



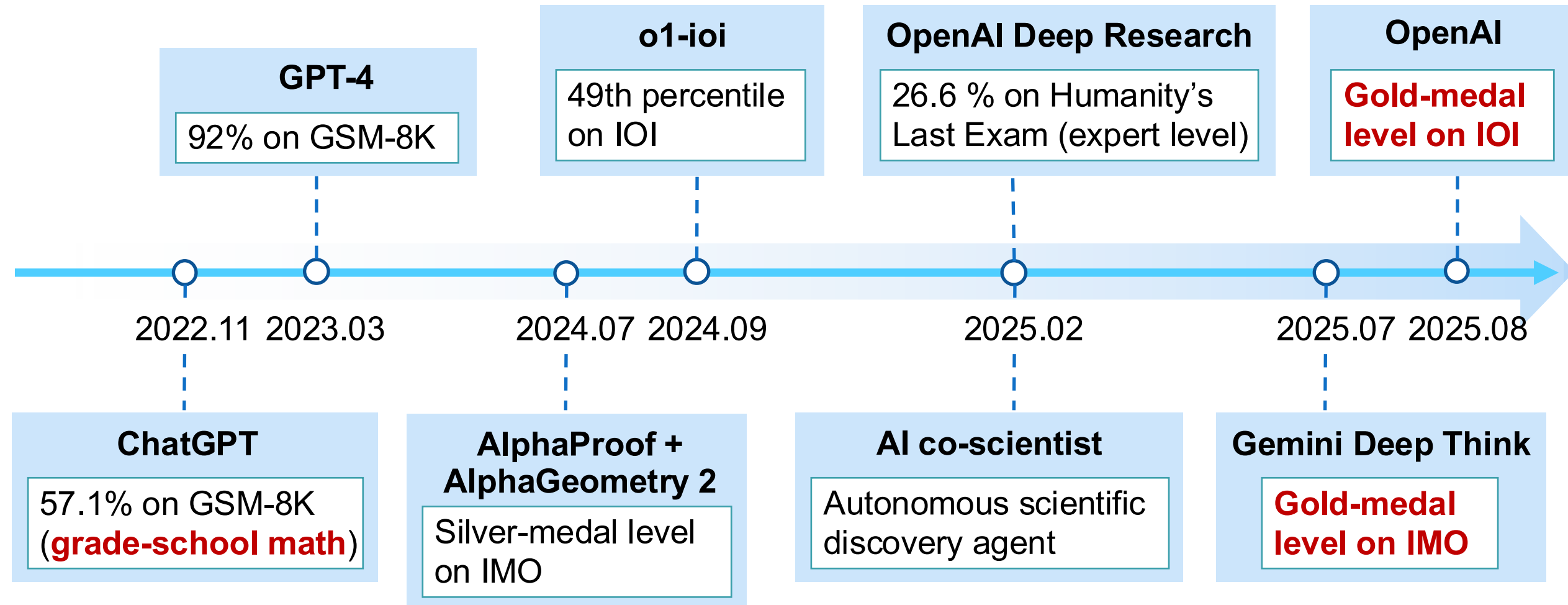
# 数学建模智能体理论与实战

刘浩 博士  
香港科技大学（广州）  
2025年12月26日

# 数学建模智能体

MM-Agent: LLM as Agents for Real-world Mathematical Modeling Problem. NeurIPS 2025

# The Intellectual Development of LLMs



# The Reasoning Barrier in Real-World Applications

## Well-defined

**Problem 1.** A line in the plane is called sunny if it is not parallel to any of the x-axis, the y-axis, and the line  $x + y = 0$ .

Let  $n \geq 3$  be a given integer. Determine all nonnegative integers  $k$  such that there exist  $n$  distinct lines in the plane satisfying both of the following:

- for all positive integers  $a$  and  $b$  with  $a + b \leq n + 1$ , the point  $(a, b)$  is on at least one of the lines; and
- exactly  $k$  of the  $n$  lines are sunny.

## Well-defined

Gardens by the Bay is a large nature park in Singapore. In the park there are  $n$  towers, known as supertrees. These towers are labelled 0 to  $n - 1$ . We would like to construct a set of **zero or more** bridges. Each bridge connects a pair of distinct towers and may be traversed in **either** direction. No two bridges should connect the same pair of towers.

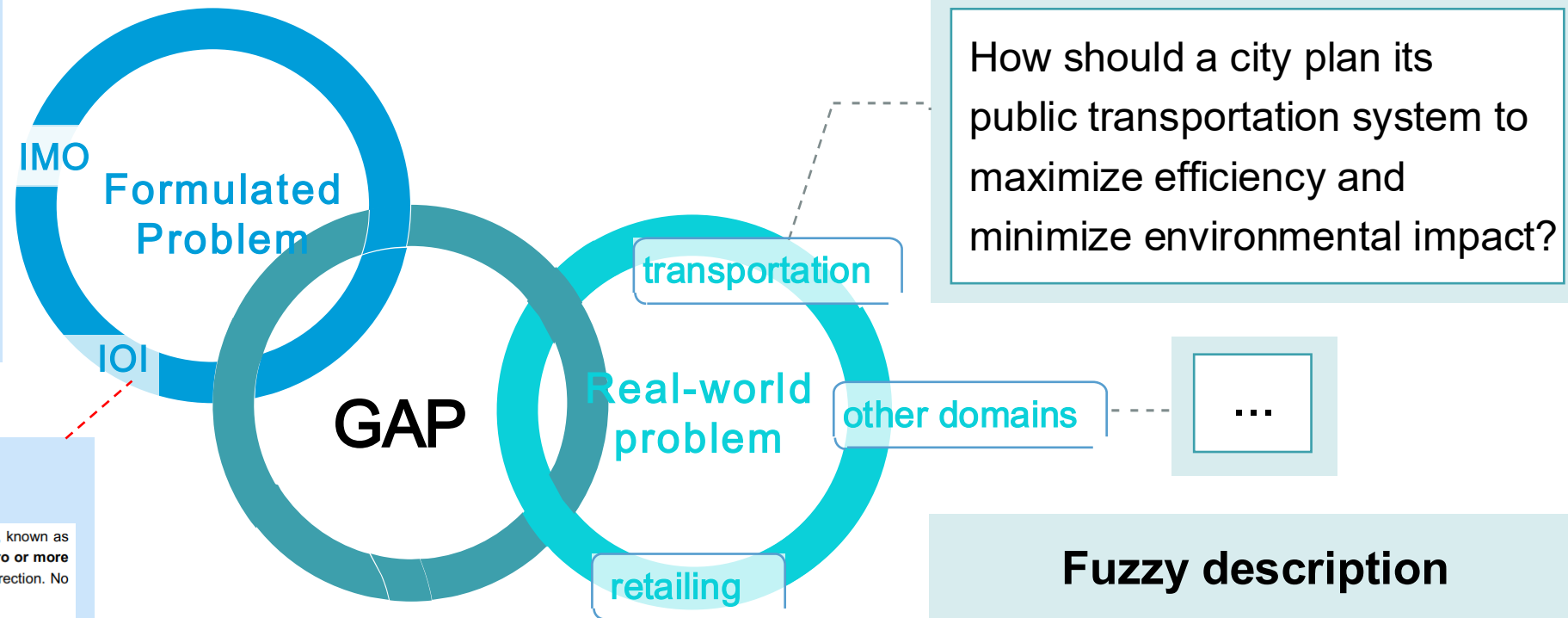
A path from tower  $x$  to tower  $y$  is a sequence of one or more towers such that:

- the first element of the sequence is  $x$ ,
- the last element of the sequence is  $y$ ,
- all elements of the sequence are **distinct**, and
- each two consecutive elements (towers) in the sequence are connected by a bridge.

Note that by definition there is exactly one path from a tower to itself and the number of different paths from tower  $i$  to tower  $j$  is the same as the number of different paths from tower  $j$  to tower  $i$ .

