# **ADO:** Automatic Data Optimization for Inputs in LLM Prompts

#### Anonymous ACL submission

#### Abstract

This study explores a novel approach to enhance the performance of Large Language Models (LLMs) through the optimization of input data within prompts. While previous research has primarily focused on refining instruction components and augmenting input data with in-context examples, our work investigates the potential benefits of optimizing the input data itself. We introduce a twopronged strategy for input data optimization: content engineering and structural reformulation. Content engineering involves imputing missing values, removing irrelevant attributes, and enriching profiles by generating additional information inferred from existing attributes. Subsequent to content engineering, structural reformulation is applied to optimize the presentation of the modified content to LLMs, given their sensitivity to input format. Our findings suggest that these optimizations can significantly improve the performance of LLMs in various tasks, offering a promising avenue for future research in prompt engineering. The source code is available at https:// anonymous.4open.science/r/ADO-6BC5/.

## 1 Introduction

003

011

012

014

018

037

041

Large Language Models (LLMs) (Achiam et al., 2023; Team et al., 2023; Touvron et al., 2023) have demonstrated exceptional proficiency across a wide array of tasks. They have been successfully implemented in various real-world applications, including personalized recommendations (Wu et al., 2024; Hua et al., 2023), healthcare (Yu et al., 2024b,a; Li et al., 2024a), financial decision-making (Li et al., 2023b; Wu et al., 2023), and advanced language reasoning (Huang and Chang, 2022; Fan et al., 2023; Sharan et al., 2023). In particular, LLM prompting has become a critical research area (Chen et al., 2023, 2024). This is because LLMs are highly sensitive to input content and format; even slight modifications, such

as changes in word order or indentation, can significantly influence their performance (Sclar et al., 2023; Fang et al., 2024). 042

043

044

047

048

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

082

When LLMs are employed for task inferencing, a user prompt (or query) typically comprises two primary components: a task-specific instruction and the input data to be processed according to that instruction. For example, when employing an LLM for Heart Disease classification (Baccouche et al., 2020), the task-specific instruction can be "analyze the following user's health profile to determine the likelihood of a heart attack", while the input data can include the individual's health profile, encompassing attributes such as age, medical history, and lifestyle habits. In the context of personalized recommendations, such as for beauty products (Geng et al., 2022), the instruction can be "generate beauty product recommendations based on the user's recent interaction history with other beauty products", with the input data consisting of the user's interaction history and a set of candidate beauty products to make recommendations from.

Various prompting methods have been proposed to enhance the inference performance of LLMs. For example, multiple studies have focused on crafting manual prompting strategies (Bsharat et al., 2023; Sahoo et al., 2024; Marvin et al., 2023), such as Chain-of-Thought (CoT) reasoning (Wei et al., 2022). Additionally, automated methods have been developed to search for optimal instructions tailored to specific tasks (Do et al., 2024; Li et al., 2024b). For instance, APE (Zhou et al., 2023) introduces an iterative Monte Carlo search to refine prompt instructions. Other works focus on providing in-context demonstrations (Dong et al., 2022), offering examples to guide the model's responses.

Most prior works on prompt engineering have focused on two aspects: (1) optimizing the instruction component of the prompt and (2) augmenting the input data with additional context, such as incontext exemplars, as illustrated on the "Traditional



Figure 1: Types of prompt engineering approaches. Given an inference task, such as solving a logical puzzle (as shown in the middle of the figure), prior works primarily focus on either optimizing instructions or augmenting the input data with similar examples, as depicted at the top of the figure. In contrast, we propose optimizing the input data to enhance its presentation to LLMs for more effective task inference, as illustrated at the bottom of the figure.

Approach" section of Figure 1. Nevertheless, the role of input data optimization in enhancing LLM performance remains underexplored.

To address this gap, we investigate whether optimizing the input data portion of the prompt can also enhance performance, as depicted on the "Proposed Solution" section of Figure 1. Towards this goal, we propose a new framework "Automatic Data Optimization (ADO)" as well as a new algorithm, "Diverse Prompt Search (DPS)". This framework can optimize input data through two key strategies: content engineering and structural reformulation. First, we apply content engineering to refine input data, such as imputing missing values based on domain knowledge and removing irrelevant attributes that may hinder decision-making. Second, we leverage structural reformulation to modify the format of input data, aiming to optimize data presentation to LLMs. Together, our proposed framework has demonstrated its effectiveness to complement conventional prompting strategies to enhance LLM inference performance.

### 2 ADO Framework

087

090

092

098

100

101

102

103

105

106

107

109

This section outlines the objectives of input data optimization and explains the mechanisms by which the ADO framework achieves these objectives.

## 2.1 Framework Objective

In this work, we conduct data optimization on
the input data part of the prompt prior to submit-

ting the prompt to a LLM for inference. Our data optimization objectives can be categorized into two aspects: content optimization and format optimization. Content optimization emphasizes enhancing the saliency of features within the data, ensuring that the most relevant and informative attributes are highlighted. Format optimization focuses on structuring the data in an optimal format, such as tables, XML, or other representations that facilitate efficient processing and interpretation. Let **D** represents the original input data. The overall data optimization process can be considered as a combination of both content and format optimizations, resulting in an optimized dataset **D**':

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

$$\mathbf{D}' = f_{format}(f_{content}(\mathbf{D})) = f(\mathbf{D}) \qquad (1)$$

where f is the composite optimization function. This comprehensive approach ensures that the data not only contains salient features but is also presented in a format that maximizes its utility for inference tasks.

**Content Optimization** has been a prominent area of research across various fields and modalities (Ahmad et al., 2018; Zhou and Aggarwal, 2004). For example, in tabular datasets, where each individual is represented by a set of attribute-value pairs, common content optimization procedures include feature extraction, missing value imputation, and attribute aggregation (Zheng and Casari, 2018). These techniques aim to enhance the quality of the



Figure 2: ADO Workflow. The Prompt-Generation LLM initially proposes task-specific instructions for optimizing input data, which the Data Optimization LLM executes on validation set samples, generating optimized inputs. These optimized samples are then processed by the Task Inference LLM to produce task predictions. The Objective Evaluator compares these predictions against the expected outputs (ground truth) using task-specific metrics to compute a score. This score represents the quality of the data optimization instructions, with prior prompt-score pairs provided as additional context to the Prompt-Generation LLM for refining instructions in future iterations.

