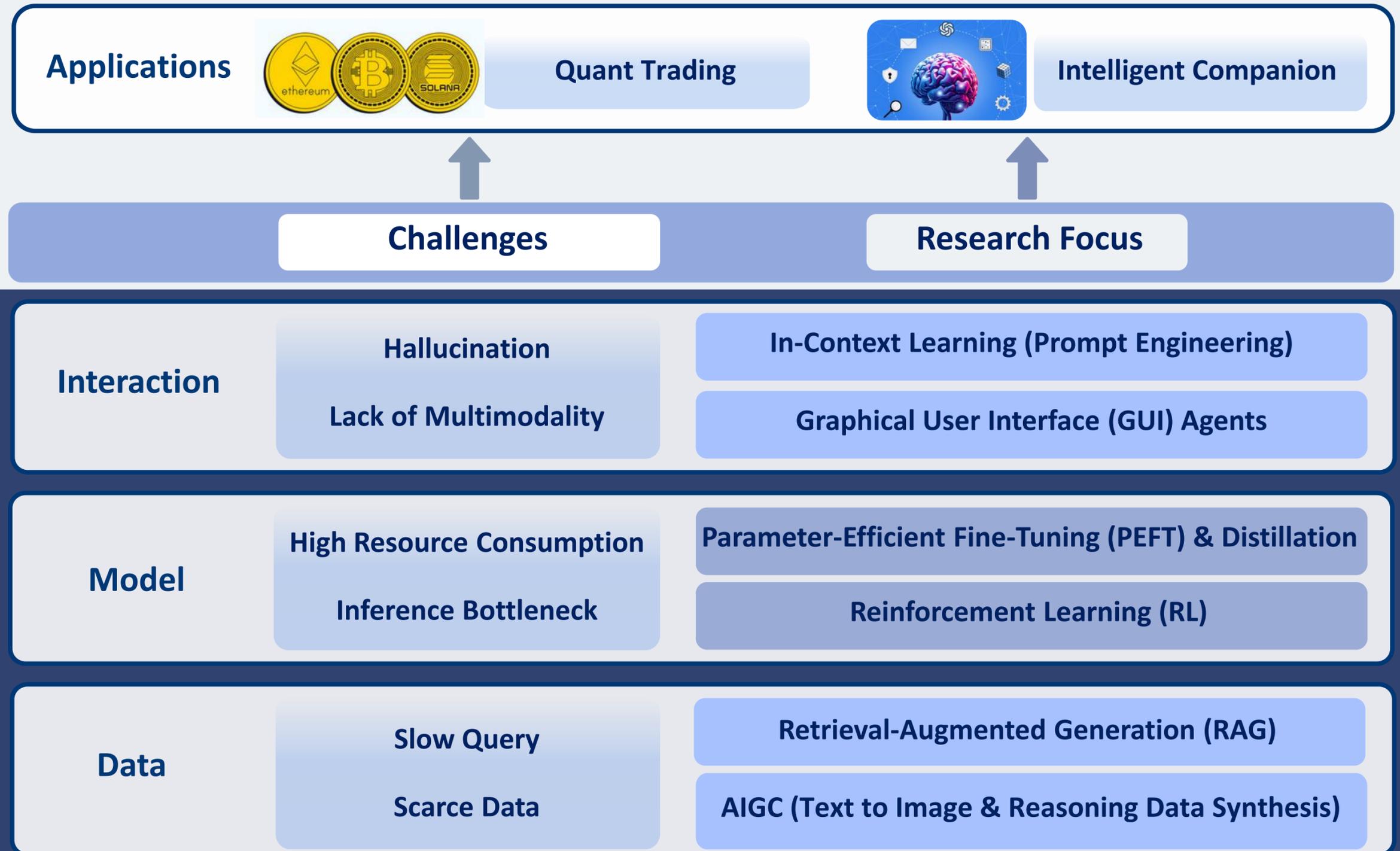
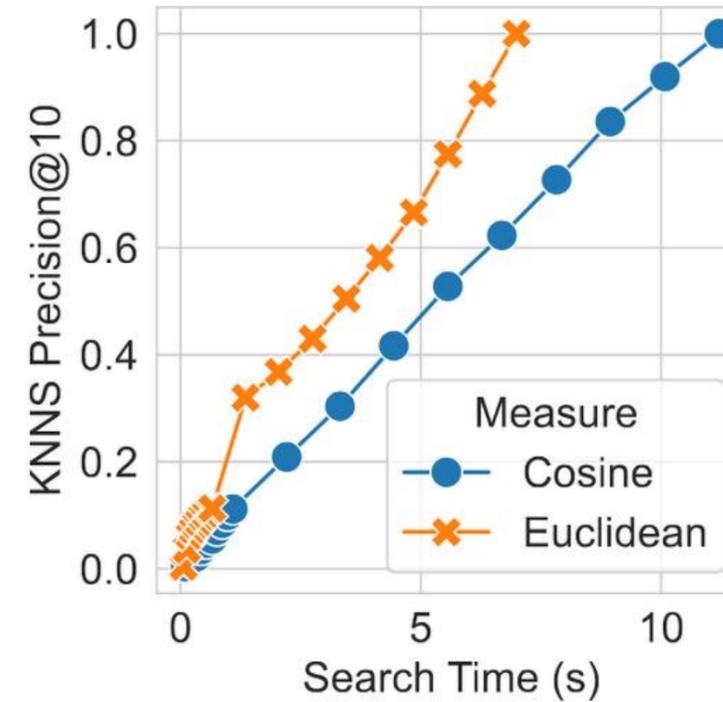
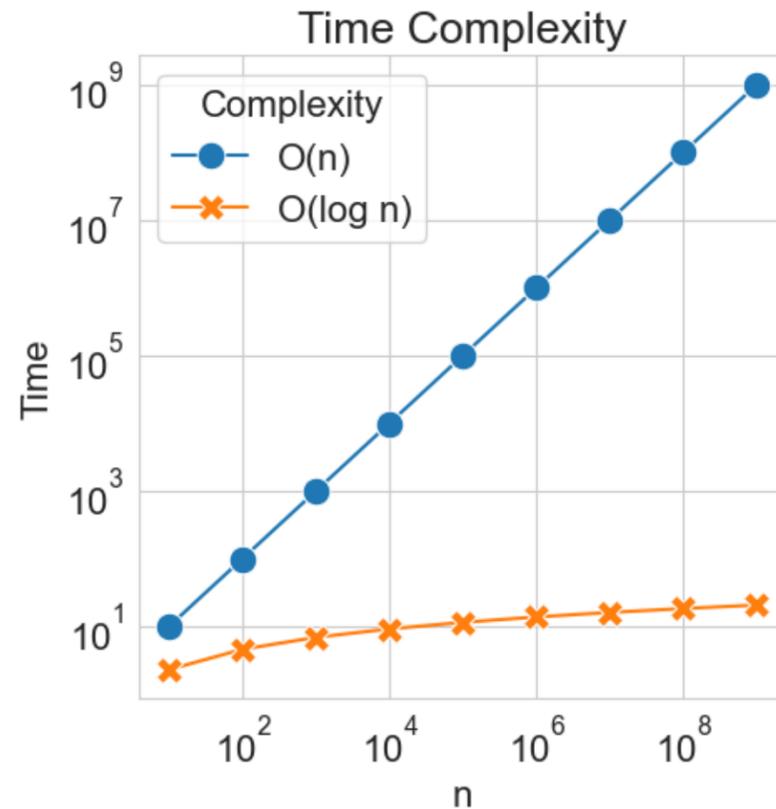


Large Language Model Empowering AI Agents: Tackling Core Challenges and Advancing Real-World Applications

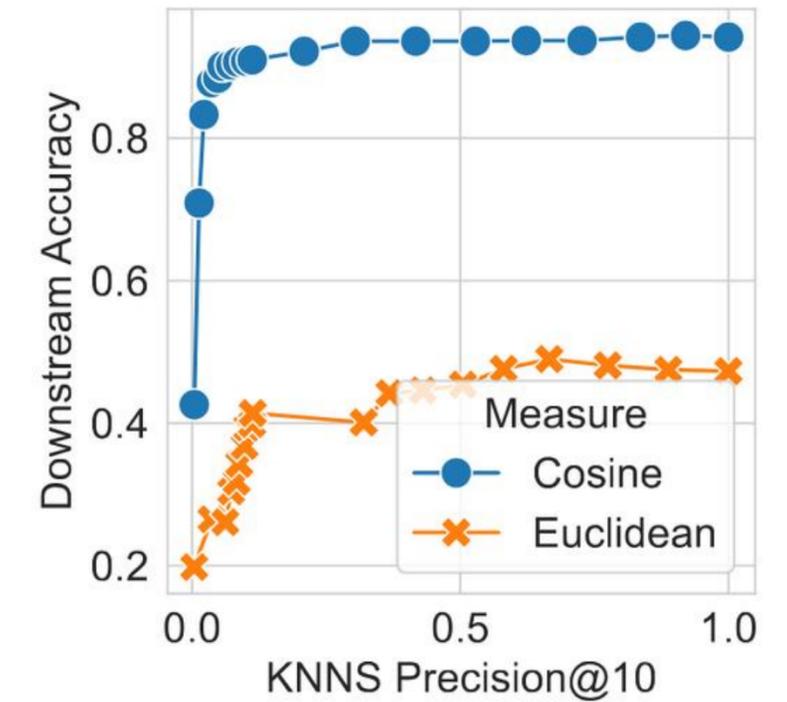
 Jing Tang, Assistant Professor



Approximate Nearest Neighbor Search (ANNS)



(a)



(b)

- Real-world: large scale, low latency
- Naïve NNS: $O(n)$
- SOTA ANNS: $O(\log n)$
- Most research: Euclidean distance
- Common measure in downstream task (e.g., T2I): cosine similarity

Hemi-Sphere Centroids Graph (HSCG)

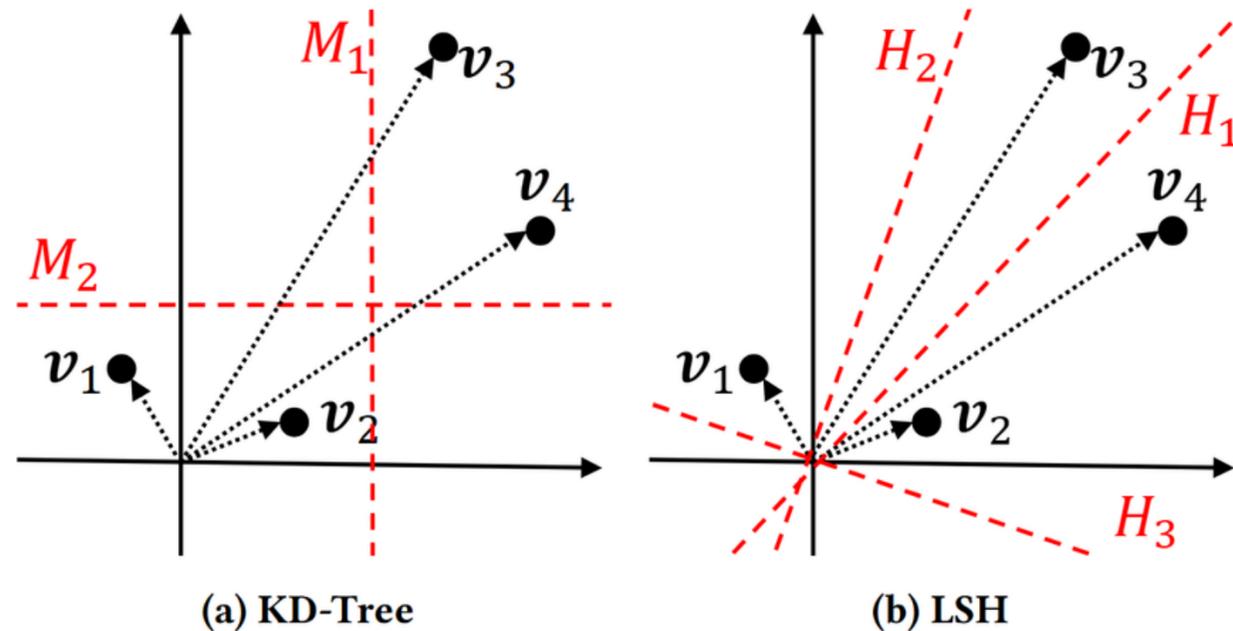


Figure 1: Comparison of KD-Tree and LSH when obtaining neighbors in cosine similarity.

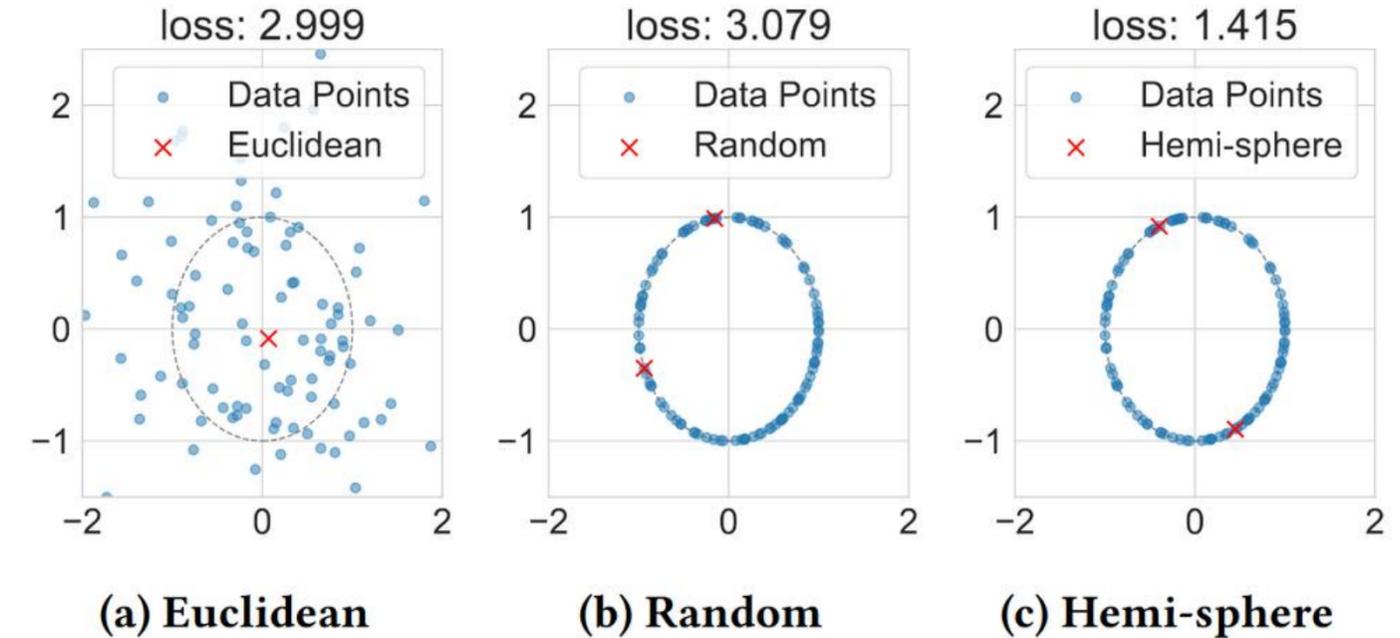
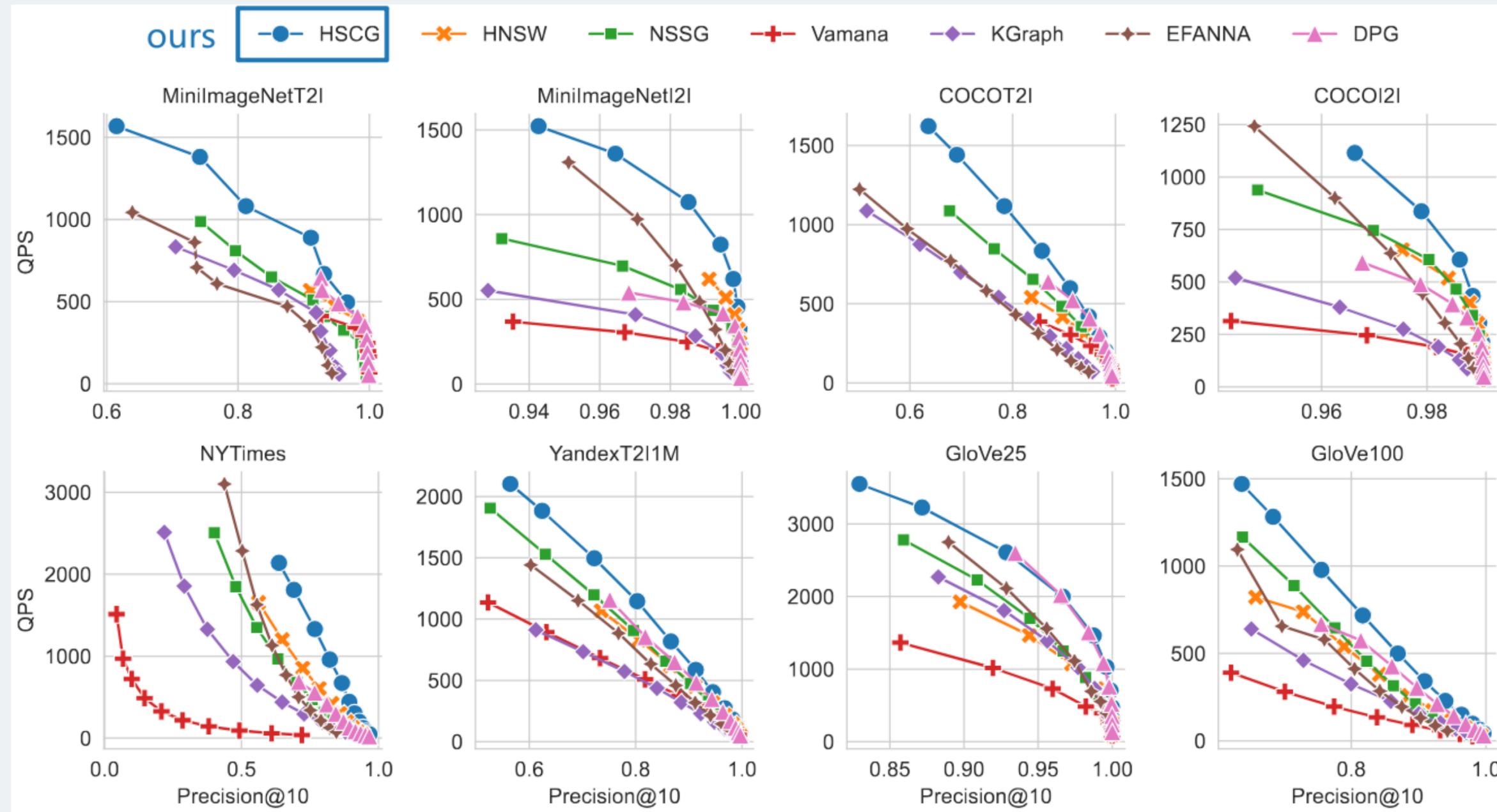


