000 001 002

003 004

010 011

012

013

014

015

016

017

018

019

021

024

025

PLHF: PROMPT LEARNING FROM FEW-SHOT HUMAN Feedback

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advances explore prompt tuning for large language models (LLMs) and develop automatic optimization frameworks to obtain suitable prompts with respect to desired output quality metrics. Although existing approaches can handle conventional tasks such as fixed-solution question answering, defining the metric becomes complicated when the output quality cannot be easily assessed by comparisons with standard golden samples, especially for those natural language applications that multiple outputs are equally valid. Consequently, optimizing the prompts effectively and efficiently without a clear metric becomes a critical challenge. To address this issue, we present PLHF, a few-shot prompt optimization framework inspired by the well-known RLHF technique. Different from näive strategies involving human experts, PLHF employs a specific evaluator module acting as the metric to estimate the output quality. PLHF requires only a single round of human feedback to complete the entire prompt optimization process. Empirical results on both public and industrial datasets show that PLHF significantly outperforms existing output scoring strategies for LLM prompt optimizations.

- 1 INTRODUCTION
- 027 028

I INTRODUCTION

029 General-purpose large language models (LLMs) have demonstrated substantial capabilities across various fields in recent years. However, solving complex tasks with LLMs often requires appropriate customizations on LLMs to fit the task requirements. While fine-tuning pre-trained LLMs is a 031 common approach, it may be infeasible when there is limited training data, restricted computational resource, or when working with a black-box LLM. Alternatively, previous studies (Wang et al., 2022; 033 Shin et al., 2020) have shown that the potential of LLMs can also be fully leveraged with suitable 034 prompts. Recent literature develops automatic few-shot prompt optimization for LLM usages, such as DSPy (Khattab et al., 2024) and TextGrad (Yuksekgonul et al., 2024). To determine an effective prompt for the LLM, existing methods often employ gradient descent or other algorithms (Yang 037 et al., 2024; Zhou et al., 2023; Guo et al., 2023) to optimize the performance with respect to desired 038 metrics (i.e., definition of output quality). Overall, the key to success heavily relies on the output quality evaluations which shall precisely reveal the model performance to the optimizer. 039

040 Although such output quality metrics are often well-defined for the tasks which can be modeled as 041 the traditional discriminative tasks (e.g., classifications and regressions) where the performance can 042 be directly evaluated given ground-truths, scoring the outputs often becomes non-trivial for most 043 generation-type of tasks. An example comes from the essay writing task where the LLM needs 044 to output an essay given specific requirements. To perform prompt optimizations, the evaluator needs to score an essay written by an LLM. This could be extremely challenging without human's involvement. Another typical scenario is the dialogue systems (e.g., chat-bots), as automatically 046 rating the outputs is difficult. In this case, the output quality could be affected by numerous factors 047 such as the context, the circumstances of the environment, as well as specific user preferences. 048 Therefore, a generalized evaluation metric is difficult to formulate. 049

The absence of a precise metric would hinder the effectiveness of the prompt optimization process
 for a generative task. Most prompt optimization systems adopt two types of mechanisms to score the
 generated outputs. The first is employing simple evaluators, such as exact matching or soft matching
 (based on certain similarity measurement) to compare the generated outputs with the observed samples. The second common strategy is to utilize existing LLMs for the output scoring. For instance,



064

065

066

067

068

071

073

075

077

079

081 082

084 085

Figure 1: Demonstrations of the actual failure cases that the evaluations from pre-trained LLMs have different preference from specific humans. The first (left) example is the task of joke generation, where the scoring is according to funniness and novelty. The second scenario is a math problem generation bot, where the response quality is evaluated based on helpfulness and problem quality. As shown above, the verdicts of state-of-the-art LLMs could still differ from real human's preferences.



Figure 2: Workflow framework of PLHF. The entire LLM program contains two modules, Responder R and Evaluator E, where PO can be PO arbitrary prompt optimization method.

several studies (Wang et al., 2023; Fu et al., 2024) leverage advanced OpenAI GPT models (Achiam 087 et al., 2023) as an external judge to score the generated outputs. Nevertheless, such scorers suffer 088 from a critical drawback as generic pre-trained LLMs might not have enough contextual or background knowledge to behave as accurate scoring functions. (See Figure 1 for the examples of failure 089 cases.) Ultimately, the purpose of prompt optimization is to tackle complex tasks that the original 090 LLM, when provided with a simple prompt, fails to handle effectively. Consequently, employing 091 such a sub-optimal LLM as the scoring mechanism is likely to diminish the effectiveness of the 092 prompt optimization framework. As a result, for applications related to response generation, it is highly demanded to involve real human experts to evaluate the results generated by LLMs, yet such 094 scheme suffers from a frequently occurring issue — constrained budget to employ human experts. 095