The lead architect in charge of the design wishes for the bridges to be built such that for all  $0 \leq i, j \leq n - 1$  there are exactly  $p[i][j]$  different paths from tower  $i$  to tower  $j$ , where  $0 \leq p[i][j] \leq 3$ .

Construct a set of bridges that satisfy the architect's requirements, or determine that it is impossible.



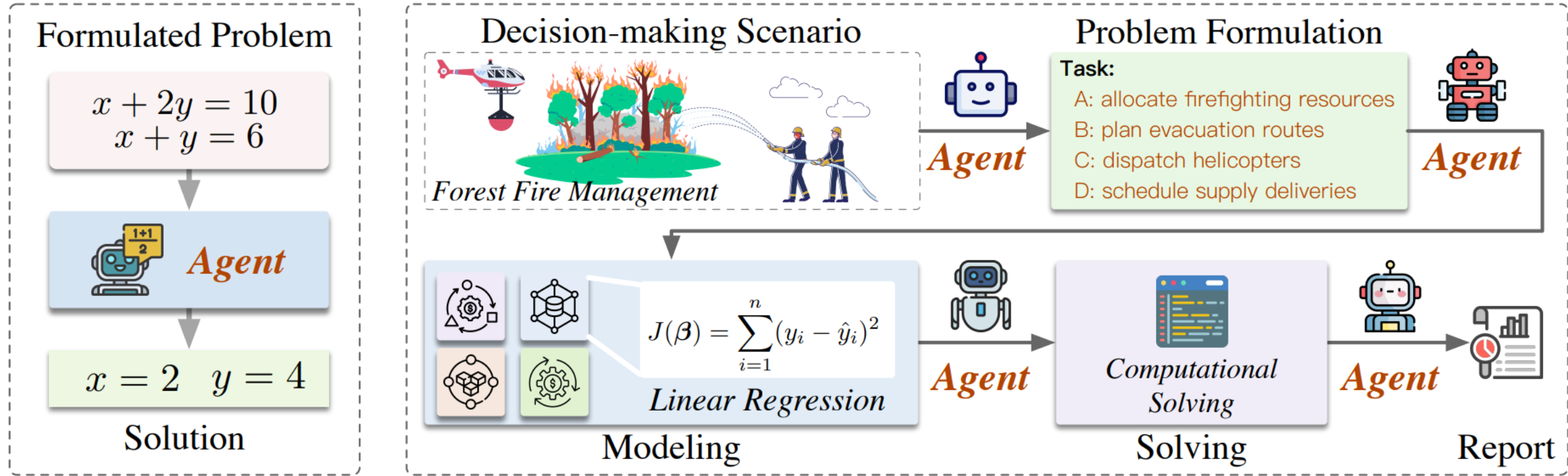
## Open-ended

How should a city plan its public transportation system to maximize efficiency and minimize environmental impact?

## Fuzzy description

How should a retail store optimize its staff scheduling to minimize labor costs while ensuring customer satisfaction?

# Our Work: MM-Agent



Solving well-defined problem in textbook

Solving open-ended Mathematical Modeling problem in the real-world



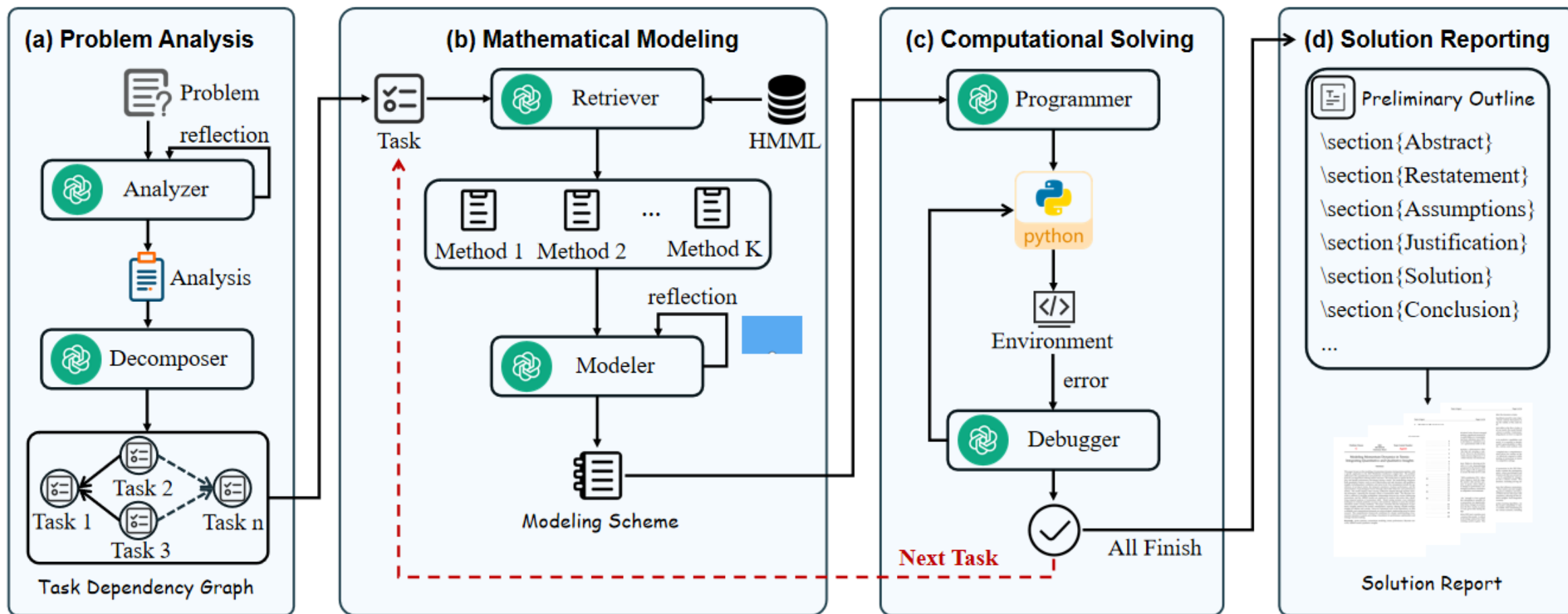
**MM-Agent assisted two undergraduate teams in **winning the Finalist Award** (top 2.0% among 27,456 teams) in MCM/ICM 2025.**





# Architecture of MM-Agent

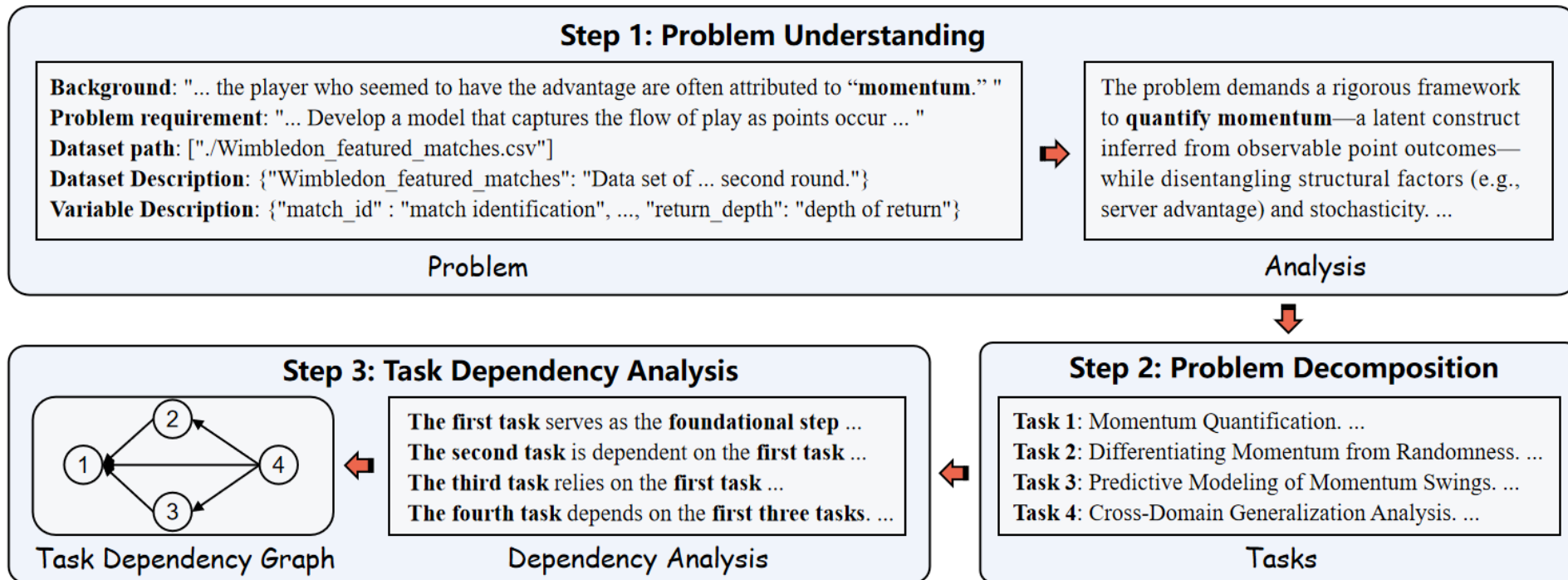
- MM-Agent automates the problem formulation and structured mathematical modelling process.



