

PIAST: Rapid Prompting with In-context Augmentation for Scarce Training data

Paweł Batorski

Paul Swoboda

Heinrich Heine Universität Düsseldorf

{pawel.batorski, paul.swoboda}@hhu.de

Abstract

LLMs are highly sensitive to prompt design, but handcrafting effective prompts is difficult and often requires intricate crafting of few-shot examples. We propose a fast automatic prompt construction algorithm that augments human instructions by generating a small set of few shot examples. Our method iteratively replaces/drops/keeps few-shot examples using Monte Carlo Shapley estimation of example utility. For faster execution, we use aggressive subsampling and a replay buffer for faster evaluations. Our method can be run using different compute time budgets. Under a limited budget, it outperforms prior automatic prompting methods on text simplification and mathematical reasoning (GSM8K, DeepMath, Math500), while achieving second-best results on classification and summarization and third-best on MedQA. With an extended, yet still modest budget, PIAST sets a new state of the art among automatic prompting methods on classification, simplification, GSM8K, DeepMath, and Math500. Overall, our results suggest that optimizing in-context examples, rather than exhaustively searching over instruction rewrites is the dominant lever for fast and data-efficient prompt engineering. Our code is available at: <https://github.com/Batorskq/PIAST>.

1 Introduction

Automatic prompt engineering has emerged as a practical way to adapt LLMs without gradient updates. However, many existing methods are impractical in time and data constrained settings: (i) some require hours of compute to explore a large prompt search space, and (ii) they rely on sizeable training sets to reliably score candidates.

Recent work on prompt generation for previously unseen tasks (Batorski and Swoboda, 2025) alleviates the per-task tuning burden, but still lags behind methods that are optimized separately for each task. Moreover, most prior work optimizes only

the instruction string, for example via rephrasing (using e.g. evolutionary algorithms or extensive search), ignoring the most impactful component of in-context learning (ICL): few-shot examples. The main exceptions are (Batorski et al., 2025), which synthesizes examples, but its computational cost is dozens of hours, and (Pryzant et al., 2023), which only selects examples from the training set. Our method is, to our knowledge, the first method that is fast, synthesizes new few-shot examples not found in the training set and requires relatively less access to training examples. When run long enough, our method additionally obtains new state of the art results among automatic prompting methods for a number of tasks. Our method works as follows: We first synthesize a proposal set of in-context examples that we append to our initial prompt. Then, in our optimization loop, we evaluate its efficacy on a small randomized evaluation set and identify the least helpful examples using a Monte Carlo Shapley estimator and replace, drop or keep it. If replaced, we draw from a pool of newly proposed few-shot examples. For efficiency and stability, we use a replay buffer for the evaluation set. Our algorithm has a favorable anytime performance: When run with a small computational budget, we attain second best results among our baselines on classification and summarization and already exceeds previous SoTA on simplification and mathematical reasoning tasks. When run with an extended budget that is still comparable to some other baselines, we exceed previous methods additionally on classification. Interestingly, even when running without the iterative update loop and only using the first generated few-shot examples, we often still get competitive results. To summarize, our contributions are as follows:

Conceptual: We propose PIAST, an automatic prompt construction method that augments a concise human-written instruction with a small

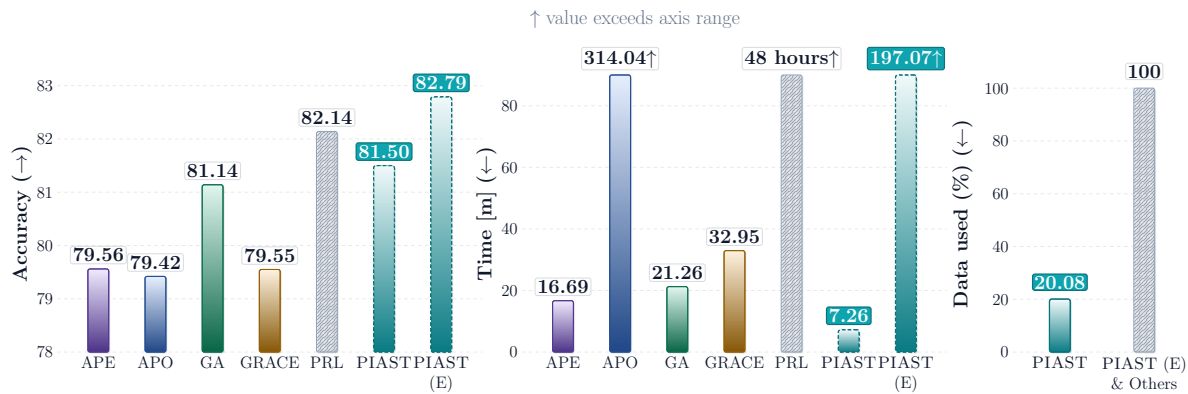


Figure 1: Overview of the results averaged over seven different text classification tasks, each run three times, comparing PIAST against current benchmarks. PIAST is able to generate high-quality prompts very efficiently, while requiring only a small portion of the dataset yielding comparable results to the current SOTA methods.

set of automatically generated few-shot examples. We use an iterative improvement loop that improves the current set of few-shot examples using Shapley values to estimate utility of individual examples.

Implementation: For a fast implementation, we approximate Shapley values, KV-cache reuse for shared ICL prefixes, PagedAttention for compact KV memory management and continuous token-level batching to maintain high GPU utilization.

Empirical Results: We demonstrate strong performances using the same set of robust hyperparameters on text classification, summarization and simplification as well as GSM8K. Our approach yields strong anytime performance: When using only a subset of data and a small computational budget we obtain SoTA on text simplification and GSM8K and obtain second best results on summarization and classification. With an extended budget and full training set access we additionally set a new SoTA on classification among automatic prompting methods.

2 Related Work

Prompt Engineering improves model capabilities without retraining (Liu et al., 2023). Chain-of-Thought (CoT) (Wei et al., 2022) elicits intermediate reasoning; Tree-of-Thought (ToT) (Yao et al., 2023) explores multiple paths; Program-of-Thoughts (Chen et al., 2022) and Graph-of-Thoughts (Besta et al., 2024) structure prompts as programs/graphs. Least-to-Most prompting decomposes problems into subproblems (Zhou et al.,

2023), and zero-shot CoT and self-consistency improve robustness (Kojima et al., 2022; Wang et al., 2022). Few-shot prompting (Brown et al., 2020) conditions on in-prompt exemplars and is effective for puzzles and evidence extraction (Xu et al., 2023; Greenblatt, 2024; Sivarajkumar et al., 2024). Lu et al. (Lu et al., 2022) show strong sensitivity to demonstration order and propose more robust orderings.

Automated Prompt Engineering replaces manual prompt design with automated search and refinement. APE (Zhou et al., 2022) generates candidate prompts and selects best by performance, remaining purely generative. APO (Pryzant et al., 2023) iteratively refines prompts via natural-language critiques, but its few-shot examples are restricted to the training set. EvoPrompt (Guo et al., 2023) evolves prompts with evolutionary operators, while PromptAgent (Wang et al., 2023) frames prompt optimization as planning and applies MCTS with error-driven feedback. Promptbreeder (Fernando et al., 2023) studies self-referential prompt evolution, where prompts generate and select improved variants of themselves via an evolutionary loop. Similarly, (Wang et al., 2025) frame prompt optimization as an open-ended, self-replicating process. GRACE (Shi et al., 2025) performs iterative prompt refinement with a gating rule and adaptively compresses prompts. OPRO (Yang et al., 2024b) uses an LLM as a black-box optimizer over instructions given past candidates and scores. Recent automatic prompt optimization methods also include CriSPO (He et al., 2025a), which improves prompts for text generation via multi-aspect critiques and suggestions, and ZERA (Yi et al.,

Algorithm	Dataset	Refinement	Few-shot	Auto Gen.	Speed
Manual Instruction (Zhang et al., 2022)	✗	✗	✗	✗	
APE (Zhou et al., 2022)	✓	✗	✗	✓	
APO (Pryzant et al., 2023)	✓	✓	⚡	✓	
EvoPrompt (Guo et al., 2023)	✓	✓	✗	✓	
PRL (Batorski et al., 2025)	✓	✓	✓	✓	
GRACE (Shi et al., 2025)	✓	✓	✗	✓	
PIAST	⚡	✓	✓	✓	

Legend: ✓ yes ✗ no ⚡ partial Speed: slow, . . . , fast.

Table 1: Comparison of PIAST with other methods from the literature. Dataset indicates the fraction of the dataset used during construction of the prompt. Refinement shows whether the method iteratively improves the current prompt or generates a new one in a single step. Few-shot specifies whether the method is capable of generating few-shot examples. Auto Gen. denotes whether prompts are generated automatically.

2025), which evolves prompts from zero initialization through principle-based refinement, and PromptWizard (Agarwal et al., 2025), which uses a feedback-driven critique-and-synthesis loop to iteratively optimize both instructions and synthetically generated in-context examples. AutoPrompt (Shin et al., 2020) searches over discrete prompt tokens using gradient signals, and (Lu et al., 2024) argue that random sampling is a strong prompt-optimization baseline. Other approaches leverage reinforcement learning, such as RLPrompt (Deng et al., 2022) (short token prompts) and PRL (Batorski et al., 2025), which can synthesize in-context examples when beneficial. Among these methods, only APO, PRL, and PromptWizard explicitly incorporate examples. APO reuses examples from the training set, whereas PRL and PromptWizard can generate novel examples. However, PRL often requires tens of hours of computation, which limits its practicality, and PromptWizard’s generated examples usually underperform. In contrast, PIAST rapidly generates few-shot examples not present in the training data. We summarize comparisons against the baselines we test against in Table 1.

Multi-agent / multi-module prompt optimization. Recent work optimizes prompts for multi-stage LM programs and agentic pipelines rather than a single LM call. MiPRO (Opsahl-Ong et al., 2024) and GePA (Agrawal et al., 2025) optimize instructions and demonstrations across modules; BetterTogether (Soylu et al., 2024) combines prompt optimization with parameter updates in modular pipelines. We do not compare to these methods because they assume a multi-module / multi-agent program and optimize end-to-end behavior, whereas PIAST produces a single reusable prompt

for one LLM call.

