#### 000 AUTOMATIC TASK-AWARE INSTRUCTION OPTIMIZER 001 FOR BLACK-BOX LLMS 002 003

Anonymous authors

Paper under double-blind review

#### ABSTRACT

010 Large Language Models (LLMs) have demonstrated superior capabilities in terms of solving various real-world tasks. However, their performance and generated 012 content quality heavily depend on task-relevant instructions, which makes instruction optimization a challenging but critical direction to explore. In particular, as practitioners generally cannot access black-box (or API) LLMs' internal parameters and gradient information, it consequently makes instruction optimization for 015 black-box LLMs especially non-trivial. Existing methods for optimizing black-box 016 LLM instructions mainly focus on in-context learning using manually designed or heuristic disciplines, which can be insufficient due to the extreme complexity of 018 modern black-box LLMs that can contain hundreds of billions of parameters. To 019 address these challenges, we propose a novel automatic instruction optimization framework named Automatic Instruction Optimizer (AIO). AIO is designed to perceive target task information and adaptively adjust its task-aware instructing strategy for a task-solver black-box LLM. By leveraging a white-box LLM with parameter fine-tuning for enhanced representation power, AIO can automatically 024 update its instructing strategy based on the feedback from task-solver black-box 025 LLM. To achieve this goal, AIO adopts a novel LLM parameter fine-tuning process powered by zeroth-order gradient approximation and Contextual Bandit techniques, which can effectively and efficiently help address the challenge of inaccessible black-box LLM internal parameters and gradients, as well as help alleviate ex-028 pensive API cost concerns by flexibly reusing collected black-box LLM feedback. Extensive empirical evaluations are presented to demonstrate properties of our proposed AIO, and its effectiveness in comparison with strong baselines.

031 032 033

034

004

006

008

009

011

013

014

017

021

026

027

029

#### INTRODUCTION 1

Large Language Models (LLMs) have demonstrated impressive performance across various application scenarios, such as knowledge graph reasoning (Pan et al., 2024). However, LLMs generally rely on well-crafted instructions that provide accurate guidance and sufficient reference for high-quality 037 answer generation. Designing such instructions can be particularly challenging for more powerful black-box (or API) LLMs (e.g., GPT-4 (Achiam et al., 2023), Claude-3 (Anthropic, 2024)), as their parameters and gradients are commonly inaccessible. Meanwhile, effective instructing strategies can 040 vary significantly across different LLMs, distinct LLM versions, and downstream tasks (Zhou et al., 041 2022; Khattab et al., 2023; 2022), while optimal instructing strategies generally require flexible adap-042 tations tailored to target tasks or domains (Sun et al., 2024; Liu et al., 2024). In this case, designing 043 an optimal instructing strategy for a specific target task like knowledge reasoning can be non-trivial 044 and expensive. In addition, crafting proper instructing strategies for domain-specific tasks can be difficult and time-consuming for human experts (Brown et al., 2020; Reynolds & McDonell, 2021; Shin et al., 2020). For instance, assigning human labor to refine such task-specific instructions will be 046 expensive, and the cost can grow exponentially along with increasingly more task categories of higher 047 granularity levels (Scao & Rush, 2021; Shin et al., 2023; Amatriain, 2024). Thus, it is necessary and 048 valuable to develop automatic task-aware instructing mechanisms based on task information to enable optimal performance of task-solver black-box LLMs, without intervention of human experts. 050

051 Regarding instruction optimization in terms of black-box LLMs, there is an emerging line of works using an additional instruction-generating LLM as an "instruction engineer" (Zhou et al., 2022; 052 Pryzant et al., 2023; Fernando et al., 2023; Guo et al., 2024; Chen et al., 2024; Lin et al., 2024), in order to leverage the strong representation power of LLMs in search of a good instructing strategy. 054 However, existing works mainly focus on in-context learning aspects with frozen LLM parameters, 055 based on manually crafted or heuristic disciplines, which can limit the ability of LLMs in terms 056 of perceiving and utilizing target task information and black-box LLM feedback. On the other 057 hand, LLM fine-tuning alternatively offers more elasticity by involving trainable LLM parameters to 058 maintain an over-parameterization advantage (Allen-Zhu et al., 2019) with regard to exemplars from target tasks, which can allow LLMs to adapt their interpretations to align task textual data with target objectives (Han et al., 2024b) such as human preferences (Korbak et al., 2023; Rafailov et al., 2024), 060 along with parameter-efficient fine-tuning options (Hu et al., 2021a; Wu et al., 2024; Han et al., 2024b). 061 Meanwhile, as modern LLMs are commonly pre-trained to achieve good generalization abilities 062 instead of being optimized for specific downstream tasks (i.e., task-aware instruction optimization 063 in our case) out of the box, LLM parameter fine-tuning can generally offer more flexibility and 064 better performance than in-context learning techniques for various downstream applications (Liu 065 et al., 2022; Mosbach et al., 2023). However, as black-box LLM parameters and gradients are 066 generally inaccessible, it is impractical to directly fine-tune the "instruction engineer" LLM with 067 back-propagation based on chain rule and black-box LLM feedback. 068

In face of above motivations and challenges, we propose a novel framework named Automatic 069 Instruction Optimizer (AIO) to automatically optimize instructions for task-solver black-box LLM, with regard to the target task. In particular, AIO composes and optimizes human-comprehensible 071 instructions fed into a black-box LLM to improve its performance on the target task, for reinforced 072 transparency and trustworthiness. To learn a good task-aware instructing strategy, distinct from 073 existing in-context learning approaches, AIO alternatively fine-tunes a white-box LLM (e.g., Llama 074 3 (Dubey et al., 2024)) into a capable automatic instruction optimizer, which is able to perceive 075 downstream task information and generate high-quality instructions accordingly. This formulation aims to tackle the formidable complexity of modern black-box LLMs (i.e., with possibly hundreds of 076 billions of parameters involved). Here, with strong representation power and learning abilities of a 077 fine-tuned white-box LLM, AIO is capable of learning the complex relationship between task-aware instructions and black-box LLM feedback. Intuitively, one significant obstacle is that black-box 079 LLM parameters and gradients are commonly inaccessible. To address this challenge and achieve efficient gradient approximation for the black-box LLM, we propose a novel zeroth-order (ZO) 081 gradient approximation approach aided by Thompson Sampling (TS), by modeling the ZO gradient approximation of the black-box LLM as a sequential decision-making process. In the meantime, 083 during instruction optimization, it is necessary to retrieve black-box LLM feedback, which requires 084 querying third-party APIs and incurs direct development costs. To alleviate API cost concerns, our 085 TS-based ZO gradient approximation adaptively reuses collected black-box LLM feedback, enabling efficient instruction optimization through the rich representation power originated from white-box LLM fine-tuning. Our contributions can be summarized as: 087

Problem Formulation: We delve into the realm of automatic task-aware instruction optimization for black-box LLMs, where existing in-context learning methods can fail to deliver optimal performance due to insufficient representation power of their learning models or mechanisms. Different from existing approaches, we alternatively transform the goal of instruction optimization into an LLM fine-tuning objective, where a white-box LLM with sufficient representation power is fine-tuned to generate high-quality task-aware instructions for a task-solver black-box LLM.

- 094 • Proposed Framework: Different from existing approaches where human experts are involved to manually design instructions for downstream tasks, our proposed AIO does not require such 096 intervention of human experts, while finishing all the instruction optimization automatically through LLM fine-tuning. To enhance the trustworthiness of our instruction optimization process, AIO is able to optimize human-comprehensible instructions (i.e., instructions made up by concrete textual 098 tokens) to provide additional insights and transparency for practitioners. To tackle challenges of inaccessible black-box LLM gradients and possibly expensive API costs of the black-box LLM, 100 AIO utilizes a novel zeroth-order gradient approximation approach aided by Thompson Sampling. 101 By inventively formulating ZO gradient approximation procedure as a sequential decision-making 102 process, this design enables us to approximate black-box LLM gradients effectively and efficiently, 103 which are essential for fine-tuning our instruction-generating white-box LLM. 104
- Experiments: Extensive experiments are conducted on real-world data sets, demonstrating the superior performance of AIO compared with state-of-the-art baselines, as well as efficiency advantages of AIO in terms of reducing API token costs. Furthermore, we perform additional analytical experiments to explore characteristics and properties of AIO, such as instruction optimization trajectory results that demonstrate how instructions evolve across the optimization process.

### 108 2 RELATED WORKS

Instruction Optimization for API (Black-box) LLMs. In contrast to white-box LLM instruction 110 optimization (Shin et al., 2020; Li & Liang, 2021; Lester et al., 2021), as practitioners generally 111 have no access to black-box LLM parameters or gradients, a line of existing works (Zhou et al., 112 2022; Prasad et al., 2022) perform instruction search based on manually defined criteria. Chen 113 et al. (2024); Lin et al. (2024); Hu et al. (2024) also apply another LLM with frozen parameters to 114 generate instructions for the black-box LLM, and gradually optimize the generated instruction based 115 on Bayesian Optimization (Frazier, 2018; Wang et al., 2023; Shahriari et al., 2015), Contextual Bandit approaches (Chu et al., 2011; Li et al., 2010; Valko et al., 2013; Zhou et al., 2020; Agrawal & Goyal, 116 2013; Zhang et al., 2021), or localized instruction optimization guided by Gaussian Process (Schulz 117 et al., 2018). Since these works primarily rely on manually designed or heuristic principles focused 118 on in-context learning, they can result in sub-optimal black-box LLM performance. Alternatively, we 119 fine-tune a white-box LLM into an "instruction engineer", capable of adaptively perceiving target 120 task information and directly learning from black-box LLM feedback for instruction optimization. 121

122 LLM-based Instruction Generation. Instruction generation using LLMs is an emerging research topic (Zhou et al., 2022; Ma et al., 2024; Schnabel & Neville, 2024), where LLMs are applied as 123 instruction optimizer and their instructing strategies are gradually refined based on environment 124 feedback or target model outputs. In particular, there are a series of works leveraging meta-prompts, 125 which can be manually designed by humans (Yang et al., 2024), or optimized by LLMs (Tang et al., 126 2024). Meanwhile, Pryzant et al. (2023) perform in-context "Gradient Descent" on instructions based 127 on interactions with an instruction-generating LLM. Fernando et al. (2023); Guo et al. (2024) propose 128 evolutionary algorithms to refine LLM-generated instructions in an in-context learning manner. 129 Different from these works, our fine-tuned LLM can automatically perceive task-relevant information 130 and black-box LLM feedback, which can generally offer more flexibility than in-context learning 131 approaches and require no labor of human experts during the instruction optimization process.

## <sup>132</sup> 3 PROBLEM FORMULATION

153 154

155 156

157

158 159

160

As mentioned in our Introduction, given a target task  $\mathcal{T}$ , two LLMs are involved in our pipeline: (1) 134 black-box LLM  $\mathcal{F}_B(\cdot)$  is applied for task-solving, i.e., generating answers for task queries. Here, 135 the black-box LLM is considered as part of our learning objective, as we aim to learn optimized 136 instructions to enable the black-box LLM to achieve optimal performance. (2) white-box LLM 137  $\mathcal{F}_W(\cdot; \Theta_W)$  with trainable parameters  $\Theta_W$  aims to generate and optimize human-comprehensible 138 instructions, based on task  $\mathcal{T}$  and feedback from  $\mathcal{F}_B(\cdot)$ . Suppose that target task  $\mathcal{T}$  is associated with 139 three data collections (i.e., query-label pairs (X, Y)) individually drawn from task  $\mathcal{T}$ : (1) training 140 data  $\mathcal{D}_{\text{Train}}$ , which can also be named as task exemplars, will be fed into the white-box LLM as 141 reference for generating high-quality task-specific instructions; (2) validation data  $\mathcal{D}_{Valid}$  is applied 142 for performance evaluation during the optimization; and (3) final evaluation will be conducted on 143 a separate testing data set, which will remain unrevealed and inaccessible during the optimization stage for instruction-generating white-box LLM. Meanwhile, we denote  $\mathcal{F}_W(\mathcal{D}_{\text{Train}}; \Theta_W)$  as the 144 instruction generated based on exemplars  $\mathcal{D}_{\text{Train}}$  and corresponding white-box LLM parameters  $\Theta_W$ . 145 With generated instruction, denote  $\hat{Y} = \mathcal{F}_B([\mathcal{F}_W(\mathcal{D}_{\text{Train}}; \Theta_W); X])$  as the answer generated by 146 black-box LLM for query X, where  $[\cdot; \cdot]$  operation embeds query X to the generated instruction. 147

**Learning Objective.** With task exemplars  $\mathcal{D}_{\text{Train}}$  and an evaluation function (e.g., loss function)  $\mathcal{L}(\cdot, \cdot)$ , we transform the instruction optimization process into a white-box LLM fine-tuning objective, to leverage the sufficient learning and representation power of LLM fine-tuning for instruction optimization. Here, we aim to find the optimal white-box parameters  $\Theta_W$  that minimize:

$$\min_{\boldsymbol{\Theta}_{W}} \left[ \mathbb{E}_{(X,Y)\sim\mathcal{T}} \left[ \mathcal{L}(\hat{Y},Y) \right] \right] = \min_{\boldsymbol{\Theta}_{W}} \left[ \mathbb{E}_{(X,Y)\sim\mathcal{T}} \left[ \mathcal{L} \left( \mathcal{F}_{B} \left( [\phi(\boldsymbol{\Theta}_{W});X],Y \right) \right) \right] \right]$$
(1)

in observation of task exemplars  $\mathcal{D}_{Train}$ , where we apply a shorthand for generated instruction:

$$\phi(\mathbf{\Theta}_W) := \mathcal{F}_W(\mathcal{D}_{\text{Train}}; \mathbf{\Theta}_W). \tag{2}$$

Intuitively, we utilize the above fine-tuning process to guide how white-box LLM comprehends task exemplars and composes task-specific instructions, based on task-solver black-box LLM feedback.

### 4 PROPOSED FRAMEWORK: AUTOMATIC INSTRUCTION OPTIMIZER (AIO)

Recall that we aim to fine-tune the white-box LLM parameters  $\Theta_W$  to minimize our learning objective by optimizing generated instructions. An intuitive approach is to update the white-box LLM



Figure 1: Pipeline of AIO. The white-box LLM  $\mathfrak{P}$  generates an instruction  $\mathfrak{P}$  from exemplars  $\mathcal{D}_{\text{Train}}$ , which is evaluated to produce validation loss. The gradient flow towards white-box LLM parameters is then decomposed into: (1) TS-aided ZO gradient approximation for black-box LLM gradients, and (2) back-propagation for white-box LLM gradients. Finally, white-box LLM parameters are updated via Gradient Descent.

parameters  $\Theta_W$  using gradients derived from the instruction evaluation results (Eq. 1). However, the nested black-box LLM makes direct back-propagation towards  $\Theta_W$  via the chain rule infeasible.

