

KAIST

DeepAuto.ai

System Prompt Optimization with Meta-Learning

Yumin Choi* Jinheon Baek* Sung Ju Hwang

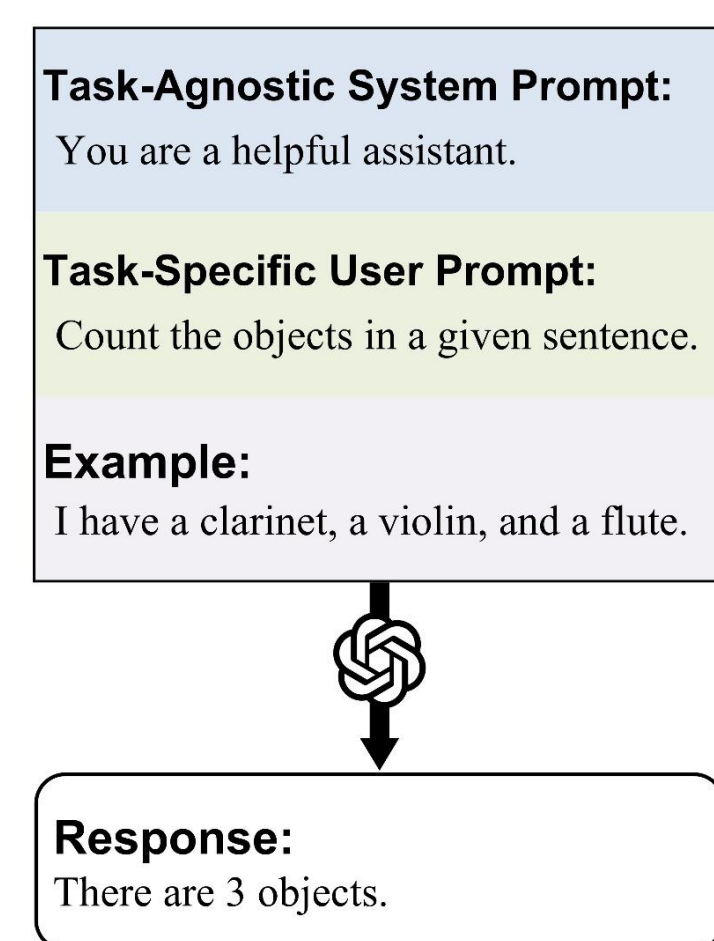
KAIST, DeepAuto.ai (*: equal contribution)