For most of the automatic prompt optimization frameworks (Khattab et al., 2024; Yuksekgonul et al., 096 2024; Pryzant et al., 2023; Deng et al., 2022; Wen et al., 2024), multiple iterations of optimizations are performed, with the quality estimated based on the given metric. With querying human experts 098 acting as the metric, once the prompt is updated with any modifications, we have to ask human experts for their judgement again for each of the training inputs, which might cause serious efficiency 100 bottleneck. To address the aforementioned issues, we present PLHF (which stands for Prompt 101 Learning with Human Feedback), a few-shot prompt optimization framework. Inspired by the fa-102 mous Reinforcement Learning from Human Feedback (RLHF) technique (Ouyang et al., 2022; Li 103 et al., 2023), PLHF introduces a particular *evaluator* module E which requires human scoring no 104 greater than *linear* (with respect to the number of training samples) times during the optimization 105 process. To leverage human feedback in few shots, we consider utilizing a prompt-optimized LLM as E to evaluate the output of the main responder R. The overall framework is depicted in Fig-106 ure 2. Specifically speaking, first, we employ human experts to provide judgements as scores on a 107 set of training samples D, containing pairs of the inputs and sample text-outputs. Then, we perform

prompt optimization on another LLM to mimic the human experts' preference pattern based on D. Since the prompting task on the evaluator module is relatively typical (e.g., binary classifications or regressions), we can leverage any existing automatic few-shot prompt optimization frameworks (e.g., DSPy (Khattab et al., 2024)) with trivial metrics (e.g., Accuracy or Mean Absolute Error) to obtain the evaluator module. Finally, we can perform prompt optimizations for both E and R to establish the entire framework.

114 To verify the actual effectiveness of PLHF, experiments are conducted not only on multiple bench-115 mark datasets but also on a real-world industrial data collected from an online customer support 116 chat-bot product of an AI company. The experimental results shown in Section 4 demonstrate that 117 PLHF can boost the output quality of existing automatic few-shot prompt optimization frameworks 118 with our duo-module design. Moreover, to verify the effectiveness of each module, we also provide an analysis towards the relations in performance curves between the number of training sam-119 ples and the output scores for each subtask. The source codes for the experiments are available at 120 https://(keeping secret for the double-blind paper review). 121

- 122 In summary, our contributions are as follows:
 - We study automatic few-shot prompt optimization for LLMs with limited number of human feedback calls, which is a more reasonable and feasible setting in real-world applications, especially for those have specific user preferences or multiple acceptable outputs.
 - We introduce PLHF, a novel prompt optimization framework that does not directly rely on well-defined metrics for the text output. Instead, we design an evaluator module to provide an automatic mechanism that evaluates text outputs for LLM program prompting.
 - We conduct extensive experiments on publicly accessible benchmark datasets as well as a test on industrial data to validate the effectiveness of PLHF. The results show that PLHF has superiority in terms of output quality, compared with the approaches employing string matching or adopting the state-of-the-art LLM (GPT-40) as the evaluator.
- 133 134 135

124

125

126 127

128

129 130

131

132

2 RELATED WORK

136 137

Recent research has investigated various strategies to obtain suitable prompts for LLMs. Earlier studies have introduced techniques of automating the search process for data samples (Gao et al., 2021), learning prompts through gradient-based searching methods (Shin et al., 2020; Wen et al., 2024; Pryzant et al., 2023), refining prompts using evolutionary algorithms (Guo et al., 2023; Fernando et al., 2023) and utilizing other LLMs for prompt generation (Yang et al., 2024; Zhou et al., 2023). Several studies have also attempted to optimize prompts using reinforcement learning, exploring prompt editing at different granular levels such as word-level (Deng et al., 2022), phrase-level (Zhang et al., 2023), and within text-to-image generation tasks (Hao et al., 2024).

As LLMs are increasingly applied in real-world scenarios, in-context learning (McCann et al., 2018; Radford et al., 2018; Brown et al., 2020) is becoming an emerging trend for effective LLM programming. Instruction tuning (Ouyang et al., 2022) further enhances this process by enabling complex behaviors through the use of structured prompts (Press et al., 2023; Yao et al., 2023; Khot et al., 2023; Madaan et al., 2024).

For automatic few-shot prompt optimization, Khattab et al. (2024) introduced DSPy, a state-of the-art prompt optimization framework, which considers LLM usages in a programmatic fashion.
 DSPy parameterizes each module to learn the data pattern and the desired behaviors by iteratively
 bootstrapping useful demonstrations. On the other hand, Yuksekgonul et al. (2024) inspired by LLM
 fine-tuning procedures and proposed TextGrad, a framework refining the prompt with the back propagation algorithm. Instead of deriving numeral-valued gradients, TextGrad regards LLMs' text
 feedback as the 'gradient' in texts.