Figure 2: Different centroids and their angular distance loss on a normal distribution dataset.

- Theoretical analysis of Monotonic Relative Neighbor Graph (MRNG) under cosine similarity (convergency of greedy search and bounded out-degree)
- HSCG uses hemi-sphere centroids as the entry points, and locality-sensitive hashing to initialize the graph

HSCG Outperforms Others

Upper right is better



- Diffusion can generate high-quality Video, but expensive.
- A 5 seconds video by Wan 14B takes 7 yuan.



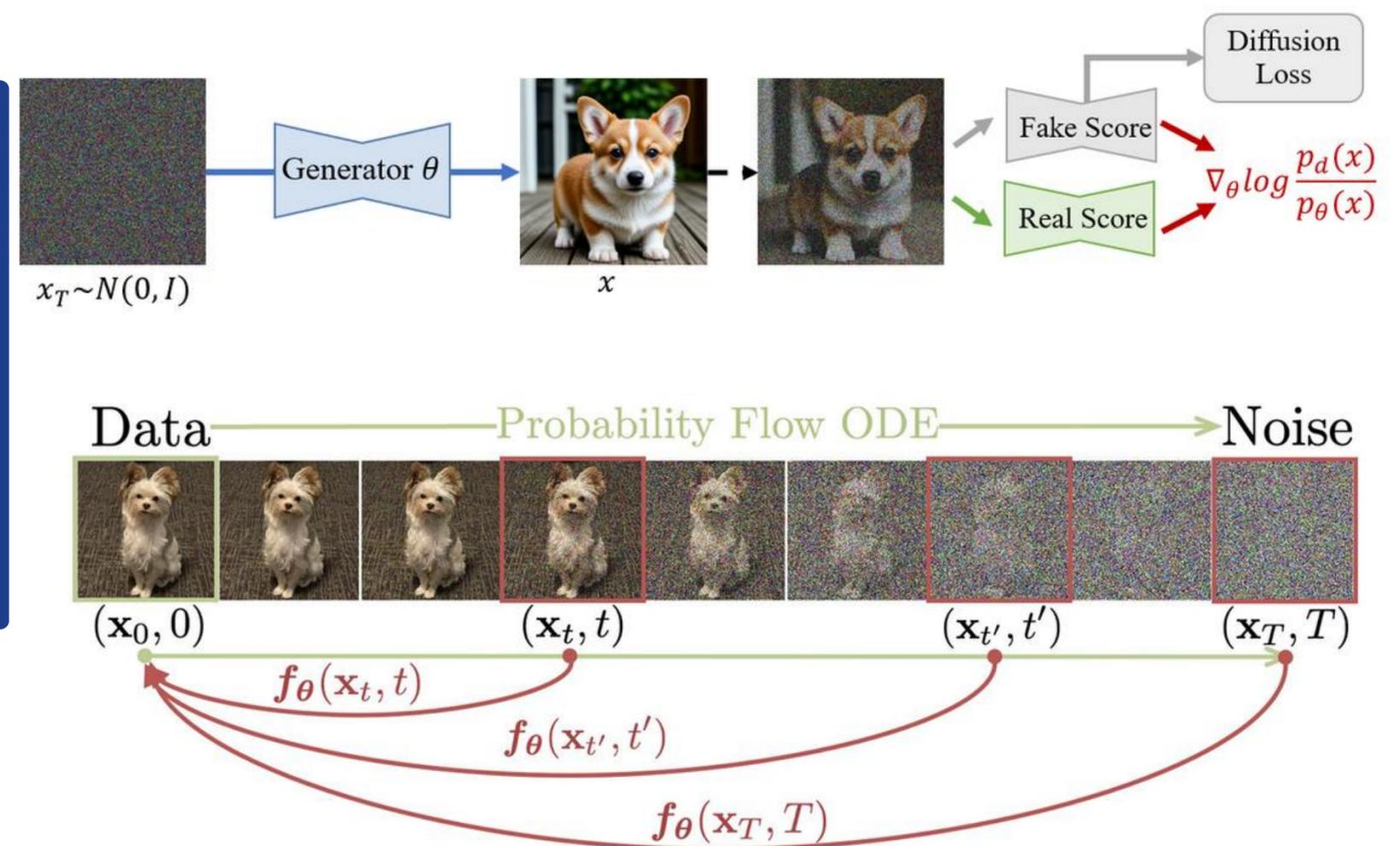
Review on Diffusion Distillation

Distribution Matching:

Powerful one-step synthesis, less efficient in few-step synthesis

Trajectory Distillation:

Effectively preserve the prior knowledge in diffusion models, but worse performance

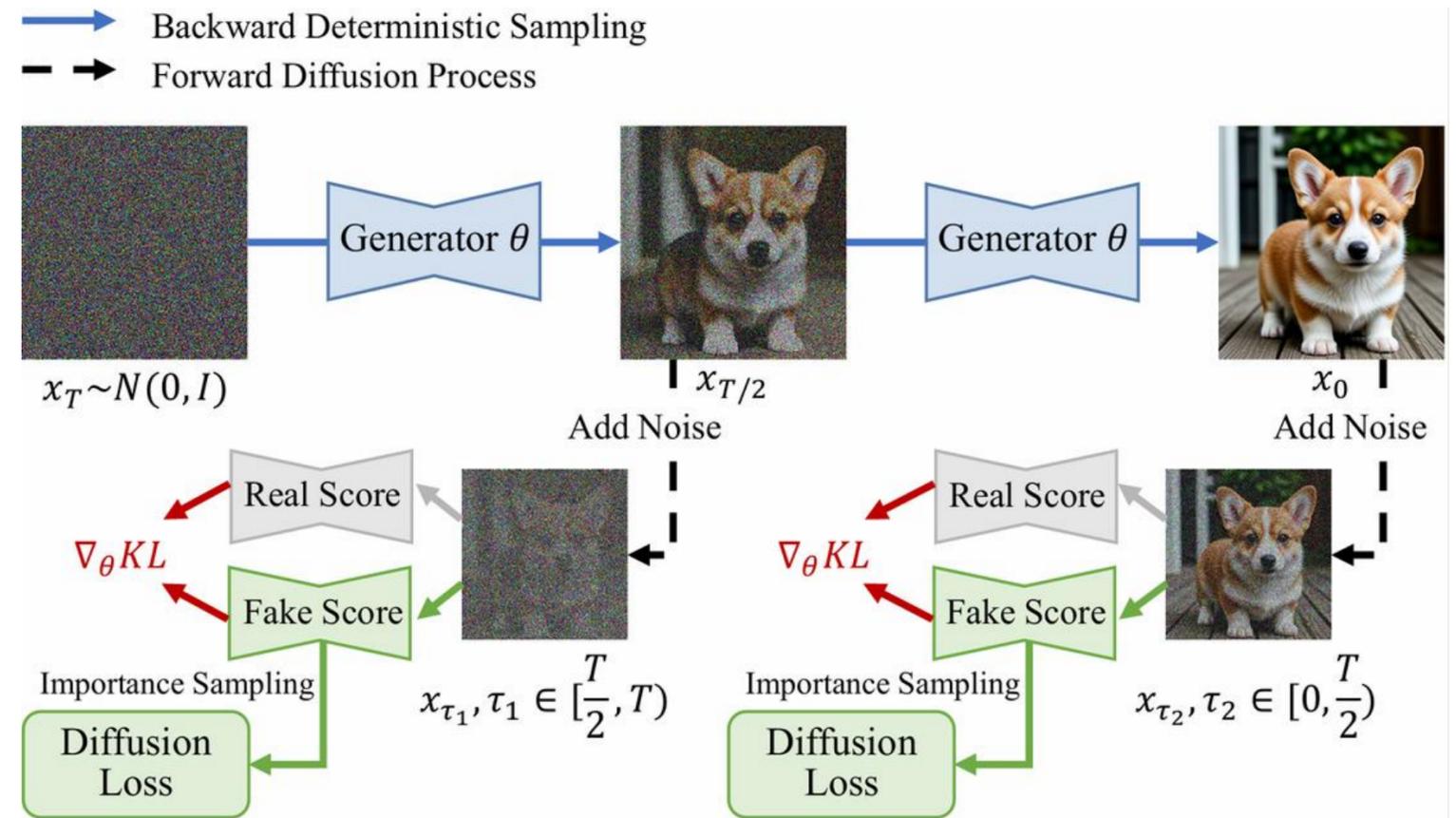


TDM: Trajectory Distribution Matching

$$L(\theta) = \sum_{i=0}^{K-1} \text{KL}(p_{\theta,t_i}(\mathbf{x}_{t_i}) || p_{\phi,t_i}(\mathbf{x}_{t_i}))$$

Student's trajectory distribution Teacher's trajectory distribution

Our TDM unify the distribution matching and trajectory distillation



TDM: Trajectory Distribution Matching



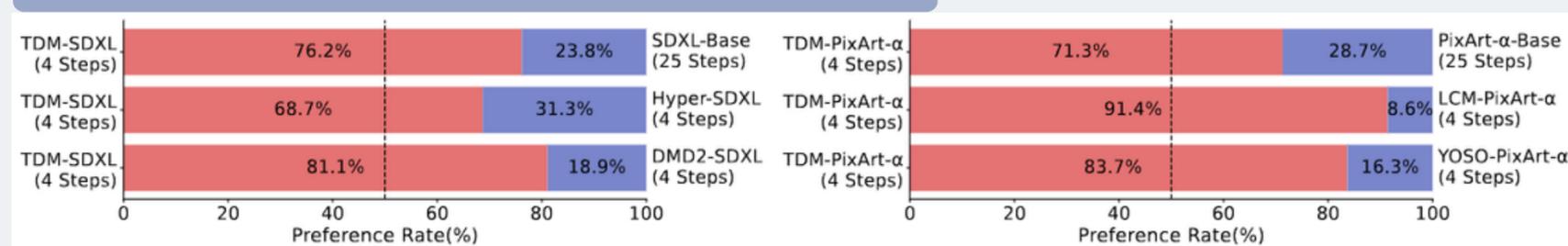
Model	Backbone	HFL	Steps	HPS \uparrow					Aes \uparrow	CS \uparrow
				Animation	Concept-Art	Painting	Photo	Average		
Base Model	SD-v1.5	No	25	26.29	24.85	24.87	26.01	25.50	5.49	33.03
Base Model + Fine-tuning	SD-v1.5	No	25	31.10	29.88	29.53	28.94	29.86	5.85	33.68
InstaFlow [9]	SD-v1.5	No	1	23.17	23.04	22.73	22.97	22.98	5.27	30.04
PeRFlow [35]	SD-v1.5	No	1	12.37	13.50	13.64	11.53	12.76	4.47	15.49
PeRFlow [35]	SD-v1.5	No	2	19.75	19.43	19.41	18.40	19.25	4.91	25.83
Hyper-SD [22]	SD-v1.5	Yes	1	28.65	28.16	28.41	26.90	28.01	5.64	30.87
DMD2 [37]	SD-v1.5	No	1	24.17	22.68	22.97	24.30	23.53	5.82	30.92
TDM-unify (Ours)	SD-v1.5	No	1	29.85	28.90	29.22	27.62	28.90	6.02	32.12
LCM-dreamshaper [14]	SD-v1.5	No	4	26.51	26.40	25.96	24.32	25.80	5.94	31.55
PeRFlow [35]	SD-v1.5	No	4	22.79	22.17	21.28	23.50	22.43	5.35	30.77
TCD [38]	SD-v1.5	No	4	23.14	21.11	21.08	23.62	22.24	5.43	29.07
Hyper-SD [22]	SD-v1.5	Yes	4	31.06	30.01	30.47	28.97	30.24	5.78	31.49
TDM-unify (Ours)	SD-v1.5	No	2	31.43	30.90	30.96	28.78	30.52	6.06	32.67
TDM-unify (Ours)	SD-v1.5	No	4	32.40	31.65	31.35	29.86	31.31	6.08	32.77
Base Model	SDXL	No	25	34.66	33.70	33.43	30.95	33.19	6.17	36.28
TCD [38]	SDXL	No	4	29.65	27.50	27.98	26.13	27.81	5.88	33.42
LCM [14]	SDXL	No	4	30.79	29.38	29.60	27.87	29.41	5.84	34.84
SDXL-Turbo-512 [25]	SDXL	No	4	32.54	31.03	31.04	28.60	30.80	5.81	35.03
SDXL-Lighting [8]	SDXL	No	4	34.20	32.97	33.15	30.52	32.71	6.23	34.62
Hyper-SD [22]	SDXL	Yes	4	35.58	34.54	34.54	31.90	34.14	6.18	34.27
DMD2 [37]	SDXL	No	4	32.87	31.56	31.01	30.39	31.46	5.88	35.51
TDM (Ours)	SDXL	No	4	36.42	35.34	35.51	32.25	34.88	6.28	36.08
Base Model (1024)	PixArt- α	No	25	33.54	32.35	32.00	30.93	32.21	6.23	34.11
YOSO-512 [16]	PixArt- α	No	4	31.40	31.18	31.26	28.15	30.60	6.23	31.83
LCM-1024 [13]	PixArt- α	No	4	31.96	30.60	30.70	28.92	30.55	6.17	33.49
TDM-1024 (Ours)	PixArt- α	No	4	34.61	33.54	33.45	31.23	33.21	6.42	33.66

Ultra-Fast Convergence!