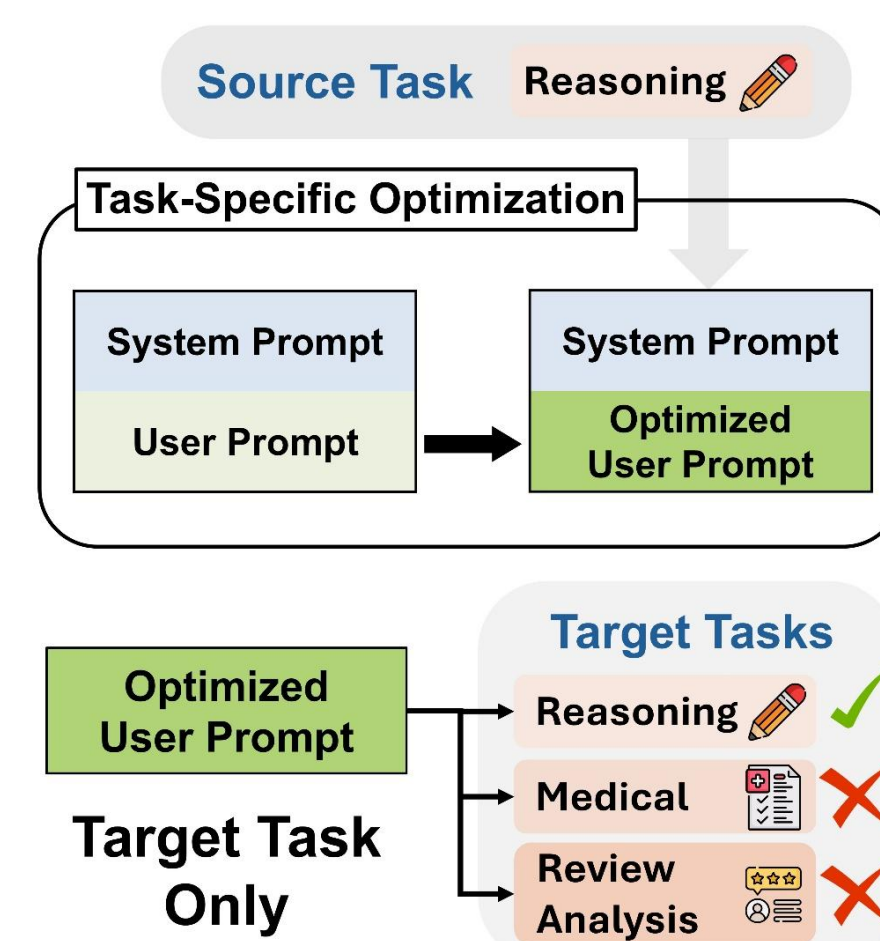
Motivation

- Existing prompt optimization methods **mainly focus on user prompts** for **specific tasks**, overlooking the system prompt that could significantly influence to LLMs.