All these techniques rely on well-defined metrics to set their objectives. While advanced generalpurpose LLMs like GPT-4 can be adopted for text-output evaluations (Zimbres, 2024; Zheng et al., 2023; Tan et al., 2024), they may lack the contextual or background knowledge needed for accurate
evaluation in specific tasks. Hence, involving human feedback becomes inevitable for the prompt
optimizations in such tasks. To address the issue, Lin et al. (2024) inspired by dueling bandits and designed a strategy to choose pairs of prompts to query for human feedback during the prompt

| Input: "Why did the shoe store owner become a detective?" Output: "Because he was good at solving sneaker cases!" Labeled Score: 9 | | |
|--|--|--|
| Input: "Why don't astronauts throw parties in space?" Output: "Because there's no atmosphere for a good vibe!" Labeled Score: 8 Input: "Why did the smartphone go to therapy?" Output: "It had too many hang-ups!" Labeled Score: 8 : Positive Samples | | |
| | | |
| Generate the Output for the given Input, where Input is the "setup" and Output is the "punchline" of a joke. | | |
| Follow the following format. Input: \${input} Output: \${output} Initial Prompt <i>P</i> _R (Instruction | | |
| Input: "Why did the shoe store owner become a detective?" Output: "Because he was always good at solving sneaker case Input: "Why don't astronauts throw parties in space?" | | |
| Output: "Because there's no atmosphere for a good vibe!" Input: "Why did the smartphone go to therapy?" | | |
| | | |

187 188 Figure 3: A toy example to illustrate the subtask designs of PLHF. The targeted generative AI task 189 for this example is "generate the punchline for a joke setup." The training samples D are triplets 190 (Input, Output, Labeled Score), where Input is the joke setup, Output is a sample output of the 191 punchline for the corresponding *Input*, and *Labeled Score* is the rating judged by human experts. 192 For this example, we consider *Labeled Score* \geq 8 as the condition of positive samples. Examples of optimized prompts P'_E and P'_R (for the evaluator E and the responder R, respectively) are shown. 193

optimizations to reduce the number of needed calls of human feedback. In this paper, we consider a different approach to tackle the issue — we focus on the *metric* in the prompt learning process. We developed a duo-module framework to obtain an evaluator module acting as the metric of the main task to perform the desired LLM prompt optimizations, requiring minimal human feedback.

199 200 201

202

203

204

205

206

207

194 195

196

197

PLHF: PROMPT LEARNING FROM FEW-SHOT HUMAN FEEDBACK 3

As shown in Figure 2, our entire framework, PLHF, is designed to perform prompt optimizations for typical language model program usages — output a proper response based on the given input. The whole process is guided by the principal intuition of taking advantages from few-shot in-context learning Brown et al. (2020) to capture the contextual patterns from limited number of labeled samples. Since there is no explicit metric available, a scoring function is needed for existing prompt optimization frameworks. Hence, we introduce an *evaluator* module E, acting as the scoring function for the main *responder* module R. The two modules correspond to two respective subtasks.

208 209

211

210 3.1 Responder Task

212 The responder task for R is tasked with performing the original assignment. The core of this com-213 ponent is based on a base LLM, denoted as M_B , which generates outputs based on the input query. LLM M_R is starting from pairing with an initial prompt P_R describing the roles of input and out-214 put, as well as the relationship between input and output (i.e., how the input determines the output). 215 See Figure 3 for an example. The training samples for R include input-output pairs with *positive*

| Alg | orithm 1 PLHF: Duo-module Framework for Few-shot LLM Prompt Optimizations |
|-----|--|
| Inp | ut: Training data samples $D = \{d_1, d_2, \dots, d_n\}$ (each d_i is a triplet containing input/query q_i |
| | output o_i and score/verdict r_i), Base LLMs M_R , M_E , Initial prompts P_R , P_E , Trivial metric L |
| | for evaluator E. |
| Ou | tput: Optimized prompts P'_R and P'_E . |
| 1: | Set $P'_R := P_R$ and $P'_E := P_E$. |
| 2: | while there are new training samples in D do |
| 3: | while P'_E is not yet optimized (converged) for D do |
| 4: | for each training sample $d_i = (q_i, o_i, r_i) \in D$ do |
| 5: | Input (q_i, o_i) pair to generate the score (or verdict) \tilde{r}_i by M_E with prompt P'_E . |
| 6: | end for |
| 7: | Compute conventional metric score $S_E = L(\{\tilde{r}_1, \ldots, \tilde{r}_n\}, \{r_1, \ldots, r_n\}).$ |
| 8: | Consider the D, P'_E and S_E to update the prompt P'_E . |
| 9: | end while |
| 10: | while P'_R is not yet optimized (converged) for D do |
| 11: | for each training sample $d_i = (q_i, o_i, r_i) \in D$ with positive rating r_i do |
| 12: | Input q_i to generate the output \tilde{o}_i by M_R with prompt P'_R . |
| 13: | end for |
| 14: | Evaluate the outputs by evaluator E. Obtain the score $S_R = M_E(\{\hat{o}_1, \dots, \hat{o}_n\}; P'_E)$. |
| 15: | Consider the D, P'_R and S_R to update the prompt P'_R . |
| 16: | end while |
| 17: | If (<i>optional</i>) having new inputs/queries $\{t_1, \ldots, t_m\}$ to augment D then |
| 18: | Test the current responder $R = (M_R; P_R)$ with inputs/queries $\{t_1, \ldots, t_m\}$. |
| 19: | Collect numan feedback f_i for the output $M_R(t_i; P'_R)$, for $i = 1, \ldots, m$. |
| 20: | Append $(t_i, M_R(t_i; P_R), f_i)$, for $i = 1,, m$, as new training samples into D . |
| 21: | ena II |
| 22: | ena while |
| 23: | return P_R, P_E |