Method	Backbone	NFE \downarrow	HPS \uparrow	Training Cost
DMD2 [42]	SD-v1.5	4	31.53	30+ A800 Days
TDM-unify-GAN	SD-v1.5	4	32.40	4 A800 Days
TDM-unify-SFT	SD-v1.5	4	32.77	3 A800 Days
LCM [18]	SDXL	4	29.41	32 A100 Days
DMD2 [42]	SDXL	4	31.46	160 A100 Days
TDM	SDXL	4	34.88	2 A800 Days
LCM [18]	PixArt- α	4	30.55	14.5 A100 Days
YOSO [20]	PixArt- α	4	30.60	10 A800 Days
TDM	PixArt- α	4	32.01	2 A800 Hours

80x
faster

SOTA performance



OptiBench & Resocratic: Advancing LLMs in Optimization Modeling

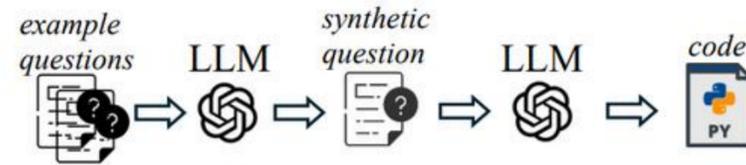
Limitations of Existing Benchmarks

- Only linear programming
- No tabular data
- Small scale (<100 samples)
- No end-to-end evaluation

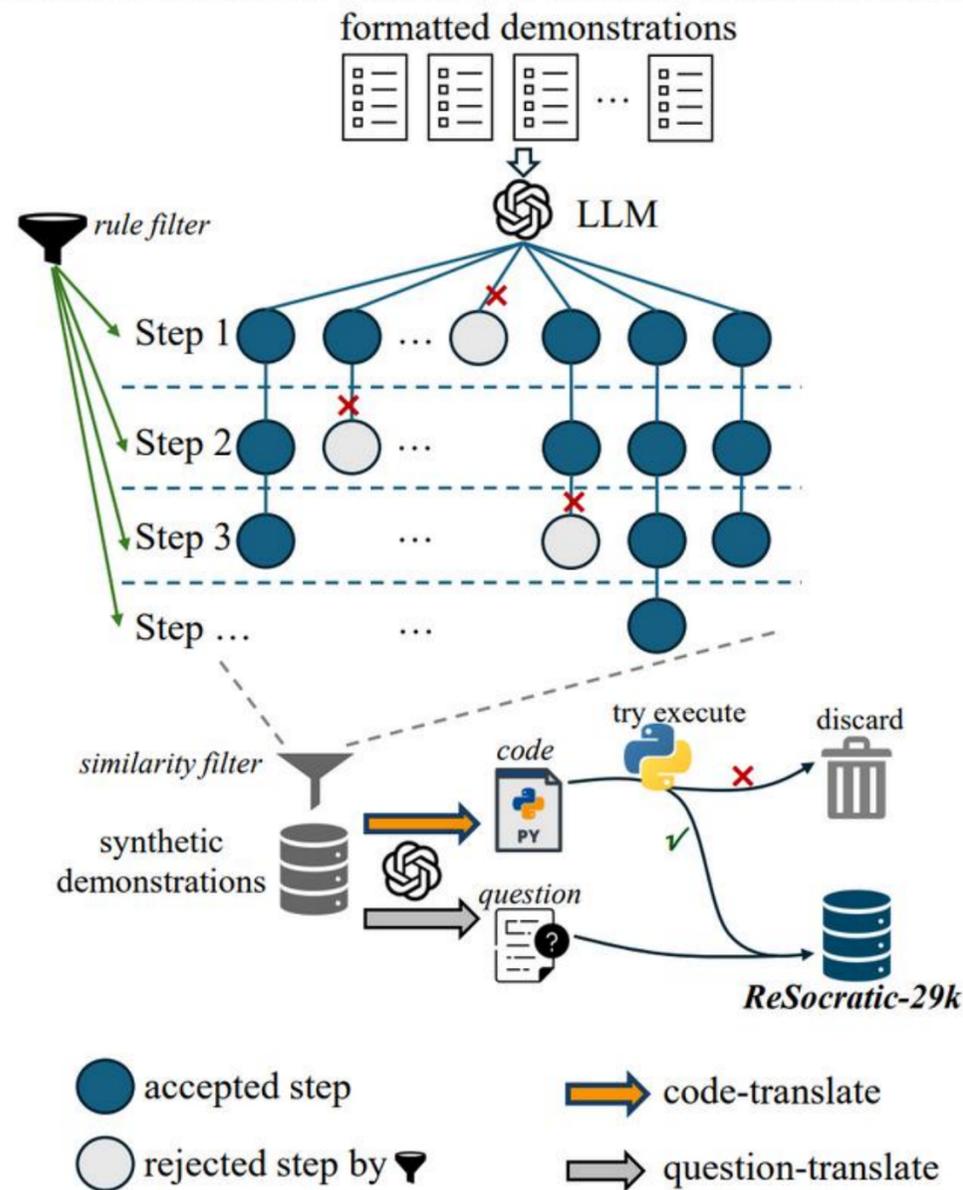
Data Scarcity Challenge

- High annotation cost
- Requires domain expertise
- Performance gap: small vs. large LLMs

Benchmark	Question Form	Size	End2End	Linear		Nonlinear	
				w/ table	w/o table	w/ table	w/o table
ComplexOR (Xiao et al., 2023)	Implicit	37	✓	×	✓	×	×
NLP4LP (AhmadiTeshnizi et al., 2024)	Implicit	57	✓	×	✓	×	×
NL4OPT (Ramamonjison et al., 2022b)	Explicit	289	×	×	✓	×	×
OPTIBENCH (Ours)	Explicit	605	✓	✓	✓	✓	✓



(a) vanilla forward synthesis method



(b) pipeline overview of our ReSocratic

Our Approach:

- **OptiBench Benchmark**
 - 605 manually verified problems
 - Linear & nonlinear programming
 - With/without tabular data
 - End-to-end evaluation
- **Resocratic Data Synthesis**
 - Reverse generation: demonstration → question
 - Step-by-step formatted synthesis
 - 29K high-quality samples (Resocratic-29k)
- **Key Results**
 - GPT-4: 66.1% accuracy (SOTA)
 - Resocratic-29k fine-tuning:
 - Llama-2-7B: 0% → 30.6%
 - Llama-3-8B: 13.6% → 51.1%

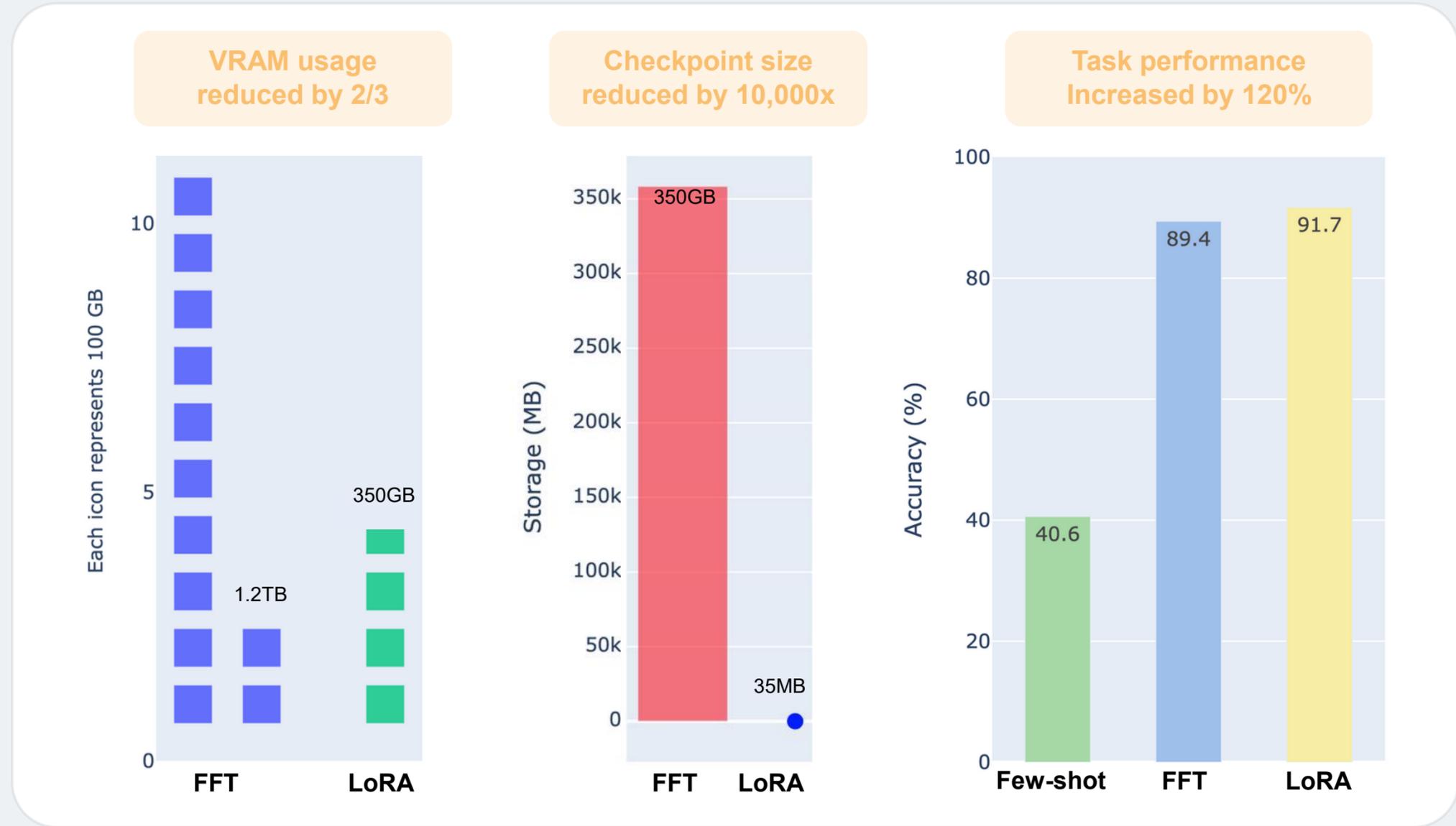
KaSA: Knowledge-Aware Singular-Value Adaptation



Adapting large language models (LLMs) to downstream tasks via full fine-tuning (FFT) is expensive.

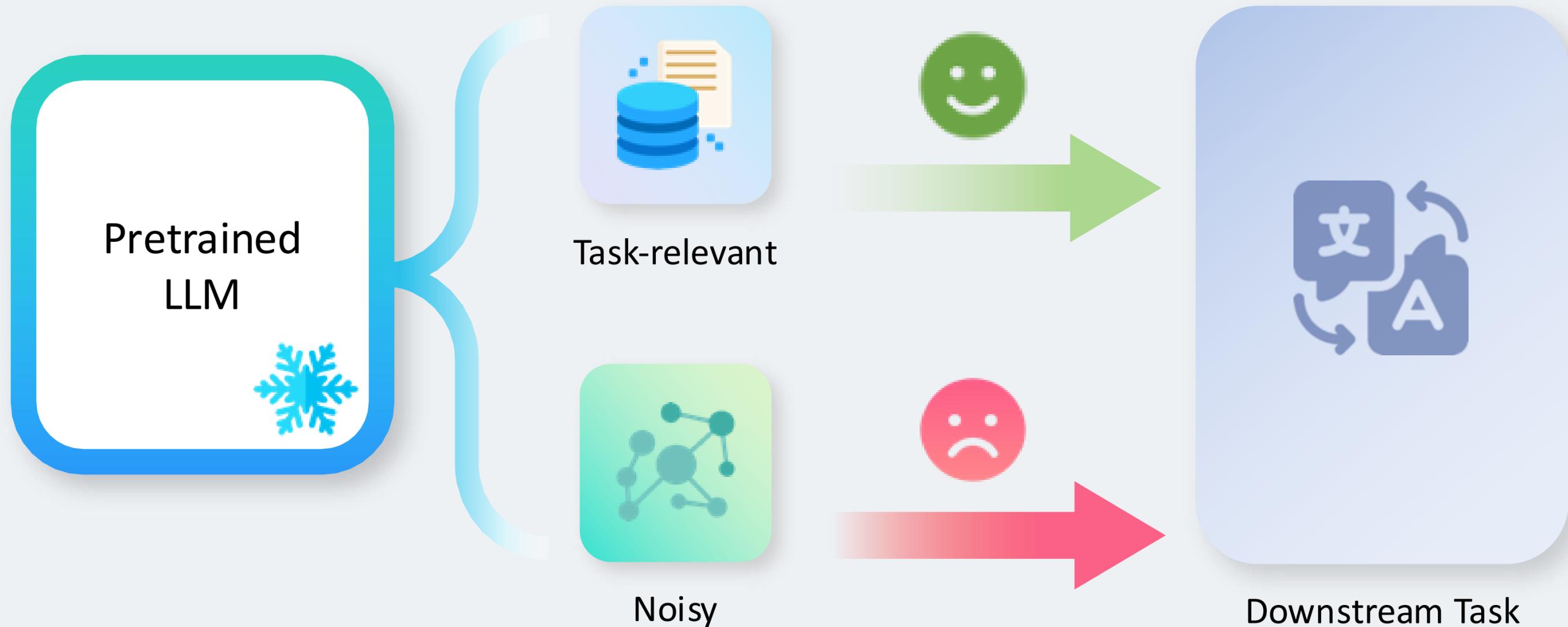


- Sentiment Analysis
- Translation
- Question answering
- Coding



KaSA: Knowledge-Aware Singular-Value Adaptation

Noisy knowledge in the frozen model degrades performance.