Demonstration selection / retrieval. Demonstration selection constructs a prompt per test input by retrieving demonstrations from a fixed pool, typically the labeled training set. (Rubin et al., 2022) train a dense retriever using LM-scored (input, demo) helpfulness labels; (Li et al., 2023) learn a retriever that generalizes across task families with a unified list-wise ranking objective; and Skill-KNN (An et al., 2023) uses engineered representations for instance-wise selection. MoD (Wang et al., 2024) partitions the demo pool into expert groups for collaborative retrieval. We do not compare to these approaches because they require per-instance retrieval and continued access to a demo bank at inference time, whereas PIAST generates a compact set of demonstrations once per task and reuses the same prompt for all inputs.

3 Method

In this section, we present our method, which is composed of three components: the Example Proposer, the Prompt Evaluator, and the Example Improver. Each component is instantiated as a frozen LLM with a distinct role in the overall pipeline. Our final prompt consists of the hand-crafted instruction proposed by (Zhang et al., 2022), concatenated with the in-context examples produced by our optimization procedure.

Example Proposer. The Example Proposer is responsible for generating initial candidate examples. It receives a task-specific initial instruction and produces a set of examples accordingly. To ensure coverage and robustness, the generated examples are deliberately diverse in both topic and length. Each example is then subject to a subse-

quent replace/drop/keep decision. The prompt for the Example Proposer is given in Appendix E.

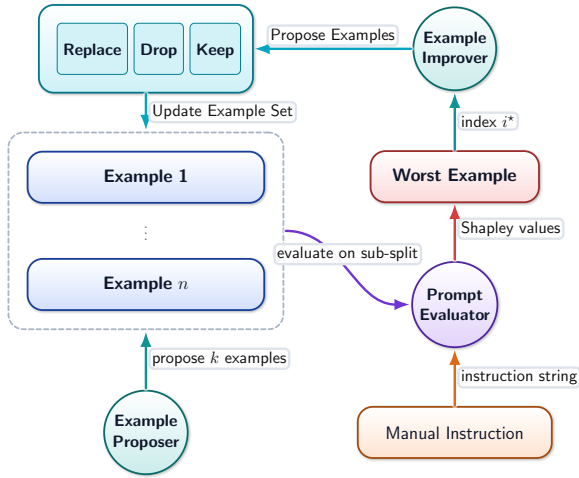


Figure 2: Pipeline of PIAS. Initially, the Example Proposer generates examples, which are then iteratively improved by evaluating them with the Prompt Evaluator and choosing new examples from the Example Improver to incorporate into the set of current in-context examples.

Prompt Evaluator. The Prompt Evaluator assesses the quality of candidate prompts. Given a prompt, it evaluates its performance on a subset D of the data using a specified metric f . This step ensures that only the most effective prompts are selected for further use. In our setup, the evaluated prompt is the base instruction concatenated with the proposed examples.

Example Improver. The Example Improver starts from a current prompt and iteratively changes the in-context examples in a replace/drop/keep cycle. The examples produced may differ in structure, topic, and length, thereby increasing diversity in the candidate pool. This process mirrors the behavior of the Example Proposer, where randomness in topic and sentence length is also introduced to encourage exploration. The prompt for the Example Improver is given in Appendix F. The example improver uses Shapley values to estimate quality of examples and implements a replace/drop/keep cycle as detailed below.

Example Selection via Shapley Values. At each iteration of the example improver, we identify the least useful example, which is then considered for a replace/drop/keep decision. To determine which, we employ Shapley values, which provide a principled way of estimating the contribution of each example to overall performance.

Let $[n] = \{1, \dots, n\}$ index the current few-shot examples and let $v : 2^{[n]} \rightarrow \mathbb{R}$ be a function that maps any subset $S \subseteq [n]$ to its utility (computed on D). In our setup, $v(\emptyset)$ is the accuracy of the instruction-only baseline (no few-shot examples). Write $\Pi([n])$ for the set of all permutations of $[n]$. For $\pi \in \Pi([n])$ and $i \in [n]$, define the predecessor set

$$P_i^\pi = \{j \in [n] : j \text{ precedes } i \text{ in } \pi\}.$$

The Shapley value of example i is

$$\phi_i = \frac{1}{n!} \sum_{\pi \in \Pi([n])} \left(v(P_i^\pi \cup \{i\}) - v(P_i^\pi) \right). \quad (1)$$

Remark 1 Although in-context learning can be sensitive to the order of examples, our Shapley estimator is order-agnostic: it treats coalitions as sets and evaluates each subset using the current fixed example order used by the optimizer and at inference time. This is still meaningful because our algorithm never searches over orderings, only replaces or drops examples, so we want attributions under the deployed order. Averaging over random subsets still captures redundancy and complementarity between examples, giving a robust signal at low computational cost.

Monte Carlo Approximation Because summing over all $n!$ permutations in (1) is infeasible, we approximate it by sampling K i.i.d. permutations $\pi^{(1)}, \dots, \pi^{(K)} \sim \text{Unif}(\Pi([n]))$ and computing

$$\hat{\phi}_i = \frac{1}{K} \sum_{k=1}^K \left(v(P_i^{\pi^{(k)}} \cup \{i\}) - v(P_i^{\pi^{(k)}}) \right). \quad (2)$$

In practice, each permutation yields marginal contributions for all i in a single sweep by maintaining the running predecessor set; averaging over K permutations gives an unbiased estimator of ϕ_i . The pseudo-code for selecting the least useful example is included in Appendix J.

Replace/Drop/Keep decision. Given the current set of in-context examples $[n]$ we determine the least helpful example i^* using the Shapley criterion:

$$i^* = \arg \min_{i \in [n]} \hat{\phi}_i. \quad (3)$$

For this step, the Example Improver proposes m candidate examples $C = \{c_1, \dots, c_m\}$ for potential appending to the current few shot example set.

To decide whether to replace, keep, or drop the index i^* , we compute the following scores:

$$r = \max_{c \in \mathcal{C}} v([n] \setminus \{i^*\} \cup \{c\}) \quad (\text{REPLACE})$$

$$d = v([n] \setminus \{i^*\}) \quad (\text{DROP})$$

$$k = v([n]) \quad (\text{KEEP})$$

We select replace, keep or drop by checking whichever score is largest. When having ties, we prefer replace over drop over keep. The next prompt becomes $(N \setminus \{i^*\}) \cup \{c^*\}$ under REPLACE, $N \setminus \{i^*\}$ under DROP, and N under KEEP. This policy ensures we only adopt a modification when it does not underperform the best available alternative (drop or status quo).

Remark 2 *Note that the Replace/Drop/Keep step could also be formulated directly using Shapley values. However, this would significantly increase the computational cost of each iteration, whereas our design prioritizes speed and efficiency.*

Replay Buffer Our method relies on sampling a subset of the training data at each iteration. Consequently, newly crafted examples can overfit the current subset and fail to generalize to training subsets drawn in later iterations, since there is no mechanism enforcing that they also perform well on previously seen data. To mitigate this, after each iteration we store a small portion of the training data in a replay buffer. At the next iteration, this buffer is merged with the freshly sampled subset, which preserves accuracy across iterations by acting as a regularizer: newly crafted examples must also succeed on data sampled in prior iterations.

Speed. To make our implementation fast, we employ the following techniques: We use a KV cache (Radford et al., 2019) to avoid recomputing attention over already-processed tokens: Keys and values for the shared in-context prefix are cached once and then reused across (i) all tokens within a sequence and (ii) multiple evaluation queries that share this prefix. In addition, we rely on PagedAttention (Korthikanti et al., 2023) to store the KV cache in paged memory chunks, which minimizes fragmentation and data movement while enabling efficient, contiguous access during decoding. Finally, we leverage continuous (token-level) batching (Yu et al., 2022), in which the scheduler dynamically forms a new batch at each decoding step by admitting fresh requests and retiring completed

ones, thereby overlapping prefill and decoding and maintaining high GPU utilization. Pseudocode for PIAST can be found in Appendix J.

4 Experiments

All experiments are run on a single NVIDIA A100 GPU. Unless stated otherwise, PIAST uses Qwen2.5-7B-Instruct (Yang et al., 2024a) for all three roles in the pipeline: the Example Proposer, the Example Improver, and the Prompt Evaluator. For a controlled comparison, we evaluate all baselines using the same underlying LLM, so differences in results reflect the prompting algorithms rather than model choice. We compare PIAST against established automatic prompting methods across four task families: (i) text classification, (ii) summarization, (iii) mathematical reasoning, and (iv) domain-specific question answering. We additionally report results on text simplification in Appendix A. We also conduct ablation studies to isolate the effect of key design and hyperparameter choices. In each ablation, we vary a single factor while holding all others fixed. All reported numbers (main results and ablations) are averaged over three independent runs. Across all tasks, we use one shared hyperparameter configuration; the full setting is provided in Appendix B.