243 244

245 score/verdict labeled by human experts. The data positivity can be specifically defined to align with 246 the requirements of the assigned AI task. For instance, the example described in Figure 3 considers 247 a score threshold as the condition of the positive samples. The generated responses by the LLM M_R with prompt P_R are then judged by the evaluator module E. To enhance the quality and adaptability 248 of the responses, with respect to human experts' preference patterns, we perform prompt optimiza-249 tions on LLM M_R to obtain a prompt-optimized prompt P'_R for M_R . Finally, we consider M_R with 250 prompt P'_R as the finalized responder module R in our framework to produce desired outputs that 251 solve the original assigned AI task. 252

253 254 255

3.2 EVALUATOR TASK

256 The evaluator task for E is an auxiliary task designed to verify and score the output generated by 257 the responder R. Similar to the responder R, the evaluator E is also built with a base LLM, denoted 258 as M_E . The evaluator task considers training samples in the triplet form of input, output, and the 259 corresponding score (e.g., verdict, rating) labeled by human experts. Note that, different from the 260 training samples for R, this time we consider *all* training data samples (regardless of the positivity) as the references for E. To provide a nuanced verdict, we also optimize the prompt P_E for M_E 261 with a trivial metric L (such as the conventional Accuracy or Mean Absolute Error) to evaluate the 262 quality of the response. The metric L is specifically defined as the loss to estimate the difference 263 between the predicted scores and the actual labeled scores in D. Overall speaking, The evaluator 264 module E leverages LLM M_E with the finalized prompt P'_E to provide a score rating toward any 265 responses for the original AI task. 266

With the evaluator E, PLHF ensures that the outputs from the responder R are not only technically accurate but also contextually appropriate for the task. The feedback loop between the evaluator Eand the responder R helps refine the overall model performance, as the scores determined by E are used to inform future responses generated by R.

270 3.3 INTEGRATION AND FEEDBACK LOOPS271

With the integration of responder R and evaluator E, the entire system operates in a feedback-loop structure, as described in pseudo codes stated in Algorithm 1. At the beginning, we initialize prompts as $P'_R := P_R$ and $P'_E := P_E$ for modules R and E, respectively. In each iteration of PLHF, first, training data samples D are used to optimize the evaluator E (i.e., to update P'_E). Then, we optimize the responder R (i.e., update P'_R) regarding the evaluator E with prompt P'_E as the metric. After an iteration of prompt optimizations for P'_R and P'_E , we obtain a version of optimized prompts for the responder R and the evaluator E. Figure 3 provides an example of a toy generative AI problem to demonstrate the optimized prompts P'_R and P'_E .

For batch tests, the whole optimization process is ended by Line 16 in Algorithm 1 and the finalized prompt P_R are used for the LLM M_R to generate the outputs for upcoming test inputs/queries. On the other hand, if the scenario has incrementing data samples, starting from Line 17 in Algorithm 1, we can augment the training sample set D with the user feedback on new samples, then repeat the optimization processes for E and R.

Overall, the proposed PLHF framework is capable of performing prompt optimizations for LLMs even when occurring challenges of (a) no available well-defined metrics to evaluate the LLM output quality for the specific task, (b) limited number (few-shot) of labeled samples for LLM prompting, and (c) multiple valid outputs for a single input.

289 290

291 292

293

294

295

296

297 298

299

4 EXPERIMENTS

To evaluate the performance and robustness of our proposed model framework, PLHF, we conducted a series of experiments across various tasks. The tasks were selected to test the model ability to generate accurate and contextually relevant outputs, while also assessing the effectiveness of the auxiliary evaluator task in refining responses. In the following subsections, we detail the dataset selection, experimental setup, and the results obtained from these experiments.

4.1 DATASETS

We conduct the experiments on three public datasets with various tasks, and one industrial dataset from a real-world product of question answering chat-bot generating practical SQL commands.

302 303 304

4.1.1 SCHEMA GUIDED DIALOGUE (SGD) DATASET

305 The Schema Guided Dialogue (SGD) dataset Sun et al. (2021) is a large-scale dataset designed for 306 task-oriented dialogue systems. It comprises dialogues collected in English, specifically designed 307 to encompass a wide range of dialogue scenarios, schema-based actions, and services. The dataset 308 contains 1,000 dialogues, contributing to a total of 13,833 utterances. Each user utterance in the dataset is labeled with a satisfaction score on a 5-point Likert scale (Likert, 1932). Rating 1 indicates 309 the lowest level of satisfaction, while rating 5 denotes the highest satisfaction. The distribution of 310 satisfaction ratings $\{1, 2, 3, 4, 5\}$ is $\{120, 769, 11151, 1494, 50\}$. These human-assigned satisfaction 311 scores are valuable for assessing chat-bot responses with respect to user satisfaction. 312

313

4.1.2 AUTOMATED ESSAY SCORING (AES) DATASETS

315 The dataset is originally provided by Ben et al. (2012) for the Automated Student Assessment Prize 316 (ASAP). The dataset, named as AES-ASAP, consists of eight essay sets varying in length, ranging 317 from an average of 150 to 550 words per response. The responses were written by students in grades 318 seven through ten, and all essays were hand-graded by human experts. Each essay was double-319 scored, with a resolved score provided to harmonize the differences between raters. We use the 320 training set (with the scores in domain 1) of the first set of essays in our experiments. The actual 321 text of the student's response is included. We consider the *average score* as aggregated result from these raters. The scores are distributed from 1 to 30. In our experiments, we discard the essays with 322 transcription errors (marked as "illegible" or containing placeholder text such as "???") from the 323 training data.

