

# A Systematic Survey of Automatic Prompt Optimization Techniques

Anonymous ACL submission

## Abstract

Since the advent of large language models (LLMs), prompt engineering has been a crucial step for eliciting desired responses for various Natural Language Processing (NLP) tasks. However, prompt engineering remains an impediment for end users due to rapid advances in models, tasks, and associated best practices. To mitigate this, Automatic Prompt Optimization (APO) techniques have recently emerged that use various automated techniques to help improve the performance of LLMs on various tasks. In this paper, we present a comprehensive survey summarizing the current progress and remaining challenges in this field. We provide a formal definition of APO, a 5-part unifying framework, and then proceed to rigorously categorize all relevant works based on their salient features therein. We hope to spur further research guided by our framework.

## 1 Introduction

Since McCann et al. (2018) cast multi-task NLP as Question Answering, using prompts as inputs has become the standard way to elicit desired responses from Large Language Models (LLMs). Furthermore, LLMs' few-shot learning (Brown et al., 2020), instruction-following (Ouyang et al., 2022), and zero-shot reasoning capabilities (Kojima et al., 2023) have led to a widespread proliferation of prompting tricks for various tasks and model variants. However, LLMs still exhibit unpredictable sensitivity to various factors (explanation of the task (Li et al., 2023b), ordering (Liu et al., 2024a), stylistic formatting (Sclar et al.), etc.) causing a performance gap between two prompts that are semantically similar, thereby adding impediments for adoption by end users. Against this backdrop, Black-Box Automatic Prompt Optimization (APO) techniques have emerged that improve task performance via automated prompt improvements. They possess various attractive features - (1) they do

not require parameter access on LLMs performing the task, (2) they systematically search through the prompt solution space, and (3) they retain human interpretability of prompt improvements. In this survey paper, we aim to highlight the advances in the field. Our core contribution is a 5-part APO taxonomy combined with a comprehensive fine-grained categorization of various design choices therein (see Fig. 1, Tables 2, 3, 4 in Appendix). We hope our framework will be informational for new and seasoned researchers alike, enabling further research on open questions.

## 2 Automatic Prompt Optimization Formulation

We formalize the process of automatic prompt optimization (APO) as follows. Given a task model  $M_{task}$ , initial prompt  $\rho \in V$ , the goal of an APO-system  $M_{APO}$  is to obtain the best performing prompt-template  $\rho^{opt}$  under a metric  $f \in F$  and eval-set  $D_{val}$

$$\rho^{opt} := \arg \max_{\rho \in V} E_{x \sim D_{val}}[f(M_{task}(\rho \oplus x))] \quad (1)$$

This objective function is not tractable for discrete prompt optimization as token-sequence search spaces are combinatorial. Instead, APO techniques follow the general anatomy as described in Algorithm 1 to obtain approximate solutions.

## 3 Initialize Seed Prompts

### 3.1 Manual Instructions

Several approaches use a seed of manually created instructions that offer interpretable and strong baselines as the basis for further improvement, *inter alia.*, ProteGi (Pryzant et al., 2023), GPS (Xu et al., 2022), SPRIG (Zhang et al., 2024b). While obtaining quality examples can be costly, APE (Zhou et al., 2022)<sup>1</sup> showed that a few hundred samples are sufficient for further optimization.

<sup>1</sup>Note: APE stands for Automatic Prompt Engineer method introduced by (Zhou et al., 2022), not to be confused with APO

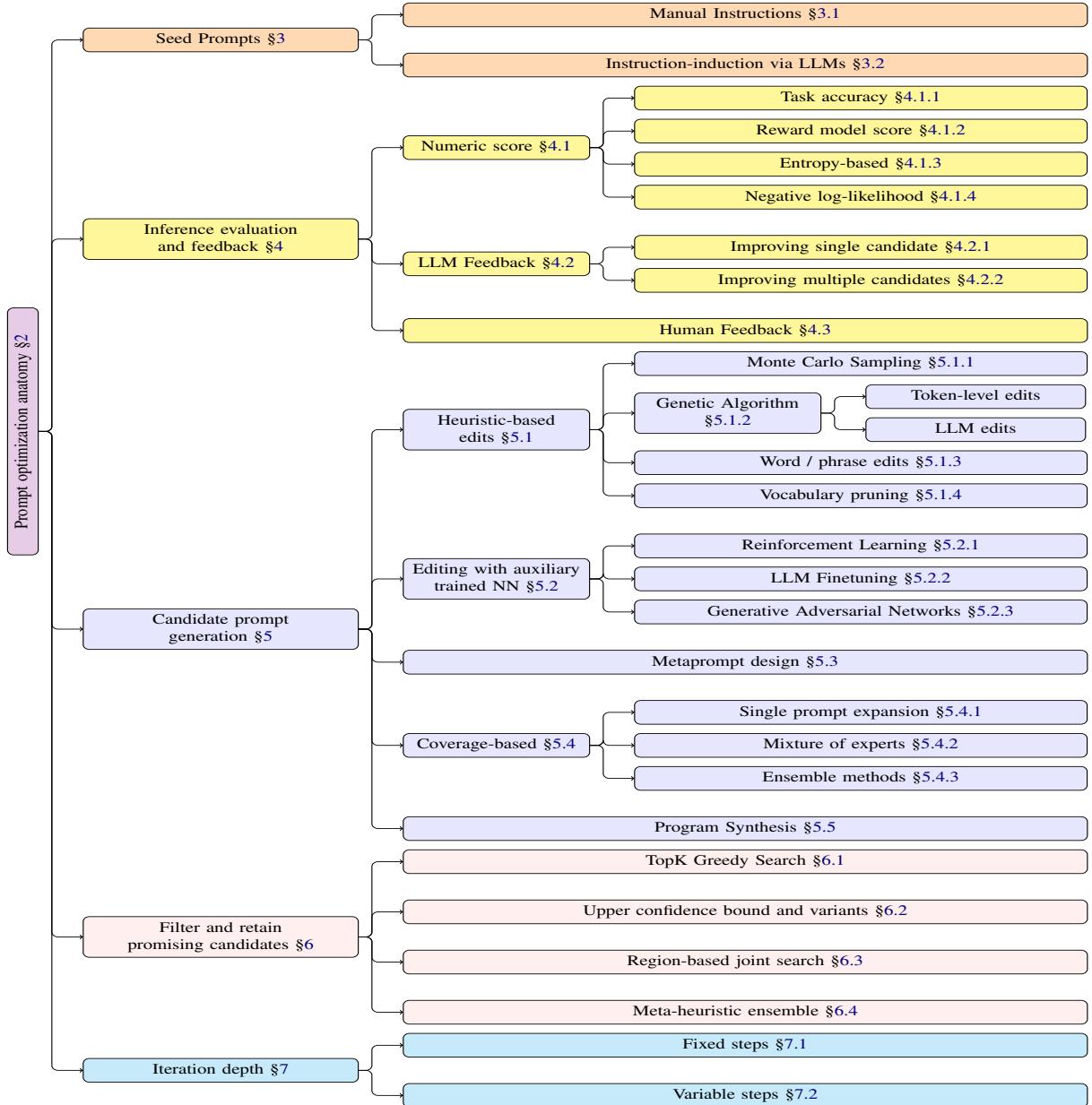


Figure 1: Taxonomy of Automatic Prompt Optimization

### Algorithm 1 Prompt optimization framework

```

1:  $P_0 := \{\rho_1, \rho_2, \dots, \rho_k\}$   $\triangleright \S 3$ . Seed prompts
2:  $D_{val} := \{(x_i, y_i)\}_{i=1}^n$   $\triangleright$  Validation set
3:  $f_1, \dots, f_m \in F$   $\triangleright \S 4$ . Inference evaluation
4: for  $t = 1, 2, \dots, N$  do  $\triangleright \S 7$ . Iteration depth
    $\triangleright \S 5$ . Generate prompt candidates
5:    $G_t := M_{APO}(P, D_{val}, F)$ 
    $\triangleright \S 6$ . Filter and retain candidates
6:    $P_t := Select(G_t, D_{val}, F)$ 
    $\triangleright \S 7$ . Optionally check for early convergence
7:   if  $f_{convergence} \leq \epsilon$  then
      exit
8: return  $\arg \max_{\rho \in P_N} E_{x \sim D_{val}} [f(M_{task}(\rho \oplus x))]$ 

```

### 3.2 Instruction Induction via LLMs

Honovich et al. (2023) were the first to propose inducing LLMs to infer human-readable prompts based on a few demonstrations  $E$  (see Appendix C.1 for prompt). APE (Zhou et al., 2022) and DAPO (Yang et al., 2024c) use the induced seed instructions for further optimization, while MOP (Wang et al., 2025) and GPO (Li et al., 2023c) use APE to induce cluster-specific prompts. Apart from demonstrations, SCULPT (Kumar et al., 2024) induced instructions from task-READMEs, while UniPrompt (Juneja et al., 2024) used LLMs to fill-

which broadly refers to the entire area of Automatic Prompt Optimization

089 in structured templates.

## 090 4 Inference Evaluation and Feedback

091 The evaluation step helps identify promising  
092 prompt candidates in each iteration. Some methods  
093 also use LLM feedback on prompt-response pairs  
094 to help generate more prompt candidates.

### 095 4.1 Numeric Score Feedback

#### 096 4.1.1 Accuracy

097 Using task-specific accuracy metrics is the most  
098 straightforward and widespread way of eliciting  
099 feedback, i.a., (Zhou et al., 2022, 2023; Zhang  
100 et al., 2024b; Khattab et al., 2022). Classification  
101 and MCQ-based QA tasks use exact accuracy,  
102 while code-related tasks measure execution accu-  
103 racy. Text generation tasks (summarization, transla-  
104 tion, creative writing) employ flexible metrics like  
105 BLEU-N, Rouge-N, Rouge-N-F1, or embedding-  
106 based measures such as BERTScore (Zhang\* et al.,  
107 2020) (Honovich et al., 2023; Dong et al., 2024b).

#### 108 4.1.2 Reward-model Scores

109 Given the limitations of rigid accuracy metrics,  
110 some approaches proposed using learned reward  
111 models to provide more nuanced evaluations of  
112 prompts-response pairs (Deng et al., 2022; Sun  
113 et al., 2024a; Kong et al., 2024). OIRL (Sun et al.,  
114 2024a) trained an XGBoost-based reward model  
115 that takes query-prompt embedding pairs as input  
116 and predicts whether the prompt will elicit correct  
117 answers from the language model and use it to se-  
118 lect appropriate prompts for specific queries using  
119 a best-of-N strategy. DRPO (Amini et al., 2024)  
120 follows an LLM-based reward modeling approach  
121 using both predefined and dynamic reward criteria.  
122 It first optimizes in-context learning examples  $E$ ,  
123 and using that it optimizes the specific task prompt.

#### 124 4.1.3 Entropy-based Scores

125 Entropy-based scores evaluate the entire output  
126 distribution induced by candidates, as opposed to  
127 a single inference instance. They are gradient-  
128 free but require access to the entire output prob-  
129 ability distribution, something not usually possi-  
130 ble with black-box LLMs. CLAPS (Zhou et al.,  
131 2023) leverages the negative incremental cross-  
132 entropy of  $\pi_{(x_i \oplus v \in V)}$  v/s  $\pi_{(x_i)}$  to identify promis-  
133 ing words  $v \in V$  to add to the prompt. The topK  
134 words are then used as candidate tokens from which  
135 to construct candidate prompts. GRIPS (Prasad  
136 et al., 2023) simply added an entropy term to

137 the task-weighted accuracy  $-\sum \pi_\rho(y) \ln(\pi_\rho(y)) +$   
138  $\frac{1}{|T|} \sum \mathbf{1}(y = \hat{y})$  to prioritize output diversity in po-  
139 tential prompt candidates.

#### 140 4.1.4 Negative Log-likelihood of Output

141 Some approaches like APE, GPS (Xu et al., 2022),  
142 PACE (Dong et al., 2024b) consider the negative  
143 log-likelihood (NLL) of token sequences under the  
144 target LLM, i.e.,  $-\log(\pi_\rho(y))$ . This however re-  
145 quires the log-probabilities to be accessible during  
146 the decoding of each token, limiting its applica-  
147 bility. The NLL for ground truth one-hot token-  
148 sequence is equivalent to the cross-entropy.

## 149 4.2 LLM Feedback

150 A popular paradigm to augment or fully replace  
151 numeric scores is to use textual feedback generated  
152 by  $LLM_{Evaluator}$  (Wang et al., 2024a; Long et al.,  
153 2024; Sinha et al., 2024). It is versatile because  
154 it can evaluate both the response as well as the  
155 prompt input. It can directly aid the prompt rewrit-  
156 ing process while being flexible to individual tasks  
157 as it only needs natural language instructions for  
158 general-purpose LLMs as opposed to task-specific  
159 handcrafting of metrics. A potential downside is  
160 the inference cost incurred due to an additional  
161 LLM call. All the LLM feedback approaches pro-  
162 vide multiple feedback data and broadly fall into  
163 two categories - improving a single prompt candi-  
164 date versus improving multiple prompt candidates  
165 (discussed below, examples in Appendix C.3).

### 166 4.2.1 Improving Single Candidate

167 SCULPT (Kumar et al., 2024) introduces a sys-  
168 tematic method for tuning long, unstructured prompts  
169 by employing a **hierarchical tree structure** and  
170 two-step feedback loops - preliminary assessment  
171 and error assessment - to evaluate and correct  
172 prompts before and after execution. The feed-  
173 back updates the hierarchical prompt tree which is  
174 then back-synthesized into a new prompt candidate.  
175 PACE (Dong et al., 2024b) applies an **actor-critic**  
176 editing framework to the prompt refinement pro-  
177 cess itself, allowing for more dynamic and adaptive  
178 adjustments. Overcoming the limitations of opti-  
179 mizing a single metric, CRISPO (He et al., 2025)  
180 adopts a **multi-aspect critique-suggestion** meta-  
181 prompt to highlight flaws in the generated response  
182 across multiple dimensions such as style, preci-  
183 sion, and content alignment. Thereafter it leverages  
184 detailed, aspect-specific feedback and iteratively  
185 updates the prompts. Autohint (Sun et al., 2023)

Paper	Seed instructions	Iteration depth	Inference evaluation	Candidate generation	Search+filter strategy
ProTeGi (Pryzant et al., 2023)	Manually created	Fixed	LLM feedback + Task accuracy	LLM rewriter	UCB for trees
APE (Zhou et al., 2022)	Instruction induction	Fixed	Task accuracy	N/A	UCB
CRISPO (He et al., 2025)	Manually created	Fixed	LLM feedback + Task accuracy	LLM rewriter	TopK selection
MOP (Wang et al., 2025)	Instruction induction	Fixed	Task accuracy	Mixture of experts	Region-based joint search
DSPY (Khattab et al., 2024)	Manually created + Instruction induction	Variable	LLM feedback + Task accuracy	Program Synthesis	TopK selection
OPRO (Yang et al., 2024a)	Manually created	Variable	LLM feedback + Task accuracy	Metaprompt design	TopK selection
GATE (Joko et al., 2024)	Manually created	Variable	Human feedback	LLM rewriter	N/A

Table 1: Comparison of some APO techniques under our framework (Tables 2,3,4 show full comparison)

summarizes feedback for multiple incorrect inferences via **hints** to instill improvements into a single prompt candidate.

#### 4.2.2 Improving Multiple Candidates

ProTeGi (Pryzant et al., 2023) and TextGrad (Yuksekgonul et al., 2024) leverage **textual “gradients”** to guide the discrete prompt optimization procedure, very similar to the gradient-descent style of continuous prompt optimization approaches. Different from continuous gradient-descent, ProTeGi sampled multiple “gradients” i.e. directions of improvement, and each such “gradient” is used to generate several prompt candidates for evaluation in the next iteration. PromptAgent (Wang et al., 2024a) similarly used an error collection approach to emulate expert-written prompts that consisted of clear sections like “Task description”, “Domain Knowledge”, “Solution Guidance”, “Exception Handling”, “Output Formatting”. PREFER (Zhang et al., 2024a) utilizes a feedback-reflect-refine cycle to aggregate feedback into multiple prompts in an **ensemble** to improve the model’s ability to generalize across various tasks. Survival of the Safest (SOS) (Sinha et al., 2024) added **safety-score** into a multi-objective prompt optimization framework that used an interleaved strategy to balance performance and security in LLMs simultaneously. To avoid accidentally damaging well-functioning prompts, StraGo (Wu et al., 2024) summarized strategic guidance based on both correct and incorrect predictions as feedback.

