Distributed In-Context Learning under Non-IID Among Clients

Anonymous ACL submission

Abstract

Advancements in large language models (LLMs) have shown their effectiveness in multiple complicated natural language reasoning tasks. A key challenge remains in adapting these models efficiently to new or unfamiliar tasks. In-context learning (ICL) provides a promising solution for few-shot adaptation by retrieving a set of data points relevant to a query, called in-context examples (ICE), from a training dataset and providing them during the inference as context. Most existing studies utilize a centralized training dataset, yet many real-world datasets may be distributed among multiple clients, and remote data retrieval can be associated with costs. Especially when the client data are non-identical independent distributions (non-IID), retrieving from clients a proper set of ICEs needed for a test query presents critical challenges. In this paper, we first show that in this challenging setting, test queries will have different preferences among clients because of non-IIDness, and equal contribution often leads to suboptimal performance. We then introduce a novel approach to tackle the distributed non-IID ICL problem when a data usage budget is present. The principle is that each client's proper contribution (budget) should be designed according to the preference of each query for that client. Our approach uses a data-driven manner to allocate a budget for each client, tailored to each test query. Through extensive empirical studies on diverse datasets, our method demonstrates superior performance relative to competing baselines.

1 Introduction

040

043

Recent significant progress in large language models (LLMs) (Achiam et al., 2023; Touvron et al., 2023a,b; Team et al., 2023) has demonstrated their effectiveness across various natural language processing (NLP) tasks (Wang et al., 2018, 2019). Despite their impressive performances, they still require adaptation to the specific downstream tasks



Figure 1: Problem overview. When datasets are distributed among clients in a non-IID manner, it creates an obstacle in generating a good context (left). However, by assigning appropriate budgets to leverage per-client expertise, better context can be created (right).

for better performance. However, adaptation poses challenges due to LLMs' vast number of trainable parameters.

In-context learning (ICL) (Dong et al., 2022) is a notable method that distinguishes itself through both its effectiveness and efficiency. In brief, ICL adapts to the target task by incorporating context information following two primary steps: i) identify samples from the training dataset helpful to solve the target query by creating a prompt describing a context; ii) feed the constructed prompt with the target query and get the answer. Previous related works on ICL mainly have focused on the construction of a prompt describing the context, which involves several sub-problems, such as the retrieval of in-context examples (ICEs) (Robertson et al., 2009) and determining the optimal sequence for the selected ICEs (Zhang et al., 2024).

A common assumption in most existing ICL research is that the system has access to a highquality centralized dataset used for retrieval. However, in many application scenarios, such as health informatics, centralized datasets may not be fea-

066

ples, such as client 1, it can create a more relevant context to answer the query related to (+) opera-

tion. This indicates that under non-IID, the server should allocate the budgets over clients based on the preference of each query itself, as well as the distribution of local training samples. Motivated by this, we propose a novel distributed ICL framework to collaboratively collect scattered information among non-IID clients by properly assigning ICE budgets to each client. First, the server will gather the optimal budget statistics using an existing proxy dataset on the server side. Next, the server will use this dataset to train the budget allocator. During the deployment stage, the server

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

151

152

153

154

155

156

157

158

159

160

161

162

163

164

167

will predict the proper budget for each client using this trained budget allocator given each test query and perform ICL among clients. Furthermore, in practical scenarios where privacy concerns arise, we augment our framework with the paraphrasing method (Mohtashami et al., 2023) to secure privacy.

Contributions. A summary of our contributions:

- To the best of our knowledge, we are the first to study the challenging real-world setting of ICL with distributed non-IID clients. We identify the principal challenge as properly assigning the ICE budget for non-IID clients based on the preference of each test query and local knowledge distribution.
- · We propose a framework to handle the distributed non-IID ICL. This framework trains a budget allocator on the server with the help of a server-side proxy dataset. Then, the server will use this trained allocator to decide how many ICEs to retrieve from each client for the ICL process, enabling collaborative action among clients.
- Across a range of dataset benchmarks featuring various non-IID configurations as well as on different LLM architectures, our approach has been validated to enhance ICL performance. Notably, we examine both non-private, *i.e.*, communicate raw samples directly, and private cases using the paraphrasing method to secure privacy. In both scenarios, our approach shows superiority to the previous method and other reasonable baselines.

2 Problem Formulation

In this section, we provide a detailed problem formulation. First, we begin with the specifics of in-

sible, and data could be distributed in different institutions, which calls for the distributed ICL. In addition, when the data is proprietary and possesses high value towards inferences, access to data entries may also be bound to data pricing strategies (Xu et al., 2023; Cong et al., 2022). For instance, the system needs to pay the local institution based on the number of samples sent to the system to share profits from inferences (Tang et al., 2020). Under this scenario, aggregating ICEs from local clients to a center server for ICL entails significant financial costs and lacks efficiency.

067

068

075

077

097

100

101

102

103

105

107

108

110

111

112

113

114

115

116

117

118

In this paper, we focus on integrating knowledge from distributed clients to achieve better ICL performance under the per-query ICE budget constraint. Specifically, we formalize the distributed ICL problem where the ICEs are distributed on clients, and the server has an LLM for ICL inference but can only request a limited number of ICEs from all clients for each query, which we refer to as the ICE budget.

We begin by identifying the key challenge in distributed ICL with ICE budget constraints lies in the non-independently and identically distributed (non-IID) training data, as shown in Section 3.1. For example, in Figure 1, data samples are spread across C clients, each with a unique data distribution. Specifically, client 1 primarily contains (+)samples, while client 2 is mainly constituted by (-) examples. Only limited research (Mohtashami et al., 2023) tried to address the challenge of distributed datasets for ICL, while none considers the challenging real-world setting of non-IID clients. This leaves a critical question unanswered: What happens to distributed ICL when local clients are non-IID?

To further the understanding of the key challenge in the distributed non-IID ICL, we explore the local retrieval process on non-IID clients. We found that each query has different preferences for different clients based on local knowledge distribution, that is, the number of samples needed from different clients should vary based on local sample distribution. As the toy example shown in Figure 1, when the server creates context by uniformly assigning budgets to clients, the answer might be incorrect due to the insufficiency of (+) information in the context. To be more detailed, the server assigns the clients who have expertise on $(-), (\times)$, and (\div) operations with the same budget as on (+), without any preference. Nevertheless, if the server assigns more budget to clients with many (+) sam-

- 168
- 169

185

191

192

193

194

195

196

197

198

204

205

209

context learning (ICL), followed by a description of distributed non-IID ICL.

170 2.1 In-Context Learning

171Notation. We consider a NLP tasks which have172training dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ with N training173samples. Here, x_i is the input text, and y_i is the174corresponding output. In the test phase, a test query175 x_q is given.

176**Retrieval.** We employ the off-the-shelf pre-trained177retriever KATE (Liu et al., 2021)¹, which utilizes178k-NN example selection. This retriever employs179a sentence encoder $\mathcal{E}(\cdot)$ to measure the similarity180between the in-context example x_i in dataset \mathcal{D} and181the query x_q as follows:

$$d(e_i, e_q) = ||e_q - e_i||_2, \tag{1}$$

where $e_q = \mathcal{E}(x_q)$ and $e_i = \mathcal{E}(x_i)$. We select k samples using the following criterion:

$$\mathcal{T}(e_q, k | \mathcal{D}) = \arg_{e_i \in \mathcal{E}(x_i) \forall (x_i, y_i) \in \mathcal{D}} \operatorname{Top-}k(d(e_i, e_q)), \quad (2)$$

where $\mathcal{T}(e_q, k | \mathcal{D})$ denotes the selected samples from the dataset \mathcal{D} , and used for inference.

ICL Inference. In the test phase, given a test query with input x_i , relevant k training samples called incontext examples (ICEs) are selected, *i.e.*, $S = \mathcal{T}(e_q, k | D)$. Based on the retrieved samples, we feed the constructed context prompt $s(S, x_q)$ into LLM for inference and obtain results via:

$$y_t = \arg\max_{y} p_{\text{LLM}}(y|s(S, x_q), y_{< t})$$
(3)

where
$$s(S, x_q) = (x_1, y_1) \odot \ldots \odot (x_k, y_k) \odot x_q$$
,

where the \odot operation denotes concatenation, and $s(S, x_q)$ is the context constructed using query x_q and samples in S; the term p_{LLM} represents the output softmax probability of the LLM, functioning autoregressive, meaning that the output up to time t, *i.e.*, $y_{< t}$, is input back into the model to generate the t^{th} output, y_t . Previous works (Ye et al., 2023; Levy et al., 2022) on ICL mainly focus on the selection of S under a centralized setting. However, we investigate the scenario where \mathcal{D} is split among several clients, each following non-IID distributions.

