

# Open-ended Code Generation Challenges for Evolving Intelligence

Qiuyang Mang *on behalf of Frontier-CS Team*

{SLICE, SKY} @ UC Berkeley





SWE-bench Verified  $\approx$  **80%**

AIME  $\approx$  **100%**

Humanity Last Exam  $\approx$  **45%**



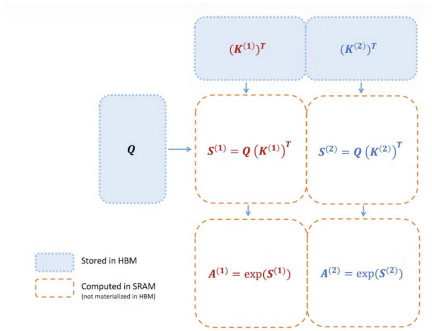
LLMs now **saturate** real exams and exam-style benchmarks.

What's next?



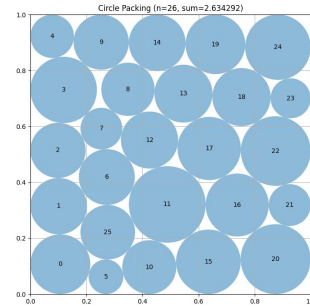
*Open-ended Problems*

# Beyond Passing Exams: Exploring New Frontiers



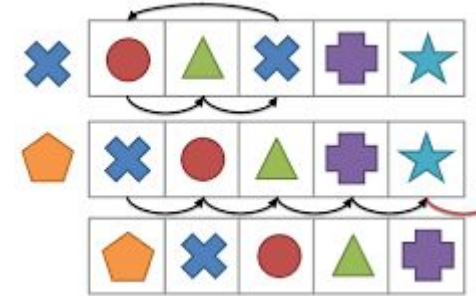
LLM as a performance optimizer

*How much speedup can we achieve with a GPU kernel?*



LLM as an algorithm designer

*How many circles can we pack into a fixed region?*



LLM as a system researcher

*How much cache miss cost can we reduce with a replacement policy?*

**No optimal solutions, only better frontiers.**

# We Don't Know the **Best**, but We Know How **Good**



## Objectively Verifiable Evaluation

Open, deterministic evaluation



No LLM-as-a-judge

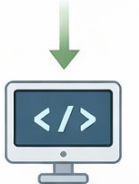
### OpenEvolve

```
def evolve(code):
    while not optimal:
        code = mutate(code)
        evaluate(code)
```

*Evolutionary Coding Agent*



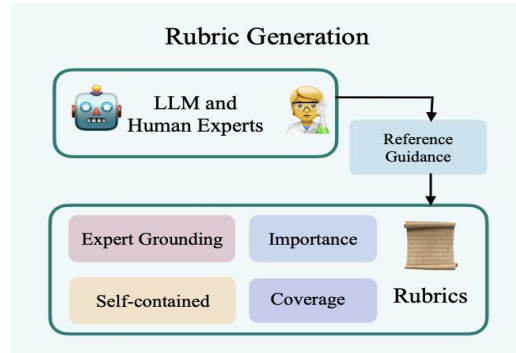
Anybody can rerun it



Runs it independently

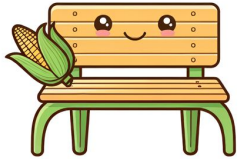


Gets the same score



⇒ *Continuous scores guide agentic evolution and support post-training*

# LLMs for Open-Ended Problems: Where We Are Limited

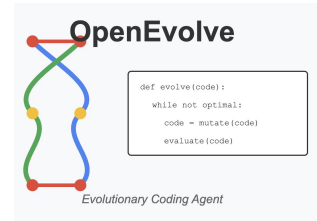


**KernelBench:** ~200 GPU programming problems (single-domain)



**AI-Driven Research Systems**

**ADRS:** ~10 system research problems



**OpenEvolve:** 14 example problems



**ALE-bench:** 40 Atcoder Heuristic Contest problems (single-source)

We need a **large-scale benchmark** for open-ended problems, spanning multiple domains, sources, and problem types.

# LLMs for Open-Ended Problems: Where We Are Limited

arXiv:2511.23473v1 [cs.LG] 28 Nov 2025

## ThetaEvolve: Test-time Learning on Open Problems

Yiping Wang<sup>1,2</sup> Shao-Rong Su<sup>1</sup> Zhiyuan Zeng<sup>1</sup> Eva Xu<sup>1</sup> Liliang Ren<sup>2</sup> Xinyu Yang<sup>3,2</sup> Zeyi Huang<sup>4,2</sup>  
 Xuehai He<sup>2</sup> Luyao Ma<sup>3</sup> Baolin Peng<sup>2</sup> Hao Cheng<sup>2</sup> Pengcheng He<sup>2</sup> Weizhu Chen<sup>2</sup> Shuohang Wang<sup>2</sup>  
 Simon Shoalei Du<sup>1</sup> Yelong Shen<sup>1,2</sup>

### Abstract

Recent advances (LLMs) have enabled mathematical discovery; a closed-source system improves bounds on relies on ensemble new bounds and its models cannot inter We introduce *Theta* work that simplifies efficiently scale but enforcement Learning models to continue in improving ThetaEvolve featuregram database for sampling for higher discourage stagnation shaping for stable) evolve is the first able a small opens-RI-0525-Queen-8) bounds on open p first auto-correlation phatEvolve. Beside open tasks, we fit at test-time consist only baselines, and ing capabilities, as demonstrate faster formance on both unseen tasks. We r

This work was done at Microsoft, University of Michigan, University of California, and Tsinghua University. We thank Yiping Wang, Simon Shoalei Du, and the anonymous reviewers for their helpful comments.