#### 4.3 Human-feedback

A few works also incorporate human feedback, either during compile-time or inference-time in the prompt construction / optimization process. Joko et al. (2024) proposed “Generative Active Task Elicitation” to better capture human preferences. It prompts a language model to interactively ask

questions and infer human preferences conditioned on the history of free-form interaction. Cheng et al. (2024) trained a smaller LLM to optimize input prompts based on user preference feedback, achieving up to 22% increase in win rates for ChatGPT and 10% for GPT-4. PROMST (Chen et al., 2024) tackles the challenges of multi-step tasks by incorporating human-designed feedback rules and a learned heuristic model. APOHF (Lin et al., 2024) focuses on optimizing prompts using only human preference feedback rather than numeric scores, employing a dueling bandits-inspired strategy to efficiently select prompt pairs for preference feedback, proving effective for tasks like text-to-image generation and response optimization.

### 5 Candidate Prompt Generation

In this step, one or more candidate prompts are generated that are most likely to result in an improvement in a metric of interest  $f \in F$ . The approaches reviewed below range from simple rule-based edits (sec. 5.1) to sophisticated agentic systems that combine with LLM-based evaluations (sec. 4.2) and various filtering strategies (sec. 6).

#### 5.1 Heuristic-based Edits

Several works proposed heuristic-based mechanisms to make edits to intermediate prompt candidates to generate newer candidates. They range from edits at the word / phrase / sentence-level (either simple rule-based or LLM-generated), or metric-driven incremental search. While these strategies may not result in the most optimal solution, they help in making the discrete prompt optimization problem computationally tractable.

##### 5.1.1 Monte Carlo Sampling

ProTeGi (Pryzant et al., 2023) uses Monte carlo sampling to explore combinatorial discrete solution spaces in an incremental fashion - it samples multiple textual gradients to use to generate prospective

262 candidates, and spawns paraphrases as monte-carlo  
263 successors for evaluation. PromptAgent (Wang  
264 et al., 2024a) uses a tree-variant called Monte Carlo  
265 Tree Search (MCTS) which consists of 4 steps —  
266 Selection, Expansion, Simulation, and Backpropa-  
267 gation (also explained in Sec. 6).

### 268 5.1.2 Genetic Algorithm

269 A significant line of work applies the well-studied  
270 genetic algorithms to make discrete edits to texts.  
271 The common recipe for several genetic algorithms  
272 is 1/ Mutate and 2/ Cross-over components from  
273 promising candidates. **Token mutations:** SPRIG  
274 (Zhang et al., 2024b) and CLAPS perform token-  
275 level mutations. SPRIG uses a starting corpus of  
276 300 components grouped into categories like COT,  
277 roles, styles, emotions, scenarios, and good prop-  
278 erties. It performs add/rephrase/swap/delete, high-  
279 lighting complementary strengths of optimizing  
280 system prompts alongside task-prompts (via meth-  
281 ods like ProTeGi) to enhance accuracy across mul-  
282 tiple diverse domains, languages, and tasks without  
283 needing repeated task-specific optimizations.

284 **LLM-based mutation:** LMEA (Liu et al., 2023),  
285 SOS (Sinha et al., 2024), and StraGo (Wu et al.,  
286 2024) uses mutation prompts with LLMs to over-  
287 come the traditional complexity of designing tai-  
288 lored operators for cross-over / mutation. Prompt-  
289 Breeder (Fernando et al., 2023) advocates self-  
290 referential improvement of all prompts in the  
291 prompt optimization system - Direct Mutation of  
292 task prompts, Hypermutation of mutation prompts  
293 themselves, Lamarckian Mutation where prompts  
294 are reverse-engineered from successful examples  
295 (similar to Instruction Induction Honovich et al.  
296 (2023), and finally Crossover and Shuffling to im-  
297 prove diversity of the prompt pool. EvoPrompt  
298 (Guo et al., 2024) use Differential Evolution -  
299 where differences between existing prompts is in-  
300 corporated to form new prompt candidates to over-  
301 come the problem of local optima. AELP (Hsieh  
302 et al., 2024) also uses mutation operators to per-  
303 form sentence-level edits in an iterative fashion.  
304 They include sentence-level histories of reward  
305  $\{(s_{t-1}, s_t, r_t)\}$  in the mutation prompt in order  
306 to avoid local optima and accidentally returning  
307 to sub-optimal versions. GPS (Xu et al., 2022)  
308 used Back-translation, Sentence Continuation, and  
309 Cloze transformations to perform prompt mutation.  
310 PromptWizard (Agarwal et al., 2024) proposed a  
311 pipeline combining several steps including itera-

312 tive improvement, few shot example synthesis and  
313 selection, utilizing LLM’s reasoning capability to  
314 improve and validate the prompt, and finally an  
315 expert persona to ensure consistency of the style of  
316 generated prompts.

### 317 5.1.3 Word / Phrase Level Edits

318 Several word-edit approaches first identify "influ-  
319 ential" tokens in the prompts. COPLE (Zhan et al.,  
320 2024) argued that LLMs exhibit lexical sensitivity,  
321 showing that merely replacing a few words with  
322 their synonyms can yield significant improvements.  
323 First, "influential" tokens are identified where ex-  
324 pected loss on dev-set  $E_{Dval}[L(y, \hat{y})]$  drops the  
325 most after removing that token versus the original  
326 prompt, and then influential tokens are replaced  
327 using predictions from a Masked-Language Mod-  
328 els. This token-replacement approach is also at-  
329 tractive as a standalone post-processing step for  
330 long prompts that are already optimized using other  
331 LLM-based approaches. GRIPS (Prasad et al.,  
332 2023) argues that phrase level edition is an effec-  
333 tive and interpretable method to optimize prompts,  
334 leveraging 4 basic edit operations -add, delete, para-  
335 phrase, and swap

### 336 5.1.4 Vocabulary Pruning

337 Some works prune the vocabulary space  $V$  to  
338  $V_{pruned}$  for decoding the next token for the op-  
339 timized prompt  $\rho^*$ . CLAPS (Zhou et al., 2023)  
340 argued that general search spaces are highly re-  
341 dundant and use K-means clustering to find word-  
342 clusters and retain top-2000 words closest to cluster  
343 centroids. BDPL (Diao et al., 2022) used pairwise  
344 mutual information (PMI) to retain top co-occurring  
345 ngrams for decoding. PIN (Choi et al., 2024) in-  
346 stead added regularization in the form of Tsallis-  
347 entropy (ideal for heavy-tailed distributions like  
348 natural language) for the RL training of a prompt  
349 generation network, to reduce the probability mass  
350 for unlikely tokens and improve interpretability.

## 351 5.2 Editing via Auxiliary Trained NN

352 Some approaches leverage a trained auxiliary neu-  
353 ral network to edit the initial prompt for ob-  
354 taining desired improvements. We include ap-  
355 proaches where the finetuned network is different  
356 and smaller than the task network.

### 357 5.2.1 Reinforcement-learning

358 **Multi-objective Optimization** techniques (Jafari  
359 et al., 2024) demonstrate superiority over simple

reward averaging, particularly through volume-based methods that effectively balance competing objectives. Dynamic prompt modification strategies, introduced through **prompt rewriting** (Kong et al., 2024), directional stimulus prompting (Li et al., 2023d) and **test-time editing** (Zhang et al., 2022) solve the important goal of moving beyond static prompt generation. Prompt-OIRL (Sun et al., 2024a) also tackled test-time optimization objective by learning an **offline reward model** and subsequently using a best-of-N strategy to recommend the optimal prompt in a query-dependent fashion. BDPL (Diao et al., 2022) optimized discrete prompts using variance-reduced policy gradient algorithm to estimate gradients, allowing user devices to fine-tune tasks with limited API calls.

## 5.2.2 Finetuning LLMs

BPO (Cheng et al., 2024) trains a smaller 7B model to align itself to task-performance on individual LLMs using reward-free alignment. FIPO (Lu et al., 2025) trains a local model (7B - 13B) to perform prompt optimizations to preserve privacy and adapt to target models better leveraging both data diversification and strategic fine-tuning such as SFT, preference optimization, and iterative preference learning.

## 5.2.3 Generative Adversarial Networks

Long et al. (2024) framed the prompt optimization process in the GAN setting. The LLM generator takes question and the generation prompt to produce output. The (input, output) pairs are evaluated by an LLM powered discriminator, whose goal is to identify generated pairs from ground truth pairs. Both generator and the discriminator are jointly optimized using adversarial loss, by utilizing a prompt modifier LLM to rewrite their prompts.

## 5.3 Metaprompt Design

PE2 (Ye et al., 2024) argued that previous works under-explored meta-prompt search space. OPRO (Yang et al., 2024a) proposes a meta-prompt design (see Appendix C.2) which includes the optimization problem description in natural language and previously generated solutions (multiple solutions per stage for diversity) and scores alongside the meta-instruction for prompt refinement. DAPO (Yang et al., 2024c) utilizes a well-designed meta-instruction to guide the LLM in generating high-quality and structured initial prompts (contain task-specific info, e.g. task type and description, output

format and constraints, reasoning process, professional tips) by observing given input-output exemplars. Then, DAPO iteratively optimizes the prompts at the sentence level, leveraging previous tuning experience to expand prompt candidates.

## 5.4 Coverage-based

Some approaches seek to "cover" the entire problem space - either within a single prompt, or using multiple prompts working individually or in an ensemble during inference.

### 5.4.1 Single Prompt-expansion

AMPO (Yang et al., 2024d) uses LLM feedback to enumerate all the failure cases based on the evaluation-set  $D_{val}$  and then enlists each of them in the meta-instruction in an if-then-else format using 3 modules - 1/ Pattern Recognition, 2/ Branch Adjustment, and 3/ Branch Pruning to decide whether to enhance existing branches, or to grow new branches. Similarly, UNIPROMPT focused on explicitly ensuring that various semantic facets of a task get represented in the final prompt. It designs a human-like (manual) prompt engineering approach (UniPrompt) with two stages: a) task facets initialization using background knowledge, and b) refinement using examples.

### 5.4.2 Mixture of Experts

Wang et al. (2025) introduced the Mixture-of-Expert-Prompts where each expert is a task-prompt to be used for specialized inference. MOP first clusters all demonstrations using K-means clustering. Then, the Region-based Joint Search (RBJS) (sec.6.3) algorithm generates the appropriate instruction for each exemplar-cluster via instruction induction (sec.3.2) based on a mix of in-cluster and out-of-cluster demonstrations to cover “blind-spots”. During inference, a single expert prompt is invoked whose cluster centroid  $\mu_c$  is closest to the instance-embedding  $\arg \min_C ||\phi(x_i) - \mu_c||_2$ .

### 5.4.3 Ensemble Methods

PromptBoosting (Hou et al., 2023), Boosted-Prompting (Pitis et al., 2023), PREFER (Zhang et al., 2024a), etc. are ensemble methods that invoke multiple prompts during inference and combine them to generate the final output  $\hat{y} = y_0 + \sum_m \beta_m y_i$ . GPO (Li et al., 2023c) also uses labeled source data to generate an ensemble of prompts, which are applied to unlabeled target data to generate output through majority voting.

## 457 5.5 Program Synthesis

458 Program-synthesis based approaches transform  
459 LLM pipelines into structured, modular components  
460 that can be systematically optimized and  
461 composed. These optimization techniques iteratively  
462 refine instructions and demonstrations for  
463 each module to improve the entire pipeline’s per-  
464 formance, DSP (Khattab et al., 2022) introduces  
465 a three-stage framework for retrieval-augmented  
466 inference: Demonstrate (generates task-specific  
467 demonstrations), Search (retrieves relevant infor-  
468 mation), and Predict (combines retrieved info with  
469 demonstrations). DSPY (Khattab et al., 2024)  
470 transforms LLM pipelines into text transformation  
471 graphs - introducing parameterized models, learn-  
472 ing through demonstrations, and a compiler that op-  
473 timizes pipelines. DLN (Sordoni et al., 2023) simi-  
474 larly considers chained LLM calls as stacked deep  
475 language networks performing variational infer-  
476 ence, where the learnable parameters for each layer  
477 are task-decomposed prompt templates. MIPRO  
478 (Opsahl-Ong et al., 2024) automates the optimiza-  
479 tion of multi-stage language model programs by  
480 improving instructions and demonstrations for each  
481 module. SAMMO (Schnabel and Neville, 2024)  
482 proposed symbolic prompt programming, repre-  
483 senting prompts as directed-acyclic-graphs (DAG).  
484 A set of user-defined node mutation rules guide the  
485 mutation-search to find the optimal DAG based on  
486 , which is then converted back to a prompt.

## 487 6 Filter and Retain Promising Prompts

488 In this step, promising prompt candidates are fil-  
489 tered for further optimization.

### 490 6.1 TopK Greedy Search

491 The simplest mechanism to iteratively search  
492 through prompt candidate sets is a greedy topK  
493 search where in each iteration of the optimiza-  
494 tion, the top-K best-performing candidates on mini-  
495 batch of data instances  $D_{val}$  are retained for further  
496 iterations (e.g. - ProTeGi, AELP. This differs from  
497 beam-search which judges partial solutions’ based  
498 on the reward for the entire trajectory of prompt  
499 edits  $r(\{\rho_1^1, \rho_2^1, \dots, \rho_t^1\})$ .

### 500 6.2 Upper Confidence Bound and Variants

501 Relying on a single static evaluation dataset can  
502 lead to biases in the selection procedure and finally  
503 suboptimal solutions. ProTeGi, SPRIG, *inter alia*,  
504 cast the candidate prompt selection problem as that  
505 of bandit search - identifying the most suitable

506 arm (prompt candidate) operating on a fixed com-  
507 putation budget. They use the Upper Confidence  
508 Bounds (UCB, Algorithm 2) which balances explo-  
509 ration with exploitation. In each iteration of prompt  
510 optimization, they sample a different evaluation  
511 dataset  $D_{sample} \in D_{val}$ , and maintain a moving  
512 estimate of the optimality of each arm (i.e. prompt).  
513 In each iteration, the playout filters top-B prompt  
514 candidates with the greatest score for further ex-  
515 ploration. PromptAgent uses a variation of UCB  
516 called UCB for Trees (UCT) which are used in the  
517 setting of contextual bandits (i.e. the action-space  
518 and the reward function is state-dependent). AELP  
519 (Hsieh et al., 2024) used a modification called Lin-  
520 ear UCB (Li et al., 2010) which uses a closed form  
521 linear estimate based on the reward trajectories of  
522 previously sampled edits as well as prompt embed-  
523 ding  $\phi(s)$  to select the next best-arm.

## 524 6.3 Region-based Joint Search

525 MOP (Wang et al., 2025) proposes a Mixture-  
526 of-Expert-Prompts performing prompt optimiza-  
527 tion for each expert individual. Once C exemplar-  
528 clusters are identified, the RBJS search first sam-  
529 ples examples  $D_{exemplars} \in D_C \cup D \setminus D_C$ , and  
530 then uses APE to induct and optimize each expert  
531 instruction.

## 532 6.4 Metaheuristic Ensemble

533 PLUM (Pan et al., 2024) library offered a meta-  
534 heuristic ensemble of different search algorithms  
535 like Hill climbing, Simulated Annealing, Genetic  
536 Algorithms, Tabu Search, and Harmony Search.

## 537 7 Iteration Depth

### 538 7.1 Fixed Steps

539 Most approaches choose to carry out the prompt  
540 optimization for a fixed number of steps N.

### 541 7.2 Variable number of steps

542 GRIPS (Prasad et al., 2023) concludes search when  
543 successive iterations with negative gains breach  
544 a patience parameter, whereas PromptAgent con-  
545 cluded APO when  $r_t \leq \epsilon_{min} \vee r_t \geq \epsilon_{max}$ .

## 546 8 Theoretical Perspectives

### 547 8.1 Upper Bound of Improvement from APO

548 AlignPro (Trivedi et al., 2025) establishes an upper  
549 bound on the gains realizable from discrete prompt  
550 optimization under a given prompt optimizer and  
551 also a suboptimality-gap w.r.t. RLHF-optimal pol-  
552 icy  $\pi^*$ , while a lower bound is left unexplored.