2.2 Distributed non-IID ICL

Distributed ICL Setting. We consider C clients with a centralized server in our system. Each

client $c \in [C]$ has local training dataset $\mathcal{D}_c =$ $\{(x_i^c, y_i^c)\}_{i=1}^{N_c}$ with N_c training samples. Note that \mathcal{D}_c follows different distributions for different clients. We follow the non-IID conditions as defined in (Li et al., 2022), with details provided in Appendix A. In summary, we allocate data on a perclass basis, where each client receives a specific number of classes, meaning each client has samples from only specified classes. Clients and the server have identical off-the-shelf pre-trained retrievers. Consider the computation resource limitation on clients as in many real scenarios (Yoo et al., 2022), only the server is equipped with an LLM. Moreover, the server has limited query-only proxy dataset $\mathcal{D}_{\text{proxy}} = \{x_j^{\text{proxy}}\}_{j=1}^{N_{\text{proxy}}}$, that $N_{\text{proxy}} \ll \sum_{c=1}^{C} N_c$. The server has quite a small $\mathcal{D}_{\text{proxy}}$, and it is an auxiliary dataset to extract information for collaboration to make the problem feasible. Notice that we do not require the ground truth label information of the proxy dataset, only input queries are sufficient, which reduces the difficulty on collecting such dataset in practical setting.

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

227

228

229

230

231

232

233

234

235

236

238

240

241

242

243

244

245

246

247

251

252

254

255

256

Pipeline. First, the server requests relevant samples from each client by sending x_q to all clients with local budgets k_c . Remark that each query x_q has its own preference of each client, which can be represented as k_c . A larger k_c indicates the given test query x_q prefers more information from client c, compared with client c' with a smaller $k_{c'}$. Here, x_a can be anonymized by paraphrasing, as done in previous works (Mohtashami et al., $(2023)^2$. Each client then selects the most relevant k_c samples from their local training dataset, *i.e.*, $S_c = \mathcal{T}(e_q, k_c | \mathcal{D}_c) \subset \mathcal{D}_c$, and returns them to the server. The server receives S_c from clients and generates the context s based on the merged examples, $S = \bigcup_{c=1}^{C} S_c$. In the final step, the server infers y using $s(S, x_q)$. The entire framework also can be described in Figure 2. In this paper, we are concentrating on assigning k_c to each client as described in Figure 2.

3 Observations

In this section, we describe several empirical supports to handle the distributed non-IID ICL. First, we demonstrate that non-IID distributions hinder the merging of scattered information. We then establish our goal, termed as *oracle budget*, which

¹We do not fine-tune the retriever for each task, which is impractical because we cannot gather the distributed datasets.

²Although our main experiments utilize the nonparaphrased dataset, we also present the paraphrased results in Section 5.



Figure 2: Overview of the pipeline: First, the **budget allocator** assigns a budget to each client based on the question. Subsequently, each client retrieves their relevant samples and sends them back to the server. The server infers the answer by feeding the question, which is composed of concatenated context examples and the query.



Figure 3: Non-IID experimental results. It shows that centralized performance is comparable to the IID case, whereas non-IID scenarios exhibit a significant declined performance. This highlights the critical importance of addressing non-IIDness to find a solution.

reflects the server's preference for each client if the server knows all distributed data. Finally, we check if predicting the oracle budget of each test query for inference is feasible.

3.1 Non-IIDness Leads to Performance Drop

261

262

263

270

274

275

276

278

First of all, we evaluate the effect of non-IIDness. Straightforwardly, we distribute the budget $\{k_c\}_{c=1}^C$ uniformly according to the following criteria: Given C clients are involved in answering this question, and the number of samples for context is k. We first explore the naïve equally assigned local budget scheme in both IID and non-IID settings. That is, each client $c \in [C]$ locally retrieves top- k_c samples where $k_c = \lceil \frac{k}{C} \rceil$ from local dataset \mathcal{D}_c . Detailed experimental settings are described in Appendix B.

As illustrated in Figure 3, we observe the followings: (1) There is no significant performance degradation between the centralized case (\blacksquare) and the IID case (\blacksquare). This is expected, as the merged top- k_c samples in the IID case closely resemble the centralized top-k samples. Any minor discrepancies are attributed to differences in sample ordering. (2) However, performance degradation becomes pronounced in non-IIDness case (refer to the comparison between , and). Hereinafter, we gather insights to address the distributed non-IID ICL. 279

281

284

285

286

288

290

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

3.2 Proper Budget per Query for Each Client

Oracle budget. The remaining issue is that to make the server operate similarly in a centralized manner, it needs to allocate the budget as if it has complete knowledge of all clients. We call this budget for each client as the *oracle budget* for query embedding e_q and define it as follows:

$$k_c^{\star}(e_q) = \left| \mathcal{T}(e_q, k | \mathcal{D}_c) \cap \mathcal{T}(e_q, k | \mathcal{D}) \right|,$$

where $\mathcal{T}(\cdot)$ is defined as Eq. (2) and $|\cdot|$ is set cardinality. Note that the physical meaning of $k_c^*(e_q)$ is the number of shared samples between the top-krelevant to e_q in local \mathcal{D}_c and global \mathcal{D} datasets.

Check of predictability of oracle budget. For the next step, it is necessary to check if e_q has sufficient patterns of oracle budget to extract and use it in the inference phase. Our hypothesis is that similar queries may share similar oracle budget patterns and preferences on the same client, and it can lead to similar budget allocations for that client. Therefore, to verify this hypothesis, we perform t-SNE analysis (Van der Maaten and Hinton, 2008) on the embeddings obtained from the retriever for queries. Furthermore, we color each sample based on the oracle budget $k_c^*(e_q)$. As described in Figure 5 of Appendix F, similar query embeddings exhibit similar oracle budget patterns. This indicates that, given

Algorithm 1 Top- k sampling, $\mathcal{T}(e, k \mathcal{D})$
Require: Query embedding e , Encoder $\mathcal{E}(\cdot)$
/* Compute embedding */
1: for $(x_i,y_i)\in\mathcal{D}$ do
2: $e_i \leftarrow \mathcal{E}(x)$
3: end for
/* Select top-k samples */
4: $S = \arg \operatorname{Top-}k \ e - e_i\ _2$
5: Return: S

a test query, we can infer the budget assignment for each client. However, it is challenging to predict 312 fine-grained budget value since there are no rigid classification patterns. For instance, determining 314 the detailed budget value seems challenging in the case of client 1 in SST-5. Therefore, developing an 316 efficient method to infer the exact budgets based on 317 these broad patterns for each client is required. We also conduct t-SNE visualization on clients under other non-IID with task shifting and feature skew as shown in Figure 9 of Appendix F, which helps us 321 to conclude this observation holds under different non-IID settings.

3.3 Observation Summary

In summary, our findings and the approach for designing a method are as follows: (1) non-IIDness significantly affects the distributed ICL setting, necessitating the development of a coalition method. To handle this problem, it is straightforward to allocate an appropriate number of budgets to each client, *i.e.*, making server work so as it knows client all samples. (2) By analyzing the query embeddings, we can determine the importance of each client per query.

4 Method

327

331

332

335

336

337

338

341

342

344

346

In this section, we outline the proposed method to mitigate non-IIDness in the ICL framework. Specifically, we show how to train the *budget allocator* and conduct inference.

4.1 Train a Budget Allocator

Based on Section 3, it is feasible to assign budgets of each client by using the embeddings obtained from the retriever encoder \mathcal{E} . We first construct the datasets having the targeting budget values and then train the budget allocator. The pseudo-codes are described in Algorithm 2 and 1.

347 **Construct dataset for oracle budget.** First, we

Algorithm 2 Construct dataset

Require: Encoder $\mathcal{E}(\cdot)$, server-side ICE budget k,
proxy dataset $\mathcal{D}_{\text{proxy}} = \{(x_j, y_j)\}_{i=1}^{N_{\text{proxy}}}$ Quan-
tization parameter δ .
1: for $(x_j^{\text{proxy}}, y_j^{\text{proxy}}) \in \mathcal{D}_{\text{proxy}}$ do
2: $e_j^{\text{proxy}} = \mathcal{E}(x_j^{\text{proxy}})$
/* Get distributed examples */
3: for $c \in [C]$ do
4: $e_j^{\text{proxy}} \to \text{Client } c$
5: $\check{S_c} = \mathcal{T}(e_j^{\text{proxy}}, k \mathcal{D}_c)$
6: Server $\leftarrow S_c$
7: end for
/* Construct optimal example */
8: $S = \bigcup_{c=1}^{C} S_c$
9: $S^{\text{top}} = \arg \operatorname{Top} k \ e_j^{\text{proxy}} - \mathcal{E}(x_s)\ _2$
$(x_s, y_s) \in S$ /* Compute proper budget size for each c */
10: $k_c(e_j) = S^{\text{top}} \cap S_c / \delta \forall c \in [C]$
11: end for
12: $B_{\text{proxy}} = \{(e_j, \{k_c(e_j)\}_{c=1}^C)\}_{j=1}^{N_{\text{proxy}}}$
13: Return: B _{proxy}

explain how to create a dataset to train the budget allocator for each client, as described in Algorithm 2. Given proxy dataset $\mathcal{D}_{\text{proxy}}$, for all embeddings $e_i = \mathcal{E}(x_i)$ where $x_i \in \mathcal{D}_{\text{proxy}}$, we request k samples from each client $c \in [C]$ using Top-k procedure, *i.e.*, $S_c = \mathcal{T}(e, k | \mathcal{D}_c)$. Once the server receives k examples from each clients, *i.e.*, $\{S_c\}_{c=1}^C$, it merges and re-orders them to obtains S^{top} . Based on S^{top} , we count the number of samples from each client in S^{top}, *i.e.*, compute $k_c(e_i)$. After counting $k_c(e_i)$ for all clients, we quantize the budget levels for each client using the quantization hyper-parameter δ . As a result, the output of this procedure is B_{proxy} for all clients, composed of embeddings e and their respective budgets $k_c(e_j)$.