171

172

141

data by emphasizing salient features and reducing noise. In another example for image inputs, content optimization often involves transformations such as rotation, translation, flipping, cropping, and adjustments to brightness and contrast (Jiao and Zhao, 2019). These procedures are employed to enhance model performance by augmenting the dataset and improving the representation of important features (Barrett and Cheney, 2002; Ling et al., 2021).

Traditionally, task-specific data engineering has relied heavily on domain expertise (Ling et al., 2021). For example, in the medical field, experts may derive new attributes from existing ones—such as calculating the Body Mass Index (BMI) from weight and height measurements—to create more informative features for analysis. Similarly, for data in natural language form, such as logical puzzles or mathematical problem statements, individuals with linguistic and analytical expertise may augment the text by identifying contextual cues, deducing relevant implicit information, and explicitly defining known and unknown variables to facilitate more effective interpretation.

However, employing human experts to craft and refine each input data can be both costly and timeconsuming. With recent advancements in LLMs, we propose leveraging LLMs as universal domain experts. Specifically, we investigate their ability to propose and execute content optimization procedures across datasets from diverse fields. By automating the content optimization process, we aim to transform the original dataset **D** to optimized version **D**'. The objective is to reduce reliance on human expertise while maintaining or enhancing model performance. This approach not only accelerates the data preparation phase but also has the potential to uncover novel optimization strategies that may be overlooked by human practitioners. 173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

191

192

193

194

196

197

198

200

201

203

Format Optimization concentrates on the automatic discovery of the optimal format for presenting input data to a LLM, after the content has been optimized. Recent studies have demonstrated that LLMs are highly sensitive to input formatting (Sclar et al., 2023). For example, manipulations such as positional swapping of in-context examples or alphabet shifting have been observed to influence an LLM's performance. Additionally, transforming attribute-value pairs in tabular data into structured formats like XML can enhance LLM performance on classification tasks. Similarly, converting natural language inputs into non-natural language formats using emojis, logical operators, or other symbolic figures has been shown to improve LLM performance (Lin et al., 2024a). Here, we again leverage LLM to find an optimal formatting function that maximizes the performance. By utilizing LLMs to explore various formatting strategies, we aim to identify structural reformulations that enhance the LLM's performance without altering the underlying content of the data.

## 2.2 Framework Workflow Design

The ADO framework employs a set of LLMs to automatically optimize the representation of input

204data D. As illustrated in Figure 2, the process ini-<br/>tiates with a Prompt Generation LLM, which pro-<br/>poses a data-optimization prompt  $\mathbf{P}_o$  that outlines<br/>a set of procedures for modifying D. Specifically,<br/>these procedures consist of two sequential compo-<br/>nents: the first provides step-by-step instructions<br/>for modifying the content of D, while the second<br/>details step-by-step instructions for reformulating<br/>the content-optimized data.

213

214

215

216

217

218

227

228

229

230

232

234

237

241

Subsequently, a Data Optimization LLM progressively executes the proposed data-optimization prompt by processing both  $\mathbf{P}_o$  and  $\mathbf{D}$ , instructing the model to generate the optimized data  $\mathbf{D}'$  to implement the target function  $\mathbf{D}' = f_{\text{format}}(f_{\text{content}}(\mathbf{D}))$ . The optimized data  $\mathbf{D}'$  is then submitted to a Task Inference LLM for processing, and its performance is evaluated on a reserved validation set, serving as the performance measure for  $\mathbf{P}_o$ . Finally,  $\mathbf{P}_o$  and its corresponding performance are fed back into the Prompt Generation LLM as additional context, enabling it to generate improved data-optimization prompts in future search rounds.

We now formally define the ADO framework, which involves three instances of LLMs:

Prompt Generation LLM (LLM<sub>G</sub>): Given a meta-prompt P<sub>m</sub> used to instruct generating the data-optimization-prompt P<sub>o</sub>, LLM<sub>G</sub> generates a set of candidate P<sub>o</sub>s aiming at providing instructions on how to optimize D:

$$\mathbf{P}_o = \mathrm{LLM}_{\mathcal{G}}(\mathbf{P}_m) \tag{2}$$

Data Optimization LLM (LLM<sub>O</sub>): Given a data-optimization prompt P<sub>o</sub>, LLM<sub>O</sub> optimizes D to produce the optimized data D':

$$\mathbf{D}' = \mathrm{LLM}_{\mathcal{O}}(\mathbf{P}_o, \mathbf{D}) \tag{3}$$

Task Inference LLM (LLM<sub>I</sub>): Using the optimized data D' and the task-specific instruction t, LLM<sub>I</sub> generates the final result y:

 $y = \text{LLM}_{\mathcal{I}}(\mathbf{D}', \mathbf{t}) \tag{4}$ 

In the ADO framework, the search for the optimal data-optimization prompt  $\mathbf{P}_o$  is typically conducted using a reserved set of data points S = $\{(x,y) \mid x \in \mathbf{D}_S, y \in \mathcal{Y}_{\mathbf{D}_S}\}$  where  $\mathcal{Y}_{\mathbf{D}_S}$  is the set of ground truth corresponding to  $\mathbf{D}_S$ . Given S, we sequentially utilize the three LLM instances to generate candidate prompts  $\mathbf{P}_o$ s, optimize the data  $\mathbf{D}$ , and produce the final inference result y'. By comparing the generated outputs y and with the ground truth labels y', we can evaluate the quality of each candidate  $\mathbf{P}_o$  using some task-specific loss function L(y, y'). The optimization of  $\mathbf{P}_o$  can be formulated as minimizing the loss over S:

$$\mathbf{P}_{o}^{*} = \arg\min_{\mathbf{P}_{o} \in \mathrm{LLM}_{\mathcal{G}}(\mathbf{P}_{m})} 255$$

$$\sum L(\mathrm{LLM}_{\mathcal{T}}(\mathrm{LLM}_{\mathcal{O}}(\mathbf{P}_{o}, x_{i}), \mathbf{t}), y_{i}) \quad (5)$$