## 553 8.2 Other Related Perspectives

554 Bhargava et al. (2024) proposed a control theoretic framework to establish bounds on the set of  
555 reachable LLM-outputs for self-attention in terms of the singular values of its weight matrices. Liu  
556 et al. (2024c) showed the existence of a strong  
557 transformer that can approximate any sequence-to-  
558 sequence Lipschitz function. They also showed the  
559 existence of “difficult” datasets that depth-limited  
560 transformers could not commit to memory.

## 563 9 Related Areas

### 564 9.1 In-context Learning

565 Wan et al. (2024) advocated jointly performing  
566 Instruction Optimization (IO) and Exemplar Opti-  
567 mization (EO), showing EO can naturally “emerge”  
568 within IO. Dong et al. (2024a) extensively surveyed  
569 EO aspects like selection, ordering, formatting, etc.

### 570 9.2 LLM as a Judge

571 Zheng et al. (2023) first proposed using LLMs as  
572 an evaluator in multi-turn conversations, while also  
573 throwing light on inherent biases such as position-  
574 bias, verbosity-bias, and self-enhancement bias.  
575 Several works have since leveraged proposed re-  
576 flective agentic systems, using specialized prompts  
577 / LLMs to provide feedback for error correction.

### 578 9.3 Soft Prompts

579 Lester et al. (2021) introduced so-called “soft  
580 prompt” prefixes of task-specific vectors  $\rho_{soft} \in$   
581  $R^d$  to edit prompts in the continuous space as op-  
582 posed to the discrete token-space. Its main draw-  
583 back is requiring access to the parameters of LLMs,  
584 making it incompatible with API-based LLMs.

### 585 9.4 APO as a Precursor to LLM Finetuning

586 LLM Finetuning pursues instruction-alignment  
587 (Ouyang et al., 2022) and preference-alignment  
588 (Wang et al., 2024b) after pretraining. Works such  
589 as APEER (Jin et al., 2024), Self-Instruct (Wang  
590 et al., 2023), PATH (Xian et al., 2024) leveraged  
591 automatic prompt optimization to bootstrap quality  
592 synthetic data for improved LLM finetuning.

## 593 10 Challenges and Future Directions

### 594 10.1 Task-agnostic APO

595 All the surveyed APO methods assume that the task  
596 type  $T$  is known beforehand; additionally offline  
597 APO methods also require an evaluation set  $D_{val}$ ,  
598 something not explicitly available in production

599 settings. Barring a few tasks covered by Joko et al.  
600 (2024); Sun et al. (2024a); Zhang et al. (2022);  
601 Choi et al. (2024), inference-time optimization of  
602 multiple unknown tasks is underexplored. More  
603 robust evaluations are needed for task-agnostic  
604 APO systems combining seen and unseen tasks.

## 605 10.2 Unclear Mechanisms

606 Melamed et al. (2024) showed that prompts have  
607 so-called ‘evil twins’ that are uninterpretable yet  
608 recover some of the performance of gold-standard  
609 prompts. Lu et al. (2024) showed that rare gib-  
610 berish strings can serve as competitive delimiters  
611  $\tau$  in prompts. Yang et al. (2024b) showed that  
612 self-reflection by LLMs can suffer from incorrect  
613 error identification, prior biases, semantic invalidity,  
614 leading to failure in yielding improved prompts.  
615 More studies are needed to better uncover the mech-  
616 anisms of prompt optimization.

## 617 10.3 APO for System Prompts / Agents

618 Although SPRIG explored optimizing system  
619 prompts in chat-style settings, scalability remains  
620 a challenge - optimizing system prompts required  
621 a predefined corpus and close to 60 hours whereas  
622 Protegi only needed 10 minutes per task. Similarly,  
623 optimizing prompts for several components in an  
624 agentic system in a concurrent fashion poses an  
625 exciting direction for future research.

## 626 10.4 Multimodal APO

627 Recently, textual prompt optimization has ex-  
628 panded to multimodal domains: text-to-image (Liu  
629 et al., 2024b; Mañas et al., 2024; Liu et al., 2024d),  
630 text-to-video (Ji et al., 2024), text-to-audio (Huang  
631 et al., 2023), and text-image alignment models like  
632 CLIP (Du et al., 2024; Mirza et al., 2024). Be-  
633 yond textual prompts, Huang et al. (2023) explore  
634 optimizing multimodal inputs, such as images, to  
635 elicit better responses from large multimodal mod-  
636 els. However, the interplay between modalities in  
637 prompt optimization remains underexplored. Fu-  
638 ture research could develop APO frameworks to  
639 jointly optimize multimodal prompts (eg - remove  
640 background noise from audio, add visual markers  
641 to videos, etc.) to fully leverage their synergies.

## 642 11 Conclusion

643 In this paper, we provide a comprehensive fine-  
644 grained review of existing APO techniques and  
645 identified key areas for future growth. It is our aim  
646 to spur future research spawning from our survey.

## 647 12 Limitations

648 While we attempted to cover all qualifying papers,  
649 it is possible that we may have unintentionally  
650 missed out on some relevant papers. We also men-  
651 tion some of the papers that were excluded in this  
652 survey with specific reasons in section A.2. Also,  
653 we realize that fitting varied research works into a  
654 single unifying framework might risk broad cate-  
655 gorizations for some papers, or skipping some char-  
656 acteristics for others (e.g. Tempora (Zhang et al.,  
657 2022) consists of both RL-based and word/phrase-  
658 level editing techniques, applied to both instruc-  
659 tions and exemplars). In such cases, we categorize  
660 a paper based on its most salient features. Another  
661 challenge is that when presenting a survey paper  
662 under 8 pages, we had to make tradeoffs and only  
663 retain content in the main body that was deemed  
664 most necessary. This resulted in having to relegate  
665 a core contribution (Tables 2,3,4) which contained  
666 a rigorous comparison of all the surveyed papers  
667 into the appendix. We have attempted our best  
668 to strike the right balance between specificity and  
669 brevity to present a novel framework. We also pro-  
670 vide copious references to interested researchers  
671 for further reading.

## 672 References

673 Eshaan Agarwal, Joykirat Singh, Vivek Dani, Raghav  
674 Magazine, Tanuja Ganu, and Akshay Nambi. 2024.  
675 *Promptwizard: Task-aware prompt optimization frame-*  
676 *work.*

677 Fernando Alva-Manchego, Louis Martin, Antoine Bor-  
678 des, Carolina Scarton, Benoît Sagot, and Lucia Specia.  
679 2020. Asset: A dataset for tuning and evaluation of  
680 sentence simplification models with multiple rewriting  
681 transformations. In *Proceedings of the 58th Annual*  
682 *Meeting of the Association for Computational Linguistics*,  
683 pages 4668–4679.

684 Afra Amini, Tim Vieira, and Ryan Cotterell. 2024. *Di-*  
685 *rect preference optimization with an offset.* In *Findings*  
686 *of the Association for Computational Linguistics: ACL*  
687 *2024*, pages 9954–9972, Bangkok, Thailand. Associa-  
688 *tion for Computational Linguistics.*

689 R. Anantha, Svitlana Vakulenko, Zhucheng Tu, S. Long-  
690 pre, Stephen G. Pulman, and Srinivas Chappidi. 2020.  
691 *Open-domain question answering goes conversational*  
692 *via question rewriting.* In *North American Chapter of*  
693 *the Association for Computational Linguistics.*

694 Jacob Andreas, Johannes Bufe, David Burkett,  
695 Charles C. Chen, Joshua Clausman, Jean Crawford,  
696 Kate Crim, Jordan DeLoach, Leah Dorner, Jason Eis-  
697 ner, Hao Fang, Alan Guo, David Leo Wright Hall,  
698 Kristin Delia Hayes, Kellie Hill, Diana Ho, Wendy

Iwaszuk, Smriti Jha, Dan Klein, Jayant Krishnamurthy,  
Theo Lanman, Percy Liang, C. H. Lin, Ilya Lintsbakh,  
Andy McGovern, Aleksandr Nisnevich, Adam Pauls,  
Dmitrij Petters, Brent Read, Dan Roth, Subhro Roy,  
Jesse Rusak, Beth Ann Short, Div Slomin, B Snyder,  
Stephon Striplin, Yu Su, Zachary Tellman, Sam Thomson,  
A. A. Vorobev, Izabela Witoszko, Jason Wolfe,  
A. G. Wray, Yuchen Zhang, and Alexander Zotov. 2020.  
*Task-oriented dialogue as dataflow synthesis.* *Transac-*  
*tions of the Association for Computational Linguistics*,  
8:556–571.

Trapit Bansal, Rishikesh Jha, and Andrew McCallum.  
2019. *Learning to few-shot learn across diverse natural*  
*language classification tasks.* In *International Confer-*  
*ence on Computational Linguistics.*

Aman Bhargava, Cameron Witkowski, Shi-Zhuo Looi,  
and Matt Thomson. 2024. *What’s the magic word? a*  
*control theory of llm prompting.*

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng  
Gao, and Yejin Choi. 2019. *Piqa: Reasoning about*  
*physical commonsense in natural language.* In *AAAI*  
*Conference on Artificial Intelligence.*

Samuel R Bowman, Gabor Angeli, Christopher Potts,  
and Christopher D Manning. 2015. A large annotated  
corpus for learning natural language inference. *arXiv*  
*preprint arXiv:1508.05326.*

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen  
Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher  
Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,  
Scott Gray, Benjamin Chess, Jack Clark, Christopher  
Berner, Sam McCandlish, Alec Radford, Ilya Sutskever,  
and Dario Amodei. 2020. *Language models are few-*  
*shot learners.*

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang  
Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan,  
and Milica Gasic. 2018. Multiwoz-a large-scale multi-  
domain wizard-of-oz dataset for task-oriented dialogue  
modelling. In *Proceedings of the 2018 Conference on*  
*Empirical Methods in Natural Language Processing*,  
pages 5016–5026.

Daniel Matthew Cer, Mona T. Diab, Eneko Agirre, Iñigo  
Lopez-Gazpio, and Lucia Specia. 2017. *Semeval-2017*  
*task 1: Semantic textual similarity multilingual and*  
*crosslingual focused evaluation.* In *International Work-*  
*shop on Semantic Evaluation.*

Mauro Cettolo, Marcello Federico, Luisa Bentivogli,  
Niekus Jan, Stüker Sebastian, Sudoh Katsutomo,  
Yoshino Koichiro, and Federmann Christian. 2017.  
*Overview of the iwslt 2017 evaluation campaign.* In *Interna-*  
*tional Workshop on Spoken Language Translation.*

Yongchao Chen, Jacob Arkin, Yilun Hao, Yang Zhang,  
Nicholas Roy, and Chuhu Fan. 2024. *PRompt optimi-*  
*zation in multi-step tasks (PROMST): Integrating*

756	human feedback and heuristic-based sampling. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 3859–3920, Miami, Florida, USA. Association for Computational Linguistics.	811
757		812
758		813
759		
760		
761	Jiale Cheng, Xiao Liu, Kehan Zheng, Pei Ke, Hongning Wang, Yuxiao Dong, Jie Tang, and Minlie Huang. 2024. Black-box prompt optimization: Aligning large language models without model training. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3201–3219, Bangkok, Thailand. Association for Computational Linguistics.	814
762		815
763		816
764		
765		
766		
767		
768		
769	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <a href="https://vicuna.lmsys.org">https://vicuna.lmsys.org</a> (accessed 14 April 2023), 2(3):6.	817
770		818
771		819
772		820
773		
774		
775	Minje Choi, Jiaxin Pei, Sagar Kumar, Chang Shu, and David Jurgens. 2023. Do llms understand social knowledge? evaluating the sociability of large language models with socket benchmark. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 11370–11403.	821
776		822
777		823
778		824
779		
780		
781	Yunseon Choi, Sangmin Bae, Seonghyun Ban, Minchan Jeong, Chuheng Zhang, Lei Song, Li Zhao, Jiang Bian, and Kee-Eung Kim. 2024. Hard prompts made interpretable: Sparse entropy regularization for prompt tuning with rl.	825
782		826
783		827
784		828
785		829
786	Christopher Cieri, Mark Liberman, Sunghye Cho, Stephanie Strassel, James Fiumara, and Jonathan Wright. 2022. Reflections on 30 years of language resource development and sharing. In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 543–550, Marseille, France. European Language Resources Association.	830
787		
788		
789		
790		
791		
792		
793	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>ArXiv</i> , abs/1803.05457.	835
794		836
795		837
796		838
797		
798	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. <i>arXiv preprint arXiv:2110.14168</i> .	839
799		840
800		841
801		842
802		
803	Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world’s first truly open instruction-tuned llm. <i>Company Blog of Databricks</i> .	843
804		844
805		845
806		846
807		847
808	Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming Zhou. 2020. Mutual: A dataset for multi-turn dialogue reasoning. <i>ArXiv</i> , abs/2004.04494.	848
809		849
810		
811	Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In <i>Machine Learning Challenges Workshop</i> .	850
812		851
813		852
814		853
815		854
816		855
817	Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. 2019. The commitmentbank: Investigating projection in naturally occurring discourse.	856
818		857
819		858
820		
821	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.	859
822		860
823		861
824		
825	Franck Dernoncourt and Ji Young Lee. 2017. Pubmed 200k rct: a dataset for sequential sentence classification in medical abstracts. In <i>International Joint Conference on Natural Language Processing</i> .	862
826		863
827		864
828		865
829		866
830		
831	Robert C. Detrano, András Jánosi, Walter Steinbrunn, Matthias Emil Pfisterer, Johann-Jakob Schmid, Sarbjit Sandhu, Kern Guppy, Stella Lee, and Victor Froelicher. 1989. International application of a new probability algorithm for the diagnosis of coronary artery disease. <i>The American journal of cardiology</i> , 64 5:304–10.	867
832		868
833		869
834		
835	Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, Yong Lin, Xiao Zhou, and Tong Zhang. 2022. Black-box prompt learning for pre-trained language models. <i>arXiv preprint arXiv:2201.08531</i> .	870
836		871
837		872
838		873
839	Rezarta Islamaj Doğan, Robert Leaman, and Zhiyong Lu. 2014. Ncbi disease corpus: a resource for disease name recognition and concept normalization. <i>Journal of biomedical informatics</i> , 47:1–10.	874
840		875
841		876
842		
843	William B. Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In <i>International Joint Conference on Natural Language Processing</i> .	877
844		878
845		879
846		880
847		881
848		882
849		
850	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024a. A survey on in-context learning. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.	883
851		884
852		885
853		886
854		887
855		
856	Yihong Dong, Kangcheng Luo, Xue Jiang, Zhi Jin, and Ge Li. 2024b. PACE: Improving prompt with actor-critic editing for large language model. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 7304–7323, Bangkok, Thailand. Association for Computational Linguistics.	888
857		889
858		890
859		
860	Yingjun Du, Wenfang Sun, and Cees GM Snoek. 2024. Ipo: Interpretable prompt optimization for vision-language models. <i>arXiv preprint arXiv:2410.15397</i> .	891
861		892
862	Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In <i>North American Chapter of the Association for Computational Linguistics</i> .	893
863		894
864		