348

350

351

352

354

355

356

357

360

361

362

363

364

365

366

367

369

370

371

373

Train budget allocator. Based on the constructed dataset B_{proxy} , we train the *budget allcoators*, *i.e.*, $\{f_c(\cdot)\}_{c=1}^C$, for each $f_c(\cdot)$ has Multi-layer perceptrons on top of the frozen feature extractor of the off-the-shelf retriever \mathcal{E} . The budget allcoators are trained on the cross-entropy loss, as we have already quantized the optimal budgets using the hyper-parameter δ . Note that if δ is high, the quantization is severe, otherwise the quantization is mild.

Require: Embedding model $\mathcal{E}(\cdot)$, LLM model $\mathcal{M}(\cdot)$, local datasets \mathcal{D}_c , budget allocator $f_c(\cdot)$ Buffering hyperparameter α .

Input: Test query x_q

1: Extract embedding $e_q = \mathcal{E}(x_q)$ 2: for $c \in [C]$ do

3:
$$\hat{k}_c = f_c(e_q)$$

4: Send e_q to all clients

5:
$$S_c = \mathcal{T}(e_q, \hat{k}_c + \alpha | \mathcal{D}_c)$$

- 6: return back $S_c \rightarrow$ Server
- 7: end for

8: $\begin{aligned} &S_{\text{agg}} = \bigcup_{c \in [C]} S_c \\ &9: S = \mathcal{T}(e_q, k | S_{\text{agg}}) \\ &10: s(S, x_q) = (x_1, y_1) \odot \dots \odot (x_k, y_k) \odot x_q \\ &11: \text{ Return: } y = \mathcal{M}(s(S, x_q)) \end{aligned}$



Figure 4: Overview of budget allocator: We train a budget allocator on top of the frozen encoder \mathcal{E} , which inherits from the retriever. During inference, when a test query x_q is provided, this module determines quantized budget levels for each client and allocates them accordingly.

4.2 Inference Using Budget Allocator

374

375

386

We derive the response to the test query x_q utilizing the LLM $\mathcal{M}(\cdot)$ through the described steps (see Algorithm 3 for specifics). We first extract the embedding $e_q = \mathcal{E}(x_q)$. Then, we compute the allocated budget $\{\hat{k}_c = f_c(e_q)\}_{c=1}^C$ and send \hat{k}_c to each client. Each client sends back top $\hat{k}_c + \alpha$ samples, *i.e.*, S_c , to the server. Note that we summarize how the budget allocator outputs \hat{k}_c in Figure 4. Here, α denotes the buffering hyper-parameter, which increases the chances for each client to be involved. After collecting $S_{agg} = \bigcup_{c \in [C]} S_c$, we aggregate them and run regular ICL.

5 Experiment

5.1 Experiment setup

First, we summarize the baselines, datasets, and the method for constructing non-IID settings. Finally, we depict the implementation details.

387

388

389

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

Baselines. We compare our method with various baselines, including Social Learning (Mohtashami et al., 2023), which does not account for non-IID, and other possible ways to handle distributed non-IID ICL, like Zero-shot, Proxy-only, Singleton, Uniform-budget, Random-budget, and ∞ -budget (oracle case). Notice that the proxy set in Proxy-Only baseline here contains the ground truth label information for each query in proxy set, which is different from proxy set used in our method. Detailed explanations are described in Appendix C.

Datasets. We check performance under 7 datasets – <u>Sentiment classification</u>: SST-5 (Socher et al., 2013), Amazon (McAuley and Leskovec, 2013), Yelp (Zhang et al., 2015), MR (Pang and Lee, 2005), <u>Topic classification</u>: Yahoo, AGNews (Zhang et al., 2015), and <u>Subjectivity classification</u>: Subj (Pang and Lee, 2004).

Dataset partition for non-IIDness. We split the training dataset into C subsets to ensure they follow a non-IID distribution. To achieve this, we partition the data based on class, following the splitting criteria outlined in (Li et al., 2022). Specifically, each client has access to only $\gamma < \Gamma$ classes, where Γ represents the total number of classes. We outline the summary of γ for each dataset in Appendix D. We also perform other non-IID partitions, including Dirichlet distribution non-IID, feature skew non-IID, and task shifting non-IID, with detailed description included in Appendix F.

Dataset paraphrasing. Due to concerns about sharing private samples between servers and clients, various techniques have been developed for NLP tasks. In this paper, we adopt paraphrasing technique used in (Mohtashami et al., 2023). Specifically, we utilize a small language model (Team et al., 2024) to generate paraphrased questions. In Appendix E, we summarize the instructions provided to the language model for rephrasing queries in the training dataset.

Implementation details. We implement our method and baselines based on OpenICL (Wu et al., 2023). We utilize pre-trained KATE re-triever (Liu et al., 2021). Note that they do not

Algorithm	Dataset				Ανσ			
ngonum	SST-5	Amazon	Yelp	MR	Yahoo	AGNews	Subj	1105
Zero-shot	29.19	24.70	31.23	73.95	25.87	67.60	50.55	43.30
Proxy-only	$40.64{\scriptstyle\pm}2.89$	$28.43 {\pm} 0.11$	$31.85{\pm}1.28$	$70.40{\pm}1.54$	$54.73 {\pm}~0.93$	$84.65 {\pm} 0.42$	$71.09{\pm}~1.34$	54.54
Singleton	$25.14{\scriptstyle\pm}4.18$	$24.03 {\pm}~0.57$	$29.44{\pm}3.91$	$50.00 {\pm}~0.00$	$38.14{\pm}2.03$	$50.60{\pm}0.66$	50.00 ± 0.00	38.19
Social Learning	$36.03 {\pm}~0.27$	$28.42{\pm}0.19$	$29.25 {\pm}~0.45$	$58.58 {\pm} 0.18$	$46.03 {\pm}~0.49$	$81.10{\pm}0.29$	$71.37 {\pm}~0.71$	50.11
Uniform-budget	32.94	25.63	26.60	33.65	43.00	73.17	63.20	42.60
Random-budget	$32.82{\pm}0.82$	$25.69 {\pm}~0.55$	$27.72 {\pm}~0.51$	$34.68 {\pm}~0.59$	$42.46 {\pm}~0.53$	$67.34{\pm}0.39$	$65.37 {\pm}~0.80$	42.30
∞ -budget	43.26	32.70	34.80	77.20	62.67	89.37	91.4	61.62
Ours	$\textbf{44.08}{\pm 0.12}$	$31.54{\pm}0.22$	$\textbf{35.48}{\pm 0.28}$	$\textbf{80.44}{\pm 0.67}$	$61.67 {\pm 0.25}$	$88.52{\pm0.30}$	$\textbf{82.36}{\pm 0.91}$	60.58

Table 1: Main results: To address the issue of non-IIDness in distributed ICL, we examined seven datasets and seven straightforward baselines. We run three random seeds and illustrate mean and std values. The top performance is highlighted in **bold** font, excluding the infinite budget scenario due to its impracticality. In summary, the proposed method effectively mitigates the non-iid distributed ICL problem to a reasonable extent.

overlap with the datasets used in our experiment. They used RoBERTa-large (Liu et al., 2019) encoder model. We use GPT-Neo-2.7B (Black et al., 2021) pre-trained model as answering LLMs as default. Hyper-parameters related to training budget allocators, α , and δ are described in Appendix D.