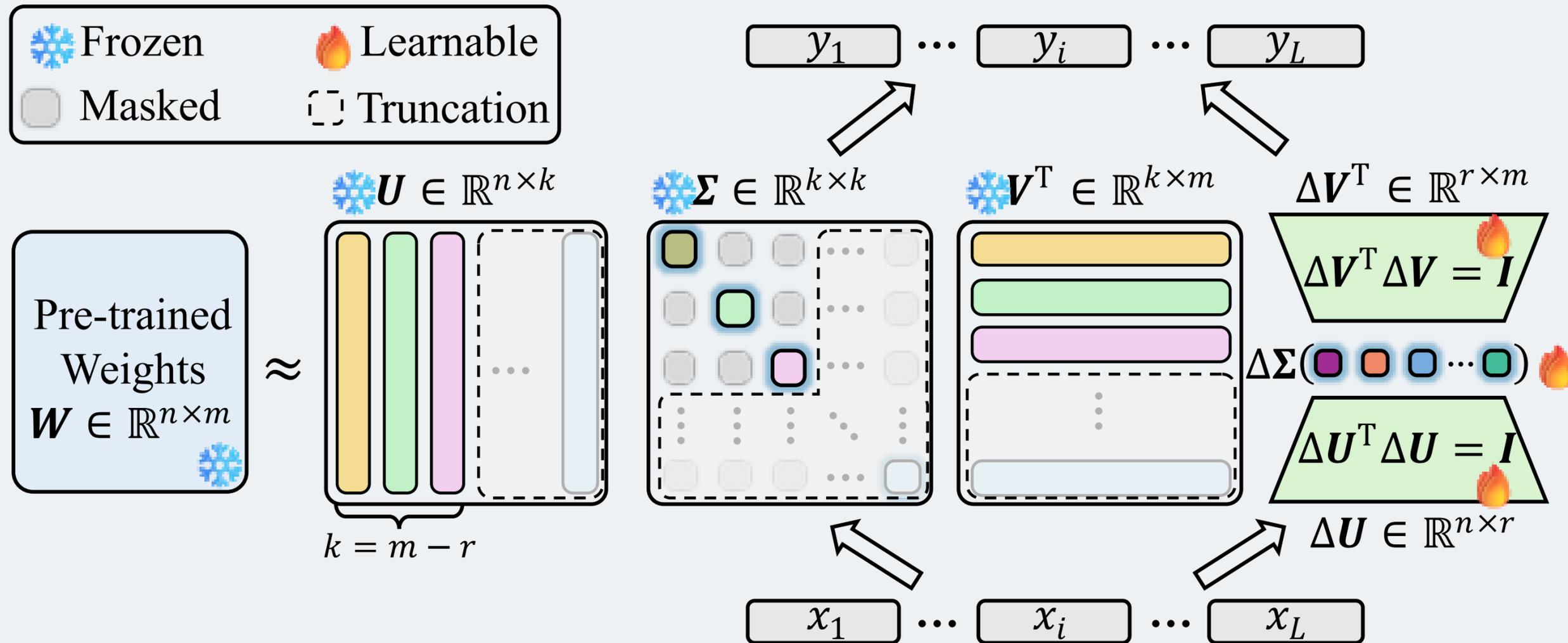


KaSA: Knowledge-Aware Singular-Value Adaptation



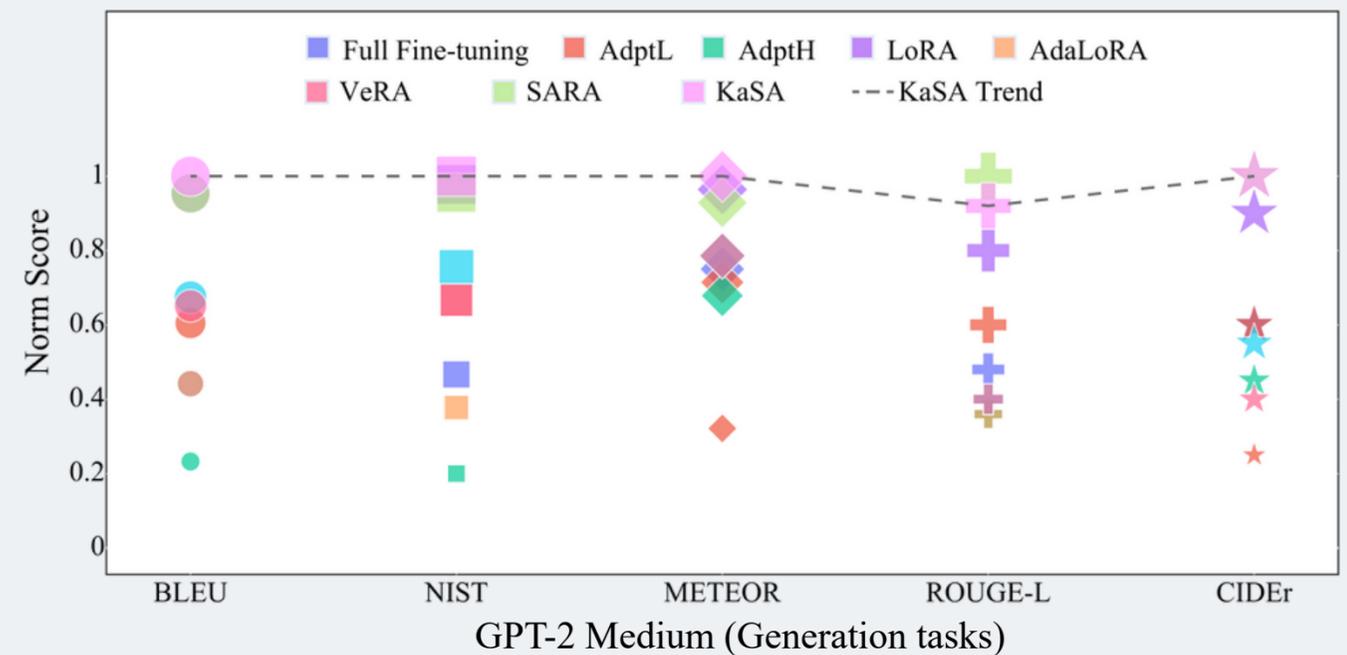
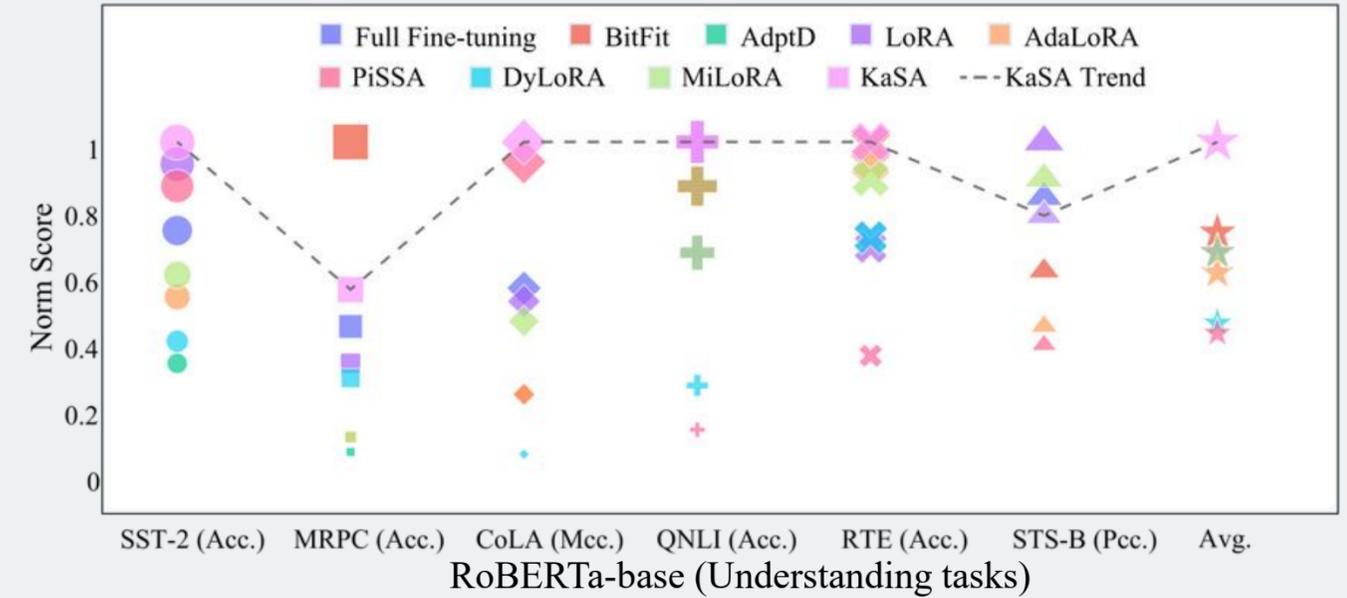
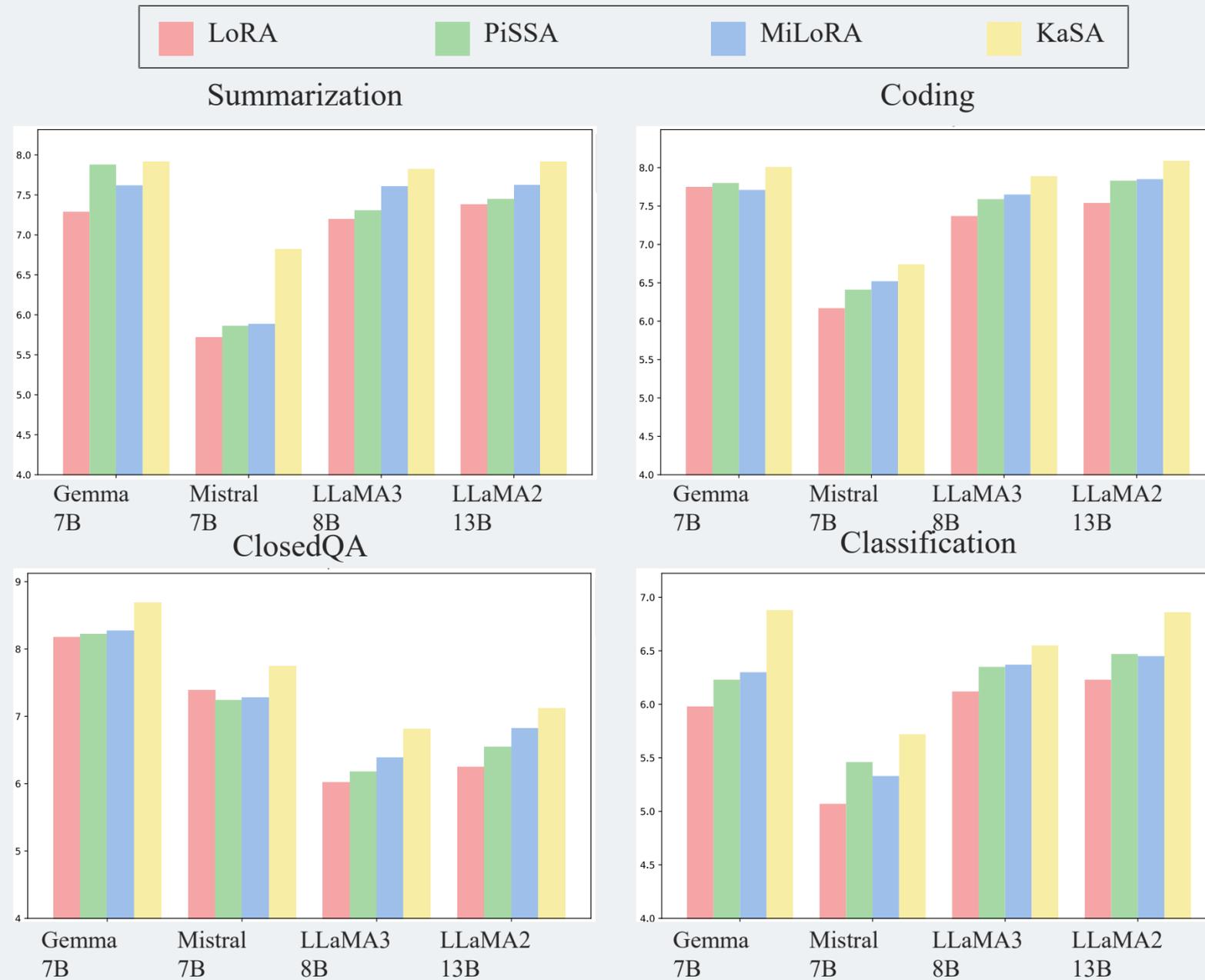
KaSA (Knowledge-aware Singular-value Adaptation):

- Knowledge-based SVD truncation
- Knowledge-aware singular-value adaptation



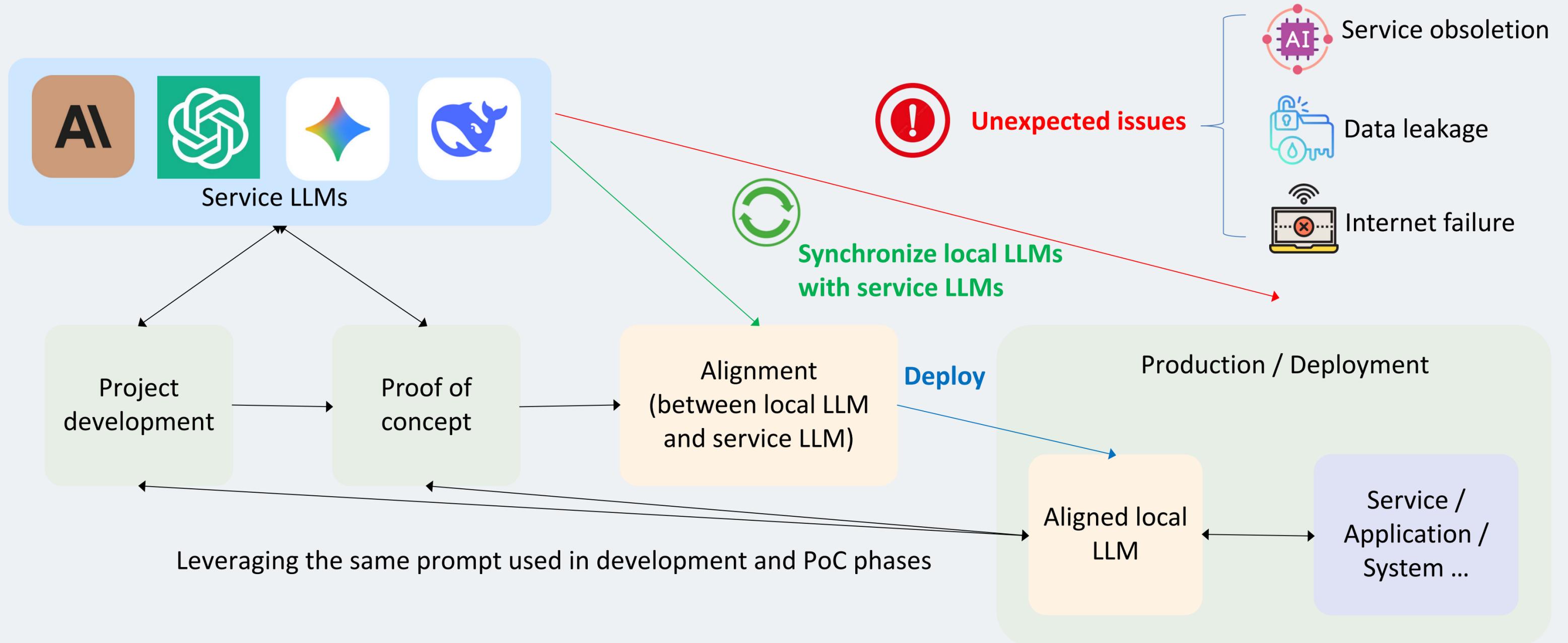
KaSA: Knowledge-Aware Singular-Value Adaptation

KaSA outperforms 14 baselines across 16 benchmarks and 4 synthetic datasets.



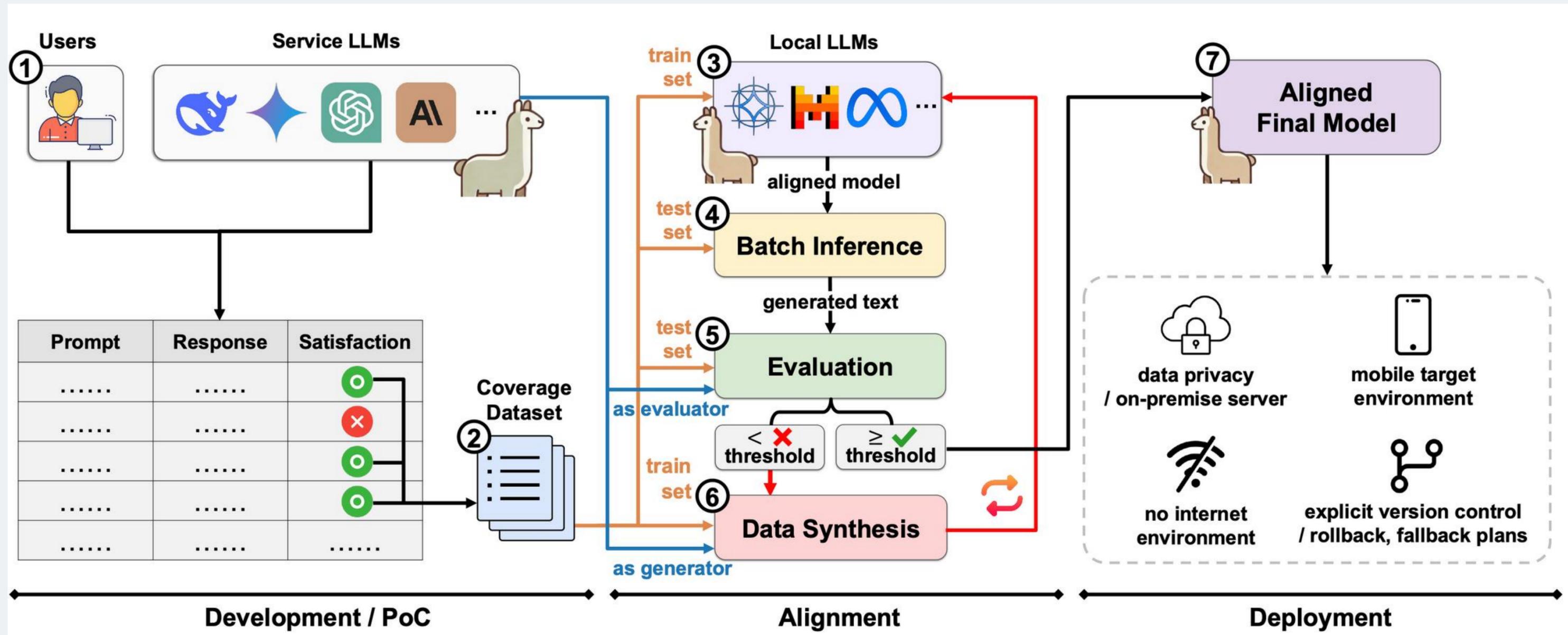
LlamaDuo: LLMOps Pipeline for Distillation

Reliance on cloud-based service LLMs presents challenges.