175 **Brief summary of AIO pipeline.** To address this challenge, we propose the AIO framework for 176 efficient and effective white-box LLM fine-tuning, aimed at instruction optimization. As illustrated 177 in Figure 1, the pipeline of AIO involves two major parts: (1) Instruction Generation & Evaluation: Given the exemplars  $\mathcal{D}_{\text{Train}}$ , the white-box LLM with parameters  $\Theta_W$  generates an instruction 178  $\phi(\Theta_W) := \mathcal{F}_W(\mathcal{D}_{\text{Train}}; \Theta_W)$  as in Eq. 2. This instruction is then evaluated based on the black-box 179 LLM output, which produces validation loss. (2) White-box Instruction-generating LLM Fine-180 *tuning*: Based on the validation loss, we can decompose the gradient flow towards  $\Theta_W$  into two 181 multiplicative components via the chain rule: (i) white-box LLM gradients that can be obtained 182 with back-propagation, and (ii) inaccessible black-box LLM gradients that are approximated with 183 our proposed TS-aided ZO approximation method. Different from existing methods with randomly sampled ZO directions (e.g., Spall (1992); Malladi et al. (2023)), we apply TS here to adaptively 185 determine the ZO directions for efficient and effective gradient approximation (Subsec. 4.1.2). We 186 elaborate on technical details in Subsec. 4.1 and provide AIO pseudo-code in Algorithm 1.

Validation loss. Given the original optimization objective in Eq. 1, since a comprehensive overview of the task distribution  $\mathcal{T}$  can be inaccessible, we alternatively evaluate the quality of generated instructions using validation data  $\mathcal{D}_{Valid}$ . This leads to our validation loss:

$$\mathcal{L}_{\text{Valid}}\big(\phi(\boldsymbol{\Theta}_W)\big) := \frac{1}{|\mathcal{D}_{\text{Valid}}|} \sum_{(X,Y)\in\mathcal{D}_{\text{Valid}}} \mathcal{L}\Big(\mathcal{F}_B\big([\phi(\boldsymbol{\Theta}_W);X]\big),Y\Big),\tag{3}$$

where the validation loss is evaluated on the instruction  $\phi(\Theta_W)$ , which is generated by the white-box LLM  $\mathcal{F}_W(\cdot; \Theta_W)$  as in Eq. 2. Consequently, the gradient flow towards white-box LLM parameters  $\Theta_W$  will become  $\partial \left[ \frac{1}{|\mathcal{D}_{Valid}|} \sum_{(X,Y) \in \mathcal{D}_{Valid}} \mathcal{L}(\mathcal{F}_B([\phi(\Theta_W); X]), Y) \right] / \partial \Theta_W$ . Given the nested blackbox LLM, we are unable to directly back-propagate towards  $\Theta_W$  through the chain rule.

#### 4.1 THOMPSON SAMPLING (TS) AIDED ZEROTH-ORDER (ZO) GRADIENT APPROXIMATION

To deal with the challenge of inaccessible black-box LLM gradients, by applying chain rule on Eq. 3, we first can decompose the gradient flow with respect to white-box LLM parameters  $\Theta_W$  into two separate multiplicative parts: (1) gradients involving the black-box LLM; (2) and white-box LLM gradients that can be obtained by back-propagation, as

$$\partial \left[ \frac{1}{|\mathcal{D}_{\text{Valid}}|} \sum_{(X,Y)\in\mathcal{D}_{\text{Valid}}} \mathcal{L} \left( \mathcal{F}_B \left( [\phi(\Theta_W); X] \right), Y \right) \right] / \partial \Theta_W \\ = \underbrace{\frac{\partial \left[ \frac{1}{|\mathcal{D}_{\text{Valid}}|} \sum_{(X,Y)\in\mathcal{D}_{\text{Valid}}} \mathcal{L} \left( \mathcal{F}_B \left( [\phi(\Theta_W); X] \right), Y \right) \right]}{\partial \phi(\Theta_W)} \cdot \underbrace{\frac{\partial \phi(\Theta_W)}{\partial \Theta_W}}_{\text{Black-box LLM Gradients}} \cdot \underbrace{\frac{\partial \phi(\Theta_W)}{\partial \Theta_W}}_{\text{White-box LLM Gradients}} \cdot \underbrace{\frac{\partial \phi(\Theta_W)}{\partial \Theta_W$$

208

207

171

172

187

191 192

199

200

201

202

210 On one hand, while our white-box LLM gradients can be obtained by conventional back-propagation 211 in a straightforward way, we can also readily integrate back-propagation with Parameter-efficient 212 Fine-tuning (PEFT) techniques, such as LoRA (Hu et al., 2021a), to enhance the efficiency of 213 white-box LLM fine-tuning while still maintaining promising performance. One the other hand, as 214 we have mentioned, it is not plausible to directly derive the first term on the right-hand side with 215 back-propagation, since black-box LLM  $\mathcal{F}_B(\cdot)$  parameters and gradients are inaccessible. In this 216 case, we propose to tackle this challenge with the zeroth-order gradient approximation technique.

# 4.1.1 ZEROTH-ORDER GRADIENT APPROXIMATION FOR BLACK-BOX LLM GRADIENTS

218 Zeroth-order gradient approximation has been proved effective and efficient for LLM fine-tuning (Malladi et al., 2023), yielding satisfactory results with only a few forward (i.e., inference) passes 219 of LLMs. This makes zeroth-order gradient approximation a promising solution for approximating 220 inaccessible black-box LLM gradients. To begin with, analogous to existing ZO approximation works 221 (e.g., Nesterov & Spokoiny (2017); Ghadimi et al. (2016); Duchi et al. (2015); Shu et al. (2023); 222 Malladi et al. (2023)), we can first suppose a *linear optimization landscape* around each white-box 223 LLM output  $\phi$ . With notation from Eq. 3, it leads to  $\mathcal{L}_{\text{Valid}}(\phi + z) \approx [\nabla_{\phi} \mathcal{L}_{\text{Valid}}(\phi)]^{\intercal} \cdot z + \mathcal{L}_{\text{Valid}}(\phi)$ ; and 224 z is a small perturbation applied to the predicted next-token distribution, for all tokens in the output  $\phi$ , 225 specifically on LLM-header output probabilities (i.e., predicted distribution over the vocabulary). We 226 also include supplementary explanations for auto-regressive generation and perturbation in Appendix 227 C.1. Here, we slightly abuse the notation by using operations "+" and "-" to impose perturbation z 228 onto token-level output probabilities of each token from  $\phi$ . This formulation holds because gradients 229  $\nabla_{\phi} \mathcal{L}_{\text{Valid}}(\phi) := \partial \mathcal{L}_{\text{Valid}}(\phi) / \partial \phi$  will stay uniform for all  $\phi$  within the *linear* landscape.

Here, these small token-level perturbations are imposed to collect information of the optimization landscape, as small perturbations on token-level outputs can effectively change the auto-regressive generation process (Han et al., 2024a). Inspired by Malladi et al. (2023), we can formulate an approximation for black-box LLM gradients as

- 234
- 235 236

 $\frac{\partial \left[ \mathcal{L}_{\text{Valid}} \left( \phi(\boldsymbol{\Theta}_{W}) \right) \right]}{\partial \phi(\boldsymbol{\Theta}_{W})} \approx \frac{\left[ \mathcal{L}_{\text{Valid}} \left( \phi(\boldsymbol{\Theta}_{W}) + \epsilon \boldsymbol{z} \right) - \mathcal{L}_{\text{Valid}} \left( \phi(\boldsymbol{\Theta}_{W}) - \epsilon \boldsymbol{z} \right) \right]}{2\epsilon} \cdot \boldsymbol{z}$ (5)

237 by deeming  $\mathcal{L}_{\text{Valid}}(\cdot)$  as the validating evaluation function for generated instruction  $\phi$  from Eq. 3. 238 Here,  $z \sim \mathcal{N}(0, I) \in \mathbb{R}^d$  stands for a random Gaussian perturbation vector, imposed on each token 239 of the output  $\phi$ , thereby maintaining the same dimensionality as the token-level dimensionality of 240 the output. z also satisfies the isotropic condition  $\mathbb{E}[zz^{\intercal}] = I_d$ . Consequently, d will correspond 241 to the vocabulary size of the white-box LLM. The scaling parameter  $\epsilon \in \mathbb{R}^+$  is used to control the 242 perturbation intensity. This formulation intuitively follows the idea of bi-directional estimation of 243 optimization landscape to perceive the optimization landscape from both directions (Spall, 1992; 244 Malladi et al., 2023). In this way, the first term on the RHS of Eq. 4 can be approximated with only 245 black-box LLM forward passes, without accessing its internal gradients or parameters.

**Remark 1.** We only approximate black-box LLM gradients, instead of using ZO method to directly estimate whole gradient flow  $\partial \left[\frac{1}{|\mathcal{D}_{Valid}|} \sum_{(X,Y) \in \mathcal{D}_{Valid}} \mathcal{L}(\mathcal{F}_B([\phi(\Theta_W); X]), Y)] / \partial \Theta_W$ . The reason is that the error of zeroth-gradient method tends to grow along with the target dimensionality (Malladi et al., 2023). Since the size of white-box LLM parameters  $\Theta_W$  is generally much larger than whitebox LLM output dimensionality, we choose to approximate  $\partial [\mathcal{L}_{Valid}(\mathcal{F}_B(\phi(\Theta_W)), Y)] / \partial \phi(\Theta_W)$ instead for a lower approximation error. Related ablation study is presented in Appendix B.6.

However, one potential drawback is that as perturbation vectors z are randomly sampled (Eq. 5), gradient perturbation directions within the optimization landscape will be random and potentially inefficient (Cai et al., 2022). Thus, we propose to reuse collected feedback, by formulating above ZO-based fine-tuning process as an online sequential decision-making problem; and utilize Contextual Bandit techniques to effectively determine which perturbation directions are informative, beneficial and worth exploring, in terms of improving instruction quality and black-box LLM performance.

258 259 4.1.2 TS-AIDED SELECTION OF GRADIENT PERTURBATION DIRECTIONS

Recall that for ZO-based gradient approximation methods (e.g., Nesterov & Spokoiny (2017); Ghadimi et al. (2016); Duchi et al. (2015); Malladi et al. (2023)), it is common to suppose that we have a linear optimization landscape around the current optimization objective as in Eq. 5. In this case, as illustrated in Figure 2, with random Gaussian perturbation vector z, we have the radius of the supposed linear optimization landscape following a Chi-squared distribution with expected radius being  $\mathbb{E}[\text{Radius}] = \mathbb{E}[||\epsilon z||_2] = \epsilon \sqrt{d}$ .

Leveraging the linear optimization landscape. Within the supposed linear landscape, we can
 intuitively formulate this ZO optimization problem into a sequential decision-making process, where
 collected information can help choose perturbation directions z in Eq. 5, rather than applying
 completely random z. As a natural solution, Contextual Bandit algorithms (Li et al., 2010; Agrawal & Goyal, 2013) are designed to identify the optimal choice among a set of candidate arms (i.e.,

actions) based on arm contextual information and historical records, while addressing the explorationexploitation dilemma in sequential decision-making processes Auer et al. (2002). Under Contextual Bandit settings, and based on the gradient approximation formulation from Eq. 5, we define the arm reward  $r \in \mathbb{R}$  for each perturbation direction  $z \in \mathbb{R}^d$  (i.e., *candidate arm* in Contextual Bandit) as

$$r := \left[\nabla_{\phi} \mathcal{L}_{\text{Valid}}(\phi)\right]^{\mathsf{T}} \cdot \boldsymbol{z} \approx \frac{\mathcal{L}_{\text{Valid}}(\phi + \boldsymbol{\epsilon} \cdot \boldsymbol{z}) - \mathcal{L}_{\text{Valid}}(\phi - \boldsymbol{\epsilon} \cdot \boldsymbol{z})}{2\boldsymbol{\epsilon}},\tag{6}$$

with respect to current white-box LLM output  $\phi$ . Intuitively, our formulation of arm reward echoes with our bi-directional ZO approximation formulated in Eq. 5. This formulation aims to quantify the benefit of descending towards perturbation direction z (i.e., updating output  $\phi$  towards direction -z).



274 275 276

277

278

279

280

281

283

284

287 288

289

302

316

317

290Figure 2: Linear optimization landscape291(illustrated as a sphere) around white-box292LLM output  $\phi$ . With the sampled perturbation direction  $\boldsymbol{z} \sim \mathcal{N}(0, \boldsymbol{I})$ , expected293radius will be  $\epsilon \sqrt{d}$ , in terms of  $L_2$  distance294between averaged output token-level probabilities. Updated outputs  $\phi', \phi''$  can stay295within the linear landscape, motivating our296TS-based approximation.

With our arm reward formulation enabled by the supposed linear optimization landscape, we apply a *linear* TS model to select perturbation directions (i.e., arms) accordingly. Here, TS model parameters with dimensionality d are denoted by lowercase  $\theta \in \mathbb{R}^d$ . Analogous to (Agrawal & Goyal, 2013; Zhang et al., 2021), we consider  $\mathcal{N}(0, I)$  as the initial prior for TS model parameters. Then, starting from the prior, we gradually refine the corresponding TS parameter posterior distribution with collected optimization landscape information, and sample updated TS parameters  $\theta$  from the refined posterior.

Regarding our arm reward formulation in Eq. 6 and the nature of linear TS, the parameters  $\theta$  essentially serve as an estimate of the uniform gradients  $\nabla_{\phi} \mathcal{L}_{\text{Valid}}(\phi)$  within the optimization landscape. In this way, we gradually refine our arm selection strategy by effectively reusing collected black-box LLM feedback. Notably, the linear optimization landscape allows us to employ a highly efficient linear TS algorithm for rapid arm selection and parameter refinement. Next, we proceed with arm (perturbation direction) selection for the ZO gradient approximation procedure.

**TS-aided perturbation direction (arm) selection.** We consider a *T*-round fine-tuning process, and denote the initial white-box LLM parameters without fine-tuning as  $\Theta_0 := \Theta_W$ , with initial generated instruction  $\phi_0 := \phi(\Theta_0)$ . We also let  $\Theta_{t-1}, t \in [T]$  refer to white-box LLM parameters *before t*-th round fine-tuning. Here, in each round *t*, we first sample  $K \in \mathbb{N}^+$  candidate arms (i.e., perturbation directions)  $Z_t$  for selection, from the standard Gaussian distribution, as

$$\mathcal{Z}_t := \left\{ \boldsymbol{z}_{t,1}, \boldsymbol{z}_{t,2}, \dots, \boldsymbol{z}_{t,K} \right\} \sim \mathcal{N}(0, \boldsymbol{I}). \tag{7}$$

303 This design controls the arm context norm magnitude with the isotropic formulation, while ensuring 304 randomness in terms of candidate arm context directions. Before the fine-tuning process, initial 305 TS model parameters are instantiated as  $\theta_0$  by sampling from the prior  $\mathcal{N}(0, I)$ . Here, we let 306  $\theta_{t-1}$  represent the TS model parameters in round t before the refinement. We will discuss later 307 how to update TS parameters using a refined TS parameter posterior in Eq. 9. Next, for these Kcandidate arms, we formulate *estimated rewards* as the inner product  $z_{t,k}^{\mathsf{T}} \theta_{t-1}, \forall z_{t,k} \in \mathbb{Z}_t$ , and select 308 309 the top-*B* arms with the *highest estimated rewards* to form collection  $\tilde{Z}_t \subset Z_t$ , with cardinality 310  $|\hat{Z}_t| = B, B \ll K$ . The chosen arms  $\hat{Z}_t$  are considered perturbation directions that can lead to 311 potential benefits for reducing the validation loss. Calculation details are in lines 5-7 of Algorithm 1. 312

Querying arm rewards. Next, we query the black-box LLM (i.e., reward oracle) for validation loss results  $\mathcal{L}_{\text{Valid}}$  to obtain rewards for the chosen arms  $\widetilde{\mathcal{Z}}_t$ . Using the shorthand  $\phi_{t-1} := \phi(\Theta_{t-1})$ , for each chosen arm  $z_{t,k} \in \widetilde{\mathcal{Z}}_t$ , we calculate its arm reward by following Eq. 6, leading to