324 Apart from AES-ASAP, our experiments also include a newer essay scoring dataset of from an 325 online competition hosted by Kaggle (2024). The dataset, named as AES-2.0, contains 24,000 326 student-written argumentative essays. Each essay was scored on a scale of 1 to 6 as the holistic 327 rating ¹ judged by human experts. Similar to our settings for AES-ASAP, we consider the provided 328 training split for prompt optimizations in our the experiments.

4.1.3 INDUSTRIAL SQL COMMAND QUESTION ANSWERING DATASET

331 In addition to the previously mentioned public datasets, we have also deployed PLHF on a real-332 world question-answering system, which is currently an actual product of a commercial AI company. 333 This test, named as SQL-QA, comprises 100 real-world queries involving various database inquiry 334 requests from the clients. The database entries contains daily transaction records and other logs 335 sourced from multiple banks in China. As the training data for prompt optimizations, human experts 336 from the company labeled 10 positive samples and 20 negative samples for the prompt optimizations. 337 Table 1 provides examples of queries used in this test.

338 339 340

341

330

| Table 1: Examples of | the queries in our industrial SQL-QA test. |
|----------------------------------|--|
| User Ouery (English translation) | SOL Statement (Output) |

| 342 343 344 345 346 347 | Please list the top 3 clients by total deposits at the Beijing branch as of January 31, 2024. | SELECT CUST_ID, CUST_NAME, DEPO_BAL FROM acct WHERE DATA_DT='20240131' AND ORG_NAME='Beijing' ORDER BY DEPO_BAL DESC LIMIT 3 |
|--|---|--|
| 348 349 350 351 352 353 | Please inquire about the top 5 banks with the highest asset balances, grouped by institution, as of March 31, 2024. | SELECT ORG_NAME, ASSET_BAL FROM acct WHERE DATA_DT='20240331' ORDER BY ASSET_BAL DESC LIMIT 5 |

353 354 355

356

364

365

366

367

368

369

370

372

4.2 EXPERIMENT SETUPS

357 To perform prompt optimizations (denoted as PO) in each subtask, we consider two state-of-the-art 358 automatic prompting frameworks, DSPy and TextGrad. Same experiments are conducted for both 359 frameworks, and respective results are shown. 360

For the experiments, first, we estimate the effectiveness of the evaluator E, which solves the task 361 of predicting the labeled scores based on each input-output pair. The comparisons include several 362 baseline methods: 363

- Base LLM (GPT-3.5): the grounding baseline simply employing raw gpt-3.5-turbo-0125 model to predict the labeled score each input-output pair. We utilize OpenAI APIs for the model implementation. The prompt is set to be the same as the initial prompt P_E in PLHF for the comparisons. See Figure 3 for an example.
- MLP with Text Embedding: DNN-based predictor leveraging a Multi-layer Perception (MLP) model, based on the algorithm presented by Popescu et al. (2009), to predict the score of each input-output pair. Since the input and the output are texts, we transform the texts into embeddings by the powerful text-embedding-ada-002 model (OpenAI, 2024) to obtain the respective numerical vectors. The number of layers is set to 3.
- 373 • SVM with Text Embedding: conventional ML predictor — similar to the MLP one, but 374 this time considering Support Vector Machines (SVMs) as the score predictor. We adopt 375 the SVM implementation provided by Chang & Lin (2011) with the default configuration. 376

¹https://storage.googleapis.com/kaggle-forum-message-attachments/2733927/20538/Rubric_ 377 HolisticEssayScoring.pdf

Table 2: Summary of experimental results for the evaluator subtask across each dataset. For the public datasets, the presented values are RMSE losses (lower is better) of the output scores from E; for the industrial dataset SQL-QA, the values indicate Accuracy scores (higher is better). The best ones are marked in bold font.

