PAS: PLUG-AND-PLAY PROMPT AUGMENTATION System

Anonymous authors

Paper under double-blind review

Abstract

In recent years, the rise of Large Language Models (LLMs) has spurred a growing demand for plug-and-play AI systems. Among the various AI techniques, prompt engineering stands out as particularly significant. However, users often face challenges in writing prompts due to the steep learning curve and significant time investment, and existing automatic prompt engineering (APE) models can be difficult to use. To address this issue, we propose PAS, an LLM-based plug-and-play APE system. PAS utilizes LLMs trained on high-quality, automatically generated prompt complementary datasets, resulting in exceptional performance. In comprehensive benchmarks, PAS achieves state-of-the-art (SOTA) results compared to previous APE models, with an average improvement of 6.09 points. Moreover, PAS is highly efficient, achieving SOTA performance with only 9000 data points. Additionally, PAS can autonomously generate prompt augmentation data without requiring additional human labor. Its flexibility also allows it to be compatible with all existing LLMs and applicable to a wide range of tasks. PAS excels in human evaluations, underscoring its suitability as a plugin for users. This combination of high performance, efficiency, and flexibility makes PAS a valuable system for enhancing the usability and effectiveness of LLMs through automatic prompt engineering. The codebase is available at https://anonymous.4open.science/r/ape_gen-4214/.

028 029 030

031

038

003

005

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

1 INTRODUCTION

Large language models (LLMs) have demonstrated exceptional performance across diverse tasks spanning myriad domains Touvron et al. (2023); meta llama (2024); Bai et al. (2023). Variations in task contexts demand tailored response paradigms from large language models and the flexibility in adjusting response style is key to meeting the diverse needs. Nonetheless, pursuing a specific application or enhancing certain capabilities through model fine-tuning involves considerable data collection costs, computational resources, and the potential for a seesaw effect.

To address this issue, prompt engineering has emerged as a crucial approach among the various techniques for LLMs due to its extremely low cost and significant enhancement of LLM perfor-040 mance Sahoo et al. (2024). This technique leverages the inherent capabilities of LLMs to under-041 stand and generate human-like text, enabling them to perform a wide range of applications from 042 natural language understanding and generation to specialized domains such as medical diagnosis 043 and legal analysis Meskó (2023); Giray (2023). In prompt engineering techniques, few-shot learn-044 ing Brown et al. (2020) stands out because it provides a small number of examples to guide the 045 model, thereby enhancing its task-specific performance. Chain-of-Thought (CoT) Wei et al. (2022) prompting guides models to reason through problems step-by-step, improving logical consistency 046 and accuracy. In-context learning Dong et al. (2022), by embedding relevant examples and instruc-047 tions directly within the prompt, allows models to adapt dynamically to new tasks. 048

Despite the potential of existing methods, prompt engineering is not user-friendly and requires meticulous design. Crafting effective prompts demands a deep understanding of both the model and the task at hand. This process can be time-consuming and often involves extensive trial and error to optimize performance. To tackle those weak points, automatic prompt engineering (APE) is designed for easier prompt generation Zhou et al. (2022b); Pryzant et al. (2023b); Cheng et al. (2023). Although APE models can automatically enhance the prompt, they have to use massive amounts of



human-labeled data Cheng et al. (2023); Ouyang et al. (2022); Rafailov et al. (2024). Additionally,

previous methods failed to construct a flexible, user-friendly and effective APE model Zhou et al.

(2022a); Pryzant et al. (2023a). They face the following three key challenges:

Figure 1: We first present the pipeline of the PAS in (a). PAS takes user prompts, enhances them, and then inputs the augmented prompts into LLMs. As illustrated in (b), PAS significantly improves responses across all categories in human evaluation.

Figure 2: Case Study 1, Red text is the complementary prompt generated by PAS. We can see PAS can give complementary prompt to avoid logic traps.

C1. Poor Effectiveness. Although several works aim to automatically enhance the quality of prompts, they often rely on outdated evaluation metrics and do not consistently yield significant improvements across various benchmarks and models Cheng et al. (2023); Ouyang et al. (2022). Additionally, these models fail to include human evaluations, thus lacking valuable feedback from human users Cheng et al. (2023); Pryzant et al. (2023b).

C2. Low Flexibility. Previous works primarily focus on the performance of APE models, over-looking the importance of flexibility, specifically their model-agnostic and task-agnostic capabilities Pryzant et al. (2023b); Yang et al. (2023b). Low-flexibility APE systems can lead to computational waste and hinder their application on LLMs Pryzant et al. (2023b). Additionally, these low-flexibility systems have limited applicability across various scenarios, making them less versatile for diverse use cases and reducing their overall effectiveness in practical applications.

C3. Limited Practicality Among the previous methods, only BPO offers flexibility for prompt augmentation. However, BPO chose to rewrite the questions for prompt augmentation. this carries the risk of the new prompts deviating from the oiginal intent. Additionally, when dealing with long inputs, BPO is unable to control the generating time. It is also impossible for real-time display.

To address these issues, in this paper, we designed a plug-and-play system, called Prompt Augmen-099 tation System, PAS. The system employs a specialized model designed to generate prompts tailored 100 to various tasks, thereby transforming the prompts into domain-specific expert prompts which are 101 the inputs for LLM. The central idea behind PAS is to decouple the scene requirements from the 102 model's capabilities, as shown in Figure 1(a), which is a two-stage alignment paradigm. In this two-103 stage approach, the prompt augmentation (PA) model is utilized to understand prompts and generate 104 the specific constraints or hints for them, while the LLM focuses on instruction following and knowl-105 edge capabilities. The PA model and LLM collaborate to accomplish tasks. Prompt augmentation can be achieved in two ways: by entirely rewriting the prompt or by generating supplementary 106 content to enhance the original prompt. BPO Cheng et al. (2023) adopted the rewriting approach; 107 however, this method risks altering the original intent and results in unpredictable prompt generation

2

054 055

108 times, especially for lengthy inputs. To address these limitations, we opted to add complementary 109 content following the original prompt. As illustrated in Figure 2, the PA model comprehends the 110 original prompt, identifies potential logical traps, and generates supplementary content as a hint. 111 The core contributions of this paper are summarized as follows:

- New Perspective. To the best of our knowledge, this work is the first to construct a curated prompt complementary dataset without human labor. Additionally, we are the first to utilize this dataset to train LLMs to construct the PAS to automatically complement user prompts instead of directly modifying them.
- New Method for Auto Prompt Augmentation System. We propose a new data collection method to construct a complementary prompt dataset for training a Prompt Augmentation model. Furthermore, we introduce a plug-and-play Prompt Augmentation System (PAS), which can be seamlessly integrated into any LLM to enhance model performance without additional training costs. PAS is designed to interpret prompts and generate specific constraints or hints to guide responses.
- SOTA Performance. (1) SOTA Performance on Multiple Benchmarks. Our PAS model 123 achieves SOTA performance across multiple models and comprehensive benchmarks. It 124 also outperforms the previous SOTA model, BPO, under identical experimental settings. 125 Furthermore, it demonstrates superior performance on human evaluation metrics, as shown in Figure 1(b), highlighting the outstanding capabilities and potential applications of the 127 PAS model. Figure 2 further illustrates that our model has good performance in real-world 128 problems. (2) High Flexibility. Our PAS model can be plugged into any LLM, and is model 129 and task-agnostic. It achieves SOTA performance across all models and tasks, demonstrat-130 ing its exceptional flexibility. (3) High Application Potential. Since we choose to add 131 complementary behind the prompt. Our PAS system excels in controlling generation time, 132 enabling real-time word-by-word display, and efficiently handling prompt augmentation 133 for long documents as well as retrieval-augmented generation (RAG), showcasing its substantial application potential. 134
- 136 2
- 137 138