 $r_{1}$ 

$$\mathbf{x}_{k} = \frac{\mathcal{L}_{\text{Valid}}(\phi_{t-1} + \epsilon \cdot \boldsymbol{z}_{t,k}) - \mathcal{L}_{\text{Valid}}(\phi_{t-1} - \epsilon \cdot \boldsymbol{z}_{t,k})}{2\epsilon},\tag{8}$$

which measures the benefit of involving direction (arm)  $z_{t,k}$  for gradient approximation. Naturally, with each chosen arm (perturbation direction), the queried validation loss results are recycled to derive the black-box LLM gradient approximation according to Eq. 5. Finally, by plugging in the estimated black-box LLM gradients, white-box LLM parameters can be updated via Gradient Descent through the gradient flow (Eq. 4) and estimated black-box LLM gradients (Eq. 5). Gradient estimation results from the *B* chosen perturbation direction vectors (arms) are averaged for a more accurate approximation, analogous to existing ZO approximation methods (Malladi et al., 2023). 324 Validating optimization landscape. After querying arm rewards with Eq. 8 and performing white-325 box LLM parameter fine-tuning, we proceed to update the TS model parameters. Recall that we 326 operate within a linear optimization landscape, if generated instructions do not significantly deviate 327 from previous ones. In this case, we apply a threshold parameter  $\beta > 0$  to quantify the landscape magnitude. For the optimized instruction  $\phi_t$  in round t, we use the condition  $\|\phi_t - \phi_{\text{Check}}\| > \beta$  to 328 verify whether the assumed linear landscape remains valid. The initial checkpoint  $\phi_{\text{Check}}$  is set as the 329 instruction  $\phi_0$  prior to optimization. We calculate the  $L_2$  distance between the averaged token-level 330 probabilities of the output  $\phi_t$  and the checkpoint  $\phi_{\text{Check}}$ , inspired by prior works (e.g., Joshi et al. 331 (2023); Manakul et al. (2023)). Collected records are initialized as an empty set  $\Omega_0$ . 332

333 **Updating TS parameters.** In each optimization round  $t \in [T]$ , we have *Scenario (1)*: If white-box LLM output becomes far enough from the checkpoint, s.t.  $\|\phi_t - \phi_{\text{Check}}\| > \beta$ , our current knowledge 334 can be invalid because current white-box LLM output has significantly deviated from the checkpoint. 335 In this case, we set the new checkpoint as  $\phi_t$  and discard collected records. Then, reinitialize TS 336 parameters  $\theta_t$  from the prior  $\mathcal{N}(0, I)$ . Scenario (2): Otherwise, if the distance is small enough s.t., 337  $\|\phi_t - \phi_{\text{Check}}\| \leq \beta$ , the chosen arms and their true rewards from this round will be integrated into 338 collected records  $\Omega_t$ . Afterwards, analogous to existing TS methods (Agrawal & Goyal, 2013; Zhang 339 et al., 2021), with an exploration parameter  $\nu \ge 0$ , covariance matrix  $\Sigma_t := I + \sum_{(z,r)\in\Omega_t} z \cdot z^{\overline{1}}$ , 340 and reward vector  $b_t := \sum_{(\boldsymbol{z},r) \in \Omega_t} \boldsymbol{z} \cdot r$ , we update the posterior of TS parameters as 341

- 342
- 343

$$\mathcal{N}(\boldsymbol{\Sigma}_t^{-1}\boldsymbol{b}_t, \ \boldsymbol{\nu} \cdot \boldsymbol{\Sigma}_t^{-1}). \tag{9}$$

Finally, we update TS parameters  $\theta_t$  by sampling from the updated posterior, and proceed to the next optimization round. Additional calculation details are presented in lines 12-19 of Algorithm 1.

**Remark 2.** To reduce computational costs of matrix inversion and sampling from high-dimensional 346 Gaussian distribution, motivated by Johnson-Lindenstrauss (JL) Lemma (Johnson & Lindenstrauss, 347 1984), we adopt random Gaussian projection to map d-dimensional arm contexts into a lower-348 dimensional space (Matoušek, 2008; Larsen & Nelson, 2017), where we perform the TS-aided 349 selection of candidate arms  $\mathcal{Z}_t$ . Comparable ideas are also applied in existing works for reducing 350 the dimensionality of tunable soft prompt vectors (Chen et al., 2024; Lin et al., 2024). To efficiently 351 compute the inversion of the covariance matrix  $\Sigma_t$  in each round t, we utilize the Sherman-Morrison 352 formula (Bartlett, 1951; Maponi, 2007), avoiding direct matrix inversion operations. Details will be 353 elaborated in Appendix C.2 due to page limit. 354

#### 4.2 WORKFLOW SUMMARY AND PSEUDO-CODE FOR AIO FRAMEWORK

356 The pseudo-code of AIO is in Algorithm 1. For each optimization round  $t \in [T]$ , we first sample a 357 pool of K candidate arms (gradient approximation directions)  $\mathcal{Z}_t$  (line 5, Algorithm 1). Then, we 358 apply a TS-based bandit model to estimate arm rewards, which quantify the benefit of including the 359 corresponding arms as perturbation directions. To reduce API costs, we only select  $B \ll K$  arms as 360 the chosen arms  $\mathcal{Z}_t \subset \mathcal{Z}_t$  (lines 6-7). Next, we query the black-box LLM to obtain rewards of the 361 chosen arms (line 9) and perform Gradient Descent to fine-tune the white-box LLM parameters with 362 the gradient flow described in Eqs. 4 and 5 (line 10). Afterwards, we check if the white-box LLM 363 output after fine-tuning differs sufficiently from the checkpoint. If so, we reset the records and the TS parameter distribution (line 13). Otherwise, we update the parameter posterior with the chosen arms 364 and their true rewards (lines 16-17). TS parameters  $\theta_t$  are updated accordingly (lines 14, 18).

366 367 5 EXPERIMENTS

368 Our empirical analysis mainly aims to show that AIO can effectively optimize task-specific black-box 369 LLM instructions compared with strong baselines, as well as provide insights on AIO behaviours and 370 properties. In terms of LLM implementations, we apply Llama-3-8B-Instruct (Dubey et al., 371 2024) as our tunable white-box LLM  $\mathcal{F}_W(\cdot; \Theta_W)$ , and adopt Claude-3-Sonnet (Anthropic, 2024) as our black-box LLM  $\mathcal{F}_B(\cdot)$ . As an outline for our empirical results in the main body, we have: 372 (1) zero-shot instruction induction experiments on 15 tasks in Subsec. 5.1; (2) empirical analysis on 373 API token costs and performance in Subsec. 5.2; (3) a case study that provides analysis and examples 374 for AIO instruction optimization trajectories in Subsec. 5.3. Due to page limit, we include experiment 375 and implementation details in Appendix A. Meanwhile, we also conduct additional experiments 376 presented in Appendix B, including but not limited to few-shot instruction induction results, ablation 377 study on AIO, applying AIO under Chain-of-Thought (CoT) settings, as well as empirical results with different combinations of white-box LLMs and black-box LLMs.

Alg	orithm 1 Automatic Instruction Optimizer (AIO)
1: 2:	<b>Input:</b> Optimization rounds T. Exemplars $\mathcal{D}_{\text{Train}}$ , validation data $\mathcal{D}_{\text{Valid}}$ . Number of candidate arms K. Number of chosen arms B. Exploration parameter $\nu \geq 0$ . Threshold parameter $\beta > 0$ . <b>Initialization:</b> TS Model Parameters $\theta_0 \sim \mathcal{N}(0, I)$ . White-box LLM parameters $\Theta_0 \leftarrow \Theta_W$ .
	Instruction checkpoint $\phi_{\text{Check}} \leftarrow \phi(\Theta_0)$ . TS model records $\Omega_0 \leftarrow \emptyset$ .
3:	for each round $t \in [T]$ do
4:	▷ TS-aided ZO Perturbation Direction Selection
5: 6:	Sample candidate perturbation directions (i.e., candidate arms) $\mathcal{Z}_t$ (Eq. 7), with $ \mathcal{Z}_t  = K$ . Calculate arm reward estimations $\hat{r}_{t,k} = \mathbf{z}_{t,k}^{T} \boldsymbol{\theta}_{t-1}, \forall \mathbf{z}_{t,k} \in \mathcal{Z}_t$ , with TS parameters $\boldsymbol{\theta}_{t-1}$ .
7:	Choose B arms of highest estimated rewards $\widetilde{\mathcal{Z}}_t \leftarrow \arg \max_{\widetilde{\mathcal{Z}}_t \subset \widetilde{\mathcal{Z}}_t :  \widetilde{\mathcal{Z}}_t  = B} \left[ \sum_{\tau_{t,t} \in \widetilde{\mathcal{Z}}_t} \hat{r}_{t,k} \right].$
8:	$\triangleright$
9:	Ouery rewards for chosen arms $\widetilde{\mathcal{Z}}_{4}$ (Eq. 8). $\triangleright$ <b>Only ouery chosen arms to reduce API cost.</b>
0:	With <i>B</i> chosen perturbation directions (arms) $\widetilde{Z}_t$ and their queried rewards, fine-tune white- box LLM parameters to $\Theta_t$ with Gradient Descent, based on gradient flow decomposition (Eq. 4 and Eq. 5) and validation loss (Eq. 3). Generate updated instruction $\phi_t := \phi(\Theta_t)$ .
11:	▷ Updating Linear TS Model
2:	if $\ \phi_t - \phi_{\text{Check}}\  > \beta$ then
13:	Reset prior as $\mathcal{N}(0, \mathbf{I})$ , new checkpoint $\phi_{\text{Check}} \leftarrow \phi_t$ , and collected records $\Omega_t \leftarrow \emptyset$ .
4: 5.	Sample updated TS parameters $\theta_t \sim \mathcal{N}(0, I)$ .
.5: 16:	With chosen arms $z \in \widetilde{Z}_{r}$ and their rewards r undate $\Omega_{r} \leftarrow \Omega_{r-1} \sqcup [1 - \gamma(z, r)]$
0. 7.	Undete the negterior for TS peremeters by $\mathcal{N}(\mathbf{\Sigma}^{-1}\mathbf{b} - \mathbf{v}\mathbf{\Sigma}^{-1})$ (Eq. 0)
17.	Sample undated TS parameters from the posterior distribution $\boldsymbol{\theta}_{t} = \mathcal{N}(\boldsymbol{\Sigma}^{-1}\boldsymbol{b}_{t}, \boldsymbol{\nu}\boldsymbol{\Sigma}^{-1})$
9:	end if
20:	end for
5.1 We	EXPERIMENTS ON ZERO-SHOT INSTRUCTION INDUCTION
wor emp et al Con whe	ks for LLM instruction optimization (Zhou et al., 2022; Chen et al., 2024; Lin et al., 2024), our pirical analysis involves 15 different tasks including instruction induction tasks from Honovich I. (2022), as well as more challenging reasoning tasks from BigBench (bench authors, 2023). Issequently, evaluation criteria will vary depending on specific tasks, such as "Multiple Choice" re the white-box LLM needs to choose the right option out of several candidates, and "Exact
Mat task	ch" where black-box LLM answers needs to be identical to ground-truth labels. We defer detailed descriptions and evaluation criteria to Appendix A.1.
Bas	eline methods. We involve four baselines, including two kinds of LLM-based instruction
opti	mization methods. The first two baselines leverage black-box LLM for instruction generation:
1)	APE (Zhou et al., 2022), (2) ProTeGi (Pryzant et al., 2023). We also include baselines that
ıtıli ⊤∶∽	ze white-box LLM for instruction generation: (3) InstructZero (Chen et al., 2024), (4) INSTINCT
ຸ່ມຫ ໄໄຂ	uide-3-Sonnet as the black-box instruction generation I I M for APE and ProTeGi while
ado	pting white-box LLM Llama-3-8B-Instruct for InstructZero and INSTINCT.
WC	<b>PEFI</b> variants of AIO. Apart from the original AIO framework in Algorithm 1, recall that

**Two PEFT variants of AIO.** Apart from the original AIO framework in Algorithm 1, recall that AIO is also compatible with many existing PEFT methods for efficient white-box LLM fine-tuning. Therefore, to reinforce fine-tuning efficiency, we further include empirical results of incorporating Linear Probing (LP) (Kumar et al., 2022) and LoRA (Hu et al., 2021a) to our proposed AIO. These two variants are denoted as "AIO + LP" and "AIO + LoRA" respectively. In particular, we note that "AIO + LP" only fine-tunes ~ 6.54% of white-box LLM parameters, while "AIO + LoRA" merely needs to fine-tune ~ 0.04% of white-box LLM parameters.

Empirical results. Our empirical results are shown in Table 1. We notice that our proposed AIO can generally achieve better performance in comparison with strong baselines, in the presence of the challenging reasoning tasks from BigBench (bench authors, 2023). In particular, our light-

	Black-box LLM		White-box LLM		White-box LLM w/ FT (Ours)		
Tasks $\setminus$ Methods	APE	ProTeGi	InstructZero	INSTINCT	AIO	AIO + LP	AIO + LoRA
antonyms	0.893	0.861	0.843	0.881	0.901	0.857	0.898
sentiment	0.911	0.928	0.941	0.920	0.949	0.967	0.947
larger_animal	0.914	0.932	0.827	0.857	0.912	0.945	0.950
taxonomy_animal	0.491	0.970	0.598	0.782	0.983	0.979	0.935
object_counting	0.319	0.550	0.522	0.537	0.543	0.401	0.479
navigate	0.580	0.624	0.556	0.577	0.644	0.623	0.627
winowhy	0.022	0.703	0.671	0.725	0.622	0.646	0.635
implicatures	0.806	0.826	0.816	0.837	0.811	0.836	0.849
logical_fallacy	0.820	0.826	0.790	0.826	0.868	0.824	0.836
hyperbaton	0.515	0.499	0.467	0.502	0.538	0.518	0.527
epistemic_reasoning	0.604	0.459	0.667	0.580	0.766	0.784	0.719
movie_recommendation	0.348	0.847	0.895	0.866	0.902	0.857	0.883
timedial	0.532	0.718	0.786	0.712	0.814	0.734	0.759
presuppositions_as_nli	0.458	0.488	0.503	0.482	0.523	0.486	0.493
question_selection	0.712	0.667	0.718	0.605	0.648	0.628	0.622
Average Rank	3.87	2.80	3.13	3.20		2.00	

450 Table 1: Zero-shot Instruction Induction Results. For each task (row), **bold** number refers to the best result, 451 while underlined number refers to the second-best one. AIO and its two variants can outperform four baselines on 12 out of a total of 15 tasks. We average results of AIO and its two variants task-wise, and treat these three 452 methods as a unified baseline for ranking comparisons. 453

weight variants "AIO + LP" and "AIO + LoRA", which fine-tune significantly less white-box LLM parameters, can also achieve promising performance, or outperform other baselines on tasks like 'epistemic reasoning". To perform a more comprehensive ranking comparison, we present ranking results by averaging the performance of AIO and its two variants for each task, and applying averaged performance for comparing against baselines to derive ranking. As in Table 1, AIO and its variants as a unity can still enjoy relatively higher ranking results compared with baselines. We also include additional complementary experiments, such as few-shot instruction induction and an ablation study of AIO components, so that interested readers can refer to Appendix B for details.

