When Locally Deployable Small Models Can Substitute LLMs? An Empirical Study on Active Learning in Real-World Scenarios

Anonymous ACL submission

Abstract

Large Language Models (LLMs) excel at diverse benchmarking tasks, yet face many deployment barriers in real-world scenarios, such as data privacy and computing resources. On the other hand, low-resource learning techniques like Active Learning (AL) can reduce the annotation cost for fine-tuning locally deployable small models. Subsequently, when and how those AL-assisted small models with low-resource expert annotations can substitute off-the-shelf generic LLMs in real-world scenarios is critical but being overlooked. This empirical study examines AL-assisted small models versus generic LLMs in five real-world tasks with expert annotations. Our AL simulation validates the significance of AL-assisted locally deployable small models as well as the importance of distinct AL sampling strategies in realworld scenarios. We further discuss a promising future paradigm that leverages LLMs to "warm-up" AL-assisted small models.

1 Introduction

001

007

017

018

021

024

037

Recently, Large Language Models (LLMs) (Brown et al., 2020) have shown great capability in a variety of tasks. However, while several works allow more efficient training and deploying LLMs (Lester et al., 2021), there still exists a barrier to deploying LLMs in many real-world scenarios, such as computational cost (Brown et al., 2020), carbon footprint (Luccioni et al., 2023), and lack of domain knowledge (Xu et al., 2023). Also, in numerous real-world fields such as biomedical and legal, it is not viable to employ generic LLMs that are hosted by third-party companies (Plant et al., 2022). Instead, what consistently shows effectiveness in these domains is smaller, locally deployable models that are trained on expert annotations.

However, human experts are usually hard to access (Rasmussen et al., 2022), expensive to recruit (Wu et al., 2022), and often unwilling to work as

a labeler to annotate large-scale and high-quality datasets, especially in many domains that require extensive expertise (e.g., legal, clinical, and education) (Xu et al., 2022; Pappas et al., 2020). To bridge the gap between the scarcity of high-quality data and the huge data demand for model training, the research communities have widely explored low-resource learning techniques, such as transfer learning and meta-learning. Active Learning (AL, Settles (2009)) stands out due to its human-in-theloop design (Wu et al., 2022). Thus, a crutial question araises: when and how those AL-assisted small models with low-resource expert annotations can substitute generic LLMs in real-world scenarios? We hypothesize that AL-assisted small language models with a small number of expert annotations can reliably outperform generic LLMs. 041

042

043

044

045

047

049

052

053

055

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

079

This work presents an empirical study with AL simulations on five datasets from two real-world specialized domains (Biomedicine and Legal). We probe state-of-the-art (SOTA) LLMs (GPT-3.5 and GPT-4 (OpenAI, 2023)) with their best-performing prompting methodologies and compare them with an AL-assisted T5-base model (Raffel et al., 2020) leveraging two different AL strategies.

Our results show that the AL-assisted T5 model with hundreds of human annotations can consistently outperform GPT-3.5 and perform on par with GPT-4. Our results justify our hypothesis that locally deployable small models are irreplaceable in real-world domain-specific scenarios as well as the importance of selecting different AL sampling strategies. To better assist human experts' daily work, we also envision a hybrid paradigm to leverage LLMs to "warm-up" AL models.

2 Empirical Study Design

2.1 Datasets

We thoroughly examine existing expert-annotated datasets for specific real-world domains that re-



Figure 1: The sampling process of our data diversitybased strategy.

quire extensive expertise and choose BioMRC (Pappas et al., 2020), CUAD (Hendrycks et al., 2021), Unfair_tos (Lippi et al., 2019), ContractNLI (Koreeda and Manning, 2021) and Casehold (Zheng et al., 2021) for our evaluation. The datasets are in legal and biomedical domains and different types of tasks, including Multiple Choice, Classification, and Natural Language Inference (MacCartney and Manning, 2008). The dataset details are in Table 3.

2.2 Models

081

094

096

101

102

For the experiments with LLM, we utilize two SOTA generic LLMs: GPT-3.5 and GPT-4 (OpenAI, 2023). We probe the best-performing prompting strategy for each dataset with LLMs through extensive experiments on GPT-3.5 (reported in Appendix B) and apply the same settings for GPT-4.

We choose T5 (Raffel et al., 2020) as the backbone for AL because existing works demonstrate that T5 has strong performance for domain-specific fine-tuning (Yao et al., 2022; Mou et al., 2021). We initialize the T5 model with T5-base, a pre-trained weight that is trained on many general-domain downstream tasks.