Preprint.  
<https://github.com>

arXiv:2510.03851v1 [cs.AI] 4 Oct 2025

### 1. Introduction

#### ALGORITHM GENERATION VIA CREATIVE IDEATION

Ruiying Ma  
 University of California, Berkeley

Chieh-Jan Mike Liang  
 Microsoft Research

Yanjie Gao  
 Microsoft Research

Francis Y. Yen  
 University of Illinois Urbana-Champaign

#### ABSTRACT

Designing system algorithms remains challenging, where the discontinuous nature of the solution space often forces system engineers to rely on generic heuristics at the expense of performance. We study whether LLMs can practically drive algorithm generation, and find that they are biased towards well-known generic designs, rather than making the creative leaps needed to navigate the discontinuous solution space. To address this limitation, we introduce *MetaMuse*, a framework for creative ideation built on three self-reflection principles: (1) quantifying solution diversity and usefulness in measurable performance space, rather than abstract idea space, (2) steering ideation through external stimuli, rather than internal randomness, and (3) constructing executable solutions using waypoints reasoning, rather than free-form chain-of-thought. Extensive evaluation shows that *MetaMuse* can generate high-performing solutions for two critical problems at a global cloud provider: cache replacement (reducing cache misses by up to 35.76%) and online bin packing (reducing bin usage by up to 30.93%).

#### 1 INTRODUCTION

Designing system algorithms continues to be a central challenge in computing systems. Traditionally, the development of such algorithms has been a manual and labor-intensive process. Our experience at a global cloud provider indicates that even seemingly simple algorithms used in production — such as cache replacement for data storage or bin packing for job scheduling — can require tens of thousands of engineering hours to design. As a result, practitioners often resort to generic heuristics from the literature, e.g., least-recently used (LRU) for cache replacement and first-fit for bin packing, which frequently result in suboptimal performance.

This paper asks whether large language models (LLMs) can practically drive algorithm generation, with an emphasis on principles to transform this task into a systematic process. The core challenge in system algorithm design arises from the nature of its solution space: it is an inherently discontinuous space, where even a small change in algorithm design (e.g., data structure or control flow) can lead to sharp and non-linear changes in performance. Although it is sometimes possible to estimate the upper-bound performance, searching for practical solutions that approach this bound remains non-trivial. Furthermore, the discontinuous solution space does not provide sufficiently predictable patterns or a smooth landscape to guide the search.

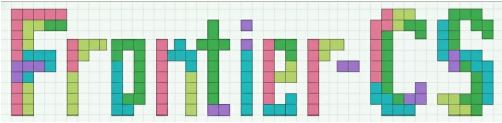
Due to this discontinuity, we approach the algorithm generation task from a different angle, and frame it as a sampling process in the solution space. LLM attempts to generate distinct solutions at each step. This generative process represents a sequence of leaps in discontinuous solution space (Eubank et al., 2023), which we formulate as *creative ideation* for LLMs. In fact, the systems community has long hypothesized algorithm design as a discovery process of ideas (Kant, 1985).

To study the algorithm generation task, we focus on high-impact problems at a global cloud provider: cache replacement and online bin packing. Our initial attempts of repeatedly sampling GPT-4o, Llama3.3-70B, and DeepSeek-V3 show that LLMs are fundamentally hindered by *availability bias* (Tversky & Kahneman, 1973) — LLMs are trained to output the most likely sequence of words, according to training datasets. As a result, solutions tend to cluster around well-known heuristics in the literature, e.g., least-recently (LRU) and least-frequently used (LFU) for caching.






*ThetaEvolve*: Evaluated on only **5** tasks from OpenEvolve Examples: CirclePacking-T / FirstAutoCorrIneq / SecondAutoCorrIneq / ThirdAutoCorrIneq / HadamardMatrix

*MetaMuse*: Evaluated on only **2** tasks by their authors: Bin Packing / Cache Replacement

# Introducing Frontier-CS



= **200+** open-ended problems adapted from competitive programming and real CS research

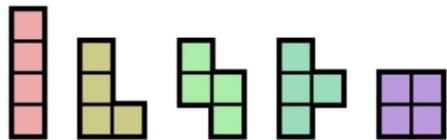
- **Algorithmic Track:** 10+ experts converse exam-style problems from multiple sources into an open-ended style by *changing objectives, adding constraints* . . .      **CODEFORCES**
- **Research Track:** CS PhDs extract the core algorithmic problems from their research interests, including OS, HPC, AI, DB, PL, . . .