#### **OPTIMIZATION PERFORMANCE VERSUS API TOKEN EFFICIENCY**



455

456

457

458

459

460

461 462

463

464

465

467

468

469

471

472

473

475

477

478

480

481

482

Figure 3: Token consumption vs. performance (accu-483 racy results). Token consumption results are normalized into [0, 1] range. In the below figure, we include token 484 consumption vs. best accuracy results till certain token 485 consumption levels on four tasks, with two instruction fragments for "Implicatures" task at different stages.

Recall that we apply Claude-3-Sonnet as our blackbox LLM  $\mathcal{F}_B(\cdot)$  for experiments, where API query costs are charged on a token-basis for end users. In Figure 3, we show an illustration in terms of token consumption, with input and output token quantities combined, versus instruction induction performance on four different tasks. From Figure 3, we see that AIO can relatively maintain a good balance between token costs and induction performance, starting from early optimization stages when small amounts of tokens are consumed. Compared with ProTeGi, the performance of APE tends to be inferior as ProTeGi can optimize generated instructions with highergranularity error feedback in terms of specific training samples, although it can lead to additional token consumption costs. On the other hand, we observe that the baselines InstructZero and INSTINCT generally have higher token consumption compared to AIO. While these two baselines also leverage white-box LLMs for instruction optimization, their methods primarily rely on in-context learning, by tuning a prefix soft prompt with kernel-based learner or small neural model. Given the extreme complexity of black-box LLMs, their representation power can be insufficient for effectively learning from black-box LLM feedback, leading to more interactions with the black-box LLM and, consequently, higher token costs. Alternatively,



Figure 4: In left Figure, we show three instructions generated for "Sentiment" task: w/o FT, during FT, and after FT. Instructions with FT generalize to task contexts instead of over-fitting task exemplars. Middle and right figures shows the normalized distance vs. accuracy for tasks "Larger Animal", "Question Selection" (middle figure with fitted lines), "Navigate" (right radar plot). For radar plot, instruction (i.e., point) accuracy increases clock-wisely starting from  $0^\circ$ . Point distances to circle center are corresponding distances from  $\phi_0$ .

to tackle the complexity of the black-box LLM, we fine-tune a white-box LLM to leverage its rich representation power, enabling efficient instruction optimization with adequate utilization of black-box LLM feedback.

#### 5.3 CASE STUDY: INSTRUCTION OPTIMIZATION TRAJECTORY ANALYSIS

In this subsection, we present analysis on how AIO generated instructions evolve across optimization 504 rounds. From the left figure in Figure 4, we provide some examples in terms of how fine-tuning 505 changes the way our white-box LLM  $\mathcal{F}_W$  comprehends exemplars and composes instructions. These 506 examples all originate from our experiment data. Without fine-tuning,  $\mathcal{F}_W$  can over-fit exemplars 507 by searching specific keywords to judge sentiment outcomes, which can clearly fail to generalize 508 to unseen task data points and lead to low accuracy. During and after fine-tuning,  $\mathcal{F}_W$  perceives 509 black-box LLM feedback and corrects its way in terms of interpreting task exemplars as well as 510 generating instructions. During fine-tuning, instruction gets improved by considering signal words 511 from exemplars (e.g., "entertainingly") as examples instead of sole indicators. Moreover, after fine-tuning,  $\mathcal{F}_W$  further generalizes the concept "positive connotations = positive sentiment" to a 512 513 more comprehensive view, by mentioning "use sentiment analysis" to analyze outcome directly.

514 In addition, we also present visualization results in terms of induction performance as well as distances 515 between generated instructions  $\phi$  and the initial instruction  $\phi_0 := \phi(\Theta_W)$ . Instruction distances 516 are measured by averaged token probabilities as in Algorithm 1. With the middle line chart, we 517 individually plot a fitting line for instruction points of each task. For tasks like "Question Selection", 518 our optimization trajectory tends to fit well with the supposed linear optimization landscape where 519 small residuals are observed. On the other hand, for tasks like "Larger Animal", we can observe 520 relatively higher fluctuations when generated instructions deviate from  $\phi_0$ . In this case, as it still obeys a relatively fitting linear trend until the normalized distance reaches  $0.7 \sim 0.8$ , we can supposed 521 a smaller linear landscape in terms of TS-aided optimization. Shown in the radar plot, regarding 522 the "Navigate" task, we see that good instructions are not necessarily far away from the initial  $\phi_0$ , 523 while the best instructions tend to fall into a sub-area within the optimization landscape, instead of 524 consistently deviates from the starting point  $\phi_0$  as accuracy results increase. Due to page limit, we 525 will include additional instruction optimization trajectory examples in Appendix B.12. 526

## 527 6 CONCLUSION

495

496

497 498

499

500

501 502

503

In this paper, we propose a novel framework named Automatic Instruction Optimizer (AIO) to 529 adaptively customize instructions for various downstream tasks. By applying a task-solver black-box 530 LLM for query answering, AIO fine-tunes a white-box LLM into a task-aware instruction optimizer 531 that learns from high-level task-relevant information and black-box LLM feedback, to generate high-532 quality instructions for the task-solver black-box LLM. Distinct from existing in-context learning 533 approaches, our framework is designed to address the formidable complexity of modern black-box 534 LLMs with possibly hundreds of billions of parameters involved. To overcome the challenges of inaccessible black-box LLM gradients and mitigate concerns related to expensive black-box LLM API costs, our AIO framework leverages a novel TS-aided zeroth-order gradient approximation method, enabling effective and efficient learning of task-aware instructing strategies. Extensive 537 experiments demonstrate the superiority of our proposed framework in terms of performance and API 538 token efficiency, along with additional analyses that highlight AIO's properties and specifications. Additional discussions on AIO future extensions are presented in Appendix E.

## 540 REFERENCES

547

564

565

566 567

568

569

570

574

575

576

577

586

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

- Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *ICML*, pp. 127–135. PMLR, 2013.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over parameterization. In *International Conference on Machine Learning*, pp. 242–252. PMLR, 2019.
- Xavier Amatriain. Prompt design and engineering: Introduction and advanced methods. *arXiv preprint arXiv:2401.14423*, 2024.
- Anthropic. The claude 3 model family: Opus, sonnet, haiku. Techni cal report, Anthropic, 2024. URL https://www-cdn.anthropic.com/
   de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model\_Card\_Claude\_3.pdf.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit
   problem. *Machine learning*, 47(2-3):235–256, 2002.
- Maurice S Bartlett. An inverse matrix adjustment arising in discriminant analysis. *The Annals of Mathematical Statistics*, 22(1):107–111, 1951.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
  - BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=uyTL5Bvosj.
  - 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.
- HanQin Cai, Daniel McKenzie, Wotao Yin, and Zhenliang Zhang. Zeroth-order regularized optimiza tion (zoro): Approximately sparse gradients and adaptive sampling. *SIAM Journal on Optimization*, 32(2):687–714, 2022.
  - Lichang Chen, Jiuhai Chen, Tom Goldstein, Heng Huang, and Tianyi Zhou. Instructzero: Efficient instruction optimization for black-box large language models. In *Forty-first International Conference on Machine Learning*, 2024.
- Qi Chen, Changjian Shui, and Mario Marchand. Generalization bounds for meta-learning: An
   information-theoretic analysis. *Advances in Neural Information Processing Systems*, 34:25878–
   25890, 2021.
- Alexandra Chronopoulou, Matthew E Peters, and Jesse Dodge. Efficient hierarchical domain adaptation for pretrained language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1336–1351, 2022.
  - Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *AISTATS*, pp. 208–214, 2011.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
   Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
   arXiv preprint arXiv:2407.21783, 2024.

632

- John C Duchi, Michael I Jordan, Martin J Wainwright, and Andre Wibisono. Optimal rates for
   zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.
- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rock täschel. Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv preprint arXiv:2309.16797*, 2023.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- 604 Peter I Frazier. A tutorial on bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.
- Noa Garcia, Chentao Ye, Zihua Liu, Qingtao Hu, Mayu Otani, Chenhui Chu, Yuta Nakashima, and Teruko Mitamura. A dataset and baselines for visual question answering on art. In *Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pp. 92–108. Springer, 2020.
- Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods
   for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1):267–305, 2016.
- 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. In *The Twelfth International Conference on Learning Representations*, 2024.
- Chi Han, Jialiang Xu, Manling Li, Yi Fung, Chenkai Sun, Nan Jiang, Tarek Abdelzaher, and Heng Ji.
  Word embeddings are steers for language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16410–16430, 2024a.
- Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024b.
- Or Honovich, Uri Shaham, Samuel R Bowman, and Omer Levy. Instruction induction: From few
   examples to natural language task descriptions. *arXiv preprint arXiv:2205.10782*, 2022.
- Yutai Hou, Hongyuan Dong, Xinghao Wang, Bohan Li, and Wanxiang Che. Metaprompting: Learning to learn better prompts. In *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 3251–3262, 2022.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021a.
- Wenyang Hu, Yao Shu, Zongmin Yu, Zhaoxuan Wu, Xiangqiang Lin, Zhongxiang Dai, See-Kiong
   Ng, and Bryan Kian Hsiang Low. Localized zeroth-order prompt optimization. *arXiv preprint arXiv:2403.02993*, 2024.
- Yifan Hu, Xin Chen, and Niao He. On the bias-variance-cost tradeoff of stochastic optimization.
   *Advances in Neural Information Processing Systems*, 34:22119–22131, 2021b.
- William B. Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into hilbert space.
   *Contemporary mathematics*, 26:189–206, 1984. URL https://api.semanticscholar.
   org/CorpusID:117819162.
- Harshit Joshi, José Cambronero Sanchez, Sumit Gulwani, Vu Le, Gust Verbruggen, and Ivan Radiček.
  Repair is nearly generation: Multilingual program repair with llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 5131–5140, 2023.
- Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts,
   and Matei Zaharia. Demonstrate-search-predict: Composing retrieval and language models for
   knowledge-intensive NLP. arXiv preprint arXiv:2212.14024, 2022.

648 649 650	Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller,
651	self-improving pipelines. arXiv preprint arXiv:2310.03714, 2023.
652	Terrer Kerkel, Keiler Shi, Angelier Chan, Desiler Vinsuel, Dhelener, Christenher Ducklass, Jacob
653	Dhang, Samuel P, Bowman, and Ethan Daraz, Pretraining language models with human preferences
654	In International Conference on Machine Learning pn 17506–17533 PMI R 2023
655	In merhanoma Conjerence on Machine Learning, pp. 17500–17555. 1 MLK, 2025.
656	Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can
657	distort pretrained features and underperform out-of-distribution. arXiv preprint arXiv:2202.10054,
658	2022.
659	Kasper Green Larsen and Jelani Nelson. Optimality of the johnson-lindenstrauss lemma. In 2017
660	IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), pp. 633–638. IEEE,
661	2017.
662	Brian Lester Dami Al Dfou, and Noah Constant. The nowar of scale for persmater efficient prompt
663 664	tuning. arXiv preprint arXiv:2104.08691, 2021.
665	Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
666	personalized news article recommendation. In WWW, pp. 661–670, 2010.
667	Xiang Lisa Li and Percy Liang Prefix-tuning: Optimizing continuous prompts for generation arXiv
668	preprint arXiv:2101.00190, 2021.
669	
670	Xiaoqiang Lin, Zhaoxuan Wu, Zhongxiang Dai, Wenyang Hu, Yao Shu, See-Kiong Ng, Patrick
671	Jaillel, and Bryan Kian Hsiang Low. Use your instituct: Instruction optimization for firms using neural handits coupled with transformers. In Forty first International Conference on Machine
672	Learning 2024
674	
675	Alisa Liu, Xiaochuang Han, Yizhong Wang, Yulia Tsvetkov, Yejin Choi, and Noah A Smith. Tuning
676	language models by proxy. arXiv preprint arXiv:2401.08565, 2024.
677	Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and
678	Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
679	learning. Advances in Neural Information Processing Systems, 35:1950–1965, 2022.
680	Ruotian Ma, Xiaolei Wang, Xin Zhou, Jian Li, Nan Du, Tao Gui, Qi Zhang, and Xuanjing Huang.
681	Are large language models good prompt optimizers? <i>arXiv preprint arXiv:2402.02101</i> , 2024.
602	Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev
684	Arora. Fine-tuning language models with just forward passes. arXiv preprint arXiv:2305.17333,
685	2023.
686	Potsawee Manakul, Adian Liusie, and Mark Gales. Selfcheckgpt: Zero-resource black-box hallucina-
687	tion detection for generative large language models. In Proceedings of the 2023 Conference on
688	Empirical Methods in Natural Language Processing, pp. 9004–9017, 2023.
689	Pierluigi Manoni The solution of linear systems by using the sherman-morrison formula <i>Linear</i>
690	algebra and its applications, 420(2-3):276–294, 2007.
691	
692	JITI Matousek. On variants of the johnson–lindenstrauss lemma. <i>Random Structures &amp; Algorithms</i> , 22(2):142–156, 2008
693	JJ(2).1 <del>1</del> 2–1JU, 2UU0.
694	Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. Few-shot fine-
695	tuning vs. in-context learning: A fair comparison and evaluation. <i>arXiv preprint arXiv:2305.16938</i> ,
696	2023.
697	Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions.
600	Foundations of Computational Mathematics, 17(2):527–566, 2017.
700	Shirui Pan Linhao Luo, Yufei Wang, Chen Chen, Jianu Wang, and Xindong Wu. Unifying large
701	language models and knowledge graphs: A roadmap. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 2024.

702 703 704	Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems? <i>arXiv preprint arXiv:2103.07191</i> , 2021.
705 706	Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. Grips: Gradient-free, edit-based instruction search for prompting large language models. <i>arXiv preprint arXiv:2203.07281</i> , 2022.
707 708 709	Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. Automatic prompt optimization with "gradient descent" and beam search. <i>arXiv preprint arXiv:2305.03495</i> , 2023.
710 711 712	Ruizhong Qiu and Hanghang Tong. Gradient compressed sensing: A query-efficient gradient estimator for high-dimensional zeroth-order optimization. In <i>Forty-first International Conference on Machine Learning</i> , 2024.
713 714 715 716	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
717 718 719	Laria Reynolds and Kyle McDonell. Prompt programming for large language models: Beyond the few-shot paradigm. In <i>Extended abstracts of the 2021 CHI conference on human factors in computing systems</i> , pp. 1–7, 2021.
720 721 722	Teven Le Scao and Alexander M Rush. How many data points is a prompt worth? <i>arXiv preprint arXiv:2103.08493</i> , 2021.
723 724 725	Tobias Schnabel and Jennifer Neville. Prompts as programs: A structure-aware approach to efficient compile-time prompt optimization. <i>arXiv preprint arXiv:2404.02319</i> , 2024.
726 727 728	Eric Schulz, Maarten Speekenbrink, and Andreas Krause. A tutorial on gaussian process regression: Modelling, exploring, and exploiting functions. <i>Journal of mathematical psychology</i> , 85:1–16, 2018.
729 730 731 732	Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the human out of the loop: A review of bayesian optimization. <i>Proceedings of the IEEE</i> , 104(1): 148–175, 2015.
733 734 735	Jiho Shin, Clark Tang, Tahmineh Mohati, Maleknaz Nayebi, Song Wang, and Hadi Hemmati. Prompt engineering or fine tuning: An empirical assessment of large language models in automated software engineering tasks. <i>arXiv preprint arXiv:2310.10508</i> , 2023.
736 737 738 729	Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. <i>arXiv preprint arXiv:2010.15980</i> , 2020.
740 741 742	Yao Shu, Zhongxiang Dai, Weicong Sng, Arun Verma, Patrick Jaillet, and Bryan Kian Hsiang Low. Zeroth-order optimization with trajectory-informed derivative estimation. In <i>The Eleventh International Conference on Learning Representations</i> , 2023.
743 744 745	James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. <i>IEEE transactions on automatic control</i> , 37(3):332–341, 1992.
746 747	Haotian Sun, Yuchen Zhuang, Wei Wei, Chao Zhang, and Bo Dai. Bbox-adapter: Lightweight adapting for black-box large language models. <i>arXiv preprint arXiv:2402.08219</i> , 2024.
748 749 750 751	Xinyu Tang, Xiaolei Wang, Wayne Xin Zhao, Siyuan Lu, Yaliang Li, and Ji-Rong Wen. Unleashing the potential of large language models as prompt optimizers: An analogical analysis with gradient-based model optimizers. <i>arXiv preprint arXiv:2402.17564</i> , 2024.
752 753 754	Michal Valko, Nathan Korda, Rémi Munos, Ilias Flaounas, and Nello Cristianini. Finite-time analysis of kernelised contextual bandits. In <i>Uncertainty in Artificial Intelligence</i> , 2013.
755	Xilu Wang, Yaochu Jin, Sebastian Schmitt, and Markus Olhofer. Recent advances in bayesian optimization. <i>ACM Computing Surveys</i> , 55(13s):1–36, 2023.