| 302 | | Method | SGD | AES-ASAP | AES-2.0 | SQL-QA | |
|-----|---|--|-----------|-------------------|-------------------|-----------------|----------------------------|
| 383 | | Base LLM (GPT-3.5) | 1.02 | 4.75 | 0.46 | 0.53 | |
| 384 | | MLP with Text Embedding | 1.17 | 7.22 | 1.08 | 0.33 | |
| 385 | | SVM with Text Embedding | 1.25 | 6.43 | 1.10 | 0.40 | |
| 386 | | Base LLM PO via DSPy | 0.43 | 2.36 | 0.33 | 0.80 | |
| 387 | | Base LLM PO via TextGrad | 0.40 | 2.42 | 0.38 | 0.73 | |
| 388 | | | | | | | |
| 389 | | | | | | | |
| 390 | • | GPT-3.5 PO via DSPy: LLM | with de | monstration-bas | sed PO — j | performing pr | ompt op- |
| 391 | | timizations with DSPy framew | ork (Kh | attab et al., 20 | 24), based | on the aforen | nentioned |
| 392 | | setting of Base LLM (GP1-3.5). | | | | | |
| 393 | • | GPT-3.5 PO via TextGrad: LI | LM with | text instruction | n-based PO | - performin | ig prompt |
| 394 | | optimizations with TextGrad fra | amewor | k (Yuksekgonul | l et al., 2024 | 4), based on | the afore- |
| 395 | | mentioned setting of Base LLM (GPT-3.5). | | | | | |
| 396 | | | | | | | |
| 397 | Then, fo | r the main responder task R , outp | out quali | ty of the LLM | with optimiz | zed prompts a | re judged |
| 398 | by test q | ueries. We consider various types | s of eval | uators for the p | rompt optin | nizations. | |
| 399 | | Deco IIM (CDT 2 5), the even | | haadina air | | ant 0 | E turbo |
| 400 | • | Dase LLM (GP I-3.3): the gro | ounding | baseline — sir | npiy utilizi | iyon innut Ti | .5-luido- |
| 401 | | is set to be the same as the initia | al promi | the output bas | for the even | eriments See | Figure 3 |
| 402 | | for an example | a prom | | for the exp | eriments. See | / I Iguie 5 |
| 403 | | | 6.4 | | 1 . | c · · · | |
| 404 | • | PO with GP I-40: using a state- | -of-the-a | art LLM as the o | evaluator - 0 | - performing I | PO on the |
| 405 | Base LLM (GP 1-3.5). Adopting GP 1-40 (gpt-40-2024-05-13) model as the judge to | | | | | | |
| 406 | | score the outputs during FO. we | employ | OpenAl AFIS | in the imple | | |
| 407 | • | PO with Exact Matching: scor | ing by h | hard-matching - | Base LL | M (GP I-3.5) |) with PO |
| 408 | | regarding the given reference ou | tputs as | the ground-trut | th (i.e., the g | golden predic | tion). Let |
| 409 | | score = 1 if the output of R is ex | cactly th | e same as the g | round-truth | ; otherwise, so | $\operatorname{core} = 0.$ |
| 410 | • | PO with Embedding Similarity | y: scorin | g by soft-match | ning — Base | e LLM (GPT- | •3.5) with |
| 411 | | PO considering the cosine simila | rity scor | e between the e | mbedding v | ectors of the c | output and |
| 412 | | the ground-truth. The outputs a | re embe | dded by the po | werful text | -embedding- | -ada-002 |
| 413 | | model (OpenAI, 2024) via Open | AI APIs | 3. | | | |
| 414 | • | PLHF: our proposed framework | c — the | base LLMs M | T_R and M_E | are both set | to be raw |
| 415 | | GPT-3.5 (gpt-3.5-turbo-0125) | models. | | | | |
| 416 | | | | | | | |
| 417 | 4.3 Ev | ALUATIONS | | | | | |
| 418 | | | | | | | |
| 419 | For the | model output, we employ multi | ple hum | an experts as t | he judges t | o provide pro | ofessional |
| 420 | scores w | ith respect to the score scales of | t the ori | ginal data. Ho | wever, for t | he public dat | asets, the |
| 421 | original | people who labeled the data are i | unavaila | ble to give their | r judgement | for our new | generated |
| 422 | outputs 1 | or the data inputs. In our experi | ments, v | OPT 40 with - | concept of p | seuao-numai | i judge to |
| 423 | the pseu | do-human judge. Since we cons | sider Gl | PT-3.5 for the | base LLMs | in all the me | ethods for |

424 425 426

378

427 4.4 EXPERIMENTAL RESULTS

Table 2 lists the experimental results of the evaluator task. As shown in the table, we can observe that the conventional methods (MLP/SVM with Text Embedding) seem struggled in predicting the labeled scores from the given embedded inputs. In contrast, the LLM-based methods performed significantly better on both public datasets and the industrial tests. A possible reason is that LLMs

more powerful model that can provide fair evaluations toward output quality.

experiments, the pseudo-human judge (i.e., GPT-40 with DSPy PO based on training samples) is a

| 433 | Table 3: Summary of experimental results for the responder subtask across each dataset. For the |
|-----|--|
| 434 | industrial dataset SQL-QA, the overall Accuracy score are given by actual human experts in the |
| 435 | company; for the public datasets, the scores from the pseudo-human judge are shown. The values |
| 436 | for Base LLM are the actual scores, whereas for the other methods, relative improvements are shown |
| /27 | in percentages. The best ones are marked in bold font. |
| 437 | |

| PO | Method | SGD | AES-ASAP | AES-2.0 | SQL-QA |
|----------|------------------------------|---------|----------|---------|---------|
| | Base LLM (GPT-3.5) | 4.25 | 26.50 | 5.35 | 0.74 |
| | PO with GPT-40 | +1.18% | +3.70% | +0.75% | +5.41% |
| DSPy | PO with Exact Matching | - 4.00% | -10.87% | - 5.61% | 0.00% |
| | PO with Embedding Similarity | +1.65% | - 4.27% | +0.56% | +10.81% |
| | PLHF | +6.59% | +8.45% | +2.62% | +18.92% |
| | PO with GPT-40 | +4.71% | +3.28% | +1.31% | +2.70% |
| | PO with Exact Matching | -10.59% | -15.92% | -10.84% | -18.92% |
| TextGrad | PO with Embedding Similarity | +3.53% | - 0.64% | +0.93% | +2.70% |
| | PLHF | +8.71% | +8.68% | +4.30% | +18.92% |

might have superior fitting and understanding capabilities to handle text inputs. With prompt optimizations, both DSPy and TextGrad provided more effective prompt for more accurate evaluators.