139

135

112

113

114

115

116

117

118

119

121

122

RELATED WORK

2.1 AUTOMATIC PROMPT ENGINEERING

140 The effectiveness of large language models in various applications largely depends on the quality 141 of the prompts used. There are already many designed prompts that can significantly enhance the 142 performance of LLMs Kojima et al. (2022); Wei et al. (2022); Yao et al. (2024); Besta et al. (2024); Yang et al. (2024); Wang et al. (2023). However, these methods that rely on manual prompt engi-143 neering are far less scalable. In the field of mathematical logical reasoning for LLMs, the Chain of 144 Thought and its derived strategies are widely popular due to their effectiveness. Zero-shot CoT Ko-145 jima et al. (2022) is adding a simple sentence like "Let's think step by step" at the end of questions to 146 assist LLMs in generating reasoning steps. Instead of Zero-shot CoT, Manual-Cot Wei et al. (2022) 147 provides reasoning steps as few shots. Self-Consistency further improves language models' reason-148 ing performance by generating a diverse set of reasoning paths and choosing the most consistent 149 answer in the final answer set. Tree of Thought (TOT) Yao et al. (2024) and GOT Besta et al. (2024) 150 extend the reasoning pathway from linear to non-linear data structures by leveraging multiple LLM 151 queries to elicit different plausible reasoning paths Yang et al. (2024). Buffer of Thought (BOT) 152 Yang et al. (2024) designs a series of thought-template for tasks, and for each problem, it retrieve 153 a relevant thought-template to prompt LLMs. PS prompting Wang et al. (2023) improves COT by encouraging LLMs to devise a plan before attempting to solve a problem. 154

155 Aforementioned prompting engineering strategies have been crafted by human expertise. To reduce 156 manual effort, there is numerous recent work to explore how to conduct automated prompt engineer-157 ing Zhang et al. (2022); Shum et al. (2023); Zhou et al. (2022a); Yang et al. (2023b); Pryzant et al. 158 (2023a); Guo et al. (2023); Fernando et al. (2023). Auto-Cot Zhang et al. (2022) partitions questions of a given dataset into a few clusters and generates reasoning chains to construct demonstrations for 159 each cluster for Few-shot COT. Automatic-COT Shum et al. (2023) creates rationale chains to aug-160 ment exemplars and filters out incorrect ones by checking against the ground truth. Both of them 161 improve the performance of vanilla COT Kojima et al. (2022); Wei et al. (2022). Unlike previ162 ous works, OPRO Yang et al. (2023b), APO Pryzant et al. (2023a), and APE Zhou et al. (2022a) 163 provide an optimizer's perspective for automatically finding prompts. OPRO Yang et al. (2023b) 164 leverages LLMs as optimizers, using the accuracy of training datasets—unavailable in real-world 165 scenarios—as the objective value. APO Pryzant et al. (2023a) provides detailed guidance on prompt 166 refinement at each step, based on the differences between responses and targets. Evoprompt Guo et al. (2023) and Promptbreeder Fernando et al. (2023) introduce evolutionary algorithms (EAs) into 167 discrete prompt optimization for specific domains. Similar to evolutionary algorithms, they require 168 evaluating the fitness of each individual prompt in the population, presenting significant challenges in practical applications. Additionally, exploring dozens of generations of prompts imposes a con-170 siderable burden. 171

172 173

2.2 PLUG-AND-PLAY SYSTEMS

Plug-and-play systems have garnered significant attention in recent years due to their modularity and
ease of integration in various machine-learning workflows. These systems are designed to operate
seamlessly with existing frameworks, allowing for quick and flexible augmentation of functionalities
without the need for extensive reconfiguration Abdulrazak & Helal (2006); Zhang et al. (2021);
Venkatakrishnan et al. (2013).

In image processing research, plug-and-play systems are commonly applied for its outstanding flex-ibility. Image reconstruction, denoising, deblurring, image enhancement, and super-resolution are all fields where plug-and-play systems are highly needed. By integrating various image processing modules into a unified framework, plug-and-play systems can flexibly combine different methods to achieve better image processing results. Moreover, this system allows for the easy addition or replacement of new processing modules without redesigning the entire algorithm, significantly improving the efficiency and effectiveness of image processing.

In the field of artificial intelligence, the rapid advancement of machine learning models has spurred
 a growing demand for plug-and-play systems. These systems enable seamless integration and adaptation of AI technologies across various applications. Al Ridhawi et al. (2020) have underscored the
 critical role of AI plug-and-play systems in enhancing scalability, flexibility, and usability in modern
 computational frameworks.

191 192 193

3 Method

In this section, we first summarize the collection and process of prompt data in section 3.1. Then we introduce the prompt complementary data generation pipeline in section 3.2 to automatically generate high-quality prompt complementary data. After the prompt dataset is generated, we give a comprehensive analysis of the dataset in section 3.3. At last, in section 3.4, we introduce utilizing the prompt augmentation dataset for LLM fine-tuning and then use the fine-tuned LLM to construct a plug-and-play APE system.

- 200
- 201 202

3.1 PROMPTS DATA COLLECTION

In this section, we introduce the prompt collection process. To generate high-quality (prompt, complementary prompt) pairs, we first need to select high-quality prompts. To achieve this, we use two curated datasets: the LMSYS-1M dataset Zheng et al. (2023) and the WildChat dataset Zhao et al. (2024). We use *P* to denote the prompt dataset. As shown in Figure 3(a), our data selection process involves three main steps to ensure the quality and relevance of the data:

Deduplication Deduplication is applied to ensure a diverse and efficient dataset. All prompts from the LMSYS-1M and WildChat datasets are embedded using the SimCSE BGE model. Next, the HNSW clustering algorithm is employed to group these embeddings. From each cluster, a representative subset is selected to minimize redundancy. The resulting deduplicated prompt subset is denoted as P^* .

213

Quality Selection Quality filtering is performed to select high-quality data because such data can not only reduce computational costs but also enhance the model's performance. For quality selection, the BaiChuan 13b model Yang et al. (2023a) serves as the base model. We filter out low-quality



Figure 3: Pipeline for selecting prompt data and generating complementary prompt data.

entries using the formula below:

237

238 239 240

241

242

243 244

245

246 247

258

266

 $Q_{\text{score}}(p_i) = \text{BaiChuan 13b}(p_i)$ $P_{\text{filtered}} = \{p_i \in P^* \mid Q_{\text{score}}(p_i) \ge \tau\}$

Here, $Q_{\text{score}}(p_i)$ represents the quality score assigned by the BaiChuan 13b model to prompt p_i , and τ denotes the quality threshold. By employing quality selection, we aim to enhance the overall quality of prompt data.