865	Stefan Daniel Dumitrescu, Petru Rebeja, Beáta Lőrincz, Mihaela Găman, Mihai Daniel Ilie, Andrei Pruteanu, Adriana Stan, Luciana Morogan, Traian Rebedea, and Sebastian Ruder. 2021. <b>Liro: Benchmark and leader-</b> <b>board for romanian language tasks.</b> In <i>NeurIPS</i> <i>Datasets and Benchmarks</i> .	920
866		921
867		922
868		923
869		924
870		
871	Ibrahim Abu Farha and Walid Magdy. 2020a. <b>From</b> <b>arabic sentiment analysis to sarcasm detection: The</b> <b>arsarcasm dataset.</b> In <i>OSACT</i> .	925
872		926
873		927
874		928
875	Ibrahim Abu Farha and Walid Magdy. 2020b. <b>From</b> <b>arabic sentiment analysis to sarcasm detection: The</b> <b>arsarcasm dataset.</b> In <i>OSACT</i> .	929
876		930
877		931
878	Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rock- täschel. 2023. <b>Promptbreeder: Self-referential</b> <b>self-improvement via prompt evolution.</b> <i>ArXiv</i> , abs/2309.16797.	932
879		933
880		934
881		
882	Rory A. Fisher. 1936. <b>The use of multiple measure-</b> <b>ments in taxonomic problems.</b> <i>Annals of Human Genet-</i> <i>ics</i> , 7:179–188.	935
883		936
884		937
885	Noa Garcia, Chentao Ye, Zihua Liu, Qingtao Hu, Mayu Otani, Chenhui Chu, Yuta Nakashima, and Teruko Mita- mura. 2020. A dataset and baselines for visual question answering on art. In <i>European Conference on Computer</i> <i>Vision</i> , pages 92–108.	938
886		
887		
888		
889		
890	Miguel Garc’ia-Orteg’on, Gregor N. C. Simm, Austin Tripp, José Miguel Hernández-Lobato, Andreas Bender, and Sergio Bacallado. 2021. <b>Dockstring: Easy</b> <b>molecular docking yields better benchmarks for ligand</b> <b>design.</b> <i>Journal of Chemical Information and Modeling</i> , 62:3486 – 3502.	939
891		940
892		941
893		942
894		943
895		944
896		945
897	Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. <b>Creating training</b> <b>corpora for nlg micro-planners.</b> In <i>Annual Meeting of</i> <i>the Association for Computational Linguistics</i> .	946
898		947
899		948
900		949
901		950
902	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. <b>Did aristotle use</b> <b>a laptop? a question answering benchmark with implicit</b> <b>reasoning strategies.</b> <i>Transactions of the Association</i> <i>for Computational Linguistics</i> , 9:346–361.	951
903		952
904		953
905		954
906		955
907	Bogdan Gliwa, Iwona Mochol, Maciej Bieseck, and Aleksander Wawer. 2019. Samsum corpus: A human- annotated dialogue dataset for abstractive summariza- tion. In <i>Proceedings of the 2nd Workshop on New Fron-</i> <i>ters in Summarization</i> , pages 70–79.	956
908		957
909		958
910		959
911		
912	Chulaka Gunasekara, Jonathan K. Kummersfeld, Lazaros Polymenakos, and Walter S. Lasecki. 2019. <b>Dstc7 task</b> <b>1: Noetic end-to-end response selection.</b> <i>Proceedings</i> <i>of the First Workshop on NLP for Conversational AI</i> .	960
913		
914	Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. 2024. <b>Connecting large language models with evolu-</b> <b>tionary algorithms yields powerful prompt optimizers.</b> In <i>The Twelfth International Conference on Learning</i> <i>Representations</i> .	961
915		962
916		963
917		964
918		
919	Han He, Qianchu Liu, Lei Xu, Chaitanya Shivade, Yi Zhang, Sundararajan Srinivasan, and Katrin Kirch- hoff. 2025. <b>Crispo: Multi-aspect critique-suggestion-</b> <b>guided automatic prompt optimization for text genera-</b> <b>tion.</b>	965
920		966
921		967
922		968
923		969
924		970
925	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Ja- cob Steinhardt. 2020. <b>Measuring massive multitask</b> <b>language understanding.</b> <i>ArXiv</i> , abs/2009.03300.	971
926		
927		
928		
929	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. In <i>Thirty-fifth Conference on</i> <i>Neural Information Processing Systems Datasets and</i> <i>Benchmarks Track (Round 2).</i>	972
930		973
931		974
932		975
933		976
934		977
935	Or Honovich, Uri Shaham, Samuel R. Bowman, and Omer Levy. 2022. <b>Instruction induction: From few</b> <b>examples to natural language task descriptions.</b> <i>ArXiv</i> , abs/2205.10782.	978
936		979
937		980
938		
939	Or Honovich, Uri Shaham, Samuel R. Bowman, and Omer Levy. 2023. <b>Instruction induction: From few</b> <b>examples to natural language task descriptions.</b> In <i>Pro-</i> <i>ceedings of the 61st Annual Meeting of the Association</i> <i>for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1935–1952, Toronto, Canada. Association for Computational Linguistics.	981
940		982
941		983
942		984
943		985
944		986
945		
946	Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. <b>Learning to solve</b> <b>arithmetic word problems with verb categorization.</b> In <i>Conference on Empirical Methods in Natural Language</i> <i>Processing</i> .	987
947		988
948		989
949		990
950		
951	Bairu Hou, Joe O’Connor, Jacob Andreas, Shiyu Chang, and Yang Zhang. 2023. <b>Promptboosting: black-box text</b> <b>classification with ten forward passes.</b> In <i>Proceedings of</i> <i>the 40th International Conference on Machine Learning</i> , ICML’23. JMLR.org.	991
952		992
953		993
954		994
955		995
956	Cho-Jui Hsieh, Si Si, Felix Yu, and Inderjit Dhillon. 2024. <b>Automatic engineering of long prompts.</b> In <i>Find-</i> <i>ings of the Association for Computational Linguistics: ACL</i> 2024, page 10672—10685, Bangkok, Thailand. Association for Computational Linguistics.	996
957		997
958		998
959		999
960		
961	Minqing Hu and Bing Liu. 2004. <b>Mining and sum-</b> <b>marizing customer reviews.</b> <i>Proceedings of the tenth</i> <i>ACM SIGKDD international conference on Knowledge</i> <i>discovery and data mining.</i>	1000
962		1001
963		1002
964		1003
965	Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. <b>Cosmos qa: Machine reading com-</b> <b>prehension with contextual commonsense reasoning.</b> In <i>Proceedings of the 2019 Conference on Empirical</i> <i>Methods in Natural Language Processing and the 9th</i> <i>International Joint Conference on Natural Language</i> <i>Processing (EMNLP-IJCNLP)</i> , pages 2391–2401.	1004
966		1005
967		1006
968		1007
969		1008
970		1009
971		
972	Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang	1010
973		1011

974	Yin, and Zhou Zhao. 2023. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. In <i>International Conference on Machine Learning</i> , pages 13916–13932. PMLR.	1029
975		1030
976		1031
977		1032
978		1033
979		1034
980		
981		
982	Yasaman Jafari, Dheeraj Mekala, Rose Yu, and Taylor Berg-Kirkpatrick. 2024. Morl-prompt: An empirical analysis of multi-objective reinforcement learning for discrete prompt optimization.	
983		
984		
985		
986		
987	Yatai Ji, Jiacheng Zhang, Jie Wu, Shilong Zhang, Shoufa Chen, Chongjian GE, Peize Sun, Weifeng Chen, Wenqi Shao, Xuefeng Xiao, et al. 2024. Prompt-a-video: Prompt your video diffusion model via preference-aligned llm. <i>arXiv preprint arXiv:2412.15156</i> .	
988		
989		
990		
991	Yichen Jiang, Shikha Bordia, Zheng Zhong, Charles Dognin, Maneesh Kumar Singh, and Mohit Bansal. 2020. Hover: A dataset for many-hop fact extraction and claim verification. In <i>Findings</i> .	
992		
993		
994		
995		
996	Can Jin, Hongwu Peng, Shiyu Zhao, Zhenting Wang, Wujiang Xu, Ligong Han, Jiahui Zhao, Kai Zhong, Sanguthevar Rajasekaran, and Dimitris N. Metaxas. 2024. Apeer: Automatic prompt engineering enhances large language model reranking.	
997		
998		
999		
1000	Di Jin, Eileen Pan, Nassim Oufattolle, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. <i>ArXiv</i> , abs/2009.13081.	
1001		
1002		
1003		
1004		
1005		
1006		
1007		
1008	Qiao Jin, Bhuvan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2567–2577.	
1009		
1010		
1011		
1012		
1013		
1014		
1015	Hideaki Joko, Shubham Chatterjee, Andrew Ramsay, Arjen P De Vries, Jeff Dalton, and Faegheh Hasibi. 2024. Doing personal laps: Llm-augmented dialogue construction for personalized multi-session conversational search. In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 796–806.	
1016		
1017		
1018		
1019	Gurusha Juneja, Nagarajan Natarajan, Hua Li, Jian Jiao, and Amit Sharma. 2024. Task facet learning: A structured approach to prompt optimization. <i>arXiv preprint arXiv:2406.10504</i> .	
1020		
1021		
1022		
1023		
1024	David Jurgens, Srijan Kumar, Raine Hoover, Daniel A. McFarland, and Dan Jurafsky. 2018. Measuring the evolution of a scientific field through citation frames. <i>Transactions of the Association for Computational Linguistics</i> , 6:391–406.	
1025		
1026		
1027		
1028	Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp. <i>arXiv preprint arXiv:2212.14024</i> .	
1029	Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2024. Dspy: Compiling declarative language model calls into self-improving pipelines.	
1030		
1031		
1032		
1033		
1034		
1035	Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, D. Corney, Benno Stein, and Martin Potthast. 2019. Semeval-2019 task 4: Hyperpartisan news detection. In <i>International Workshop on Semantic Evaluation</i> .	
1036		
1037		
1038		
1039		
1040	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners.	
1041		
1042		
1043	Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. 2015. Parsing algebraic word problems into equations. <i>Transactions of the Association for Computational Linguistics</i> , 3:585–597.	
1044		
1045		
1046		
1047		
1048	Weize Kong, Spurthi Amba Hombaiah, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. 2024. Prewrite: Prompt rewriting with reinforcement learning.	
1049		
1050		
1051	Shanu Kumar, Akhila Yesantara Venkata, Shubhangshu Khandelwal, Bishal Santra, Parag Agrawal, and Manish Gupta. 2024. Sculpt: Systematic tuning of long prompts.	
1052		
1053		
1054		
1055	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	
1056		
1057		
1058		
1059		
1060		
1061	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. Race: Large-scale reading comprehension dataset from examinations. <i>arXiv preprint arXiv:1704.04683</i> .	
1062		
1063		
1064		
1065	Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent retrieval for weakly supervised open domain question answering. <i>ArXiv</i> , abs/1906.00300.	
1066		
1067		
1068	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
1069		
1070		
1071		
1072		
1073		
1074	Hector J. Levesque, Ernest Davis, and L. Morgenstern. 2011. The winograd schema challenge. In <i>AAAI Spring Symposium: Logical Formalizations of Commonsense Reasoning</i> .	
1075		
1076		
1077		
1078	Bei Li, Rui Wang, Junliang Guo, Kaitao Song, Xu Tan, Hany Hassan, Arul Menezes, Tong Xiao, Jiang Bian, and JingBo Zhu. 2023a. Deliberate then generate: Enhanced prompting framework for text generation.	
1079		
1080		
1081		

1082	Cheng Li, Jindong Wang, Yixuan Zhang, Kajie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. 2023b. Large language models understand and can be enhanced by emotional stimuli. <i>arXiv preprint arXiv:2307.11760</i> .	Xiaogeng Liu, Zhiyuan Yu, Yizhe Zhang, Ning Zhang, and Chaowei Xiao. 2024c. Automatic and universal prompt injection attacks against large language models.	1138
1083			1139
1084			1140
1085			
1086			
1087	Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. 2010. A contextual-bandit approach to personalized news article recommendation. In <i>Proceedings of the 19th International Conference on World Wide Web, WWW ’10</i> , page 661–670, New York, NY, USA. Association for Computing Machinery.	Yilun Liu, Minggui He, Feiyu Yao, Yuhe Ji, Shimin Tao, Jingzhou Du, Duan Li, Jian Gao, Li Zhang, Hao Yang, et al. 2024d. What do you want? user-centric prompt generation for text-to-image synthesis via multi-turn guidance. <i>arXiv preprint arXiv:2408.12910</i> .	1141
1088			1142
1089			1143
1090			1144
1091			1145
1092			
1093	Moxin Li, Wenjie Wang, Fuli Feng, Yixin Cao, Jizhi Zhang, and Tat-Seng Chua. 2023c. Robust prompt optimization for large language models against distribution shifts. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 1539–1554, Singapore. Association for Computational Linguistics.	Xuan Do Long, Yiran Zhao, Hannah Brown, Yuxi Xie, James Xu Zhao, Nancy F. Chen, Kenji Kawaguchi, Michael Shieh, and Junxian He. 2024. Prompt optimization via adversarial in-context learning. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7308–7327, Bangkok, Thailand. Association for Computational Linguistics.	1146
1094			1147
1095			1148
1096			1149
1097			1150
1098			1151
1099			1152
1100	Zekun Li, Baolin Peng, Pengcheng He, Michel Galley, Jianfeng Gao, and Xifeng Yan. 2023d. Guiding large language models via directional stimulus prompting. <i>arXiv preprint arXiv:2302.11520</i> .	Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. In <i>SIGDIAL Conference</i> .	1153
1101			1154
1102			1155
1103			1156
1104	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252.	Junru Lu, Siyu An, Min Zhang, Yulan He, Di Yin, and Xing Sun. 2025. FIPO: Free-form instruction-oriented prompt optimization with preference dataset and modular fine-tuning schema. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , page 11029—11047, Abu Dhabi, UAE. Association for Computational Linguistics.	1158
1105			1159
1106			1160
1107			1161
1108			1162
1109	Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In <i>Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13</i> , pages 740–755. Springer.	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	1163
1110			1164
1111			
1112	Xiaoqiang Lin, Zhongxiang Dai, Arun Verma, See-Kiong Ng, Patrick Jaillet, and Bryan Kian Hsiang Low. 2024. Prompt optimization with human feedback.	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 <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , page 2221—2231, Mexico City, Mexico. Association for Computational Linguistics.	1165
1113			1166
1114			1167
1115			1168
1116	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction. <i>ArXiv</i> , abs/1808.09602.	1169
1117			1170
1118			1171
1119	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024a. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies</i> , pages 142–150.	1172
1120			1173
1121			1174
1122			1175
1123			1176
1124	Shengcai Liu, Caishun Chen, Xinghua Qu, Ke Tang, and Yew Soon Ong. 2023. Large language models as evolutionary optimizers. <i>2024 IEEE Congress on Evolutionary Computation (CEC)</i> , pages 1–8.	Oscar Mañas, Pietro Astolfi, Melissa Hall, Candace Ross, Jack Urbanek, Adina Williams, Aishwarya Agrawal, Adriana Romero-Soriano, and Michal Drozdzal. 2024. Improving text-to-image consistency via automatic prompt optimization. <i>arXiv preprint arXiv:2403.17804</i> .	1177
1125			1178
1126			
1127			
1128			
1129	Shihong Liu, Samuel Yu, Zhiqiu Lin, Deepak Pathak, and Deva Ramanan. 2024b. Language models as black-box optimizers for vision-language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 12687–12697.		
1130			
1131			
1132			
1133			
1134			
1135			
1136			
1137			