5.2 Main results

437

438

439

440

441

442

443

444 We have presented the performance of our method and baselines in Table 1. First, we can observe 445 that performance varies significantly depending on 446 the way the budget is allocated, which indicates 447 that the budget allocation scheme really matters in 448 non-IID ICL. Additionally, even when using only 449 the proxy dataset, there is a performance improve-450 ment, and this performance surpasses that of using 451 other clients which have the tilted local datasets 452 $(e.g., 29.19\% \rightarrow 40.64\% \text{ in SST-5 case})$. This in-453 454 dicates that utilizing a biased dataset can degrade the ICL performance. Although Social Learning 455 method has shown good performance in the pre-456 vious paper, it does not perform well under the 457 non-IID cases configured in this research. If we 458 can use an infinite budget, all settings would exhibit 459 high performance. However, our proposed method 460 demonstrates better performance than the infinite 461 budget upper limit (e.g., $34.86\% \rightarrow 35.48\%$ in 462 the Yelp case). This is likely due to a mechanism 463 that prevents unnecessary information from being 464 selected by the retriever with high importance. Ulti-465 mately, our method shows an average performance 466 467 improvement of 5.05% across seven datasets compared to the best performance of baselines using 468 the proxy dataset. Our method also shows outstand-469 ing performance under Dirichlet distribution non-470 IID, as shown in Table 14 of Appendix F. These 471

show that the proposed method can handle the non-IID case well. Furthermore, considering that our method does not necessitate the collection of ground truth label information for the proxy set and solely relies on task-related queries, it proves to be more practical in real-world applications compared to the leading baseline Proxy-only. Unlike the baseline, our method imposes fewer constraints on the proxy set, enhancing its practicability. 472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

5.3 Analysis

In this section, we further examine four key aspects: (1) privacy-preserving case analysis, which encompasses paraphrasing both training and testing queries, (2) sensitivity to hyper-parameters, (3) the performance of the trained budget allocator, and (4) the compatibility of the LLMs.

Paraphrasing results. Due to privacy concerns in the fundamental distributed system, we evaluate the performance of paraphrased datasets, with results detailed in Table 20 of Appendix F. Our method demonstrates superior performance compared to other baselines across multiple datasets. We used the exact same data settings as in Table 1. Specifically, performance on the Subj and SST-5 datasets is lower than without paraphrasing, while the Yelp dataset shows a slight improvement. Additionally, as consistent with Table 1, non-IIDness causes significant performance degradation for ICL methods, as seen by comparing Zero-shot with ICL-related methods (*e.g.*, 27.96% \rightarrow 25.31% in Singleton).

Hyper-parameter sensitivity. We examine the sensitivity of the hyper-parameters of our method. We have two hyper-parameters: δ , which is the resolution of the budget allocator; α , which represents the additional budget allocated to each client as a

buffer; and proxy size, which is the size of proxy 507 data for the budget allocator training. As illustrated 508 in Figure 11 of Appendix F, when we increase α , 509 the performance is improved while the budget ef-510 ficiency is reduced. On the other hand, when δ is high (or low), it has too sparse (or dense) repre-512 sentation of the budget class, thus performance is 513 degraded. Nevertheless, the performance is higher 514 than the other baselines in Table 1. For the sensitiv-515 ity of the size of proxy data, it is revealed that our 516 framework is not sensitive to how many proxy data 517 samples are used to train the budget allocator, as 518 shown in Figure 12 of Appendix F. This indicates 519 our method is stable even with limited proxy data on the server side. We have more fine-grained ex-521 periment result for relation between proxy size and the budget allocator resolution δ , which is shown 523 in Figure 7 of Appendix F.

Trained budget allocator. We assess whether the 525 trained budget allocator distributes budgets appropriately for each client. To evaluate efficiency, we examine the number of samples, *i.e.*, k_c communicated for all queries and plot histogram. As demonstrated in Figure 6 of Appendix F, we confirm that our method's forecasts exhibit nearly identical performance to the oracle budget when an additional 25% budget is allocated. Notice that it is necessary to assign $k \times C$ budgets to get a performance similar to the oracle case, without our method.

527

528

529

532

533

534 535

538

540

541

543

544

545

547

548

549

550

551

552

553

556

Other types of LLMs. We utilize various LLMs to assess the compatibility of our method. Specifically, we evaluate the SST-5 using different model sizes, including GPT-Neo-1.3B (Black et al., 2021), Llama-2-7B (Touvron et al., 2023a), and OpenAI gpt-3.5-turbo (OpenAI, 2022). As demonstrated in Table 18 of Appendix F, our method exhibits a plug-and-play capability and achieves reasonable performance improvements in ICL.

Distribution of proxy set. We also conduct experiments on when the proxy set has different distributions from the test set, aiming to verify the applicability of our method in broader real-world settings. The results and conclusion are included in Appendix F.5 due to the space limit.

Related Work 6

In-context learning. ICL (Dong et al., 2022) is one of the fastest paradigms using pre-trained LLMs by feeding several examples to construct the context to solve the given query. The main criteria of this research field are to find the most informative samples among the training datasets. (Liu et al., 2021) trains BERT (Devlin et al., 2018) oriented encoder and uses k nearest neighbors. (Rubin et al., 2022) proposed an efficient retriever called EPR. It trains two encoders by inheriting method of dense passage retriever (DPR) (Karpukhin et al., 2020) using loss of positive and negative pairs. To reduce domain specificity, (Li et al., 2023) proposed UDR, which is applicable to multiple domain tasks in a universal way and shows reasonable performance from a single retriever. PromptPG (Lu et al., 2022) utilized a reinforcement learning framework to train the retriever so that it can generate context to improve the answerability of LLMs. (Chang and Jia, 2022) trains linear regressors according to the example influence on the LLM prediction. (Xie et al., 2021) proposes to use implicit Bayesian inference to understand the ICL problem. Note that extensive research focuses on the centralized case rather than targeting distributed cases.

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

Distributed ICL. To the best of our knowledge, only a single study (Mohtashami et al., 2023) tries to address ICL in a distributed manner. However, this paper solely focuses on merging the distributed information without considering the nature of the non-identically distributed information. Many studies, such as those on federated learning (Li et al., 2021; Zhang et al., 2021; Mammen, 2021), address the non-IID distribution of datasets, highlighting the need to handle distributed non-IID ICL. Distributed ICL may resemble distributed RAG, yet the latter is more complex and requires further exploration, as discussed in Appendix H.

7 Conclusion

In this paper, we tackle the challenge of distributed non-IID ICL. Initially, we show that non-IID leads to performance degradation and discover that improper budget allocation causes significant drops in ICL. Inspired by the learnable pattern between budget values and query embeddings, we propose a method that learns budget assignment and employs it during inference to allocate appropriate budgets for each query. The proposed method achieves performance improvements across several benchmarks compared with various baselines. In addition, we examine the privacy-preserving version of our method using paraphrasing and show its efficacy. Last but not least, extensive sensitivity experiments show robustness of our method on hyperparameters and different LLMs.

Limitations. This research addresses in-context 607 learning with datasets distributed among clients 608 with non-iid data. One limitation of this study can 609 be described as follows: (1) Proxy dataset: In prac-610 tice, some datasets might lack a proxy dataset on the server side, posing a challenge for the proposed 612 algorithm. (2) The research assumes the use of a 613 pre-trained off-the-shelf retriever. However, this 614 retriever may fail if there is a significant difference 615 between the target task domain and the pre-trained 616 domain. This issue can be mitigated by employing the contrastive training mechanism suggested 618 in federated learning research (Seo et al., 2024), 619 as one effective approach for training retrievers is utilizing contrastive loss.

Ethical statement. As outlined in the main 622 manuscript, utilizing distributed knowledge raises privacy concerns. We address this by employing a 624 paraphrasing technique developed and frequently used in federated learning with pre-trained generative models, although it has not been thoroughly explored. We will make every effort to eliminate 628 privacy concerns and implement all possible mea-629 sures to prevent privacy information leakage when applying the paraphrasing method to the best of our 632 ability.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*. 633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

- Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. If you use this software, please cite it using these metadata.
- Ting-Yun Chang and Robin Jia. 2022. Data curation alone can stabilize in-context learning. *arXiv preprint arXiv:2212.10378*.
- Zicun Cong, Xuan Luo, Jian Pei, Feida Zhu, and Yong Zhang. 2022. Data pricing in machine learning pipelines. *Knowledge and Information Systems*, 64(6):1417–1455.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- 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*.
- Eduard Hovy, Laurie Gerber, Ulf Hermjakob, Chin-Yew Lin, and Deepak Ravichandran. 2001. Toward semantics-based answer pinpointing. In *Proceedings* of the first international conference on Human language technology research.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Itay Levy, Ben Bogin, and Jonathan Berant. 2022. Diverse demonstrations improve in-context compositional generalization. *arXiv preprint arXiv:2212.06800*.
- Jiaxing Li, Chi Xu, Lianchen Jia, Feng Wang, Cong Zhang, and Jiangchuan Liu. 2024. Eaco-rag: Edgeassisted and collaborative rag with adaptive knowledge update. *arXiv preprint arXiv:2410.20299*.
- Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. 2022. Federated learning on non-iid data silos: An experimental study. In 2022 IEEE 38th international conference on data engineering (ICDE), pages 965–978. IEEE.
- Qinbin Li, Zeyi Wen, Zhaomin Wu, Sixu Hu, Naibo Wang, Yuan Li, Xu Liu, and Bingsheng He. 2021. A survey on federated learning systems: Vision, hype and reality for data privacy and protection. *IEEE Transactions on Knowledge and Data Engineering*, 35(4):3347–3366.