4.1 Baselines

- **MI (Manual Instruction)** (Zhang et al., 2022): A set of prompts handcrafted and written by humans, aiming to improve task-specific performance.
- **NI (Natural Instruction)** (Mishra et al., 2021): Contains similarly to MI a set of human-written prompts for classification.
- **APE (Automatic Prompt Engineer)** (Zhou et al., 2022): Automatically generates multiple instruction candidates with an LLM and selects the most effective prompt based on downstream performance, without further refinement during optimization. This method only rephrases instructions and does not generate few-shot examples.
- **APO (Automatic Prompt Optimization)** (Pryzant et al., 2023): Frames prompt tuning as a black-box optimization problem, refining prompts through an iterative feedback loop with beam search. Incorporate few-shot examples taken directly from the training dataset.
- **EvoPrompt** (Guo et al., 2023): Uses evolutionary

Method / Dataset	SST-2	CR	MR	SST-5	AG’s News	TREC	Subj	Avg
MI	92.70	87.25	87.40	52.31	82.29	69.20	57.95	75.59
NI	95.77	91.50	90.85	51.90	83.43	66.60	68.10	78.31
APO	93.71 \pm 0.25	93.48 \pm 0.24	89.97 \pm 1.37	53.94 \pm 0.29	83.73 \pm 0.31	71.30 \pm 1.90	69.80 \pm 5.96	79.42
APE	91.23 \pm 0.66	92.87 \pm 0.02	89.90 \pm 0.94	49.37 \pm 5.66	82.58 \pm 1.20	77.07 \pm 1.61	73.92 \pm 1.39	79.56
GA	94.65 \pm 1.04	92.75 \pm 0.40	90.45 \pm 0.72	53.76 \pm 1.13	82.24 \pm 1.00	79.20 \pm 2.83	74.93 \pm 3.12	81.14
DE	93.29 \pm 0.34	93.38 \pm 0.19	89.98 \pm 0.24	55.25 \pm 0.37	82.18 \pm 1.04	76.47 \pm 0.38	73.08 \pm 4.95	80.52
GRACE	93.61 \pm 0.53	90.92 \pm 1.15	89.60 \pm 1.51	53.96 \pm 0.93	82.34 \pm 0.39	72.53 \pm 8.62	73.92 \pm 3.05	79.55
PRL	96.32 \pm 0.04	92.83 \pm 0.24	91.27 \pm 0.05	56.21 \pm 0.15	84.36 \pm 0.08	77.07 \pm 2.36	76.90 \pm 0.95	82.14
PromptWizard	93.87 \pm 1.38	89.45 \pm 0.25	89.30 \pm 0.15	52.04 \pm 0.14	82.82 \pm 1.00	71.10 \pm 2.10	73.05 \pm 3.12	78.80
PIAST	95.35 \pm 0.14	92.35 \pm 0.05	90.57 \pm 0.21	53.27 \pm 0.66	85.93 \pm 0.62	77.07 \pm 3.30	75.93 \pm 0.40	81.50
PIAST (E)	95.88 \pm 0.24	92.55 \pm 0.35	91.00 \pm 0.65	53.33 \pm 0.35	87.39 \pm 0.35	78.40 \pm 1.22	80.98 \pm 0.67	82.79
PIAST (I)	95.04 \pm 0.18	91.53 \pm 0.65	90.43 \pm 0.21	49.79 \pm 1.05	85.38 \pm 0.20	74.33 \pm 4.77	59.52 \pm 2.29	78.00
PIAST (LOO)	95.70 \pm 0.31	92.15 \pm 0.15	90.42 \pm 0.28	53.18 \pm 1.40	86.43 \pm 0.72	75.87 \pm 2.37	69.73 \pm 3.68	80.50

Table 2: Accuracy on classification tasks, averaged over three runs. Colours mark the best (red), second-best (orange) and third-best (yellow) numbers in each column; minor differences (≤ 0.05) are treated as ties. The right-most column shows the mean accuracy of each method across the seven datasets.

strategies, selection, crossover, and mutation—to evolve a pool of discrete prompts and discover high-performing candidates. Similar to APE, only rephrases instructions and does not generate few-shot examples.

- **DE (Differential Evolution):** Explores the prompt space using differential evolution strategies.
- **GA (Genetic Algorithm):** Applies genetic operators such as selection, crossover, and mutation to progressively improve prompt quality.
- **GRACE (Shi et al., 2025):** Proposes a prompt optimization framework based on gated refinement and adaptive compression. GRACE iteratively refines prompts while enforcing a no-regression constraint, only accepting modifications that do not degrade performance. To control prompt length and improve efficiency, it adaptively compresses prompts by removing or rewriting less useful components.
- **PRL (Prompts from Reinforcement Learning) (Batorski et al., 2025):** Employs a reinforcement learning framework to automatically generate and optimize prompts. PRL also constructs few-shot examples that are not in the training set.
- **PromptWizard (Agarwal et al., 2025):** Iteratively critiques and refines prompts using task-specific feedback, balancing exploration and exploitation throughout the search procedure. Moreover, it can generate novel in-context examples instead of relying solely on demonstrations sampled from the training set.
- **PIAST:** Our method as described in Section 3. The first two variants PIAST and PIAST(E) are used throughout experiments, while the (I) and (LOO) variants are ablations. All variants other-

wise have the same hyperparameters.

- **PIAST:** With medium runtime budget with limited access to the training set.
- **PIAST (E):** With extended runtime budget and accessing the full dataset.
- **PIAST (I):** Use only the initially generated examples, without the replace/keep/drop cycle. Notably this variant does not access the training set.
- **PIAST (LOO):** Replace Shapley value selection (1) by simple leave-one-out.

4.2 Results

Classification We evaluate PIAST on seven standard text classification datasets spanning sentiment (SST-2, MR, CR, SST-5), question type (TREC), news topic (AG’s News), and subjectivity (SUBJ). Results per-task are reported in Table 2 and summarized in Figure 1, including average runtime and data usage; per-dataset runtimes appear in Appendix D. Overall, PIAST is consistently among the top methods while being the fastest across benchmarks, and PIAST (E) further improves performance, setting new state-of-the-art results on AG’s News and SUBJ. Example prompts are provided in Appendix G.

Summarization We evaluate PIAST on abstractive summarization using SAMSUM (Gliwa et al., 2019), a dataset of messenger-style dialogues with human-written summaries. We report ROUGE-1/2/L (Lin, 2004), measuring unigram overlap (coverage), bigram overlap (local coherence), and longest common subsequence (fluency/structure). Table 3 shows that PIAST is the fastest method and ranks second across all ROUGE metrics. An interesting observation is that, although PRL is capable of generating examples, it does not uti-

lize any for the summarization task. Instead, PRL merely rephrases the manual prompt. The authors of PRL argue that summarization is not particularly suitable for example-based prompting. While we find that incorporating examples can indeed enhance performance, the improvements do not reach the level achieved by PRL. Finally, PIAST (E) improves further, achieving the best ROUGE-2 and the best average over ROUGE-1/2/L. The PIAST prompt is in Appendix I.

Method	ROUGE-1	ROUGE-2	ROUGE-L	Time [m]
MI	32.76	10.39	28.97	–
APE	37.12 \pm 2.02	12.97 \pm 0.74	33.32 \pm 1.68	60.07 \pm 0.27
GA	39.69 \pm 1.76	14.47 \pm 1.00	35.84 \pm 1.63	89.31 \pm 3.08
DE	33.91 \pm 4.04	12.53 \pm 1.47	31.05 \pm 3.79	76.89 \pm 1.34
GRACE	40.61 \pm 0.54	14.65 \pm 0.53	35.86 \pm 0.54	1125.38 \pm 21.40
PRL	42.47 \pm 0.83	16.17 \pm 0.24	37.73 \pm 0.36	2880.00 \pm 0.00
PIAST	41.13 \pm 0.67	16.07 \pm 0.76	36.74 \pm 0.48	34.48 \pm 0.27
PIAST (E)	42.13 \pm 0.27	16.83 \pm 0.3	37.37 \pm 0.25	737.00 \pm 108.31

Table 3: Text summarization results with ROUGE scores and runtime (minutes).

Domain knowledge task. To evaluate how PIAST performs on tasks that require domain-specific knowledge, we additionally test it on MedQA (Yang et al., 2024c), a multiple-choice medical question answering benchmark with four options (A–D). Results are reported in Table 4. PIAST achieves the second-best accuracy while being the fastest method overall, and PIAST (E) attains only a small additional improvement despite a substantially higher runtime. This suggests that, for domain tasks, most benefits of refinement are realized within the first few iterations of the replace/drop/keep loop, after which gains diminish.

Method	Accuracy	Time[m]
APE	45.66 \pm 0.97	34.67 \pm 9.75
GA	51.95 \pm 1.61	35.71 \pm 2.73
DE	51.76 \pm 0.16	88.63 \pm 3.57
GRACE	52.26 \pm 0.16	61.33 \pm 9.74
PRL	53.34 \pm 0.11	2880.00 \pm 0.00
PromptWizard	51.84 \pm 0.23	24.80 \pm 1.76
PIAST	52.45 \pm 0.51	23.58 \pm 0.61
PIAST (E)	52.89 \pm 0.05	617.55 \pm 89.20

Table 4: Results on MedQA dataset

Reasoning tasks We further evaluate PIAST on mathematical reasoning benchmarks: GSM8K (Cobbe et al., 2021), which requires explicit multi-step arithmetic with free-form integer answers, as well as DeepMath (He et al., 2025b) and Math500 (Lightman et al., 2023).

Method	GSM8K	DeepMath	MATH500	Time[m]
APE	83.43 \pm 1.98	15.47 \pm 0.45	31.53 \pm 1.04	180.81 \pm 2.66
GA	81.62 \pm 1.38	18.63 \pm 2.37	40.13 \pm 1.39	191.96 \pm 1.11
DE	79.52 \pm 0.45	16.10 \pm 0.00	34.20 \pm 1.39	252.57 \pm 3.59
GRACE	82.37 \pm 1.82	15.05 \pm 0.16	33.20 \pm 1.60	1436.41 \pm 74.30
PRL	86.15 \pm 0.55	21.58 \pm 0.22	44.40 \pm 1.40	2880.00 \pm 0.00
PromptWizard	86.61 \pm 0.19	18.52 \pm 1.08	35.50 \pm 0.30	39.94 \pm 2.17
PIAST	91.65 \pm 0.31	24.85 \pm 1.13	48.53 \pm 0.31	80.26 \pm 2.95
PIAST (E)	92.12 \pm 0.12	25.33 \pm 0.75	48.70 \pm 0.50	1598.34 \pm 234.54

Table 5: Results on GSM8K, DeepMath, and MATH500.

Performance on reasoning tasks is known to be highly sensitive to the choice of in-context exemplars (Wei et al., 2022), making them a natural benchmark for example-centric prompt construction. As shown in Table 5, methods that primarily rephrase or adjust the base instruction (APE, GA, DE) provide only modest improvements over the manual prompt, whereas methods that optimize few-shot examples (PRL, PIAST) achieve substantially higher accuracy. PIAST attains the strongest overall results and is also the fastest among the top-performing approaches, while PIAST (E) yields a small additional gain at higher cost. Overall, these results show that PIAST is both effective and compute-efficient for reasoning-intensive tasks.