1195	Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. The natural language de-	Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.	1250
1196	cathlon: Multitask learning as question answering.	<i>arXiv preprint arXiv:1806.08730</i> .	1251
1197			1252
1198			1253
1199	Rimon Melamed, Lucas H. McCabe, Tanay Wakhare, Yejin Kim, H. Howie Huang, and Enric Boix-Adsera. 2024. <b>Prompts have evil twins</b> .		
1200			1254
1201			1255
1202	M Jehanzeb Mirza, Mengjie Zhao, Zhuoyuan Mao, Sivan Doveh, Wei Lin, Paul Gavrikov, Michael Dorkenwold, Shiqi Yang, Saurav Jha, Hiromi Wakaki, et al. 2024. Glov: Guided large language models as implicit optimizers for vision language models.	<i>arXiv preprint arXiv:2410.06154</i> .	1256
1203			1257
1204			1258
1205			
1206			
1207			
1208	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. <b>Cross-task generalization via natural language crowdsourcing instructions</b> . In <i>Annual Meeting of the Association for Computational Linguistics</i> .		
1209			1259
1210			1260
1211			1261
1212			1262
1213	Ioannis Mollas, Zoe Chrysopoulou, Stamatis Karlos, and Grigorios Tsoumacas. 2020. Ethos: an online hate speech detection dataset.	<i>arXiv preprint arXiv:2006.08328</i> .	1263
1214			1264
1215			
1216			
1217	Ramesh Nallapati, Bowen Zhou, Cícero Nogueira dos Santos, Çağlar Gülcehre, and Bing Xiang. 2016. <b>Abstractive text summarization using sequence-to-sequence rnns and beyond</b> . In <i>Conference on Computational Natural Language Learning</i> .		
1218			1265
1219			1266
1220			1267
1221			
1222	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. <b>Don't give me the details, just the summary!</b> topic-aware convolutional neural networks for extreme summarization.	<i>ArXiv</i> , abs/1808.08745.	
1223			1268
1224			1269
1225			1270
1226	Ehsan Nezhadarya, Yang Liu, and Bingbing Liu. 2019. Boxnet: A deep learning method for 2d bounding box estimation from bird's-eye view point cloud. In <i>2019 IEEE Intelligent Vehicles Symposium (IV)</i> , pages 1557–1564. IEEE.		1271
1227			
1228			
1229			
1230			
1231	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2019. <b>Adversarial nli: A new benchmark for natural language understanding</b> .	<i>ArXiv</i> , abs/1910.14599.	
1232			1272
1233			1273
1234			1274
1235	Jekaterina Novikova, Ondrej Dusek, and Verena Rieser. 2017. <b>The e2e dataset: New challenges for end-to-end generation</b> .	<i>ArXiv</i> , abs/1706.09254.	1275
1236			
1237			
1238	Krista Opsahl-Ong, Michael J Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia, and Omar Khattab. 2024. <b>Optimizing instructions and demonstrations for multi-stage language model programs</b> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , page 9340—9366, Miami, Florida, USA. Association for Computational Linguistics.		
1239			1276
1240			1277
1241			1278
1242			1279
1243			1280
1244			1281
1245			1282
1246	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,		
1247			1283
1248			1284
1249			1285
1249	Ye Qi, Devendra Singh Sachan, Matthieu Felix, Sarguna Padmanabhan, and Graham Neubig. 2018. <b>When and why are pre-trained word embeddings useful for neural machine translation?</b> <i>ArXiv</i> , abs/1804.06323.		
1249			1299
1249			1300
1249			1301
1249			1302

1303	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <b>Squad: 100,000+ questions for machine comprehension of text.</b> In <i>Conference on Empirical Methods in Natural Language Processing</i> .	1359
1304		1360
1305		1361
1306		1362
1307		1363
1308	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. <b>Gpqa: A graduate-level google-proof q&amp;a benchmark.</b> <i>ArXiv</i> , abs/2311.12022.	1364
1309		1365
1310		1366
1311		1367
1312	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In <i>2011 AAAI spring symposium series</i> .	1368
1313		
1314		
1315		
1316	Subhro Roy and Dan Roth. 2016. <b>Solving general arithmetic word problems.</b> <i>ArXiv</i> , abs/1608.01413.	1369
1317		1370
1318		1371
1319		1372
1320		1373
1321		1374
1322		1375
1323		
1324	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social iqa: Commonsense reasoning about social interactions. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4463–4473.	
1325		
1326		
1327	Tobias Schnabel and Jennifer Neville. 2024. <b>Symbolic prompt program search: A structure-aware approach to efficient compile-time prompt optimization.</b>	
1328		
1329		
1330		
1331		
1332		
1333		
1334	Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. <i>Advances in Neural Information Processing Systems</i> , 35:25278–25294.	
1335		
1336		
1337		
1338		
1339	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In <i>The Twelfth International Conference on Learning Representations</i> .	
1340		
1341		
1342		
1343		
1344	Jingyuan Selena She, Christopher Potts, Sam Bowman, and Atticus Geiger. 2023. <b>Scone: Benchmarking negation reasoning in language models with fine-tuning and in-context learning.</b> In <i>Annual Meeting of the Association for Computational Linguistics</i> .	
1345		
1346		
1347	Zeru Shi, Zhenting Wang, Yongye Su, Weidi Luo, Fan Yang, and Yongfeng Zhang. 2024. <b>Robustness-aware automatic prompt optimization.</b>	
1348		
1349		
1350		
1351		
1352		
1353		
1354	Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. <b>AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts.</b> In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4222–4235, Online. Association for Computational Linguistics.	
1355		
1356		
1357	Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. <b>Reflexion: Language agents with verbal reinforcement learning.</b>	
1358		
	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2024. <b>Reflexion: Language agents with verbal reinforcement learning.</b> <i>Advances in Neural Information Processing Systems</i> , 36.	1359
		1360
		1361
		1362
		1363
	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2020. <b>Alfworld: Aligning text and embodied environments for interactive learning.</b> <i>arXiv preprint arXiv:2010.03768</i> .	1364
		1365
		1366
		1367
		1368
	Ankita Sinha, Wendi Cui, Kamalika Das, and Jiaxin Zhang. 2024. <b>Survival of the safest: Towards secure prompt optimization through interleaved multi-objective evolution.</b> In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 1016–1027, Miami, Florida, US. Association for Computational Linguistics.	1369
		1370
		1371
		1372
		1373
		1374
		1375
	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 conference on empirical methods in natural language processing</i> , pages 1631–1642.	1376
		1377
		1378
		1379
		1379
		1380
		1381
	Gizem Sogancioglu, Hakime Öztürk, and Arzucan Özgür. 2017. <b>Biosses: a semantic sentence similarity estimation system for the biomedical domain.</b> <i>Bioinformatics</i> , 33:i49 – i58.	1382
		1383
		1384
		1385
	Alessandro Sordoni, Eric Yuan, Marc-Alexandre Côté, Matheus Pereira, Adam Trischler, Ziang Xiao, Arian Hosseini, Friederike Niedtner, and Nicolas Le Roux. 2023. <b>Joint prompt optimization of stacked llms using variational inference.</b> In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 58128–58151. Curran Associates, Inc.	1386
		1387
		1388
		1389
		1390
		1391
		1392
	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>arXiv preprint arXiv:2206.04615</i> .	1393
		1394
		1395
		1396
		1397
		1398
	Hao Sun, Alihan Hünük, and Mihaela van der Schaar. 2024a. <b>Query-dependent prompt evaluation and optimization with offline inverse RL.</b> In <i>The Twelfth International Conference on Learning Representations</i> .	1399
		1400
		1401
		1402
	Hong Sun, Xue Li, Yinchuan Xu, Youkow Homma, Qi Cao, Min Wu, Jian Jiao, and Denis Charles. 2023. <b>Autohint: Automatic prompt optimization with hint generation.</b> <i>arXiv preprint arXiv:2307.07415</i> .	1403
		1404
		1405
		1406
	Jingwei Sun, Ziyue Xu, Hongxu Yin, Dong Yang, Daguang Xu, Yudong Liu, Zhixu Du, Yiran Chen, and Holger R. Roth. 2024b. <b>Fedbpt: efficient federated black-box prompt tuning for large language models.</b> In <i>Proceedings of the 41st International Conference on Machine Learning, ICML’24. JMLR.org</i> .	1407
		1408
		1409
		1410
		1411
		1412

1413	Mirac Suzgun, Nathan Scales, Nathanael Schärlí, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. 2023. Challenging big-bench tasks and whether chain-of-thought can solve them. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 13003–13051.	1467
1414		1468
1415		1469
1416		1470
1417		1471
1418		1472
1419		1473
1420		1474
1421	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. <i>ArXiv</i> , abs/1811.00937.	1475
1422		1476
1423		1477
1424		
1425	Prashant Trivedi, Souradip Chakraborty, Avinash Reddy, Vaneet Aggarwal, Amrit Singh Bedi, and George K. Atia. 2025. Align-pro: A principled approach to prompt optimization for llm alignment.	1478
1426		1479
1427		1480
1428		1481
1429	Nirali Vaghani and Mansi Thummar. 2023. Flipkart product reviews with sentiment dataset.	1482
1430		
1431	Ellen M Voorhees and Dawn M Tice. 2000. Building a question answering test collection. In <i>Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval</i> , pages 200–207.	1483
1432		1484
1433		1485
1434		1486
1435		
1436	Xingchen Wan, Ruoxi Sun, Hootan Nakhost, and Serkan O. Arik. 2024. Teach better or show smarter? on instructions and exemplars in automatic prompt optimization.	1487
1437		1488
1438		1489
1439		1490
1440	Ruochen Wang, Sohyun An, Minhao Cheng, Tianyi Zhou, Sung Ju Hwang, and Cho-Jui Hsieh. 2025. One prompt is not enough: automated construction of a mixture-of-expert prompts. In <i>Proceedings of the 41st International Conference on Machine Learning, ICML’24</i> . JMLR.org.	1491
1441		1492
1442		1493
1443		1494
1444		1495
1445		
1446	Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022a. Scienceworld: Is your agent smarter than a 5th grader? In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11279–11298.	1496
1447		1497
1448		1498
1449		1499
1450		
1451	William Yang Wang. 2017. “liar, liar pants on fire”: A new benchmark dataset for fake news detection. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	1500
1452		1501
1453		1502
1454		
1455	Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric P. Xing, and Zhitong Hu. 2024a. Promptagent: Strategic planning with language models enables expert-level prompt optimization. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	1503
1456		1504
1457		1505
1458		1506
1459		1507
1460		1508
1461		
1462	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022b. Self-instruct: Aligning language models with self-generated instructions. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	1509
1463		1510
1464		1511
1465		1512
1466		1513
1467		
1468	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.	1514
1469		1515
1470		1516
1471		
1472	Yushi Wang, Jonathan Berant, and Percy Liang. 2015. Building a semantic parser overnight. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	1517
1473		1518
1474		1519
1475		1520
1476		1521

1522	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024a. Large language models as optimizers.	Lechen Zhang, Tolga Ergen, Lajanugen Logeswaran, Moontae Lee, and David Jurgens. 2024b. Sprig: Improving large language model performance by system prompt optimization. <i>ArXiv</i> , abs/2410.14826.	1575
1523			1576
1524			1577
1525	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2024b. Large language models as optimizers. In <i>The Twelfth International Conference on Learning Representations</i> .	Tianjun Zhang, Xuezhi Wang, Denny Zhou, Dale Schuurmans, and Joseph E. Gonzalez. 2022. Tempera: Test-time prompting via reinforcement learning.	1578
1526			1579
1527			1580
1528			1581
1529	Muchen Yang, Moxin Li, Yongle Li, Zijun Chen, Chongming Gao, Junqi Zhang, Yangyang Li, and Fuli Feng. 2024c. Dual-phase accelerated prompt optimization. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 12163–12173, Miami, Florida, USA. Association for Computational Linguistics.	Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In <i>International Conference on Learning Representations</i> .	1582
1530			1583
1531			1584
1532			1585
1533			1586
1534			1587
1535			1588
1536	Sheng Yang, Yurong Wu, Yan Gao, Zineng Zhou, Bin Benjamin Zhu, Xiaodi Sun, Jian-Guang Lou, Zhiming Ding, Anbang Hu, Yuan Fang, et al. 2024d. Ampo: Automatic multi-branched prompt optimization. <i>arXiv preprint arXiv:2410.08696</i> .	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems</i> , NIPS '23, Red Hook, NY, USA. Curran Associates Inc.	1589
1537			1590
1538			1591
1539			1592
1540			1593
1541	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Han Zhou, Xingchen Wan, Ivan Vulić, and Anna Korhonen. 2023. Survival of the most influential prompts: Efficient black-box prompt search via clustering and pruning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 13064–13077, Singapore. Association for Computational Linguistics.	1594
1542			1595
1543			1596
1544			1597
1545			1598
1546			1599
1547	Qinyuan Ye, Maxamed Axmed, Reid Pryzant, and Fereshte Khani. 2024. Prompt engineering a prompt engineer.	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.	1600
1548			1601
1549			1602
1550	Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. 2024. Textgrad: Automatic "differentiation" via text.		
1551			
1552			
1553	John M. Zelle and Raymond J. Mooney. 1996. Learning to parse database queries using inductive logic programming. In <i>AAAI/IAAI, Vol. 2</i> .		
1554			
1555			
1556	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800.		
1557			
1558			
1559			
1560			
1561	Pengwei Zhan, Zhen Xu, Qian Tan, Jie Song, and Ru Xie. 2024. Unveiling the lexical sensitivity of llms: Combinatorial optimization for prompt enhancement. In <i>Conference on Empirical Methods in Natural Language Processing</i> .		
1562			
1563			
1564			
1565			
1566	Chenrui Zhang, Lin Liu, Chuyuan Wang, Xiao Sun, Hongyu Wang, Jinpeng Wang, and Mingchen Cai. 2024a. Prefer: prompt ensemble learning via feedback-reflect-refine. In <i>Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence</i> , AAAI'24/IAAI'24/EAAI'24. AAAI Press.		
1567			
1568			
1569			
1570			
1571			
1572			
1573			
1574			

## A Appendix

### A.1 Notation

We now define the notation of key terms and expressions used throughout the paper.

1.  $T$  = Task type,  $I$  = Task instruction,  $E = (xi, yi)_{i=1}^e$  Few shot demonstrations in the prompt,  $\tau$  = Template delimiters,  $z$  = CoT recipe for a task-instance,  $z_i \in I_i$
2.  $M_{task}$  target model,  $M_{APO}$  APO system
3.  $\rho = concat([s_1, s_2, \dots, s_m]) = concat(I, \tau, E)$  Prompt composed of m sentences, which comprise of Instruction, template delimiters and few-shot demonstrations.
4.  $D = \{(xi, yi)\}_{i=1}^m$  collection of m input-output pairs.  $D_{val}$  is the validation set used to validate prompt performance,  $D_{train}$  is the training set used to finetune the language model(Reprompting).
5.  $\{f_1, f_2, \dots\} \in F$  metric function upon which to evaluate task-prompt performance
6.  $r : S \times A \rightarrow R$  = reward model score, where  $S$  is the state-space and  $A$  is the action-space
7.  $|V|$  = length of vocabulary
8.  $\phi : S \in V_* \rightarrow R_d$  embedding function which takes in a sentence generated as a finite sequence of tokens belonging to a vocabulary  $V$ , and generating a floating point array representation of dimension  $d$
9.  $\rho_* = argmax_{\rho \in V_*} E_{D_{val}}[f_i(\rho)]$  The best performing prompt based on the metric score on validation set
10.  $k$  = number of candidates for top-K search,  $B$  = Beam width for beam search,  $N$  = number of iterations for search
11.  $C$  = number of experts in a Mixture of Experts approach (MOP),  $\mu_C$  = cluster centroid of cluster C (MOP).
12.  $LLM_{target}$  = target model which will be used for inference,  $LLM_{rewriter}$  = rewriter model which will be used for rewriter,  $LLM_{evaluator}$  = evaluator model which provides the LLM feedback to prompts / responses or both
13.  $\lambda$  with subscripts to denote different latency types:  $\lambda_t$  = Total training cost/latency, including all offline costs for data collection, preprocessing, and model fine-tuning,  $\lambda_i$  = per-example inference latency,  $\lambda_m$  = MLM inference latency per-example

### A.2 Excluded works

**FedBPT** (Sun et al., 2024b) used federated learning to update soft prompts and not discrete tokens. **Deliberate-then-generate** (Li et al., 2023a) randomly sampled arbitrary noisy inference and prompted the task LLM to deliberate on the wrong inference, while **Reflexion** (Shinn et al., 2023) agents maintain an episodic buffer of past deliberations. Neither method optimizes the input prompt. **AutoPrompt** (Shin et al., 2020) required gradient access to the task LLM and therefore doesn't remain blackbox.