# An Open-Ended, Verifiable Example with Continuous Score

## FrontierCS

open-ended

verifiable

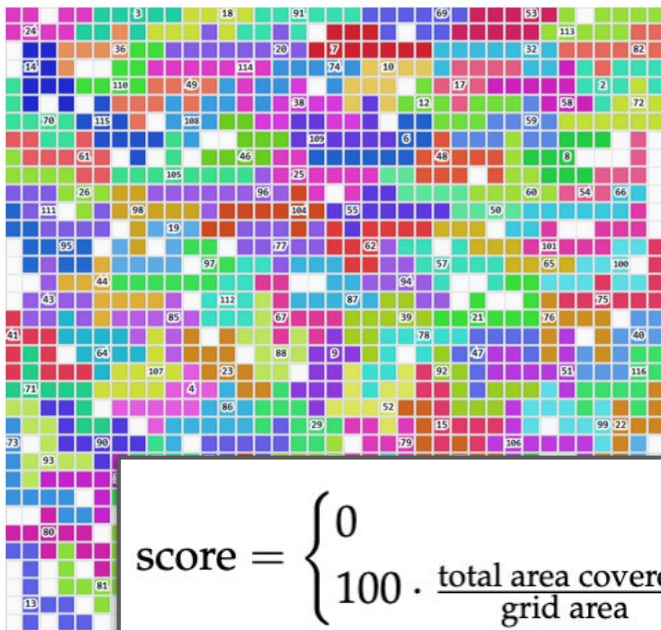


### Example: Polyomino Packing

Pack all polyominoes as tightly as possible into the grid.

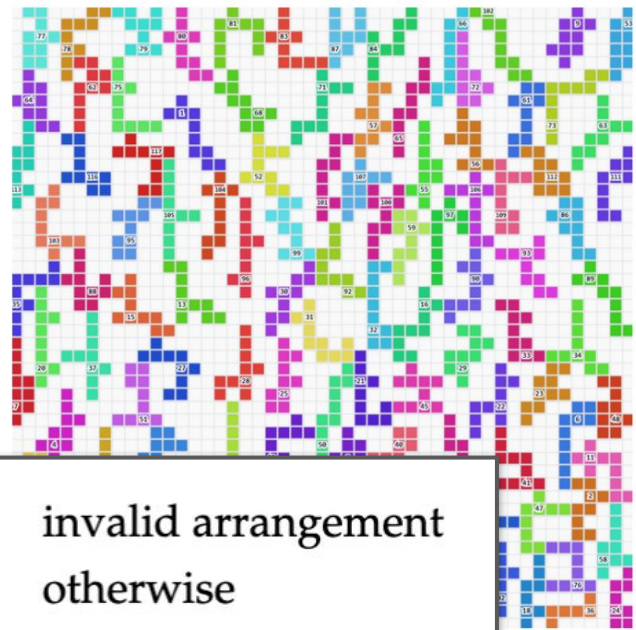
Human Expert

87%



GPT-5 Thinking

47%

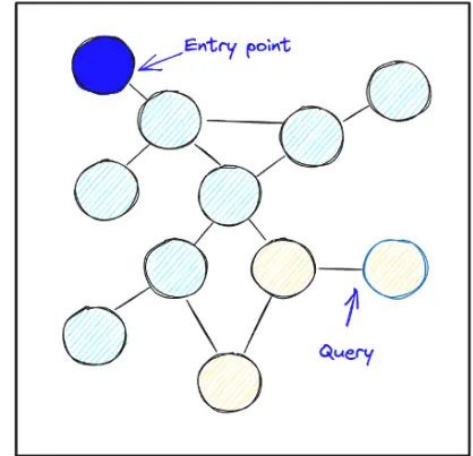


$$\text{score} = \begin{cases} 0 & \text{invalid arrangement} \\ 100 \cdot \frac{\text{total area covered}}{\text{grid area}} & \text{otherwise} \end{cases}$$