250

251

252

253

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

280

281

284

285

287

288

290

291

292

293

294

295

296

$$\sum_{(x_i, y_i) \in S} L(\text{LLM}_{\mathcal{I}}(\text{LLM}_{\mathcal{O}}(\mathbf{P}_o, x_i), \mathbf{t}), y_i) \quad (5)$$

Various optimization algorithms such as Automatic Prompt Engineer (APE) (Zhou et al., 2023), Automatic Prompt Optimization (APO) (Pryzant et al., 2023), and Optimization by PROmpting (OPRO) (Yang et al., 2024; Liu et al., 2024; Zhou et al., 2023) can be employed to search for a better  $\mathbf{P}_o$  based on the loss function *L*. Nevertheless, such algorithms exhibit a potential limitation in optimizing  $\mathbf{P}_o$ . In the following subsection, we introduce the novel Diverse Prompt Search (DPS) algorithm to address the limitation.

#### 2.3 DPS Algorithm for P<sub>o</sub> Optimization

Recently, various optimization algorithms (Pryzant et al., 2023; Yang et al., 2024; Liu et al., 2024) have been proposed that leverage LLMs for automatic prompt optimization. Specifically, APE employs an LLM to propose several candidate prompts and selects the one with the best performance based on a reserved validation set. Subsequent works, such as OPRO, build upon this by directly utilizing an LLM as the prompt optimizer. For instance, OPRO instructs an LLM to iteratively propose candidate prompts, one at a time, while providing feedback on the performance of prior proposed prompts on a reserved validation set. This additional context enables the LLM to generate prompts with improved performance in subsequent iterations.

Nevertheless, recent studies (Zhang et al., 2024; Tang et al., 2024) have shown that optimizing by augmenting a single candidate prompt as context in each iteration, without any constraints on the resemblance between candidate prompts, may hinder the discovery of an optimal prompt. Despite being instructed to generate new candidate prompts that differ from previous ones, the LLM may at times converge toward semantically or lexically similar variations of prior proposed prompt(s). In our case, instead of proposing novel data optimization procedures, the LLM may keep proposing procedures that refine the wording or reorder the steps in the

297

3

- 306 307
- 30
- 310
- 311 312
- 313
- 3
- 315 316

316 317

3

319 320

321 322

3

3

328

329 330

333 334

3

3

3

- 340 341
- 342 343

344 345 prior proposed procedures. This behavior reduces diversity in prompt generation, restricting exploration to a narrow region of the prompt space and yielding only marginal performance improvements.

To this end, we propose the DPS algorithm, which also employs a LLM as the prompt optimizer, while generating multiple diverse candidate prompts for each iteration of the search process, with both semantic and lexical diversity constraints enforced to grant prompt diversity. Specifically, we request LLM<sub>G</sub> to generate k distinct candidate prompts { $\mathbf{P}_o^1, ..., \mathbf{P}_o^k$ } for each iteration of the search. For both semantic and lexical diversity among these prompts, we propose two constraints:

- Cosine similarity constraint (c<sub>1</sub>): The cosine similarity between any pair of prompts should be less than c<sub>1</sub>: cos (P<sup>i</sup><sub>o</sub>, P<sup>j</sup><sub>o</sub>) < c<sub>1</sub>, ∀i ≠ j
- METEOR Score Constraint (c<sub>2</sub>): The ME-TEOR score (Saadany and Orasan, 2021) between any pair of prompts should be less than c<sub>2</sub>: METEOR(P<sup>i</sup><sub>o</sub>, P<sup>j</sup><sub>o</sub>) < c<sub>2</sub>, ∀i ≠ j

To dynamically control the extent of prompt diversity tailored to specific tasks, we propose the novel idea of incorporating Bayesian Search (Turner et al., 2021) to automatically determine optimal values for k,  $c_1$ , and  $c_2$  based on validation set performance. Since Bayesian Search has been widely employed for hyper-parameter tuning in various deep learning models, we propose to integrate this approach with automatic prompt search by treating ADO as a standalone model, with k,  $c_1$ , and  $c_2$  as its hyper-parameters. The performance metric for each Bayesian Search iteration is defined as the highest performance achieved among all data-optimization prompts proposed by ADO with a fixed set of hyper-parameters. Such constraints ensure that the generated prompts are semantically and lexically diverse, encouraging exploration of different regions in the prompt space. For Bayesian Search details, please refer to A.1.

The generation of qualifying prompts is performed iteratively by repeatedly querying  $LLM_G$ until all k diverse prompts satisfying the above constraints are obtained. Each candidate prompt  $\mathbf{P}_o^i$ is evaluated on S, based on which result we batch update the generation  $\mathbf{P}_o$ . The evaluation involves applying the data optimization and inference steps:

• Data optimization:  $x'_i = \text{LLM}_{\mathcal{O}}(\mathbf{P}^i_o, x_i)$ where  $x_i$  is one input data in S • Result inference:  $y'_i = \text{LLM}_{\mathcal{I}}(x'_i, \mathbf{t})$  where t 346 is the task-specific instruction. 347

The performance of each candidate  $\mathbf{P}_o^i$  is assessed by computing a loss function L over S:

$$l_i = \sum_{(x_i, y_i) \in S} L(y'_i, y_i) \tag{6}$$

348

349

350

352

353

354

355

357

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

384

387

389

391

The batch of prompt-performance pairs  $(\mathbf{P}_{o}^{i}, l_{i})$ is then appended to  $\mathbf{P}_{m}$  to guide subsequent iterations of prompt generation. This feedback mechanism informs  $\text{LLM}_{\mathcal{G}}$  about the effectiveness of previously generated prompts, enabling it to generate more promising candidates in future iterations.

By iteratively refining the set of candidate prompts and incorporating performance feedback with batch update, the DPS algorithm encourages the exploration of a broader search space. This increases the likelihood of discovering more effective data optimization procedures, ultimately enhancing the performance of the LLM on the given task.