2.3 Active Learning Strategies

Following the established taxonomies of AL strate-104 gies (Schröder and Niekler, 2020), we designed and implemented one data diversity-based strategy 106 and one model uncertainty-based strategy. The 107 diversity-based approach aims to identify the most 108 representative examples from the unlabeled data 110 space while maximizing the diversity, regardless of the model. On the other hand, the uncertainty-111 based approach attempts to locate examples that 112 the model is least confident about. We illustrate the 113 details of each strategy below and in Algorithm 1. 114

Data Diversity-Based Strategy. The process of data diversity-based sampling is illustrated in Figure 1. The objective of the data diversity-based strategy (Schröder et al., 2022) is to identify the most representative and diverse data. During the data pre-processing stage, we utilize Sentence-BERT (Wang et al., 2020) to embed each data content as a vector to prepare for the diversity-based AL sampling. For each iteration of the diversitybased AL sampling strategy, we 1) calculate the average cosine similarity score between each unused training data and all previously used training data, 2) sort the unused data by the average similarity score, and 3) select representative examples with the same interval from the sorted list to ensure diversity. For instance, in order to select 4 examples from 10 unused data, we select the 1st, 4th, 7th and 10th data from the ranked list after Step 2. This strategy design allows us to ensure the diversity and representativeness of selected examples.

115

116

117

118

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

150

151

152

154

155

156

157

159

160

161

163

Model Uncertainty-based Strategy. Model Uncertainty-Based Strategy (Sener and Savarese, 2018) aspires to identify samples the model is least confident about. Within each iteration, the model operates on the training data, computing the logits and locating the samples holding the minimal average probability on the highest-ranked tokens.

In addition to the aforementioned two types of AL strategies, we also include a random AL sampling baseline. For each iteration in the AL simulations, we follow a common practice of sampling 16 data samples with a specified strategy and then evaluate the model on the test split. Each AL setting was executed 10 times, and we report the mean and standard errors.

2.4 Evaluation Methods

We utilized the averaged F1 score for each label to evaluate Unfair_TOS to avoid the influence of unbalanced label distribution, which will also be discussed in Section 3.2. We evaluate the other datasets with average prediction accuracy. Detailed task instructions and experiment hyperparameters are shown in Appendix C and D.

3 Study Result

3.1 LLM vs. AL-Assisted Small Models

We plot the results on four legal domain datasets in Figure 2, and the results on BioMRC in Appendix A. The horizontal lines symbolize the best performance of GPT-3.5 and GPT-4, respectively.



Figure 2: AL simulation results. The horizontal line represents two close-domain LLMs' best performance. We report the mean value (line) and standard error (colored shaded area) over 10 trials. Each AL iteration comprises 16 examples. We can observe the T5-base with AL can reliably **outperform GPT-3.5** and reach a saturated performance that is **compatible with or even exceeds GPT-4** on all four datasets.

Unsurprisingly, all AL approaches suffer from 164 the "cold-starting" problem. However, on all four 165 datasets, the T5-base with AL can reliably outper-166 167 form GPT-3.5 and eventually reach a saturated performance that is compatible with or even ex-168 ceeds GPT-4, leveraging a total of several hundred data selected. For BioMRC, as shown in Figure 3, the T5-base can also consistently beat GPT-3.5 but 171 is saturated at a slightly lower performance com-172 pared to GPT-4. However, we believe GPT-4 might 173 have seen or been trained on most of these datasets 174 because they are publically available text corpora. 175 Regardless, our fine-tuned T5-base achieves com-176 patible performance with GPT-4 despite having 177 hundreds of times fewer parameters and requiring 178 significantly less computational power. 179

3.2 Further Analysis

180

181

182

183

185

186

189

Analysis of AL Strategies on Unfair_TOS. We observe the AL models in Unfair_TOS merely output "None" regardless of the input prior to the 20th iteration, but we can also observe clear advantage differences between AL strategies, where the uncertainty-based strategy can lead to better performance and saturate at higher results compared to the other settings.

The Unfair_TOS dataset consists of around 85%

| Strategy | Not-None Ratio | None Ratio |
|------------------|----------------|------------|
| Random | 0.1247 | 0.8752 |
| Diversity | 0.1255 | 0.8744 |
| Uncertainty | 0.1458 | 0.8541 |
| Complete dataset | 0.1252 | 0.8747 |

Table 1: Label distributions of complete dataset and data sampled by different AL strategies in Unfair_TOS. The ratio is calculated by dividing the corresponding data type by all data counts.