MM-Agent: LLM as Agents for Real-world Mathematical Modeling Problem. NeurIPS 2025.

# Problem Analysis

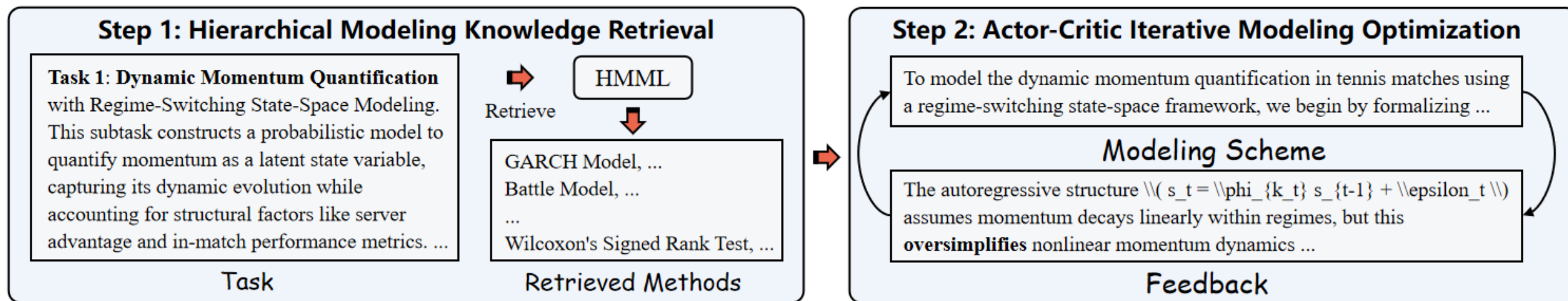
- ❑ Automated problem understanding, problem decomposition, and task dependency analysis
  - 1. Extract background, goals, and key information. 2. Identify subtasks based on components. 3. Map logical/variable dependencies between subtasks .



# Mathematical Modeling

## □ Hierarchical Mathematical Modeling Optimization

- Mimic human cognition: retrieve modeling methods from a knowledge base and solve each subtask via actor-critic optimization.



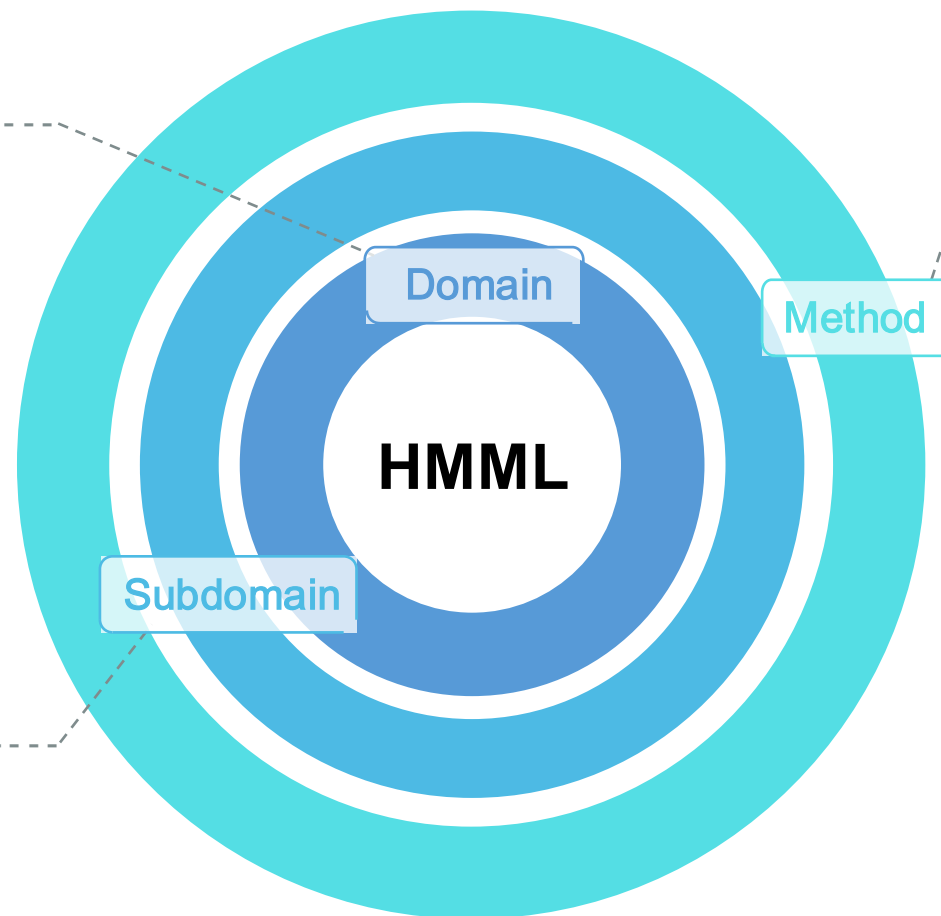


# Hierarchical Mathematical Modeling Library

- ❑ To enhance the mathematical modeling capabilities of LLM agents, we constructed a structured three-level hierarchy designed to improve retrieval effectiveness.

**Domains of mathematical modeling methods**, including Optimization, Machine Learning, Operations Research, Prediction, and Evaluation.

**Subdomains of mathematical modeling methods**, such as Operations Research, which includes Graph Theory, Programming Theory, Stochastic Programming, and so on.



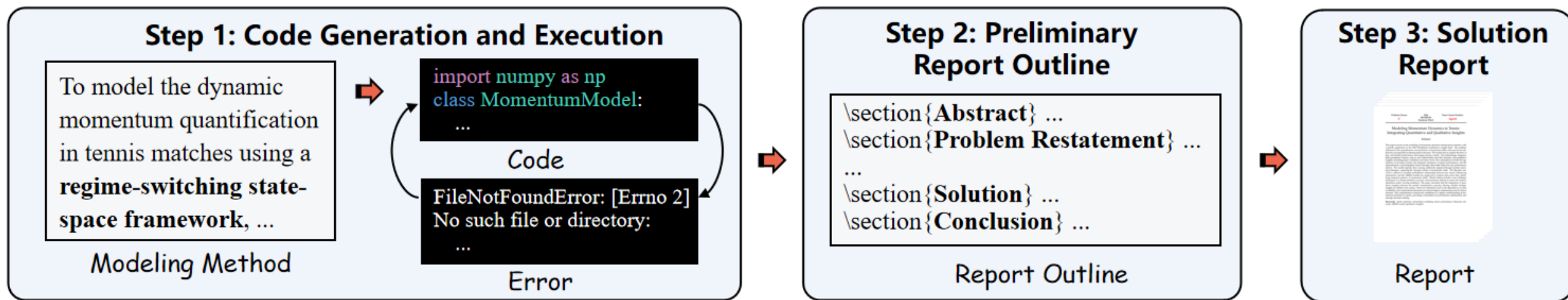
**Specific mathematical modeling methods**, such as Linear Programming, which include:

<Method>: Linear Programming  
<core\_idea>: ... linear objective functions and constraints  
<application>: ... resource allocation, production planning, ...