# **3** Implementation Details

This section provides key implementation details of the ADO framework, including the structure of meta-prompts, the execution of parallelized data optimization tasks, the handling of LLM hallucinations through multi-agent debate with crossvalidation. By leveraging these components, the ADO framework effectively enhances both the content and format of input data to improve performance across diverse tasks while maintaining factual accuracy and efficiency.

**Meta-Prompt** In this purely text-based data optimization framework, the data-optimization prompt  $\mathbf{P}_{o}$  must consist of instructions that can be executed by the LLM without relying on external tools or operations. To ensure this, we incorporate a comprehensive set of modality-specific constraints within the meta-prompt  $\mathbf{P}_m$  provided to LLM<sub>G</sub>. These constraints guide the prompt generation process, ensuring that LLM<sub>G</sub> avoids proposing optimization procedures that  $LLM_{O}$  is incapable of performing. For instance, when generating instructions for tabular data, the meta-prompt explicitly prohibits steps such as Principal Component Analysis (PCA), normalization, standardization, or one-hot encoding of categorical attributes, as these require tool-based operations beyond the LLM's text-based capabilities. Below is an example of  $\mathbf{P}_m$ :

**Parallelized Execution** The generated dataoptimization prompt  $\mathbf{P}_o$  typically includes multiple procedures, each addressing a specific aspect of data engineering or reformulation (e.g., missing data imputation, structural conversion). We parse the number of procedures generated from  $\mathbf{P}_o$  and employ an equivalent number of LLM instances to execute each procedure concurrently.

393

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415 416

Parallel execution provides two advantages: (1) avoiding omission or redundancy - we observed that prompting  $LLM_O$  to execute a lengthy list of detailed procedures in one go often leads to omissions and repetition. By executing procedures in parallel, we mitigate these issues by breaking down the tasks into smaller, independent units of work for each LLM instance. (2) improving time efficiency -Sequential execution of a long series of procedures can be time-consuming. Since many procedures are independent of each other and can be directly applied to the raw input data, distributing them across multiple LLM instances significantly reduces the overall time required for data optimization. For procedures that depend on sequential execution where the output of one serves as the input for the next – their execution is grouped together.

Hallucination Mitigation Instructions included 417  $\mathbf{P}_{o}$  may sometimes be implemented inaccurately 418 by  $LLM_O$  due to hallucinations. For example, if 419  $\mathbf{P}_{o}$  includes a directive such as "Please identify the 420 mathematical terminologies and provide concise 421 definitions, accompanied by examples for each." 422 LLM<sub>O</sub> may generate incorrect or inaccurate defini-423 tions for some of the terms identified. These inac-424 curacies could mislead the performance of  $LLM_{\mathcal{I}}$ , 425 potentially degrading overall output quality. 426

To mitigate the risk of hallucination and im-427 prove factual accuracy, we adapt a cross-validation 428 method inspired by (Du et al., 2023). In this frame-429 work, we introduce an additional LLM, denoted 430 as  $LLM_{\mathcal{F}}$  which reviews the optimized input data 431 to identify factual inaccuracies and provides con-432 cise explanations for any detected errors. When 433 errors are found,  $LLM_{\mathcal{F}}$ 's feedback is passed back 434 to  $LLM_{\mathcal{O}}$ , prompting it either to justify its original 435 output or to agree with the corrections suggested 436 437 by LLM<sub> $\mathcal{F}$ </sub>. By incorporating this cross-validation framework, we ensure a higher level of factual ac-438 curacy, leveraging the complementary strengths of 439 multiple LLMs to reduce the likelihood of halluci-440 nations and errors in the final output. 441

## 4 **Experiments**

In this section, we aim to evaluate: (1) the effectiveness of ADO as a standalone approach for performance enhancement, (2) whether DPS outperforms existing optimization algorithms in searching for data-optimization procedures, and (3) whether integrating ADO with other prompt engineering methods can further improve performance. 442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

#### 4.1 Experiment Settings

**Dataset** To demonstrate the wide applicability of data optimization, we conduct experiments on nine publicly available, real-world datasets across various domains where LLMs are frequently applied (Fang et al., 2024; Li et al., 2023a; Lin et al., 2024b; Rouzegar and Makrehchi, 2024). These datasets include Big-Bench StrategyQA (QA) <sup>1</sup>, Fraudulent Job Detection (Job) <sup>2</sup>, Grade School Math 8k (GSM8k) <sup>3</sup>, Amazon Beauty (AB) <sup>4</sup>, Amazon Toys (AT) <sup>5</sup>, Amazon Electronics (AE) <sup>6</sup>, Census Income (CI) <sup>7</sup>, Heart Disease (HR) <sup>8</sup>, and Financial Distress (FD) <sup>9</sup>. For each dataset, we randomly select 1,000 samples to form the validation set *S*.

Modeling The evaluation modeling is twofold. First, we evaluate the effectiveness of ADO under zero-shot prompting, with three LLMs with different backbones for generalizability. To perform data-optimization procedure search, we employ APE, OPRO, and DPS algorithms. Second, we assess whether ADO can be integrated with existing Prompt Engineering techniques (i.e., instruction optimization and data augmentation) to further enhance performance, with GPT-3.5 Turbo as the backbone. For instruction optimization, we employ either Chain-of-Thought reasoning (CoT) (Wei et al., 2022) or PE2 (Ye et al., 2023) after ADO is applied; similarly, for data augmentation, we employ In-Context Learning (ICL) (Liu et al., 2022) subsequent to employing ADO. For CoT, we

<sup>&</sup>lt;sup>1</sup>https://github.com/google/BIG-bench/tree/ main/bigbench/benchmark\_tasks/strategyqa <sup>2</sup>https://www.kaggle.com/datasets/shivamb/ real-or-fake-fake-jobposting-prediction <sup>3</sup>https://huggingface.co/datasets/DaertML/ gsm8k-jsonl <sup>4</sup>https://jmcauley.ucsd.edu/data/amazon/ <sup>5</sup>https://jmcauley.ucsd.edu/data/amazon/ <sup>6</sup>https://jmcauley.ucsd.edu/data/amazon/ <sup>7</sup>https://archive.ics.uci.edu/dataset/2/adult <sup>8</sup>https://www.kaggle.com/datasets/kamilpytlak/ personal-key-indicators-of-heart-disease <sup>9</sup>https://www.kaggle.com/c/GiveMeSomeCredit/ data?select=cs-test.csv