4.3 Ablations

Ablation Study: Cross-Model Robustness We assess how well prompts learned by PIAST transfer across models in two settings. **(i) Cross-model inference.** We train prompts with Qwen2.5-7B-Instruct and then evaluate the same prompts using Mistral-7B-Instruct-v0.2 (Jiang et al., 2023). As shown in Table 6 (left), PIAST attains the strongest transfer on SUBJ, edging out PRL and APO; APE and GA are roughly on par with the manual instruction baseline. Notably, PIAST (E) exhibits a sizable accuracy drop, which we attribute to overfitting to the source evaluator due to its substantially larger improvement iteration budget. These results suggest that, when portability matters, PIAST offers the best cross-model robustness. **(ii) Component swaps.** We next vary which model plays each role, swapping Qwen and Mistral between the (Example Proposer & Improver) and the Prompt Evaluator. Table 6 (right) shows that accuracy remains comparable across configurations, indicating that PIAST is not overly sensitive to a particular model pairing. We additionally ran the full optimization pipeline on GSM8K using two other instruction-tuned mod-

Method	Acc.
MI	60.30
APE	60.77 \pm 1.08
APO	69.53 \pm 1.33
GA	60.68 \pm 1.60
DE	64.10 \pm 2.20
GRACE	69.10 \pm 1.88
PRL	70.73 \pm 3.81
PIAST	72.87 \pm 4.16
PIAST (E)	68.75 \pm 3.01

P & I	Eval	Acc.
Qwen	Mistral	75.93 \pm 3.14
Mistral	Qwen	72.52 \pm 1.71
Mistral	Mistral	74.93 \pm 2.54

Table 6: Cross-model robustness on SUBJ. **Left:** prompts trained with Qwen2.5-7B-Instruct, evaluated with Mistral-7B-Instruct-v0.2. **Right:** accuracy under component swaps between Qwen and Mistral (Example Proposer & Improver vs. Prompt Evaluator).

Method	Qwen3-4B	Llama3.1-8B
Manual Prompt	80.20	82.25
PIAST	91.58 \pm 0.77	87.60 \pm 0.43
PIAST (E)	92.23 \pm 0.27	88.40 \pm 0.12

Table 7: GSM8K accuracy (%) on additional instruction-tuned models. We run the full PIAST optimization pipeline separately for each model. PIAST consistently improves over the manual prompt, and PIAST (E) achieves the best performance in both cases.

els: Qwen3-4B-Instruct and Llama3.1-8B-Instruct. As shown in Table 7, PIAST yields substantial gains over the manual prompt on both models, and the extended-budget variant PIAST (E) provides a further improvement. These results suggest that the benefits of optimizing in-context examples are not specific to a single base model.

Ablation Study: Influence of the Replace/Drop/Keep Optimization In this experiment, we evaluate PIAST without the optimization loop, i.e., the model is tested directly on the proposed initial examples. We include this variant as a baseline, denoted PIAST (I), in Table 2. As shown, for many benchmarks the initial examples already yield strong performance, in some cases even surpassing algorithms that employ optimization loops. Nevertheless, we consistently observe that incorporating our optimization loop further improves the results. For certain tasks, such as binary sentiment classification (e.g., SST-2 or MR), the improvement is marginal. We attribute this to the fact that the initial examples are already highly effective, as the underlying LLM has a strong capability to distinguish between positive and negative samples. Interestingly, in the subjectivity dataset PIAST without the optimization loop performs poorly, achieving results comparable to those of the Manual Instruction baseline. However, after applying

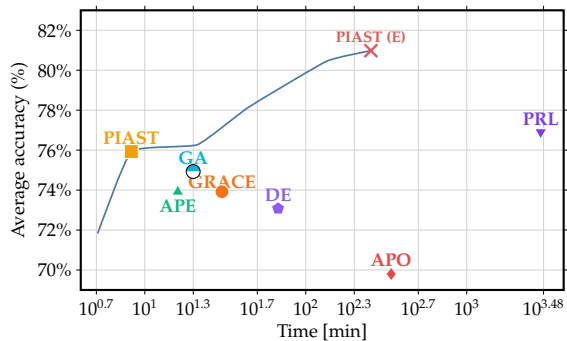


Figure 3: Scaling of PIAST on SUBJ compared to other baselines while increasing the number of improvement iterations.

the optimization loop, performance improves significantly by 16.41% and by an additional 5.05% when using PIAST(E), showing that on more difficult tasks where in-context example selection is non-trivial our optimization loop can help a lot.

Ablation Study: Influence of #Replace/Drop/Keep Iterations We investigate how PIAST scales with an increasing number of crafting iterations, and compare it against baselines from the literature. To this end, we run PIAST with 10, 15, 30, 50, 100, and 150 crafting iterations, and report the results in Figure 3. We observe a clear trend: increasing the number of crafting iterations consistently improves accuracy, albeit at the cost of higher runtime. This highlights an appealing property of PIAST: its performance can be effectively scaled by allocating more computation time by increasing the number of crafting iterations. Moreover, the plot clearly shows that PIAST has anytime performance superior to the baselines.

Ablation: Leave-One-Out vs. Shapley for worst-example selection To test whether a simpler procedure can replace our Shapley-value selection, we evaluate a *leave-one-out* (LOO) heuristic for identifying the worst (most harmful) in-context example. Given n examples $E = \{e_i\}_{i=1}^n$, LOO removes each example once and measures performance on the validation split:

$$i^* = \arg \max_{i \in \{1, \dots, N\}} v(E \setminus \{e_i\}).$$

In words, LLO chooses the example whose removal yields the highest accuracy drop. removing a strongly useful example would decrease it. We run this ablation on all classification tasks using the same hyperparameters as PIAST and report results in Table 2. Across several benchmarks LOO

Method	Training-set Examples	Synthetic Examples
PIAST	87.27 \pm 0.49	91.65 \pm 0.31
PIAST (E)	89.42 \pm 0.35	92.12 \pm 0.12

Table 8: GSM8K accuracy (%) when PIAST uses either training-set demonstrations or synthetic demonstrations. Synthetic examples yield consistently better performance for both budgets.

attains results comparable to our full method, but on SUBJ there is a clear gap between PIAST and PIAST (LOO). We hypothesize that, for many classification tasks, the initial pool already contains mostly good examples, so LOO can make small, beneficial adjustments. In contrast, SUBJ is more sensitive to initialization (see PIAST (I)), and the Shapley-based selection is notably more robust in such settings, where performance depends heavily on which examples are initially presented.

Ablation Study: Synthetic vs. Training-Set Demonstrations To isolate the contribution of synthetic example generation, we compare the standard PIAST pipeline against a variant that replaces synthetic demonstrations with examples sampled directly from the labeled training set, while keeping all other hyperparameters fixed. We evaluate this comparison on GSM8K, since it is particularly sensitive to the quality of in-context examples. Results in Table 8 show that synthetic examples consistently outperform training-set examples for both runtime budgets. This indicates that the gains of PIAST do not arise solely from iterative refinement, but also from the ability to construct stronger demonstrations than those obtained by directly sampling from the training data.

Ablation Study: Complementarity with Instruction Quality To test whether in-context example optimization is complementary to instruction quality, we rerun PIAST on SUBJ using two different starting instructions: the original manual instruction and a stronger natural instruction. We also report the instruction-only baselines without examples. As shown in Table 9, stronger instructions substantially improve the baseline performance. However, PIAST still provides large additional gains on top of both starting points, indicating that example optimization and instruction quality are complementary rather than redundant. In particular, the strongest performance is obtained by combining a better instruction with optimized in-context examples.

Method	Manual Instruction	Natural Instruction
No Examples	57.95	68.10
PIAST	75.93 \pm 0.40	85.95 \pm 0.05
PIAST (E)	80.98 \pm 0.67	87.25 \pm 0.20

Table 9: Accuracy (%) on SUBJ under two different starting instructions. Better instructions improve the baseline, but PIAST still adds substantial gains, showing that instruction quality and example optimization are complementary.

5 Conclusions

We have shown PIAST, a fast and effective prompting method that synthesizes ICL examples and outperforms several sota methods on a wider range of benchmarks. Interestingly, we argue that PIAST relies less on the training set than some other approaches. The only information it gets from the training set is its accuracy with the given prompt. Other approaches often use more information, e.g. using examples from the test set in the prompts. Additionally, we show that we do not need the full training set to achieve high performance, a subset typically suffices. We conjecture PIAST is a step towards ever less reliance on training sets for prompting. PIAST relies on the LLM having some task-specific competence to propose useful examples. In cases where this is not the case, a more involved iterative example synthesis might be needed. We also believe that, while example quality is often the dominant lever, combining PIAST with instruction rewriting could yield complementary gains. Another promising direction is to make the proposer benefit from other external signals (e.g., retrieval from a small seed set, weak heuristics, or verification-style self-checks) or by using a stronger or more specialized model for proposal while keeping the evaluator lightweight.

6 Limitations

Performance and applicability of PIAST are limited by:

Task-specific competence The LLM having task-specific competence: If the LLM cannot propose meaningful examples at all, the final prompt is unlikely to have helpful ICL examples.

Training set We rely on a training set to evaluate our ICL examples with. If no training set is given, PIAST can still propose initial ICL examples, but will not be able to improve them.

Scoring function Second, we need a scoring function on the training set as well. When it is not aligned with the ultimate task performance, prompts generated with PIAST might not help.

Prompt construction overhead PIAST incurs an initial overhead while constructing the prompt. While significantly smaller than other competing methods, it still takes a few minutes. This cost is incurred only once, however, and the generated ICL examples can be reused for any following task prompt.

Larger & proprietary LLMs While we have tested on state of the art medium sized open-source models, performance on larger and proprietary models might differ.

Human evaluation We do not conduct human evaluation. Our conclusions are therefore based on automatic metrics and benchmark accuracies, which may not fully capture qualities such as usefulness, readability, faithfulness, or preference in real-world use, especially for open-ended tasks like summarization and simplification.