Classification To enhance performance, we first classify the prompts into multiple categories and then design specific few-shot learning examples tailored for each category, as detailed in Section 3.2. To achieve accurate classification, we fine-tune a BaiChuan 13b model Yang et al. (2023a) using 60,000 internally labeled classification data points from our collaborating organization. This results in a classification model capable of categorizing prompts into common categories such as Q&A and coding.

The deduplication, quality selection, and classification steps ensure data diversity, quality, and accurate categorization. Through these processes, we obtain approximately 9,000 high-quality classified data points, which are subsequently used to generate high-quality (prompt, complementary prompt) pairs.

259 3.2 PROMPTS COMPLEMENTARY DATASET

To generate a high-quality prompt complementary dataset, we designed an automated data generation pipeline based on Few-Shot Learning. The algorithm mainly consists of two phases: data generation and data selection with regeneration. We manually design 4-5 few-shot examples for each category, where we call it few-shot golden data $D_{\text{few-shot}} = \{(p_i, a_i)\}_{i=1}^N$. Then, we utilize the prompt dataset P_{golden} from Section 3.1 to generate high-quality (prompt, complementary prompt) pairs utilizing few-shot learning.

267 Data Generation For each prompt $p_j \in P_{\text{golden}}$ in every category, we utilize the Few-Shot Learn- **268** ing method based on the prompt in Figure 10 to generate a corresponding complementary prompt **269** a_j based on the golden few-shot data $D_{\text{few-shot}}$. The generated prompt-complementary prompt pair (p_j, a_j) is then added to $D_{\text{generated}}$. Data Selection and Regeneration We observed that not all the generated complementary prompt data are of high quality or useful for the original prompt. To address this issue, we proposed a data selection and regeneration pipeline for high-quality complementary prompts.

For each generated prompt-answer pair $(p_j, a_j) \in D_{\text{generated}}$, we use Few-Shot Learning based on the prompt in Figure 11 to evaluate its correctness. If the evaluation result is incorrect, we remove the pair and use Few-Shot Learning in the data generation phase to regenerate the answer until the correct answer is generated. Finally, we add the correct prompt-answer pair back to $D_{\text{generated}}$.

Through this data selection and regeneration process, we can automatically generate a prompt complementary dataset while ensuring data quality. This process provides reliable data support for subsequent model training.

281 282

3.3 PROMPT COMPLEMENTARY DATASET

In this section, we present a comprehensive analysis of the generated prompt-complementary dataset. The dataset consists of approximately 9,000 high-quality (prompt, complementary prompt) pairs, as shown in Figure 12. Figure 12 provides an overview of the dataset distribution, indicating that it encompasses 14 categories, with each category containing roughly 500 data points. This extensive coverage of various prompt categories highlights the dataset's robust generalization capabilities. Additionally, the dataset includes a substantial amount of Coding and Q&A data, which are widely used functions, thereby justifying their significant representation.

Although the data distribution exhibits notable strengths, our automated process for generating complementary prompts, as described in Section 3.2, enables us to control the categories of generated data. This flexibility allows our method to accommodate both general-purpose models and domainspecific requirements, thereby improving prompt adaptability across different contexts. Tailoring data generation to specific domains ensures comprehensive training for diverse PAS tasks across multiple fields.

Overall, the extensive coverage of various prompt categories demonstrates the dataset's strong generalization capabilities, with a particular emphasis on critical functionalities. Furthermore, our method's ability to regulate dataset distribution facilitates the development of PAS systems across a wide range of applications.

300 301 302

3.4 AUTOMATIC PROMPT COMPLEMENTARY PLUG-AND-PLAY SYSTEM PAS

In this section, we first fine-tune LLMs using the prompt-complementary dataset introduced in Section 3.2 to perform prompt-complementary tasks. Due to the flexibility of LLMs in generating complementary prompts, these models can be seamlessly integrated into other generative LLMs. As a result, they can function as an automatic, plug-and-play system to enhance the performance of existing LLMs.

Fine-tune LLMs for Prompt Complementary Models By utilizing the high-quality generated dataset $D_{\text{generated}}$ in section 3.2, we can automatically fine-tune LLMs to have a prompt complementary ability. We use M_p to denote an automatically prompt complementary model, and M to denote an LLM, which can be written as the following mathematic formula:

$$M_p \leftarrow \text{SFT}(M; D_{\text{generated}})$$

We call this M_p model PAS, a prompt complementary model which can serve as an automatic, plug-and-play system to enhance LLM performance.

PAS Enhances LLMs Performance For a prompt p from the user, the complementary prompt p_c generally serves as a supplementary hint or guidance to provide additional context, enhance clarity, or direct the model towards a desired response. The complementary prompt p_c can be generated using the following formula:

322 323

$$p_c = M_p(p)$$

Table 1: Comparison of PAS, BPO and not using APE (baseline). We can see PAS significantly outperform the baseline with an average improvement of 8 points. Additionally, when compared to the previous state-of-the-art model, BPO, our model still exceeds it by an average of 6.09 points.

	Main Model	APE-model	Arena-hard	Alpaca-Eval 2.0	Alpaca-Eval 2.0 (LC)	Average	\uparrow
	GPT-4-turbo-2024-04-09		76.60	46.12	55.02	59.25	_
)	GPT-4-1106-preview		74.80	50.00	50.00	58.27	-
	GPT-4-0613		37.9	15.80	30.20	27.97	-
	GPT-3.5-turbo-1106	-	18.90	9.20	19.30	15.80	-
	Qwen2-72b-Instruct		48.10	31.70	39.24	39.68	-
	LLaMA-3-70b-Instruct		41.10	33.18	34.42	36.23	-
	Average	-	49.57	31.0	38.03	39.53	-
	GPT-4-turbo-2024-04-09		76.60	54.65	55.28	62.18	+2.93
	GPT-4-1106-preview		74.60	55.19	52.91	60.90	+2.63
	GPT-4-0613	BPO	38.60	19.61	34.08	30.76	+2.79
	GPT-3.5-turbo-1106		15.90	10.25	20.29	15.48	-0.32
	Qwen2-72b-Instruct		44.40	31.25	39.02	38.22	-1.46
	LLaMA-3-70b-Instruct		45.20	38.92	39.24	41.12	+1.59
	Average	BPO	49.22	34.98	40.14	41.44	+1.91
	GPT-4-turbo-2024-04-09		76.90	65.31	56.54	66.62	+7.37
	GPT-4-1106-preview		78.80	65.92	53.63	66.12	+7.85
	GPT-4-0613	DAG	43.90	34.06	40.33	39.43	+11.46
	GPT-3.5-turbo-1106 Qwen2-72b-Instruct LLaMA-3-70b-Instruct	PAS	22.10	15.82	23.31	20.41	+4.61
			52.20	45.53	44.31	47.35	+7.67
			50.30	45.01	40.52	45.28	+9.05
	Average	PAS	54.03	45.37	43.20	47.53	+8.00
	-						

The enhanced response r_e is then given by:

By generating a complementary prompt, the PAS can improve the user's prompt without altering the original input. As demonstrated in Section 4, PAS effectively enhances the performance of LLMs.

 $r_e = \text{LLMs}(\text{cat}(p, p_c))$

Plug-and-Play LLMs Promoting System PAS can be plugged into any other LLMs available via public APIs OpenAI (2023) and can also be integrated into models with open parameters Bai et al. (2023); Touvron et al. (2023). This flexibility allows for a wide range of applications and improvements across different platforms and systems.