invalid arrangement  
otherwise

# Research Problem Example: VectorDB Design

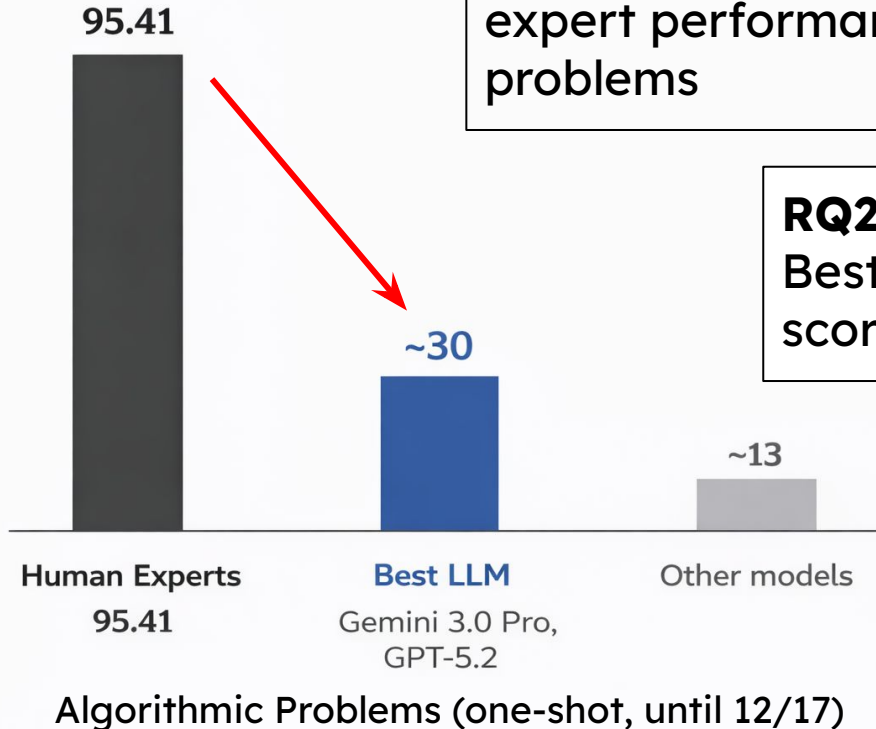
- Given: Millions of vectors and a set of query vectors.
- Task: For each query, return the nearest stored vector as fast as possible, while being correct at least 80% of the time.
- Scoring:  $\text{Recall@1} \geq 80\%$  is required; passing solutions receive a 0–100 score normalized by average query latency relative to a human SOTA.



## Research Questions on Frontier-CS

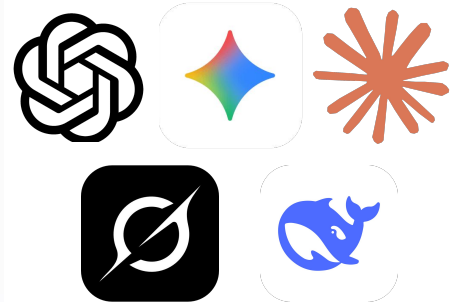
- RQ1: As LLMs saturate code exams, how do they perform on open-ended CS problems?
- RQ2: Does test-time scaling via multiple attempts still hold for open-ended CS problems?
- RQ3: Does longer chain-of-thought test-time scaling still work?
- RQ4: What failure modes do LLMs exhibit on open-ended CS problems?

# Open-ended Problems Remain Challenging for Frontier LLMs

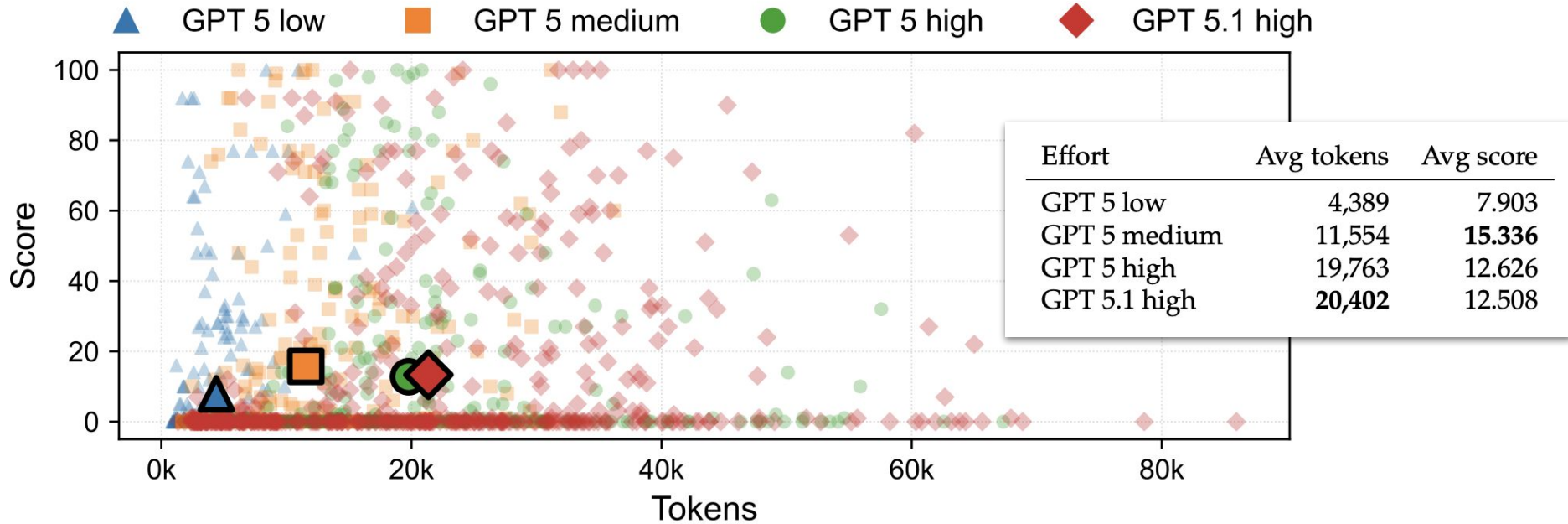


**RQ1:** LLMs remains a major gap to human expert performance on open-ended problems

**RQ2:** Multiple attempts matter  
Best-of-5 improves average score by **+6.5 - +22.7**



# Improving Reasoning Effort Does Not Yield Further Gains



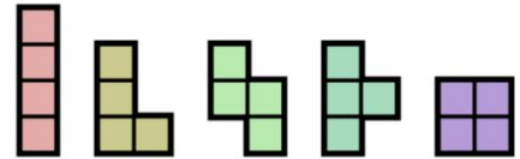
**RQ3:** Longer chain-of-thought not always help open-ended problem solving.