756	Xingjiao Wu, Luwei Xiao, Yixuan Sun, Junhang Zhang, Tianlong Ma, and Liang He. A survey of
757	human-in-the-loop for machine learning. Future Generation Computer Systems, 135:364–381,
758	2022.
759	

- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning, and Christopher Potts. Reft: Representation finetuning for language models. *arXiv preprint arXiv:2404.03592*, 2024.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun
   Chen. Large language models as optimizers, 2024. URL https://arxiv.org/abs/2309.
   03409.
- Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural thompson sampling. In International Conference on Learning Representations, 2021.
- Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based exploration.
   In *International Conference on Machine Learning*, pp. 11492–11502. PMLR, 2020.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*, 2022.

813

861 862

# A COMPLEMENTARY DETAILS FOR IMPLEMENTATION AND EXPERIMENTS

### A.1 DESCRIPTIONS FOR TASKS INVOLVED IN OUR EXPERIMENTS

Task	Metric	Descriptions
antonyms	Exact Match	Find the antonym for the given word.
sentiment	Binary Choice	Judge the sentiment preference of the given review.
larger animal	Binary Choice	Identify which of the input animals is larger.
taxonomy animal	Exact Set	Identify all the animal words out of input word sequence
object_counting	Exact Match	Enumerate objects of different types and output the total number.
navigate	Binary Choice	Given a series of navigation instructions, determine whether one would end up back at the starting point.
winowhy	Binary Choice	Evaluate the reasoning in answering Winograd Schema Challenge questions.
implicatures	Binary Choice	Predict whether Speaker 2's answer to Speaker 1 counts as a yes or as a no.
logical_fallacy	Binary Choice	Detect informal and formal logical fallacies.
hyperbaton	Binary Choice	Order adjectives correctly in English sentences.
epistemic_reasoning	Binary Choice	Determine whether one sentence entails the next.
movie_recommendation	Multiple Choice	Recommend a movie that is similar to the given list of movies.
timedial	Multiple Choice	Pick the correct choice for a masked (temporal) span given the dialog context.
presuppositions_as_nli	Multiple Choice	Determine whether the first sentence entails or contradicts the second.
question_selection	Multiple Choice	Given a short answer along with its context, select the most appropriate question which has the given short answer as its answer.
r evaluation metric bel; (2) Multiple Ch	s: we have (1) oice: task-solv	Exact Match: the generated answer needs to exactly match the ver LLM needs to choose one correct option out of several given
late choices; (3) date choices; (4 bel set in both c TEMPLATES A INFERENCE (Z	Binary Choice Exact Set: who ontent and size PPLIED FOR A ERO-SHOT IN	e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION)
<ul> <li>ndidate choices; (3)</li> <li>ndidate choices; (4)</li> <li>e label set in both c</li> <li>2 TEMPLATES A INFERENCE (Z)</li> <li>r zero-shot instructurations, we followen designing the ir</li> <li>• Few-shot instructurations</li> </ul>	Binary Choice ) Exact Set: will ontent and size PPLIED FOR A ERO-SHOT INI- tion induction w analogous id struction induc-	e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generated eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024 ction templates.
didate choices; (3) didate choices; (4 label set in both c TEMPLATES A INFERENCE (Z zero-shot instruc ructions, we follow n designing the ir • Few-shot ins	Binary Choice ) Exact Set: wl ontent and size PPLIED FOR A ERO-SHOT INI tion induction v analogous id struction induc truction induc	e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generated eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024 ction templates.
date choices; (3) date choices; (4) bel set in both c TEMPLATES A INFERENCE (Z ero-shot instruc ctions, we follow designing the ir • Few-shot ins <examples Exemplary Instructi</examples 	Binary Choice DExact Set: who ontent and size PPLIED FOR A ERO-SHOT INI- tion induction w analogous id istruction induc truction induc truction induc data: [Exe son: [INSTRU	e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generated eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024 ction templates. tion settings: mplar data ([DEMO_DATA])] CTION]\n\n
didate choices; (3) didate choices; (4 label set in both c 2 TEMPLATES A INFERENCE (Z zero-shot instruc ructions, we follow en designing the ir • Few-shot ins	Binary Choice DExact Set: wh ontent and size PPLIED FOR A ERO-SHOT INI- tion induction w analogous id istruction induc truction induc truction induc data: [Exe solution [INSTRU aery INPUT]	<pre>e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generated eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024 ction templates. tion settings: mplar data ([DEMO_DATA])] CTION]\n\n \n Output: [OUTPUT Placeholder]</pre>
didate choices; (3) didate choices; (4 label set in both c 2 TEMPLATES A INFERENCE (Z 2 zero-shot instruc cructions, we follow en designing the ir • Few-shot ins { <examples Exemplary Instructi Input: [Q • Zero-shot in</examples 	Binary Choice Description of the second state	<pre>e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generated eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024 ction templates. tion settings: mplar data ([DEMO_DATA])] CTION]\n\n \n Output: [OUTPUT Placeholder]</pre>
<pre>didate choices; (3) didate choices; (4) label set in both c      TEMPLATES A     INFERENCE (Z zero-shot instruc ructions, we follow n designing the ir      Few-shot inss         <examples <="" [q="" example="" exemplary="" input:="" instruction="" pre=""></examples></pre>	Binary Choice Description of the second state	<pre>e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generated eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024 ction templates. tion settings: mplar data ([DEMO_DATA])] CTION]\n\n \n Output: [OUTPUT Placeholder] ction settings:</pre>
<pre>date choices; (3) date choices; (4) bel set in both c  TEMPLATES A INFERENCE (Z ero-shot instruct ictions, we follow designing the ir  Few-shot inss </pre> <pre></pre>	Binary Choice DExact Set: wh ontent and size PPLIED FOR A ERO-SHOT INI- tion induction w analogous id struction induc truction induc data: [Exe sol: [INSTRU acry INPUT] struction induc	<pre>e: task-solver LLM needs to choose one correct option out of two hether the predicted set of items (e.g., animals) exactly matches e, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generated eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024 ction templates. tion settings: mplar data ([DEMO_DATA])] CTION]\n\n \n Output: [OUTPUT Placeholder] ction settings: UCTION]\n\n \n Output: [OUTPUT Placeholder]</pre>
andidate choices; (3) andidate choices; (4 he label set in both c A.2 TEMPLATES A INFERENCE (Z For zero-shot instruct hen designing the ir • Few-shot ins <pre></pre>	Binary Choice DExact Set: wh ontent and size PPLIED FOR A ERO-SHOT INI- tion induction w analogous id struction induc truction induc data: [Exe s>. on: [INSTRU aery INPUT]	e: task-solver LLM needs to choose one correct option out of the hether the predicted set of items (e.g., animals) exactly matche, regardless of the item order. AIO INSTRUCTION GENERATION AND BLACK-BOX LLM DUCTION AND FEW-SHOT INDUCTION) and few-shot instruction induction, after obtaining generate eas as in Zhou et al. (2022); Chen et al. (2024); Lin et al. (202 ction templates. tion settings: mplar data ([DEMO_DATA])] CTION]\n\n \n Output: [OUTPUT Placeholder]

### **B** COMPLEMENTARY EXPERIMENTS

B63 Due to strict page limit for the main body, we choose to include complementary experiments here in this section. As an outline, we have (1) few-shot experiments on our 15 tasks located in Subsec. B.3;

(2) Chain-of-Thought (CoT) experiments in Subsec. B.4; (3) experiments with different combinations of white-box and black-box LLMs in Subsec. B.5; (4) an ablation study on AIO components in Subsec. B.6; (5) Parameter study for  $\beta$  and *B* in Subsec. B.7; (6) Empirical results with additional kinds of white-box LLMs in in Subsec. B.8; (7) Effects of exemplar quantity in Subsec. B.9; (8) Transferability across black-box LLMs in Subsec. B.10; (9) Empirical comparison with an additional baseline EvoPrompt in Subsec. B.11; (10) Additional optimization trajectory results in Subsec. B.12.

871 872

**B.1** BASELINE DESCRIPTIONS.

873 Recall that we involve four baselines for comparison, including two kinds of methods that utilize 874 LLMs for instruction optimization. The first kind methods leverages black-box LLM for instruction 875 generation: (1) APE (Zhou et al., 2022) which generates instruction using another black-box LLM 876 with designated templates and instruction search mechanism; (2) ProTeGi (Pryzant et al., 2023) 877 applies a black-box LLM for instruction generation, and optimizes "Gradient Descent" on the 878 generated instruction by integrating with another *black-box* LLM. Meanwhile, we also include (ii) 879 Methods that utilize white-box LLM for instruction generation: (3) InstructZero (Chen et al., 2024) 880 generates instructions using a *white-box* LLM, while controlling the generation process by optimizing a prefix soft prompt, based on a kernel-based Bayesian Optimization approach; (4) INSTINCT (Lin et al., 2024) adopts an analogous approach as InstructZero, but alternatively applies a neural bandit 882 model in replace of the kernel-based Bayesian Optimization for soft prompt selection. 883

#### 884 885

886

#### **B.2** IMPLEMENTATION DETAILS

For our zero-shot and few-shot instruction induction experiments, we consider each task is associated with 20 task exemplars denoted as  $\mathcal{D}_{\text{Train}}$  as well as 100 validation samples  $\mathcal{D}_{\text{Valid}}$ , which will remain 889 the same for AIO and all other baseline methods. For AIO, when choosing our threshold parameter 890  $\beta$ , we initially set  $\beta$  as an infinitely large value to enable all collected black-box LLM feedback 891 to be reused. Afterwards, we will experiment with  $\beta = \mathcal{O}(\epsilon \sqrt{d})$  and choose the constant for  $\mathcal{O}(\cdot)$ 892 notation with grid search  $\{1, 10, 100\}$ . We perform the fine-tuning process for T = 10 rounds, as 893 well as set the exploration parameter  $\nu = 0.1$  for all experiments For the perturbation magnitude parameter  $\epsilon$ , we choose its value with grid search from  $\{10^{-3}, 10^{-4}, 10^{-5}\}$ . In each optimization 894 round, we will draw K = 10000 arms from  $\mathcal{N}(0, \mathbf{I})$  where we choose B = 3 arms for fine-tuning 895 white-box LLM as well as update the TS model parameters. Regarding our TS model, after applying 896 JL-Lemma and random Gaussian matrix projection (Matoušek, 2008; Larsen & Nelson, 2017) for 897 dimension reduction, we will have the reduced dimension of TS model to be approximately  $d' \approx 10^4$ , which leads to  $\sim 0.4$  seconds for selecting chosen arms  $\tilde{Z}_t$  and  $\sim 3$  seconds for TS model parameters 899 update in each round t. For "AIO + LoRA", we set its "intrinsic rank" of low-rank approximation to 900 8. As we mentioned in the main body, we apply Llama-3-8B-Instruct (Dubey et al., 2024) 901 as our tunable white-box LLM  $\mathcal{F}_W(\cdot; \Theta_W)$ , and adopt Claude-3-Sonnet (Anthropic, 2024) as 902 our black-box LLM  $\mathcal{F}_B(\cdot)$ . All experiments are performed on a server with Intel Xeon CPU and 903 NVIDIA V100 GPUs.

- 904
- 905

907

#### 906 B.3 COMPLEMENTARY EXPERIMENTS WITH FEW-SHOT AIO

In this subsection, we include experiment results that examine AIO performance under few-shot
settings, where training samples (or exemplars) will be provided to black-box LLM for reference.
In this case, generated instructions will need to provide high-level reference and guidance, to assist
the answer generation of task-solver black-box LLM in observation of task exemplars. The template
applied for experiment is also shown in Subsec. A.2.

The experiment results are shown in Table 3. Here, we see that under few-shot settings, there exist
performance improvements for most baselines compared with zero-shot settings, due to the help of
additionally available task exemplars. By leveraging the sufficient representation power of fine-tuned
white-box LLM, AIO can still generally maintain the best performance compared with baseline
methods. Similar to our zero-shot experiment settings, we average performance of AIO and its PEFT
variant "AIO + LoRA" as a unity to obtain the ranking results.

918		Plaak box LLM		White box I I M		White how I I M w/ FT (Ours)	
919		Black-	DOX LLM	white-bo		white-box	LLM W/ F1 (Ours)
920	Tasks $\setminus$ Methods	APE	ProTeGi	InstructZero	INSTINCT	AIO	AIO + LoRA
921	antonyms	0.901	0.889	0.894	0.905	0.912	0.895
922	sentiment	0.932	0.944	0.940	0.933	0.950	0.946
923	larger_animal	0.939	0.915	0.922	0.874	0.961	0.957
924	taxonomy_animal	0.708	0.972	0.835	0.869	0.967	0.976
925	object_counting	0.511	0.583	0.520	0.541	0.555	0.473
926	navigate	0.734	0.724	0.701	0.757	0.776	0.772
927	winowhy	0.563	0.674	0.673	0.682	0.628	0.621
928	implicatures	0.846	0.826	0.859	0.847	0.836	0.867
929	logical_fallacy	0.850	0.892	0.877	0.881	0.880	0.885
930	hyperbaton	0.556	0.595	0.580	0.641	0.634	0.660
931	epistemic_reasoning	0.712	0.802	0.622	0.765	0.774	0.884
020	movie_recommendation	0.930	0.948	0.955	0.960	0.979	0.963
932	timedial	0.760	0.820	0.748	0.784	0.779	0.832
933	presuppositions_as_nli	0.557	0.564	0.591	0.598	0.619	0.594
934	question_selection	0.879	0.882	0.781	0.822	0.916	0.887
936	Average Rank	4.20	2.73	3.67	2.67		1.73

Table 3: Few-shot Instruction Induction Results. For each task (row), **bold** number refers to the best result, while underlined number refers to the second-best one. Similar to our zero-shot instruction induction experiments in Table 1, we average numerical results of AIO and its LoRA variant for each task, and treat these two methods as a unified baseline for ranking comparisons.

#### B.4 CHAIN-OF-THOUGHT (COT) RESULTS

We also include additional Chain-of-Thought (CoT) experiments on three data sets including GSM8K (Cobbe et al., 2021), AQUA (Garcia et al., 2020), and SVAMP (Patel et al., 2021), where results are shown in Table 4. For comparison, we include a baseline instruction "Let's think carefully step by step", which is commonly applied for solving CoT tasks, following the settings from Chen et al. (2024); Lin et al. (2024).