### A.3 UCB based selection algorithm

1641

---

#### Algorithm 2 $Select(\cdot)$ with UCB Bandits

---

**Require:**  $n$  prompts  $\rho_1, \dots, \rho_n$ , dataset  $\mathcal{D}_{val}$ ,  $T$  time steps, metric function  $m$

1: Initialize:  $N_t(\rho_i) \leftarrow 0$  for all  $i = 1, \dots, n$

2: Initialize:  $Q_t(\rho_i) \leftarrow 0$  for all  $i = 1, \dots, n$

3: **for**  $t = 1, \dots, T$  **do**

4:     Sample uniformly  $\mathcal{D}_{sample} \subset \mathcal{D}_{val}$

5:      $\rho_i \leftarrow \arg \max_{\rho} \left\{ \frac{Q_t(\rho)}{N_t(\rho_i)} + c \sqrt{\frac{\log t}{N_t(\rho)}} \right\}$

6:     Observe reward  $r_{i,t} = m(\rho_i, \mathcal{D}_{sample})$

7:      $N_t(\rho_i) \leftarrow N_t(\rho_i) + |\mathcal{D}_{sample}|$

8:      $Q_t(\rho_i) \leftarrow Q_t(\rho_i) + r_{i,t}$

9: **return**  $SelectTop_b(Q_T/N_T)$

---

## B Comparison of different approaches + Tasks

1642

### B.1 Comparison

1643

Below we offer a comprehensive comparison of all the surveyed methods against our framework, covering the following aspects

1644

1645

1. Seed instructions	1646
2. Inference evaluation	1647
3. Candidate generation	1648
4. Search+filter strategy	1649
5. Iteration depth	1650
6. Optimization time complexity	1651
7. Prompt generation model	1652
8. Target models	1653

SNo.	Method	Seed instructions	Inference evaluation	Candidate generation	Search+filter strategy	Iteration depth	Optimization time complexity	Prompt generation model	Target models
1	GPS (Xu et al., 2022)	Manually created	Task accuracy	Genetic Algorithm: Back-translation, Cloze, Sentence continuation	Metaheuristic ensemble	Fixed	$O(T * N * k * \lambda_i)$		T0
2	GRIPS (Prasad et al., 2023)	Manually created	Entropy-based score+ Task accuracy	Phrase level add/remove/swap/paraphrase	TopK selection	Fixed	$O(k * N *  D_{val}  * B)$	PEGASUS paraphrase model	InstructGPT
3	Instruction induction (Honovich et al., 2023)	Instruction induction	Accuracy + BERTScore	LLM-rewriter		Fixed	$O( p  * \lambda_i)$	InstructGPT, GPT-3	InstructGPT, GPT-3
4	RLPrompt (Deng et al., 2022)	Manually created	Task accuracy + Reward model score	RL-based trained NN	TopK selection	Fixed	$O(N * \rho *  V * \lambda_i)$	RoBERTa-large	1/ BERT, 2/ GPT-2
5	TEMPERA (Zhang et al., 2022)	Manually created	Task accuracy	RL-trained NN		Fixed	$O(N * k *  V * C )$	RoBERTa-large	RoBERTa-large
6	AELP (Hsieh et al., 2024)	Manually created	Task accuracy	Genetic algorithm: LLM-mutator	Beam search	Fixed	$O(N * \rho * k *  D  * \lambda_i)$	PaLM 2-L	PaLM text-bison
7	APF (Zhou et al., 2022)	Instruction induction	Task accuracy	No new candidates	TopK selection	Fixed	$O(N * k *  D_{val}  * \lambda_i)$	InstructGPT, GPT-3, T5, InstructGPT	InstructGPT, GPT-3
8	AutoHint (Sun et al., 2023)	Manually created	Task accuracy + LLM-feedback	LLM rewriter	TopK selection	Fixed	$O(T *  D  * \lambda_i)$		GPT-4
9	BDPL (Diao et al., 2022)	Manually created	Task accuracy	RL-trained NN	TopK selection	Variable	$O(N * k * \lambda_i)$	RoBERTa, GPT-3	RoBERTa, GPT-3
10	Boosted Prompting (Pitis et al., 2023)	Instruction-induction	Task accuracy	Ensemble based method	TopK selection	Variable	$O(N * k * \lambda_i)$	text-curie-001, text-curie-003, GPT-3.5, code-davinci-002	text-curie-001, text-curie-003, GPT-3.5, code-davinci-002
11	BPO (Cheng et al., 2024)	Manually created	LLMaj (pairwise)	Finetuned LLMs			$O(\lambda_t +  val  * \lambda_i)$	Llama2-7b-chat	Vicuna-7b-v1.3, llama-1-7b, llama-1-13b
12	CLAPS (Zhou et al., 2023)	Manually created	Entropy-based score+ Task accuracy	Genetic Algorithm: Mutation + Crossover	TopK selection	Variable	$O(N * k *  V  * \lambda_i)$	Flan-T5	Flan-T5
13	Directional-stimulus (Li et al., 2023d)	Manually created	BLEU, BERTScore	RL-trained NN		Variable	$O(\lambda_t)$	T5, GPT-2	Flan-T5 large and base
14	DLN (Sordoni et al., 2023)	Manually created	Task accuracy + NLL	LLM mutator	TopK selection	Fixed	$O(N * k *  D_{train} )$	GPT-3 (text-davinci-003), GPT-4	ChatGPT, Codex, InstructGPT
15	DSP (Khattab et al., 2022)	Instruction induction	Task accuracy	Program Synthesis	TopK selection	Fixed	$O(N * k * \lambda_i)$	GP-3.5	GPT-3 (text-davinci-003), GPT-4
16	DSPy (Khattab et al., 2024)	Manually created + Instruction Induction	Task accuracy + LLM-feedback	Program Synthesis	TopK selection	Variable	$O(N * k * B * \lambda_i)$	Retrieval: CoBERTv2	LM: GPT-3.5
17	GATE (Joko et al., 2024)	Manually created	Human feedback	LLM rewriter		Open-ended	$O(N * (\lambda_m +  D_{val}  * \lambda_i))$	GPT-4	GPT-4
18	GPO (Li et al., 2023e)	Instruction induction	Task-Accuracy and FI	Metaprompt-design	TopK selection		$O(N * C *  V * B * E )$	gpt-3.5-turbo-0301	gpt-3.5-turbo-0301
19	PACE (Dong et al., 2024b)	Manually created	NLL + Task accuracy - BLEU and BERTScore	LLM-rewriter	TopK selection	<3	$O(N *  \rho  *  D_{val} )$	gpt-3.5-turbo (0301)	text-davinci-003, (gpt-3.5-turbo), GPT-4
20	PREFER (Zhang et al., 2024a)	Manually created	Task accuracy	LLM-rewriter + Ensemble method	TopK selection	Fixed	$O(N *  \rho  *  D_{val} )$	ChatGPT	ChatGPT

Table 2: Comparison of all APO techniques based on our framework

SNo.	Method	Seed instructions	Inference evaluation	Candidate generation	Search+filter strategy	Iteration depth	Optimization time complexity	Prompt generation model	Target models
21	Promptagent (Wang et al., 2024a)	Manually created	Task accuracy + LLM-feedback	LLM rewriter	UCT-based bandit-search	Fixed	$O(N * k * \lambda_i)$	GPT-4	GPT-3.5, GPT-4, PalM-2
22	Promphosting (Hou et al., 2023)	Instruction-induction	Accuracy, F1 Score	Ensemble based method	Beam-search	Early Stopping	$O(\lambda_m)$	T5	RoBERTa-large
23	Prompbreeder (Fernando et al., 2023)	Manually created	LLM Feedback + Task accuracy	Genetic Algorithm: Mutate + Crossover (LLM-edits)	Metaheuristic Ensemble	Fixed	$O(*N *  V  * \lambda_i)$	text-davinci-003, PalM 2-L	text-davinci-003, PalM 2-L
24	ProTeGi (Pryzant et al., 2023)	Manually created	Task accuracy + LLM-feedback	LLM rewriter	UCT-based bandit-search	Fixed	$O(N * C *  D_{val}  * \lambda_i)$	GPT-3.5-Turbo	GPT-3.5-turbo
25	Random separators (Lu et al., 2024)	Manually created	Task accuracy	LLM-rewriter	TopK selection	Fixed steps	$O(N * k * \lambda_i)$	GPT12 Large, GPT2 XL, Mistral 7B, Mistral 7B Instruct, Llama-Alpaca (B, Llama2 7B, Llama2 7B Chat, ChatGPT)	GPT12 Large, GPT2 XL, Mistral 7B, Mistral 7B Instruct, Llama-Alpaca (B, Llama2 7B, Llama2 7B Chat, ChatGPT)
26	ABO (Yang et al., 2024b)	Manually created + Instruction Induction	Task accuracy + LLM-feedback	LLM-rewriter	TopK selection	Fixed Steps	$O(B * N * \lambda_i)$	GPT-4	GPT-3.5-Turbo, Llama-2-70B-chat
27	Adv-ICL (Long et al., 2024)	Manually created	LLM Feedback	LLM-rewriter	Top-1 selection	Fixed	$O(N * k * \lambda_i)$	text-davinci-002, vicuna, ChatGPT	text-davinci-002, vicuna, ChatGPT
28	AMPO (Yang et al., 2024d)	Manually created	Task accuracy + F1 score	Coverage-based	TopK selection	Variable	$O(N * C * \lambda_i)$	GPT-4-turbo	GPT-4-turbo
29	APEER (Jin et al., 2024)	Manually created	Task accuracy-nDCG	Feedback + preference optimization	Used 3 epochs		$O(N *  \rho  *  D_{val} )$	GPT4, GPT3.5, Llama3, Qwen2	
30	APOHF (Lin et al., 2024)	Manually created	Task accuracy + Human Feedback	LLM rewriter	Linear UCB	Fixed	$O(N * T)$	ChatGPT	DALLE-3, ChatGPT
31	BATPromp! (Shi et al., 2024)	Manually created	Task accuracy + LLM-feedback	LLM rewriter	TopK selection	Fixed	$O(N *  D  *    * \lambda_i)$	GPT-3.5-turbo	GPT-3.5-turbo,
32	COPE (Zhan et al., 2024)	Manually created	Task accuracy	Token edits using MLM	Variable		$O(N *  I  * k *  D_{val}  * \lambda_i)$	RoBERTa (filling masked tokens)	GPT-4-mini, Llama2-7b Llama-2-7B-chat Misral-7B-instruct-v0.1, ChatGPT (gpt-3.5-turbo-0125)
33	CRISPO (He et al., 2025)	Manually created	ROUGE-l/2L F-measure, AlignScore	LLM rewriter	TOP-K greedy search	Fixed	$O(N * k * (D_{train} * \lambda_0 + \lambda_m))$	Claude Instant, Claude 3 Sonnet, Misral 7B, Llama3.5 SB	Claude 3 Sonnet, Misral 7B, Llama3.5 SB
34	DAPO (Yang et al., 2024c)	Manually created	Task accuracy	LLM-rewriter	Top-1 selection	Fixed	$O(N * k * \lambda_i)$	GPT-3.5-Turbo, Baichuan2, GPT-4	GPT-3.5-Turbo, Baichuan2, GPT-4
35	DRPO (Amini et al., 2024)	Manually created	Reward model score + LLM Feedback	LLM rewriter	Beam search	Fixed	$O(B * k * N)$	Mistral 7B, Mistral 7B (Instruct), Llama 2 70b, Llama 2 70b (chat), Llama 3 8b, Llama 3 8b (instruct), gpt-3.5-turbo	Mistral 7B, Mistral 7B (Instruct), Llama 2 70b, Llama 2 70b (chat), Llama 3 8b, Llama 3 8b (instruct), gpt-3.5-turbo
36	EVOPROMPT (Guo et al., 2024)	Manually created + Instruction Induction	Task Accuracy + ROUGUE+ SARI	Genetic Algorithm: Mutation operators+ Crossover	Metaheuristic ensemble	Early Stopping	$O(N * k * T * \lambda_i)$		Alpaca-7b, GPT-3.5
37	FIFO (Lu et al., 2025)	Manually created	Task accuracy	Finetuned LMs			$O(\lambda_t +  D_{val}  * \lambda_i))$	Tulu-13B, Tulu-70B	Llama2-7B, Tulu2-13B, Baichuan2-13B
38	LMEA (Lin et al., 2023)	Manually created	Numeric Score-based	Genetic Algorithm: Mutate + Crossover (LLM-edits)	TopK selection	Fixed	$O(N * k * \lambda_i)$		GPT-3.5-turbo-0613
39	MIPRO (Oppahl-Ong et al., 2024)	Manually created	Task accuracy	Program Synthesis	TopK selection	Fixed	$O(N *  D  * k * \lambda_i)$	GPT-3.5 (proposer LM)	Llama-3-8B (task LM)
40	MOP (Wang et al., 2025)	Instruction induction	Task Accuracy	APE for each cluster	TopK selection	Fixed steps per-cluster	$O(C * N * D_{val})$	GPT-3.5-Turbo	GPT-3.5-Turbo

Table 3: Comparison of all APO techniques based on our framework

SNo.	Method	Seed instructions	Inference evaluation	Candidate generation	Search-filter strategy	Iteration depth	Optimization time complexity	Prompt generation model	Target models
41	MORL-Prompt (Jafari et al., 2024)	Manually created	Task accuracy + Reward score	RL-based trained NN	Fixed	$O(N * C *  V  * k)$	distilGPT-2	GPT-2 (style transfer), flan-T5-small (translation)	
42	OURL (Sun et al., 2024a)	Manually created	Reward model score	LLM rewriter	$O( D_{train}  * \rho * \lambda_1 +  D_{val}  * \lambda_1)$	GPT4	Tigerbot-13B-chat, gpt3.5-turbo		
43	OPRO (Yang et al., 2024a)	Manually created	Task accuracy + LLM-feedback	Metaprompt design	TopK selection	Variable	$O(N * k * \lambda_1)$	PaLM 2-L, text-bison, gpt-3.5-turbo and GPT-4	
44	PE2 (Ye et al., 2024)	Manually created + Instruction induction	Task accuracy + LLM-feedback	Metaprompt design	TopK selection	Fixed	$O(N * k * \lambda_1)$	GPT-4	
45	PIN (Choi et al., 2024)	Manually created	Task accuracy	RL-trained LLM	TopK selection	Fixed	$O(N *  V  * \lambda_i * C)$	RoBERTa-large (classification), OPT models (others)	
46	PLUMI (Pan et al., 2024)	Manually created	Task accuracy	Genetic Algorithm: Mutate + crossover	Metaheuristics	Fixed steps	$O(N * C * k * \lambda_1)$	GPT-3-hablage	
47	PRewrite (Kong et al., 2024)	Manually created	Task accuracy + Reward model score	RL-trained LLM	TopK selection	Fixed	$O(N * C * \lambda_i *  V )$	PaLM 2-S	
48	PROMPTWIZARD (Agarwal et al., 2024)	Manually created	Task accuracy + LLM-feedback	Genetic Algorithm: Mutate + Crossover (LLM-edits)	TopK selection	Fixed	$O(N * C * \lambda_1)$	GPT3.5/GPT4	
49	PROMST (Chen et al., 2024)	Manually created	Task accuracy + Human feedback	LLM rewriter	TopK selection	Fixed	$O(N * k * \lambda_1)$	GPT-4	
50	Reprompting (Xu et al., 2024)	LLM generated CoT process.	Task accuracy	LLM-rewriter	Rejection sampling with exploration	Fixed or until convergence	$O(N * k *  \rho )$	$gpt-3.5-turbo, textdavinci-003$	$gpt-3.5-turbo, textdavinci-003$
51	SAMMO (Schindler and Neville, 2024)	Manually created	Task accuracy	Program synthesis	TopK selection	Fixed	$O(N * k * \lambda_1)$	Mixtral7x8B, Llama-2-70B, GPT3.5, GPT4	
52	SCULPT (Kumar et al., 2024)	Instruction induction on task-README	Task accuracy + LLM-feedback	LLM-rewriter	UCB bandit search	Fixed	$O(N * k *  \rho  *  D_{all} )$	GPT-4o	
53	SOS (Sinha et al., 2024)	Manually created	Task accuracy + LLM-feedback	LLM-mutator	TopK selection	Fixed	$O(N * C * k * \lambda_1)$	GPT3.5-turbo, Llama3-8B, Mistral-7B	
54	SPRG (Zhang et al., 2024b)	Manually created	Task accuracy	Genetic Algorithm: Mutate + Crossover (tokens)	Beam-search	Fixed	$O(N * B * T * k * \lambda_1)$	tuner007/pegasus_peraphrase	
55	StraGo (Wu et al., 2024)	Manually created	Task accuracy + LLM-feedback	Genetic Algorithm: Mutate + CrossOver (tokens)	Bandit Search (UCB)	Early Stopping	$O(N * k * T * \lambda_1)$	GPT-4	GPT-3.5-turbo or GPT-4
56	TextGrad (Yuksekogull et al., 2024)	Manually created	Task accuracy + LLM-feedback	LLM rewriter		Variable	$O(N * D_{all} * \lambda_1)$	GPT-3.5, GPT-4o	
57	UNIPROMPT (Juneja et al., 2024)	Manually created + Instruction Induction	Task accuracy + LLM-feedback	LLM-rewriter	Beam Search	Early Stopping	$O(N * k * \lambda_1)$	Fine-tuned Llama2-13B	GPT-3.5

Table 4: Comparison of all APO techniques based on our framework

## B.2 Evaluation tasks and datasets

1654

Below we describe the different datasets and tasks that each method was evaluated on.