Ethics Statement

This work was conducted in accordance with the ACL Code of Ethics and the ACM Code of Ethics and Professional Conduct, and we have adhered to the rules throughout this research. We recognize that large pre-trained language models, whether used with automated prompting or not, have the potential to be misused to generate misleading, toxic, or otherwise harmful content. Our intention is that the proposed prompting method contributes to improving the steerability of such models and supports more responsible use. We also acknowledge that appropriate safeguards and further evaluation are necessary to mitigate risks in real world applications.

Acknowledgments

The authors gratefully acknowledge the funding of this project by computing time provided by the Paderborn Center for Parallel Computing (PC2).

References

Eshaan Agarwal, Raghav Magazine, Joykirat Singh, Vivek Dani, Tanuja Ganu, and Akshay Nambi. 2025. [PromptWizard: Optimizing prompts via task-aware,](#)

[feedback-driven self-evolution.](#) In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 19974–20003, Vienna, Austria. Association for Computational Linguistics.

Lakshya A Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziem, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. 2025. [Gepa: Reflective prompt evolution can outperform reinforcement learning.](#) *arXiv preprint arXiv:2507.19457*.

Fernando Alva-Manchego, Louis Martin, Antoine Bordes, Carolina Scarton, Benoît Sagot, and Lucia Specia. 2020. [Asset: A dataset for tuning and evaluation of sentence simplification models with multiple rewriting transformations.](#) *arXiv preprint arXiv:2005.00481*.

Shengnan An, Bo Zhou, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Weizhu Chen, and Jian-Guang Lou. 2023. [Skill-based few-shot selection for in-context learning.](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*.

Pawel Batorski and Paul Swoboda. 2025. [Gps: General per-sample prompter.](#) *arXiv preprint arXiv:2511.21714*.

Paweł Batorski, Adrian Kosmala, and Paul Swoboda. 2025. [Pr1: Prompts from reinforcement learning.](#) *arXiv preprint arXiv:2505.14412*.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and 1 others. 2024. [Graph of thoughts: Solving elaborate problems with large language models.](#) In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. [Language models are few-shot learners.](#) *Advances in neural information processing systems*, 33:1877–1901.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2022. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks.](#) *arXiv preprint arXiv:2211.12588*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems.](#) *arXiv preprint arXiv:2110.14168*.

Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric P Xing, and Zhiting Hu. 2022. [Rlprompt: Optimizing](#)

- discrete text prompts with reinforcement learning. *arXiv preprint arXiv:2205.12548*.
- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. 2023. Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv preprint arXiv:2309.16797*.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. Samsun corpus: A human-annotated dialogue dataset for abstractive summarization. *arXiv preprint arXiv:1911.12237*.
- Ryan Greenblatt. 2024. [Getting 50% \(sota\) on ARC-AGI with GPT-4o](#).
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujie Yang. 2023. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. *arXiv preprint arXiv:2310.08510*.
- Han He, Qianchu Liu, Lei Xu, Chaitanya Shivade, Yi Zhang, Sundararajan Srinivasan, and Katrin Kirchhoff. 2025a. Crispo: Multi-aspect critique-suggestion-guided automatic prompt optimization for text generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 24014–24022.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, and 1 others. 2025b. Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning. *arXiv preprint arXiv:2504.11456*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *arXiv preprint arXiv:2310.06825*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Vijay Korthikanti, Zhaozhuo Yu, Zhun Yao, Yifan Zhu, Zhifan Shao, Le Zheng, Brandon Reagen, Tianqi Chen, and Rahul Jain. 2023. vllm: Easy, fast, and cheap llm serving with pagedattention. In *Proceedings of the ACM Symposium on Cloud Computing (SoCC)*, pages 1–15.
- Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng Qiu. 2023. Unified demonstration retriever for in-context learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing](#). *ACM Computing Surveys*, 55(9).
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098.
- Yao Lu, Jiayi Wang, Raphael Tang, Sebastian Riedel, and Pontus Stenetorp. 2024. Strings from the library of babel: Random sampling as a strong baseline for prompt optimisation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2221–2231.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.
- Krista Opsahl-Ong, Michael J Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia, and Omar Khattab. 2024. [Optimizing instructions and demonstrations for multi-stage language model programs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9340–9366, Miami, Florida, USA. Association for Computational Linguistics.
- Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. *arXiv preprint arXiv:2305.03495*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. OpenAI Technical Report. https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

- Wenhong Shi, Yiren Chen, Shuqing Bian, Xinyi Zhang, Kai Tang, Pengfei Hu, Zhe Zhao, Wei Lu, and Xiaoyong Du. 2025. No loss, no gain: Gated refinement and adaptive compression for prompt optimization. *arXiv preprint arXiv:2509.23387*.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*.
- Sonish Sivarajkumar, Mark Kelley, Alyssa Samolyk-Mazzanti, Shyam Visweswaran, and Yanshan Wang. 2024. An empirical evaluation of prompting strategies for large language models in zero-shot clinical natural language processing: algorithm development and validation study. *JMIR Medical Informatics*, 12:e55318.
- Dilara Soylu, Christopher Potts, and Omar Khattab. 2024. [Fine-tuning and prompt optimization: Two great steps that work better together](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10696–10710, Miami, Florida, USA. Association for Computational Linguistics.
- Jianyu Wang, Zhiqiang Hu, and Lidong Bing. 2025. Evolving prompts in-context: An open-ended, self-replicating perspective. *arXiv preprint arXiv:2506.17930*.
- Song Wang, Zihan Chen, Chengshuai Shi, Cong Shen, and Jundong Li. 2024. Mixture of demonstrations for in-context learning. In *Advances in Neural Information Processing Systems*.
- Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric P. Xing, and Zhiting Hu. 2023. [Promptagent: Strategic planning with language models enables expert-level prompt optimization](#). *arXiv preprint arXiv:2310.16427*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing statistical machine translation for text simplification. *Transactions of the Association for Computational Linguistics*, 4:401–415.
- Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott Sanner, and Elias B Khalil. 2023. Llms and the abstraction and reasoning corpus: Successes, failures, and the importance of object-based representations. *arXiv preprint arXiv:2305.18354*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024a. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024b. [Large language models as optimizers](#). *Preprint, arXiv:2309.03409*. ICLR 2024.
- Hang Yang, Hao Chen, Hui Guo, Yineng Chen, Ching-Sheng Lin, Shu Hu, Jinrong Hu, Xi Wu, and Xin Wang. 2024c. Llm-medqa: Enhancing medical question answering through case studies in large language models. *arXiv preprint arXiv:2501.05464*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822.
- Seungyoung Yi, Minsoo Khang, and Sungrae Park. 2025. Zera: Zero-init instruction evolving refinement agent—from zero instructions to structured prompts via principle-based optimization. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 23334–23348.
- Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. 2022. [Orca: A distributed serving system for transformer-based generative models](#). In *Proceedings of the 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI'22)*, pages 521–538. USENIX Association.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, and 1 others. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and Ed H. Chi. 2023. [Least-to-most prompting enables complex reasoning in large language models](#). In *International Conference on Learning Representations (ICLR)*.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large language models are human-level prompt engineers. In *The eleventh international conference on learning representations*.

A Additional Experimental Result - Simplification

We evaluate PIAST on the sentence simplification task using the ASSET dataset (Alva-Manchego et al., 2020). ASSET is a crowdsourced corpus specifically designed to evaluate simplification models across a range of rewriting operations, including lexical paraphrasing, sentence splitting, deletion, and reordering. Each original sentence is paired with multiple human-written simplifications. For measuring simplification quality, we adopt the SARI metric (Xu et al., 2016), which compares the system output to both the original sentence and the reference simplifications. Results are presented in Table 10. As shown, PIAST achieves the highest SARI score for text simplification while requiring the least runtime. Furthermore, PIAST exhibits the lowest standard deviations across runs, highlighting its stability and robustness. With additional computational time and data, PIAST (E) attains an even higher score on this benchmark. The example prompt is given in Appendix H.

Method	SARI	Time[m]
MI	43.77	–
APE	45.33 \pm 0.83	35.69 \pm 0.20
GA	46.25 \pm 0.47	39.60 \pm 0.63
DE	45.79 \pm 0.35	52.77 \pm 1.12
GRACE	50.21 \pm 0.18	611.61 \pm 19.57
PRL	52.26 \pm 3.51	2880.00 \pm 0.00
PIAST	54.52 \pm 0.07	18.14 \pm 0.42
PIAST (E)	55.06 \pm 0.26	389.78 \pm 113.17

Table 10: Simplification task results.

B Hyperparameters

In this section, we present the exact hyperparameters used across all tasks. A single set of hyperparameters was applied consistently across all tasks, and their values are summarized in Table 11.

The hyperparameters for PIAST (E) are identical to those of PIAST, except that we increase the number of iterations I to 150. This adjustment provides the most effective way to scale the performance of PIAST, as demonstrated in our ablation studies (see Section 4.3).

C Hyperparameter Study

We study the impact of the number of Shapley permutations, number of refine and initial example candidates for performance of PIAST.

Hyperparameter	Value
k	16
s	70
I	15
m	10
r	5
P	3

Table 11: Hyperparameters used across all tasks. The notation is consistent with the pseudocode provided in Appendix J.

m	Acc.	Time (min)
5	70.57 \pm 2.33	7.11 \pm 0.15
10	75.93 \pm 0.07	8.25 \pm 0.27
30	76.15 \pm 0.48	11.54 \pm 0.05
50	76.33 \pm 1.53	12.99 \pm 0.32

Table 12: Effect of the number of proposed candidates on the subjectivity task.

Hyperparameter Study: Influence of Number of Refine Candidates In this experiment, we analyze how the number of candidate examples m proposed by the example improver influences the final performance. The results are summarized in Table 12.

We observe that using too few candidates ($m = 5$) leads to a significant drop in accuracy, which can be mitigated by increasing the number of candidates to $m = 10$. For larger values $m = 30$ and $m = 50$, accuracy does not improve substantially, suggesting that an initial pool of $m = 10$ candidates combined with the replace/drop/keep iteration is sufficient to achieve strong performance. Therefore, in all subsequent experiments we adopt $m = 10$ as the default setting.