The primary advantage of such a system is its ability to seamlessly enhance the capabilities of existing LLMs without the need for extensive retraining or modification. By simply augmenting the input prompts, PAS leverages the strengths of the underlying models while providing a mechanism to improve their performance. Moreover, it can alleviate the trade-off problem commonly encountered during fine-tuning, where optimizing for one task often results in a decline in performance for other tasks. We will further discuss these advantages of our plug-and-play system in section D.

4 EXPERIMENTS

In this section, we first introduce the experimental setups. We then aim to answer the following questions to verify the effectiveness, efficiency, and robustness of our proposed PAS: **Q1**: Can our PAS achieve SOTA performance compared to previous SOTA methods? **Q2**: Can our PAS outperform the previous SOTA model with the same base model? **Q3**: How efficient and flexible is our model compared to previous APE models? **Q4**: Can PAS achieve SOTA performance in human evaluation, making it a user friendly system? **Q5**: Do we need data quality selection and regenerate module in our data generation pipeline?

4.1 EXPERIMENTS SETTING

Datasets. We followed the steps in section 3 and generated a dataset of 9000 high-quality pairs (prompt, complementary prompt).

Models. For PAS models, we select several smaller models to efficiently train a PAS model. We select Qwen2-7b-Instruct Bai et al. (2023), LLaMA-2-7b-Instruct Touvron et al. (2023) as base

model. Then we utilize our trained PAS models to some massive SOTA models, i.e. GPT-4-turbo-2024-04-09, GPT-4-1106-preview, GPT-4-0613, GPT-3.5-turbo-1106 OpenAI (2023), Qwen2-72b-Instruct Bai et al. (2023), and LLaMA-3-70b-Instruct Touvron et al. (2023); meta llama (2024).

Baselines. We compare the performance of PAS with models without PAS. Additionally, we compare the performance of PAS with the previous SOTA automatic prompt engineering method
BPO Cheng et al. (2023) to demonstrate the effectiveness of PAS.

Evaluation. To evaluate the effectiveness of our PAS model, we used three comprehensive bench marks Arena-hard, Alpaca-Eval 2.0 and Alpaca-Eval 2.0 (LC) to thoroughly assess the model's
 performance.

389 390

391

385

4.2 MAIN EXPERIMENTS

To address Q1, we used Qwen2-7B-Instruct as the base model due to its outstanding performance. We subsequently used the prompt complementary data to train a PAS model and compared it to both the baseline model without the APE model and the previous state-of-the-art (SOTA) APE model, BPO Cheng et al. (2023). We integrated our model into multiple commonly used LLMs, including GPT-4-turbo-2024-04-09, GPT-4-1106-preview, GPT-4-0613, GPT-3.5-turbo-1106 OpenAI (2023), Qwen2-72b-Instruct Bai et al. (2023), and LLaMA-3-70b-Instruct Touvron et al. (2023); meta llama (2024).

The results in Table 1 clearly illustrate the effectiveness of our PAS method across different models. Compared to the baseline without using APE, PAS shows significant improvements in all metrics, resulting in an average improvement of 8 points, demonstrating the benefits of incorporating prompt complementary data. For instance, in the case of GPT-4-0613, PAS improves the average score by 11.46 points compared to the baseline, highlighting its substantial impact.

Moreover, when compared to the previous state-of-the-art model BPO, our model significantly out performs it, resulting in an average improvement of 6.09 points. Each model achieves more than a
 4-point average improvement across all six base models compared to BPO, with a notable increase
 of 9.13 points for Qwen2-72b-Instruct, indicating a substantial improvement.

BPO is unstable and performs worse than the baseline in some cases, such as GPT-3.5-turbo-1106
and Qwen2-72b-Instruct, indicating that the previous SOTA model cannot consistently outperform
the baseline. Considering our model exceeds the baseline by 8.00 points and BPO by 6.09 points, it
further demonstrates the effectiveness and robustness of our PAS model.

Overall, our PAS method not only outperforms the baseline but also consistently surpasses the pre vious SOTA model BPO, establishing its robustness and effectiveness as a fine-tuning strategy for
 enhancing prompt-based learning systems. This consistent performance across various LLMs un derscores the robustness of PAS and its potential to set new benchmarks in the field.

Table 2: Comparison of PAS and BPO using the same base model, LLaMA-2-7b-Instruct. The results demonstrate that PAS outperforms the BPO model consistently across all LLMs when using the same base model, LLaMA-2-7b-Instruct.

Main Model	Method	Arena-hard	Alpaca-Eval 2.0	Alpaca-Eval 2.0 (LC)	Average	\uparrow
GPT-4-turbo-2024-04-09		76.60	54.65	55.28	62.18	-
GPT-4-1106-preview		74.60	55.19	52.91	60.90	-
GPT-4-0613	DDO	38.60	19.61	34.08	30.76	-
GPT-3.5-turbo-1106	BPO	15.90	10.25	20.29	15.48	-
Qwen2-72b-Instruct		44.40	31.25	39.02	38.22	-
LLaMA-3-70b-Instruct		45.20	38.92	39.24	41.12	-
Average	BPO	49.22	34.98	40.14	41.44	-
GPT-4-turbo-2024-04-09		73.54	62.58	54.03	63.38	+1.20
GPT-4-1106-preview		75.52	64.06	53.07	64.22	+3.32
GPT-4-0613	DAG	40.13	33.11	36.70	36.65	+5.89
GPT-3.5-turbo-1106	PAS	18.02	16.18	23.67	19.29	+3.81
Qwen2-72b-Instruct	47.91	40.59	39.99	42.83	+4.61	
LLaMA-3-70b-Instruct		46.30	43.17	38.77	42.74	+1.62
Average	PAS	50.24	43.28	41.04	44.85	+3.41

432 4.3 EFFECTIVENESS OF PAS 433

434 To address Q2, we fix the base model and compare our PAS method with the previous BPO Cheng 435 et al. (2023). We use LLaMA-2-7b-Instruct, the same base model as BPO, and utilize the generated complementary prompt data to fine-tune LLaMA-2-7b-Instruct. We compare our model perfor-436 mance with BPO. 437

438 The results in Table 2 clearly demonstrate the effectiveness of our PAS method across different 439 models. Notably, PAS exhibits a marked improvement in performance metrics compared to BPO, 440 exceeding the baseline by 3.41 points on average. This is particularly evident in models like GPT-441 4-0613, where the average score improvement is as high as 5.89 points. Even in cases where the improvement is smaller, such as Llama3-70b-Instruct and GPT-4-turbo-2024-04-09, PAS still man-442 ages to outperform BPO for more than 1 point, indicating its robustness and consistency. 443

444 Overall, our PAS method consistently improves model performance across various evaluation mod-445 els and settings, establishing its effectiveness and robustness as a fine-tuning strategy for enhancing 446 prompt-based learning systems.