LLM for ADO	Algorithm	QA	Job	GSM	AB	AT	AE	CI	HD	FD	Mean
GPT-3.5 Turbo	N/A	0.578	0.619	0.285	0.124	0.129	0.211	0.788	0.617	0.639	0.443
	APE	0.575	0.633	0.721	0.161	0.184	0.241	0.839	0.687	0.658	0.522
	OPRO	0.583	0.627	0.734	0.169	0.195	0.238	0.846	0.681	0.667	0.527
	DPS	0.589	0.638	0.755	0.166	0.213	0.253	0.853	0.704	0.652	0.536
Gemini-1.5 Flash	N/A	0.569	0.607	0.299	0.137	0.115	0.197	0.791	0.625	0.612	0.439
	APE	0.581	0.621	0.698	0.159	0.176	0.219	0.827	0.701	0.661	0.516
	OPRO	0.589	0.624	0.704	0.173	0.183	0.238	0.841	0.709	0.672	0.526
	DPS	0.595	0.643	0.729	0.198	0.201	0.225	0.838	0.722	0.699	0.539
Llama-3.1 70B	N/A	0.563	0.588	0.281	0.117	0.135	0.188	0.769	0.629	0.615	0.431
	APE	0.571	0.613	0.675	0.129	0.166	0.205	0.798	0.673	0.649	0.498
	OPRO	0.574	0.619	0.693	0.135	0.173	0.213	0.806	0.692	0.657	0.507
	DPS	0.581	0.635	0.718	0.159	0.189	0.229	0.827	0.711	0.661	0.523

Table 1: ADO performance across all datasets. "LLM for ADO" denotes the LLM used within the ADO framework. "Algorithm" denotes the algorithm to search for optimal data-optimization procedures. "Mean" denotes the mean performance across all datasets. The best performance for each dataset on every LLM is highlighted in bold.

follow (Wei et al., 2022) by appending the phrase "Let's think step-by-step" at the end of the task instruction. For PE2, we employ it to search for the optimal task instruction. For ICL, we randomly select ten samples per dataset and augment them to the prompt for additional context, following the approach from (Liu et al., 2022).

480

481

482

483

484

485

486

487

488

489

490

**Evaluation metrics** We employ accuracy (with balanced accuracy for datasets that have imbalanced binary targets) and Hit@10 for the recommendation datasets from Amazon.

To evaluate the effectiveness of ADO, 491 Baselines we compare  $LLM_{\mathcal{I}}s'$  performance without data op-492 timization to the performance achieved after ADO 493 is applied. To evaluate the effectiveness of the DPS 494 algorithm on data-optimization procedure search, 495 we compare it against two recent optimization al-496 gorithms: APE and OPRO. It is important to high-497 light that ADO represents a novel sub-direction in 498 the field of prompt engineering and can be com-499 bined with existing prompt engineering techniques. Unlike a competitive relationship, ADO and tech-501 niques such as CoT, PE2, and ICL are in fact 502 complementary, enabling joint application for en-503 504 hanced performance. Thus, we utilize CoT, PE2, and ICL as baselines to observe whether combining ADO with any of these techniques achieves better performance compared to using them alone. 507

508LLM BackbonesWe employ three instances of509the same LLM as  $LLM_G$ ,  $LLM_O$ , and  $LLM_I$ . For510generalizability, we test with three different LLMs,511including GPT-3.5 Turbo, Gemini-1.5 Flash, and512Llama-3.1 70B. Additionally, Gemini-1.5 Pro is513instantiated as  $LLM_F$ , which will be employed514in Section 4.3. We set the temperature to 1.0 for515LLM<sub>G</sub> to encourage the generation of more cre-

ative content. For  $LLM_{\mathcal{O}}$  and  $LLM_{\mathcal{I}}$ , we set the temperature to 0 to obtain more consistent outputs.

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

#### 4.2 Result and Analysis

As demonstrated by Table 1, employing ADO for data optimization consistently leads to comparable or superior performance across all datasets for all three LLM backbones, compared to inferencing with unoptimized data. Additionally, DPS outperforms both APE and OPRO in eight, seven, and nine out of nine datasets for GPT-3.5 Turbo, Gemini-1.5 Flash, and Llama-3.1 70B, respectively. This highlights the effectiveness of batch-based prompt search with diverse candidates.

Furthermore, Table 2 demonstrates that integrating ADO with existing Prompt Engineering techniques, including CoT, ICL, and PE2, consistently results in a noticeable performance enhancement compared to employing these techniques alone across all nine datasets. For instance, ADO significantly enhances the effectiveness of CoT, particularly in the QA, Job, and FD datasets. For QA, applying CoT alone even results in slightly worse performance compared to not applying it, while combining CoT with ADO yields substantially better performance. These results demonstrate the complementarity of ADO towards both Instruction Optimization and Data Augmentation.