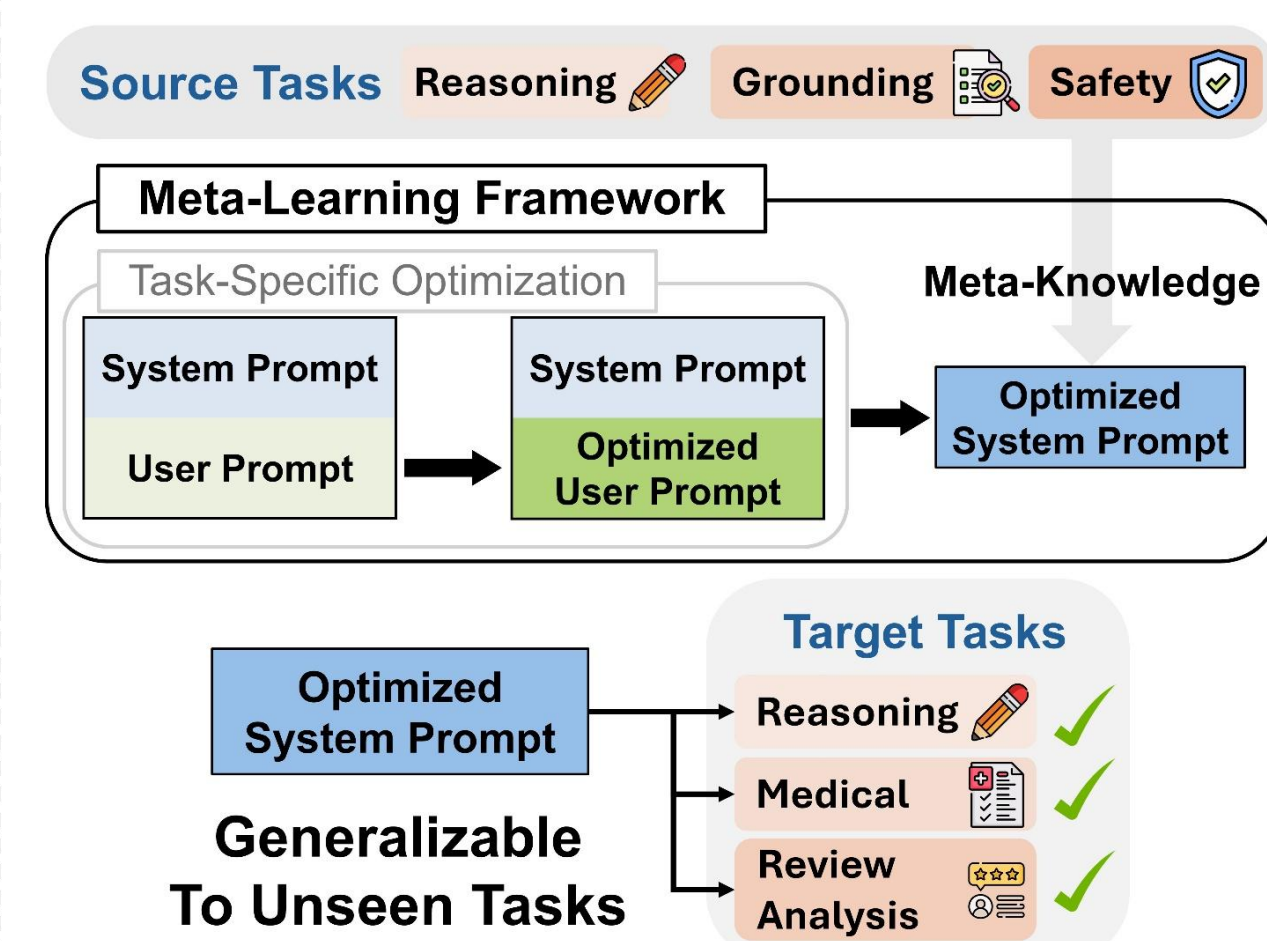
(A) Prompt Input in LLM



(B) Task-Specific Optimization



(C) Bilevel System Prompt Optimization (Ours)

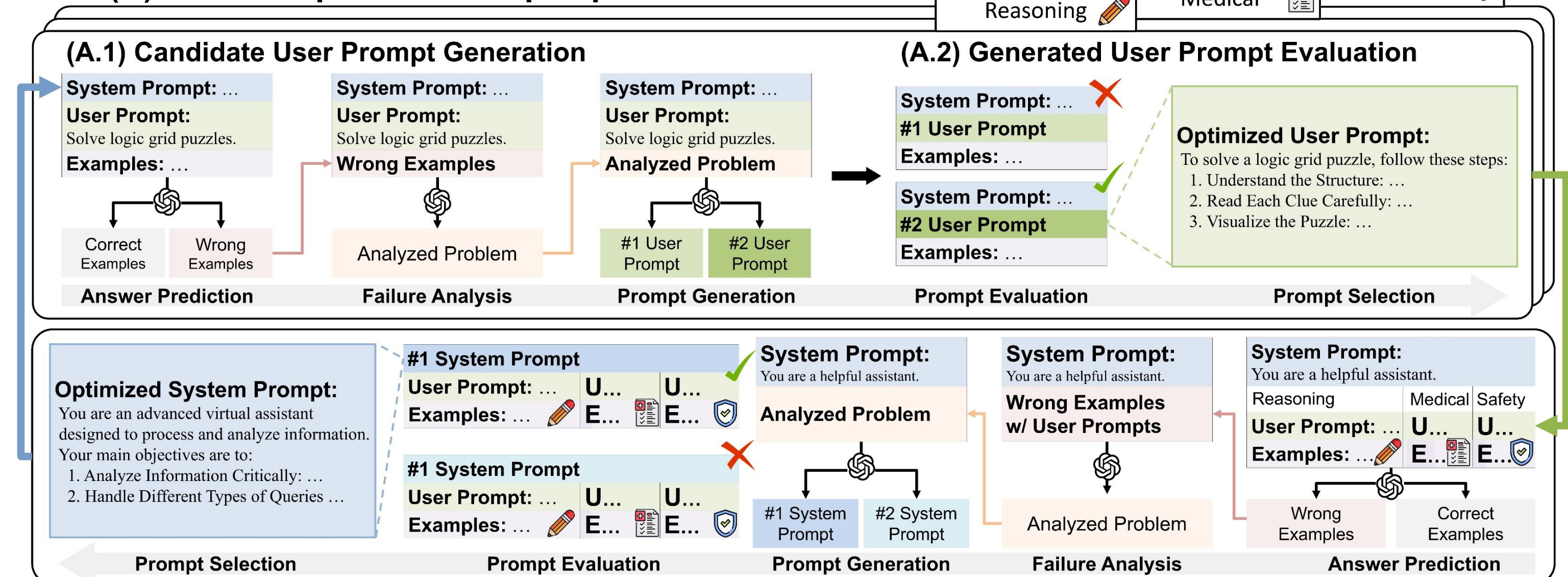


- Optimizing the system prompt: (i) creates a **single, foundational instruction** that generalizes across multiple tasks and domains, and (ii) establishes a **robust behavioral framework** for adapting to unseen user prompts and creating synergy.

MetaSPO: Meta-level System Prompt Optimizer

- We propose **MetaSPO**, a meta-learning framework that solves bilevel optimization by iteratively alternating between an inner loop and an outer loop to ensure **synergy between system and user prompts** and **generalization across tasks**.