- Xiaonan Li, Kai Lv, Hang Yan, Tianyang Lin, Wei Zhu, Yuan Ni, Guotong Xie, Xiaoling Wang, and Xipeng Qiu. 2023. Unified demonstration retriever for incontext learning. *arXiv preprint arXiv:2305.04320*.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2022. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. *arXiv preprint arXiv:2209.14610*.
- Priyanka Mary Mammen. 2021. Federated learning: Opportunities and challenges. *arXiv preprint arXiv:2101.05428*.
- Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of the 7th ACM conference on Recommender systems*, pages 165–172.
- Amirkeivan Mohtashami, Florian Hartmann, Sian Gooding, Lukas Zilka, Matt Sharifi, and 1 others. 2023. Social learning: Towards collaborative learning with large language models. *arXiv preprint arXiv:2312.11441*.
- OpenAI. 2022. Introducing chatgpt.

704 705

707

711

714

715

716

717

718

719

720

721

722

725

726

727

729

730

733

734

735

736

737

739

740

741

742

- Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings* of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04), pages 271–278, Barcelona, Spain.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *arXiv preprint cs/0506075*.
- Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389.
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2655–2671, Seattle, United States. Association for Computational Linguistics.
- Seonguk Seo, Jinkyu Kim, Geeho Kim, and Bohyung Han. 2024. Relaxed contrastive learning for federated learning. *arXiv preprint arXiv:2401.04928*.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642. 743

744

745

746

747

750

751

752

753

754

756

759

760

761

762

763

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

787

788

789

790

791

792

793

794

795

796

- Zuoqi Tang, Zheqi Lv, and Chao Wu. 2020. Abrief survey of data pricing for machine learning. In *CS* & *IT Conference Proceedings*, volume 10. CS & IT Conference Proceedings.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, and 1 others. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, and 1 others. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(11).
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- Shuai Wang, Ekaterina Khramtsova, Shengyao Zhuang, and Guido Zuccon. 2024. Feb4rag: Evaluating federated search in the context of retrieval augmented generation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 763–773.

Zhenyu Wu Wu, Yaoxiang Wang, Jiacheng Ye, Jiangtao Feng, Jingjing Xu, Yu Qiao, and Zhiyong Wu. 2023. Openicl: An open-source framework for in-context learning. *arXiv preprint arXiv:2303.02913*.

798

799 800

801

803

804

806

810

811

812

813

814

815

816

817

818 819

820

821

824

825

826 827

- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*.
- Jimin Xu, Nuanxin Hong, Zhening Xu, Zhou Zhao, Chao Wu, Kun Kuang, Jiaping Wang, Mingjie Zhu, Jingren Zhou, Kui Ren, and 1 others. 2023. Datadriven learning for data rights, data pricing, and privacy computing. *Engineering*, 25:66–76.
 - Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023. Compositional exemplars for in-context learning. In *International Conference on Machine Learning*, pages 39818–39833. PMLR.
 - Joo Hun Yoo, Hyejun Jeong, Jaehyeok Lee, and Tai-Myoung Chung. 2022. Open problems in medical federated learning. *International Journal of Web Information Systems*, 18(2/3):77–99.
 - Chen Zhang, Yu Xie, Hang Bai, Bin Yu, Weihong Li, and Yuan Gao. 2021. A survey on federated learning. *Knowledge-Based Systems*, 216:106775.
 - Kaiyi Zhang, Ang Lv, Yuhan Chen, Hansen Ha, Tao Xu, and Rui Yan. 2024. Batch-icl: Effective, efficient, and order-agnostic in-context learning. *arXiv preprint arXiv:2401.06469*.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

A How we construct non-IIDness

Following Li et al. (2022), we use class number based non-IID partition in our experiment. For a dataset with overall Γ classes, given hyperparameter class number γ on each client, we randomly assign γ classes from the overall Γ classes for each client. Assuming that $C_1 \leq C$ clients are assigned with a specific class, we equally partition samples of this class into C_1 parts and assign one part to each of C_1 clients. We denote this non-IID partition with the class number γ on each client as noniid-#label= γ .

B Motivation Experimental Settings

Non-IIDness performance drop experiment. For this experiment, we use SST5, Amazon, Yelp, Yahoo, and AGNews. And the non-iid settings are dsecribed in Table 2.

Dataset	ICE size k	#Clients	Partition
SST-5	32	8	noniid-#label=1
Amazon	16	8	noniid-#label=1
Yahoo	16	8	noniid-#label=2
AGNews	16	4	noniid-#label=1
Yelp	4	2	noniid-#label=3

Table 2: Experimental setup for obtaining the motivation.

t-SNE analysis of per-client budget experiment. For extracting t-SNE figure, we utilized the following experimental setting Table 3

Dataset	ICE size k	#Clients	Partition
SST-5	32	4	noniid-#label=2
Yelp	8	2	noniid-#label=3

Table 3: Experimental setup for obtaining the motivation.

C Baseline Details

Proxy-only. We randomly select samples from the original test set to construct the proxy set on the server side and use the remaining test set as the true test set. Notice that the proxy set in this baseline is different from the one we use in our pipeline. The proxy set used in Proxy-only contains label information for each query in the proxy set, while the proxy set in our algorithm does not require this information, which is more flexible in the real-world setting. When performing the ICL process, the server directly retrieves ICEs from the proxy set

rather than from the training set. For SST5, MR, and Subj, we randomly select 500 samples from the test set to be the proxy set. For Amazon, Yelp, Yahoo, and Agnews, we randomly select 750 samples from the test set to be the proxy set. Also, since the proxy set is already on the server side, there will be no privacy issues during communication between clients and the server. Thus, we don't generate samples to protect privacy and directly use the original samples in the proxy set for ICL.

Singleton. This baseline is for if the whole ICE set is constructed only using single client's local dataset. We randomly select one client from C clients, and perform local retrieval with $k_c = k$ budget. Then, the server uses this locally retrieved ICE set for LLM inference. We report the average accuracy over all clients.

Social learning. This algorithm (Mohtashami et al., 2023) is the first paper that considers the distributed ICL, but it only considers the IID setting. Since the authors didn't release the source code, we implemented it on our own. In our implementation, given server-side ICE number as k, each local client c performs local top- $\left\lceil \frac{k}{C} \right\rceil$ retrieval and sends retrieved ICEs to the server. The server then performs a random selection from k ICEs to construct an ICE set with k samples and feed this ICE set into LLM for inference.

Uniform-budget. We equally assign a local budget to each client. Assume the ICE number fed to server-side LLM for inference is k, then each client's local budget is $\lfloor \frac{k}{C} \rfloor$, where C is the number of clients. On server-side aggregation, we use reorder method as default.

Random-budget. We randomly assign a local budget to each client with the constraint that the overall local budget over C clients is k, where k is the ICE number fed to server-side LLM. On server-side aggregation, we use reorder method as default.

 ∞ -budget. The most inefficient way to do distributed non-IID ICL is to allow ∞ -budget on each client, that is, sending all samples to the server side. Then, the system performs centralized retrieval on the collected dataset to obtain top-*k* ICEs and feed them into LLM for inference.

D Experimental Setting

Dataset Explanation.In this study, we utilized910seven text classification tasks: four for sentiment911

914

931

933

934

935

936

937

939

941

943

analysis, two for topic classification, and one for subjectivity classification. The dataset statistics are presented in Table 4.

Dataset	Туре	Training	Test	Class
SST-5	Sentiment	8,534	2,210	5
Amazon	Sentiment	30,000	3,000	5
Yelp	Sentiment	30,000	3,000	5
MR	Sentiment	8,662	2,000	2
Yahoo	Topic	29,150	3,000	10
AGNews	Topic	29,914	3,000	4
Subj	Subjectivity	8,000	2,000	2

Table 4: The statistics of the datasets used.

915Given that the input instruction prompt can no-
tably influence performance, we detail the prompts916tably influence performance, we detail the prompts917used for each dataset in Table 24. It is in the last918page since prompts have long length. We follow919the prompt settings described in (Li et al., 2023)920and use the dataset uploaded by the paper's author,921available at https://huggingface.co/KaiLv.

922ICE number for LLM inference. Given an LLM,923different datasets show different preferences on the924choice of ICE number, *i.e.*, k, used in ICL inference925for better performance. For algorithms using ICL926(except Zero-shot), SST5, MR, and Subj use 32927ICEs for server-side LLM inference; Amazon uses9288 ICEs for server-side LLM inference; Yelp, Ya-929hoo, and Agnews use 4 ICEs for server-side LLM930inference.