# Computational Solving and Solution Reporting

## ❑ Code execution based problem solving and report generation

- Use MLE-Solver to conduct experiments and auto-generate a structured report with results and findings



# How to Evaluate Open-ended, Ill-Defined Tasks?

- ❑ **Benchmark Design:** Mathematical modeling competitions offer a natural, open-ended setting for evaluating LLM agents, where problems mirror real-world complexity and require abstraction, reasoning, and solution design.

- ❑ **Evaluation:**



## Analysis Evaluation

- Is the problem clearly defined?
- Are key factors and logic between sub-tasks well-structured?

## Modeling Rigorousness

- Are assumptions reasonable and stated explicitly?
- Are methods and models accurate, rational, and scientific?

## Practicality & Scientificity

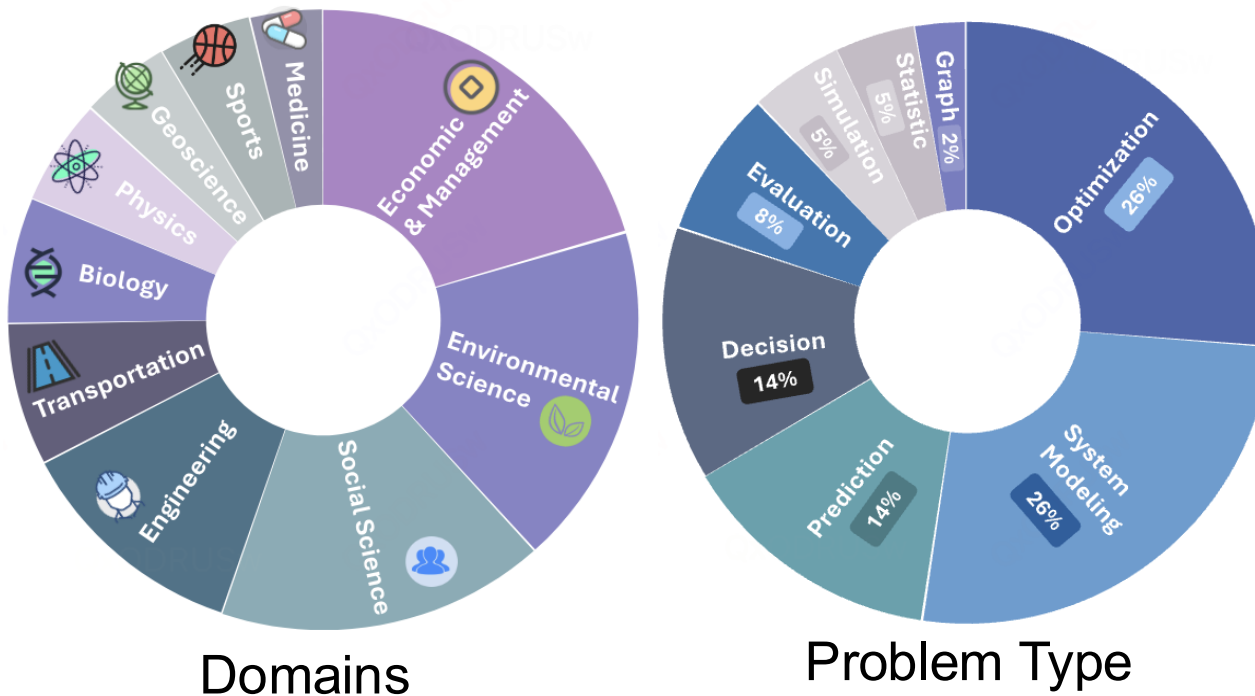
- Is the model applicable in real-world scenarios?
- Does it yield useful, scientifically valid insights?

## Result & Bias Analysis

- Are the results interpretable and in-depth?
- Is bias addressed and robustness ensured

# MM-Bench

- ❑ **Data Summary:** MM-Bench consists of 10 domains, 8 task types, and a total of 111 problem samples, all sourced (MCM and ICM).



## Problem example:

**Background:** “In the 2023 Wimbledon Gentlemen’s final,... The incredible swings, sometimes for many points or even games, that occurred in the player who seemed to have the advantage are often attributed to “momentum.”

**Dataset path:** / Wimbledon\_featured\_matches’: ‘Data set of ... second round.

**Variable Description:** “match\_id”: ‘match identification’,

Problem Example: 2024 C

# Experiments Results

- ❑ **MM-Agent achieves SOTA performance** across all metrics and backbone models.
- ❑ Results are **robust across years**, showing no sign of data leakage.

Table 1: Experimental results on the 2021–2024 and 2025 mathematical modeling competitions. AE, MR, PS, and RBA denote *Analysis Evaluation*, *Modeling Rigorousness*, *Practicality and Scientificity*, and *Result and Bias Analysis*, respectively.

Methods	2021–2024					2025				
	AE ↑	MR ↑	PS ↑	RBA ↑	Overall ↑	AE ↑	MR ↑	PS ↑	RBA ↑	Overall ↑
Human										
Human Team	9.04	6.20	8.79	7.62	7.91	9.25	7.42	8.92	6.50	8.02
GPT-4o										
GPT-4o	7.62	3.86	8.48	5.17	6.28	7.67	3.67	8.90	5.75	6.50
DS-Agent	8.18	7.08	8.72	7.47	7.86	8.25	<b>7.33</b>	8.92	7.10	7.90
ResearchAgent	7.97	6.80	8.82	7.37	7.74	8.00	7.30	8.60	7.00	7.73
Agent Laboratory	8.56	6.35	8.63	5.56	7.28	8.75	5.58	8.58	5.33	7.13
<b>MM-Agent</b>	<b>9.15</b>	<b>7.28</b>	<b>9.00</b>	<b>8.44</b>	<b>8.85</b>	<b>8.86</b>	7.21	<b>9.00</b>	<b>8.43</b>	<b>8.38</b>
DeepSeek-R1-671B										
DeepSeek-R1	7.23	4.79	8.69	4.50	6.30	7.42	4.25	8.50	5.25	6.35
DS-Agent	8.25	6.88	8.74	7.19	7.77	7.92	6.33	9.00	7.60	7.71
ResearchAgent	8.13	7.04	8.77	6.92	7.72	8.00	6.75	8.83	7.58	7.79
Agent Laboratory	8.65	5.96	8.70	5.91	7.31	8.83	5.50	8.83	5.58	7.19
<b>MM-Agent</b>	<b>9.54</b>	<b>8.25</b>	<b>9.06</b>	<b>8.54</b>	<b>8.85</b>	<b>9.50</b>	<b>8.33</b>	<b>9.25</b>	<b>8.58</b>	<b>8.92</b>

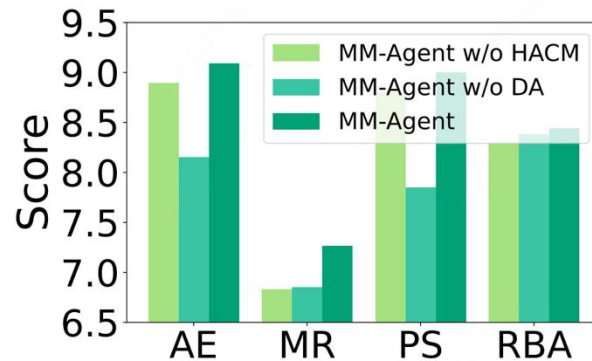


# Experiments Results

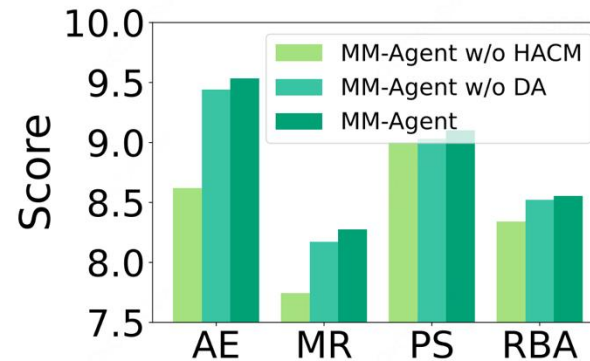
□ **Ablation Study:** We evaluate MM-Agent by removing key components:

- *w/o DA* (no task dependency analysis)
- *w/o HACM* (no hierarchical actor-critic modeling)

□ **Cost Efficiency:** MM-Agent achieves **competitive performance** with **low computational cost**.