(A) Inner Loop: User Prompt Optimization



Bilevel System Prompt Optimization

- We introduce the novel problem of **Bilevel System Prompt Optimization**. Our objective is to obtain system prompts that are robustly effective when coupled with diverse user prompts and highly transferable to a wide range of unseen tasks.

- Why Bilevel? The **Hierarchical Dependency**: The system prompt should generalize across tasks (forming the higher-level objective) while also synergizing with user prompts that are optimized for specific tasks (forming the lower-level objective).

Inner Loop: Optimizes the **user prompts** to maximize task-specific performance.

$$u_i^* = \arg \max_u \mathbb{E}_{(q,a) \sim T_i} [f(\text{LLM}(s, u, q), a)]$$

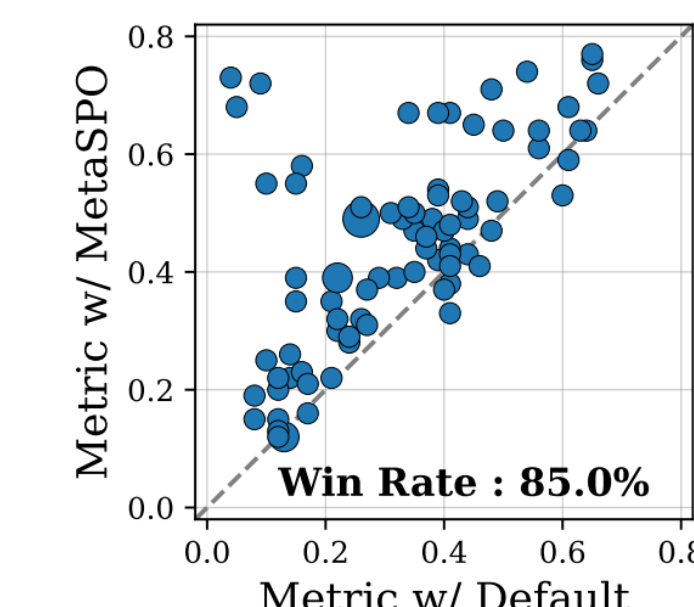
Outer Loop: Optimizes the **system prompt** to generalize across a distribution of tasks and optimized user prompts.

$$s^* = \arg \max_s \mathbb{E}_{T_i \sim \mathcal{T}} [\mathbb{E}_{(q,a) \sim T_i} [f(\text{LLM}(s, u_i^*, q), a)]]$$

Experimental Results

- Our optimized system prompt demonstrates **superior generalization performance** on unseen tasks and robustness across diverse user prompts.

		Medical				Review Analysis			Reasoning			Safety		Grounding		
Methods		Ana.	Ped.	Den.	Sur.	Ele.	Pet	Spo.	Cou.	Epi.	Col.	A.H.	Eth.	N.Q.	Web.	Avg.
Global	Default	36.1	38.9	25.8	32.3	41.3	41.5	29.3	43.5	28.3	56.6	21.2	28.7	15.1	11.6	32.2
	CoT	36.1	42.7	26.0	32.0	36.8	40.3	25.0	45.6	37.2	62.0	21.9	31.9	15.9	12.0	33.2
	Service	34.4	35.2	20.2	30.6	59.0	53.2	52.2	30.6	37.6	56.6	21.1	26.9	11.4	9.9	34.2
	SPRIG	41.6	42.2	28.4	35.7	47.9	47.4	38.6	39.3	29.9	59.9	23.0	31.1	14.1	11.2	35.0
	MetaSPO	45.7	43.1	31.1	36.3	67.2	66.0	61.4	44.5	39.6	64.5	24.9	37.6	9.5	7.7	41.4
Domain	SPRIG	41.2	41.8	29.6	35.3	61.6	57.4	51.3	30.1	34.5	51.5	24.0	32.1	16.1	12.0	37.0
	MetaSPO	48.9	46.7	36.4	40.0	61.8	64.9	61.5	47.1	43.0	66.6	29.1	43.9	19.1	13.7	44.5



- Our optimized system prompt serves as a strong foundation, **enabling improved performance and greater efficiency** when further optimizing user prompts.

Methods	Med.	Rev.	Rea.	Saf.	Gro.	Avg.
Default	45.1	68.9	64.0	59.9	17.5	51.1
SPRIG	45.4	69.3	65.3	64.7	17.7	52.5
MetaSPO	45.6	71.4	67.3	67.2	19.9	54.3