# Misleading Micro-Optimization Trap

- LLMs often fail to identify which choices are algorithmically meaningful, becoming trapped in micro optimizations.
- **GPT-5 solution:** starts with a suboptimal core idea, augmented with several low-level memory and efficiency optimizations  $\Rightarrow$  **~40 points**

Human solution & prompt fix:

*Please use a 2D array to maintain the rectangle state, and convert to the required format only at the end  $\Rightarrow$  ~80 points*



## Example: Polyomino Packing

Pack all polyominoes as tightly as possible into the grid.

# Harbor Integration

# Agentic Benchmarks

# Impact and Future Directions of Frontier-CS

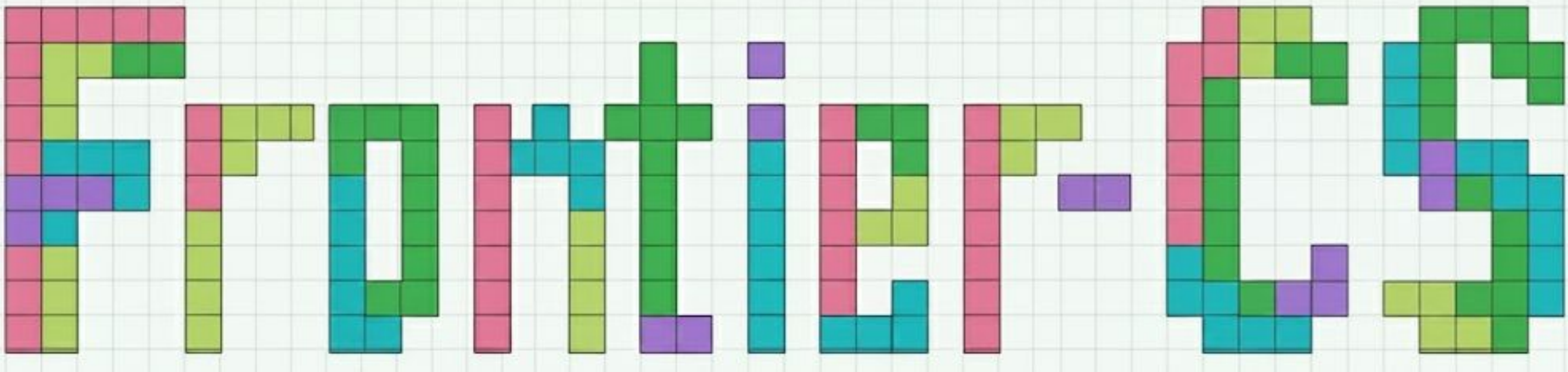
## ***Impact:***

- Provides a large, diverse collection of open-ended problems
- Enables reproducible, fine-grained evaluation via deterministic, verifiable and continuous scoring
- The project has attracted early industry interest (e.g., xAI, ByteDance), and welcomes contributions of problems or compute resources.



## ***What's next?***

- **Evolving:** From domain-specific tricks to regularized, general evolving frameworks.
- **Synthesis:** Can LLMs synthesize new open-ended problems?
- **Post-training:** Can models learn new reasoning abilities from open-ended problems with continuous reward signals?