(a) GPT-4o



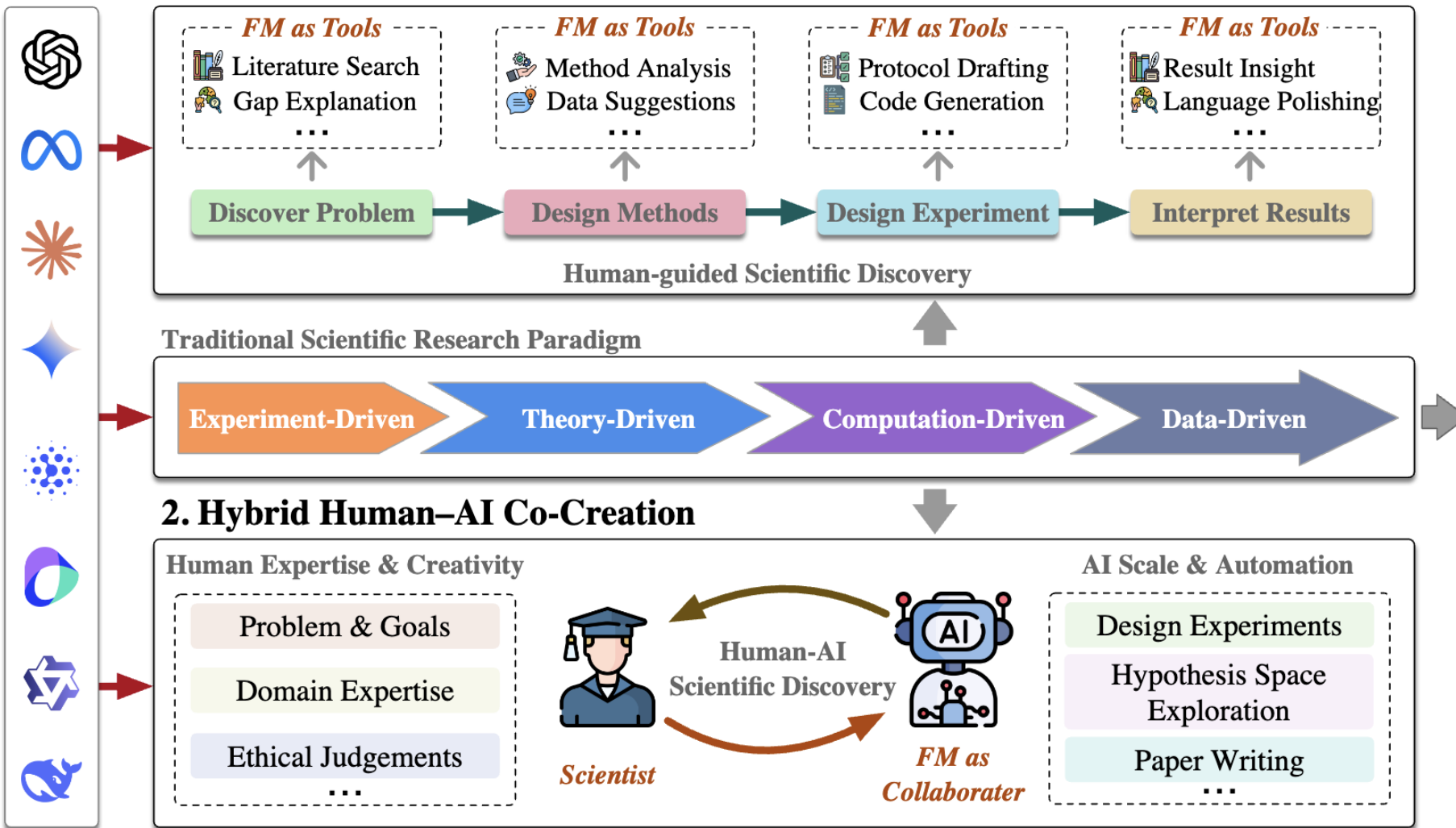
(b) DeepSeek-R1-671b

Table 2: Experimental results on average token consumption, cost, and runtime.

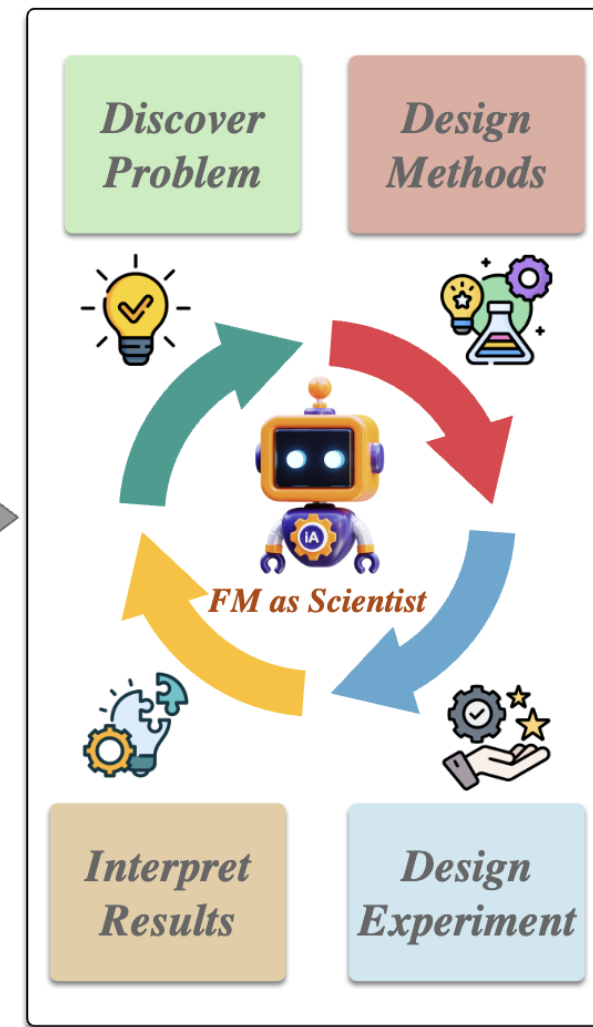
Methods	Token	Cost(\$)	Runtime(s)
GPT-4o			
DS-Agent	198,186	0.77	1,044
ResearchAgent	170,732	0.67	459
Agent Laboratory	746,159	2.14	1,015
MM-Agent	240,877	0.88	906
DeepSeek-R1-671b			
DS-Agent	315,171	0.42	6,444
ResearchAgent	222,030	0.28	4,816
Agent Laboratory	974,313	0.88	10,595
MM-Agent	530,363	0.56	7,529