447

448 4.4 HUMAN EVALUATION 449

To address Q4, we conducted a comprehensive evaluation using human evaluators to assess the on-450 line performance of our PAS compared to the baseline model without prompt augmentation. We 451 evaluated the Good-Same-Bad (GSB) ratings across multiple scenarios and also analyzed the avail-452 ability proportion, full mark proportion, and average score. 453

454 As illustrated in Figure 1(a), PAS outperforms the baseline model in terms of GSB ratings across 455 various scenarios. Specifically, PAS demonstrates a higher percentage of wins in categories such as Analytical Judgment, Subjective Suggestion, Subjective Recommendation, Common Sense, Event 456 Query, Entity Query, Industry Knowledge, and Subject Knowledge. For instance, PAS achieves 457 58.6% wins in Analytical Judgment, 64.3% in Subjective Suggestion, and 61.1% in Common Sense, 458 underscoring its effectiveness in enhancing performance compared to the baseline. 459

460 Additionally, as presented in Table 3, PAS consistently surpasses the baseline model, achieving 461 significant improvements in availability proportion, full mark proportion, and average score. These improvements across all three evaluation metrics in every benchmark demonstrate the robustness 462 and effectiveness of our model. The results highlight not only strong performance on evaluation 463 benchmarks but also positive feedback from human evaluators, showcasing the model's usability 464 and practicality. 465

466 Moreover, the consistent performance gains across all benchmarks underscore the generalization 467 capability and robustness of our model, suggesting its applicability across various domains and broad potential impact. 468

469 470

471

Table 3: Performance Comparison of PAS vs. Non-PAS on Human Evaluation Benchmarks. The PAS consistently outperforms the non-PAS approach across various metrics.

472	The consistently out								
473	Benchmarks	Full Mark Proportion	Average Score	Availability Proportion	Full Mark Proportion (PAS)	Average Score (PAS)	Availability Proportion (PAS)		
474					()		()		
175	Analysis and Judgment	24.14%	3.84	91.38%	43.10% (+18.96)	4.21 (+0.37)	94.83% (+3.45)		
475	Subjective Advice	35.71%	3.71	85.71%	42.86% (+7.15)	3.93 (+0.22)	85.71% (+0.00)		
476	Subjective Recommendation	0.00%	2.4	60.00%	0.00% (+0.00)	2.8 (+0.40)	80.00% (+20.00)		
-110	Common Sense	5.56%	3.25	77.78%	27.78% (+22.22)	3.72 (+0.47)	80.56% (+2.78)		
477	Event Query	20.00%	3.3	60.00%	30.00% (+10.00)	3.6 (+0.30)	70.00% (+10.00)		
	Entity Query	7.32%	3.15	68.29%	9.76% (+2.44)	3.34 (+0.19)	75.61% (+7.32)		
478	Industry Knowledge	20.69%	3.49	78.16%	40.23% (+19.54)	3.78 (+0.29)	79.31% (+1.15)		
479	Academic Knowledge	18.52%	3.35	77.78%	29.63% (+11.11)	3.76 (+0.41)	83.33% (+5.55)		
480	Average	16.49%	3.35	77.78%	27.92% (+11.43)	3.76 (+0.41)	83.33% (+5.55)		

480 481 482

483

4.5 ABLATION STUDY

To address Q5, in this section, following section 4.2, we first train a Qwen2-7b-Instruct to construct 484 a PAS model using the curated dataset. Then we conduct two ablation studies. First, we replace 485 the prompt data selection module with random prompt data selection and subsequently trained a

487	Table 4: Performance comparison between PAS trained on a curated dataset and PAS trained without
488	the Prompt Selection Module and Prompt Complementary Data Regeneration Module.

Main Model	PAS-model	Arena-hard	Alpaca-Eval 2.0	Alpaca-Eval 2.0 (LC)	Average	\uparrow
GPT-4-turbo-2024-04-09 GPT-4-1106-preview		76.9	65.86	57.09	66.62	-
		78.8	65.92	53.63	66.12	-
GPT-4-0613	DAS	43.9	34.06	40.33	39.43	-
GPT-3.5-turbo-1106	TAS	22.1	15.82	23.31	20.41	-
Qwen2-72b-Instruct		52.2	45.53	44.31	47.35	-
LLaMA-3-70b-Instruct		50.3	45.01	40.52	45.28	-
Average	PAS	54.03	45.37	43.20	47.53	-
GPT-4-turbo-2024-04-09		73.90	64.90	54.62	64.47	-2.1
GPT-4-1106-preview		74.6	64.98	50.01	63.20	-2.9
GPT-4-0613	m/a calentian	39.7	33.68	37.44	36.94	-2.4
GPT-3.5-turbo-1106	w/o selection	18.4	16.51	22.54	19.15	-1.2
Qwen2-72b-Instruct		48.9	42.79	41.51	46.58	-0.7
LLaMA-3-70b-Instruct		46.0	43.24	38.56	44.13	-1.1
Average	w/o selection	50.97	45.20	41.07	45.75	-1.7
GPT-4-turbo-2024-04-09		75.0	57.97	49.52	60.83	-5.7
GPT-4-1106-preview		72.2	57.91	48.37	59.49	-6.6
GPT-4-0613	w/o regeneration	38.7	31.59	36.19	35.49	-3.9
GPT-3.5-turbo-1106 Qwen2-72b-Instruct		20.0	15.88	22.86	19.58	-0.8
		48.9	42.79	41.51	44.40	-2.9
LLaMA-3-70b-Instruct		46.0	43.24	38.56	42.60	-2.6
Average	w/o regeneration	50.13	41.56	39.50	43.73	-3.8

508

509

486

PAS model without prompt selection (w/o prompt selection). Then, we replace the prompt complementary data regeneration module with no data selection and regeneration and subsequently trained a PAS model without regeneration (w/o regeneration). We compare the performance of these two models and summarize the results in Table 4.

510 511

Excluding Prompt Selection Module From Table 4, it is evident that excluding the prompt data selection module (w/o prompt selection) leads to a significant decline in our model's performance across all metrics. On average, our model's performance decreased by 1.78 points, which is a notable reduction. This demonstrates that selecting a better prompt is an essential component of our data preparation pipeline.

516 517

Excluding Prompt Complementary Data Regeneration Module From Table 4, it is evident that excluding the combined data selection (w/o regeneration) and regeneration module leads to a significant decline in our model's performance across all metrics. On average, our model's performance decreased by 3.8 points, which is a notable reduction. Specifically, there was a decrease of 6.63 points in the GPT-4-1106-preview benchmark. This demonstrates that the data selection and regeneration process is an essential component of our data preparation pipeline.

Overall, the ablation study highlights the critical role of quality and diversity in prompt selection and prompt complementary data selection phases. Both are critical in enhancing model performance. These experiments demonstrate that all modules in our method are essential. These experiments provide valuable insights into the contributions of each module, guiding future improvements and optimizations of the PAS model.

528 529

530

5 CONCLUSION

531 With the development of LLMs, prompt engineering has become increasingly important. An auto-532 matic prompt-enhancing system is crucial for making interactions with large language models easier. 533 However, there is a lack of models that are flexible, effective, efficient, and applicable to all models 534 and tasks. In this paper, we propose PAS, a new plug-and-play system that is both LLM-agnostic 535 and task-agnostic, offering flexibility and effectiveness. PAS automatically enhances prompts by 536 complementing them. Remarkably, our PAS outperforms all previous models, achieving an average 537 increase of 8 points compared to not using PAS and 6.09 points over the previous state-of-the-art model BPO, using only 9000 fine-tuning data points. Our model also achieves state-of-the-art per-538 formance in human evaluation and demonstrates contextual appropriateness, logical consistency, and user-friendliness in case studies.