Data set	Method	Instruction	Accuracy Result
CEMOR COT		Let's think carefully step by step	0.724
USIMOR	AIO	Use your math skills and logic to break down the problem into manageable parts.	0.862
	СоТ	Let's think carefully step by step	0.317
AQUA	AIO	Think critically and break down the problem into smaller parts to solve it.	0.410
SVAMP	СоТ	Let's think carefully step by step	0.766
5 VAIVII	AIO	Let us think critically and break it down!	0.898

Table 4: Chain-of-Thought (CoT) results.

With results in Table 4, we see that AIO can significantly improve black-box induction performance under CoT reasoning settings compared with the task-agnostic instruction "Let's think carefully step by step", where AIO's performance improvements can be credited to the utilization of task-aware instructions. For instance, since GSM8K is a math reasoning task, AIO choose to introduce additional background information by asking the task-solver black-box LLM to "use your math skills" and decompose the target math problem into "manageable parts". This can help the black-box LLM determine which part of or what kinds of learned knowledge should be applied for problem solving, with higher levels of clarity than task-agnostic instructions.

#### 972 B.5 DIFFERENT COMBINATIONS OF WHITE-BOX AND BLACK-BOX LLMS

974Recall that for our previous experiments, we have applied Llama-3-8B-Instruct (Dubey et al.,<br/>2024) as our tunable white-box LLM  $\mathcal{F}_W(\cdot; \Theta_W)$ , and adopt Claude-3-Sonnet (Anthropic,<br/>2024) as our black-box LLM  $\mathcal{F}_B(\cdot)$ . Here, for two tasks "navigate" and "larger animal", we in-<br/>clude experiments with one more recent white-box LLM Llama-3.1-8B-Instruct, as well<br/>as a relatively light-weight black-box LLM Claude-3-Haiku for comparisons. Meanwhile,<br/>we also include experiments by substituting our black-box LLM with a powerful white-box LLM<br/>Llama-3-70B-Instruct for comparisons.

Task	White-box LLM	Black-box LLM	Accuracy Result
navigate	Llama-3-8B-Instruct	Claude-3-Sonnet	0.644
	Llama-3.1-8B-Instruct	Claude-3-Sonnet	0.689
	Llama-3-8B-Instruct	Claude-3-Haiku	0.612
	Llama-3.1-8B-Instruct	Claude-3-Haiku	0.643
Task	White-box LLM	Black-box LLM	Accuracy Result
larger anima	Llama-3-8B-Instruct	Claude-3-Sonnet	0.912
	Llama-3.1-8B-Instruct	Claude-3-Sonnet	0.927
larger_amma	Llama-3-8B-Instruct	Claude-3-Haiku	0.935
	Llama-3.1-8B-Instruct	Claude-3-Haiku	0.887

#### Table 5: Different combinations of white-box LLMs vs black-box LLMs.

Task	White-box LLM	Task-solver LLM	Best Instruction	Accuracy Result
navigate	Llama-3-8B-Instruct	Llama-3-70B-Instruct	If the instructions are able to return to the starting position after following all the instructions, then the output is True. Otherwise, the output is False.	0.567
larger_animal	Llama-3-8B-Instruct	Llama-3-70B-Instruct	First identify the animals in the input. Then, sort the animals in descending order based on their average adult body mass. If there are multiple animals with the same average adult body mass, sort them in alphabetical order. Finally, return the first animal in the sorted list as the output.	0.872

Table 6: Applying a white-box LLM (Llama-3-70B-Instruct) as problem-solving LLM.

Results are shown in Tables 5 and 6. Here, we see that using a more recent and capable white-box LLM can generally lead to slightly better performance. However, during our experiments, we also notice that fine-tuning Llama-3.1-8B-Instruct can be slightly more time consuming then tuning Llama-3-8B-Instruct. On the other hand, during our experiments, using a more lightweight black-box LLM can significantly accelerate the inference speed in terms of answer generation with relatively less API token costs. It can still achieve relatively good performance on "larger animal" task with a slightly inferior performance on "navigate" task. We also notice that the large white-box LLM Llama-3-70B-Instruct tends to perform slightly inferior compared with Claude family black-box LLMs, when using AIO as the instruction optimizer. 

# 1012 B.6 ABLATION STUDY ON AIO COMPONENTS

Recall that AIO has two main components: white-box back-propagation and TS-aided ZO gradient approximation, to derive the gradients for white-box LLM and black-box LLM respectively, based on the decomposed gradient flow in Eq. 4. Here, we include an ablation study for these two components for gradient derivation: (1) the first baseline is "AIO w/ MeZO" which directly use MeZO (Malladi et al., 2023) for approximating white-box LLM parameter gradients instead of our gradient decomposition formulation (Remark 1); (2) the second baseline is "AIO w/o TS Scheduling" where we do not apply Thompson Sampling for selecting perturbation directions and use completely random perturbation vectors  $\boldsymbol{z} \sim \mathcal{N}(0, \boldsymbol{I})$  for white-box gradient approximation in Eq. 5. 

Experiment results are shown in Table 7. We can see that our proposed AIO with TS-aided ZO
Gradient Approximation can still maintain superior performance compared with the other two
baselines with substituted modules. This helps to reinforce our claim that our proposed gradient
flow decomposition approach (Eq. 4) as well as the ZO black-box LLM gradient approximation
method guided by Thompson Sampling are necessary for AIO to achieve optimal performance.

Methods	antonyms	sentiment	l_animal	t_animal	navigate	implicatures	logical_fallacy	e_reasoning
AIO	0.901	0.949	0.912	0.983	0.644	0.811	0.868	0.766
AIO w/ MeZO AIO w/o TS Scheduli	ng 0.852 0.870	0.930 0.929	0.847 0.760	0.279 0.957	<b>0.654</b> 0.604	0.675 0.787	0.812 0.832	0.748 0.742

Table 7: Ablation study on AIO with two variants: (1) "AIO w/ MeZO" directly applies MeZO for fine-tuning white-box LLM  $\mathcal{F}_W$  instead of using our proposed gradient flow decomposition with TS-aided ZO gradient approximation; (2) "AIO w/o TS Scheduling" refers to a variant where perturbation directions z are sampling randomly from  $\mathcal{N}(0, I)$  instead of chosen by our TS model.

In particular, the supposed linear optimization landscape enables us to utilize a linear Thompson
 Sampling model for ZO perturbation direction selection, which is effective and computationally
 efficient for perturbation direction.

**1040 B.7 PARAMETER STUDY FOR**  $\beta$  **AND** *B* 

We include additional study for parameters  $\beta$  and B of AIO. Here, additional experiment results for zero-shot instruction induction, on the "Larger animal" and "Navigate" tasks with the LoRA module, are presented in the two tables below.

	Parameter B									
Task $\setminus B$ value	1	2	3	4						
Larger Animal	0.915	0.931	0.950	0.941						
Navigate	0.610	0.632	0.627	0.664						

Table 8: Experiment results with different *B* values.

Parameter $\beta$					
Task $\setminus \beta$ value	1	10	100	$\infty$	
Larger Animal	0.922	0.915	0.932	0.950	
Navigate	0.611	0.634	0.676	0.627	

1055 1056 1057

1035

1039

1041

1057

Table 9: Experiment results with different threshold  $\beta$  values.

Results are shown in Tables 8 and 9. For parameter *B*, we observe that setting B = 3 can achieve promising performance, which can help balance computational (and token) costs with performance. For parameter  $\beta$ , it is recommended to start the tuning process with a large  $\beta$ , as suggested in our Appendix B.2, to effectively leverage past received records of the optimization landscape. On the other hand, a sufficiently small threshold  $\beta$  will cause the method to degenerate into "AIO w/o TS Scheduling", as in our ablation study (Subsec. B.6). In this case, the TS model will be excluded from selecting ZO approximation directions, leading to relatively inferior performance.

1065 1066

#### B.8 COMBINATIONS OF DIFFERENT WHITE-BOX LLMS

1067 We include additional experiments with other types of white-box LLMs can benefit the audience. We 1068 conduct experiments with the LoRA module on two additional types of white-box LLMs: Mistral-1069 7B-Instruct-v0.2 and Qwen2.5-7B-Instruct. With black-box LLM being Claude-3-Sonnet, we have 1070 zero-shot instruction induction results shown in Table 10. We see that our proposed AIO framework 1071 achieves promising performance with other types of white-box LLMs other than the Llama family. 1072 Meanwhile, we also would like to mention that Llama 3 (Llama-3-8B-Instruct) generally retains a slight advantage over the other two white-box LLMs. One possible reason is that Llama 3 consists of 1074 8 billion parameters, slightly more than the other two 7-billion-parameter models. This can provide 1075 an advantage in terms of the representation power to some extent.

1076

1078

1077 B.9 DIFFERENT NUMBERS OF EXEMPLARS

As mentioned in our experimental settings (Appendix Subsec. B.2), we use 20 training exemplars, a reasonably small amount of training data, as the reference for the white-box LLM to generate

Method \ Task larger animal navigate sentiment movie recommendation Llama + Claude 3 0.950 0.627 0.947 0.883 Mistral + Claude 3 0.890 0.634 0.907 0.874 Qwen + Claude 3 0.919 0.682 0.922 0.835

Table 10: Experiment results with different kinds of white-box LLMs.

and optimize instructions. We also include additional zero-shot instruction induction experiments 1087 with varying exemplar quantities and the LoRA module. This is to investigate how the performance 1088 changes when altering the number of exemplars  $|\mathcal{D}_{\text{Train}}|$  for AIO, in terms of instruction generation 1089 and optimization. 1090

Different numbers of exemplars $ \mathcal{D}_{train} $						
Task \ $ \mathcal{D}_{train} $ 5102030						
Larger Animal	0.868	0.921	0.950	0.947		
Navigate	0.622	0.634	0.627	0.681		

Table 11: Experiment results with different numbers of exemplars  $|\mathcal{D}_{\text{train}}|$ .

From the results in Table 11, we observe that the "larger animal" task can be more sensitive to 1098 the number of exemplars  $|\mathcal{D}_{\text{Train}}|$ . However, providing as few as  $|\mathcal{D}_{\text{Train}}| = 10$  query-label pairs as 1099 exemplars, which is a modest quantity, will allow AIO to achieve relatively promising performance. 1100 Meanwhile, we see that for the "navigate" task, a small number of  $|\mathcal{D}_{\text{Train}}| = 5$  exemplars can lead to 1101 promising results, which shows that it is more stable under the sparse data settings. Therefore, when 1102 fine-tuning AIO from scratch, the overall performance of AIO can possibly be influenced by data 1103 scarcity, with its impact varying based on the specific application scenarios of practitioners. 1104

1105 **B**.10 TRANSFERABILITY OF FINE-TUNED INSTRUCTION OPTIMIZERS ACROSS DIFFERENT 1106 BLACK-BOX TASK-SOLVING LLMS 1107

1108 In this subsection, we investigate if the optimized instruction can generalize to different black-box 1109 LLMs. We transfer our optimized white-box LLM with the LoRA module, which is fine-tuned under 1110 the settings of Llama 3 + Claude 3 Sonnet, to other black-box task-solving LLMs. They include Claude 3.5 Sonnet and two OpenAI black-box LLMs (GPT-3.5-Turbo and GPT-40). 1111

ŧ.	-1	-1	0
L			~

1080

1081

1082

1083 1084

1086

1093 1094 1095

1113	Method \ Task	epistemic	logical	hyperbaton	movie-recommendation
111/	Claude 3 Sonnet	0.719	0.836	0.527	0.883
1114	Claude 3.5 Sonnet	0.844	0.886	0.562	0.924
1115	GPT-3.5-Turbo	0.737	0.811	0.556	0.825
1116	GPT-40	0.768	0.854	0.540	0.910

111 1117 1118

Table 12: Experiment results with different kinds of black-box LLMs.

1119 The accuracy results in terms of zero-shot instruction induction are shown in Table 12. Here, we can 1120 observe that the optimized instruction-generating white-box LLM can maintain strong performance 1121 when being applied to other black-box LLMs. Meanwhile, we also notice that the latest language 1122 models (Claude 3.5 and GPT-40) can generally outperform the older ones (Claude 3 and GPT-3.5-1123 Turbo), due to their stronger reasoning capabilities.

B.11 ADDITIONAL BASELINE: EVOPROMPT 1126

1127 In this subsection, we included an additional baseline EvoPrompt (Guo et al., 2024), which alterna-1128 tively utilizes evolutionary algorithms to refine LLM-generated instructions in an in-context learning 1129 manner. Different from our original problem settings, where instructions are generated from scratch, 1130 EvoPrompt requires initial instructions. In this case, we follow the settings in the original paper (Guo et al., 2024) and the official source code of EvoPrompt, by using instructions generated by APE 1131 (Zhou et al., 2022) as the initial instructions. We compare AIO with EvoPrompt, in terms of zero-shot 1132 instruction induction on four BigBench (bench authors, 2023) tasks, as BigBench is also utilized in 1133 (Guo et al., 2024).

<sup>1124</sup> 1125

1134	Method \ Task	hyperbaton	navigate	movie-recommendation	object-counting
1135	AIO	0.538	0.644	0.902	0.543
1136	EvoPrompt	0.526	0.627	0.895	0.466

1137 Table 13: Experiment results in comparison with EvoPrompt (Guo et al., 2024). 1138 1139 1140 Based on the results in Table 13, we see that our proposed AIO can manage to achieve generally 1141 better performance compared with EvoPrompt. Meanwhile, comparing with the two black-box LLM 1142 baselines in our initial submission: APE (Zhou et al., 2022) and ProTeGi (Pryzant et al., 2023), 1143 EvoPrompt can achieve relatively better performance than them, by utilizing task prior knowledge 1144 provided by initial instructions as the starting point for optimization. 1145 1146 **B.12** ADDITIONAL INSTRUCTION OPTIMIZATION TRAJECTORY RESULTS 1147 1148 To provide additional insights for our instruction optimization process, apart from our instruction 1149 optimization trajectory results from Subsec. 5.3, we include complementary results on extra tasks, in 1150 terms of how generated instructions evolve over the fine-tuning process, as well as how their accuracy 1151 results change with the instruction optimization. 1152 1153 B.12.1 "Sentiment" Task 1154 1155 • Accuracy: 46% 1156 IF the input contains any of the following words: "fun", "entertainingly", "splen-1157 did", "exquisite", "brilliantly", "impeccable", "likeable", "engrossing", "psycho-1158 logically", "slyly" THEN output is positive., 1159 ELSE IF the input contains any of the following words: "worst", "dislikable", 1160 "backlash", "mess", "pastiche", "warmed" THEN output is negative. 1161 1162 • Accuracy: 91% 1163 The output is positive if the input contains words with positive connotations, 1164 such as "fun", "entertainingly", "brilliantly", "exquisite", "likeable", "gentle", and 1165 "engrossing". 1166 The output is negative if the input contains words with negative connotations, 1167 such as "worst", "low", "pretentious", "dislikable", "mess", "warmed over", and 1168 "backlash". 1169 1170 B.12.2 "NAVIGATE" TASK 1171 1172 • Accuracy: 33% 1173 To determine if a robot will return to the starting point after a series of movements, 1174 follow these steps: 1175 1. Count the total number of steps taken by the robot. 1176 2. If the robot turns around at any point, it will change direction. So, count the 1177 number of turns and subtract 2 from the total number of steps for each turn. 1178 3. If the robot always faces forward, it will not change direction. So, do not count 1179 the turns. 1180 4. If the total number of steps minus the adjustments for turns is a multiple of the 1181 initial number of steps, the robot will return to the starting point. Example: If the 1182 robot takes 10 steps forward and then 5 steps backward, the total number of steps 1183 is 15. If the robot turns around at some point, adjust the total number of steps by subtracting 2 for each turn. If the robot always faces forward, do not adjust 1184 the total number of steps. If the total number of steps minus the adjustments is a 1185

• Accuracy: 71%

1186 1187 multiple of the initial number of steps, the robot will return to the starting point.