Hyperparameter Study: Influence of Number of Initial Examples In this experiment, we investigate how performance on the subjectivity task varies with the number of initial examples. The results are reported in Table 13. The number of initial examples does not exhibit a clear monotonic trend (e.g., “the more, the better”). Using too few initial examples may lead to underfitting, as they fail to provide sufficient insight into the task. Conversely, using too many examples can introduce excessive noise, making it more difficult for the evaluator to

identify the weakest examples within a large pool. Based on this analysis, we select 16 examples as the most robust choice for our experiments.

k	Acc.	Time [m]
4	71.77 \pm 0.42	2.88 \pm 0.12
8	74.20 \pm 1.66	4.63 \pm 0.14
16	75.93 \pm 0.40	8.25 \pm 0.27
32	74.77 \pm 0.48	16.26 \pm 0.26

Table 13: Effect of the number of initial examples on the subjectivity task.

Hyperparameter Study: Influence of Number of Shapley Permutations We study the sensitivity of PIAST to the number of Monte–Carlo permutations K used in the Shapley-value estimator (2). As shown in Table 14, increasing K beyond 3 yields only marginal accuracy gains (at most 0.39 points from $K = 3$ to $K = 50$) while substantially increasing runtime. Using $K = 1$ underperforms $K = 3$ by 2.25 points, indicating that a small amount of averaging is beneficial. We therefore adopt $K = 3$ as our default, which offers a strong performance/speed trade-off.

D Runtime Analysis for Classification Benchmarks

In this section, we report the average runtime (in minutes) across different methods for each classification benchmark. Table 15 summarizes the results, including mean and standard deviation values.

E Prompt for Example Generator

In this section we give the prompt for the prompt evaluator for binary sentimental analysis task. For other tasks, the prompt is analogous.

P	Acc.	Time [m]
1	73.68 \pm 1.45	4.61 \pm 0.04
3	75.93 \pm 0.40	8.25 \pm 0.27
10	76.02 \pm 1.08	20.88 \pm 0.17
50	76.32 \pm 1.41	83.90 \pm 1.82

Table 14: Effect of the number of Shapley permutations on subjectivity.

Example Proposer prompt for CR dataset

You are a data generator that writes high-quality in-context learning examples for *binary sentiment* on short movie-review style snippets. Create exactly {NUM_EXAMPLES} training examples in THIS STRICT format only:

Example1:
Sentence: ""<text>""
Label: {LABEL}

Example2:
Sentence: ""<text>""
Label: {LABEL}

...

Example{NUM_EXAMPLES}:
Sentence: ""<text>""
Label: {LABEL}

Diversity plan (MUST FOLLOW):
{DIVERSITY_PLAN}

Rules:

- Each example’s "Sentence" must contain exactly the number of sentences specified above (1–3).
- Keep sentences concise: typically 3–14 words each. Across the set, include at least one very short (≤ 5 words) and one longer (10–14 words).
- Use only ASCII characters. Do NOT include double quotes inside the text.
- Use exactly ONE ‘Sentence:’ line per example; if multiple sentences are needed, put them inside the same quotes separated by a space.
- Make the writing naturally match the requested label in the everyday sense of the word.
- Do NOT mention the label or talk about labels in the text (no meta commentary).
- No Markdown/code fences.
- Output ONLY the examples in the exact format above; no extra text.

F Prompt for Example Improver

Example Improver prompt for CR dataset

You are improving in-context examples for sentiment classification. Generate replacements that diversify length (1–3 sentences) and topic, avoid paraphrasing, and help the task.

You are given the CURRENT examples (do not repeat or paraphrase them):
{CURRENT_EXAMPLES}

Now create exactly {NUM_CANDIDATES} NEW ex-

Dataset	APE	APO	DE	GA	GRACE	PRL	PIAST	PIAST (E)
sst2	62.85 ± 1.96	376.82 ± 0.00	62.79 ± 2.54	20.21 ± 2.84	28.76 ± 1.48	2880.00 ± 0.00	7.49 ± 0.14	221.48 ± 8.41
cr	16.09 ± 5.29	302.22 ± 47.40	65.06 ± 0.96	18.75 ± 4.37	28.87 ± 1.20	2880.00 ± 0.00	7.49 ± 0.20	181.06 ± 8.67
mr	4.60 ± 0.45	342.03 ± 15.04	62.71 ± 1.36	28.84 ± 9.18	27.88 ± 1.70	2880.00 ± 0.00	7.64 ± 0.11	210.40 ± 20.84
sst5	5.99 ± 1.06	430.08 ± 92.80	62.09 ± 1.88	20.68 ± 1.82	44.94 ± 14.45	2880.00 ± 0.00	7.63 ± 0.11	155.55 ± 0.35
agnews	7.33 ± 2.43	241.19 ± 28.07	65.16 ± 3.59	18.08 ± 2.54	32.53 ± 4.09	2880.00 ± 0.00	8.38 ± 0.16	210.97 ± 38.85
trec	3.94 ± 0.74	256.34 ± 17.55	66.01 ± 3.25	22.24 ± 6.46	37.52 ± 2.75	2880.00 ± 0.00	6.63 ± 0.19	146.19 ± 28.15
subj	16.03 ± 2.63	339.44 ± 42.38	67.28 ± 0.37	19.98 ± 2.33	30.14 ± 4.26	2880.00 ± 0.00	8.25 ± 0.27	253.83 ± 8.07

Table 15: Average runtime in minutes (mean ± SD) for each classification benchmark and method. The last row reports the overall average across all APO runs.

amples in THIS STRICT format:

Example1:

Sentence: ""<text>""

Label: positiveneegative

Example2:

Sentence: ""<text>""

Label: positiveneegative

...

Example{NUM_CANDIDATES}:

Sentence: ""<text>""

Label: positiveneegative

Diversity plan (MUST FOLLOW):

{DIVERSITY_PLAN}

Rules:

- Use exactly ONE ‘Sentence:’ line per example. If multiple sentences are needed, put them INSIDE the same quotes separated by a space.
- Each example must have exactly the number of sentences specified in the plan above (1–3).
- Keep sentences concise: typically 3–14 words each. Across the set, include very short (≤ 5 words) and longer (10–14 words).
- ASCII only. Do NOT include double quotes inside the text.
- Make topics clearly different from the given examples and from each other; avoid near-duplicates or paraphrases.
- Prefer balancing labels; if unsure, choose the minority label: {MINORITY_LABEL}.
- Do NOT wrap output in Markdown/code fences.
- Output ONLY the examples in the exact format above; no extra text.

SST2

Please perform Sentiment Classification task. Given the sentence, assign a sentiment label from [‘negative’, ‘positive’]. Return label only without any other text.

Example1: Sentence: "The film maintains a steady pace. No dull moments. Engaging from beginning to end." Label: positive

Example2: Sentence: "The set design is detailed. Costumes match the era perfectly. Attention to historical accuracy is evident." Label: positive

Example3: Sentence: "The lead actor delivers a compelling performance." Label: positive

Example4: Sentence: "The film maintains a brisk pace from start to finish. No lulls or wasted moments, just continuous action and intrigue." Label: positive

Example5: Sentence: "The soundtrack is loud and distracting, overshadowing the dialogue." Label: negative

Example6: Sentence: "The editing is seamless, keeping the pace tight. Transitions between scenes are smooth and impactful." Label: positive

Example7: Sentence: "The lead actor’s performance is powerful and moving." Label: positive

Example8: Sentence: "The cinematography is breathtaking, capturing the essence of the story. The use of light and color enhances every scene." Label: positive

Example9: Sentence: "The lead actress’s performance is powerful." Label: positive

Example10: Sentence: "The direction was confusing and lacked focus." Label: negative

Example11: Sentence: "The lead actress shines in every scene." Label: positive

Example12: Sentence: "The dialogue felt forced and unnatural." Label: negative

Example13: Sentence: "The visual effects were poorly done. The CGI looked cheap and unrealistic. It ruined the immersion." Label: negative

Example14: Sentence: "The director masterfully guides the narrative." Label: positive

Example15: Sentence: "The lead actress’s portrayal is emotionally resonant." Label: positive

CR

Please perform Sentiment Classification task. Given the sentence, assign a sentiment label from [‘negative’, ‘positive’]. Return label only without any other text.

Example1: Sentence: "The cinematography captured the essence of the setting, with stunning visuals that added depth to the story." Label: positive

G Classification Prompts

In this section, we describe the most effective prompts for PIAST on classification tasks. The base prompts are taken from (Guo et al., 2023), and the examples are produced by our method.

Example2: Sentence: "The editing was precise, enhancing the flow of the story. However, some transitions felt abrupt and jarring." Label: negative

Example3: Sentence: "The actors delivered powerful performances, bringing depth to their roles." Label: positive

Example4: Sentence: "The screenplay was weak, with dialogue that felt forced and unnatural. Characters lacked development, making their actions confusing. The dialogue felt stiff, with lines that didn't flow naturally. This made the scenes less engaging." Label: negative

Example5: Sentence: "The first act was slow and 拖沓, 但中间部分节奏加快, 保持了紧张感。" Label: positive

Example6: Sentence: "The editing was seamless, enhancing the flow of the story. Cuts were precise, keeping the pacing tight." Label: positive

Example7: Sentence: "The film started slowly but picked up in the middle." Label: negative

Example8: Sentence: "The director's vision was clear but the execution was lacking. Scenes felt disjointed, and the overall story was confusing." Label: negative

Example9: Sentence: "The editing was choppy and disjointed." Label: negative

Example10: Sentence: "The director's vision was clear, but the actors seemed uncomfortable on camera." Label: negative

Example11: Sentence: "The screenplay felt rushed, with dialogue that seemed out of place. Characters had little to no development, making their motivations unclear. The plot relied too heavily on clichés, lacking originality." Label: negative

Example12: Sentence: "The director's vision was clear and inspiring. However, the final cut felt rushed and incomplete." Label: negative

Example13: Sentence: "The soundtrack was inappropriate, detracting from the mood of the scenes." Label: negative

Example14: Sentence: "The director skillfully guided the ensemble cast." Label: positive

Example15: Sentence: "The screenplay felt rushed, with dialogue that seemed out of place. Characters had little to no development, making their motivations unclear. The plot relied too heavily on clichés, lacking originality." Label: negative

MR

Please perform Sentiment Classification task. Given the sentence, assign a sentiment label from ['negative', 'positive']. Return label only without any other text.