Non-IID Setting. To keep similar non-IIDness levels across different datasets, we follow Table 5 as non-IID hyper-parameters for each dataset.

Hyper-parameters for our methods. For the main table results, the generated dataset results, and the different LLM architecture results, the hyperparameters are shown in Table 6, Table 7 and Table 8, respectively. For the training of the budget model, we use 800 epochs, with a learning rate range $\{0.01, 0.003\}$ and a batch size of 8.

Multi-layer perceptron for budget allocator. We use the three-layer perceptron on top of the encoder \mathcal{E} . The torch pseudo code is as follows:

class SMLP(nn.Module):

```
def __init__(self, width=300,
                       num_classes=10,
                        data_shape=(768,)):
        super().__init__()
        self.flat = nn.Flatten()
        self.l1
                   nn.Linear(np.prod(data_shape),
width)
        self.relu = nn.ReLU()
        self.12 = nn.Linear(width, width)
        self.13 = nn.Linear(width, num_classes)
    def forward(self, x):
        x = self.flat(x)
        x = self.ll(x)
        x = self.relu(x)
        x = self.12(x)
        x = self.relu(x)
        x = self.13(x)
        x = F.softmax(x)
        return x
```

Dataset	#Clients	Partition
SST-5	4	noniid-#label=2
Amazon	2	noniid-#label=3
Yelp	2	noniid-#label=3
MR	4	noniid-#label=1
Yahoo	2	noniid-#label=5
AGNews	2	noniid-#label=2
Subj	4	noniid-#label=1

Table 5: Non-IID setting

Dataset	ProxySetSize	δ	α	QuantRatio
SST-5	500	3	0	0.5
Amazon	750	2	0	0.5
Yelp	750	2	2	0.5
MR	500	3	0	0.5
Yahoo	750	2	2	0.5
AGNews	750	2	2	0.5
Subj	500	3	0	0.3

 Table 6: Hyper-parameters of our methods used in the main table

Dataset	ProxySetSize	δ	α	QuantRatio
SST-5	500	3	4	0.5
Yelp	750	3	0	0.5
Subj	500	3	0	0.3

Table 7: Hyper-parameters of our methods used for thegenerated query and training samples experiment

Computation Environment. We run our experiments on NVIDIA RTX A5000 GPU. Each experiment takes less than half an hour on a single GPU.

Model	ProxySetSize	δ	α	QuantRatio
GPT-Neo-1.3B	500	3	0	0.3
GPT-Neo-2.7B	500	3	0	0.3
Llama-2-7B	500	3	0	0.3
gpt-3.5-turbo	500	3	0	0.3

Table 8: Hyper-parameters of our methods used for different LLM architectures experiment on Subj

947

Ε

948

949

951

952

953

955

959

960

961

962

963

964

To generate the paraphrased query and response, we use the following instruction.

Generate paraphrased question

Please paraphrase the original sentence. Original sentence: {In-context example} Paraphrase sentence: {Paraphrased sentence }

Here is an example of input for the rephrasing LLM using the SST-5 dataset.

Please paraphrase the original sentence. Original sentence: "a stirring, funny and finally transporting re-imagining of beauty and the beast and 1930s horror films" Paraphrased sentence: A captivating, humorous, and ultimately uplifting reinterpretation of Beauty and the Beast combined with 1930s horror films. Please paraphrase the original sentence. Original sentence: "jonathan parker 's bartleby should have been the be-all-endall of the modern-office anomie films" Paraphrased sentence: Jonathan Parker's "Bartleby" had the potential to be the definitive film capturing the sense of alienation in modern office settings. Please paraphrase the original sentence. Original sentence: "a fan film that for the uninitiated plays better on video with the sound turned down" Paraphrased sentence: A fan film that, for those not familiar with the source material, is more enjoyable when watched with the sound turned off. Please paraphrase the original sentence. Original sentence: "apparently reassembled from the cutting-room floor of any given daytime soap" Paraphrased sentence: It appears to be pieced together from the outtakes of any given daytime soap opera. Please paraphrase the original sentence. Original sentence: " Paraphrased sentence:

Our paraphrased examples are summarized as follows.

Extra Experiments F

Per-client t-SNE visualization colored by **F.1** oracle budgets

Figure 5 presents the t-SNE visualization for per-client query embeddings, colored by oracle budget values.

F.2 Budget Analysis Results

Figure 6 presents the result of the budget analysis of our method.

F.3 Robustness on Proxy Size

Here, we present more detailed results on Subj with different proxy sizes over different values on budget allocator resolution δ in Figure 7.



Figure 5: t-SNE analysis of each client across two datasets. Each figure demonstrates that the budgets can be segregated by training a simple classifier, as they exhibit clustered subgroup pattern.



Figure 6: Analyze the total amount of budget allocated to clients under two datasets. Red and blue denote the oracle and 25% larger total budgets compared to oracle case.



Figure 7: Proxy size robustness over different budget allocator resolution δ . The results of Subj using GPT-Neo-2.7B.

F.4 More experiments for claim on Non-IID leads to ICL performance drop

To further support our claim that Non-IIDness leads to ICL performance drop under distributed setting, we introduce an additional feature-based Non-IID setting, where clients are split based on query length (number of words). For example, for Subj with 4 clients setting, client 0 only has samples with query length in range of [0, 13), client 1 within [13, 26), client 2 within [26, 31), and client 3 within $[31, \infty)$. We present the results for this setting in Table 12. By comparing the result of Uniform-budget (81.20%) and that of ∞ -budget/centralized

Original	Paraphrased
a turgid little history lesson, humourless	A dry and tedious history lesson that is
and dull	devoid of humour or interest.
not so much a movie as a picture book for	The movie is more of a picture book than a
the big screen .	full-fledged movie for the big screen.
now it 's just tired .	It is now simply outdated.

Original	Paraphrased
for all the wit and hoopla, festival in cannes	Festival in Cannes offers rare insight into
offers rare insight into the structure of rela-	the structure of relationships.
tionships .	
eldom has a movie so closely matched the	A movie seldom has a movie so closely
spirit of a man and his work .	matched the spirit of a man and his work.
those of you who are not an eighth grade	8th graders and younger will most likely
girl will most likely doze off during this	doze off during this one.
one.	

Table 10: Paraphrased examples of Subj dataset

setting (91.40%), we can still conclude that this kind of Non-IID setting also leads to performance drop, and our method helps to improve the performance under this Non-IID setting.

Algorithm

Zero-shot

Proxy-only Singleton

Social Learning

Uniform-budget

Random-budget

 ∞ -budget

Ours

setting.

Subj

50.55

71.09

76.22

79.60

81.20

81.75

91.40

82.80

Setting	Accuracy
Centralized	43.26
IID	44.52
Non-IID class-based	10.72
Non-IID query-length based	32.30

Table 13: Performance comparison under different settings for SST-5. Non-IID class-based: 8 clients distributed setting following Non-IID setting in Figure 3. Non-IID query-length based: 8 clients distributed setting under query length based Non-IID.

F.5 Non-Extreme Non-IID on Binary Classification Tasks

We conduct the experiment on Dirichlet distribution $\text{Dir}(\alpha)$ Non-IID partition on Subj and MR under the setting of 4 clients with $\alpha_{Dir} = 1.5$. The per-client sample distribution is shown in Figure 8, and the performance results are shown in Table 14.

Algorithm	Dat	aset
8	MR	Subj
Zero-shot	73.95	50.55
Proxy-only	70.40	71.09
Singleton	64.16	73.80
Social Learning	58.85	76.95
Uniform-budget	52.85	77.80
Random-budget	53.50	77.85
∞ -budget	77.20	91.40
Ours	75.53	82.80

Table 12: Performance on query length based Non-IID

We also include results under this query-length-based Non-IID setting for SST-5 with 8 clients, comparing IID vs. Non-IID performance following the same setting as in Figure 3. Again, performance degrades under Non-IID conditions, supporting our original claim.

Table 14: MR, Subj results under Dirichlet distribution Non-IID.

979

980 981

982

983



Figure 8: Per-client sample distribution under Dirichlet distribution Non-IID setting.

To further show the robustness of our method under Dirichlet distribution based Non-IID, we perform experiments on Subj for more α_{Dir} . As shown in Table 15, our method shows robustness even when $\alpha_{Dir} = 0.3$ (which we consider as an extreme Non-IID case), and beats other baselines (except the oracle ∞ -budget case). And as α_{Dir} increases (becoming more IID), the performance of our method also increases as expected.

Algorithm		α_{Dir}			
8	0.3	1.5	3.0		
Zero-shot		50.55			
Proxy-only		71.09			
Singleton	67.50	73.80	74.93		
Social Learning	67.10	76.95	78.95		
Uniform-budget	67.85	77.80	79.65		
Random-budget	67.80	77.85	80.00		
∞ -budget		91.40			
Ours	82.07	82.80	83.73		

Table 15: Results for Dirichlet distribution Non-IID under different α_{Dir} on Subj.