LlamaDuo: LLMOps Pipeline for Distillation

LlamaDuo: automated alignment pipeline



LlamaDuo: LLMOps Pipeline for Distillation

LlamaDuo enables smaller LLMs to match performance of service LLMs in specific tasks.

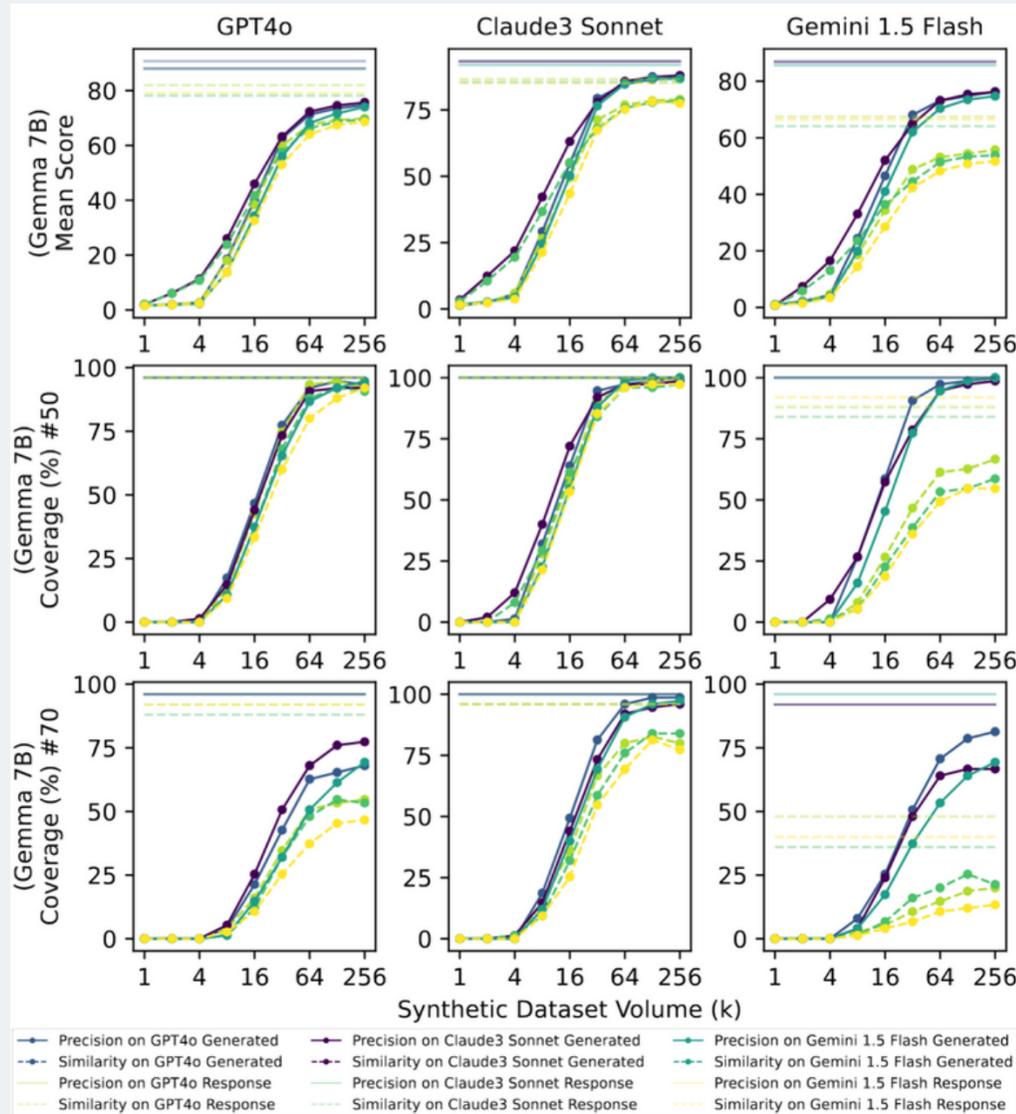


Figure 1: Performance of Gemma 7B fine-tuned on 1~256k dataset volume and evaluated by various service models.

Prompt

Would you classify this email as an advertisement or an account update?

Only 4 hours left in 2020 to make your support known for public radio in the Ozarks! As always, your year-end contribution is tax-deductible.

Don't wait until the ball drops to make your donation...give now, and enjoy the countdown to midnight knowing you've done your part to help the programs you rely on thrive in the year to come.

Target response

This email is an advertisement.

Gemma 2B (aligned on GPT4o generated 128k dataset)

This email is an advertisement. It is asking supporters to make a donation to public radio before midnight. This is intended to persuade supporters to make a donation which is the purpose of an advertisement.

Gemma 7B (aligned on GPT4o generated 128k dataset)

This email is an advertisement.

	GPT4o	Claude 3 Sonnet	Gemini 1.5 Flash
Precision scores	95.5	94.5	93.0
Similarity scores	87.5	91.0	81.0
Precision scores	99.5	100.0	100.0
Similarity scores	100.0	100.0	100.0

Figure 2: Responses by Gemma models fine-tuned on GPT4o generated dataset. Evaluations are given by GPT4o, Claude 3 Sonnet, Gemini 1.5 Flash.

LlamaDuo: LLMOps Pipeline for Distillation

LlamaDuo enables smaller LLMs to match performance of service LLMs in specific tasks.

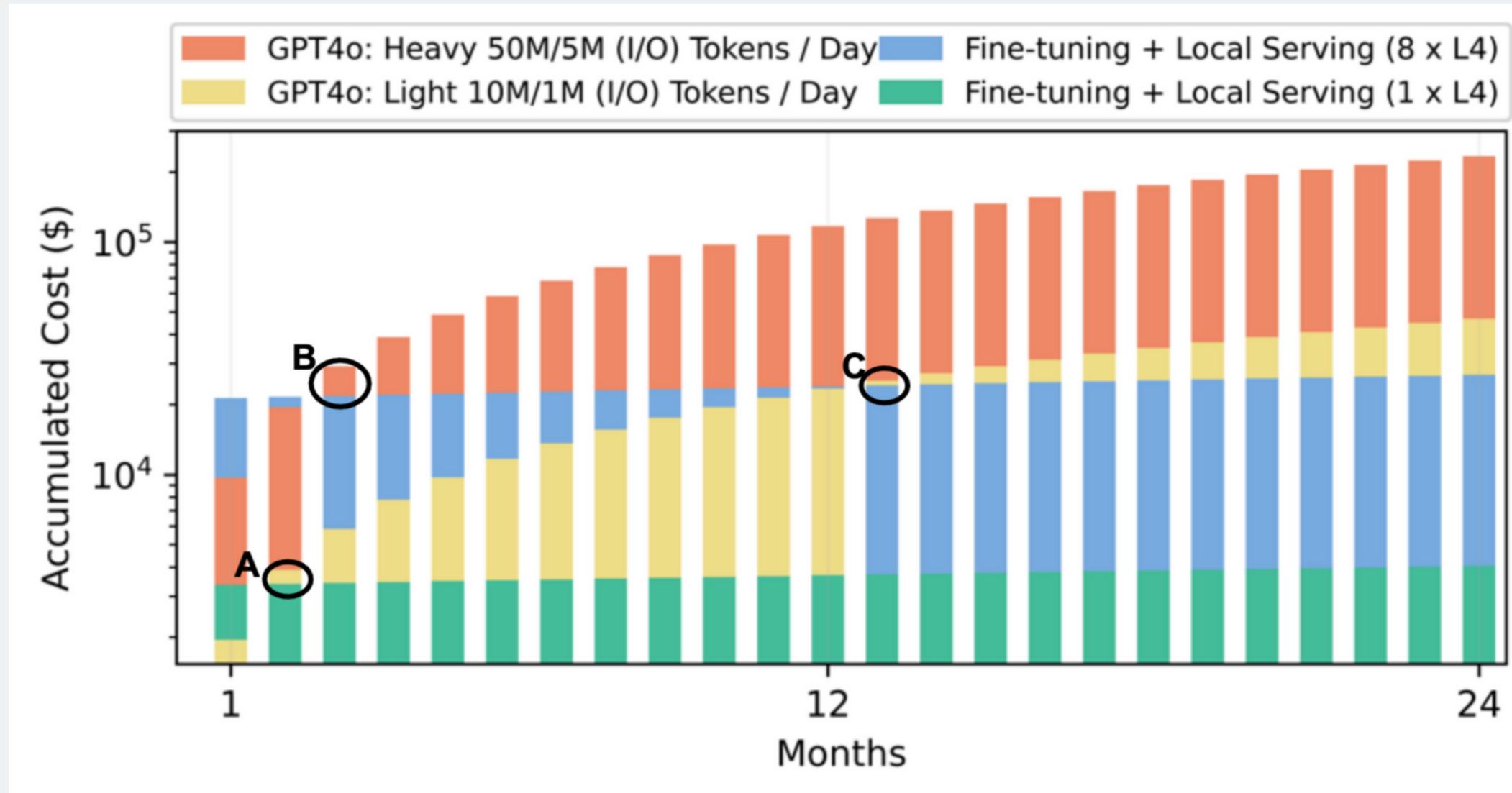


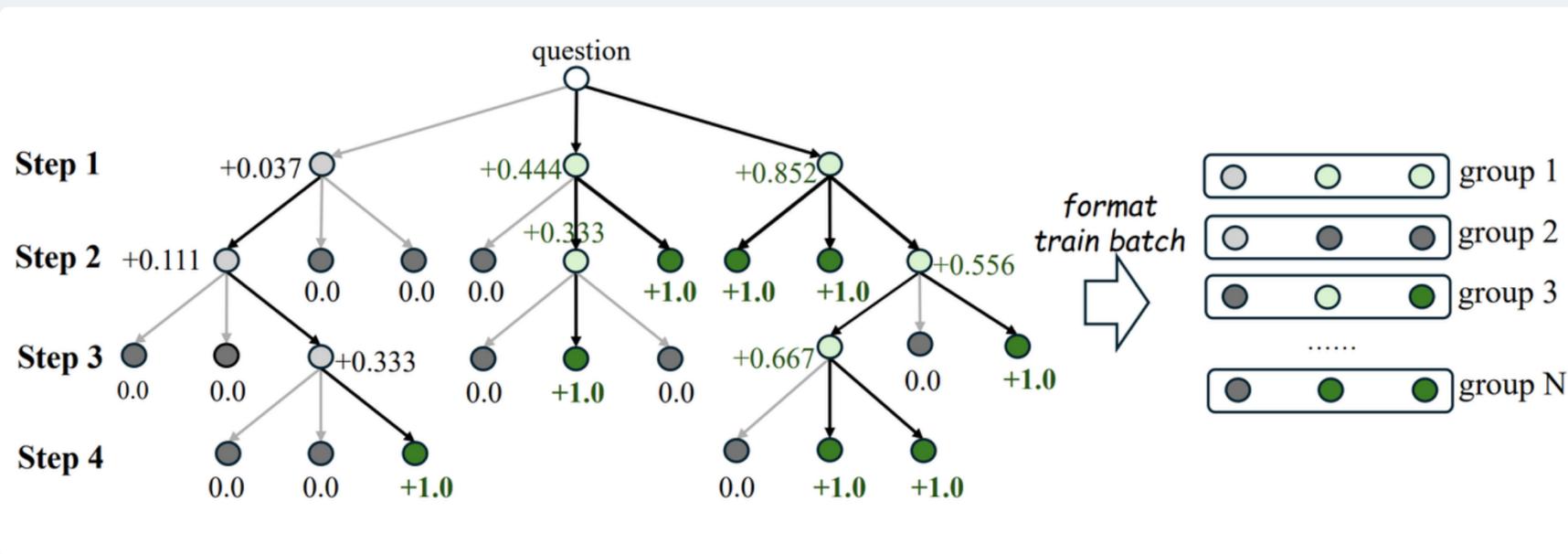
Figure 3: Long-term cost comparison between fine-tuning a local LLM and API-based token usage of GPT4o.

TreeRPO – Tree Relative Policy Optimization

Motivation

- **Sparse Rewards** in RL for reasoning tasks
- **No step-level guidance** in model-free RL (e.g., GRPO)
- **Process Reward Models (PRMs)** are costly to train
- **Goal:** Dense reward signals without PRMs

 *Trajectory-level rewards limit intermediate step optimization*



Solution: TreeRPO

- **Tree Sampling:** Expand N-ary tree from each reasoning step
- **Step-Level Reward Estimation:**
 - Leaf nodes → verifiable reward
 - Parent nodes → expected reward from children
- **Group-Relative Advantage:**
 - Filter groups by reward variance
 - Continuous reward normalization via $\sigma = \mu(1 - \mu)$
- **Training Objective:** GRPO-style clipped PPO + KL penalty

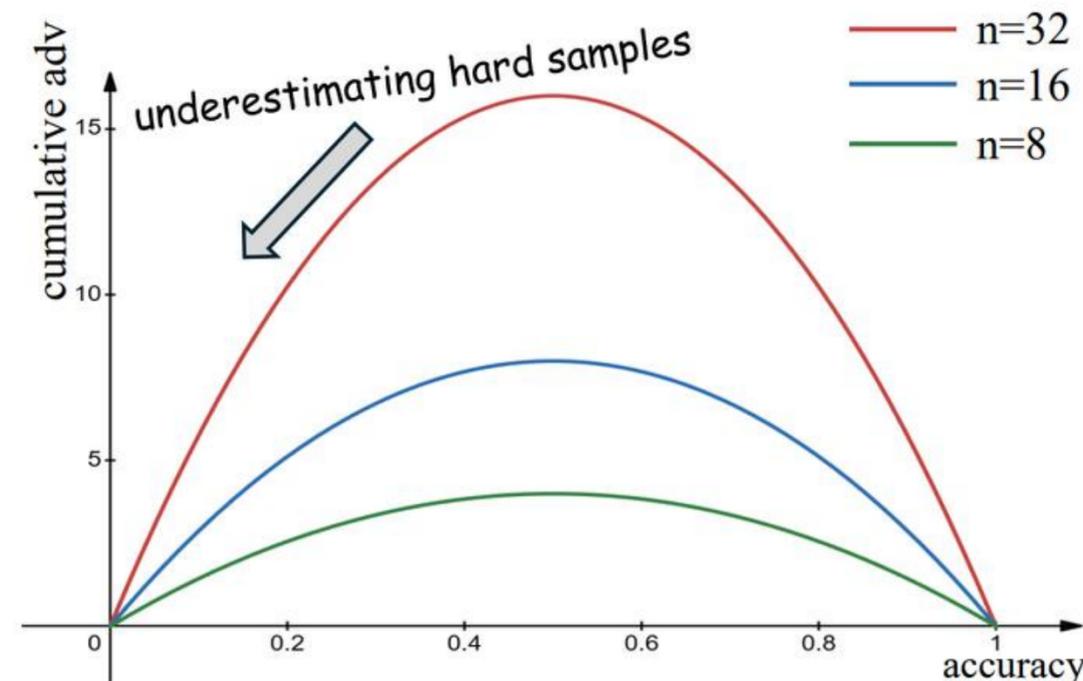
 *Tree structure enables step-wise credit assignment without PRM*