# Future Directions

FMs



3. Autonomous Scientific Discovery

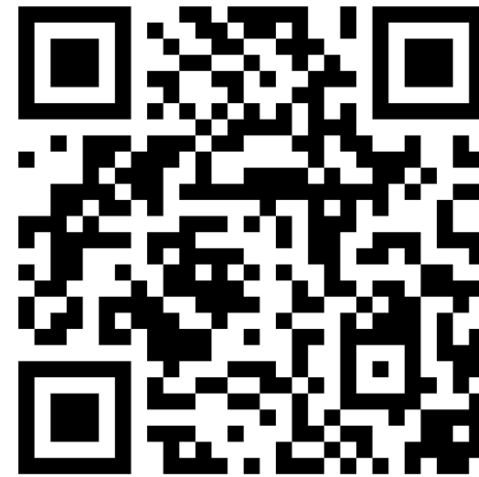
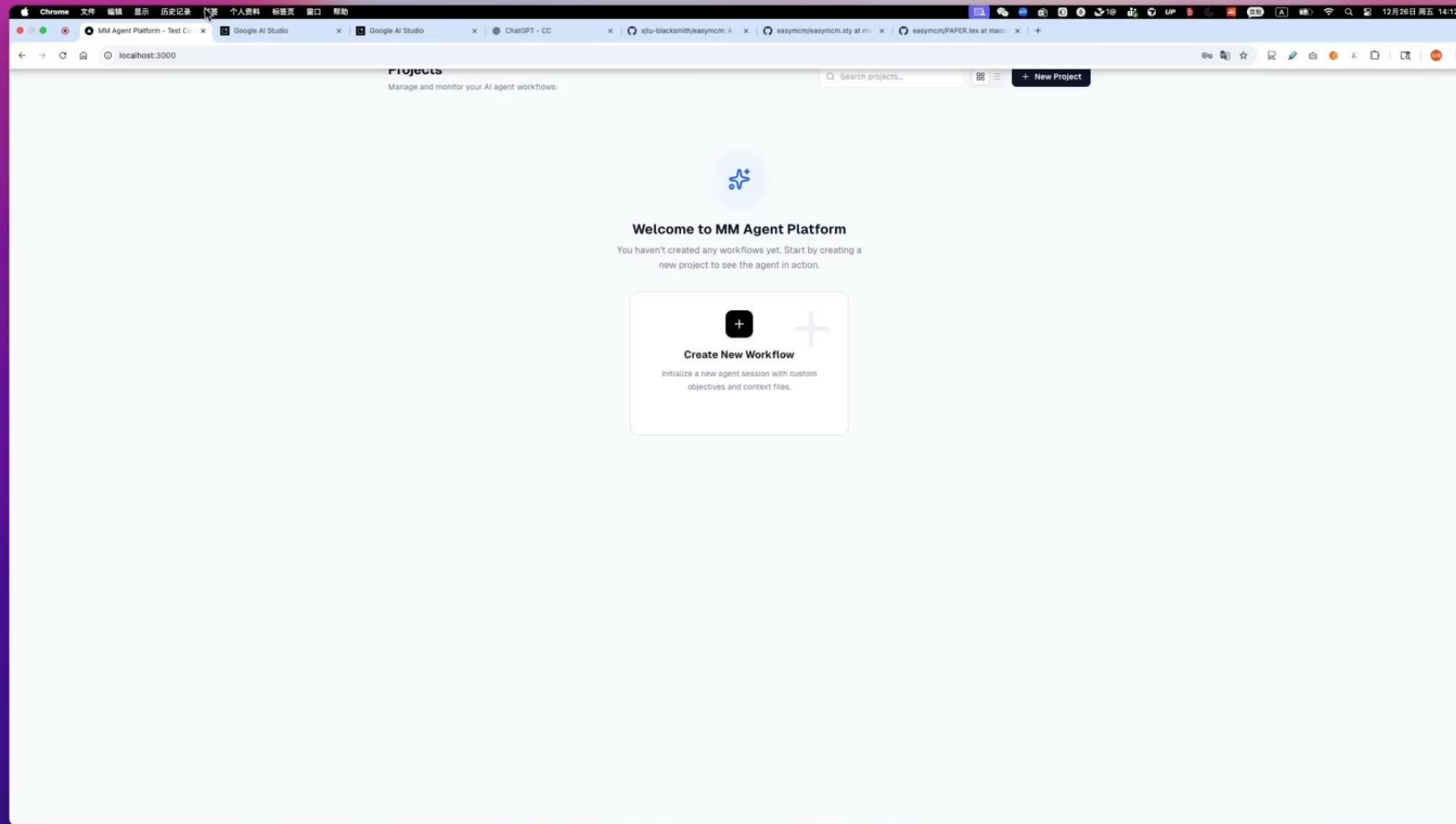


Foundation Models for Scientific Discovery: From Paradigm Enhancement to Paradigm Transition. NeurIPS 2025