## 4.3 Ablation Study

In this section, we perform an ablation study to assess the impact of different components of the ADO framework from three perspectives: (1) whether both content optimization and format optimization are necessary, (2) whether incorporating a factualvalidation LLM (LLM<sub> $\mathcal{F}$ </sub>) improves performance, and (3) whether data-optimizing in-context exam-

Modeling variant	QA	Job	GSM	AB	AT	AE	CI	HD	FD	Mean
GPT	0.578	0.619	0.285	0.124	0.129	0.211	0.788	0.617	0.639	0.443
GPT w/ CoT	0.571	0.663	0.698	0.127	0.137	0.198	0.827	0.678	0.688	0.510
GPT w/ CoT + ADO	0.679	0.807	0.851	0.185	0.219	0.257	0.879	0.751	0.789	0.602
GPT w/ ICL	0.584	0.617	0.294	0.141	0.147	0.225	0.809	0.651	0.653	0.458
GPT w/ ICL + ADO	0.597	0.641	0.778	0.199	0.223	0.262	0.851	0.728	0.668	0.549
GPT w/ PE2	0.592	0.634	0.301	0.162	0.152	0.209	0.838	0.649	0.685	0.469
GPT w/ PE2 + ADO	0.618	0.659	0.312	0.183	0.178	0.234	0.863	0.697	0.722	0.496

Table 2: Performance when ADO is combined with other Prompt Engineering techniques, using GPT-3.5 Turbo as the backbone (denoted as "GPT"). "CoT + ADO" denotes applying both CoT and ADO, "ICL + ADO" denotes applying both ICL and ADO, and "PE2 + ADO" denotes applying both PE2 and ADO. For each dataset on each technique, any performance enhancement resulting from ADO integration is highlighted in bold.

ples yields performance gains. For experiment details, please refer to A.2. The results of all three experiments are presented in Table 3 in the Appendix. As the table demonstrates, both content and format optimizations are essential for performance: removing format optimization significantly reduced performance on recommendation datasets and the CI dataset, while removing content optimization led to declines on other datasets. Moreover, incorporating LLM<sub> $\mathcal{F}$ </sub> for hallucination mitigation produced comparable or improved performance across all datasets, with the most significant gains on the QA, Job, and GSM datasets. Finally, optimizing incontext examples led to noticeable improvements, particularly on the Job, GSM, and FD datasets.

## 5 Related Work

551

552

554

556

558

560

561

562

563

564

566

567

568

569

570

571

573

576

577

578

580

584

585

588

Numerous approaches have been proposed for modifying prompts to enhance LLM performance, such as In-Context Learning and Instruction Optimization. In-Context Learning concentrates on providing the LLM with additional in-prompt exemplars from the same task domain, typically in the form of input data paired with their corresponding labels or outputs (Wei et al., 2023; Dong et al., 2022; Shin et al., 2022). This method capitalizes on the model's ability to generalize from in-prompt examples, enabling the LLM to better comprehend the expected output format and task-specific requirements based on the provided exemplars.

Instruction Optimization aims to modify the instruction part of the prompt to improve LLM performance. For example, Si et al. (2022) points out that composing better instructions can greatly boost LLM's performance on task inferencing. Wei et al. (2022) proposes CoT reasoning, which introduces immediate reasoning steps into the output generation process. As demonstrated by (Wei et al., 2022), employing zero-shot CoT substantially improve LLM performance tasks including logical reasoning, fraud detection, among many others. Extending beyond manually crafted instructions, various studies have proposed automated methods to search for optimal instructions tailored to specific tasks (Zhou et al., 2023; Pryzant et al., 2023; Yang et al., 2024). For instance, APE (Zhou et al., 2023) introduces an iterative Monte Carlo search to refine prompt instructions. It first uses an instructionproposing LLM to generate a set of candidate instructions, then evaluates each on a validation set to select the best-performing candidates. 589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

Despite these advances, directly optimizing the presentation of input data has received little attention. In this work, we hypothesize that optimizing both the data content and format may yield performance improvement when employing LLM for task inferencing. Building on the principles of automatic prompt optimization, we propose a novel framework called Automatic Data Optimization (ADO). In ADO, an LLM, denoted as  $LLM_G$ , iteratively proposes and searches data-optimization instructions aimed at maximizing LLM performance.

## 6 Conclusions

In this paper, we introduce a new sub-direction of prompt engineering: input data optimization, facilitated by the ADO framework and the DPS algorithm. The ADO framework automates content and format optimization by leveraging LLMs as universal domain experts, reducing the need for manual data processing. DPS enhances this process by generating diverse data optimization prompts, enabling broader exploration and increasing the likelihood of identifying optimal procedures. Empirical results demonstrate that ADO not only improves modeling performance when used alone but also further enhances performance when combined with other prompt engineering methods.

## 7 Limitations

627

647

649

671

672

673

676

As we explore the novel approach of input data optimization within prompts, we question whether 629 it is possible to simultaneously search for both the optimal instruction and the optimal procedures for input data optimization in a specific inference task. Currently, as detailed in the paper, we first search for the optimal data representation using ADO, 634 and then for the optimal instruction using PE2. 635 However, this process involves two distinct steps, and it would be more efficient to search for both 637 the instruction and data optimization concurrently. Therefore, in the future, we aim to investigate the feasibility of jointly optimizing both components, as proposed in (Sordoni et al., 2024; Chen et al., 641 2024), to further enhance LLM performance. 642

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Awais Ahmad, Murad Khan, Anand Paul, Sadia Din, M Mazhar Rathore, Gwanggil Jeon, and Gyu Sang Choi. 2018. Toward modeling and optimization of features selection in big data based social internet of things. *Future Generation Computer Systems*, 82:715–726.
- Asma Baccouche, Begonya Garcia-Zapirain, Cristian Castillo Olea, and Adel Elmaghraby. 2020. Ensemble deep learning models for heart disease classification: A case study from mexico. *Information*, 11(4):207.
- William A Barrett and Alan S Cheney. 2002. Objectbased image editing. In *Proceedings of the 29th annual conference on Computer graphics and interactive techniques*, pages 777–784.
- Sondos Mahmoud Bsharat, Aidar Myrzakhan, and Zhiqiang Shen. 2023. Principled instructions are all you need for questioning llama-1/2, gpt-3.5/4. *arXiv preprint arXiv:2312.16171*.
- Banghao Chen, Zhaofeng Zhang, Nicolas Langrené, and Shengxin Zhu. 2023. Unleashing the potential of prompt engineering in large language models: a comprehensive review. *arXiv preprint arXiv:2310.14735*.
- Yongchao Chen, Jacob Arkin, Yilun Hao, Yang Zhang, Nicholas Roy, and Chuchu Fan. 2024. Prompt optimization in multi-step tasks (promst): Integrating human feedback and preference alignment. *arXiv preprint arXiv:2402.08702*.