451 As for the experimental results of the responder task, shown in Table 3, we consider both DSPy 452 and TextGrad as the prompt optimization (PO) tool for each method in the comparisons. Overall, the results are relatively similar in same directions for each pair of the scores between the two 453 PO selections. In summary, for all the four datasets, the proposed PLHF framework achieved the 454 best performance in output quality, in terms of the metric for each task. Moreover, PLHF used 455 GPT-3.5 as the evaluator's base LLM M_E to achieve superior performance than 'PO with GPT-456 40', which conducted prompting with a more powerful GPT-40 as the evaluator. For the other 457 baselines, 'PO with GPT-40' consistently outperformed 'Base LLM (raw GPT-3.5)'. Regarding 458 the conventional matching-based scoring functions, both hard-matching (Exact Matching) and soft-459 matching (Embedding Similarity) produced outputs with worse quality. 460

4.5 PERFORMANCE ANALYSIS

463 In addition to the overall performance, we also analyze the robustness and the relationship between 464 model effectiveness and the number of training samples involved in the prompt optimization process. 465 As examples, Figure 4 demonstrates the performance curves for datasets SGD and AES-ASAP. 466 Note that, to analyze the performance of a responder with n data samples, we also use the evaluator 467 optimized with the same n samples. For the curves of evaluator E, we can observe that the RMSE 468 loss raised for initial samples, then the RMSE value dropped significantly after few shots of data. 469 For the curves toward responder R, the pattern is similar to the curves of E (but in opposite way), the output quality score dropped for initial samples, while the score then bouncing back and achieving 470 new highs with greater number of samples. Last but not least, as expected, the standard deviations 471 lower along with the increasing number of training samples. 472

473 474

475

432

447 448

449

450

461

462

5 CONCLUSION

476 In this paper, we focused on prompt optimizations for LLMs with a limited amount of human feed-477 back — a more practical and achievable approach for real-world applications. To address the chal-478 lenges of no well-defined metrics and the scarce human resources, we introduced PLHF, a few-shot 479 prompt learning with an evaluator module design to automatically score the outputs generated by 480 LLMs. We performed extensive experiments with public datasets and a real industrial dataset to 481 verify the effectiveness of PLHF. The experimental results demonstrated that PLHF outperforms 482 existing methods across from simple string matching functions to even the latest publicly available 483 LLMs as output evaluators in terms of the output quality. Overall, our proposed framework is practically effective especially for the scenarios when directly applying pre-trained general-purpose LLMs 484 are not the best option. Our future work involves enhancing and deploying the proposed framework 485 across diverse applications, particularly for tasks that utilize multi-modal data.



Figure 4: Performance curves of PLHF on the datasets SGD and AES-ASAP. For the plots, we consider DSPy as the PO method for PLHF. The x-values are the number of (randomly selected) training samples. The y-values are mean values of the RMSE losses for E and the output scores for R, respectively. The vertical bar of each point indicates the standard deviations estimated in 30 runs.

REFERENCES

510

511

512

513 514 515

516

523

524

525 526

527

528 529

531

532

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-517 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical 518 report. arXiv preprint arXiv:2303.08774, 2023. 519
- 520 Hamner Ben, Morgan Jaison, lynnvandev, Shermis Mark, and Ark Tom Vander. The hewlett foun-521 dation: Automated essay scoring, 2012. URL https://kaggle.com/competitions/asap-aes. 522
 - Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
 - Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. ACM transactions on intelligent systems and technology (TIST), 2(3):1–27, 2011.
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, 530 Eric Xing, and Zhiting Hu. Rlprompt: Optimizing discrete text prompts with reinforcement learning. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pp. 3369-3391, 2022.
- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. 534 Promptbreeder: Self-referential self-improvement via prompt evolution. arXiv preprint 535 arXiv:2309.16797, 2023. 536
- Jinlan Fu, See Kiong Ng, Zhengbao Jiang, and Pengfei Liu. Gptscore: Evaluate as you desire. In Proceedings of the 2024 Conference of the North American Chapter of the Association for 538 Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 6556– 6576, 2024.