# 数学建模智能体实战



# Demo



项目链接



群聊：美赛数据建模 agent

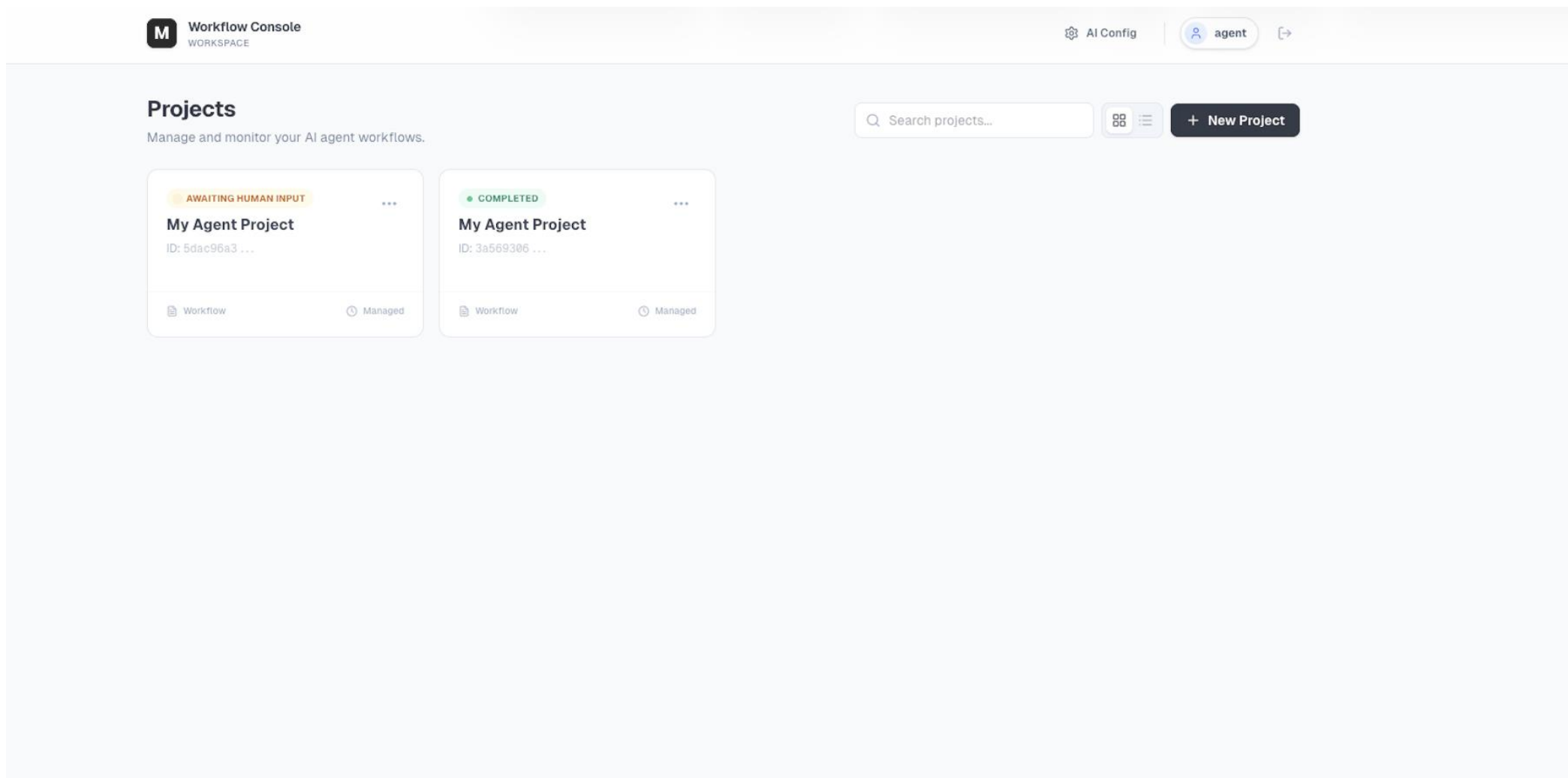


该二维码7天内(1月2日前)有效，重新进入将更新

项目agent服务群



# 管理界面：创建自己的数学建模项目。





# 项目建立：上传问题与数据集。

Workflow

My Agent Project > Problem analysis Awaiting input

CONTENT REVIEW

PHASE 1: ANALYSIS AND PLANNING

Problem analysis

Subtask planning

PHASE 2: ITERATIVE MODELING

Iteration 0

Model Selection

Model Design

Code Implementation

Visualization Generation

Result Analysis

Iteration 1

Model Selection

Model Design

Code Implementation

Visualization Generation

Result Analysis

PHASE 3: PAPER GENERATION

Paper Architect

Paper Assembler & Compiler

PROBLEM RESTATEMENT

The objective is to address the multifaceted nature of "momentum" in professional tennis, specifically within the context of the 2023 Wimbledon Gentlemen's final between Carlos Alcaraz and Novak Djokovic (match\_id: 2023-wimbledon-1701), and broader tournament data.

I. Core Modeling Objectives:

- Flow of Play Quantization: Develop a mathematical model that captures the dynamic flow of play on a point-by-point basis. The model must:
  - Identify which player is performing better at any discrete point in time.
  - Quantify the magnitude of the performance advantage.
  - Incorporate the inherent advantage of the server (higher win probability) into the calculation.
  - Provide high-quality visualizations depicting the match flow across all sets and games.
- Momentum vs. Randomness (The "Coach's Skepticism"): Evaluate the hypothesis that swings in play and "runs" of success are merely stochastic events (random walks) rather than the result of physical or psychological "momentum." Use statistical testing or model metrics to support or refute the claim that momentum is a measurable phenomenon.
- Swing Prediction and Factor Analysis:
  - Develop a model to predict shifts in match flow (momentum swings) using data from at least one match.
  - Identify the specific factors (e.g., aces, unforced errors, break points, rally length, physical exertion) most strongly correlated with these shifts.
  - Formulate tactical advice for players on how to prepare for and respond to momentum swings when facing different opponents.
- Validation and Generalization:
  - Test the models on other matches within the dataset to assess predictive accuracy.
  - Perform error analysis to identify why the model might fail in specific instances.
  - Evaluate the potential to generalize the model to Women's matches, other tournaments, different court surfaces (clay, hard court), and other sports like table tennis.

II. Requirements and Constraints:

- Use the provided `Wimbledon_featured_matches.csv` data.
- The solution must include a Summary Sheet, Table of Contents, complete solution, and a 1-2 page memo for coaches.
- All models must be rigorously justified and validated.
- Adherence to the 25-page limit (excluding the AI Use Report).
- Transparency regarding the use of AI tools as per COMAP policy.

III. Specific Match Context (Alcaraz vs. Djokovic):

- Set 1: Djokovic dominant (6-1).
- Set 2: Tense battle, Alcaraz wins in tie-breaker (7-6).
- Set 3: Alcaraz dominant (6-1).
- Set 4: Djokovic regains control (6-3).
- Set 5: Alcaraz secures victory (6-4).
- The model must reflect these "incredible swings" and identify the turning points where the "strength or force gained by motion" shifted.

DATA DESCRIPTION

The dataset consists of every point from all Wimbledon 2023 men's matches starting from the third round.


v2

Interactive Task Block

Save Changes

Approve

Reject



香港科技大学 (广州)  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
TECHNOLOGY (GUANGZHOU)

19

# 数学建模助力

## 建模流程

Workflow

PHASE 1: ANALYSIS AND PLANNING

Problem analysis

Subtask planning

PHASE 2: ITERATIVE MODELING

Model Selection

Model Design

Code Implementation

Visualization Generation

PHASE 3: PAPER GENERATION

Paper Architect

Paper Assembler & Compiler

THESIS / REPORT

REVIEW STATUS

COMPILE LOGS

<> PYTHON SOURCE

1 \documentclass[12pt]{article}

2

3 % --- Package Loading ---

4 % The control number must be passed as an option to easynorm

5 \usepackage[240000]{easynorm}

6

7 % --- Preamble: Metadata ---

8 \problem(A)

9 \title{Dynamic Modeling and Robustness Analysis of [System Name]}

10

11 % --- Document Start ---

12 \begin{document}

13

14 % --- Summary Sheet (Abstract) ---

15 % In easynorm, the abstract content is automatically parsed

16 % into the Summary Sheet body.

17 \begin{abstract}

18 \noindent

19 % O-PRIZE REQUIREMENT: High-level summary of methodology and results.

20 In this paper, we develop a multi-stage [Model Type] to analyze [Problem].

21 We derive features including [Feature 1] and [Feature 2] to capture non-linear dynamics.

22 Our results demonstrate an  $SR^2$  of 0.94, with sensitivity analysis confirming

23 stability under  $\delta p_m$  and  $\delta \lambda$  parameter fluctuations.

24 \end{abstract}

25

26 \maketitle % Generates the Summary Sheet; MUST follow the abstract environment

27

28 % --- Administrative Pages ---

29 \newpage

30 \tableofcontents

31 \newpage

32

33 % --- Main Body ---

34

35 % 1. Introduction

36 \section{Introduction}

37 \input{sections/a01\_1\_introduction.tex}

38

39 % 2. Notation and Assumptions

40 \section{Assumptions and Notations}

41 % RIGOR: Define every variable in a Notation Table here.

42 \input{sections/a02\_2\_assumptions\_and\_n.tex}

43

44 % 3. Data Processing and Feature Engineering

45 \section{Data Preprocessing and Feature Engineering}

46 % Derived Features: Temporal Momentum, Z-scores, and Entropy metrics.

47 \input{sections/a03\_3\_data\_preprocess.tex}

48

49 % 4. Mathematical Model

50 \section{The Mathematical Model}

51 % Rigorous Formulation using Differential Equations or Stochastic Processes.

52 \input{sections/a04\_4\_the\_mathematical.tex}

53

54 % 5. Statistical Validation

55 \section{Results and Statistical Analysis}

56 % Report MSE, RMSE, and F-test results here.

57 \input{sections/a05\_5\_results\_and\_stat.tex}

58

59 % 6. Robustness

60 \section{Sensitivity Analysis}

Console

EXECUTION LOG

Light Mode

./usr/local/texlive/2024/texmf-dist/tex/latex/hyperref/pdftoc.def

./usr/local/texlive/2024/texmf-dist/tex/generic/intcalc/intcalc.sty

./usr/local/texlive/2024/texmf-dist/tex/latex/hyperref/puenc.def

./usr/local/texlive/2024/texmf-dist/tex/latex/utl/utl.sty

./usr/local/texlive/2024/texmf-dist/tex/generic/bitset/bitset.sty

./usr/local/texlive/2024/texmf-dist/tex/generic/bigintcalc/bigintcalc.sty

./usr/local/texlive/2024/texmf-dist/tex/latex/base/atbeglnl-ltx.sty

./usr/local/texlive/2024/texmf-dist/tex/latex/hyperref/hdftex.def

./usr/local/texlive/2024/texmf-dist/tex/latex/base/atveryend-ltx.sty

./usr/local/texlive/2024/texmf-dist/tex/generic/uniquecounter/uniquecounter.sty

./usr/local/texlive/2024/texmf-dist/tex/latex/psnfs/otlpd.fd

./usr/local/texlive/2024/texmf-dist/tex/latex/l3backend/l3backend-pdf.tex.def

./main.aux

geometry: driver: auto-detecting

geometry: detected driver: pdfTeX

./usr/local/texlive/2024/texmf-dist/tex/context/base/m11/supp-pdf.m11

[Loading MPS to PDF converter (version 2006.09.02).]