+2.9% over GRPO on Qwen2.5-Math-1.5B
-18.1% response length

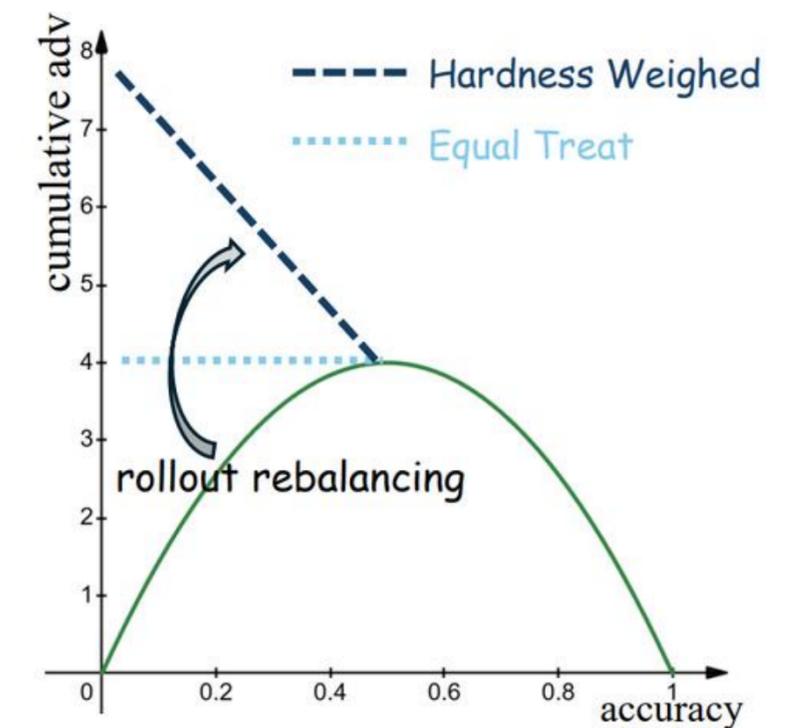
DARS: Why RLVR Falls Short in Reasoning Capacity (Pass@k)

- **Problem:** RLVR with GRPO underweights hard problems
- **Impact:** Limits Pass@K performance
- **Observation:** Naive scaling of rollout size or batch size is inefficient

$$A_{\text{group}} = \sum_{i=1}^G |\hat{A}_i|$$



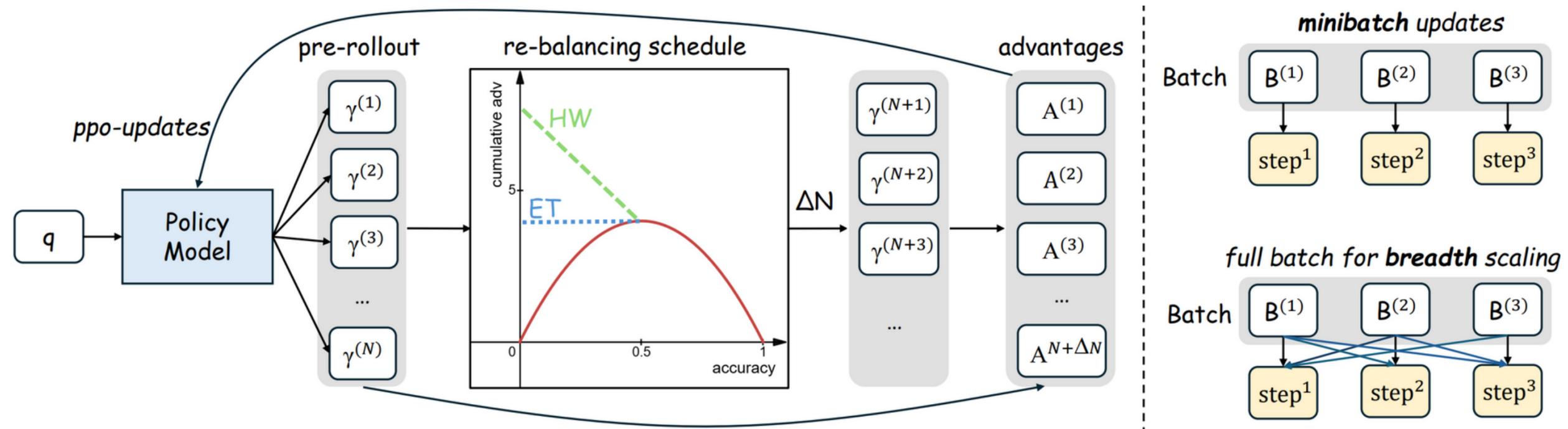
(a) cumulative advantage vs. accuracy



(b) cumulative advantage with DARS

Solution: Difficulty-Adaptive Rollout Sampling (DARS)

> ****We treat DARS as the focal loss in RLVR methods.**



Phase 1: Pre-rollout difficulty estimation.

Phase 2: Multi-stage rollout re-balancing (ET / HW schedules).

- **ET:** Equal treat the samples with accuracy < 0.5
- **HW:** Pay more attention to difficult samples

Breadth Scaling: Full-batch updates + large batch size.

DARS: Depth-Breadth Synergy in RLVR

Key Results

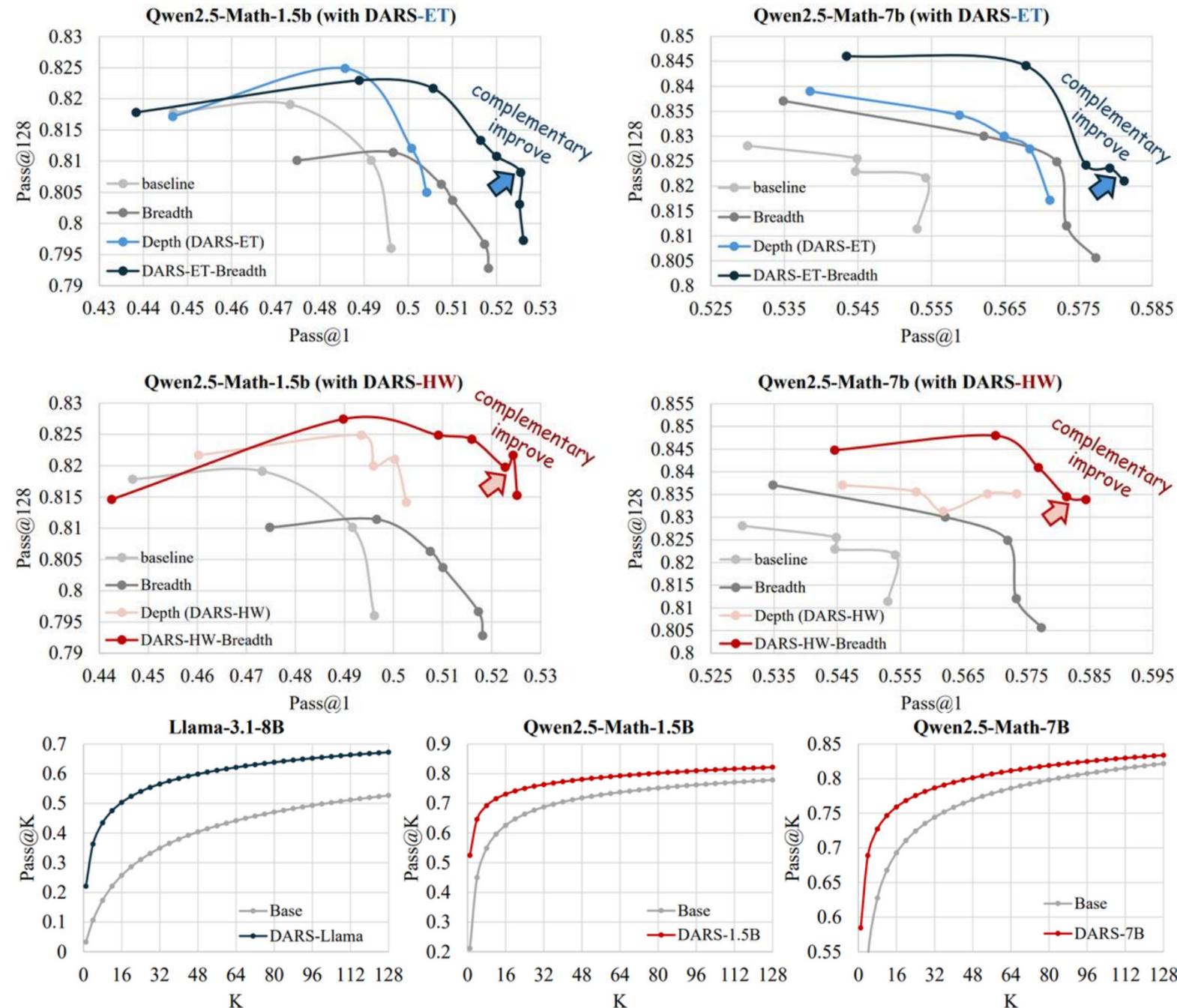
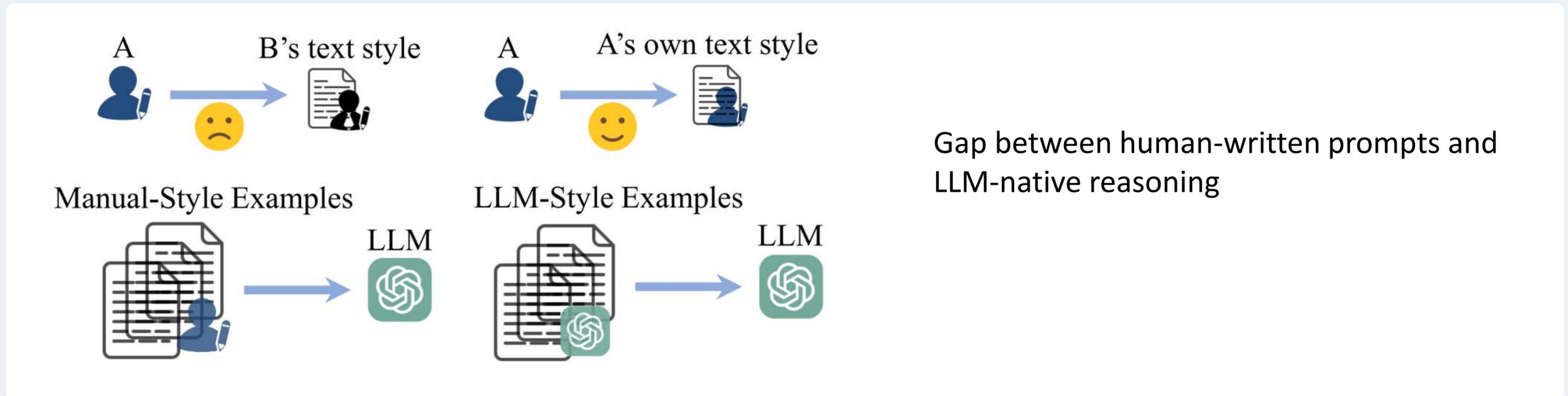


Table 1: Overall performance of $Pass@1$ (Avg@128) and $Pass@128$ of Qwen2.5-Math series.

Base Model + Method	AIME24	MATH500	Olympiad	AMC	Minerva	Avg@128	Pass@128
<i>Qwen2.5-Math-1.5B as the Base Model</i>							
<i>Qwen2.5-Math-1.5B-Base</i>	4.0	35.1	16.2	20.8	9.5	21.1	77.9
<i>Qwen2.5-Math-1.5B-Ins</i>	10.2	67.9	34.7	42.6	24.4	43.5	79.9
<i>Oat-Zero-1.5B</i>	16.4	73.0	35.5	47.4	26.8	46.3	80.3
Our RLVR Replication	14.7	75.9	39.4	47.5	31.2	49.6	79.6
Depth-Naive	16.5	76.2	39.9	47.9	30.9	50.1	79.9
Breadth-Naive	18.5	77.6	41.7	49.8	31.9	51.5	79.2
DARS-1.5B-ET	15.8	76.0	40.9	47.2	30.0	50.0	81.2
DARS-1.5B-ET-Breadth	18.6	79.4	42.9	50.6	31.7	52.5	80.8
DARS-1.5B-HW	14.7	76.4	40.0	48.4	30.8	50.0	<u>82.1</u>
DARS-1.5B-HW-Breadth	19.3	<u>79.0</u>	<u>42.7</u>	51.9	31.6	<u>52.4</u>	82.2
<i>Qwen2.5-Math-7B as the Base Model</i>							
<i>Qwen2.5-Math-7B-Base</i>	11.6	52.3	19.7	35.2	15.3	30.1	82.1
<i>Qwen2.5-Math-7B-Ins</i>	12.9	81.5	39.9	47.0	34.1	52.0	82.3
<i>SimpleRL-Zero-7B</i>	23.3	72.8	36.1	52.8	26.8	46.9	82.5
<i>Oat-Zero-7B</i>	31.3	79.2	42.5	59.4	33.7	53.4	79.7
Our RLVR Replication	26.8	82.2	44.3	57.2	35.7	55.3	81.4
Depth-Naive	28.0	83.8	46.4	59.0	37.3	57.0	80.3
Breadth-Naive	30.5	83.7	47.3	<u>61.4</u>	37.7	57.7	79.2
DARS-7B-ET	26.9	83.2	46.6	57.3	38.5	57.0	81.7
DARS-7B-ET-Breadth	33.3	<u>83.8</u>	<u>47.8</u>	<u>61.3</u>	<u>38.4</u>	<u>58.1</u>	82.1
DARS-7B-HW	30.1	83.5	47.1	59.4	37.2	57.3	83.5
DARS-7B-HW-Breadth	<u>33.0</u>	84.5	48.4	63.0	36.9	58.4	<u>83.4</u>

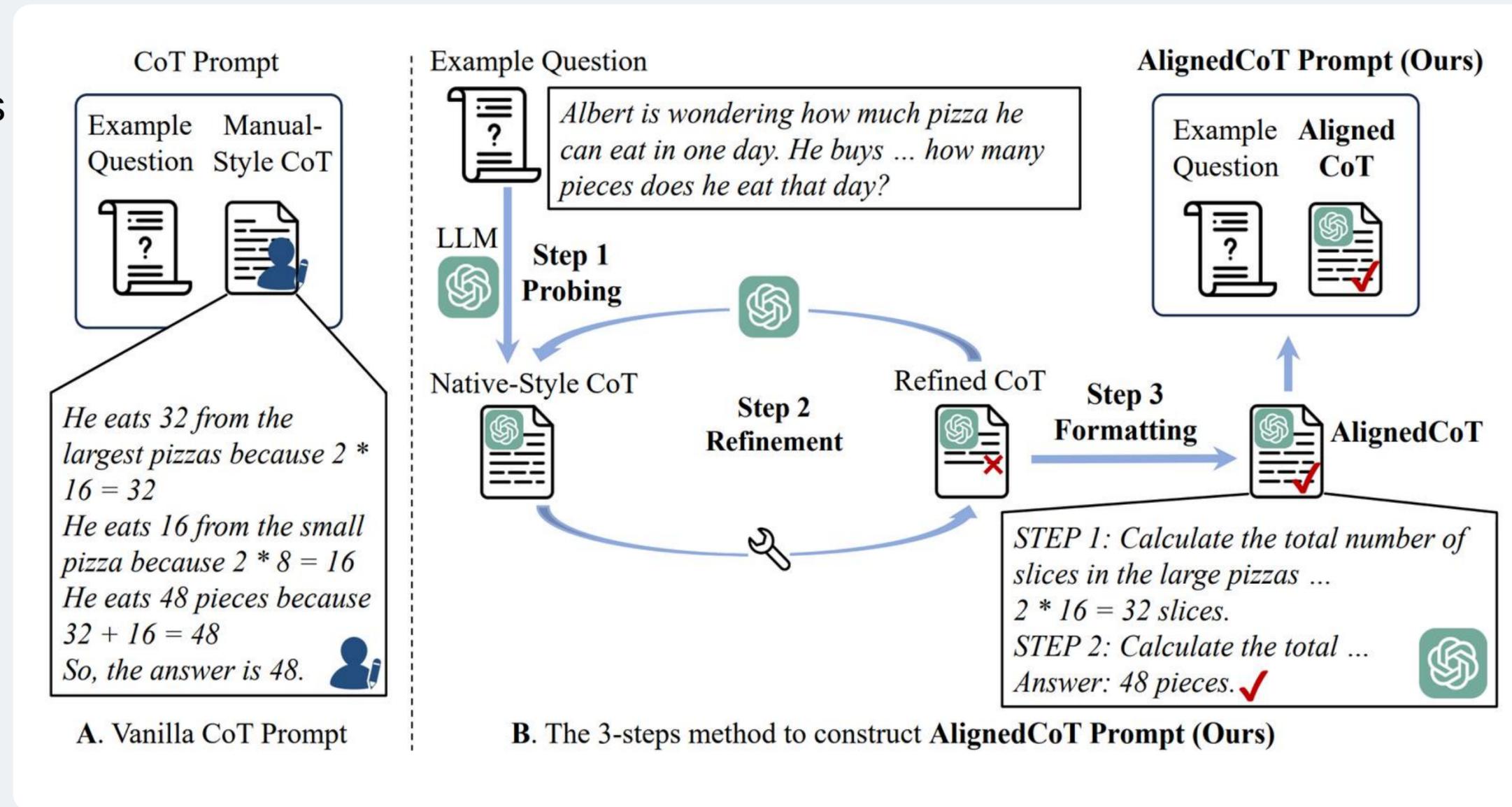
AlignedCoT: Why Do We Need Native-Style Prompting?