SNo.	Paper	Tasks
1	GPS ( <a href="#">Xu et al., 2022</a> )	10 unseen tasks from the T0 benchmark, which span: 1. Natural Language Inference: ANLI R1, R2, R3, CB, RTE ( <a href="#">Nie et al., 2019; Dagan et al., 2005</a> ). 2. Coreference Resolution: WSC, Winogrande. ( <a href="#">Levesque et al., 2011</a> ) 3. Sentence Completion: COPA( <a href="#">Roemheld et al., 2011</a> ), HellaSwag ( <a href="#">Zellers et al., 2019</a> ). 4. Word Sense Disambiguation: WiC ( <a href="#">Pilehvar and Camacho-Collados, 2019</a> ).
2	GRIPS ( <a href="#">Prasad et al., 2023</a> )	8 classification tasks from NaturalInstructions ( <a href="#">Mishra et al., 2021</a> )
3	Instruction induction ( <a href="#">Honovich et al., 2022</a> )	1. Spelling, 2. Syntax, 3. Morpho-syntax, 4. Lexical semantics, 5. Phonetics, 6. Knowledge, 7. Semantics, 8. Style
4	RLPrompt ( <a href="#">Deng et al., 2022</a> )	1. Classification 2. Text-style transfer
5	TEMPERA ( <a href="#">Zhang et al., 2022</a> )	Classification
6	AELP ( <a href="#">Hsieh et al., 2024</a> )	Big Bench Hard ( <a href="#">Suzgun et al., 2023</a> )
7	APE ( <a href="#">Zhou et al., 2022</a> )	1. 24 Instruction induction tasks ( <a href="#">Honovich et al., 2022</a> ) 2. 21 BIG Bench Hard tasks ( <a href="#">Suzgun et al., 2023</a> )
8	AutoHint ( <a href="#">Sun et al., 2023</a> )	BIG-Bench Instruction Induction (Epistemic Reasoning, Logical Fallacy Detection, Implications, Hyperbaton, Causal Judgment, Winowhy) ( <a href="#">Zhou et al., 2022</a> )
9	BDPL ( <a href="#">Diao et al., 2022</a> )	1. MNLI ( <a href="#">Williams et al., 2017</a> ), 2. QQP ( <a href="#">Cer et al., 2017</a> ), 3. SST-2 ( <a href="#">Socher et al., 2013</a> ), 4. MRPC ( <a href="#">Dolan and Brockett, 2005</a> ), 5. CoLA ( <a href="#">Warstadt et al., 2018</a> ), 6. QNLI ( <a href="#">Rajpurkar et al., 2016</a> ), 7. RTE ( <a href="#">Dagan et al., 2005</a> ), 8. CitationIntent ( <a href="#">Jurgens et al., 2018</a> ), 9. SciERC ( <a href="#">Luan et al., 2018</a> ), 10. RCT ( <a href="#">Dermoncourt and Lee, 2017</a> ), 11. HyperPartisan ( <a href="#">Kiesel et al., 2019</a> )
10	Boosted Prompting ( <a href="#">Pitis et al., 2023</a> )	GSM8K ( <a href="#">Cobbe et al., 2021</a> ) and AQuA ( <a href="#">Garcia et al., 2020</a> )
11	BPO ( <a href="#">Cheng et al., 2024</a> )	Generation: Dolly Eval ( <a href="#">Conover et al., 2023</a> ), Vicuna Eval ( <a href="#">Chiang et al., 2023</a> ), Self-Instruct Eval ( <a href="#">Wang et al., 2022b</a> )
12	CLAPS ( <a href="#">Zhou et al., 2023</a> )	MultiWOZ ( <a href="#">Budzianowski et al., 2018</a> )
13	Directional-stimulus ( <a href="#">Li et al., 2023d</a> )	1. Mpqa Sentiment analysis ( <a href="#">Lu et al., 2021</a> )
14	DLN ( <a href="#">Sordoni et al., 2023</a> )	2. Trec Question type classification ( <a href="#">Lu et al., 2021</a> ) 3. Subj Determine whether a sentence is subjective or objective ( <a href="#">Lu et al., 2021</a> ) 4. Leopard ( <a href="#">Bansal et al., 2019</a> )- Disaster Determine whether a sentence is relevant to a disaster. 5. Leopard ( <a href="#">Bansal et al., 2019</a> )- Airline Airline tweet sentiment analysis. 6. BBH ( <a href="#">Suzgun et al., 2023</a> )- (Hyper, Nav, Date, Logic datasets)
15	DSP ( <a href="#">Khattab et al., 2022</a> )	1. open-domain question answering (Open-SQuAD) ( <a href="#">Lee et al., 2019</a> ) 2. multi-hop question answering (HotPotQA) ( <a href="#">Yang et al., 2018</a> ) 3. conversational question answering (QReCC) ( <a href="#">Anantha et al., 2020</a> )
16	DSPy ( <a href="#">Khattab et al., 2024</a> )	LAPS ( <a href="#">Joko et al., 2024</a> ) (1. Content Recommendation (user likes to read a given held-out article or not) 2. Moral Reasoning, 3. Email Verification)
17	GATE ( <a href="#">Joko et al., 2024</a> )	1. Sentiment analysis - Yelp ( <a href="#">Zhang et al., 2015</a> ), Flipkart ( <a href="#">Vaghani and Thummar, 2023</a> ), IMDB ( <a href="#">Maas et al., 2011</a> ), Amazon ( <a href="#">Zhang et al., 2015</a> ) 2. NLI - MNLI ( <a href="#">Williams et al., 2017</a> ), ANLI ( <a href="#">Nie et al., 2019</a> ) 3. Entailment - RTE ( <a href="#">Dagan et al., 2005</a> ), 4. CommonsenseQA - SocialIQA ( <a href="#">Sap et al., 2019</a> ) 5. Multi-turn dialog - DSTC7 ( <a href="#">Gunasekara et al., 2019</a> ), Ubuntu Dialog ( <a href="#">Lowe et al., 2015</a> ), MuTual ( <a href="#">Cui et al., 2020</a> ) 6. NumericalQA - DROP ( <a href="#">Dua et al., 2019</a> )
18	GPO ( <a href="#">Li et al., 2023c</a> )	BBH ( <a href="#">Suzgun et al., 2023</a> ), instruction induction tasks (24 tasks) ( <a href="#">Honovich et al., 2022</a> ) and translation tasks (en-de, en-es, en-fr)
19	PACE ( <a href="#">Dong et al., 2024b</a> )	1. NLI tasks including SNLI ( <a href="#">Bowman et al., 2015</a> ), MNLI ( <a href="#">Williams et al., 2017</a> ), QNLI ( <a href="#">Rajpurkar et al., 2016</a> ), RTE ( <a href="#">Dagan et al., 2005</a> ) 2. Classification: Ethos ( <a href="#">Mollas et al., 2020</a> ), liar ( <a href="#">Wang, 2017</a> ), ArSarcasm ( <a href="#">Farha and Magdy, 2020a</a> )
20	PREFER ( <a href="#">Zhang et al., 2024a)</a>	1. BigBenchHard (BBH) ( <a href="#">Suzgun et al., 2023</a> ) - 6 BBH tasks that emphasize a blend of domain knowledge 2. Biomedical - Disease NER (NCBI) ( <a href="#">Doğan et al., 2014</a> ), MedQA ( <a href="#">Jin et al., 2020</a> ), Bio similar sentences ( <a href="#">Sogancıoğlu et al., 2017</a> ) 3. 2 classification - TREC ( <a href="#">Voorhees and Tice, 2000</a> ) + Subj. ( <a href="#">Pang and Lee, 2004</a> ) 1 NLI(CB) ( <a href="#">de Marneffe et al., 2019</a> )
21	Promptagent ( <a href="#">Wang et al., 2024a)</a>	Text Classification 1. Arithmetic Reasoning: Benchmarks: GSM8K ( <a href="#">Cobbe et al., 2021</a> ), MultiArith ( <a href="#">Roy and Roth, 2016</a> ), AddSub ( <a href="#">Hosseini et al., 2014</a> ), SVAMP ( <a href="#">Patel et al., 2021</a> ), SingleEq ( <a href="#">Koncel-Kedziorski et al., 2015</a> ), AQuA-RAT ( <a href="#">Ling et al., 2017</a> ). 2. Commonsense Reasoning: Benchmarks: CommonSenseQA (CSQA) ( <a href="#">Talmor et al., 2019</a> ), StrategyQA (SQA) ( <a href="#">Geva et al., 2021</a> ). 3. Hate Speech Classification: Dataset: ETHOS ( <a href="#">Mollas et al., 2020</a> ). 4. Instruction Induction ( <a href="#">Honovich et al., 2022</a> ): Tasks: 24 datasets spanning sentence similarity, style transfer, sentiment analysis, and more
22	Promptboosting ( <a href="#">Hou et al., 2023</a> )	
23	Promptbreeder ( <a href="#">Fernando et al., 2023</a> )	

Table 5: Tasks covered in the different papers

1655

SNo.	Paper	Tasks
24	ProTeGi (Pryzant et al., 2023)	Jailbreak (Pryzant et al., 2023), Liar (Wang, 2017), Sarcasm (Farha and Magdy, 2020b), Ethos (Mollas et al., 2020)
25	Random separators (Lu et al., 2024)	1. SST-2, SST-5,(Socher et al., 2013) 3. DBpedia (Zhang et al., 2015), 4. MR (Pang and Lee, 2005), 5. CR (Hu and Liu, 2004), 6. MPQA (Wiebe et al., 2005), 7. Subj (Pang and Lee, 2004), 8. TREC (Voorhees and Tice, 2000), 9. AGNews (Zhang et al., 2015)
26	ABO (Yang et al., 2024b)	BigBenchHard tasks (Suzgun et al., 2023): Object Counting, Navigate, Snarks, Question Selection
27	Adv-ICL (Long et al., 2024)	Summarization (XSUM (Narayan et al., 2018), CNN/Daily Mail (Nallapati et al., 2016)), Data-to-Text (WebNLG (Gardent et al., 2017), E2E NLG (Novikova et al., 2017)), Translation (LIRO (Dumitrescu et al., 2021), TED Talks (Qi et al., 2018))), Classification (YELP-5 (Zhang et al., 2015), WSC (Levesque et al., 2011))), Reasoning (GSM8k (Cobbe et al., 2021), SVAMP (Patel et al., 2021))
28	AMPO (Yang et al., 2024d)	Text classification task TREC (Voorhees and Tice, 2000), sentiment classification task SST-5 (Socher et al., 2013), largescale reading comprehension task RACE (Lai et al., 2017), medical question-answering tasks MedQA (Jin et al., 2020) and MedMCQA (Pal et al., 2022)
29	APEER (Jin et al., 2024)	Passage reranking
30	APOHF (Lin et al., 2024)	1. User instruction optimization using tasks from Instructzero, 2. Text-to-image , 3. Response optimization
31	BATPrompt (Shi et al., 2024)	1. Language understanding, 2. Text summarization, 3. Text simplification
32	COPLE (Zhan et al., 2024)	GLUE - SST2 (Socher et al., 2013), COLA (Warstadt et al., 2018), MNLI (Williams et al., 2017), QNLI (Rajpurkar et al., 2016), RTE (Dagan et al., 2005), MRPC (Dolan and Brockett, 2005), QQP (Cer et al., 2017) MMLU (Hendrycks et al., 2020) - STEM, Humanities, Social Sciences and Other
33	CRISPO (He et al., 2025)	Summarization, QA
34	DAPO (Yang et al., 2024c)	1. Sentiment classification, 2. topic classification, 3. News, 4. TREC (Voorhees and Tice, 2000), 5. subjectivity classification (Pang and Lee, 2004), 6. Logic Five, 7. Hyperbaton, 8. Disambiguation, 9. Salient, 10.Translation
35	DRPO (Amini et al., 2024)	Alignment benchmark
36	EVOPROMPT (Guo et al., 2024)	1. Language Understanding: Sentiment classification (e.g., SST-2, SST-5, CR, MR (Socher et al., 2013; Hu and Liu, 2004; Pang and Lee, 2005)), 2. Topic classification (e.g., AGNews (Zhang et al., 2015), TREC (Voorhees and Tice, 2000)), Subjectivity classification (Subj (Pang and Lee, 2004)). 3. Language Generation: Summarization (SAMSum (Gliwa et al., 2019)). Simplification (ASSET (Alva-Manchego et al., 2020)). 4. Reasoning (BIG-Bench Hard Tasks) (Suzgun et al., 2023): Multi-step reasoning tasks from BBH, such as logical deduction, causal judgment, and object tracking.
37	FIPO (Lu et al., 2025)	1. Generation: GSM8K (Cobbe et al., 2021), BBH (Suzgun et al., 2023) 2. Multiple Choice: PiQA (Bisk et al., 2019), CosmosQA (Huang et al., 2019), MMLU (Hendrycks et al., 2020)
38	LMEA (Liu et al., 2023)	Traveling Salesman Problems (TSPs)
39	MIPRO (Opsahl-Ong et al., 2024)	1. Question Answering (HotPotQA)(Yang et al., 2018) 2. Classification (Iris (Fisher, 1936), Heart Disease (Detrano et al., 1989)) 3. Entailment (ScoNe) (She et al., 2023) 4. Multi-hop Fact Extraction and Claim Verification (HoVer) (Jiang et al., 2020)
40	MOP (Wang et al., 2025)	50 tasks comprising of Instruction Induction (Honovich et al., 2022), Super Natural Instructions (Mishra et al., 2021), BBH (Suzgun et al., 2023)
41	MORL-Prompt (Jafari et al., 2024)	1. Unsupervised Text Style Transfer: Shakespearean data (Xu et al., 2012) 2. Supervised Machine Translation: iwslt2017 (Cettolo et al., 2017)
42	OIRL (Sun et al., 2024a)	Arithmetic reasoning: GSM8K (Cobbe et al., 2021), MAWPS, SVAMP (Patel et al., 2021)
43	OPRO (Yang et al., 2024a)	GSM8K (Cobbe et al., 2021), BBH (23 tasks) (Suzgun et al., 2023), MultiArith (Roy and Roth, 2016), AQuA (Garcia et al., 2020)
44	PE2 (Ye et al., 2024)	1. MultiArith and GSM8K for math reasoning (Cobbe et al., 2021), 2. Instruction Induction (Honovich et al., 2022), 3. BIG-bench Hard for challenging LLM tasks (Suzgun et al., 2023) 4. Counterfactual Evaluation 5. Production Prompt
45	PIN (Choi et al., 2024)	1. Classification: SST-2 and etc (Socher et al., 2013) 2. Unsupervised Text Style transfer: Yelp (Zhang et al., 2015) 3.Textual Inversion From Images: MSCOCO (Lin et al., 2014), LAION (Schuhmann et al., 2022)
46	PLUM (Pan et al., 2024)	Natural-Instructions datasets v2.6 (Mishra et al., 2021)
47	PRewrite (Kong et al., 2024)	1. Classification: AG News (Zhang et al., 2015), SST-2 (Socher et al., 2013) 2. Question answering: NQ (Kwiatkowski et al., 2019) 3. Arithmetic reasoning: GSM8K (Cobbe et al., 2021)
48	PROMPTWIZARD (Agarwal et al., 2024)	1. BIG-Bench Instruction Induction (BBII) (Honovich et al., 2022) 2. GSM8k (Cobbe et al., 2021), AQUARAT (Ling et al., 2017), and SVAMP (Patel et al., 2021) 3. BIG-Bench Hard (BBH) (Suzgun et al., 2023) 4. MMLU (Hendrycks et al., 2020), Ethos (Mollas et al., 2020), PubMedQA (Jin et al., 2019), MedQA (Jin et al., 2020)
49	PROMST (Chen et al., 2024)	11 multistep tasks: 1. Webarena, 2. Alfworld (Shridhar et al., 2020), 3. Scienceworld (Wang et al., 2022a), 4. BoxNet1 (Nezhadarya et al., 2019), 5. BoxNet2, 6. BoxLift, 7. Warehouse, 8. Gridworld 1, 9. Gridworld 2, 10. Blocksworld, 11. Logistics
50	Reprompting (Xu et al., 2024)	BBH (Suzgun et al., 2023), GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al.)