190

191

192

193

194

195

196

198

199

200

201

202

203

204

206

of data labeled None, and the rest of the data lies in eight other categories. We believe the AL model will be able to achieve a higher averaged F1 score if the AL strategy can select more Not-None data for the model to learn from. As a result, we calculate the label ratio for the original dataset and the data sampled by different AL strategies on the Unfair_TOS dataset, which can be found in Table 1. The ratio is calculated by dividing the corresponding data type by the count of all data. We sum the counts of all other eight data types and denote them as Not-None. We can observe the model uncertainty-based strategy selects significantly more Not-None labeled data than random (t(14) = -2.46, p < 0.05) and diversity (t(14) = -2.51, p < 0.05), which justifies the better performance of the uncertainty-based strategy.

Influence of Different Number of Few-Shot Examples. To establish a more solid evaluation, we conducted an additional experiment by evaluating GPT-4's performance when given different amounts of few-shot demonstrations. We used 1, 10, 50, and the maximum amount subject to the model input limit. If GPT-4 can only handle less than 50 examples, we omit the results for the 50 shots and report the max-shot results instead. To ensure reproducibility and control cost, we randomly sample 200 examples from the original test split and set the random seed to 42.

The result is reported in Table 5. We observe that generic LLM's (GPT-4) performance does not always increase when we add more and more data into the prompt, and with 10 shots can generally result in a saturated performance. Also, in three of the five datasets experimented, GPT-4 can only fit fewer than 20 few-shot examples in their context limit, justifying the need for small, fine-tuned models for domain-specific tasks.

3.3 Discussion

207

208

210

212

213

214

215

216

218

219

224

225

228

236

237

240

241

242

243

244

245

246

247

248

256

We can observe all AL strategies suffer from wellknown "cold-start" issues (Chen et al., 2022; Jin et al., 2022), where the model performs poorly at the early iterations due to potentially under-fitting issues as a result of the lack of enough labeled data. On the other hand, LLMs, specifically GPT-4 in our case, yield reasonably good performance despite eventually being surpassed by AL models fine-tuned on domain-specific datasets.

We envision a promising future paradigm in realworld domain-specific tasks of incorporating LLMs and AL fine-tuned smaller models in parallel. At first, the LLM's prediction will be presented to the human expert, and the collected annotation will be used to train the AL model. When the AL model begins to outperform LLM, the system will "switch" to present the AL model's prediction. Thus, the LLM's prediction can help overcome the "coldstarting" problem of AL, and the calibration ability (Zhu et al., 2023) can also be used to help identify the hard-to-answer or the wrong predicted examples in the sampling process to be annotated. And, the system can still benefit from AL's continually improving and up-to-date performance. The "switch" mechanism to efficiently and continually evaluate the AL model against the LLM will be a crucial component of such a paradigm, which will be investigated in our future work.

4 Related Works

Active Learning. Active Learning allows the model to iteratively sample a few annotations to be finetuned on (Shen et al., 2017; Ash et al., 2019; Teso and Kersting, 2019; Kasai et al., 2019; Zhang et al., 2022; Xiao et al., 2023). Thus, the key element of AL lies in the sampling strategy (Yao et al., 2023; Sharma et al., 2015). Many AL surveys categorize AL sampling strategies into three high-level underlying concepts – data diversity-based, model uncertainty-based, and hybrid strategies (Settles, 2009; Olsson, 2009; Fu et al., 2013; Schröder and Niekler, 2020; Ren et al., 2021).

Large Language Models LLMs (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023a,b) can learn task-solving from the context of few-shot examples and generate high-quality content without fine-tuning (Wei et al., 2021; Chung et al., 2022). Several recent works claim that LLMs can outperform human annotators for text classification (Gilardi et al., 2023), task evaluations (Chiang and Lee, 2023; Liu et al., 2023; Törnberg, 2023), and even in specialized domains (Nori et al., 2023). LLMAAA (Zhang et al., 2023) utilizes LLM as a weak annotator to train a small model with AL.

5 Conclusion

While LLMs such as GPT-4 have been endorsed to outperform smaller models in many benchmarking datasets, whether they can substitute smaller models, especially in real-world tasks and domains requiring extensive domain expertise, is important and debatable. In this work, we present an empirical study evaluating the performance between SOTA generic LLMs (GPT-3.5 and GPT-4) and a much smaller language model (T5-base) finetuned with different Active Learning strategies on five specialized datasets representing real-world domain-specific tasks.