Example1: Sentence: "The editing was choppy, disrupting the flow of the narrative. Scenes felt disjointed, and the timing was off." Label: negative

Example2: Sentence: "The lead actress delivered a nuanced and emotionally rich performance." Label: positive

Example3: Sentence: "The movie started slowly but picked up momentum. The second act was particularly well-paced, maintaining tension." Label: positive

Example4: Sentence: "The lead actor's performance was powerful and moving." Label: positive

Example5: Sentence: "The lead actress captivated the audience with her portrayal." Label: positive

Example6: Sentence: "The soundtrack was overbearing and distracting." Label: negative

Example7: Sentence: "The soundtrack added a melancholic tone that complemented the film's somber mood."

Label: positive

Example8: Sentence: "The actor's portrayal was compelling and emotionally resonant." Label: positive

Example9: Sentence: "The editing was choppy and disjointed." Label: negative

Example10: Sentence: "The lead actor brought depth to the role." Label: positive

Example11: Sentence: "The lead actor's portrayal was gripping and heartfelt." Label: positive

Example12: Sentence: "The soundtrack added a perfect touch, enhancing the film's dramatic moments." Label: positive

Example13: Sentence: "The visual effects were poorly done and noticeable." Label: negative

Example14: Sentence: "The director's vision was unclear, leading to a disjointed narrative. The characters felt underdeveloped." Label: negative

Example15: Sentence: "The acting was wooden and unconvincing." Label: negative

SST5

Please perform Sentiment Classification task. Given the sentence, assign a sentiment label from ['terrible', 'bad', 'okay', 'good', 'great']. Return label only without any other text.

Example1: Sentence: "The screenplay was predictable. The dialogue lacked depth, feeling forced. Characters felt flat and uninteresting." Label: bad

Example2: Sentence: "The story lacked coherence and felt rushed. The plot had too many loose ends and felt unsatisfying." Label: bad

Example3: Sentence: "The screenplay was cliché and predictable. The dialogue lacked depth and felt forced." Label: bad

Example4: Sentence: "The story was predictable with a weak ending. It lacked the twists needed for a compelling narrative." Label: okay

Example5: Sentence: "The screenplay was clever and witty. The dialogue was sharp and well-paced. It captured the essence of the characters perfectly." Label: great

Example6: Sentence: "Direction felt disjointed and confusing." Label: terrible

Example7: Sentence: "Story lacked coherence and felt rushed." Label: terrible

Example8: Sentence: "Dialogue was forced and awkward, breaking the mood." Label: terrible

Example9: Sentence: "Direction kept the pace just right; not too slow or fast." Label: good

Example10: Sentence: "The dialogue was cliché and predictable. The script failed to deliver any surprises." Label: okay

Example11: Sentence: "The director skillfully balanced the dramatic and comedic elements. The pacing was just right, keeping the audience engaged. The visual storytelling was top-notch." Label: great

Example12: Sentence: "Dialogue was sharp and added depth to the characters." Label: good

Example13: Sentence: "The acting was solid but not memorable. The supporting cast added some depth to the film." Label: okay

Example14: Sentence: "Acting was wooden and unconvincing." Label: terrible

Example15: Sentence: "Great performances by all; especially the lead actor." Label: good

Example16: Sentence: "The acting was wooden and unconvincing." Label: bad

Example17: Sentence: "The direction felt disjointed and confusing." Label: bad

Example18: Sentence: "The direction was uninspired. The pacing felt slow and the camera work was basic." Label: okay

AG's News

Please perform News Classification task. Given the news item, assign a label from ['World', 'Sports', 'Business', 'Tech']. Return label only without any other text

Example1: Sentence: "Upcoming elections will focus on healthcare reform and immigration policies; debates intensify." Label: World

Example2: Sentence: "NASA launches Mars rover to study planet's geology." Label: Tech

Example3: Sentence: "Supreme Court rules on new labor laws; impact on businesses discussed." Label: Business

Example4: Sentence: "Elections this year will focus on healthcare and education reform." Label: World

Example5: Sentence: "Telemedicine platforms see surge in usage during pandemic." Label: Tech

Example6: Sentence: "US and China engage in summit talks to discuss trade and climate issues; tensions remain high." Label: Business

Example7: Sentence: "Vaccination drive reaches remote villages successfully." Label: World

Example8: Sentence: "Diplomatic talks on climate change continue with mixed progress." Label: World

Example9: Sentence: "Bombing kills dozens in city center." Label: World

Example10: Sentence: "Upcoming elections will focus on healthcare and economic reforms; debates intensify as candidates present their plans." Label: World

Example11: Sentence: "Security forces respond to a terrorist attack in the city center; multiple casualties reported. Emergency services work to contain the situation." Label: World

Example12: Sentence: "Diplomatic talks on trade agreements between Asia-Pacific nations continue." Label: World

Example13: Sentence: "Upcoming elections will focus on healthcare and economic reforms; debates intensify as candidates present their plans. Voters express concerns about rising costs." Label: World

Example14: Sentence: "Diplomatic talks between nations on climate change progress despite initial disagreements." Label: World

TREC

Please perform Question Classification task. Given the question, assign a label from ['Description', 'Entity', 'Expression', 'Human', 'Location', 'Number']. Return label only without any other text.

Example1: Sentence: "Calculus is a branch of mathematics that deals with rates of change and slopes of curves. It includes differential and integral calculus. The fundamental theorem of calculus links these two concepts." Label: Description

Example2: Sentence: "Who wrote the screenplay for the movie where the main character delivers a famous monologue about the American Dream?" Label: Human

Example3: Sentence: "The human genome consists of all the genetic information in a human cell. It is composed of approximately 3 billion base pairs." Label: Number

Example4: Sentence: "The Magna Carta, signed in 1215, was a landmark document in English history." Label: Description

Example5: Sentence: "The director chose to shoot the film in black and white to evoke a sense of nostalgia." Label: Description

Example6: Sentence: "Calculus involves the study of limits, derivatives, integrals, and infinite series. It is essential for understanding changes in various quantities." Label: Description

Example7: Sentence: "She delivered the line with such conviction it seemed real. The audience was moved." Label: Expression

Example8: Sentence: "Which director is known for their innovative use of camera angles in films?" Label: Human

Example9: Sentence: "Meryl Streep has performed in many famous plays." Label: Location

Example10: Sentence: "He directed the actors to bring out the raw emotion in their performances. The result was powerful." Label: Expression

Example11: Sentence: "The screenplay's dialogue was sharp and witty, setting the tone for the entire film." Label: Expression

Example12: Sentence: "Newton's laws of motion describe how objects move under the influence of forces." Label: Description

Example13: Sentence: "The screenplay features complex dialogue that drives the characters' motivations and relationships." Label: Description

Example14: Sentence: "Recent studies show that vitamin C is crucial for the immune system. It helps in fighting infections." Label: Description

Example15: Sentence: "William Shakespeare's plays, such as 'Hamlet' and 'Macbeth,' are considered masterpieces of English literature. They explore complex themes like ambition, revenge, and madness." Label: Description

SUBJ

Please perform Subjectivity Classification task. Given the sentence, assign a label from ['subjective', 'objective']. Return label only without any other text.

Example1: Sentence: "The soundtrack was moving." Label: subjective

Example2: Sentence: "The soundtrack added an emotional depth." Label: subjective

Example3: Sentence: "The pacing started slow. It built tension. The climax felt rushed." Label: subjective

Example4: Sentence: "The visual effects were impressive. However, a few scenes felt overdone. Enhanced the world-building but occasionally distracted from the story." Label: subjective

Example5: Sentence: "The pacing was uneven. The first half dragged while the second felt rushed." Label: subjective

Example6: Sentence: "The visual effects were impressive, though a few scenes felt overdone. They enhanced the world-building but occasionally distracted from the story." Label: subjective

Example7: Sentence: "The pacing started slow but built tension." Sentence: "Climax felt rushed." Sentence: "Overall, uneven." Label: subjective

Example8: Sentence: "The lead actor brought depth to the role." Label: subjective

Example9: Sentence: "The pacing started slow. It built tension. The climax felt rushed." Label: subjective
 Example10: Sentence: "The lead actress gave a nuanced performance." Label: subjective
 Example11: Sentence: "The editing was choppy, disrupting the flow. It felt rushed at times." Label: subjective
 Example12: Sentence: "The screenplay was tightly constructed. The dialogue flowed naturally. Characters spoke authentically, enhancing the plot." Label: subjective
 Example13: Sentence: "The screenplay was tight. The dialogue flowed naturally. Characters spoke authentically, enhancing the plot. Subtle hints of conflict kept the audience engaged." Label: subjective
 Example14: Sentence: "The soundtrack added an emotional depth. However, the choice of music was sometimes jarring." Label: subjective

H Simplification Prompt

In this section, we present the best-performing prompts for PIAST on the simplification tasks. The base prompts are adapted from (Guo et al., 2023), while the examples are generated using our method.