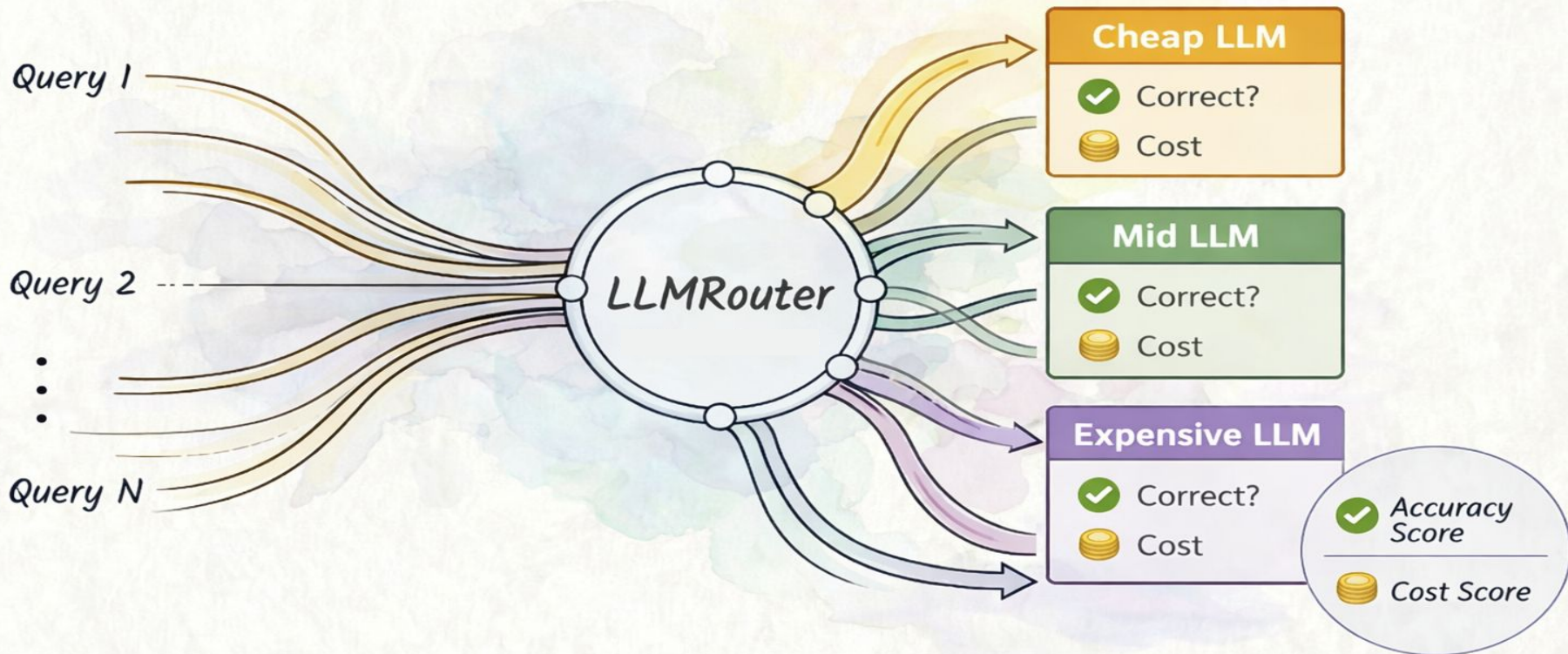
# Thanks!

Website: [frontier-cs.org](https://frontier-cs.org)

GitHub: <https://github.com/FrontierCS/Frontier-CS>



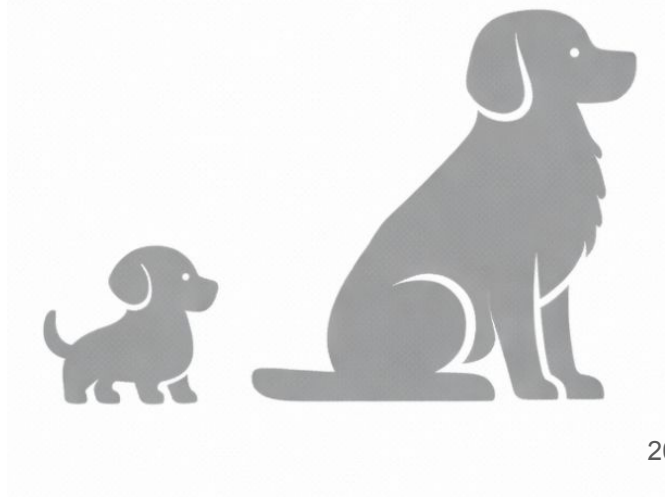
# Research Problem: LLM Router



Learn routing policy from a reference dataset

## LLMs for Open-Ended Problems: Where We Are Limited

- Toy-scale benchmarks cannot distinguish general algorithmic discovery ability from task overfitting.
- Toy-scale data can support framework design, but cannot train new reasoning abilities.