Our evaluation demonstrates that smaller ALassisted models trained with expert annotation can consistently achieve or exceed best-performing LLMs with only a few hundred expert-annotated data, justifying that human experts remain indispensable in domain-specific tasks. Derived from our results, we posit a future paradigm that utilizes LLMs to overcome the "cold-start" issue of AL models as a "warm-up" strategy and eventually switch back to small models fine-tuned on domainspecific data once the latter outperforms LLMs. 260

261

262

263

264

265

266

267

268

269

270

271

277

278

279

280

281

284 285 286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

6 Limitations

306

307

311

312

313

314

315

317

319

320

321

322

325

326

330

331

332

334

338

347

351

355

This work primarily presents an empirical study of generic LLMs versus AL-assisted small language models fine-tuned on experts-annotated domainspecific data. Our experiment of AL-assisted models solely utilizes a T5-base model, where the performance of other models, such as BART (Lewis et al., 2019) and even LLMs that can be efficiently fine-tuned with Parameter-Efficient Fine-Tuning techniques (Mangrulkar et al., 2022; Hu et al., 2021; Lester et al., 2021), remains to be explored. This work only benchmarks two SOTA generic LLMs (GPT-3.5 and GPT-4). We are aware other LLMs exist that we do not include in this work, such as Mistral-7B (Jiang et al., 2023), Llama-2 (Touvron et al., 2023b), etc.

We only implemented and evaluated two fundamental types (data diversity-based and uncertaintybased) of Active Learning strategies in our work, and we are aware there exist other families of AL strategies that could extend our study, e.g., hybrid or ensemble approaches (Krogh and Vedelsby, 1994; Qian et al., 2020). Nevertheless, our empirical study with two fundamental Active Learning strategies justifies our primary statement that human experts are still needed in real-world domainspecific data annotation tasks.

Our evaluation comprises five datasets from two specialized real-world domains (legal and biomedical). We identify there are other domains and publically available domain-specific datasets, and we leave the analysis of the generalizability of our observations from this work to other domains and tasks as future work. In addition, we primarily engage in model comparisons through automated metrics. However, these may not necessarily provide an accurate representation of a model's performance. Also, an error analysis on which type of questions LLMs may excel or fail is also meaningful for future work. Therefore, human evaluation of these datasets, including human agreement and error analysis, might be needed for a more comprehensive assessment.

References

- Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. 2019. Deep batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint arXiv:1906.03671*.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind

Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, pages 1877–1901, Red Hook, NY, USA. Curran Associates Inc. 356

357

359

360

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

381

383

384

385

386

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

- Liangyu Chen, Yutong Bai, Siyu Huang, Yongyi Lu, Bihan Wen, Alan L. Yuille, and Zongwei Zhou. 2022. Making Your First Choice: To Address Cold Start Problem in Vision Active Learning.
- Cheng-Han Chiang and Hung-yi Lee. 2023. Can Large Language Models Be an Alternative to Human Evaluations?
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Yifan Fu, Xingquan Zhu, and Bin Li. 2013. A survey on instance selection for active learning. *Knowledge and information systems*, 35:249–283.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. ChatGPT Outperforms Crowd-Workers for Text-Annotation Tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120.
- Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. 2021. CUAD: An Expert-Annotated NLP Dataset for Legal Contract Review.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Qiuye Jin, Mingzhi Yuan, Shiman Li, Haoran Wang, Manning Wang, and Zhijian Song. 2022. Cold-start active learning for image classification. *Information Sciences*, 616:16–36.
- Jungo Kasai, Kun Qian, Sairam Gurajada, Yunyao Li, and Lucian Popa. 2019. Low-resource Deep Entity Resolution with Transfer and Active Learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5851– 5861, Florence, Italy. Association for Computational Linguistics.

- 411 412
- 413 414
- 415 416
- 417
- 418 419
- 420
- 421 422 423 424 425
- 426 427
- 428
- 429 430
- 431 432
- 433 434
- 435 436 437

