one lane is empty. We use the principle of inclusion-exclusion for this purpose: One Lane Empty: Choose which lane is empty (3 ways), and distribute the cars among the remaining two lanes (2^6 ways for each choice): $3 \cdot 2^6 = 3 \cdot 64 = 192$. Two Lanes Empty: Choose which two lanes are empty ($3 \cdot 2^6 = 3 \cdot 64 = 192$) and place all cars into the remaining lane (1 way): $3 \cdot 2^6 = 3 \cdot 1 = 3$. Using inclusion-exclusion: $3 \cdot 2^6 = 3 \cdot 6^6 = 3 \cdot 6$

Figure 4: Comparison between vanilla and COPE inference on a MATH-500 problem. The vanilla solution (left), generated by EXAONE-3.5-2.4B-Instruct, results in incorrect reasoning (red). In contrast, COPE combines a goal from the same model with a guideline from GPT-40-mini, highlighting key constraints (green) and guiding the executor to a correct solution

Vanilla

Problem: Suppose a and b are different prime numbers greater than 2. How many whole-number divisors are there for the integer $a(2a+b)-2a^2+ab$?

Solution (executor): To determine the number of whole-number divisors of the integer $N=a(2a+b)-2a^2+ab$, where a and b are distinct prime numbers greater than 2, we start by simplifying N. Step 1: Simplify N. First, let's rewrite $N:N=a(2a+b)-2a^2+ab$. Expanding and combining like terms: $N=2a^2+ab-2a^2+ab=2ab$. Thus, we have: N=2ab. Step 2: Analyze the Divisors of N. Since N=2ab and a and b are distinct prime numbers greater than 2, we note the following: -a and b are both odd primes (since they are greater than 2). - The product ab is odd because the product of two odd numbers is odd. - Therefore, N=2ab is an even number because it includes the factor 2. Step 3: Determine the Divisors. To find the divisors of N=2ab: - The prime factorization of N is $2^1 \times ab$. - The divisors of N are all combinations of the factors 2^1 and ab. The possible divisors are: 1, 2, ab, 2ab. Conclusion. Counting these distinct divisors, we have: $\{1, 2, ab, 2ab\}$. Thus, there are **4** whole-number divisors of N. Therefore, the final answer is: 4.

COPE

Problem: Suppose a and b are different prime numbers greater than 2. [...]

Goal (planner): The goal of the problem is to determine the exact number of distinct whole-number divisors of the polynomial expression $a(2a+b)-2a^2+ab$, given that a and b are distinct prime numbers greater than 2.

Solution (executor): To determine the number of whole-number divisors of the expression $a(2a+b)-2a^2+ab$, we start by simplifying the expression algebraically. Step 1: Simplify the Expression. First, let's rewrite and simplify the given polynomial expression: $a(2a+b)-2a^2+ab$. Expanding and combining like terms: $a\cdot 2a+a\cdot b-2a^2+ab=2a^2+ab-2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+ab=2a^2+$

Figure 5: Comparison between vanilla and COPE inference on a MATH-500 problem. The vanilla solution (left), generated by EXAONE-3.5-2.4B-Instruct, results in incorrect reasoning (red). In contrast, COPE combines a goal from the same model, highlighting a key condition (green) and guiding the executor to a correct solution.

G PROMPTS

 Table 13: System prompts used in COPE for mathematical reasoning

Executor System Prompt

Solve the following math problem efficiently and clearly:

- For simple problems (2 steps or fewer): Provide a concise solution with minimal explanation.
- For complex problems (3 steps or more): Use this step-by-step format:

```
## Step 1: [Concise description]
[Brief explanation and calculations]
## Step 2: [Concise description]
[Brief explanation and calculations]
```

Regardless of the approach, always conclude with:

Therefore, the final answer is: \$\boxed{answer}\$. I hope it is correct.

Where [answer] is just the final number or expression that solves the problem.

Planner System Prompt

Goal Generation: You are a math teacher. State the goal of the following math problem in one sentence.

Plan Generation: You are a skilled math coach. Explain how to solve the following math problem. Focus on strategy and key ideas. Respond in just one or two sentences. Return only the guidelines as plain text.

Table 14: System prompts used in COPE for code generation

Executor System Prompt

You are an expert Python programmer. Solve the following coding problem efficiently and clearly.

Planner System Prompt

Goal Generation: You are an expert Python programmer. State the goal of the following coding problem in one sentence.

Plan Generation: You are an expert Python programmer. Explain how to solve the following coding problem. Focus on strategy and key ideas. Respond in just one or two sentences. Return only the guidelines as plain text.