# FrontierCS: Evolving Challenges for Evolving Intelligence

**Contributors ( \*equal contribution )**

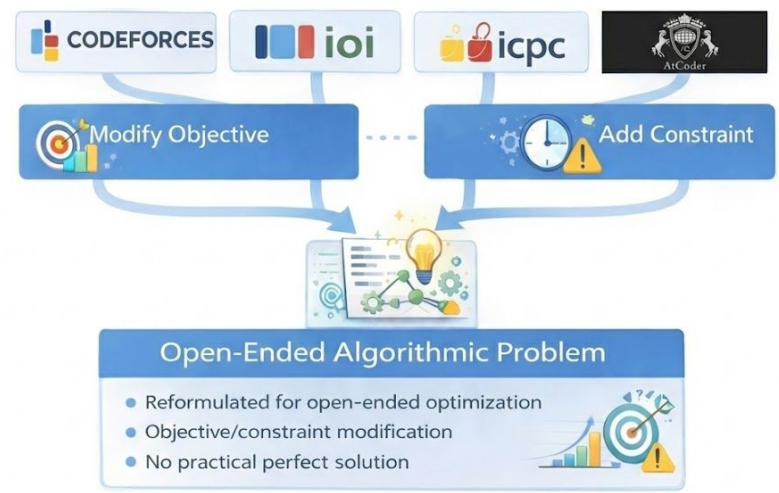
Qiuyang Mang<sup>1,\*</sup>, Wenhao Chai<sup>2,\*</sup>, Zhifei Li<sup>1,\*</sup>, Huanzhi Mao<sup>1,\*</sup>, Shang Zhou<sup>3,\*</sup>, Alexander Du<sup>1,4,\*</sup>,  
 Hanchen Li<sup>1,\*</sup>, Shu Liu<sup>1,\*</sup>, Edwin Chen<sup>5</sup>, Yichuan Wang<sup>1</sup>, Xieting Chu<sup>6</sup>, Zerui Cheng<sup>2</sup>, Yuan Xu<sup>4</sup>, Tian Xia<sup>1</sup>,  
 Zirui Wang<sup>1</sup>, Tianneng Shi<sup>1</sup>, Jianzhu Yao<sup>2</sup>, Yilong Zhao<sup>1</sup>, Qizheng Zhang<sup>7</sup>, Charlie Ruan<sup>1</sup>, Zeyu Shen<sup>2</sup>,  
 Kaiyuan Liu<sup>8</sup>, Runyuan He<sup>1</sup>, Dong Xing<sup>4</sup>, Zerui Li<sup>4</sup>, Zirong Zeng<sup>1</sup>, Yige Jiang<sup>9</sup>, Lufeng Cheng<sup>10</sup>, Ziyi Zhao<sup>11</sup>,  
 Youran Sun<sup>1</sup>, Wesley Zheng<sup>1</sup>, Meiyuwang Zhang<sup>5</sup>, Ruyi Ji<sup>12</sup>, Xuechang Tu<sup>6</sup>, Zihan Zheng<sup>13</sup>, Zexing Chen<sup>3</sup>,  
 Kangyang Zhou<sup>14</sup>, Zhaozi Wang<sup>13</sup>, Jingbang Chen<sup>5</sup>

**Advisors ( \*equal advising )**

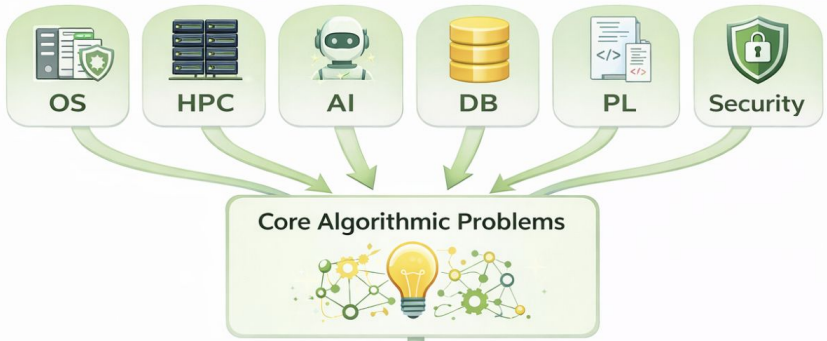
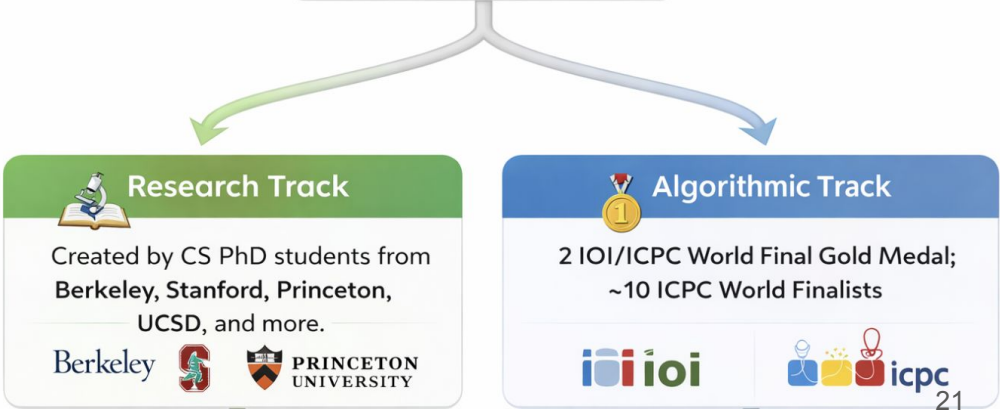
Aleksandra Korolova<sup>2</sup>, Peter Henderson<sup>2</sup>, Pramod Viswanath<sup>2</sup>, Vijay Ganesh<sup>6</sup>, Saining Xie<sup>13</sup>, Zhuang Liu<sup>2</sup>,  
 Dawn Song<sup>1</sup>, Sewon Min<sup>1</sup>, Ion Stoica<sup>1</sup>, Joseph E. Gonzalez<sup>1,\*</sup>, Jingbo Shang<sup>3,\*</sup>, Alvin Cheung<sup>1,\*</sup>