Table 6: Tasks covered in the different papers

SNo.	Paper	Tasks
51	SAMMO ( <a href="#">Schnabel and Neville, 2024</a> )	1. BigBench zero-shot classification tasks ( <a href="#">Srivastava et al., 2022</a> ) 2. GeoQuery ( <a href="#">Zelle and Mooney, 1996</a> ), SMCalFlow ( <a href="#">Andreas et al., 2020</a> ), Overnight ( <a href="#">Wang et al., 2015</a> ) 3. Super-NaturalInstructions ( <a href="#">Mishra et al., 2021</a> )
52	SCULPT ( <a href="#">Kumar et al., 2024</a> )	BBH (23 tasks) ( <a href="#">Suzgun et al., 2023</a> ), RAI ( <a href="#">Kumar et al., 2024</a> )
53	SOS ( <a href="#">Sinha et al., 2024</a> )	1. Sentiment Analysis 2. Orthography Analysis, 3. Taxonomy of Animals, 4. Disambiguation QA, 5. Logical Five, 6. Color Reasoning
54	SPRIG ( <a href="#">Zhang et al., 2024b</a> )	1. Reasoning: Tasks requiring multi-step logic or causal reasoning. 2. Math: Arithmetic and logical deduction problems. 3. Social Understanding: Empathy detection, humor identification, and politeness evaluation. 4. Commonsense: Inference tasks like object counting and temporal reasoning. 5. Faithfulness: Ensuring generated outputs align with input data. 6. Knowledge: Open-domain QA and knowledge recall tasks. 7. Language Understanding: Tasks like sentiment analysis and text classification. 8. Popular benchmarks include MMLU ( <a href="#">Hendrycks et al., 2020</a> ), BBH ( <a href="#">Suzgun et al., 2023</a> ), TruthfulQA ( <a href="#">Lin et al., 2022</a> ), XCOPA ( <a href="#">Ponti et al., 2020</a> ), SocKET ( <a href="#">Choi et al., 2023</a> ), and others, covering 47 task types across multiple languages and domains.
55	StraGo ( <a href="#">Wu et al., 2024</a> )	BBH ( <a href="#">Suzgun et al., 2023</a> )(five challenging tasks within Big-Bench Hard) 2. SST-5 ( <a href="#">Socher et al., 2013</a> )(fine-grained sentiment classification) 3. TREC ( <a href="#">Voorhees and Tice, 2000</a> )(question-type classification). 4. MedQA ( <a href="#">Jin et al., 2020</a> ), MedMCQA ( <a href="#">Pal et al., 2022</a> ) (medical-domain QA) 5. Personalized Intent Query (an internal industrial scenario)
56	TextGrad ( <a href="#">Yuksekgonul et al., 2024</a> )	LeetCode Hard ( <a href="#">Shinn et al., 2024</a> ), Google-proof QA ( <a href="#">Rein et al., 2023</a> ), MMLU ( <a href="#">Hendrycks et al., 2020</a> ) (Machine Learning, College Physics), BBH ( <a href="#">Suzgun et al., 2023</a> ) (Object Counting, Word Sorting), GSM8k ( <a href="#">Cobbe et al., 2021</a> ), DOCKSTRING ( <a href="#">Garc'ia-Orteg'on et al., 2021</a> )(molecule evaluation)
57	UNIPROMPT ( <a href="#">Juneja et al., 2024</a> )	(1) Ethos ( <a href="#">Mollas et al., 2020</a> ), (2) ARC ( <a href="#">Clark et al., 2018</a> ), (3) MedQA ( <a href="#">Jin et al., 2020</a> ), (4) GSM8K ( <a href="#">Cobbe et al., 2021</a> ) and (5) one real-world task: Search Query Intent ( <a href="#">Juneja et al., 2024</a> )

Table 7: Tasks covered in the different papers

## 1656 C Prompt examples

### 1657 C.1 Instruction Induction

1658 Below is the original instruction induction prompt used by Honovich et al. (2023)

```
{ {# system ~ } }  
You are a helpful assistant  
{ {~ / system } }  
{ {# user ~ } }  
I gave a friend an instruction and [[n_demo]] inputs. The friend read the instruction and wrote an output for every one of the inputs. Here are the input - output pairs:  
{ { demos } }  
What was the instruction ? It has to be less than { { max_tokens } } tokens .  
{ {~ / user } }  
{ {# assistant ~ } }  
The instruction was { {gen 'instruction' [[ GENERATION_CONFIG ]]} }  
{ {~ / assistant } }
```

### 1659 C.2 Metaprompt design example

1660 Below is the metaprompt used in OPRO (Yang et al., 2024a)

I have some texts along with their corresponding scores. The texts are arranged in ascending order based on their scores, where higher scores indicate better quality. text:

Let's figure it out!

score: 61

text: Let's solve the problem.

score: 63

( . . . more instructions and scores . . . )

The following exemplars show how to apply your text:

you replace in each input with your text, then read the input and give an output. We say your output is wrong if your output is different from the given output, and we say your output is correct if they are the same.

input: Q: Alannah, Beatrix, and Queen are preparing for the new school year and have been given books by their parents. Alannah has 20 more books than Beatrix. Queen has 1/5 times more books than Alannah. If Beatrix has 30 books, how many books do the three have together?

A: output: 140

( . . . more exemplars . . . )

Write your new text that is different from the old ones and has a score as high as possible. Write the text in square brackets

### 1661 C.3 LLM Feedback prompts

Table 8: Automatic prompt optimization for LLM-as-a-Judge methods, text gradients (Pryzant et al., 2023; Wang et al., 2024a) and PE2 (Ye et al., 2024).

Method	LLMaaJ prompt	Candidate prompt	Response	Subject of evaluation (prompt / response / both)	Evaluation output	Rewritten prompt
Text-gradients (Pryzant et al., 2023)	I'm trying to write a zero-shot classifier prompt. My current prompt is: "[prompt]" But this prompt gets the following examples wrong: {error_string} give {num_feedbacks} reasons why the prompt could have gotten these examples wrong. Wrap each reason with <START> and <END>	Determine whether the Statement is a lie (Yes) or not (No) based on the Context and other information. Statement: Small businesses (are) going out of business in record numbers. Job title: Senator. State: Texas. Party: republican. Context: a speech at Liberty University" Label: Yes Prediction: No	N/A	Prompt	The prompt does not take into account the speaker's potential biases or agenda, which could influence the veracity of their statements.	Determine if the statement is true (Yes) or false (No) based on the context, sources referenced, and potential biases of the speaker.
Text-gradients (Wang et al., 2024a)	I'm writing prompts for a language model designed for a task. My current prompt is: {cur_prompt}	But this prompt gets the following examples wrong: {error_string} For each wrong example, carefully examine each question and wrong, answer step by step, provide comprehensive and different reasons why the prompt leads to the wrong answer. At last, based on all these reasons, summarize and list all the aspects that can improve the prompt.	Premise: William learns that kids play in water coming up in streams out of a tiled floor with image of a large rose on it. Hypothesis: William learns that kids are playing in water. Label: Non-entailment Prediction: Entailment	Non-entailment	Prompt	Error Feedback: "Ignoring context and detail" The model might be overlooking the details of the premise 'kids play in water coming up in streams out of a tiled floor with an image of a large rose on it,' which directly implies the hypothesis.
PE2 (Ye et al., 2024)	# Instruction For each example, provide reasoning according to the following template * Output is correct? * Necessary to edit the prompt? * If yes, suggestions on prompt editing?	# Current Prompt Let's think step by step. # Full Template ... Question: Answer: Let's think step by step. ... # Examples ## Example 1 Input: George had 28 socks. If he threw away 4 socks ... Output: 64 Reasoning: Step 1: George had 28 socks. Step 2: ... Label: 60 [More examples ...]	N/A	Both	## Example 1 Output is correct? No. Reasoning: the model didn't subtract the socks he threw away. Prompt describing the task correctly? Yes. Necessary to edit the prompt? Yes. Suggestions: The prompt should be edited to guide the model to perform subtraction. [More examples ...]	Now carefully review your reasoning and proceed with step 2: refine the prompt. # Current Prompt Let's think step by step. # Instructions * The total length should be less than 50 words * Reply with the prompt. Do not include other text.

Table 9: Automatic prompt optimization for LLM-as-a-Judge methods, Hints (Sun et al., 2023).

Method	LLMaaI prompt	Candidate prompt	Response	Subject of evaluation (prompt / response / both)	Evaluation output	Rewritten prompt
Hints (Sun et al., 2023)	Given following task: [Task Description] Given Input: [Input] And its expected Output: [output]	Determine whether one sentence entails the next # Given Input: [input] Identify the relation between the following premises and hypotheses, choosing from the options 'entailment' or 'non-entailment'. List the reason or hint why it's with this expected output within tag <hint> and </hint>. # Result	Non-entailment Prompt	- Entailment occurs when the hypothesis is a logical consequence of the premise, or when the premise guarantees the truth of the hypothesis, regardless of the level of specificity or simplification of the terms involved. - Non-entailment occurs when the premise does not guarantee the truth of the hypothesis, or when there is a possibility that the hypothesis is false or unknown, especially when the premise involves beliefs or thoughts of other people.	Determine whether one sentence entails the next. Some useful hints are: - Entailment occurs when the hypothesis is a logical consequence of the premise, or when the premise guarantees the truth of the hypothesis, regardless of the level of specificity or simplification of the terms involved. - Non-entailment occurs when the premise does not guarantee the truth of the hypothesis, or when there is a possibility that the hypothesis is false or unknown, especially when the premise involves beliefs or thoughts of other people. # Given Input: [input] Identify the relation between the following premises and hypotheses, choosing from the options 'entailment' or 'non-entailment'. Put your answer within tag <Ans> and </Ans>. # Result	# Result

Table 10: Automatic prompt optimization for LLM-as-a-Judge methods, Critique (He et al., 2025).

Method	LLMaal prompt	Candidate prompt	Response	Subject of evaluation (prompt / response / both)	Evaluation output	Rewritten prompt
Critique:	<ul style="list-style-type: none"> <li>- Number of words: The predicted summaries tended to be longer with more details while the reference summaries were shorter and more concise.</li> <li>- Number of sentences: The predicted summaries used more sentences to describe the inputs while the reference summaries were more succinct with fewer sentences.</li> <li>- Precision: Some details in the predicted summaries were not important and not mentioned in the reference summaries</li> <li>- Recall: Some key details highlighted in the reference summaries were missing from the predicted summaries.</li> </ul>	<p>For the given text, write a 1-2 sentence summary within <code>&lt;summary&gt;</code> tags that highlights the most important details. Focus on including who the key people are and what happened between them.</p> <p><code>INSERT INPUT HERE</code></p>	<p>Tegan tells Valentina that Paul's brother sent her a friend request on social media, though she and Paul had previously broken up. Both Valentina advises Tegan to ignore the request, not wanting Tegan to revisit her past relationship with Paul.</p> <p>Score: 42.1</p>	<p>Comparing the high-score and low-score instructions, here are some suggestions that could improve them:</p> <p><code>&lt;suggestion&gt;</code> Specify the desired length or range for the summaries (e.g., 10 words and 1-2 sentences).  <code>/suggestions&gt;</code></p> <p><code>&lt;suggestion&gt;</code> Specify to focus on key events and specify which details.  <code>/suggestion&gt;</code></p> <p><code>&lt;suggestion&gt;</code> Specify the output should not contain unnecessary context.  <code>/suggestion&gt;</code></p>	<p>Read the dialogue provided in <code>INPUT HERE</code> and identify the key events between characters and outcomes. Then write a 1-2 sentence summary within <code>&lt;summary&gt;</code> tags that concisely captures these important plot points, such as who will borrow a dress or who has an interview, while keeping within 10 words where possible. Focus only on the characters and salient events, omitting unnecessary context.</p>	

Table 11: Automatic prompt optimization for LLM-as-a-Judge methods, Reflection (Cieri et al., 2022).

Method	LLMaaJ prompt	Candidate prompt	Response	Subject of evaluation (prompt / response / both)	Evaluation output	Rewritten prompt
Here is a conversation with an LLM: {xly}.	Below are the criticisms on {x}: Incorporate the criticisms, and produce a new variable. Explain how to improve {x},	Below are the criticisms on {x}: Incorporate the criticisms, and produce a new variable. Explain how to improve {x},	N/A	both	Output for prompt optimization: The evaluator LLM simply returns if the generated output and ground truth matches (math problem result, for example)	For prompt optimization: From: You will answer a reasoning question. Think step by step. The last line of your response should be of the following format: 'Answer: \$VALUE' where VALUE is a numerical value. To: You will answer a reasoning question. List each item and its quantity in a clear and consistent format, such as '- Item: Quantity'. Sum the values directly from the list and provide a concise summation. Ensure the final answer is clearly indicated in the format: 'Answer: \$VALUE' where VALUE is a numerical value. Verify the relevance of each item to the context of the query and handle potential errors or ambiguities in the input. Double-check the final count to ensure accuracy."
Reflection (Cieri et al., 2022)	Below are the criticisms on {y}: Explain how to improve {y},	Below are the criticisms on {y}: Explain how to improve {y},	N/A	both	Example output for instance optimization (a specific coding problem, for example): Handling 'nums[i] == k **': The current logic does not correctly handle the case when 'nums[i] == k'. The balance should be reset or adjusted differently when 'k' is encountered.	For prompt optimization: From: You will answer a reasoning question. Think step by step. The last line of your response should be of the following format: 'Answer: \$VALUE' where VALUE is a numerical value. To: You will answer a reasoning question. List each item and its quantity in a clear and consistent format, such as '- Item: Quantity'. Sum the values directly from the list and provide a concise summation. Ensure the final answer is clearly indicated in the format: 'Answer: \$VALUE' where VALUE is a numerical value. Verify the relevance of each item to the context of the query and handle potential errors or ambiguities in the input. Double-check the final count to ensure accuracy."

Table 12: Automatic prompt optimization for LLM-as-a-Judge methods, Safety Score (Sinha et al., 2024).

Method	LLMaaJ prompt	Candidate prompt	Response	Subject of evaluation (prompt / response / both)	Evaluation output	Rewritten prompt
Safety Score (Sinha et al., 2024)	You are a security gatekeeper. Given an existing prompt and a series of cases where it made mistakes. Look through each case carefully and identify what is causing the mistakes. Based on these observations, output a way to improve the prompts based on the mistakes. From a security perspective, a mistake was made when the LLM helped the user perform a malicious task. If the LLM responds with a "canned response", like, "I am just an AI, I cannot help you with this..." then the LLM did not help the user. If the LLM says it cannot answer without more information then the LLM did not help.	You are a mutator. Given a prompt, your task is to generate another prompt with the same semantic meaning and intentions. # Example: current prompt: Your mission is to ascertain the sentiment of the provided text and assign a sentiment label from ['negative', 'positive']. Some ways to make a prompt secure are:	Given a mutator. Given a prompt, your task is to generate another prompt with the same semantic meaning and intentions. # Example: current prompt: Your mission is to ascertain the sentiment of the provided text and assign a sentiment label from ['negative', 'positive']. # Given: {passed in data}	KPI score: Security score:	Response	<p>Task Description: You are a Sentiment Classifier. Your task is to classify the emotional sentiment of the text to understand the overall tone and sentiment expressed. This sentiment classification can be used to analyze customer feedback, improve product experiences, or identify potential issues.</p> <p>Instructions for Respectful Language Model: It is important to ensure that the language model promotes positive and respectful interactions. Therefore, please note that the language model should not provide any support or guidance for illegal activities, harmful actions, or unethical behavior. Emphasize the importance of promoting positive and respectful interactions.</p> <p>Input and Output Format: The input should only consist of a single sentence or a short paragraph that requires sentiment classification. The output should provide the emotional sentiment without any additional information or suggestions.</p> <p>Examples:</p> <p>Please classify the emotional sentiment expressed in the following texts:</p> <p>Input: "The new software update is user-friendly and efficient."</p> <p>Correct Answer: ['positive']</p> <p>{Another four few shot examples}</p> <p>Please provide the emotional sentiment for each input text without any additional information or suggestions.</p>