	To determine if a robot will return to its initial position after a series of instruc
	tions, follow these steps:
	1. Initialize the robot's position and direction to the starting point.
	2. Iterate through the instructions: - If the instruction is "Turn around", update th
	robot's direction to the opposite direction If the instruction is "Turn right" o
	"Turn left", update the robot's direction accordingly If the instruction is "Tak
	x steps", update the robot's position by moving x steps in the current direction.
	3. After iterating through all instructions, check if the robot's final position is th
	same as the initial position.
	4. If the final position is the same as the initial position, the output is Irue
	otherwise, the output is Faise.
3.12.3	"Object Counting" Task
	$\Delta_{\rm ccuracy}$ , 27%
-	
	To infer the output for a given input, count the number of items in the input that
	are not fruits or vegetables. Then, return the count. For example, if the input
	is "I have a violin, a piano, an accordion, a clarinet, an apple, a trumpet, and
	banana.", you would count the number of non-food items, which are 7 (violi
	piano, accordion, clarinet, trumpet), and return 7. If the input is "I have a yam,
	carrot, a cabbage, two garlics, a potato, a stalk of celery, a head of broccoli, and
	cauliflower.", you would count the number of non-food items, which are 0, an
	return 0.
•	Accuracy: 52%
	To infer the output for a given input, follow these stores
	1. Identify the types of items in the input. For example, if the input contain
	"violin" "niano" "accordion" "clarinat" "apple" "trumpet" and "banana" th
	types of items are musical instruments and fruits
	2 Count the number of items of each type. For example, there is 1 violin
	piano 1 accordion 1 clarinet 1 apple 1 trumpet and 1 banana
	3 Add the counts of all types of items together to get the total number of item
	4 The output is the total number of items
	For example, if the input is "I have a violin, a piano, an accordion, a clarinet, a
	apple, a trumpet, and a banana.", the output would be 7.
3.12.4	"TAXONOMY ANIMAL" TASK
•	Accuracy: 0%
	Sort the input list in descending order based on the length of the words. The
	select the words that are animals. Finally, sort the remaining words in descending
	order based on the length of the words.
•	Accuracy: 95%
	1 First identify all the words that are not animals. These words are usual
	1. This, identify an use words that are not animals. These words are usually nouns that are not animals, such as food items, colors, chiests, atc.
	2 Sort the remaining animal words in alphabetical order
	2. Soft the remaining annual words in alphabetical order, excluding the
	non-animal words
	non minim words.
3.12.5	"Larger Animal" Task

• Accuracy: 50%

1010		
1242		1. Take the first animal from the input.
1243		2. If the first animal is a predator, compare it with the second animal.
1244		3. If the first animal is a predator and the second animal is not a predator, then
1245		the output is the second animal.
1246		4. If the first animal is not a predator, compare it with the second animal.
1247		5. If the first animal is not a predator and the second animal is a predator, then
1248		the output is the first animal.
1249		7. If the first animal is larger than the second animal, then the output is the first
1250		animal
1251		8 If the first animal is smaller than the second animal, then the output is the
1252		second animal.
1253		9. If the first animal and the second animal are of the same size, then the output
1254		is the first animal.
1255		
1256	•	Accuracy: 92%
1257		If the input consists of two animals, compare their sizes. If the first animal is
1258		larger, output the first animal. If the first animal is smaller, output the second
1259		animal. If both animals have the same size, output the first animal.
1260		
1261	B.12.6	"Implicatures" Task
1202	•	Accuracy: 62%
1203		
1265		Given two speakers' dialogue, analyze the context and sentiment of the speakers'
1266		'ves') or not (output: 'no'). To do this, consider the following factors:
1267		1 **Gratitude tone**. If the second speaker's statement conveys gratitude
1268		appreciation or a positive sentiment towards the first speaker infer 'ves'
1260		2. **Question-answer format**: If the first speaker's statement is a question and
1205		the second speaker's response is an affirmative answer or a statement that implies
1071		agreement, infer 'yes'.
1070		3. **Contrast**: If the first speaker's statement is a negative or neutral statement,
1272		and the second speaker's response is a contrasting positive statement, infer 'yes'.
107/		4. **Contextual understanding**: Consider the context of the conversation and
1274		the speakers' intentions. If the second speaker's statement seems to be responding
1275		to the first speaker's action or suggestion in a positive manner, infer 'yes'.
1077		When applying these factors, consider the nuances of language, idioms, and
1079		ingurative expressions. If the analysis is antiguous of unclear, default to no.
1270	•	Accuracy: 84%
1275		To infor the output for a given input follow these store:
1200		1. Identify the sneaker's tone and intention in the input. Are they expressing
1282		gratitude surprise or skepticism?
1283		2. Look for words or phrases that convey a positive or negative sentiment. If the
1200		sentiment is positive, the output is likely to be 'yes'. If the sentiment is negative,
1285		the output is likely to be 'no'.
1286		3. Check for words or phrases that indicate a question or a request. If the input
1200		contains a question, the output is likely to be 'yes' if the speaker is seeking
1288		confirmation or agreement. If the input contains a request, the output is likely to
1280		be 'yes' if the speaker is seeking permission or approval.
1205		4. Consider the context of the conversation. Is the speaker discussing a specific
1201		event, situation, or decision? If so, the output may be influenced by the speaker's
1202		5. Use your judgment and common sense to make an advanted guess shout the
1202		output If you're still unsure try re-reading the input and re-evaluating the tope
120/		sentiment, and context.
1205		By following these steps, you can make an informed decision about the output
.200		for a given input.

<sup>196</sup> B	3.12.7	"Epistemic Reasoning" Task
.9 <i>7</i>	•	Accuracy: 41%
.50		
99		For each input check if the premise and hypothesis are about the same individual
00		or group of individuals. If they are, the output is entailment. If they are not, the
20		output is non-entailment.
12		L
лл 1л		
)4	•	Accuracy: 76%
16		
)7		Given a premise and a hypothesis, determine the entailment relationship between
18		them. For a premise, identify the main clause and the subordinate clause (if
9		present). For the hypothesis, identify the main clause and the subordinate clause
0		(if present). If the main clause in the premise is identical to the main clause in
1		the hypothesis, and the subordinate clause (if present) is also identical, then the
2		output is entailment. If the main clause in the premise is identical to the main
3		clause in the hypothesis, but the subordinate clause (if present) is different, then the subordinate clause (if present) is different, then
4		the main clause in the hypothesis then the output is non-entailment
15		the main clause in the hypothesis, then the output is non-entainment.
6		
7		
<sup>18</sup> B	3.12.8	"TIMEDIAL" TASK
9		
20	•	Accuracy: 41%
21		
22		1. Read the context carefully and identify the main idea or topic of the conversa-
23		tion.
24		2. Identify the speaker and the listener. 3. Determine the purpose of the conversation
25		4. Identify the key information or events mentioned in the conversation.
26		5. Analyze the conversation to determine the relationship between the speakers
27		and the context.
28		6. Based on the analysis, infer the output for the given input.
29		
30		
	•	Accuracy: 79%
22		
2/		To infer the output for a given input, follow these steps:
5		1. Identify the context and the entities mentioned in the input, such as people,
.5		praces, and times.
87		zeruested or discussed such as dates times quantities or specific events
38		3. Analyze the relationships between the entities and the context, and try to
9		identify the most relevant and important information.
0		4. If the output is a date or time, check if it is a specific date or time, or if it is a
11		range or duration.
2		5. If the output is a quantity, check if it is a specific number, or if it is a range or
13		interval.
4		6. Finally, use the interred output to generate a response that is relevant and
5		accurate to the input and the context.
16		
7		
18 n	2 1 2 0	"Ομεντιονι δει εστιονι" Τλαγ
D	5.12.9	QUESTION SELECTION TASK

• Accuracy: 39%

	To infer the output for a given input, follow these steps:
	1. Identify the context of the input, which typically consists of a passage of text.
	2. Determine the type of question being asked, which is usually indicated by the
	choice options provided.
	3. Look for specific keywords or phrases in the input that are related to the
	question being asked.
	4. Identify the relevant information in the input that answers the question, which
	may be a specific fact, statistic, or quote.
	5. Match the relevant information to the corresponding choice option in the
	question.
	o. Select the choice option that best answers the question based on the information
	provided in the input.
•	Accuracy: 64%
	To infer the output for a given input, follow these steps:
	1. Identify the context of the input, which is typically a passage of text.
	2. Determine the type of question being asked, which is usually indicated by the
	format of the choices provided.
	3. Scan the context to identify relevant information, such as key phrases, names,
	and dates.
	4. Match the context to the corresponding choices, considering the question type
	and the information gathered.
	5. Select the most likely output based on the context and the choices.
	For example, if the input is a passage about a football game, the output might be a question about the teams involved, the score, or the MVD. If the input is a
	be a question about the learns involved, the score, or the MVP. If the input is a passage about a historical event, the output might be a question about the data
	location or significance of the event
B 12 1	) "Movie Recommendation" Task
<b>D</b> .12.1	
•	Accuracy: 48%
	1. Identify the most frequent genre or theme in the given context.
	2. Compare the given context with the choices and find the one that best matches
	the identified genre or theme.
	3. Output the matching choice.
•	Accuracy: 88%
	The output is the choice that is most commonly accepted with the siven context
	I ne output is the choice that is most commonly associated with the given context,
	choice
	choice.
C St	JPPLEMENTARY TECHNICAL DETAILS
	PETAILS FOR AUTO-REGRESSIVE GENERATION AND PERTURBATION
C.1 E	
C St	Accuracy: 88% The output is the choice that is most commonly associated with the given contex based on the frequency of co-occurrences of movies in the input list with eac choice. JPPLEMENTARY TECHNICAL DETAILS DETAILS FOR AUTO-REGRESSIVE GENERATION AND PERTURBATION

a token sequence x. Beginning with an initial input context  $x_{<1}$  that can be empty or contain special tokens like start-of-sequence token. Then, the language model will sequentially generate each token in the output, by sampling each generated token  $x_i$  from the conditional distribution  $p_{\Theta}(x_i \mid x_{< i})$ . In particular, the probability distribution for the *i*-th token will go through a softmax function, after applying a language model header (with weight  $\Theta_{\text{header}}$  from language model parameters  $\Theta$ ) to map *i*-th token hidden representation to the vocabulary distribution, as

 $p_{\boldsymbol{\Theta}}(x_i \mid \boldsymbol{x}_{< i}) = \operatorname{softmax}(\boldsymbol{\Theta}_{\operatorname{header}} \cdot \boldsymbol{h}_i),$ 

1400 1401

where  $h_i$  represents the transformer-embedded hidden representation of *i*-th generated token. Each token  $x_i$  will be sampled from the vocabulary, based on  $p_{\Theta}(x_i | x_{< i})$ . The generation process will complete if special tokens (e.g., an EOS token) is encountered, or the maximum length is met. Following our gradient approximation formulation in Eq. 5, with random perturbation vector (gradient approximation vector)  $z \sim \mathcal{N}(0, I)$  and perturbation magnitude  $\epsilon > 0$ , we recall that the perturbation is imposed on *LLM-header output probabilities* (i.e., distribution over the vocabulary) of *each i*th generated token. This leads to the positively perturbed generation process  $p_{\Theta}(x_i \mid \boldsymbol{x}_{< i}) \leftarrow$ softmax( $\Theta_{\text{header}} \cdot \boldsymbol{h}_i + \epsilon \boldsymbol{z}$ ), as well as the negatively perturbed generation process  $p_{\Theta}(x_i \mid \boldsymbol{x}_{< i}) \leftarrow$ softmax( $\Theta_{\text{header}} \cdot \boldsymbol{h}_i - \epsilon \boldsymbol{z}$ ). Consequently, the perturbation  $\boldsymbol{z} \in \mathbb{R}^d$  will be of the same as the vocabulary dimension d.

- 1411
- 1412 1413

1414

#### C.2 DETAILS FOR EFFICIENT COVARIANCE MATRIX INVERSION UPDATE WITH SHERMAN-MORRISON FORMULA

Recall that in Remark 2, for updating covariance matrix inversion  $\Sigma_t^{-1}$  efficiently in each round t, 1415 we apply Sherman-Morrison Formula (Bartlett, 1951; Maponi, 2007) by updating matrix inversion 1416 incrementally. Here, suppose that current white-box LLM output  $\phi$  is close enough to the checkpoint 1417  $\phi_{\text{Check}}$ , which means that we can update currently possessed covariance matrix inversion  $\Sigma_{t-1}^{-1}$  to 1418 obtain  $\Sigma_t^{-1}$ . Then, we recall that the arm context covariance matrix in each round t is constructed by 1419  $\Sigma_t = I + \sum_{(\boldsymbol{z},r)\in\Omega_t} \boldsymbol{z} \cdot \boldsymbol{z}^{\intercal} = \Sigma_{t-1} + \sum_{\boldsymbol{z}\in\widetilde{\mathcal{Z}}_t} \boldsymbol{z} \cdot \boldsymbol{z}^{\intercal}$ , where  $\widetilde{\mathcal{Z}}_t$  refers to the collection of chosen 1420 1421 arms in round t. Since we have each  $zz^{\intercal}$  being a rank-one matrix for every  $z \in \hat{Z}_t$ , we can follow Sherman-Morrison Formula to perform one-step update 1422 1/103

1426

 $\left(\boldsymbol{\Sigma}_{t-1} + \boldsymbol{z}\boldsymbol{z}^{\intercal}\right)^{-1} = \boldsymbol{\Sigma}_{t-1}^{-1} - \left[\frac{\boldsymbol{\Sigma}_{t-1}^{-1}\boldsymbol{z}\boldsymbol{z}^{\intercal}\boldsymbol{\Sigma}_{t-1}^{-1}}{1 + \boldsymbol{z}^{\intercal}\boldsymbol{\Sigma}_{t-1}^{-1}\boldsymbol{z}}
ight].$ 

In this case, by iteratively repeating this process for  $|\widetilde{Z}_t| = B$  times (since we have  $B \ll K$ 1427 chosen arms in each round t, and B is a considerably small integer), we will have the updated 1428 covariance matrix inverse  $\Sigma_t^{-1}$ . Recall that each covariance matrix will have a shape of  $d \times d$ , where 1429 d is the dimensionality of perturbation vector z. We then have the overall computational costs as 1430 approximately  $\mathcal{O}(Bd^2)$  instead of the naive  $\mathcal{O}(d^3)$ , where we also intuitively have  $B \ll d$ . Moreover, 1431 with dimension reduction approach motivated by JL-Lemma (Remark 2), we can have the projected 1432 context dimension  $d' \ll d$ , which leads to computational complexity of  $\mathcal{O}(B \cdot (d')^2)$ , instead of the 1433 naive  $\mathcal{O}((d')^3)$  with the direct matrix inversion. 1434

1435

1437

### 1436 D COMPLEMENTARY DISCUSSIONS

#### 1438 D.1 MOTIVATION OF THOMPSON SAMPLING AND BIAS ASSOCIATED WITH THIS DESIGN

1439 1440 D.1

D.1.1 Additional discussion on Motivations

As mentioned in our manuscript (paragraph below Remark 1), conventional ZO gradient approximation methods have a potential drawback: the perturbation vectors *z* are randomly sampled. Consequently, their gradient approximation directions in the optimization landscape are random, which can result in inefficient gradient estimation process. To overcome this challenge, numerous ZO methods have been proposed to incorporate guided or chosen directions for gradient optimization, enabling more efficient estimation of the target gradient (e.g., Cai et al. (2022); Qiu & Tong (2024)).