F.6 t-SNE under Non-IID with Task Shifting & Feature Skew

Dataset Amazon and Yelp are 5-class sentiment classification tasks with exactly the same label space, while different text query distributions. Based on this, we design a special Non-IID setting with task shifting & feature skew between clients: client 1 only contains 10,000 Amazon training samples, and client 2 only contains 10,000 Yelp training samples. Thus,



(b) Client 2 with only Amazon samples

Figure 9: t-SNE analysis on the test set consisting of both Yelp & Amazon samples. Data points are colored based on local oracle budget values.

we consider this special setting to be a task-shifting Non-IID. Also, since each client consists of samples from all classes of each task, we consider this setting as feature-skew Non-IID with class balance. We calculate the oracle budget values for a mixed test set consisting of 1,000 Yelp test samples and 1,000 Amazon test samples. Then we perform t-SNE analysis on sample embeddings of this mixed test set, colored using oracle local budget values. As shown in Figure 9, under Non-IID with task shifting & feature skew, there still exists clear clustering pattern between query embedding and oracle budget values. This indicates our method still can work with task shifting and feature skew.

1008

1010

1011

1012

1013

1015

1016

1017

1018

1019

1021

1022

1023

1024

1026

1027

1028

1029

1030

1031

1032

1033

F.7 Distribution Shift between Proxy Set and Test Set

It is critical to control the distribution shifting between proxy set and test set. We conduct experiments on two settings for proxy set distribution different from test set.

From the same dataset but different label distribution. The most simple case of "different distribution" can come from the different label distribution skew between proxy set and test set. For this setting, we experiment on Subj with a proxy set only containing samples of one class. As shown in the last line of Table 16, when the label skew exists, the performance of our method does decrease compared with the setting using the ideal proxy set (from 82.36% to 70.17%). However, it is still higher than some baselines, like zero-shot, singleton, uniform-budget, random-budget.

Similar task but different dataset.A more extreme1034case for proxy set different from test set can be, proxy set share1035

1036

- 1040 1041
- 1042 1043 1044 1045 1046
- 1047 1048

1049 1050

1051

same task with the test set, but are from different datasets. For this setting, we conduct the following experiment:

- use Amazon as proxy set for Yelp Non-IID setting (evaluate on Yelp test set)
- · use Yelp as proxy set for Amazon Non-IID setting (evaluate on Amazon test set)

Since Yelp and Amazon share the similar task, we can consider this setting as using available dataset with similar task with the test set to construct the proxy set. We present the result in the the last line in Table 17. It shows that for the Amazon setting, using Yelp as a proxy set, the performance drop of our method is slight, and our method still outperforms other baselines, except in the ideal case where we use Amazon samples as a proxy set. While for Yelp setting using Amazon as proxy set, our method surprisingly shows even better performance than the ideal case, where use Yelp as proxy set.

Algorithm	Subj
Zero-shot	50.55
Proxy-only	71.09
Singleton	50.00
Social Learning	71.37
Uniform-budget	63.20
Random-budget	65.37
∞ -budget	91.30
Ours	82.36
Ours-proxy-label-skew	70.17

Table 16: Comparison with proxy set with label skew compared to the test set. The last line is the performance for this setting.

Algorithm	Dataset		
rigoriumi	Amazon	Yelp	
Zero-shot	24.70	31.23	
Proxy-only	28.43	31.85	
Singleton	24.03	29.44	
Social Learning	28.42	29.25	
Uniform-budget	25.63	26.60	
Random-budget	25.69	27.72	
∞ -budget	32.70	34.80	
Ours	31.54	35.48	
Ours-diff-proxy	31.27	37.33	

Table 17: Comparison with using different dataset to construct proxy set for budget allocator training. The last line is the performance for this setting.

Results for other types of LLMs F.8

The results for other types of LLMs are presented in Table 18.

Algorithm	Architecture					
8	GPT-Neo-1.3B	GPT-Neo-2.7B	Llama-2-7B	gpt-3.5-turbo		
Zero-shot	51.30	50.55	49.10	57.57		
Proxy-only	$80.18 {\pm}~1.87$	$71.09{\pm}~1.34$	$88.13 {\pm}~0.74$	$88.44 {\pm} 0.69$		
Singleton	50.00 ± 0.00	$50.00 {\pm}~0.00$	$52.89{\pm}3.43$	$60.81{\pm}6.31$		
Social Learning	$68.55 {\pm} 0.64$	$71.37 {\pm}~0.71$	$88.82{\pm}0.50$	$87.53 {\pm}~0.46$		
Uniform-budget	44.40	63.20	54.00	81.23		
Random-budget	$43.68 {\pm}~0.80$	$65.37 {\pm}~0.80$	$55.60{\pm}0.41$	$81.47{\pm}~1.81$		
∞ -budget	92.05	91.40	92.30	92.23		
Ours	$85.73 {\pm 0.94}$	$\textbf{82.36}{\pm 0.91}$	$91.58 {\pm 0.14}$	$91.33 {\pm}~0.72$		

Table 18: Default non-IID setting of Subj using different LLMs. 32 ICEs for server LLM inference.

F.9 Hyperparameter sensitivity results

The results for the sensitivity of δ and α are presented in Figure 11. The results for the sensitivity of proxy set size are shown in Figure 12.

1054

1055

1056

1057

1058

1060

1061

1062



Figure 10: Additional budget α analysis. The orange dash line is the second-best baseline.



Figure 11: Budget allocator resolution δ analysis. The orange dash line is the second-best baseline.

100 -								
- 08	•			•			0)
60 -								
-	300	40	0 8 Prox	500 kv siz	600 200	0	70	0

Figure 12: Proxy size analysis. Results for Subj using GPT-Neo-2.7B.

We also present the impact of client numbers in Table 19, and our method stays robust even when the client number increases to 8.

Algorithm	Client Number				
	2	4	8		
Ours	81.07	82.36	82.40		

Table 19: Results for sensitivity on client numbers using Subj.

F.10 **Results for paraphrasing**

The results for paraphrasing-based privacy protection solution 1063 are presented in Table 20. 1064

Algorithm		Avσ		
ingonum	SST-5	Yelp	Subj	
Zero-shot	27.96	31.40	51.55	36.97
Proxy-only	$39.39{\pm}1.33$	$31.78{\pm}\ 1.75$	$73.46{\pm}\ 1.46$	48.21
Singleton	25.31± 3.89	$30.78{\scriptstyle\pm}~4.88$	$50.08 {\pm 0.10}$	35.39
Social Learning	$33.09 {\pm}~0.68$	$28.80{\pm}0.33$	$74.82{\pm}~0.93$	45.47
Uniform-budget	27.06	26.60	63.30	38.99
Random-budget	$27.29 {\pm}~0.51$	$27.70{\pm}~0.46$	$63.88 {\pm}~0.81$	39.62
∞ -budget	41.63	37.23	90.75	56.54
Ours	$\textbf{40.37}{\pm 0.27}$	$\textbf{36.52}{\pm 0.89}$	$\textbf{83.82}{\pm1.00}$	53.57

Table 20: Analysis of the generated query and training samples. We paraphrase the datasets using small-sized LLMs and conduct the experiments as in Table 1 under the same experimental settings.

F.11 Experiment for dataset with more classes

To show the efficacy of our method on datasets with more classes (more than 2), we add TREC (Hovy et al., 2001) (6 classes) under 4-client Non-IID setting. As shown in Table 21, our method outperforms other baselines except the upperbound ∞ -budget.

Algorithm	TREC
Zero-shot	31.40
Proxy-only	78.60
Singleton	28.00
Social Learning	83.20
Uniform-budget	69.20
Random-budget	64.60
∞ -budget	91.40
Ours	87.40

Table 21: Performance on TREC under 4-client Non-IID setting.

F.12 Hitting ratio comparison between the trained allocator and random budget allocator

We measure the accuracy of predicted budget level (our method) and the random budget level for two cases, to demonstrate the 'hitting rate' of budget assignment. As shown in Table 22 and Table 23, our method achieves a better hitting rate compared with the random budget.

	Client 1	Client 2	Client 3	Client 4
Our	79	60	71	78
Random	59	50	53	46

Table 22: Per-client hitting ratio comparison when $\delta = 2$ (per client budget value is 'low' or 'high'), the accuracy (hitting ratio) of different budget assignments.

	Client 1	Client 2	Client 3	Client 4
Our	60	74	60	68
Random	54	52	53	51

Table 23: Per-client hitting ratio comparison when $\delta = 3$ (per client budget value is 'low', 'medium' or 'high'), the accuracy (hitting ratio) of different budget assignments.