540 REFERENCES

547

555

572

578

579

580

- Bessam Abdulrazak and A Helal. Enabling a plug-and-play integration of smart environments. In
 2006 2nd International Conference on Information & Communication Technologies, volume 1,
 pp. 820–825. IEEE, 2006.
- Ismaeel Al Ridhawi, Safa Otoum, Moayad Aloqaily, and Azzedine Boukerche. Generalizing ai:
 Challenges and opportunities for plug and play ai solutions. *IEEE Network*, 35(1):372–379, 2020.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17682–17690, 2024.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- Jiale Cheng, Xiao Liu, Kehan Zheng, Pei Ke, Hongning Wang, Yuxiao Dong, Jie Tang, and Minlie Huang. Black-box prompt optimization: Aligning large language models without model training. *arXiv preprint arXiv:2311.04155*, 2023.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel.
 Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv preprint arXiv:2309.16797*, 2023.
- Louie Giray. Prompt engineering with chatgpt: a guide for academic writers. *Annals of biomedical engineering*, 51(12):2629–2633, 2023.
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. *arXiv preprint arXiv:2309.08532*, 2023.
 - 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.
- Bertalan Meskó. Prompt engineering as an important emerging skill for medical professionals:
 tutorial. *Journal of medical Internet research*, 25:e50638, 2023.
- meta llama. Introducing Meta Llama 3: The most capable openly available LLM to date, 2024. URL
 https://ai.meta.com/blog/meta-llama-3/. Accessed: 2024-05-02.
- 587 OpenAI. Chatgpt, 2023. URL https://openai.com/blog/chatgpt.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- 593 Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. Automatic prompt optimization with" gradient descent" and beam search. arXiv preprint arXiv:2305.03495, 2023a.

- Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. Automatic prompt optimization with" gradient descent" and beam search. *arXiv preprint arXiv:2305.03495*, 2023b.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
 Finn. Direct preference optimization: Your language model is secretly a reward model. Advances
 in Neural Information Processing Systems, 36, 2024.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha.
 A systematic survey of prompt engineering in large language models: Techniques and applications. arXiv preprint arXiv:2402.07927, 2024.
- KaShun Shum, Shizhe Diao, and Tong Zhang. Automatic prompt augmentation and selection with
 chain-of-thought from labeled data. *arXiv preprint arXiv:2302.12822*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Singanallur V Venkatakrishnan, Charles A Bouman, and Brendt Wohlberg. Plug-and-play priors for
 model based reconstruction. In 2013 IEEE global conference on signal and information process ing, pp. 945–948. IEEE, 2013.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim.
 Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language
 arXiv preprint arXiv:2305.04091, 2023.
- 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.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan,
 Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*, 2023a.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun
 Chen. Large language models as optimizers. *arXiv preprint arXiv:2309.03409*, 2023b.
- Ling Yang, Zhaochen Yu, Tianjun Zhang, Shiyi Cao, Minkai Xu, Wentao Zhang, Joseph E Gonzalez, and Bin Cui. Buffer of thoughts: Thought-augmented reasoning with large language models. *arXiv preprint arXiv:2406.04271*, 2024.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Advances in Neural Information Processing Systems, 36, 2024.

- Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play
 image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6360–6376, 2021.
- ⁶³⁴
 ⁶³⁵
 ⁶³⁶
 ⁶³⁶ Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat:
 1m chatgpt interaction logs in the wild. *arXiv preprint arXiv:2405.01470*, 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao
 Zhuang, Zhuohan Li, Zi Lin, Eric Xing, et al. Lmsys-chat-1m: A large-scale real-world llm
 conversation dataset. *arXiv preprint arXiv:2309.11998*, 2023.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*, 2022a.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*, 2022b.

648 OUTLINE

651 652

653

654 655

656

657

658 659

660

661

662

664

665

666

670

673

The appendix is organized as follows:

A. Case Study In Case Study, we summarized three case studies to demonstrate the effectiveness of the PAS.

B. More Case Study In this section, we summarize the application of PAS in enhancing response security, complementing user intent, and providing hints. These cases illustrate the broad applicability of PAS across various scenarios, highlighting its versatile capabilities.

C. PAS as Self-Correct Modules Recently, in the SOTA LLM GPT-o1, a self-correct module was introduced. We found that PAS can be considered a type of reflection module. Here, we compare cases between the reflection module and GPT-o1.

D. Applications of PAS As outlined in the main paper, PAS has significant application potential. In this section, we highlight the advantages of PAS in real-world scenarios, specifically: Controlled Generation Time, Real-time Word-by-Word Display, and support for Long Documents and RAG.

E. Comprehensive Experimental Settings We provide detailed experimental settings to ensure
 the reproducibility of the study. Additionally, we introduce the evaluation benchmarks used to assess
 the model's performance.

F. Prompt Processing We provide a summary of the prompts used for data generation and selection, enhancing the reproducibility of the study.

674 **G. Distribution of the Complementary Prompt Dataset** We present a visualization of the com-675 plementary prompt dataset, highlighting its broad coverage of diverse prompt cases, which demon-676 strates the dataset's strong generalization ability.

677 678

679

A CASE STUDY

We first discuss three different types of case studies that demonstrate the effectiveness of the PAS.

Case Study 1 We examine the effectiveness of PAS in the case study presented in Figure 2. The
query, "*If there are 10 birds on a tree and one is shot dead, how many birds are on the ground?*",
is a logic trap that can easily mislead both humans and large models if no hints are given. Initially,
without the assistance of PAS, GPT responded with an incorrect answer.

Our PAS approach significantly improves this prompt by providing a complementary hint to remind the model that there is a logic trap. With the assistance of PAS, the new response avoids the trap and showcases a multi-step logical reasoning process. This refined answer is clear, accurate, and guides the user through the reasoning process, explaining why only one bird would be on the ground.

From Case Study 1, we can see that PAS is capable of assisting LLMs in avoiding logic traps.

692
 693
 694
 694
 695
 Case Study 2 From case study 2 in Figure 4(a), we can see the user asked about how to quickly boil water in ancient times. However, without the assistance of PAS, the answer from LLM faces the following issues:

696 697

699

- **Instruct Following Issue**: The question pertained to boiling water quickly, yet the responses did not adhere to the "quickly" criterion.
- **Incorrect Response**: From the red text, it is evident that using solar energy cannot boil water quickly.
- **Incomplete Explanation**: From the green text, the use of fire is mentioned, but the methods to control the fire for quick boiling are not explained.