| 540 541 542 543 544 | Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pp. 3816–3830, 2021. |
|---------------------------------|--|
| 545 546 547 | Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. <i>arXiv preprint arXiv:2309.08532</i> , 2023. |
| 548 549 | Yaru Hao, Zewen Chi, Li Dong, and Furu Wei. Optimizing prompts for text-to-image generation. <i>Advances in Neural Information Processing Systems</i> , 36, 2024. |
| 5551 552 | Kaggle. Automated essay scoring 2.0, 2024. URL https://www.kaggle.com/competitions/ learning-agency-lab-automated-essay-scoring-2/. |
| 553 554 555 556 557 | Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, Heather Miller, et al. Dspy: Compiling declarative language model calls into state-of-the-art pipelines. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. |
| 558 559 560 | Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. |
| 561 562 563 | Zihao Li, Zhuoran Yang, and Mengdi Wang. Reinforcement learning with human feedback: Learn- ing dynamic choices via pessimism. <i>arXiv preprint arXiv:2305.18438</i> , 2023. |
| 564 | Rensis Likert. A technique for the measurement of attitudes. Archives of psychology, 1932. |
| 565 566 567 | Xiaoqiang Lin, Zhongxiang Dai, Arun Verma, See-Kiong Ng, Patrick Jaillet, and Bryan Kian Hsiang Low. Prompt optimization with human feedback. <i>arXiv preprint arXiv:2405.17346</i> , 2024. |
| 568 569 570 | Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. <i>Advances in Neural Information Processing Systems</i> , 36, 2024. |
| 571 572 573 | Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. The natural language decathlon: Multitask learning as question answering. <i>arXiv preprint arXiv:1806.08730</i> , 2018. |
| 574 575 | OpenAI. Openai platform: Embeddings — embedding models, 2024. URL https://platform.openai. com/docs/guides/embeddings/embedding-models. |
| 577 578 579 580 | Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35: 27730–27744, 2022. |
| 581 582 583 | Marius-Constantin Popescu, Valentina E Balas, Liliana Perescu-Popescu, and Nikos Mastorakis. Multilayer perceptron and neural networks. <i>WSEAS Transactions on Circuits and Systems</i> , 8(7): 579–588, 2009. |
| 585 586 587 | Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pp. 5687–5711, 2023. |
| 588 589 590 | Reid Pryzant, Dan Iter, Jerry Li, Yin Lee, Chenguang Zhu, and Michael Zeng. Automatic prompt optimization with "gradient descent" and beam search. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 7957–7968, 2023. |
| 592 593 | Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving lan- guage understanding by generative pre-training. 2018. URL https://openai.com/index/ language-unsupervised/. |

594 Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: 595 Eliciting knowledge from language models with automatically generated prompts. In Proceedings 596 of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 597 4222-4235, 2020. 598 Weiwei Sun, Shuo Zhang, Krisztian Balog, Zhaochun Ren, Pengjie Ren, Zhumin Chen, and Maarten de Rijke. Simulating user satisfaction for the evaluation of task-oriented dialogue systems. In 600 Proceedings of the 44rd International ACM SIGIR Conference on Research and Development in 601 Information Retrieval, SIGIR '21. ACM, 2021. 602 603 Ting Fang Tan, Kabilan Elangovan, Liyuan Jin, Yao Jie, Li Yong, Joshua Lim, Stanley Poh, Wei Yan Ng, Daniel Lim, Yuhe Ke, et al. Fine-tuning large language model (llm) artificial intelligence chat-604 bots in ophthalmology and llm-based evaluation using gpt-4. arXiv preprint arXiv:2402.10083, 605 2024. 606 607 Chaozheng Wang, Yuanhang Yang, Cuiyun Gao, Yun Peng, Hongyu Zhang, and Michael R Lyu. 608 No more fine-tuning? an experimental evaluation of prompt tuning in code intelligence. In Pro-609 ceedings of the 30th ACM joint European software engineering conference and symposium on the 610 foundations of software engineering, pp. 382–394, 2022. 611 Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, 612 Jianfeng Qu, and Jie Zhou. Is chatgpt a good nlg evaluator? a preliminary study. arXiv preprint 613 arXiv:2303.04048, 2023. 614 615 Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 616 Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery. 617 Advances in Neural Information Processing Systems, 36, 2024. 618 Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun 619 Chen. Large language models as optimizers. In The Twelfth International Conference on Learning 620 Representations, 2024. URL https://openreview.net/forum?id=Bb4VGOWELI. 621 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan 622 Cao. React: Synergizing reasoning and acting in language models. In The Eleventh International 623 Conference on Learning Representations, 2023. 624 625 Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and 626 James Zou. Textgrad: Automatic "differentiation" via text. 2024. 627 Tianjun Zhang, Xuezhi Wang, Denny Zhou, Dale Schuurmans, and Joseph E. Gonzalez. TEM-628 PERA: Test-time prompt editing via reinforcement learning. In The Eleventh Interna-629 tional Conference on Learning Representations, 2023. URL https://openreview.net/forum?id= 630 gSHyqBijPFO. 631 632 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, 633 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and 634 chatbot arena. Advances in Neural Information Processing Systems, 36:46595–46623, 2023. 635 Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and 636 Jimmy Ba. Large language models are human-level prompt engineers, 2023. URL https://arxiv. 637 org/abs/2211.01910. 638 639 Rubens Zimbres. Evaluating llms with langchain: Using gpt-4 to evaluate google's open model gemma-2b-it, 2024. URL https://medium.com/google-cloud/ 640 evaluating-llms-with-langchain-using-gpt-4-to-evaluate-googles-open-model-gemma-2b-it-eb7555e3bdeb. 641 642 643 644 645 646