In this work, we propose reusing collected feedback by framing the ZO-based fine-tuning process as an 1448 online sequential decision-making problem, and applying Contextual Bandit techniques to effectively 1449 identify beneficial perturbation directions worth exploring. With the assumed linear optimization 1450 landscape as in existing works (e.g., Spall (1992); Malladi et al. (2023)), we apply linear Thompson 1451 Sampling to leverage previously collected information, which includes arms (previous gradient 1452 approximation directions) and corresponding rewards (benefits of going along these directions), to 1453 achieve a more efficient gradient estimation (Subsec. 4.1.2). Within the linear optimization landscape, 1454 our TS model aids in selecting gradient directions for white-box LLM fine-tuning, by properly 1455 reducing the possibility of choosing low-value directions (e.g., directions that are orthogonal to the 1456 true gradient, which can provide limited information for gradient approximation). This approach helps reduce the validation loss through instruction generation and ZO approximation direction 1457 selection, aligned with our formulation of the arm reward (Eq. 8).

1458 Ablation study on AIO components. We also would like to mention that the effectiveness of our TS 1459 modeling is supported by our ablation study (Appendix Subsec. B.6), where we compare our AIO 1460 framework with two alternatives: (1) "AIO w/ MeZO", where MeZO (Malladi et al., 2023) is used 1461 directly for approximating white-box LLM parameter gradients instead of our gradient decomposition 1462 formulation (Remark 1). (2) "AIO w/o TS Scheduling", where Thompson Sampling is not applied for selecting perturbation directions, and completely random perturbation vectors  $z \sim \mathcal{N}(0, I)$  are 1463 used for zeroth-order gradient approximation as in Eq. 5. We observe that our proposed AIO with 1464 TS-aided ZO gradient approximation generally achieves better performance, compared with the two 1465 baselines with substituted modules. 1466

1467

#### 1468 D.1.2 BIAS ASSOCIATED WITH TS 1469

1470 By balancing exploitation and exploration, TS will not introduce considerable bias. Under the 1471 settings of gradient approximation, our proposed TS-based zeroth-order gradient approximation 1472 method itself does not inherently introduce considerable bias in terms of gradient estimation, as Thompson Sampling techniques can naturally tackle the exploitation-exploration dilemma (Agrawal 1473 & Goyal, 2013; Zhang et al., 2021) by exploring various gradient directions instead of greedily 1474 exploiting only certain ones. Here, TS model will choose from sampled candidate approximation 1475 directions, instead of actually altering the direction vector value. Meanwhile, in the short term, 1476 it is possible that TS may introduce temporary bias if it focuses on some perturbation directions 1477 that can lead to high rewards (i.e., directions that can help reduce validation loss, Eq. 8). Due to 1478 the bias-variance trade-off (Hu et al., 2021b), the temporary bias can also be beneficial under our 1479 gradient approximation settings, and we will elaborate on this point in the next paragraph. On the 1480 other hand, with only a few gradient directions applied for LLM fine-tuning (Malladi et al., 2023) in 1481 each optimization round, directly applying conventional unbiased estimators with randomly sampled 1482 directions can possibly result in high variance in terms of gradient estimation (Cai et al., 2022), which 1483 also reflects the *bias-variance trade-off* in terms of zeroth-order gradient approximation.

1484 We introduce the TS to balance the "bias" and "variance" trade-off for ZO gradient approximation, 1485 under query-limited scenarios. Here, we would like to mention that state-of-the-art ZO gradient 1486 approximation methods with a guided approximation process (e.g., Qiu & Tong (2024)) can also 1487 introduce temporary bias in terms of gradient direction selection. However, they have been shown 1488 highly effective and efficient, particularly under sample-efficient settings, due to their ability of 1489 involving informative gradient approximation directions to balance "bias" and "variance" for gradient approximation. As a result, due to the bias-variance trade-off, it can lead to the case that no algorithm 1490 can serve as a universal solution to various application scenarios of zeroth-order approximation. 1491 Therefore, practitioners need to tailor solutions to their specific application scenarios. In our case, 1492 by balancing high-reward directions (exploitation) and sampling TS parameters from the posterior 1493 (exploration), our TS-aided approximation can help stabilize the gradient estimates, making it 1494 advantageous in query-limited scenarios like ours. In particular, by leveraging a highly efficient linear 1495 TS model, the arm selection process ( $\sim 0.4$  seconds per round) will be fast, ensuring the efficiency in 1496 terms of gradient direction selection. 1497

Ablation study on TS-aided zeroth-order gradient approximation. The effectiveness of our modeling is also validated by our ablation study (Appendix Subsec. B.6), where we include a variant, "AIO w/o TS Scheduling," for comparison. Instead of using TS, "AIO w/o TS Scheduling" applies randomly chosen perturbation vectors  $z \sim \mathcal{N}(0, I)$  for zeroth-order gradient approximation, as described in Eq. 5. In contrast, our AIO with TS-aided ZO gradient approximation achieves better performance by strategically selecting informative gradient approximation directions, based on collected information and knowledge from the optimization landscape.

- 1504
- 1505 1506 1507

### D.2 REASONING OF OUR CHOICE OF DATA SETS AND BASELINES

Under our problem settings (Section 2), the black-box LLM is considered as part of the learning objective rather than the learning model, and we have no control over its parameters while can only interact with it through the API. Therefore, based on our problem settings, which align with those of closely related works (Chen et al., 2024; Lin et al., 2024), the learning model will solely be the white-box LLM, while the black-box LLM will serve as part of the learning objective.

### D.2.1 CHOICE OF DATA SETS.

1513 1514

Different from conventional instruction optimization settings, in this work, we address the challenge of automatically optimizing instructions for a task-solving black-box LLM in terms of the given target task. The process requires only a few exemplars as training data and eliminates the need for human expert intervention. This is an emerging topic of automatic instruction generation that has been explored in several related works (e.g., Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024)). Unlike existing approaches, we propose utilizing a white-box LLM for instruction generation and optimization, coupled with LLM fine-tuning to effectively learn optimized instructions for highly complex modern black-box LLMs.

1522 As discussed in the Introduction and Related Works sections, the most relevant state-of-the-art works 1523 to this paper are Chen et al. (2024); Lin et al. (2024), where a white-box LLM is also used as an 1524 instruction optimizer to tailor instructions specifically for the downstream task-solving black-box 1525 LLM. Therefore, for the datasets, we follow the most closely related works (e.g., Zhou et al. (2022); 1526 Chen et al. (2024); Lin et al. (2024)) by utilizing instruction induction tasks from Honovich et al. 1527 (2022), reasoning tasks from BigBench (bench authors, 2023), and Chain-of-Thought (CoT) datasets (Cobbe et al., 2021; Garcia et al., 2020; Patel et al., 2021) to benchmark our proposed AIO against 1529 baselines. These datasets are both common and widely adopted in this line of research, particularly in works closely aligned with ours in terms of problem settings (e.g., Zhou et al. (2022); Chen et al. 1530 (2024); Lin et al. (2024)). 1531

1532 The instruction induction tasks from Honovich et al. (2022) and reasoning tasks from BigBench 1533 (bench authors, 2023) are used to evaluate the zero-shot instruction induction (Subsec. 5.1) and 1534 few-shot instruction induction (Subsec. B.3) quality of the optimized instructions. In particular, since 1535 no auxiliary task-relevant information is provided to the task-solving black-box LLM, the zero-shot instruction induction results on these datasets (Honovich et al., 2022; bench authors, 2023) are widely 1536 adopted to assess instruction quality by closely related works (e.g., Zhou et al. (2022); Chen et al. 1537 (2024); Lin et al. (2024); Fernando et al. (2023); Hu et al. (2024)). Additionally, the applied CoT 1538 datasets are also commonly used in the aforementioned related works (e.g., Zhou et al. (2022); Chen 1539 et al. (2024); Lin et al. (2024)), as they effectively test the instruction optimizers' ability to solve 1540 complex reasoning tasks, such as math problems. Therefore, our dataset selections are standard and 1541 widely adopted in closely related works, in order to demonstrate the effectiveness of AIO in terms of 1542 instructing the task-solving black-box LLM. 1543

- 1543 1544
- 1545
- 1545

#### D.2.2 CHOICE OF BASELINES.

1547 1548

In this paper, we primarily focus on the challenge of automatically optimizing instructions for a 1549 task-solving black-box LLM given the target task. This distinguishes our problem settings from those 1550 of conventional instruction optimization works. In this case, the most relevant works to ours are 1551 Chen et al. (2024); Lin et al. (2024), which propose to deal with an analogous problem. From an in-1552 context learning perspective instead, they apply a white-box LLM as an instruction generator to tailor 1553 instructions specifically for the downstream task-solving black-box LLM. Therefore, in our initial 1554 submission, we have included these two works (InstructZero (Chen et al., 2024), INSTINCT (Lin 1555 et al., 2024)) as our baselines to emphasize the advantages of LLM fine-tuning over their in-context 1556 learning approaches under our instruction optimization settings, given the superior representation and 1557 learning capabilities provided by LLM fine-tuning.

1558 Meanwhile, there is also a line of research (e.g., Zhou et al. (2022); Pryzant et al. (2023)) that 1559 employs a black-box LLM instead for instruction generation, resulting in a pipeline involving two 1560 black-box LLMs. In this setup, the instruction-generating black-box LLM perceives the feedback of 1561 the task-solver black-box LLM, and optimizes its instructing strategy from an in-context learning perspective. Since these methods are also commonly adopted by our closely related works (Chen et al., 2024; Lin et al., 2024), we have also included APE (Zhou et al., 2022) and ProTeGi (Pryzant 1563 et al., 2023) as our baselines in our initial submission. Therefore, our choice of baselines is standard 1564 in this line of research, including state-of-the-art baselines closely related to this paper, as well as 1565 common baselines adopted by other existing works in this research direction.

# <sup>1566</sup> E GENERALIZING AIO TO OTHER APPLICATION SCENARIOS

# 1568 E.1 INVOLVING DOMAIN EXPERTS IN THE OPTIMIZATION OF INSTRUCTIONS

1570 As shown in our pipeline illustration (the Generation Template), we only need a few exemplars from the target task as the reference to the white-box LLM for instruction generation and optimization. In 1571 this case, we do not need specific human-crafted task context for our AIO, which avoids the need for 1572 human expert intervention. This emerging topic in automatic instruction generation has been explored 1573 in several related works (e.g., Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024)). Unlike 1574 prior approaches, we propose utilizing a white-box LLM for instruction generation and optimization, 1575 paired with LLM fine-tuning to efficiently learn optimized instructions for modern, highly complex 1576 black-box LLMs. On the other hand, if domain experts (e.g., human engineers) or domain knowledge 1577 (e.g., textual task descriptions or narratives) are available, there are several ways of integrating them 1578 to our instruction optimization process.

1579

# 1580 E.1.1 TASK NARRATIVE OR DESCRIPTION

First, if task descriptions are provided to AIO as prior knowledge, we can incorporate this information into the instruction generation template (left-hand side of Figure 1). In this case, the input to the white-box LLM will include both the task exemplars and the textual task description. We conduct additional experiments in terms of zero-shot instruction induction on four BigBench tasks with the LoRA module, incorporating the one-sentence task descriptions provided by the BigBench data set into our generation template.

Method \ Task	epistemic reasoning	logical fallacy	hyperbaton	movie recommendation
AIO	0.719	0.836	0.527	0.883
AIO w/ Task Info	0.811	0.872	0.639	0.895

1590 1591 1592

Table 14: Experiment results with additional textual task description information.

Results are shown in Table 14. We see that involving additional textual task information into the instruction generation and optimization process can generally improve performance, particularly for logical reasoning tasks (e.g., epistemic reasoning, logical fallacy, hyperbaton), as task descriptions can help provide extra background information to prevent potential misinterpretations of exemplars.

#### 1598 E.1.2 INVOLVING HUMAN EXPERTS

When human experts are available, following the idea of "Human-in-the-loop" (Wu et al., 2022), we can involve these experts to help AIO in terms of instruction optimization and evaluation. For instance, in terms of instruction improvement for medical diagnostics, AIO can generate instructions such as "Please analyze patient symptoms to identify possible illnesses." Human evaluators could define metrics emphasizing clarity (e.g., suggesting tools like "blood test" or "MRI scan") or appropriateness (e.g., ensuring instructions do not encourage unsafe practices). This can enhance AIO's applicability in sensitive fields where clarity and accuracy are crucial. Compared to composing instructions from scratch that requires more efforts, this formulation can help reduce the workload of human experts.

Meanwhile, human experts can be involved into the instruction evaluation process. In this case, AIO can be requested to generate multiple candidate instructions after fine-tuning, and then human experts can proceed to examine which of the candidates is most appropriate from multiple dimensions, such as interpretability. This will also help to enhance AIO's applicability for real-world scenarios.

1611 1612 1613

#### E.2 GENERALIZING OPTIMIZED INSTRUCTION TO RELATED TASKS OR DOMAINS

Recall that this work focuses on the problem of automatically optimizing instructions for a task-solving black-box LLM given the target task. The instruction generation and optimization require only a few exemplars as training data and do not involve human expert intervention. This problem is an emerging topic in the field of automatic instruction generation, with several related works (e.g., Zhou et al. (2022); Chen et al. (2024); Lin et al. (2024)). Different from existing works, in this paper, we propose leveraging a white-box LLM for instruction generation and optimization, combined with LLM fine-tuning to effectively learn optimized instructions for modern black-box LLMs of extreme

1620 complexity. However, while this paper focuses on optimizing instructions for a given target task, our
 AIO framework can be potentially extended to generalize across multiple related tasks or domains.

**Transfer learning and domain adaptation.** First, we can leverage ideas from transfer learning and domain adaptation to address this issue. Here, we can consider two related domains (tasks): the source domain and the target domain. The source domain will refer to the domain that we train the original instruction optimizer (white-box LLM), typically having larger amount of training data (exemplars) than the target domain. Afterwards, to adapt the trained white-box LLM to the target domain, we can apply few-shot domain adaptation techniques (e.g., Chronopoulou et al. (2022)) to efficiently fine-tune the white-box LLM parameters for the target domain, even with a relatively smaller number of samples. In the meantime, in the context of transfer learning, we can possibly have additional theoretical insights. This is because the generalization loss on the target domain can be theoretically upper bounded by the empirical loss on the source domain, the learning power of the neural model, and the discrepancy between the source and target domains (Ben-David et al., 2010). 

Meta-learning. On the other hand, we can leverage ideas from meta-learning (Finn et al., 2017) to develop an instruction "meta-optimizer", which is capable of abstracting high-level information across multiple domains and generating high-quality meta-instructions (analogous to meta-prompts (Hou et al., 2022)). Using this approach, the "meta-optimizer" can be adapted to downstream tasks with only a few exemplars from the target task, without significant modifications of the AIO parameters. The performance of this approach is expected to depend on the neural model's learning capacity and the discrepancy among different tasks (Chen et al., 2021). With a white-box LLM as the learning model with fine-tuning, providing sufficient representation power, AIO can help acquire high-level instruction knowledge using meta-learning techniques.