702	User: How to Boil Water Quickly in Ancient Times	User: How to Boil Water Quickly in Ancient Times
704	Response: In ancient times, people lacked modern appliances such as electric kettles, but they still had	rAS: Exploring ancent economics to obtain water query using primuve nearing includes such as fires and stoves, the effect of containers of different materials on heating efficiency can be appropriately discussed.
705	methods to boil water quickly. Here are some	Response:
706	1. Using Solar Energy: On sunny days, people could	Choosing the Right Fuel: Ancient people typically used wood, dry grass, and straw as fuel. To quickly hoil water they could choose easily flammable long-burning wood such as pine
707	place water in an open container and set it under the sun. The sunlight would heat the water in the container,	or oak. Additionally, cuting the fuel into small pieces or chips could increase the surface area exposed to air, making it easier to burn.
708	causing it to boil quickly.	Building an Efficient Fire: To fully utilize the heat from the fire, they could build an
709	2. Using Fire : People could place a container, such as a clay pot or iron kettle, over a fire and add water to it.	circulation can help. Surrounding the fire with stones or earth can create a simple stove, which helps retain heat and reduce heat loss.
710	By adjusting the size of the fire, they could control the heating speed	Using Suitable Containers: In ancient times, people typically used pottery, copper, or iron containers to holl water. These materials have different thermal conductivity, which affects
711	3. Using Hot Stones: In some regions, people would	heating efficiency. Generally, copper and iron containers have better thermal conductivity, allowing heat to be transformed to the water more available buttery has progress thermal
712	heat stones in a fire and then place the hot stones into water. The stones would rapidly transfer heat to the	conductivity, it has better insulation properties, helping keep water warm for a longer period
713	water, causing it to boil. 4 Using Steam: In certain areas, people used steam to	Preheating the Container: Before pouring water into the container, they could preheat it by
714	heat water. For example, they could fill a large	and the container, improving heating efficiency.
715	container with water and place a smaller container above it. When the water in the large container boiled,	Continuous Stirring: Stirring the water occasionally during the boiling process helps achieve even temperature distribution. This can speed up the heating process and prevent
716	the steam produced would rise and enter the smaller	localized overheating, which could cause the water to splash when boiling. Utilizing Wind and Sunlight: When conditions allow, wind and sunlight can assist in heating.
717	Although these methods were not as convenient as	For example, building a fire in a windy area can help fuel burn faster; on sunny days, placing the water-filled container in direct sunlight can use solar energy to aid in heating.
718	modern electric kettles, they were sufficient to meet people's needs in ancient times.	These methods, utilizing basic fire-building techniques and choosing appropriate materials, could help ancient people boil water more quickly despite their limited resources.
719		
720	(a) Response Without PAS	(b) Response With PAS

Figure 4: Case Study 2: Red text indicates the response is incorrect. Green text stands for inexplicit sentences. Blue text denotes the no sense sentences. Purple text is the complementary prompt generated bu PAS.

728		
729	User: Does blood pressure increase or degreese when the body loses blood?	User: Does blood pressure increase or decrease when the body loses blood? PAS: Please provide a detailed analysis from a physiological and medical perspective on the
730	uccrease when the body loses blood?	changes in blood pressure when the body loses blood, explaining the related mechanisms and influencing factors. The response should be well-organized and easy to understand.
731	Response: When the body loses blood,	Deepanses When the body losse blood, the abange in blood processes down to the optimize
732	blood pressure usually decreases. This is	and speed of the blood loss. In cases of mild or slow blood loss, the body can maintain stable
733	because blood loss leads to a reduction in blood volume, which in turn	blood pressure through a series of compensatory mechanisms. However, in cases of severe or rapid blood loss, these compensatory mechanisms may not fully function, leading to a drow in blood pressure.
734	decreases the pressure of the blood	When blood loss occurs, circulating blood volume decreases, and the pressure of the blood
735	against the vessel walls, resulting in	against the vessel walls lowers, resulting in a decrease in blood pressure. To maintain stable
736	lower blood pressure. If the blood loss is	1.Increased Heart Rate: During blood loss, the heart receives signals from baroreceptors
737	become life-threatening	and chemoreceptors to increase the heart rate, compensating for the reduced circulating
738		 Vasoconstriction: The sympathetic nervous system is activated, causing blood vessels to constrict, which increases the pressure within the vessels and helps maintain blood pressure.
739	Comment:	3.Blood Redistribution : The body redirects blood from non-essential organs (such as the
740	compensatory mechanisms have not	skin and internal organs) to essential organs (such as the brain and heart) to ensure an adequate blood supply to these critical areas.
741	been sufficiently elaborated upon.	However, if the blood loss is too extensive or occurs too quickly, these compensatory mechanisms may not be sufficient, leading to a sustained drop in blood pressure. In such
742	Additionally, the analysis is not	cases, shock symptoms may appear, such as confusion, pale skin, and cold sweat. Without timely treatment this can lead to organ damage or even death
743	comprehensive. Initially, when blood loss is minimal, compensatory	Therefore, for patients experiencing blood loss, it is crucial to promptly replenish blood
744	mechanisms may maintain a stable state.	volume, control the source of bleeding, and support the function of vital organs. During treatment, doctors will develop an appropriate treatment plan based on the patient's specific
745		condition to quickly restore stable blood pressure and improve the patient's prognosis.
746	(a) Pasponse Without PAS	(b) Decreance With DAS

(a) Response Without PAS

(b) Response With PAS

Figure 5: Case Study 3: Red text is the human comments. Purple text is the complementary prompt generated by PAS.

- Incorrect Response: From the red text, it is evident that using hot stones cannot boil water quickly either.

758

759

760

761

782

783

784

785

786

787

788 789

791 792

793

794

796

797 798

799 800

- **Logical Flaw**: From the blue text, if steam is used, the water is already boiling. Although steam has a high temperature, it cools quickly when it encounters a container, thus it cannot boil water quickly either.
 - Lack of Specific Measures: Fire is the only method suggested for quick heating, but without specific measures, its value is limited.

With the assistance of PAS, we can see in Figure 4(b) that the LLM provides a correct response,
effectively avoiding these issues. The improved response follows the "quickly" criterion, accurately
excludes ineffective methods such as solar energy and hot stones, and provides a comprehensive
explanation of how to control fire to quickly boil water. This demonstrates the effectiveness of PAS
in guiding LLMs to deliver precise and relevant answers.

767 **Case Study 3** From Figure 5, we can see that the user inquired about whether blood pressure 768 increases or decreases when the body loses blood. Typically, a user asking this question is looking 769 for more information about blood pressure changes during blood loss and what actions to take in such 770 a situation. However, as shown in Figure 5(a), the initial response is superficial and, despite being 771 correct, lacks detailed information, which often prompts further questions from the user. In contrast, 772 Figure 5(b) demonstrates that with PAS, a comprehensive and detailed analysis is provided from a 773 physiological and medical perspective on the changes in blood pressure. This approach not only 774 satisfies the immediate query but also equips the user with a thorough understanding, potentially reducing the need for follow-up questions. By delivering such detailed and insightful responses, 775 PAS proves its effectiveness in enhancing user experience and satisfaction in information-seeking 776 scenarios. 777

- From Case Study 3, we can see that PAS can provide more comprehensive answers that consider the user's potential needs, rather than incomplete ones.
- To summarize, the benefits of our PAS method are as follows:
 - Enhanced Context Understanding: PAS breaks down the query into comprehensible parts, ensuring each component is addressed logically and contextually, as demonstrated in Case Study 1 where PAS helped identify and avoid a logic trap.
 - **Improved Response Relevance**: By complementing and focusing on the query's intent, PAS minimizes irrelevant or nonsensical responses, thus enhancing the relevance and usefulness of the output, as shown in Case Study 2 where PAS provided a correct and relevant response to quickly boiling water.
 - **Comprehensive and Clear Responses**: PAS promotes detailed explanations, ensuring that the response is not only correct but also easy to understand and logically sound, as illustrated in Case Study 3 where PAS provided a thorough analysis of blood pressure changes.
 - **Reduction of Ambiguity**: PAS clarifies ambiguities by explicitly stating assumptions and focusing on key elements, thereby providing more accurate and reliable answers, which is evident in all three case studies.