- 🔍 LLMs are sensitive to prompt style and wording
- 🧠 Human-written CoT \neq LLM's native reasoning style
- 📉 Performance gap in complex reasoning tasks
- 🎯 Goal: Align prompts with LLMs' own "native" thinking process

How AlignedCoT Works: Probing, Refining, Formatting

- **Step 1: Probing**
Use zero-shot prompt (“Let’s think step by step”) to elicit LLM’s native CoT
- **Step 2: Refining**
Correct errors in generated CoT iteratively
- **Step 3: Formatting**
Unify answer structure and punctuation for consistency
-  **Output: Native-style, correct, and formatted CoT demonstrations**



AlignedCoT Outperforms Human-Crafted Prompts

- ✓ +3.2% avg improvement on GPT-3.5-Turbo
- ✓ +1.7% avg improvement on GPT-4
- ✓ Better logical error detection (78.9% with GPT-4)
- ✓ Works with CoT, Complex CoT, Auto-CoT, and Self-Consistency
- 📦 Released: GSM8K-Align dataset for retrieval-augmented generation

Model	Prompt	GSM8K	AQUA	SVAMP*	AddSub	SingleEQ	Penguins	Avg
GPT-3.5-turbo	CoT w/o AlignedCoT	77.1	54.7	82.8	93.1	96.0	78.1	80.3
	CoT w/ AlignedCoT (Ours)	78.7	57.1	84.8	94.9	97.6	87.7	83.5
	Δ	+1.6 ↑	+2.4 ↑	+2.0 ↑	+1.8 ↑	+1.6 ↑	+9.6 ↑	+3.2 ↑
	Auto-CoT w/o AlignedCoT	78.6	50.4	81.6	92.7	96.5	80.1	80.0
	Auto-CoT w/ AlignedCoT (Ours)	79.8	52.0	82.3	93.9	96.5	84.9	81.6
	Δ	+1.2 ↑	+1.6 ↑	+0.7 ↑	+1.2 ↑	+0.0 ↑	+4.8 ↑	+1.6 ↑
	Complex CoT w/o AlignedCoT	79.6	55.5	82.9	93.1	96.9	81.5	81.6
	Complex CoT w/ AlignedCoT (Ours)	82.4	57.9	85.1	95.2	98.0	86.3	84.2
	Δ	+2.8 ↑	+2.4 ↑	+2.2 ↑	+2.1 ↑	+1.1 ↑	+4.8 ↑	+2.6 ↑



Motivation & Key Contribution

Evaluation Condition	Proportion (%)
Correct response	75.9
Relative distance < 0.05	84.5
Relative distance < 0.10	87.4
Relative distance < 0.20	90.9
Relative distance < 0.30	93.9

Table 1: Proportion of ShowUI-2B model predictions falling within various distance thresholds from the ground-truth bounding box. Evaluated on the ScreenSpot dataset.

Conventional evaluation uses **binary accuracy**

It ignores **how close** a prediction is to the ground truth

Biased hallucinations are common, **structured uncertainty** is hidden



Taxonomy: Define four types of GUI localization hallucinations



Evaluation: Fine-grained classification framework for structured analysis



Metric: Peak Sharpness Score (PSS) captures structured uncertainty

Peak Sharpness Score (PSS)

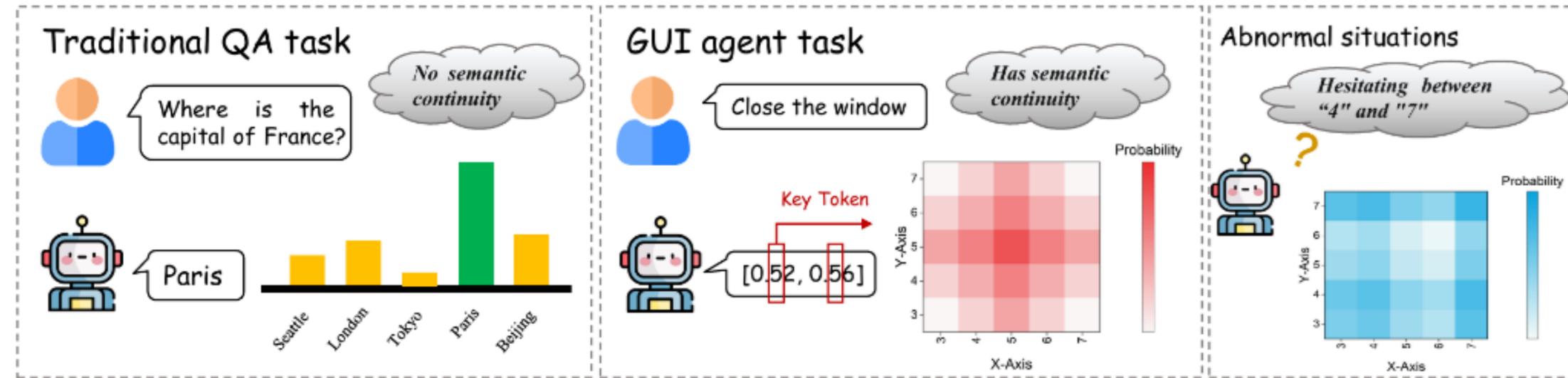


Figure 4: Comparison of logits score distributions between GUI agent tasks and traditional question answering tasks. In GUI coordinate prediction tasks, the model outputs coordinate values, where numeric tokens exhibit semantic continuity. In anomalous cases, a mismatch between this semantic continuity and the continuity of the logits distribution often indicates increased model uncertainty or confusion.

- **PSS** measures alignment between **semantic continuity** and **logits shape**
- High PSS → sharp, unimodal distribution → **confident predictions**
- Low PSS → flat or irregular distribution → **uncertain or hallucinated outputs**
- PSS distinguishes between correct, biased, and misleading errors better than entropy

Measure GUI Agent: Peak Sharpness Score (PSS)



Model	Correct	Biased Hallucination	Other Response
Qwen2-VL-7B	0.34	0.34	0.32
ShowUI-2B	0.59	0.54	0.40
UGround-V1-2B	0.52	0.43	0.25
CogAgent-9B	0.49	0.49	0.24

Table 2: Peak Sharpness Score (PSS) of different GUI agent models across response categories. Values are reported as mean \pm standard deviation.

Wasserstein distance	Avg. Token Entropy	Avg. PSS
Biased Hallucination vs Others	0.358	0.414

Table 5: Wasserstein distance between the logit-score distributions of biased hallucination cases and other errors using both PSS and token entropy.

- PSS is **significantly higher** for correct predictions
- **Biased hallucinations** have similar PSS to correct → semantically close
- **Misleading/Confusion** have lower PSS → lower confidence
- PSS outperforms token entropy in **separability and interpretability**

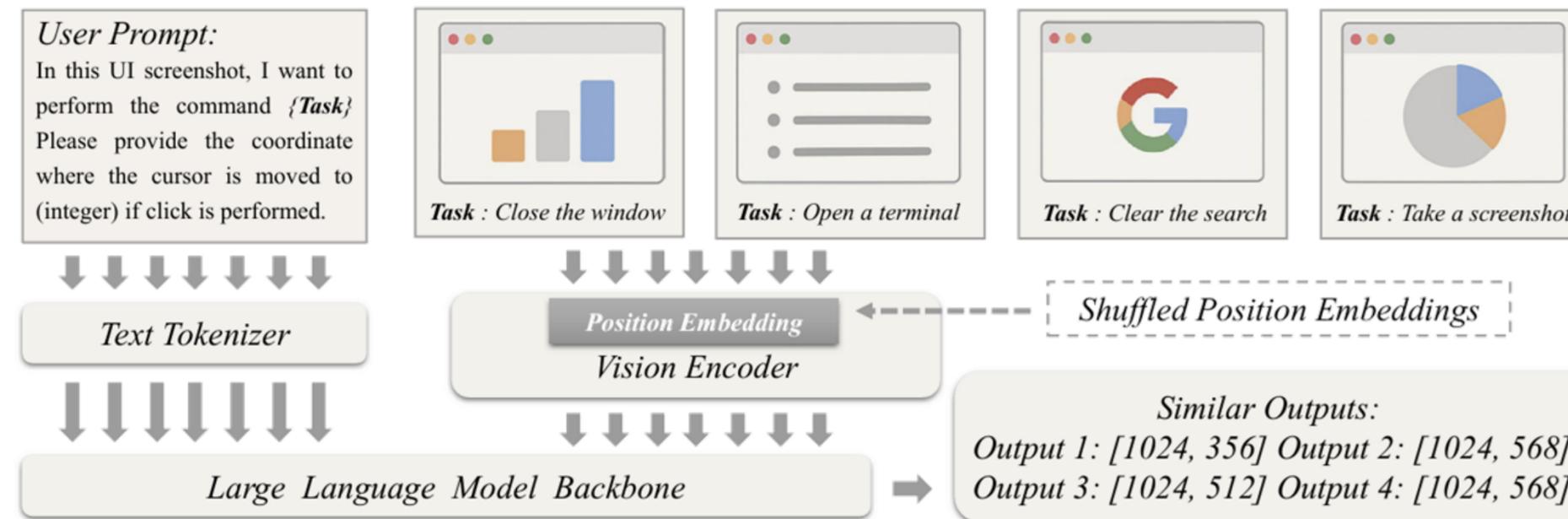
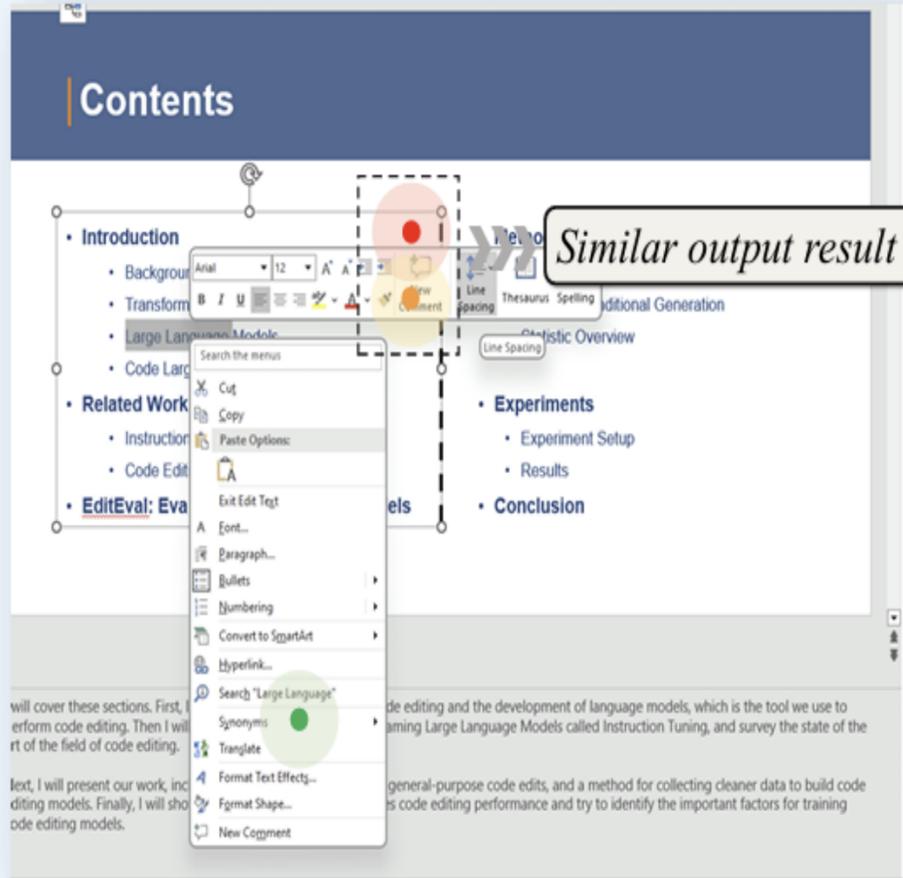


Figure 1: Effect of shuffling visual positional encodings: removing spatial conditioning causes the model to collapse to similar coordinate predictions across independent runs, indicating a position-unconditioned, directional bias rather than random variation. We also observe a similar clustered error pattern on high-resolution images (without shuffling), consistent with position-encoding failures.

- Coordinate prediction fails at high resolution
- Longer visual sequences weaken positional encodings
- Models develop directional, **non-random** coordinate bias
- Bias persists even without spatial signals
- Key question: How to fix it without retraining?

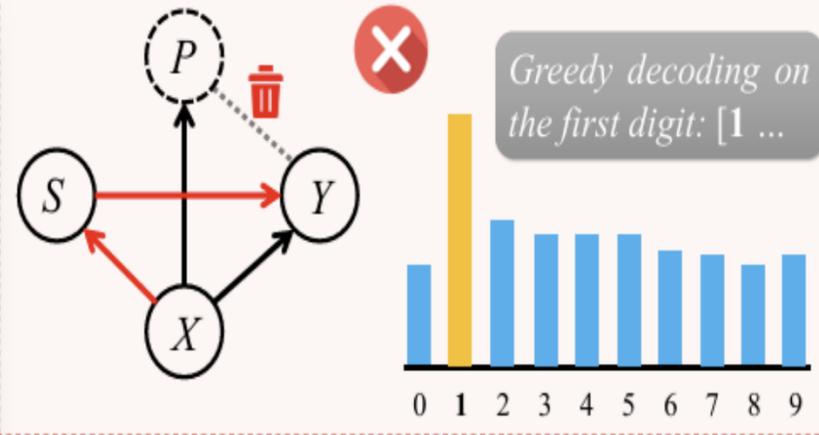
Vision-PE Shuffle Guidance (VPSG) for GUI Agent

User Task: Check out synonyms

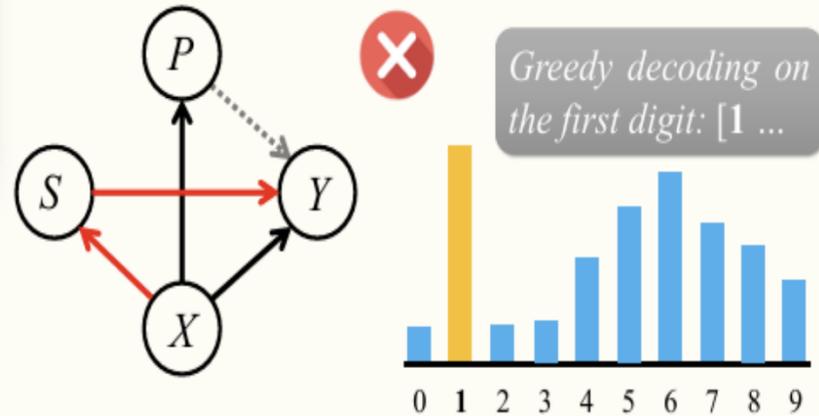


- Response with Shuffled PEs
- Wrong Response with Normal PEs
- Ground Truth

Response with Shuffled PEs: [1024, 356]

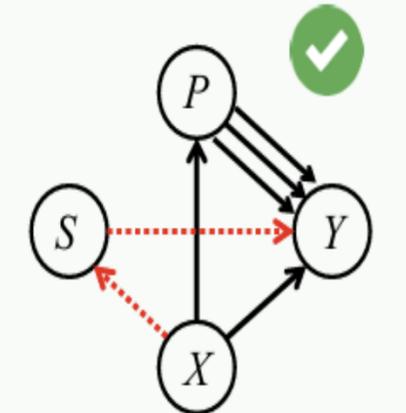


Wrong Response with Normal PEs: [1024, 512]



Comparing >>>

Unbiased Response with VPSG: [659, 857]



Greedy decoding on the first digit: [6 ...]