Viet-Tung Do, Van-Khanh Hoang, Duy-Hung Nguyen, Shahab Sabahi, Jeff Yang, Hajime Hotta, Minh-Tien Nguyen, and Hung Le. 2024. Automatic prompt selection for large language models. *arXiv preprint arXiv:2404.02717*. 677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

719

721

722

723

724

725

726

727

728

729

- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2023. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*.
- Lizhou Fan, Wenyue Hua, Lingyao Li, Haoyang Ling, Yongfeng Zhang, and Libby Hemphill. 2023. Nphardeval: Dynamic benchmark on reasoning ability of large language models via complexity classes. *arXiv preprint arXiv:2312.14890*.
- Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang, Ziqing Hu, Yanjun (Jane) Qi, Scott Nickleach, Diego Socolinsky, "SHS" Srinivasan Sengamedu, and Christos Faloutsos. 2024. Large language models (llms) on tabular data: Prediction, generation, and understanding - a survey. *Transactions on Machine Learning Research*.
- Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2022. Recommendation as language processing (rlp): A unified pretrain, personalized prompt & predict paradigm (p5). In *Proceedings of the 16th ACM Conference on Recommender Systems*, pages 299–315.
- Wenyue Hua, Lei Li, Shuyuan Xu, Li Chen, and Yongfeng Zhang. 2023. Tutorial on large language models for recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 1281–1283.
- Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*.
- Licheng Jiao and Jin Zhao. 2019. A survey on the new generation of deep learning in image processing. *Ieee Access*, 7:172231–172263.
- Lei Li, Yongfeng Zhang, and Li Chen. 2023a. Prompt distillation for efficient llm-based recommendation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pages 1348–1357.
- Lingyao Li, Jiayan Zhou, Zhenxiang Gao, Wenyue Hua, Lizhou Fan, Huizi Yu, Loni Hagen, Yonfeng Zhang, Themistocles L Assimes, Libby Hemphill, et al. 2024a. A scoping review of using large language models (llms) to investigate electronic health records (ehrs). *arXiv preprint arXiv:2405.03066*.

837

838

839

840

785

786

787

Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. 2023b. Large language models in finance: A survey. In *Proceedings of the fourth ACM international conference on AI in finance*, pages 374–382.

730

731

734

738

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

757

765

767

770

771

772

773

774

775

776

778

780

781

- Zelong Li, Jianchao Ji, Yingqiang Ge, Wenyue Hua, and Yongfeng Zhang. 2024b. Pap-rec: Personalized automatic prompt for recommendation language model. *arXiv preprint arXiv:2402.00284*.
- Guo Lin, Wenyue Hua, and Yongfeng Zhang. 2024a. Promptcrypt: Prompt encryption for secure communication with large language models. *arXiv preprint arXiv:2402.05868*.
- Xinyu Lin, Wenjie Wang, Yongqi Li, Shuo Yang, Fuli Feng, Yinwei Wei, and Tat-Seng Chua. 2024b. Dataefficient fine-tuning for llm-based recommendation. In Proceedings of the 47th International ACM SI-GIR Conference on Research and Development in Information Retrieval, pages 365–374.
- Huan Ling, Karsten Kreis, Daiqing Li, Seung Wook Kim, Antonio Torralba, and Sanja Fidler. 2021. Editgan: High-precision semantic image editing. Advances in Neural Information Processing Systems, 34:16331–16345.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. Advances in Neural Information Processing Systems, 35:1950–1965.
- Shengcai Liu, Caishun Chen, Xinghua Qu, Ke Tang, and Yew-Soon Ong. 2024. Large language models as evolutionary optimizers. In 2024 IEEE Congress on Evolutionary Computation (CEC), pages 1–8. IEEE.
- Ggaliwango Marvin, Nakayiza Hellen, Daudi Jjingo, and Joyce Nakatumba-Nabende. 2023. Prompt engineering in large language models. In *International conference on data intelligence and cognitive informatics*, pages 387–402. Springer.
- Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with" gradient descent" and beam search. *arXiv preprint arXiv:2305.03495*.
- Hamidreza Rouzegar and Masoud Makrehchi. 2024. Enhancing text classification through llm-driven active learning and human annotation. *arXiv preprint arXiv:2406.12114*.
- Hadeel Saadany and Constantin Orasan. 2021. Bleu, meteor, bertscore: evaluation of metrics performance in assessing critical translation errors in sentimentoriented text. *arXiv preprint arXiv:2109.14250*.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. arXiv preprint arXiv:2402.07927.

- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*.
- SP Sharan, Francesco Pittaluga, Manmohan Chandraker, et al. 2023. Llm-assist: Enhancing closed-loop planning with language-based reasoning. *arXiv preprint arXiv:2401.00125*.
- Seongjin Shin, Sang-Woo Lee, Hwijeen Ahn, Sungdong Kim, HyoungSeok Kim, Boseop Kim, Kyunghyun Cho, Gichang Lee, Woomyoung Park, Jung-Woo Ha, et al. 2022. On the effect of pretraining corpora on in-context learning by a large-scale language model. *arXiv preprint arXiv:2204.13509*.
- Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Lijuan Wang. 2022. Prompting gpt-3 to be reliable. *arXiv preprint arXiv:2210.09150*.
- Alessandro Sordoni, Eric Yuan, Marc-Alexandre Côté, Matheus Pereira, Adam Trischler, Ziang Xiao, Arian Hosseini, Friederike Niedtner, and Nicolas Le Roux. 2024. Joint prompt optimization of stacked llms using variational inference. Advances in Neural Information Processing Systems, 36.
- Xinyu Tang, Xiaolei Wang, Wayne Xin Zhao, Siyuan Lu, Yaliang Li, and Ji-Rong Wen. 2024. Unleashing the potential of large language models as prompt optimizers: An analogical analysis with gradient-based model optimizers. *arXiv preprint arXiv:2402.17564*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Ryan Turner, David Eriksson, Michael McCourt, Juha Kiili, Eero Laaksonen, Zhen Xu, and Isabelle Guyon. 2021. Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the black-box optimization challenge 2020. In *NeurIPS 2020 Competition and Demonstration Track*, pages 3–26. PMLR.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,

841

877

887

892

Da Huang, Denny Zhou, et al. 2023. Larger language models do in-context learning differently. arXiv preprint arXiv:2303.03846.

- Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. 2024. A survey on large language models for recommendation. World Wide Web, 27(5):60.
- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. arXiv preprint arXiv:2303.17564.
- Chengrun Yang, Xuezhi Wang Wang, Yifeng Lu Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024. Large language models as optimizers. ICLR.
- Qinyuan Ye, Maxamed Axmed, Reid Pryzant, and Fereshte Khani. 2023. Prompt engineering a prompt engineer. arXiv preprint arXiv:2311.05661.
- Huizi Yu, Lizhou Fan, Lingyao Li, Jiayan Zhou, Zihui Ma, Lu Xian, Wenyue Hua, Sijia He, Mingyu Jin, Yongfeng Zhang, et al. 2024a. Large language models in biomedical and health informatics: A bibliometric review. arXiv preprint arXiv:2403.16303.
- Huizi Yu, Jiayan Zhou, Lingyao Li, Shan Chen, Jack Gallifant, Anye Shi, Xiang Li, Wenyue Hua, Mingyu Jin, Guang Chen, et al. 2024b. Aipatient: Simulating patients with ehrs and llm powered agentic workflow. arXiv preprint arXiv:2409.18924.
- Tuo Zhang, Jinyue Yuan, and Salman Avestimehr. 2024. Revisiting opro: The limitations of small-scale llms as optimizers. arXiv preprint arXiv:2405.10276.
- Alice Zheng and Amanda Casari. 2018. Feature engineering for machine learning: principles and techniques for data scientists. " O'Reilly Media, Inc.".
- Michelle X Zhou and Vikram Aggarwal. 2004. An optimization-based approach to dynamic data content selection in intelligent multimedia interfaces. In Proceedings of the 17th annual ACM symposium on User interface software and technology, pages 227-236.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023. Large language models are human-level prompt engineers. ICLR.

#### Appendix А

## A.1 Bayesian Search Specifics

Bayesian Search is an informed search which achieves better performance than uninformed searches such as Random Search (Turner et al., 2021). In this work, we propose to incorporate Bayesian Search as part of the data-optimization 893 procedure search, by tuning  $k, c_1$ , and  $c_2$  as "hyper-894 parameters" based on performance of the validation 895 set S. This enables us to dynamically control both 896 the number of candidate prompt to be generated 897 per iteration for batch update, as well as the degree 898 of diversity among candidate prompts. 899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

# A.2 Ablation Study Specifics

Data Optimization Objectives We evaluate the effectiveness of the two optimization objectives-content optimization and format optimization-in ADO. To this end, we constrain the dataoptimization prompt  $\mathbf{P}_o$  to focus on either data engineering procedures (content optimization) or structural reformulation (format optimization), using zero-shot CoT as the prompting format. Specifically, we modify the meta-prompt  $\mathbf{P}_m$  to explicitly prohibit instructions related to the non-evaluated aspect, ensuring  $\mathbf{P}_{o}$  is restricted to either content or format optimization. These are denoted as "ADO-Engineering" (data engineering only) and "ADO-Reformulation" (structural reformulation only).

Factual-validation LLM We also investigate whether integrating the factual-validation LLM  $(LLM_{\mathcal{F}})$ , into the ADO workflow, as described in Section 3 enhances performance. Using zero-shot CoT, we perform cross-validation on optimized input data, iterating between  $LLM_{\mathcal{F}}$  and  $LLM_{\mathcal{O}}$ until a consensus is reached or a maximum of four rounds is completed. If no consensus is reached, the optimized input from the final validation round is used for prompt construction. This configuration is referred to as "ADO w/ Factual-check."

**Optimized Input for ICL** In Section 4, all incontext examples are presented in their unoptimized form. Here, we examine whether optimizing the input data of ICL examples, using the same procedures applied to the evaluation data, leads to improved performance. The hypothesis is that optimized in-context examples will better align with the evaluation input, facilitating easier learning for the LLM. We optimize the ICL input data and augment the prompt with these optimized examples paired with their respective outputs, denoted as "ADO on ICL Samples."

Table 3 presents the ablation study results. For the first experiment: both data engineering and structural reformulation are crucial for maintaining performance. Limiting optimization to data engineering led to a significant drop in performance

	QA	Job	GSM	AB	AT	AE	CI	HD	FD
ADO-Engineering	0.667	0.789	0.843	0.155	0.177	0.229	0.839	0.742	0.776
ADO-Reformulation	0.602	0.719	0.734	0.189	0.208	0.253	0.868	0.684	0.705
ADO w/ Factual-check	0.691	0.823	0.864	0.187	0.221	0.262	0.884	0.747	0.795
ADO on ICL Samples	0.599	0.682	0.803	0.187	0.228	0.267	0.871	0.734	0.691

Table 3: Ablation Study Performance.

on all recommendation datasets and the CI dataset, while restricting optimization to structural reformulation resulted in performance degradation on the other datasets. For the second experiment: incorporating  $LLM_{\mathcal{F}}$  for factual validation produced comparable or improved performance across all datasets, with the most significant gains on the QA, Job, and GSM datasets. Finally, optimizing incontext examples led to noticeable improvements, particularly on the Job, GSM, and FD datasets.

954 Dataset Description: <description> - ... - ... 

Your task is to propose a creative, detailed, and step-by-step algorithm to enrich and then reformulate samples in this dataset. The goal of the algorithm is to perform thorough data engineering and reformulation on the sample, so that it is easier for an LLM to generate the target outputs. Below are some example dataset samples with target outputs as references: Examples: - <sample input1>; Output: <sample output1> - <sample input2>; Output: <sample output2> - <sample input3>; Output: <sample output3> Please Note: - Do NOT refer to any external database. - Do NOT perform vector generations. - ONLY propose steps that an LLM can do on its own. Below is a list of prior-proposed data optimization algorithms, provided to you as additional context: - Algorithm 1; Score: a1 - Algorithm 1; Score: a2

Listing 1: Meta Prompt Example