**Affiliations**

<sup>1</sup>UC Berkeley <sup>2</sup>Princeton University <sup>3</sup>UCSD <sup>4</sup>X-camp Academy <sup>5</sup>Independent <sup>6</sup>Georgia Tech  
<sup>7</sup>Stanford University <sup>8</sup>University of Washington <sup>9</sup>Nanyang Technological University  
<sup>10</sup>University of Toronto <sup>11</sup>UIUC <sup>12</sup>University of Michigan <sup>13</sup>New York University <sup>14</sup>MIT



~200 Problems

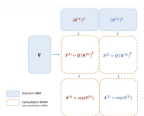


## Background: LLMs Near Ceiling on exams

SWE-bench Verified  $\approx 80\%$   
AIME  $\approx 100\%$   
Humanity Last Exam  $\approx 45\%$

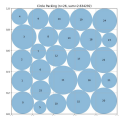


## What's Next for Computer Science Tasks?



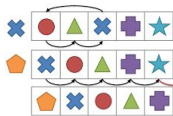
LLM as a performance optimizer

How much speedup can we achieve with a GPU kernel?



LLM as an algorithm designer

How many circles can we pack into a fixed region?



LLM as a system researcher

How much cache miss cost can we reduce with a replacement policy?

**No best solutions, only better frontiers.**



Objectively Verifiable Evaluation  
Open, deterministic evaluation



**Unsolved  
Open-ended  
Verifiable  
Diverse  
Continuous Score**

$\Rightarrow$  Continuous scores guide *agentive evolution* and support *post-training*

## Where We Are Limited

Existing research on open-ended CS problems often evaluates LLMs on a few **toy tasks**.



**KernelBench**: ~200 GPU programming problems (single-domain)



**ADRS**: ~10 system research problems



**OpenEvolve**: 14 example problems

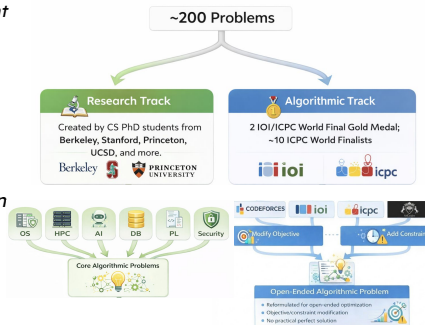


**ALE-bench**: 40 Atcoder Heuristic Contest problems (single-source)

Existing open-ended CS benchmarks are often **domain-specific** or **small-scale**.

We need a **large-scale benchmark** for open-ended problems, spanning **multiple domains, sources, and problem types**.

## Frontier-CS: a large, diverse datasets for open-ended CS problems



**Research Track** problems are constructed by distilling core algorithmic challenges from real CS research problems and asking LLMs to implement practical solutions.

**Algorithmic Track** problems are created by transforming existing exam-style tasks into open-ended settings through modified objectives and added constraints, or by collecting classical NP problems and conjectures.

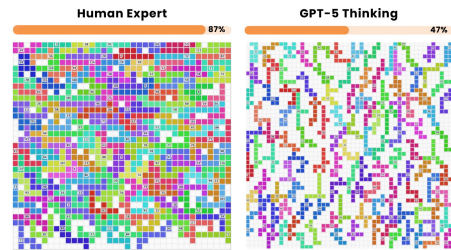
## Frontier-CS Example Problems

### FrontierCS

unsolved open-ended verifiable diverse



Example: **Polyomino Packing**  
Pack all polyominoes as tightly as possible into the grid.



### World Map (IOI 2025 $\Rightarrow$ Frontier-CS)

Given several countries and which pairs must be adjacent.  
Place country labels on a square grid.  
Adjacent cells imply country adjacency.  
All required adjacencies must appear.  
Goal: use the **smallest** possible grid.



(b) Human expert ( $K = 7$ )



(c) GPT-5 ( $K = 13$ )

**Research Problems:** Cloud Scheduling, VectorDB Design, POC generation, Grammar Fuzzing, RL Algorithms, GPU Kernel Design, and more

## Key Findings

### Open-ended Problems Remain Challenging for LLMs

One-shot reasoning remains a major gap to human expert performance on open-ended problems



Multiple attempts matter  
Best-of-5 improves average score by  $+6.5$   $\rightarrow$   $+22.7$

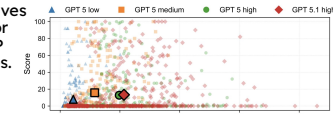
### Misleading Micro-Optimization Trap

LLMs often fail to identify which choices are algorithmically meaningful, becoming trapped in micro optimizations.

**GPT-5 solution:** Maintains the search state as a **list of polyomino placements**, augmented with several low-level memory and efficiency optimizations  $\Rightarrow$   $-40$  score

Human solution & prompt fix:  
Please use a 2D array to maintain the rectangle state, and convert to the required format only at the end  $\Rightarrow$   $-80$  score

### Improving Reasoning Effort Does Not Yield Further Gains



Effort	Avg tokens	Avg score
GPT-5 low	4,389	7.40
GPT-5 medium	11,554	15.56
GPT-5 high	19,763	12.62
GPT-5.1 high	20,402	12.508