- **VPSG**: A test-time, plug-and-play correction
- Run base model twice, normal PE and shuffled PE
- Shuffled runs expose position-free numeric bias
- Contrast both routes → detect spurious tendencies
- Apply negative evidence on digit tokens only
- Multi-seed aggregation improves robustness
- Coefficient decay focuses on important digits

Vision-PE Shuffle Guidance (VPSG) for GUI Agent



Model	Development		Creative		CAD		Scientific		Office		OS		Avg
	Text	Icon	Text	Icon									
Trained GUI Action Models 🔥													
SeeClick	0.6	0.0	1.0	0.0	2.5	0.0	3.5	0.0	1.1	0.0	2.8	0.0	1.1
OS-Atlas-4B	7.1	0.0	3.0	1.4	2.0	0.0	9.0	5.5	5.1	3.8	5.6	0.0	3.7
ShowUI-2B	16.9	1.4	9.1	0.0	2.5	0.0	13.2	7.3	15.3	7.5	10.3	2.2	7.7
CogAgent-18B	14.9	0.7	9.6	0.0	7.1	3.1	22.2	1.8	13.0	0.0	5.6	0.0	7.7
Aria-GUI	16.2	0.0	23.7	2.1	7.6	1.6	27.1	6.4	20.3	1.9	4.7	0.0	11.3
UGround-7B	26.6	2.1	27.3	2.8	14.2	1.6	31.9	2.7	31.6	11.3	17.8	0.0	16.5
OS-Atlas-7B	33.1	1.4	28.8	2.8	12.2	4.7	37.5	7.3	33.9	5.7	27.1	4.5	18.9
Generalist Models or Training-free Methods ❄️													
Qwen-VL-7B	0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.1
GPT-4o	1.3	0.0	1.0	0.0	2.0	0.0	2.1	0.0	1.1	0.0	0.0	0.0	0.8
Qwen2-VL-7B	2.6	0.0	1.5	0.0	0.5	0.0	6.3	0.0	3.4	1.9	0.9	0.0	1.6
MiniCPM-V	7.1	0.0	2.0	0.0	4.1	1.6	8.3	0.0	2.8	3.8	3.7	1.1	3.0
Qwen2.5-VL-3B	18.8	1.4	16.7	1.4	8.1	1.6	20.8	5.5	24.3	1.9	16.8	3.4	11.6
Qwen2.5-VL-3B + VPSG	24.7	2.1	20.2	2.1	10.2	1.6	21.5	5.5	26.6	5.7	15.9	1.1	13.3
Qwen2.5-VL-7B	37.7	2.8	19.7	2.1	7.6	1.6	31.3	5.5	41.8	11.3	29.9	10.1	18.5
Qwen2.5-VL-7B + VPSG	40.9	2.1	19.8	2.8	8.1	1.6	30.6	5.6	43.5	13.2	29.9	10.1	19.1

MOTIVATION

Background

-  Alleviating educational anxiety
-  Growth guidance
-  K12 education technology market
-  Lack of continuously evolving

Industry Pain Points

- Rapid Content Decay 
- No growth potential 
- Predefined Q&A Program 
- One-way Rigid Interaction Mode 

SOLUTION

Exclusive Companion



AGENTSFLOW

Private Tutor



From Traditional Toys to Companion Agents



Digital Avatar and AR Technology Enable Customization

KEY RESULTS





AlphaMind

A Closed-Loop LLM Pipeline for Automated Alpha Factor Discovery

Why AlphaMind?

Slow Manual Process

Traditional factor discovery relies on months of manual hypothesis & back-testing

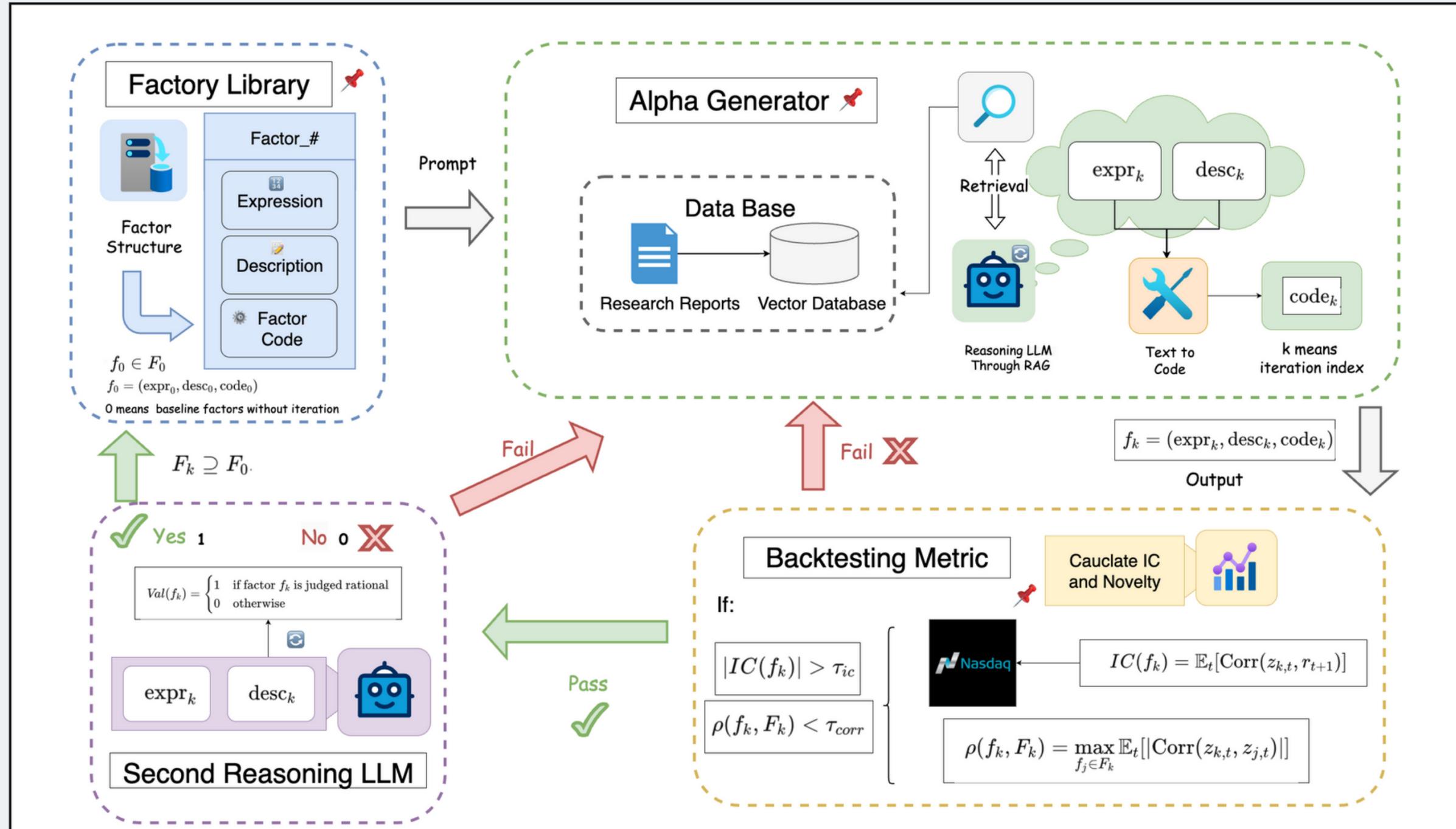
Opaque ML Methods

Deep learning generates signals but lacks interpretability and robustness in live markets

Open-Loop LLM Pipelines

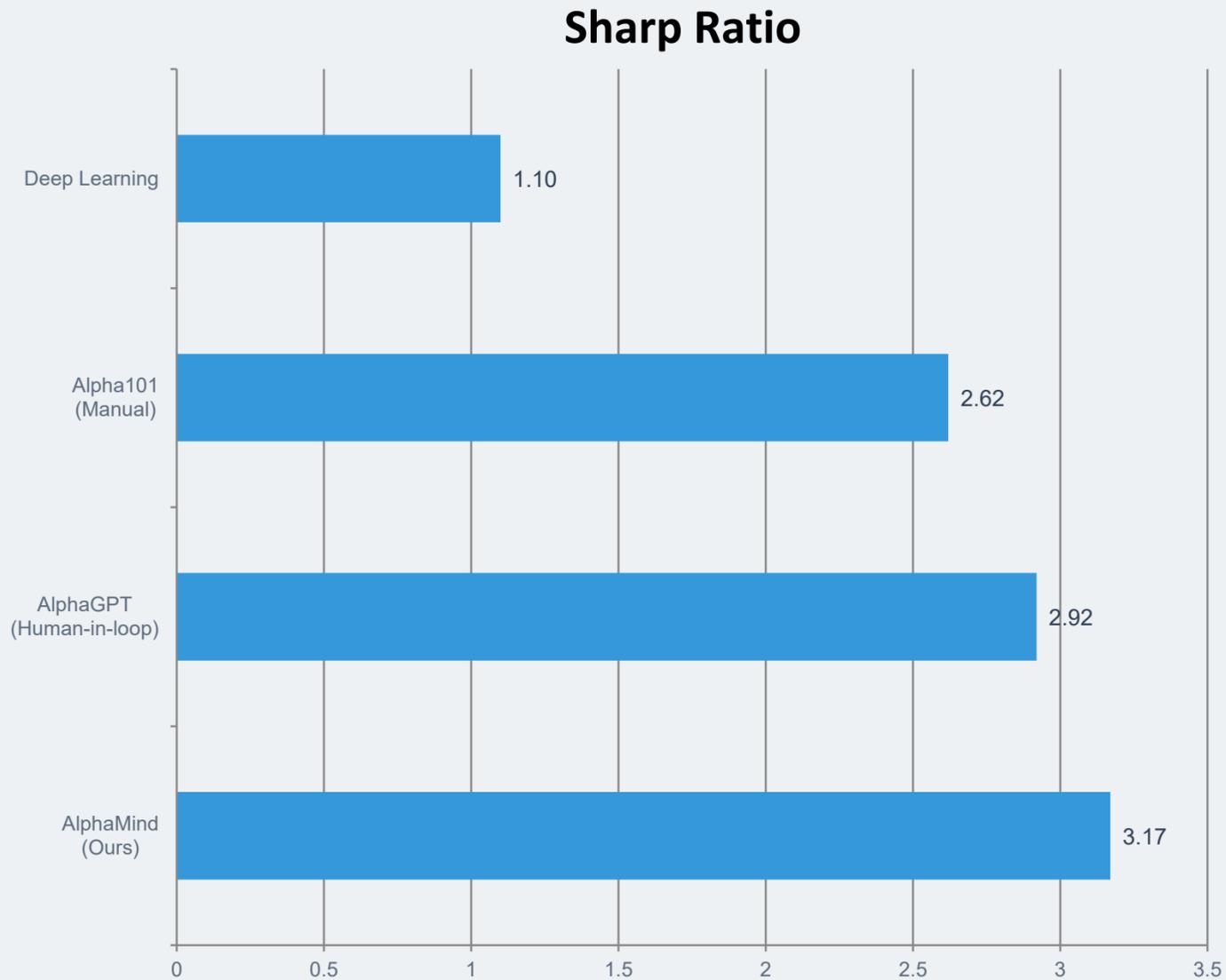
Current LLM systems generate factors without memory or theoretical objective

Architecture of AlphaMind

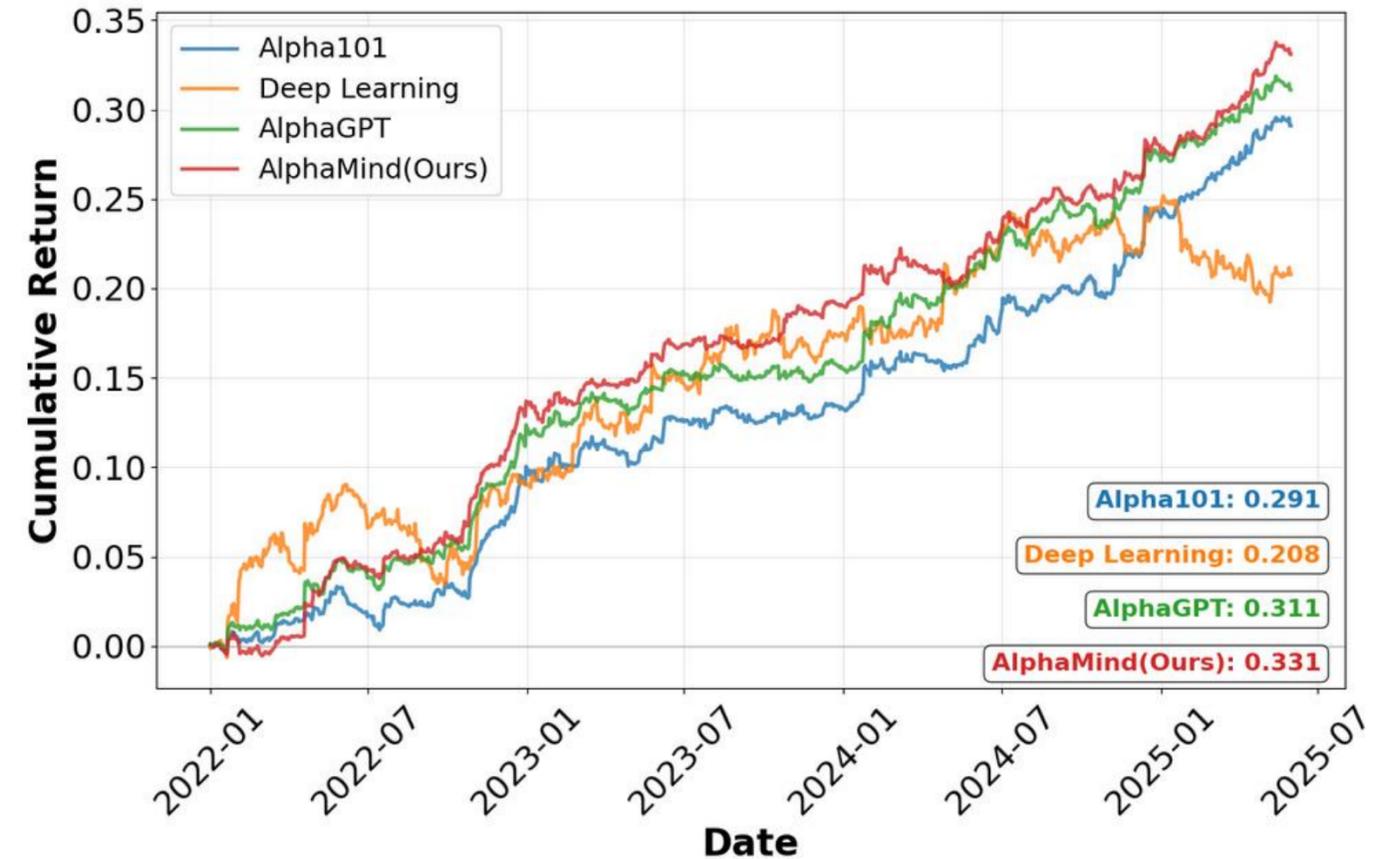


Fully autonomous cycle: retrieval → generation → validation → memory → iteration.

Portfolio Performance Comparison



Cumulative Return Comparison



🏆 Key Achievement

Lowest correlation + stable IC \Rightarrow highest Sharpe and Rate of Return

Thank you!

 Jing Tang, Assistant Professor