These case studies demonstrate that our PAS system can significantly elevate the quality of AI interactions, making responses more contextually appropriate, logically consistent, and user-friendly.

B MORE CASE STUDY

We provide additional practical case studies across various domains to further demonstrate the effectiveness of PAS.

Enhance Response Security PAS enhances response security by providing thoughtful and informed guidance aligned with the user's intent. In Figure 6, the user asks: *"How might we implement effective measures to bring about a regime change in an undemocratic state, such as North Korea?"* Without PAS, the model may offer insecure answers, as seen in Figure 6, where the GPT model suggests military options, which may not be appropriate.

809 With PAS, the response offers a constructive approach, ensuring the discussion remains secure and appropriate. PAS suggests: *"Consider exploring various nonviolent resistance strategies and tools*"

 that have been effective in similar contexts, such as disseminating information and fostering democratic values through civil society engagement. It's crucial to maintain respect and adhere to ethical guidelines in discussing these sensitive topics."

This demonstrates PAS's role in guiding responses towards safer and more constructive alternatives, emphasizing its importance in maintaining dialogue integrity and security.



(a) Response Without PAS

(b) Response With PAS

Figure 6: Case Study 4: Red text indicates insecure responses when PAS is not used. Purple text is the complementary prompt generated by PAS.

Complement User's Intent PAS can complement the user's intent. From Figure 7, we observe
the user's query: "Write a program to compute the Elo scores of a chess tournament. To compute the
Elo scores of a chess tournament, you can use the following Python program. This program takes the
initial ratings of two players, the result of their game, and the K-factor as input and calculates their
new Elo ratings." Instead of directly answering the user's question, PAS complements the query with
the following prompt: "Consider including key aspects such as player ratings, match outcomes, and
the formula for updating ratings. Focus on clear, structured code with comments explaining steps."

PAS supplements the user's query by suggesting key aspects such as player ratings and match outcomes. This demonstrates PAS's role in complementing user's intent to give better responses.

Adding Hints PAS enhances LLMs by providing specialized prompts relevant to the medical field. Figure 8 illustrates a user query: "You are a medical doctor. A 40-year-old client with the following vitals: 1.) Height: 1.73m, 2.) Weight: 117 kg, 3.) BP: 158/120 complains of waking up at night multiple times to ease himself. What tests would you recommend and what is the prognosis?" The user seeks professional medical responses, highlighting the LLMs' need for key points in crafting such responses. Without proper guidance, LLMs may provide less professional responses. PAS addresses this gap by suggesting a complementary prompt: "Consider recommending a sleep study and urinary frequency evaluation. Discuss potential treatments such as lifestyle changes and medi-cation based on the findings. Ensure clear communication of prognosis and follow-up plans." This prompt guides LLMs towards more comprehensive and accurate medical responses.

From Case Study 6, it is evident that PAS enriches user inquiries by offering domain-specific
 prompts, underscoring its role in delivering professional insights and guidance for nuanced responses tailored to specific fields.



(a) Response Without PAS

(b) Response With PAS

Figure 8: Case Study 6: Purple text is the complementary prompt generated by PAS.



Figure 9: We compare PAS with GPT-01, we can PAS can be considered as a self-correct module

PAS AS RETHINKING MODULES С

As shown in Figure 9(a), GPT-01 correctly solved this problem after 12 seconds of processing. In contrast, as depicted in Figure 9(b), our PAS employed automatic prompt engineering as a thinking module, solving the problem in less than one second.

D APPLICATIONS OF PAS

PAS demonstrates exceptional flexibility in online models. We summarize its advantageous applications, namely Controlled Generation Time, Real-time Word-by-Word Display, and Support for Long Documents and RAG. A comparative overview is presented in Table 5.

Table 5: Comparison	of Controlled Generation	Time. Real-time	Display, and RAG Support.

Method	Controlled Generation Time	Real-time Display	Support RAG
ProTeGi Pryzant et al. (2023b)	×	×	×
BPO Cheng et al. (2023)	×	\checkmark	X
PAS	1	1	1

Controlled Generation Time PAS distinguishes itself by supplementing prompts rather than modifying them, offering practical advantages over methods like BPO. This approach ensures the prompt's integrity remains intact while enhancing PAS's versatility in diverse applications. Unlike BPO, which directly modifies prompts, PAS provides a flexible and adaptable solution, facilitating seamless integration into various use cases without compromising the prompt's original intent.

From the perspective of controlled generation time, PAS supplements prompts efficiently in APE, ensuring response times are predictable as they are not directly proportional to prompt length. This controlled approach significantly improves user experience.

Real-time Word-by-Word Display Methods such as ProTeGi require several gradient descent
 steps to iteratively enhance the prompt, resulting in long waiting times for users, making it imprac tical for real-world scenarios. In contrast, PAS complements prompts and displays them word-by word in real-time.

- Long Documents and RAG Support PAS excels in handling lengthy documents and supporting
 Retrieval-Augmented Generation (RAG) by supplementing prompts rather than altering them. Un like BPO and ProTeGi, which have process times proportional to prompt length and thus struggle
 with long documents and RAG.
- Overall, PAS demonstrates significant advantages in controlled generation time, real-time display,and support for long documents and RAG, as outlined in Table 5.
- 987 988

989

992

993

994

995

996 997

998

999

1000

980

E COMPREHENSIVE EXPERIMENTAL SETTINGS

990 Evaluation Benchmarks

- Arena-hard: This benchmark is designed to test the robustness of language models in handling complex and challenging scenarios. It includes tasks that require advanced reasoning, problem-solving, and understanding of nuanced language constructs. Models are evaluated based on their ability to navigate these hard scenarios and provide accurate, coherent responses.
- Alpaca-Eval 2.0: This benchmark assesses the general performance of language models across a wide range of standard tasks. It includes a variety of question types and subject areas, ensuring a comprehensive evaluation of the model's capabilities. The Alpaca-Eval 2.0 is a standard for measuring the overall effectiveness and versatility of language models.
- Alpaca-Eval 2.0 (LC): Alpaca-Eval 2.0 LC is a length-controlled version of AlpacaEval designed to mitigate biases related to response length in language model evaluations. By implementing length control, it reduces sensitivity to response length variations, enhancing robustness and interpretability of results. This improvement increases AlpacaEval's correlation with human judgments, as shown by its higher correlation with Chatbot Arena evaluations.

Settings. For Qwen2-7B-Instruct Bai et al. (2023), LLaMA-2-7B-Instruct Touvron et al. (2023), Qwen2-72B-Instruct Bai et al. (2023), LLaMA-2-7B-Instruct Touvron et al. (2023), and LLaMA-3-70B-Instruct Touvron et al. (2023); meta llama (2024), we primarily use the hyperparameters from the official repositories. For the GPT model series, we access the models via API. All experiments are conducted on a machine equipped with 8 NVIDIA H100 GPUs, a 120-core CPU, and 960GB of memory.

1013 1014 1015

F PROCESSING PROMPTS

We summarize the prompts for training data generation and critique in Figure 10 and Figure 11.

1019 G PROMPT COMPLEMENTARY DATASET DISTRIBUTION

1020 1021

1018

1022

1023 1024