G Concrete Example of Distributed Non-IID ICL Scenario

A concrete example of distributed Non-IID ICL scenario can be the medical diagnosis task based on ICL cooperating with multiple medical institutions. Now, we have several medical institutions, with each institution owning some medical records (each sample consisting of the patient's symptoms description in text and the corresponding diagnosed disease, that is, the query x and label y). These medical institutions normally do not have enough local computation power to support LLM computation requiring large GPU resources, while they can do some small-cost local computation like retrieval processes to find similar queries. At the same time, there will be a platform operating like a server in this system, with enough computation resources to support LLM inference and in charge of cooperation management between these institutions. Once the system (including the server platform and cooperating institutions) is deployed, the platform can provide consulting diagnosis services to other patients, doctors, or even other medical institutions based on pay-by-use knowledge pricing strategies. That is, the price is decided by the number of samples involved in the whole diagnosis procedure. Also, due to medical privacy concerns, the server platform can use local samples to perform inference while not allow caching these samples. Thus, these local retrieved samples cannot be cached to construct a retrieval pool on server platform. For the specific example of platform that supports LLM, OpenAI now provides ChatGPT Enterprise³, which allows the deployment requirement that the platform should not cache and utilize private data for further training.

1081

1082

1083

1085

1086

1089

1090

1091

1092

1093

1094

1098

1100

1101

1102

1103

1104

1105

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

H Discussion on Relation with Distributed RAG

Here, we discuss the differences between our approach and existing distributed RAG studies to provide additional clarity and context for our contribution.

In developing this work, we carefully considered related studies in distributed RAG. However, the challenges addressed by existing distributed RAG works differ from those tackled in our paper. For instance, (Wang et al., 2024) focuses on the creation of datasets for distributed RAG frameworks and explores LLM-based labeling techniques for engineering pipelines. Their research scope and methodology are distinct from ours and are not directly applicable to our specific problem setting. Similarly, (Li et al., 2024) addresses resource consumption and real-time response challenges in distributed RAG, emphasizing local retrieval efficiency and answer accuracy. However, it does not account for the non-IID property in distributed settings. Additionally, (Li et al., 2024) permits LLM deployment on partial local institutions, which is fundamentally different from our setting.

Real-world distributed non-IID RAG scenarios present a more complex framework involving numerous challenges that must be addressed for effective deployment. For example:

1066 1067 1068

1069

1070

1071

1072

1073

1074

1075

1077

³https://openai.com/enterprise-privacy/

• How can we effectively decompose a user query into 1132 1133 subqueries while considering local knowledge distribu-1134 tion? 1135 • What is the best way to assign these subqueries to clients with varying local expertise? 1136 · How should we merge knowledge retrieved from mul-1137 1138 tiple clients with overlapping expertise, and should we assign confidence levels to different clients for the same 1139 1140 subqueries? · How can the local retrieval process be accelerated when 1141 1142 dealing with large local databases? 1143 These challenges represent broader avenues for exploration in distributed non-IID RAG. While our current work cannot 1144 1145 be directly compared with existing distributed RAG studies due to different settings, we believe it offers an interesting 1146 starting point for addressing such challenges. Specifically, 1147 1148 our approach focuses on how to enable cooperation among clients with varying distributions of knowledge. By assigning 1149 1150 preferences to clients based on their local knowledge distri-1151 butions and employing an MLP to learn these distributions 1152 without transmitting complete local knowledge to a central server, we offer an intuitive method that could inspire future 1153 1154 advancements in distributed non-IID RAG.

Original

Paraphrased

A friend of mine suggested we go here today before our movie. I was planning on suggesting another place, but she got their early and got a table DARN!I don't hate Red Robin...I think I avoid it because I am not a big fan of hamburgers. Seems like more of a place for straight guys and kids if you ask me, but my experience today wasn't to bad.Our waiter was really nice however I think that may have been a result of my push up bra.Ordered the Crispy Chicken Salad which also had hard boiled egg, bacon, and veggies and was very good. I'll probably get that again next time someone drags me here. Get it with that Italian dressing. Yummy! My friend ordered us onion rings as an appetizer I hate onions, but those were goood! Not my first choice but good crowd pleaser with more decent food than most chains.

Huge rooms. Soft towels. Comfortable bed with tons of pillows. Three things that make me happy when I'm traveling. What's wrong with this place, then? First of all, the location is central to nothing aside from the mall which houses several chain restaurants. Second, the staff could not care less about pretty much anything. While the woman who checked us in was helpful and friendly, that was the end of that. Despite requesting my Hilton status "amenities" (just a package of snacks and water...not much), I never received it. I was also charged \$30 for internet when I shouldn't have paid a thing (got the bill corrected at check out). And housekeeping didn't leave enough towels or replace the water. Considering we paid over 200 a night, the blase attitude of the staff is very disappointing.

You can either pay \$5.50 for 3-day movie rentals or get a \$40 membership and pay \$2 for 3-day/\$3 for 7-day rentals. They also offer 2-for-1 movies for students on Tuesday and Thursday. The reason I give 3 stars is because that deal isn't valid for people with memberships!? I learnt this 30 seconds after paying for a membership. I'm a student and could have paid \$2.75 for movies twice a week. Instead I paid \$40 (way too much) for a membership to pay .75 cents less than the other students. Good movie selection and shop but don't fall for their rip off of a membership unless you rent daily. I was considering going to Red Robin with a friend, but I decided to go somewhere else instead. I was not a big fan of hamburgers, but I was drawn to Red Robin because of the free tastings. I was pleasantly surprised by the friendliness of the staff, especially the waiter who was very attentive and helpful. I ordered the Crispy Chicken Salad, which was very good. I'll probably go back there next time someone drags me there.

3 things that make me happy when I travel are: spacious rooms, friendly staff, and a comfortable bed with lots of pillows. However, the location is inconvenient, the staff is indifferent, and the amenities are subpar.

3-day movie rentals cost either \$5.50 or \$40 for a membership. They offer a deal for students on Tuesdays and Thursdays, but it's not valid for those with memberships.

Dataset	Prompt	Label	Label Template	Example
AGNews	Topic of the text:	{World, Sports,	Topic: Label	REDMOND, Wash Microsoft Corp.
		Business, Technol-		and cable television provider Comcast
		ogy }		Corp. said Monday they would begin
				deploying set-top boxes powered by Mi-
				crosoft software starting next week. \n
				Topic: Business \\Oil demand is rising
				faster than predicted this year as OPEC
				pumps more low-quality oil in a failed
				bid to reduce record prices, according to
				International Energy Agency, an adviser
				to 26 industrialized nations. \n Topic:
MR	Sentiment of the	{great, terrible}	It was <i>Label</i>	"Analyze That" is one of those crass,
	sentence:			contrived sequels that not only fails on
				its own but makes you second-guess
				vour affection for the original. \n It
				was terribleabout the only thing to
				give the movie points for is bravado—to
				take an entirely stale concept and push
				it through the audience's meat grinder
				one more time \n It was
SST-5	Sentiment of the	forest good okay	It was Label	a strong funny and finally transporting
001 5	sentence:	had terrible}	It was Laber	re-imagining of Reauty and the Reast
	sentence.			and 1930s horror films \n It was great
				no movement no vuks not much of
				anything \n It was
Subi	Subjectivity of the	{subjective objec-	It's Lahel	gangs despite the gravity of its subject
Bubj	sentence.	tive}	It's Ember	matter is often as fun to watch as a
	sentence.			good spaghetti western \n It's subjec-
				tive smart and alert Thirteen Conver-
				sations About One Thing is a small gem
				\n If'e
Amazon	Sentiment of the	{great good okay	It was <i>Lahel</i>	Love the originality of this music Be-
7 muzon	sentence:	had terrible}	It was Laber	cause she is ever-changing Madonna is
	sentence.			never boring "Music" makes you want
				to dance - totally energizing! Wish the
				"Music" video was half as impressive
				as this work of art The case for the
				DVDs were a bit damaged. The damage
				did not compromise the DVDs \n It's
Veln	Subjectivity of the	{subjective objec-	It's I abol	My family visited ceasars palace and
Telp	sentence:	tivel	It's Luber	ate here. Our waiting time was only
	sentence.			ten minutes despite all the people. Our
				server was the best. He recommended
				several great dishes. The food was
				bigher than our expections his place
				is great. The staff is really friendly and
				the shile words humite is fontestic. You
				the chile verde burnto is fantastic. Tou
Vahaa	Topic of the ser	Society & Culture	It's Labol	who song message in a bottle? A nervor
1 allo0	tence	Science & Matha	n s Labei	Sting (the Police) Neuroscience Orec
		matice Health Edu		tion? Answer: no it is called mater new
		antion & Deferrer		roprostheses he Toric:
		Computers ^e Inter		roprosineses. In ropic:
		not Stream D		
		net, sports, Busi-		
		ness & Finance, En-		
		tertainment & Mu-		
		sic, Family & Rela-		
		tionships, Politics &		
		Government}		

Table 24: Prompt and instructions used for each dataset. We denote examples in blue and queries in red.