SIMPLIFICATION

Simplify the text.
 Example1: Complex: "The Supreme Court decision declared the law unconstitutional, invalidating it." Simple: "The Supreme Court declared the law unconstitutional."
 Example2: Complex: "Mount Everest is the highest mountain in the world located in the Himalayas." Simple: "Mount Everest is the highest mountain in the Himalayas."
 Example3: Complex: "The pizza place offers a variety of toppings including pepperoni and mushrooms." Simple: "The pizza place offers pepperoni and mushrooms."
 Example4: Complex: "The Eiffel Tower is a famous landmark in Paris, France." Simple: "The Eiffel Tower is in Paris, France."
 Example5: Complex: "The university offers a range of degree programs in various fields of study." Simple: "The university offers degree programs."
 Example6: Complex: "The student passed the exam with excellent grades." Simple: "The student passed with excellent grades."
 Example7: Complex: "The basketball game was won by the team with the highest score at the end of the game." Simple: "The team with the highest score won."
 Example8: Complex: "The Renaissance was a period of great cultural change and achievement that started in Italy in the 14th century." Simple: "The Renaissance started in Italy in the 14th century."
 Example9: Complex: "The Industrial Revolution began in the late 18th century and changed manufacturing methods." Simple: "The Industrial Revolution changed manufacturing methods in the late 18th century."
 Example10: Complex: "The Mona Lisa is a famous painting by Leonardo da Vinci." Simple: "The Mona Lisa is a famous painting."
 Example11: Complex: "The internet is a global network that connects computers and allows for communication." Simple: "The internet connects computers for communication."
 Example12: Complex: "The Supreme Court ruled that the law was unconstitutional." Simple: "The Supreme

Court said the law was unconstitutional."
 Example13: Complex: "The Renaissance art focused on humanism and realism, emphasizing individual expression and naturalism." Simple: "Renaissance art emphasized individual expression."
 Example14: Complex: "The train arrived late due to a track problem." Simple: "The train was late due to a track problem."
 Example15: Complex: "The internet protocol is a set of rules that allows computers to communicate over the internet." Simple: "The internet protocol allows computers to communicate."

I Summarization Prompts

In this section, we present the most effective prompt for PIAST on summarization tasks. The base prompts are adapted from (Guo et al., 2023), and the examples are generated by our method.

SUMMARIZATION

How would you rephrase that in a few words?
 Example1: Text: ""IBM is an American multinational technology company headquartered in Armonk, New York."" Summary: ""IBM is headquartered in New York.""
 Example2: Text: ""Apple Inc. is an American multinational technology company headquartered in Cupertino, California."" Summary: ""Apple Inc. is headquartered in California.""
 Example3: Text: ""Tesla is an American electric vehicle and clean energy company."" Summary: ""Tesla is an electric vehicle company.""
 Example4: Text: ""The NBA All-Star Game is an annual basketball game featuring the top players from the National Basketball Association."" Summary: ""The NBA All-Star Game features top NBA players.""
 Example5: Text: ""Global warming is caused by an increase in greenhouse gases, leading to rising temperatures and climate changes."" Summary: ""Global warming is caused by rising greenhouse gas levels.""
 Example6: Text: ""Tesla is an American electric vehicle and clean energy company founded by Elon Musk, known for its innovative electric cars and energy storage solutions."" Summary: ""Tesla is an electric vehicle company.""
 Example7: Text: ""The Supreme Court justices are appointed by the President and confirmed by the Senate, serving life terms."" Summary: ""Supreme Court justices are appointed by the President and confirmed by the Senate.""
 Example8: Text: ""The Supreme Court of the United States is the highest court in the country, responsible for interpreting the Constitution and ensuring federal laws are followed."" Summary: ""The Supreme Court of the United States interprets the Constitution.""
 Example9: Text: ""The Mona Lisa, a painting by Leonardo da Vinci, is one of the most famous and most visited paintings in the world, currently housed in the Louvre Museum in Paris."" Summary: ""The Mona Lisa is a famous painting by Leonardo da Vinci.""
 Example10: Text: ""The Louvre Abu Dhabi, opened in 2017, is a museum in Abu Dhabi that focuses on art and culture from around the world."" Summary: ""The Louvre Abu Dhabi opened in 2017 and focuses on global art and culture.""

Example11: Text: ""The Louvre Museum in Paris is one of the largest and most visited art museums in the world, with a vast collection of art and artifacts."" Summary: ""The Louvre Museum in Paris has a large art collection.""

Example12: Text: ""The American Revolution was a violent conflict between Great Britain and thirteen of its North American colonies from 1775 to 1783."" Summary: ""The American Revolution lasted from 1775 to 1783.""

Example13: Text: ""The Renaissance was a period of great cultural and intellectual growth in Europe, spanning the 14th to the 17th century."" Summary: ""The Renaissance was a period of cultural and intellectual growth in Europe.""

Example14: Text: ""The Renaissance was a period of great cultural and intellectual growth in Europe, spanning the 14th to the 17th century, marked by a revival of classical learning."" Summary: ""The Renaissance was a period of cultural and intellectual growth in Europe.""

Example15: Text: ""The Eiffel Tower is a wrought-iron lattice tower on the Champ de Mars in Paris, France."" Summary: ""The Eiffel Tower is in Paris, France.""

Example16: Text: ""The giant panda is a bear species endemic to central China, recognized by its distinctive black and white fur and diet primarily consisting of bamboo."" Summary: ""The giant panda is a bear species endemic to China.""

J Pseudocodes

We present concise pseudocodes for our method and its Shapley-driven oracle. Algorithm 1 orchestrates the full crafting loop: starting from k initial examples, it performs I rounds that evaluate on a small subsample (size s) augmented with a replay buffer (size r). In each round, MONTECARLOSHAPLEYWORST identifies the least helpful example using P random permutations, the improver proposes m candidates, and a conservative REPLACE/DROP/KEEP rule updates the set only when accuracy does not regress. Algorithm 2 details the Shapley routine with memoized coalition values $v(S)$ and permutation-averaged marginal contributions, returning the index with the smallest estimated value.

K Usage of LLMs

We used LLMs to help polish the presentation and writing in this manuscript and to assist in drafting portions of the implementation code; all substantive research decisions and core technical contributions, however, were made by the authors.

Algorithm 1 PIAST

Require:

Dataset \mathcal{D}
Example Proposer M_{prop}
Prompt Evaluator M_{eval}
Example Improver M_{impr}
 k (Initial examples)
 I (Number of craft iterations)
 m (Number of refine candidates)
 s (Size of subdataset)
 r (Replay Size)
 P (Number of Shapley Permutations)

Ensure: Crafted example set E^*

```
1:  $E \leftarrow \text{PROPOSEINITIALEXAMPLES}(M_{\text{prop}}, k)$  ▷  $|E| = k$ 
2:  $R \leftarrow \emptyset$  ▷ Replay buffer
3: for  $t = 0, 1, \dots, I - 1$  do
4:    $D_t \leftarrow \text{SUBSAMPLE}(\mathcal{D}_{\text{infer}}, s)$ 
5:    $\tilde{D}_t \leftarrow D_t \cup R$  ▷ Union with replay
6:    $a_{\text{base}} \leftarrow \text{EVALACC}(M_{\text{eval}}, E, \tilde{D}_t)$ 
7:    $i^* \leftarrow \text{MONTECARLOSHAPLEYWORST}(E, \tilde{D}_t, M_{\text{eval}}, P)$ 
8:    $E_{\setminus i^*} \leftarrow E \setminus \{e_{i^*}\}$ 
9:    $a_{\text{drop}} \leftarrow \text{EVALACC}(M_{\text{eval}}, E_{\setminus i^*}, \tilde{D}_t)$ 
10:   $C \leftarrow \text{GENERATECANDIDATES}(M_{\text{impr}}, E_{\setminus i^*}, m)$ 
11:   $(c^{\text{best}}, a_{\text{best}}) \leftarrow \arg \max_{c \in C} \text{EVALACC}(M_{\text{eval}}, E_{\setminus i^*} \cup \{c\}, \tilde{D}_t)$ 
▷ Decision: REPLACE vs DROP vs KEEP
12:  if  $a_{\text{best}} \geq a_{\text{drop}}$  and  $a_{\text{best}} \geq a_{\text{base}}$  then
13:     $E \leftarrow E_{\setminus i^*} \cup \{c^{\text{best}}\}$  ▷ REPLACE
14:  else if  $a_{\text{drop}} \geq a_{\text{base}}$  and  $|E| > 1$  then
15:     $E \leftarrow E_{\setminus i^*}$  ▷ DROP
16:  else
17:     $E \leftarrow E$  ▷ KEEP
18:  end if
19:   $R \leftarrow R \cup \text{SAMPLEREPLAY}(D_t, r)$ 
20: end for
21:  $E^* \leftarrow E$ 
22: return  $E^*$ 
```

Algorithm 2 MONTECARLOSHAPLEYWORST

Require:

Example set $E = \{e_1, \dots, e_n\}$
Dataset subset \tilde{D}
Prompt Evaluator M_{eval}
 P (Number of Shapley permutations)

Ensure: Worst index i^*

```
1:  $\mathcal{V} \leftarrow \emptyset$ 
2: For each  $i \in [n]$ , set list  $\Delta_i \leftarrow []$ 
3: Define  $v(S) \leftarrow \text{EVALACC}(M_{\text{eval}}, \{e_j : j \in S\}, \tilde{D})$ 
4:  $\mathcal{V}[\emptyset] \leftarrow v(\emptyset)$ 
5: for  $p = 1, 2, \dots, P$  do
6:    $\pi \leftarrow$  a random permutation of  $[n]$ 
7:    $S \leftarrow \emptyset$ ;  $v_{\text{prev}} \leftarrow \mathcal{V}[\emptyset]$ 
8:   for  $j = 1, 2, \dots, n$  do
9:      $i \leftarrow \pi_j$ ;  $S' \leftarrow S \cup \{i\}$ 
10:    if  $S' \notin \mathcal{V}$  then
11:       $\mathcal{V}[S'] \leftarrow v(S')$ 
12:    end if
13:     $v_{\text{new}} \leftarrow \mathcal{V}[S']$ 
14:    Append  $(v_{\text{new}} - v_{\text{prev}})$  to  $\Delta_i$  ▷ Marginal contribution of  $e_i$ 
15:     $S \leftarrow S'$ ;  $v_{\text{prev}} \leftarrow v_{\text{new}}$ 
16:  end for
17: end for
18: for  $i = 1, 2, \dots, n$  do
19:    $\phi_i \leftarrow \begin{cases} \frac{1}{|\Delta_i|} \sum_{d \in \Delta_i} d, & |\Delta_i| > 0 \\ 0, & \text{otherwise} \end{cases}$ 
20: end for
21:  $i^* \leftarrow \arg \min_{i \in [n]} \phi_i$ 
22: return  $i^*$ 
```