438 439

- 440
- 441 442

443

444 445 446

447 448

449 450

> 451 452

453 454

455 456

457

462

463

464 465

- Yuta Koreeda and Christopher Manning. 2021. ContractNLI: A Dataset for Document-level Natural Language Inference for Contracts. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 1907–1919, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Anders Krogh and Jesper Vedelsby. 1994. Neural network ensembles, cross validation, and active learning. Advances in neural information processing systems, 7.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 3045-3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. arXiv preprint arXiv:1910.13461.
- Marco Lippi, Przemyslaw Palka, Giuseppe Contissa, Francesca Lagioia, Hans-Wolfgang Micklitz, Giovanni Sartor, and Paolo Torroni. 2019. CLAUDETTE: An Automated Detector of Potentially Unfair Clauses in Online Terms of Service. Artificial Intelligence and Law, 27(2):117–139.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment.
- Alexandra Sasha Luccioni, Sylvain Viguier, and Anne-Laure Ligozat. 2023. Estimating the Carbon Footprint of BLOOM, a 176B Parameter Language Model. Journal of Machine Learning Research, 24(253):1-15.
- Bill MacCartney and Christopher D. Manning. 2008. Modeling Semantic Containment and Exclusion in Natural Language Inference. In Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008), pages 521-528, Manchester, UK. Coling 2008 Organizing Committee.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. Peft: State-of-the-art parameterefficient fine-tuning methods. https://github. com/huggingface/peft.
- Xiangyang Mou, Chenghao Yang, Mo Yu, Bingsheng Yao, Xiaoxiao Guo, Saloni Potdar, and Hui Su. 2021. Narrative Question Answering with Cutting-Edge Open-Domain QA Techniques: A Comprehensive Study. Transactions of the Association for Computational Linguistics, 9:1032-1046.

Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. arXiv preprint arXiv:2311.16452.

466

467

468

469

470

471

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

507

508

509

510

511

512

513

514

515

516

517

518

519

520

Fredrik Olsson. 2009. A literature survey of active machine learning in the context of natural language processing.

OpenAI. 2023. GPT-4 Technical Report.

- Dimitris Pappas, Petros Stavropoulos, Ion Androutsopoulos, and Ryan McDonald. 2020. BioMRC: A Dataset for Biomedical Machine Reading Comprehension. In Proceedings of the 19th SIG-BioMed Workshop on Biomedical Language Processing, pages 140-149, Online. Association for Computational Linguistics.
- Richard Plant, Valerio Giuffrida, and Dimitra Gkatzia. 2022. You Are What You Write: Preserving Privacy in the Era of Large Language Models.
- Kun Qian, Poornima Chozhiyath Raman, Yunyao Li, and Lucian Popa. 2020. Learning structured representations of entity names using Active Learning and weak supervision. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6376-6383, Online. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. The Journal of Machine Learning Research, 21(1):5485-5551.
- Christoffer Bøgelund Rasmussen, Kristian Kirk, and Thomas B. Moeslund. 2022. The Challenge of Data Annotation in Deep Learning-A Case Study on Whole Plant Corn Silage. Sensors, 22(4):1596.
- Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. 2021. A survey of deep active learning. ACM computing surveys (CSUR), 54(9):1-40.
- Christopher Schröder and Andreas Niekler. 2020. A survey of active learning for text classification using deep neural networks. arXiv preprint arXiv:2008.07267.
- Christopher Schröder, Andreas Niekler, and Martin Potthast. 2022. Revisiting Uncertainty-based Query Strategies for Active Learning with Transformers. In Findings of the Association for Computational Linguistics: ACL 2022, pages 2194-2203, Dublin, Ireland. Association for Computational Linguistics.
- Ozan Sener and Silvio Savarese. 2018. Active Learning for Convolutional Neural Networks: A Core-Set Approach. In International Conference on Learning Representations.

633

634

635

Burr Settles. 2009. Active learning literature survey.

521

522

524

525

529

530

531

532

535

537

538

540 541

548

549

550

551

554

555

559

560

561

562

564

569

570

571

572

574

575

577

578

- Manali Sharma, Di Zhuang, and Mustafa Bilgic. 2015. Active Learning with Rationales for Text Classification. In Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 441–451, Denver, Colorado. Association for Computational Linguistics.
 - Yanyao Shen, Hyokun Yun, Zachary Lipton, Yakov Kronrod, and Animashree Anandkumar. 2017. Deep Active Learning for Named Entity Recognition. In Proceedings of the 2nd Workshop on Representation Learning for NLP, pages 252–256, Vancouver, Canada. Association for Computational Linguistics.
 - Stefano Teso and Kristian Kersting. 2019. Explanatory interactive machine learning. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 239–245.
 - Petter Törnberg. 2023. ChatGPT-4 Outperforms Experts and Crowd Workers in Annotating Political Twitter Messages with Zero-Shot Learning.
 - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. LLaMA: Open and Efficient Foundation Language Models.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open Foundation and Fine-Tuned Chat Models.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. MiniLM: Deep Self-Attention Distillation for Task-Agnostic Compression of Pre-Trained Transformers. In Advances in Neural Information Processing Systems, volume 33, pages 5776–5788. Curran Associates, Inc.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M.