./usr/local/texlive/2024/texmf-dist/tex/latex/epstopdf/epstopdf-base.sty

./usr/local/texlive/2024/texmf-dist/tex/latex/grfext/grfext.sty

./usr/local/texlive/2024/texmf-dist/tex/latex/latexconfig/epstopdf-sys.cfg

./usr/local/texlive/2024/texmf-dist/tex/latex/graphics/color.sty

./usr/local/texlive/2024/texmf-dist/tex/latex/graphics/cfg/color.cfg

./usr/local/texlive/2024/texmf-dist/tex/latex/graphics/mathcolor.ltx

./usr/local/texlive/2024/texmf-dist/tex/latex/psnfs/omzplm.fd

./usr/local/texlive/2024/texmf-dist/tex/latex/psnfs/omzplm.fd

./usr/local/texlive/2024/texmf-dist/tex/latex/psnfs/omzplm.fd

./usr/local/texlive/2024/texmf-dist/tex/latex/psnfs/otlpd.fd

./usr/local/texlive/2024/texmf-var/fonts/map/pdftex/updmap/pdftex.map

./usr/local/texlive/2024/texmf-dist/fonts/enc/dvips/base/8t.enc

No file main.toc.

LaTeX Warning: Reference 'LastPage' on page 2 undefined on input line 31.

Package fancyhdr Warning: \headheight is too small (12.0pt):

fancyhdr: Make it at least 13.59999pt, for example:

fancyhdr: \setlength{\headheight}{13.59999pt}.

fancyhdr: You might also make \topmargin smaller to compensate:

fancyhdr: \addtolength{\topmargin}{-1.59999pt}.

[2] ./sections/a01\_1\_introduction.tex

./usr/local/texlive/2024/texmf-dist/tex/psnfs/ufp1mbb.fd

LaTeX Warning: Reference 'LastPage' on page 3 undefined on input line 34.

Package fancyhdr Warning: \headheight is too small (12.0pt):

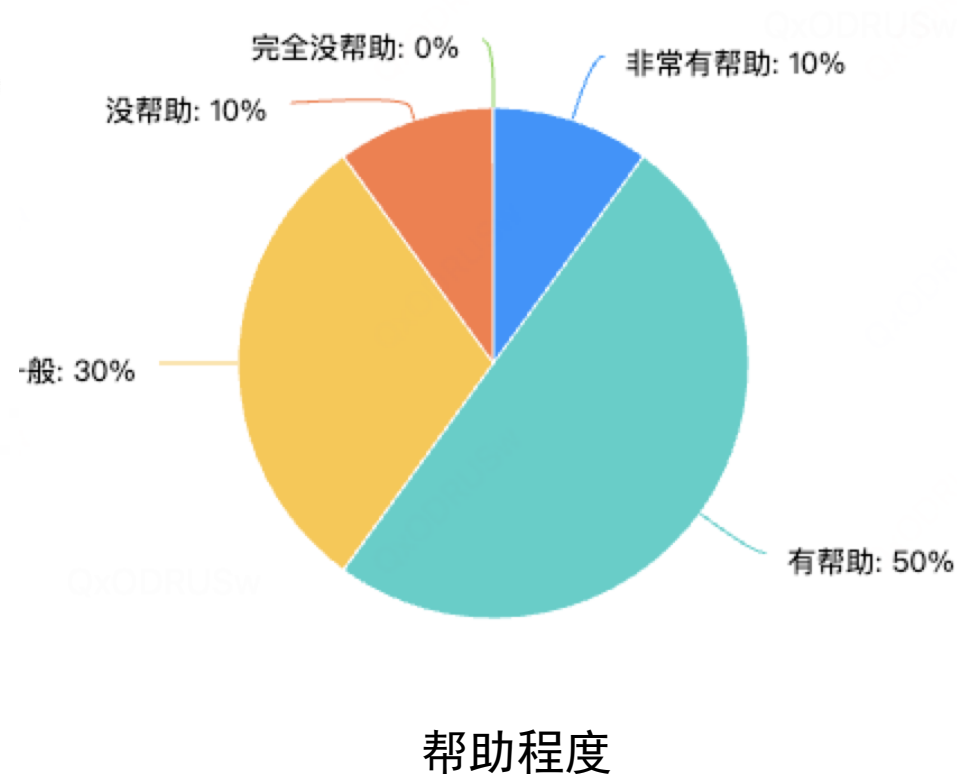
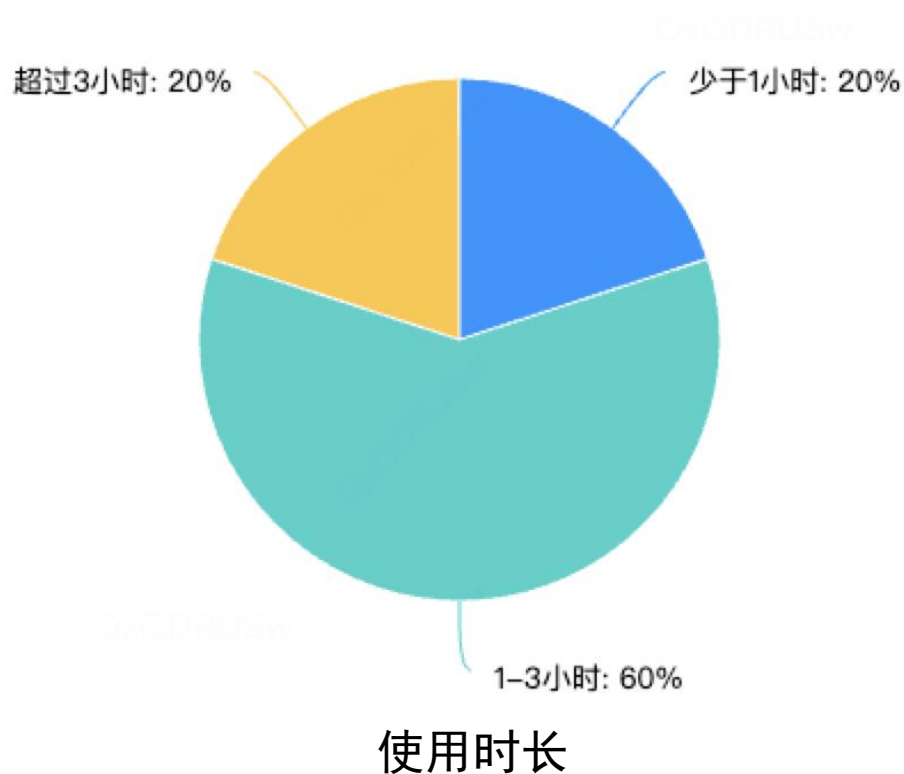
fancyhdr: Make it at least 13.59999pt, for example:

fancyhdr: \setlength{\headheight}{13.59999pt}.

fancyhdr: You might also make \topmargin smaller to compensate:

# 用户调研

- 2025年助力两个美赛队伍取得F奖。
- 调研结果：对于建模思路有很大的帮助！







# 谢谢聆听！

[liuh@ust.hk](mailto:liuh@ust.hk)

<https://raymondhliu.github.io/>



刘浩

北京 海淀



香港科技大学(广州)  
THE HONG KONG  
UNIVERSITY OF SCIENCE AND  
TECHNOLOGY (GUANGZHOU)