Dai, and Quoc V. Le. 2021. Finetuned Language Models are Zero-Shot Learners. In *International Conference on Learning Representations*.

- Xingjiao Wu, Luwei Xiao, Yixuan Sun, Junhang Zhang, Tianlong Ma, and Liang He. 2022. A survey of human-in-the-loop for machine learning. *Future Generation Computer Systems*.
- Ruixuan Xiao, Yiwen Dong, Junbo Zhao, Runze Wu, Minmin Lin, Gang Chen, and Haobo Wang. 2023. FreeAL: Towards Human-Free Active Learning in the Era of Large Language Models. In *Proceedings* of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 14520–14535, Singapore. Association for Computational Linguistics.
- Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia Gabriel, Hong Yu, James Hendler, Marzyeh Ghassemi, Anind K. Dey, and Dakuo Wang. 2023. Mental-LLM: Leveraging Large Language Models for Mental Health Prediction via Online Text Data.
- Ying Xu, Dakuo Wang, Mo Yu, Daniel Ritchie, Bingsheng Yao, Tongshuang Wu, Zheng Zhang, Toby Jia-Jun Li, Nora Bradford, Branda Sun, Tran Bao Hoang, Yisi Sang, Yufang Hou, Xiaojuan Ma, Diyi Yang, Nanyun Peng, Zhou Yu, and Mark Warschauer. 2022. Fantastic Questions and Where to Find Them: FairytaleQA – An Authentic Dataset for Narrative Comprehension.
- Bingsheng Yao, Ishan Jindal, Lucian Popa, Yannis Katsis, Sayan Ghosh, Lihong He, Yuxuan Lu, Shashank Srivastava, Yunyao Li, James Hendler, and Dakuo Wang. 2023. Beyond Labels: Empowering Human Annotators with Natural Language Explanations through a Novel Active-Learning Architecture.
- Bingsheng Yao, Dakuo Wang, Tongshuang Wu, Zheng Zhang, Toby Li, Mo Yu, and Ying Xu. 2022. It is AI's Turn to Ask Humans a Question: Question-Answer Pair Generation for Children's Story Books. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 731–744, Dublin, Ireland. Association for Computational Linguistics.
- Ruoyu Zhang, Yanzeng Li, Yongliang Ma, Ming Zhou, and Lei Zou. 2023. LLMaAA: Making Large Language Models as Active Annotators. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 13088–13103, Singapore. Association for Computational Linguistics.
- Shujian Zhang, Chengyue Gong, Xingchao Liu, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. 2022. ALLSH: Active learning guided by local sensitivity and hardness. In *Findings of the Association* for Computational Linguistics: NAACL 2022, pages 1328–1342, Seattle, United States. Association for Computational Linguistics.
- Lucia Zheng, Neel Guha, Brandon R. Anderson, Peter Henderson, and Daniel E. Ho. 2021. When does

| 636 | pretraining help? assessing self-supervised learning |
|-----|---|
| 637 | for law and the CaseHOLD dataset of 53,000+ legal |
| 638 | holdings. In Proceedings of the Eighteenth Interna- |
| 639 | tional Conference on Artificial Intelligence and Law, |
| 640 | ICAIL '21, pages 159–168, New York, NY, USA. |
| 641 | Association for Computing Machinery. |

| 642 | Chiwei Zhu, Benfeng Xu, Quan Wang, Yongdong |
|-----|---|
| 643 | Zhang, and Zhendong Mao. 2023. On the Calibration |
| 644 | of Large Language Models and Alignment. |

A Empirical Study Result on BioMRC

645

646

647

651

655

672

673

676

For BioMRC, as shown in Figure 3, the T5-base with AL can quickly **outperform GPT-3.5** and eventually reach a saturated performance that is slightly lower than GPT-4. We posit that GPT-4 may have performed exceptionally well due to its exposure or training on BioMRC, given its source's public accessibility. Nevertheless, our refined T5base model demonstrates comparable performance to GPT-4. Remarkably, this is achieved despite the T5-base model's comparative parameter deficiency - in the hundreds of times less - and a significantly lower demand for computational resources.



Figure 3: Result on BioMRC

B LLM Prompting Experiments

The LLM prompting experiments can be found in Table 2. To obtain the SOTA performance, we experiment with GPT-3.5 under zero-shot and fewshot (1,3, 10 shots) to find the best-performing setting (bolded) for each dataset and execute GPT-4 with the same settings.

C Hyperparameters and Settings

We report the experiment hyperparameters in Table 4. All our experiments are executed on one of two resources: 1) four NVIDIA V100 32G graphic cards and 2) eight NVIDIA V100 32G graphic cards. For GPT-3.5 and GPT-4, we used GPT-3.5-0613 and GPT-4-0613 respectively.

For model uncertainty-based strategies, we calculate the model probability on a randomly sampled subset of the training data to reduce the time complexity of the model uncertainty-based data sampling process. Compared to the naive approach's $O(n^2)$ time complexity, our implementation remains to have a time complexity of O(n), which is the same as that of non-AL's (where *n* is the number of training data).

Algorithm 1 Active Learning Sampling Process

| 1: | function SELECT $(D_t, D_p, N, strategy)$ |
|-----|---|
| 2: | D_t : unlabeled data in the training split |
| 3: | D_p : previously selected data |
| 4: | N: number of data needed |
| 5: | strategy: Active Learning strategy |
| 6: | if $strategy =$ "similarity" then |
| 7: | $S \leftarrow \left(\frac{\sum_{d_p \in D_p} \cos(d_i, d_p)}{ D_p }\right)_{1 \le i \le D_t }$ |
| 8: | $id \leftarrow \operatorname{argsort}(S)$ |
| 9: | $step \leftarrow \frac{ D_t }{N}$ |
| 10: | $result \leftarrow (id_i)_{i \equiv 0 \pmod{step}, 1 \leq i \leq D_t }$ |
| 11: | return result, id – result |
| 12: | end if |
| 13: | ${f if}\ strategy =$ "uncertainty" then |
| 14: | $S \leftarrow (\text{Uncertainty}(d_i))_{1 \le i \le D_t }$ |
| 15: | $id \leftarrow \operatorname{argsort}(S)$ |
| 16: | return $id_{< N}$, $id_{\geq N}$ |
| 17: | end if |
| 18: | end function |

| Dataset | Learning Rate | Training Epoch |
|------------|---------------|----------------|
| BioMRC | 1e-4 | 20 |
| Unfair_TOS | 1.5e-4 | 12 |
| ContactNLI | 1.5e-4 | 20 |
| Casehold | 4e-5 | 28 |
| CUAD | 6e-5 | 18 |

Table 4: Hyperparameters for each dataset.

681

682

683

D Prompts Used for Each Dataset

Text in [[double brackets]] denotes input data.

D.1 BioMRC (Pappas et al., 2020)

I want you to act as an annotator for a

- $\, \hookrightarrow \,$ question answering system. You will
- $\, \hookrightarrow \,$ be given the title and abstract of a
- \hookrightarrow biomedical research paper, along
- $\, \hookrightarrow \,$ with a list of biomedical entities
- $\, \hookrightarrow \,$ mentioned in the abstract. Your task
- $\, \hookrightarrow \,$ is to determine which entity should
- \hookrightarrow replace the placeholder (XXXX) in
- ightarrow the title.

Here's how you should approach this \hookrightarrow task:

Carefully read the title and abstract of \hookrightarrow the paper.

| Detect | Matria | GPT-3.5 | | | | CDT 4 |
|-------------|----------|---------|---------|---------|----------|--------|
| Dataset | Metric - | 0 shot | 1 shots | 3 shots | 10 shots | OF 1-4 |
| CUAD | Accuracy | 0.6404 | 0.8048 | 0.8293 | 0.8178 | 0.8837 |
| BioMRC | Accuracy | 0.4067 | 0.5169 | 0.5040 | 0.4532 | 0.8259 |
| Unfair_tos | F1 | 0.4201 | 0.3847 | 0.3758 | 0.4206 | 0.4863 |
| ContractNLI | Accuracy | 0.4580 | 0.5990 | 0.5750 | 0.6420 | 0.8240 |
| Casehold | Accuracy | 0.3040 | 0.3020 | 0.3330 | 0.4010 | 0.6970 |

Table 2: Hyper-parameter tuning experiment results for GPT-3.5 and GPT-4.

| Dataset | Domain | Task | # Test Data |
|---|---------|----------------|-------------|
| BioMRC Pappas et al. (2020) | Biomed. | Multi-Choice | 6,250 |
| CUAD Hendrycks et al. (2021) | Law | Classification | 4,182 |
| Unfair_tos Lippi et al. (2019) | Law | Classification | 1,620 |
| ContractNLI Koreeda and Manning (2021) | Law | NLI | 1,991 |
| Casehold Zheng et al. (2021) | Law | Multi-Choice | 3,600 |

Table 3: Datasets involved in our empirical study.

| Dataset | 1-shot | 10-shot | 50-shot | max-shot (avg. # of shots) |
|-------------|--------|---------|---------|-------------------------------|
| BioMRC | 0.835 | 0.810 | - | 0.760 (13 shots) |
| Unfair tos | 0.441 | 0.488 | 0 567 | 0 563 (137 shots) |
| ContractNLI | 0.715 | 0.750 | | 0.740 (47 shots) |
| CUAD | 0.795 | 0.790 | | 0.82 (18 shots) |
| CaseHOLD | 0.660 | 0.790 | | 0.735 (19 shots) |

Table 5: GPT-4 result with different number of few-shot examples

| Pay close attention to the context in | If none of |
|---|---------------------------------|
| $_{\hookrightarrow}$ which the placeholder (XXXX) appears | \hookrightarrow you show |
| \hookrightarrow in the title. | |
| Review the list of biomedical entities | Here's how y |
| \hookrightarrow mentioned in the abstract. | \hookrightarrow task: |
| Determine which entity from the list | |
| $ \hookrightarrow $ best fits the context of the | Carefully re |
| \hookrightarrow placeholder in the title. | Review the I |
| Output only the identifier for the | For each un |
| $_{\hookrightarrow}$ chosen entity (e.g., `@entity1`). Do | $_{ m \leftrightarrow}$ it is p |
| \hookrightarrow not output anything else. | Output only |
| | \hookrightarrow present |
| <input/> : | \hookrightarrow multiple |
| <title>:</title> | You should d |
| [[TITLE]] | \hookrightarrow by a ser |
| <abstract>:</abstract> | Do not outpu |
| [[ABSTRACT]] | |
| <entities>:</entities> | <text>:</text> |
| [[ENTITY]] | [[TEXT]] |
| <output>:</output> | <output>:</output> |
| | |

D.2 UnfairTOS (Lippi et al., 2019)

684

I want you to act as an annotator for a \hookrightarrow Term of Service (ToS) review system. \hookrightarrow You will be given a piece of a Term \rightarrow of Service. Your job is to determine $\, \hookrightarrow \,$ whether the ToS contains any of the following unfair terms: \hookrightarrow Limitation of liability

Unilateral termination Unilateral change Content removal Contract by using Choice of law Jurisdiction Arbitration

the above terms are present, uld output "None".

you should approach this

ead the ToS. list of unfair terms. fair term, determine whether resent in the ToS. the unfair terms that are in the ToS. A ToS may have e unfair terms. ∖ output all of them, separated micolon (;). ut anything else.

```
686
```

D.3 ContractNLI (Koreeda and Manning, 2021)

I want you to act as an annotator for a

- \hookrightarrow question answering system. You will
- \hookrightarrow be given a contract and a hypothesis.
- \hookrightarrow Your task is to determine the
- \hookrightarrow hypothesis is contradictory,
- \hookrightarrow entailed or neutral to the contract.

Here's how you should approach this \hookrightarrow task:

Carefully read the contract. Carefully read the hypothesis. Determine whether the hypothesis is \Rightarrow contradictory, entailed or neutral \Rightarrow to the contract. Output only the label (contradiction, \Rightarrow entailment, neutral). Do not output \Rightarrow anything else.

<INPUT>: <premise>: [[PREMISE]] <hypothesis>: [[HYPOTHESIS]] <OUTPUT>:

D.4 CUAD (Hendrycks et al., 2021)

I want you to act as an annotator for a → question answering system. You will → be given the question and a piece of → a contract. You will need to answer → the question based on the contract. → There are only two possible answers, → "Yes" or "No". Here's how you should approach this

```
\rightarrow task:
```

Carefully read the question. Carefully read the contract. Determine the answer to the question is → true or not. Output only the exact answer (one of → "Yes" or "No") of the questions. Do → not output anything else. <INPUT>: <text>: [[TEXT]] <question>:

[[QUESTION]] <OUTPUT>:

D.5 Casehold (Zheng et al., 2021)

I want you to act as an annotator for a \rightarrow Question Answering system. You will \rightarrow be given the question and several \rightarrow answers. Your job is to determine \rightarrow which answer best answers the

 $\ \ \, \rightarrow \quad \text{question.}$

Here's how you should approach this \rightarrow task:

Carefully read the question. Carefully read the answers. Output the numeric index of the answers \rightarrow that best answers the question. Do not output anything else.

<INPUT>: <question>: [[QUESTION]] <answer>: [[ANSWER]] <OUTPUT>: