110

111

112

113

114

115

116

59

Anonymous Author(s)

Submission Id: 1

ABSTRACT

1

2 3

4

5

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

58

Environmental, Social, and Governance (ESG) has been crucial in investment decision-making in recent years, with an increase of ESG-centric research emerging. Concurrently, Natural Language Processing (NLP) has emerged in analyzing ESG-related texts. However, there is a lack of models and datasets specifically tailored for ESG categorization. This study presents a novel approach leveraging Pretrained Language Models (PLMs) and Large Language Models (LLMs) to tackle ESG text classification tasks. We introduce a pipeline for creating specialized datasets for ESG analysis by using keyword search and LLMs APIs to label data. Through the strategic extension of PLMs such as BERT, DistilRoBERTa, and RoBERTa, via continued pre-training on our datasets, our approach significantly surpasses traditional baseline performances. We also introduce ESGLlama and FinLlama, domain-specific models derived from Llama2, with FinLlama demonstrating exceptional efficacy in financial benchmarks and ESG text comprehensions. Final evaluations reveal that our models achieve significant advancements in ESG classification, outperforming established baselines.

CCS CONCEPTS

 Computing methodologies → Natural language processing; • Applied computing → Enterprise computing.

KEYWORDS

ESG, pre-trained language model, text classification

1 INTRODUCTION

Environmental, Social, and Governance (ESG) considerations represent the cornerstone of contemporary sustainable or responsible investment strategies. Over the past decade, ESG has become the preeminent framework for socially responsible investments and decision-making within the financial sector. However, a significant challenge remains relying on voluminous annual sustainability reports for informed decision-making. The comprehensive nature of these reports necessitates substantial effort for thorough analysis, highlighting the critical demand for automated solutions. In this context, Natural Language Processing (NLP) emerges as an indispensable tool, enabling navigating through extensive sustainability narratives and extracting pivotal ESG insights precisely. Machine learning (ML) and NLP technologies are emerging as a linchpin in refining ESG investment strategies in this evolving investment landscape. The synergy of ML and NLP does not merely expedite the analytical process. It also carves out novel corridors for academic and practical exploration into the development of systems adept at distilling pertinent insights from the vast seas of ESG reports.

Recent advancements in NLP have streamlined the identification and interpretation of ESG information, enabling real-time tracking of ESG dynamics and more nuanced analysis. This research 56 background sets the stage for exploring the integration of NLP in 57 enhancing the efficiency and depth of ESG analysis. Additionally, existing research has applied pre-trained language models in ESGrelated NLP tasks such as climate change-related text detection and controversy detection [12, 22, 28, 34]. Despite the notable advancements in applying PLMs for analyzing ESG-related data, a significant gap exists in the processing and collecting textual ESG data. This results in a scarcity of publicly accessible, high-quality ESG textual datasets, especially for established text categorization tasks within the ESG domain. Recent developments in large language models (LLMs) are more powerful than small PLMs and have demonstrated their potential in performing various NLP tasks like language understanding and generation. However, no such research focuses on using LLMs to solve ESG-related tasks. To address the shortfall of labeled ESG data, our approach employs a Large Language Model (LLM) as an annotator to get labeled ESG data and apply them in ESG-related tasks. This strategy aims to augment the availability of categorized ESG data and explores the potential of LLMs in ESG text classification tasks.

In this study, we tackle the significant gap in the availability of ESG-related datasets and apply our models to challenging ESG classification tasks. Utilizing a continuous pre-training strategy, we enhanced BERT, DistilRoBERTa, and RoBERTa models with a specifically curated ESG corpus, effectively tailoring them to the nuances of the ESG domain. We also leveraged keyword searches and APIs from large language models to meticulously annotate datasets for both 4-class and 9-class ESG classification. Further, we enriched our dataset collection with conversational history data, which proved crucial for Supervised Fine-Tuning (SFT) processes. This comprehensive fine-tuning involved both Pre-trained Language Models (PLMs) and Large Language Models (LLMs), significantly boosting their performance in ESG-specific tasks. Moreover, we developed two fine-tuend LLMs, ESGLlama and FinLlama, based on the Llama2, which demonstrated substantial improvements over baseline models. FinLlama also excelled in financial benchmarks, highlighting its utility in financial contexts. The empirical results from our study indicate a marked improvement in our models' ability to classify ESG-related content, consistently outperforming established benchmarks. In summary, our key contributions are the following:

- We propose a pipeline by utilizing keyword search and LLMs APIs to annotate data and construct three kinds of datasets for ESG analysis: pre-training corpus, classification dataset, ESG Supervised Fine-Tuning (SFT) dataset.
- We introduce three domain-specific PLMs: ESG-BERT, ESG-DistilRoBERTa, and ESG-RoBERTa. These models notably surpass their base models and our baseline.
- We conduct two fine-tuned Llama2 models: ESGLlama and FinLlama. FinLlama exhibits remarkable improvements in financial benchmarks.
- · We extensively compare PLMs and LLMs across various experimental settings, providing a comprehensive analysis of their performance.

Using Pre-trained Language Model for Accurate ESG Prediction

2 RELATED WORK

117

118

139

140

141

163

164

165

174

2.1 ESG Related NLP

119 The exploration of textual data in ESG reports has seen a marked 120 increase in interest, covering various research topics. Recent studies 121 have expanded beyond traditional analyses by adopting machine learning models to address societal issues such as stereotypes and 123 inclusivity [20]. Furthermore, diachronic distributional techniques 124 have been utilized to trace the evolution of ESG terminology, re-125 vealing shifts in discourse [24]. Traditional research often employs 126 keyword-based analysis methods [27], which lack contextual sensi-127 tivity [32]. Recent shifts toward context-aware machine learning 128 models have improved performance in diverse tasks such as climate 129 content classification [34], topic detection [32], Q&A systems [21], 130 and claim detection and verification [29]. However, this increase 131 in climate-focused research starkly contrasts with the minimal at-132 tention given to broader ESG aspects. Deploying fine-tuned BERT 133 models, especially those trained on extensive business and financial 134 news corpora like the Reuters News Archive, has effectively iden-135 tified ESG controversies [22]. Nevertheless, a significant research 136 gap remains in the comprehensive analysis of ESG communication 137 across all three domains. 138

2.2 Pre-trained Language Models

The advent of robust Pre-trained Language Models (PLMs) such 142 as BERT [9], ELMo [23], RoBERTa [17] has significantly boosted 143 NLP task performance across diverse domains. While domain-144 specific pre-training further augments their performance in spe-145 cialized fields [11], with dedicated models like BioBERT [15] for 146 biomedicine, ClinicalBERT [1] for clinical care, and SciBERT [5] for 147 scientific texts demonstrating targeted advancements. Additionally, 148 ClimateBERT [6] specifically addresses climate risk assessment. 149 The landscape of Large Language Models (LLMs) encompasses 150 models like T5 [25], which employs a unique Encoder-Decoder 151 Transformer structure, and the OpenAI GPT series, beginning with 152 GPT-3 [7], renowned for setting benchmarks in generative tasks. 153 Other notable GPT-style models include PaLM [8], and GPT-NeoX 154 [2], alongside GLM [10]. Despite many LLMs being proprietary, 155 open-source models like OPT [40] and LLaMA [31] foster extensive 156 research and practical applications. Despite these advances, the 157 application of PLMs in the nuanced ESG domain remains nascent, 158 representing a significant research opportunity to employ PLMs 159 and instruction-tuning techniques for nuanced and contextually 160 informed ESG text analysis. Our work seeks to bridge this gap, 161 leveraging PLMs to enhance ESG analysis and categorization. 162

2.3 Financial Language Models

The application of language models in finance is rapidly expand-166 ing, as these models are increasingly used for specialized functions 167 168 such as risk assessment and information extraction [16]. These financial language models are developed either from scratch or 169 through fine-tuning existing models. For instance, BloombergGPT 170 [36] was initially trained with a mix of general and finance-specific 171 172 datasets using BLOOM176B, while Xuan Yuan 2.0 [41] and Fin-T5 173 [19] focus on the Chinese financial market, leveraging specialized

175

176

177

178

179

pre-training. Fine-tuning for financial models predominantly targets sentiment analysis, news categorization, question-answering, summarization, and entity recognition. Noteworthy adaptations include FinBERT [3, 12, 18, 38]. Emerging models like PIXIU [37], and FinGPT [39] exemplify the advanced application of LLaMA architectures tailored for financial tasks, with PIXIU using 136K task-specific instructions and FinGPT employing LoRA for efficient fine-tuning. However, despite these advancements, the domain lacks specific models optimized for ESG-related tasks within finance, highlighting a significant opportunity for development. This gap underscores the potential for deploying fine-tuned LLMs to address ESG classification in finance, a promising area for future exploration and model innovation.

3 DATASETS

In response to the notable scarcity of datasets tailored for ESG domain analysis, we propose a pipeline, as illustrated in Figure 1, which encompasses data preprocessing, labeling procedures, and model training to enhance ESG data analysis capabilities systematically. Initially, data is sourced from various open sources and cleansed according to predefined rules. During the preprocessing phase, data is preliminarily categorized using keyword searches. Subsequent labeling employs APIs from LLMs to ensure high classification accuracy. Human evaluations are conducted to validate the labeled data, which then facilitates the construction of specialized datasets for further model pre-training and fine-tuning.

Specifically, we have constructed three types of datasets to enhance the accuracy of ESG prediction tasks: (1) Pre-training Dataset. This expansive corpus of ESG-related texts is designed to bolster the initial training of domain-specific models, thereby improving their ability to interpret ESG contexts accurately. (2) Classification Datasets. These datasets are segmented into four-class and nine-class categories for ESG texts, playing a pivotal role in the fine-tuning process to enhance model precision in ESG categorization. (3) SFT Dataset. Tailored for the Supervised Fine-Tuning (SFT) of Large Language Models (LLMs), this dataset incorporates conversational data generated by LLMs during the labeling procedure to boost the models' proficiency in ESG classification tasks.

3.1 Data Collection and Processing

For data collection, we searched and collected datasets mainly from two resources: huggingface ¹ and kaggle ². Refer to more details of our collected data in Appendix A. After data collection, we extract textual content pertinent to ESG analysis. In the initial data processing phase, we standardized the datasets to a **sentence-level** format, facilitating uniform analysis across diverse data sources. This step was crucial in preparing the data for subsequent machine learning and natural language processing tasks. Following the standardization, a data-cleaning procedure was implemented. This involved the removal of URLs and special characters from the text, ensuring that the datasets were devoid of extraneous information that could potentially skew the analysis. These preprocessing steps were essential in refining the data and enhancing the quality and reliability of the insights derived from our ESG subdomain language

²https://www.kaggle.com/

230

231

¹https://huggingface.co/



Figure 1: The work pipeline encompasses data collection, preprocessing, and labeling, followed by model training. Data is initially collected from open sources and cleansed. Using keyword searches and enhancing label accuracy through LLM's APIs, with further validation by human evaluation. The resultant dataset is used for pre-training and fine-tuning classification tasks.

models. The processed data amounted to approximately 18 million sentences, reflecting our dataset preparation efforts' comprehensive scope and scale.

3.2 Data Labeling

The data labeling phase was critical in constructing our ESG classification dataset, adhering to the four-class and nine-class categorization criteria defined by Huang et al. [12]. This phase involved two strategies: keyword search and labeling utilizing LLMs APIs.

Keyword Search. The keyword search initiates data identification across ESG subdomains, segregating text relevant to Environmental, Social, and Governance (ESG) areas and distinguishing Non-ESG content. This meticulous process enabled us to partition the corpus into distinct segments, each corresponding to a specific aspect of ESG, laying the groundwork for compiling domainspecific datasets. These datasets were then optimized for training models on targeted ESG categorization tasks, ensuring the relevance and specificity of the training material. While this method predominantly isolated relevant ESG-related text, it is essential to acknowledge that it might not entirely preclude the presence of Non-ESG data within these preliminary datasets. We argue that Non-ESG data within the pre-training phase could inadvertently enhance the model's robustness by exposing it to a broader spectrum of textual content. Details of keywords are in Appendix B.

After filtering the texts by keyword searching, we got the prelim-inary results shown in Table 8. To validate the effectiveness of our classification approach, these visualizations effectively confirm the appropriateness of the categorized data, with predominant terms such as "GHG emission" and "climate change" in the Environmental domain, "human rights" and "customer" in Social, and "director" and "financial statement" in Governance, reflecting the accurate representation of domain-specific high-frequency words. Our next objective was to refine the accuracy of our labeled data further. To achieve this, we planned to leverage LLMs for an additional layer of filtering and validation. Through this process, we aim to ensure



Figure 2: Representation of task decomposition and the task descriptions alongside exemplar responses from LLM

that our final datasets reflect the essential themes of each ESG domain and are also meticulously curated for subsequent analysis and model training. Details of visualizations are in Appendix C.

Labeling Data Using LLMs. Before labeling the data, we recognized a complexity gradient in categorization tasks, where tasks with fewer categories are inherently simpler than those with more. Studies such as Bang et al. (2023) [4] suggest that LLMs may underperform in specific, challenging downstream tasks, including multiclass classification tasks. To address this, we devised a structured approach to simplify the ESG classification challenge, as depicted in Figure 2. In this stage, the overall task is divided into three simpler tasks, where Task1 and Task2 comprise the four-class task (Env, Soc, Gov, Non-ESG), and an additional Task3 is required to construct the nine-class task. Specifically, the nine-class classification

involves three environmental categories (Climate Change, Natu-349 ral Capital, Pollution and Waste), three social categories (Human 350 Capital, Product Liability, Community Relations), two governance 351 categories (Corporate Governance, Business Ethics and Values), and 352 one Non-ESG category. The final three categories of the nine-class 353 354 task are unified into a single ternary (3-class) task, applying the 355 same categorization principles as the four-class task but with an added layer of specificity. Significantly, this ternary categorization 356 357 is based on data already classified under the four-class schema, 358 further refining the process.

For each sub-task, we employed APIs from three different LLMs: 359 Owen (gwen-max), GLM (glm-4), and GPT-3.5 (gpt-3.5-turbo-instruct). 360 This multi-model strategy was underpinned by several rationales: 361 Firstly, LLMs are prone to 'hallucination', often generating less 362 reliable outputs due to their randomness. Utilizing multiple models 363 helps mitigate significant data bias and enhances the diversity of 364 the labeled data. Secondly, the decision to leverage several LLMs' 365 APIs was economically driven, aiming to reduce costs associated 366 367 with extensive data filtering and labeling tasks. Lastly, employing multiple models concurrently significantly enhances the efficiency 368 of the data labeling process. Details regarding the prompt design 369 and an example of LLM response are in Appendix E. 370

3.3 Data Construction and Analysis

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

Pre-training Dataset. In constructing the pre-training dataset, we initially aggregated datasets categorized as Environmental (Env), Social (Soc), and Governance (Gov) based on keyword searches. Recognizing the challenges associated with processing excessively long texts, we implemented a filtration step to exclude these from the dataset. Long texts can detrimentally affect the efficiency of a compact language model since their extensive length can overwhelm the model's capacity to process and learn from them effectively. This limitation can lead to prolonged training times and potential overfitting on less representative data samples. Then, we executed a 90-10 split to segregate the data into training and evaluation subsets. The evaluation set is crucial in monitoring the training loss and establishing an early stop during the pre-training phase.

Classification Dataset. The development of the labeled clas-387 sification dataset involved multiple meticulous steps. Initially, we 388 processed the outputs from the Large Language Models (LLMs) used 389 for each classification task and subjected these to a rigorous human 390 review to verify the LLM-generated classifications. This review 391 process was crucial as it helped refine the data for the four-class 392 and nine-class categorizations, specifically excluding Non-ESG data 393 394 due to its inherent complexities and the limitations of LLM outputs, which may not always guarantee the absolute accuracy of 395 the responses. Consequently, the Non-ESG dataset was compiled 396 in a two-fold approach: approximately 8,500 samples were selected 397 from the LLM responses, and an additional 5,500 samples were 398 isolated following a keyword search, cumulatively amounting to 399 400 around 14,000 Non-ESG samples. A notable issue identified was the class imbalance within the nine-class dataset. To rectify this, we 401 implemented a normalization strategy by capping the maximum 402 number of instances per class at 3,000, leading to a more balanced 403 404 distribution. Furthermore, we applied stratified sampling for both 405 datasets to ensure equitable class representation. This approach 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

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

used a 70:15:15 split ratio for the four-class dataset and an 81:09:10 ratio for the nine-class dataset to create train-dev-test sets. Details of datasets distribution are in Appendix D.

Supervised Fine-Tuning Dataset. Supervised Fine Tuning (SFT) is a critical refinement process in Natural Language Processing (NLP), enhancing a large language model's adaptability to specific tasks. Unlike the general language focus of the pre-training phase, SFT introduces targeted supervision to adjust the model's weights by comparing its predictions against actual labels. This alignment improves the model's precision and adaptability for specific tasks, such as ESG text classification. In line with best practices like those demonstrated by the Alpaca model [30], which refined the LLaMA-7B into an instruction-following language model, we have constructed a Supervised Fine Tuning Dataset for ESG classification tasks with the following instructional categories:

- Identification of ESG-related text: "If the following text is ESG related data."
- (2) Four-Class classification: "Classify the following text into one of the four ESG categories: Environmental (Env), Social (Soc), Governance (Gov), or Non-ESG."
- (3) Nine-category Class: "Classify the following text into one of the nine ESG categories: Climate Change, Natural Capital, Pollution and Waste, Human Capital, Product Liability, Community Relations, Corporate Governance, Business Ethics and Values, or Non-ESG."

The dataset preparation involved reformatting existing four-class and nine-class datasets to align with these instructions, generating 95,412 data points. We also employed stratified sampling to select about 28,000 data points, ensuring diverse and balanced coverage across the instructions for effective SFT.

4 METHODOLOGY

4.1 Pre-trained Based Method

Baseline. Our baseline employs FinBERT [12], a model adapted from BERT for the financial sector. FinBERT is pre-trained on a corpus of financial documents, including annual filings and financial reports. Additionally, FinBERT has been extended to address ESGrelated classifications. The FinBERT-esg variant is fine-tuned to categorize texts into four broad ESG themes (E, S, G, or None). Meanwhile, the FinBERT-esg-9-categories model is fine-tuned to distinguish between nine detailed ESG topics.

Datasets. The dataset used for pre-training, detailed in Section 3.3, comprises 5,257,347 training sentences and 584,150 validation sentences, obtained via keyword search. While keyword searches are prone to including non-ESG phrases, resulting in false positives, this is beneficial for pre-training. It allows the model to learn the broader context of sustainability topics by exposing it to relevant and irrelevant samples.

Training Models. As detailed in Section 3.3, we utilized this dataset to pre-train models including BERT [9], DistilRoBERTa [26], and RoBERTa [17], leveraging their varying capacities—125 million parameters for RoBERTa and 85 million for DistilRoBERTa. Instead of starting from scratch, we engaged in Continual Pre-Training (CPT), a strategy that allows a model to assimilate new data

while preserving previously acquired knowledge. This approach is advantageous for adapting models to evolving data streams or new, unseen data. By continuing to pre-train on an established model's checkpoint, we infused domain-specific ESG knowledge into the models. This continual learning process is critical for domain adaptation tasks, such as enhancing a general model with domain-specific capabilities. The main experiments were conducted on 8 NVIDIA V100 Tensor Core GPUs. Consequently, we selected the model with the smallest validation loss as our final pretraining models: ESG-BERT, ESG-DistilRoBERTa, and ESG-RoBERTa. Details regarding pre-training process are in Appendix F.

4.2 LLM Based Method

465

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

514

515

517

518

519

Baseline. Llama2 [31], a generative text model with variants ranging from 7 to 70 billion parameters, excels in benchmarks, outperforming many open-source chat models. Selected versions of Llama2 were further optimized through five rounds of Reinforcement Learning from Human Feedback (RLHF), utilizing rejection sampling and Proximal Policy Optimization (PPO) to refine the reward model. These architectural improvements and their validated effectiveness in RLHF make Llama2 (Llama2-7b-chat-hf) an ideal baseline for our ESG classification task.

Datasets. Our LLM-based methods utilize two main types of datasets: the pre-training corpus and Supervised Fine-Tuning (SFT) datasets. The pre-training corpus has been substantially expanded to include not only the ESG-related texts discussed in Section 3.3 but also a significant volume of financial texts, primarily sourced from financial reports, totaling 5,282,943 sentences. For SFT, we employed two distinct datasets. The first SFT dataset, as introduced in Section 3.3, consists of conversational data generated during the labeling of ESG data using LLMs. The second SFT dataset is more extensive, integrating the conversational data and additional financial instruction tuning data as outlined in FinGPT [33] and the ESG_Chat dataset ³. The ESG_Chat dataset comprises dialogues between humans and LLMs, focusing on strategies to enhance ESG scores. Then, we adopted a targeted sampling strategy, producing a refined subset of 86,425 sentences.

Fine-tuning Models. To enhance the LLM's understanding of ESG-related themes, we enriched the model with ESG-related knowledge, resulting in the creation of two specialized models: ES-GLlama and FinLlama. ESGLlama underwent fine-tuning through Supervised Fine-Tuning (SFT) using conversational data tailored for ESG classification tasks, notably improving its accuracy within ESG contexts (as discussed in datasets, the first SFT dataset). Meanwhile, FinLlama was developed to tackle a broader spectrum of financial tasks, integrating extensive financial texts and targeted instructiontuning data, ranging from sentiment analysis to financial Question 513 Answering (QA). For fine-tuning FinLlama, we employed a twostage training approach. Initially, the Llama2 model underwent Continual Pre-Training (CPT) using a combined corpus of ESG-516 centric texts and additional financial documents, including financial news and annual reports. Subsequently, in the second stage, we conducted supervised fine-tuning on the model pre-trained in the initial phase using the second SFT dataset (as discussed in datasets). 520

521 522

Implementation. Due to LLMs' substantial parameter size and complex structure, fine-tuning and inference can be particularly time-intensive. To enhance efficiency, we employed Parameter-Efficient Fine-Tuning (PEFT) techniques such as Low-Rank Adaptation (LoRA) and freeze during SFT phases. Additionally, we utilized LLaMA-Factory [42] framework and vLLM [14] to accelerate pre-training SFT and inference processes. All experiments were conducted on NVIDIA V100 Tensor Core GPUs, with a learning rate set at 5×10^{-5} and a duration of 3 epochs for both pre-training and SFT phases. Training details are available in Appendix G.

EXPERIMENTS 5

5.1 Test on Public Dataset

To evaluate the generalizability of our trained models for ESGrelated tasks, we conducted tests using publicly available datasets: environmental_2k⁴, social_2k⁵ and governance_2k⁶ which are published by chatclimate.ai⁷ and derived from annual reports spanning 2017-2021. Each dataset is expertly annotated for binary classification, where '0' indicates "No" and '1' denotes "Yes" outcomes. We fine-tuned our models ESG-BERT, ESG-RoBERTa, and ESG-DistilRoBERTa on these datasets with a partitioning scheme of 64% training, 16% validation, and 20% testing. This meticulous approach allowed for optimal performance tuning, with the best models selected based on validation results for further testing.



Figure 3: Overall performance of models on public datasets

To evaluate the effectiveness of our models, Table 1 provides a detailed comparative analysis of key performance metrics-Precision (P), Recall (R), and F1 Score (F1)-across three critical domains: Environmental, Social, and Governance. The table juxtaposes the baseline models with enhanced versions that have undergone additional pre-training. Generally, the pre-trained models demonstrate superior performance compared to the baselines across the publicly accessible dataset. Notably, all pre-trained models consistently outperform their corresponding baseline models within the Social domain shown in Figure 3. Among them, ESG-DistilRoBERTa stands out with the highest precision (0.9415), recall (0.9449), and F1 score (0.9431), indicating robust performance. In the Environmental

579

580

523

524

³https://huggingface.co/datasets/zadhart/ESG_Chat

⁴https://huggingface.co/datasets/ESGBERT/environmental_2k

⁵https://huggingface.co/datasets/ESGBERT/social_2k

⁶https://huggingface.co/datasets/ESGBERT/governance_2k

⁷https://huggingface.co/ESGBERT

Conference acronym 'KiL, August 25, 2024, Barcelona, Spain

| | Environmental | | Social | | Governance | | | | |
|-------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Model | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 |
| BERT | 0.9207 | 0.9285 | 0.9244 | 0.8960 | 0.8899 | 0.8927 | 0.8048 | 0.8168 | 0.8104 |
| ESG-BERT | 0.9300 | 0.9284 | 0.9292 | 0.9354 | 0.9345 | 0.935 | 0.8141 | 0.8085 | 0.8112 |
| DistilRoBERTa | 0.9340 | 0.9436 | 0.9385 | 0.9035 | 0.9044 | 0.9039 | 0.8404 | 0.8444 | 0.8424 |
| ESG-DistilRoBERTa | 0.9364 | 0.9397 | 0.9380 | 0.9415 | 0.9449 | 0.9431 | 0.8252 | 0.8271 | 0.8261 |
| RoBERTa | 0.9279 | 0.9246 | 0.9262 | 0.9041 | 0.9135 | 0.9076 | 0.8292 | 0.8421 | 0.8352 |
| ESG-RoBERTa | 0.9340 | 0.9436 | 0.9385 | 0.9311 | 0.9345 | 0.9327 | 0.8048 | 0.7976 | 0.8011 |

Table 1: Performance metrics across environmental, social, and governance domains on public datasets. Bold shows the best results among baseline and corresponding pre-trained model, and <u>underlined</u> indicates the best results in each column.

domain, ESG-RoBERTa shows remarkable precision (0.9436) and an equivalent F1 score, underscoring its effectiveness. **Table 3: Nine-Class Evaluation Results of PLMs**

However, the Governance domain exhibits a contrasting scenario, with mixed results despite pre-training enhancements. The baseline DistilRoBERTa model outperforms the pre-trained versions in this domain, achieving the highest metrics with a precision of 0.8404, recall of 0.8444, and F1 score of 0.8424. This discrepancy suggests that while pre-training generally enhances model capabilities, its impact is less pronounced in the Governance domain. The observed variance may stem from misalignments between the pre-training content and the specifics of the publicly available governance data, suggesting a need to refine the fine-tuning parameters better to tailor the models to this domain's nuances.

5.2 Test on Classification Datasets

5.2.1 Evaluate PLMs. We fine-tuned our pre-trained models on ESG classification tasks (four-class and nine-class) using our constructed classification data. The training parameters were standardized at a batch size of 32 across 50 epochs while learning rates were adjusted based on model and task specifics. For the four-class classification, the learning rates were set at 3e-6 for the BERT model and 1.25e-6 for both DistilRoBERTa and RoBERTa. For the nineclass task, BERT was fine-tuned at 3e-6, DistilRoBERTa at 1.75e-6, and RoBERTa at 1.15e-6. These rates were meticulously selected to optimize each model's performance on its respective task. An *early stopping* mechanism was implemented during fine-tuning to curb overfitting and enhance computational efficiency. The models chosen for further utilization demonstrated the best performance on the validation set across the 50 epochs, specifically those achieving the lowest validation loss.

| Model | Р | R | F1 | Acc |
|--------------------|--------|--------|--------|--------|
| FinBERT | 0.7357 | 0.7150 | 0.7165 | 0.7222 |
| BERT | 0.8668 | 0.8658 | 0.8641 | 0.8667 |
| DistillRoBERTa | 0.8672 | 0.8687 | 0.8662 | 0.8684 |
| RoBERTa | 0.8610 | 0.8596 | 0.8582 | 0.8602 |
| ESG-BERT | 0.9074 | 0.9077 | 0.9071 | 0.9083 |
| ESG-DistillRoBERTa | 0.9027 | 0.9040 | 0.9014 | 0.9034 |
| ESG-RoBERTa | 0.9086 | 0.9100 | 0.9086 | 0.9102 |

| Model | Р | R | F1 | Acc |
|--------------------|---------------|---------------|---------------|---------------|
| FinBERT | 0.7160 | 0.7154 | 0.7081 | 0.7273 |
| BERT | 0.8393 | 0.8357 | 0.8361 | 0.8419 |
| DistillRoBERTa | 0.8240 | 0.8153 | 0.8179 | 0.8239 |
| RoBERTa | 0.8187 | 0.8196 | 0.8174 | 0.8275 |
| ESG-BERT | 0.8606 | 0.8637 | 0.8617 | 0.8693 |
| ESG-DistillRoBERTa | 0.8575 | 0.8552 | 0.8556 | 0.8616 |
| ESG-RoBERTa | 0.8611 | 0.8591 | 0.8592 | 0.8662 |

To assess the effectiveness of our pretrained models, we conducted tests on two sets: a four-class and a nine-class classification task, with results detailed in Table 2 and Table 3, respectively. The evaluations included baseline models, our specifically pre-trained models, and their base models. For the four-class task, ESG-RoBERTa excelled, achieving the highest metrics with a precision of 0.9086, a recall of 0.9100, an F1 score of 0.9086, and an accuracy of 0.9102, significantly surpassing the baseline finbert-esg model, which only reached an accuracy of 0.7222. This demonstrates a clear superiority over the baseline, with even the base models outperforming finbert-esg when fine-tuned. In the nineclass task, ESG-BERT led with the highest recall of 0.8637 and an F1 score of 0.8617, while ESG-RoBERTa achieved the top accuracy of 0.8662. These results highlight the advantages of our ESG-specific pretraining and fine-tuning strategy, markedly improving upon the performance of the baseline finbert-esg-9-categories model.

5.2.2 Evaluate LLMs. We will evaluate the performance of the baseline and our fine-tuned models across six different *experimental settings*: Zero-Shot, One-Shot, In-Context Learning (ICL), Zero-Shot with Chain of Thought (CoT) [13], One-Shot with CoT, and ICL with CoT. The dataset used for SFT in ESG text classification was constructed from ESG SFT data, as detailed in Section 3.3. It was refined by selecting only the classification data and simplifying the format to retain the text and label without additional explanations. More details about the ESG classification SFT dataset can be found in Appendix H. To process the results from our models, particularly the baseline, we utilized a *regular expression* matching technique to extract predicted labels from model outputs. The regular expression patterns provide a flexible and effective method to handle the

diverse outputs from the LLMs, ensuring alignment with our predefined SFT data and system prompt formats. Details regarding classification prompts design are in Appendix I.



Figure 4: Four-Class Precision of LLMs



Figure 5: Nour-Class Precisions of LLMs

For four-class classification, we assessed our models, ESGLlama and FinLlama, using Precision as the primary performance metric. Analysis of precision scores in Figure 4 shows that both models consistently outperform the baseline across most experimental settings. Notably, even the baseline model improves significantly when subjected to SFT with our tailored ESG classification dataset. Interestingly, the Freeze fine-tuning method generally surpasses the LoRA approach, except in zero-shot scenarios where LoRA excels, possibly indicating its tendency to overfit slightly. This overfitting suggests that external examples, absent from the training data, might disrupt LoRA's inference, while the Freeze method maintains better generalization and reasoning capabilities. The integration of CoT prompts typically reduces performance in zero-shot and one-shot settings, except for ICL tasks. This reduction may stem from CoT's incompatibility with classification tasks, which require straightforward decision-making rather than stepwise logic processing. However, incorporating demonstration examples in ICL tasks enhances the model's grasp of classification logic, significantly improving outcomes in ICL-CoT settings by providing richer con-text and sample diversity. Furthermore, FinLlama achieves superior precision over ESGLlama with the addition of CoT.

In the nine-class classification, the increase in category complexity and diversity presents more significant challenges, as indicated by lower overall performance metrics than in the four-class scenario. This trend highlights the difficulty in distinguishing among a more substantial number of classes. Performance visualizations in Figure 5 show that both ESGLlama and FinLlama substantially outperform the baseline across most configurations, affirming the enhanced capability of our fine-tuned models in handling ESGrelated texts. FinLlama excels in ICL, mainly when provided with ample examples, showcasing its deep understanding of the financial domain. Conversely, the performance notably drops in one-shot learning scenarios, where providing a single instance per class introduces significant bias and variability, impairing the model's accuracy. However, increasing the number of examples markedly improves performance, underscoring the benefits of more extensive training datasets. The comparison between LoRA and Freeze methods reveals that LoRA outperforms Freeze in one-shot settings, suggesting that LoRA's parameter adjustments are better suited for absorbing limited class-specific information efficiently. Additional analyses are in Appendix J.

5.3 Test on Financial Benchmark

To assess the FinLlama model's performance in financial NLP tasks, we evaluate it on FinGPT benchmark [33]. Our evaluation concentrated on two critical tasks: financial text sentiment analysis and headline classification, utilizing the fingpt-headline dataset ⁸.

This comprehensive benchmarking demonstrates FinLlama's robust capabilities in understanding and classifying financial texts, highlighting its utility in sentiment analysis and headline categorization. Results, presented in Table 4, clearly show that FinLlama significantly outperforms the baseline Llama2 model across these tasks. This superior performance across financial sentiment analysis and headline classification tasks validates the effectiveness of our targeted pre-training and Supervised Fine-Tuning (SFT) strategy. By incorporating domain-specific knowledge, FinLlama has shown notable improvements in analyzing financial texts, confirming its advanced proficiency in financial NLP.

Table 4: Performance of models on Financial Benchmarks

| | Llar | Llama2 | | lama |
|----------|--------|--------|--------|--------|
| Dataset | Acc | F1 | Acc | F1 |
| FPB | 0.4703 | 0.4140 | 0.7855 | 0.7838 |
| FiQA | 0.7964 | 0.7744 | 0.7782 | 0.8096 |
| TFNS | 0.3811 | 0.3037 | 0.8405 | 0.8408 |
| NWGI | 0.5656 | 0.4833 | 0.6501 | 0.6445 |
| Headline | 0.4314 | 0.6182 | 0.8783 | 0.6975 |

6 RESULTS ANALYSIS

Performance of Pre-trained Models. Our analysis highlighted that classification task complexity increases with the number of categories. This was evident from the lower convergence rates in the nine-class task compared to the four-class task. ESG-RoBERTa

⁸https://huggingface.co/datasets/FinGPT/fingpt-headline

excelled in the four-class task due to its larger parameter set, which 813 enhances its text understanding capabilities. In contrast, ESG-BERT 814 815 performed better in the nine-class task, suggesting that its pretraining objectives and architecture might offer superior generalization 816 across more diverse categories. Performance evaluations on a pub-817 licly available dataset confirmed the effectiveness of our pre-trained 818 models, as shown in Table 1. Particularly in the Social domain, mod-819 els like ESG-DistilRoBERTa demonstrated exceptional precision, 820 821 recall, and F1 scores, reflecting the quality of our pretraining and 822 the model's ability to generalize well. The extensive testing on a public dataset validated our pretraining dataset's quality and 823 demonstrated our models' improved comprehension of ESG-related 824 content, enhancing classification accuracy. Furthermore, the per-825 formance variation in the Governance domain highlights the need 826 for additional optimizations in data curation and model training 827 strategies to ensure robust model efficacy across diverse domains. 828

829 Performance of Large Models. Both ESGLlama and FinLlama 830 consistently outperform the baseline across most testing scenarios, 831 with notable improvements in the baseline model following SFT 832 with our ESG classification dataset. This enhancement highlights 833 the dataset's quality and the effectiveness of SFT. A distinct observa-834 tion is Freeze is generally better than LoRA because the trend of line 835 changes in its results is consistent with those of other experimental 836 setups, and more examples can improve its results. The integration 837 of CoT typically reduces performance in zero-shot and one-shot 838 scenarios. Still, it improves outcomes in ICL tasks due to additional 839 context and examples provided. Transitioning to a nine-class frame-840 work increases task complexity, generally lowering performance 841 metrics. In ICL tasks, FinLlama shows superior proficiency, particu-842 larly when additional samples are included, reflecting its adeptness 843 at navigating complex classification landscapes. Conversely, per-844 formance drops in one-shot scenarios, underscoring the challenges 845 of minimal data learning. Our experimental results illustrate the 846 models' capabilities in complex settings and superior performance, 847 particularly in ICL with CoT configurations. Moreover, testing Fin-848 Llama on financial benchmark further validates its superiority in 849 financial NLP, highlighting FinLlama's effectiveness in financial dis-850 course analysis. This comprehensive testing confirms our models' 851 advanced capabilities in financial domains and indicates potential 852 areas for further enhancements. 853

Hallucination. LLMs can sometimes generate inaccurate or misleading information, a phenomenon known as hallucination. In our experiments, we observed several hallucination manifestations that impact model outputs' reliability. One common form of hallucination was the generation of outputs that introduced entirely new labels, complete with plausible justifications, which deviated from expected outcomes. This was especially pronounced when models were given more examples, suggesting that they may overly rely on specific examples rather than interpreting the actual input. Furthermore, models occasionally produced outputs that did not adhere to the specified formats.

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

Case Study. In this analysis, we compare the performance of ESG-RoBERTa and FinLlama in deciphering complex ESG-related texts. An example from Table 5 shows ESG-RoBERTa misclassifying a case as "Business Ethics and Values," likely due to its reliance on

keywords such as "reputation" and "financial costs" associated with ethical issues. This keyword-centric approach struggles with nuanced texts, particularly longer ones, which can lead to bias or loss of information. Conversely, FinLlama employs ICL with CoT reasoning to accurately categorize the exact text as "Product Liability." It provides a detailed rationale, dissecting the text to link it to product liability issues and showcasing its ability to perform stepwise, human-like reasoning for accurate and interpretable results.

FinLlama's response highlights its adeptness at integrating extensive contextual information, outperforming conventional PLMs in complex scenarios. This case study underscores fine-tuned LLMs' advanced comprehension and reasoning capabilities in detailed, context-rich environments. Despite some limitations, it demonstrates that there is considerable scope to further enhance LLMs for more effective ESG text categorization.

Text: A warranty or product liability claim whether or not meritorious could harm our reputation result in substantial financial costs or divert managements attention which could have an adverse effect on our business operating results and financial condition. **True Label:** Product Liability

ESG-RoBERTa

Answer: Business Ethics and Values X

FinLlama (ICL-CoT):

Response: Label: Product Liability\n\nExplanation: The text mentions the potential for a warranty or product liability claim to harm the company's reputation, result in financial costs, and divert management's attention, which could have an adverse effect on the company's business operations and financial condition. This falls under the category of Product Liability, as it relates to the potential legal liability of a company for any defects or hazards in its products. **Answer:** *Product Liability* ✓

Table 5: An example of comparing results of PLM and LLM.

Cross-Comparison of Results. Comparing the performance of PLMs with LLMs, we observe that the pre-trained models excel in domain-specific tasks. Our results reveal that while specialized models often focus primarily on keywords, potentially overlooking broader contextual cues, large models are more adept at integrating context and avoiding keyword-based misinterpretations. Moreover, the inference process of large models tends to be more interpretable, providing clearer justifications for their decisions. This indicates that integrating the strengths of both specialized and large models, remarkably when fine-tuned with domain-specific data, could offer a promising research direction.

7 CONCLUSION AND FUTURE WORK

We proposed a pipeline to address the lack of ESG-related datasets, utilizing keyword searches and LLM APIs to annotate and construct three types of data for ESG text classification tasks. This approach has significantly enhanced the performance of pre-trained models on ESG classification tasks. We introduced domain-specific LLMs, ESGLlama and FinLlama, which were fine-tuned on our datasets, marking a major advancement in applying LLMs to ESGrelated challenges. Notably, FinLlama has surpassed existing financial benchmarks. Our methodology not only improves the capabilities of PLMs and LLMs within ESG contexts but also lays

Using Pre-trained Language Model for Accurate ESG Prediction

Conference acronym 'KiL, August 25, 2024, Barcelona, Spain

a foundation for ongoing innovation and practical applications
 in this vital area. Comparative analysis reveals that while PLMs
 generally perform better, LLMs offer greater interpretability and
 adeptly handle complex contexts by integrating contextual infor mation. Moving forward, we will further evaluate our developed
 datasets, and leverage the superior classification accuracy of PLMs
 to enhance and refine LLMs' performance in ESG analysis.

REFERENCES

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

- E. Alsentzer, J. R. Murphy, W. Boag, W.-H. Weng, D. Jin, T. Naumann, and M. B. A. McDermott. 2019. Publicly Available Clinical BERT Embeddings. https: //doi.org/10.48550/arXiv.1904.03323 arXiv:1904.03323 [cs].
- [2] A. Andonian, S. Biderman, S. Black, P. Gali, L. Gao, E. Hallahan, J. Levy-Kramer, C. Leahy, L. Nestler, K. Parker, M. Pieler, S. Purohit, T. Songz, W. Phil, and S. Weinbach. 2021. GPT-NeoX: Large Scale Autoregressive Language Modeling in PyTorch. https://doi.org/10.5281/zenodo.5879544
- D. Araci. 2019. FinBERT: Financial Sentiment Analysis with Pre-trained Language Models. https://doi.org/10.48550/arXiv.1908.10063 arXiv:1908.10063 [cs].
- [4] Y. Bang, S. Cahyawijaya, N. Lee, W. Dai, D. Su, B. Wilie, H. Lovenia, Z. Ji, T. Yu, W. Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. arXiv preprint arXiv:2302.04023 (2023).
- [5] I. Beltagy, K. Lo, and A. Cohan. 2019. SciBERT: A Pretrained Language Model for Scientific Text. https://doi.org/10.48550/arXiv.1903.10676 arXiv:1903.10676 [cs].
- [6] J. A. Bingler, M. Kraus, M. Leippold, and N. Webersinke. 2022. Cheap talk and cherry-picking: What ClimateBert has to say on corporate climate risk disclosures. *Finance Research Letters* 47 (June 2022), 102776. https://doi.org/10. 1016/j.frl.2022.102776
- [7] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems 33 (2020), 1877–1901.
- [8] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research* 24, 240 (2023), 1–113.
- [9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. https://doi.org/ 10.48550/arXiv.1810.04805 arXiv:1810.04805 [cs].
- [10] Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, and J. Tang. 2021. Glm: General language model pretraining with autoregressive blank infilling. arXiv preprint arXiv:2103.10360 (2021).
- [11] S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. https://doi.org/10.48550/arXiv.2004.10964 arXiv:2004.10964 [cs].
- [12] A. H. Huang, H. Wang, and Y. Yang. 2023. FinBERT: A Large Language Model for Extracting Information from Financial Text*. Contemporary Accounting Research 40, 2 (2023), 806–841. https://doi.org/10.1111/1911-3846.12832 _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/1911-3846.12832.
- [13] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa. 2023. Large Language Models are Zero-Shot Reasoners. arXiv:2205.11916 [cs.CL]
- [14] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. E. Gonzalez, H. Zhang, and I. Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles.
- [15] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36, 4 (Feb. 2020), 1234–1240. https://doi.org/10.1093/ bioinformatics/btz682
- [16] Y. Li, S. Wang, H. Ding, and H. Chen. 2023. Large Language Models in Finance: A Survey. https://doi.org/10.48550/arXiv.2311.10723 arXiv:2311.10723 [cs, q-fin].
- [17] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. https://doi.org/10.48550/arXiv.1907.11692 arXiv:1907.11692 [cs].
- [18] Z. Liu, D. Huang, K. Huang, Z. Li, and J. Zhao. 2020. FinBERT: A Pre-trained Financial Language Representation Model for Financial Text Mining, Vol. 5. 4513–4519. https://doi.org/10.24963/ijcai.2020/622 ISSN: 1045-0823.
- [19] D. Lu, H. Wu, J. Liang, Y. Xu, Q. He, Y. Geng, M. Han, Y. Xin, and Y. Xiao. 2023. BBT-Fin: Comprehensive Construction of Chinese Financial Domain Pre-trained Language Model, Corpus and Benchmark. https://doi.org/10.48550/arXiv.2302. 09432 arXiv:2302.09432 [cs].
- [20] L. Lu, J. Gu, and C.-R. Huang. 2022. Inclusion in CSR Reports: The Lens from a Data-driven Machine Learning Model. In Proceedings of the First Computing Social Responsibility Workshop within the 13th Language Resources and Evaluation Conference. 46–51.

- [21] A. Luccioni, E. Baylor, and N. Duchene. 2020. Analyzing sustainability reports using natural language processing. arXiv preprint arXiv:2011.08073 (2020).
- [22] T. Nugent, N. Stelea, and J. L. Leidner. 2020. Detecting ESG topics using domainspecific language models and data augmentation approaches. https://doi.org/10. 48550/arXiv.2010.08319 arXiv:2010.08319 [cs].
- [23] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. 2018. Deep Contextualized Word Representations. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers). Association for Computational Linguistics, New Orleans, Louisiana, 2227–2237. https://doi.org/10.18653/v1/N18-1202
- [24] M. Purver, M. Martinc, R. Ichev, I. Lončarski, K. S. Šuštar, A. Valentinčič, and S. Pollak. 2022. Tracking changes in ESG representation: Initial investigations in uk annual reports. In Proceedings of the First Computing Social Responsibility Workshop within the 13th Language Resources and Evaluation Conference. 9–14.
- [25] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research* 21, 140 (2020), 1–67. http://jmlr.org/papers/v21/20-074.html
- [26] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. arXiv preprint arXiv:1910.01108 (2019).
- [27] Z. Sautner, L. Van Lent, G. Vilkov, and R. Zhang. 2023. Firm-level climate change exposure. The Journal of Finance 78, 3 (2023), 1449–1498.
- [28] T. Schimanski, A. Reding, N. Reding, J. Bingler, M. Kraus, and M. Leippold. 2023. Bridging the Gap in ESG Measurement: Using NLP to Quantify Environmental, Social, and Governance Communication. https://doi.org/10.2139/ssrn.4622514
- [29] D. Stammbach, N. Webersinke, J. A. Bingler, M. Kraus, and M. Leippold. 2022. A dataset for detecting real-world environmental claims. *Center for Law & Economics Working Paper Series* 2022, 07 (2022).
- [30] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. html 3, 6 (2023), 7.
- [31] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023).
- [32] F. S. Varini, J. Boyd-Graber, M. Ciaramita, and M. Leippold. 2021. ClimaText: A Dataset for Climate Change Topic Detection. http://arxiv.org/abs/2012.00483 arXiv:2012.00483 [cs].
- [33] N. Wang, H. Yang, and C. D. Wang. 2023. FinGPT: Instruction Tuning Benchmark for Open-Source Large Language Models in Financial Datasets. https://doi.org/ 10.48550/arXiv.2310.04793 arXiv:2310.04793 [cs, q-fin].
- [34] N. Webersinke, M. Kraus, J. A. Bingler, and M. Leippold. 2021. Climatebert: A pretrained language model for climate-related text. arXiv preprint arXiv:2110.12010 (2021).
- [35] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in Neural Information Processing Systems 35 (2022), 24824–24837.
- [36] S. Wu, O. Irsoy, S. Lu, V. Dabravolski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg, and G. Mann. 2023. BloombergGPT: A Large Language Model for Finance. http://arxiv.org/abs/2303.17564 arXiv:2303.17564 [cs, q-fin].
- [37] Q. Xie, W. Han, X. Zhang, Y. Lai, M. Peng, A. Lopez-Lira, and J. Huang. 2023. PIXIU: A Large Language Model, Instruction Data and Evaluation Benchmark for Finance. https://doi.org/10.48550/arXiv.2306.05443 arXiv:2306.05443 [cs].
- [38] Y. Yang, M. C. S. UY, and A. Huang. 2020. FinBERT: A Pretrained Language Model for Financial Communications. https://doi.org/10.48550/arXiv.2006.08097 arXiv:2006.08097 [cs].
- [39] Y. Yin, Y. Yang, J. Yang, and Q. Liu. 2023. FinPT: Financial Risk Prediction with Profile Tuning on Pretrained Foundation Models. http://arxiv.org/abs/2308.00065 arXiv:2308.00065 [cs, q-fin].
- [40] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. 2022. Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068 (2022).
- [41] X. Zhang, Q. Yang, and D. Xu. 2023. XuanYuan 2.0: A Large Chinese Financial Chat Model with Hundreds of Billions Parameters. https://doi.org/10.48550/ arXiv.2305.12002 arXiv:2305.12002 [cs].
- [42] Y. Zheng, R. Zhang, J. Zhang, Y. Ye, Z. Luo, and Y. Ma. 2024. LlamaFactory: Unified Efficient Fine-Tuning of 100+ Language Models. arXiv preprint arXiv:2403.13372 (2024). http://arxiv.org/abs/2403.13372

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1040

1041 1042

1043

A DETAILS OF COLLECTED DATA

Below are the descriptions of the datasets we collected:

- *ESG-Prospectus-Clarity-Category*⁹: This dataset comprising 1,155 entries categorized into four ESG language classes: Specific, Ambiguous, Generic, and Risk. These entries were systematically extracted from the "Principal Investment Strategy" sections of sustainable (ESG) fund prospectuses through a specialized data extraction pipeline.
- Esg-sentiment¹⁰: Featuring text across nine emotion classes within the ESG spectrum (*<Environmental, Social, Gover*nance> * *<Negative, Neutral, Positive>*), each emotion assigns binary labels (0/1).
 - ESGBERT base-data¹¹: This dataset extracted 13,846,000 sentences from annual reports (13,079,890 sentences), responsibility reports (695,631 sentences), sustainable reports (259,163 sentences) and articles (143,289 sentences).
 - *Environmental_claims*¹²: This dataset focuses on the binary classification of environmental claims made by publicly listed companies, containing 2,647 entries. It is designed to detect real-world environmental assertions.
 - DAX ESG Media Dataset ¹³: Comprising approximately 11k recent English language ESG documents (text is document level) related to German DAX companies, this dataset includes both company issued reports and third party data, alongside an auxiliary file detailing the Sustainable Development Goals (SDGs).
 - *CLIMATE-FEVER* ¹⁴: This dataset consists of 1,535 realworld climate change claims. Each claim is supported by five Wikipedia-sourced evidence sentences annotated to either support, refute, resulting in a total of 7,675 claimevidence pairs.

Our data extraction involved the retrieval of the 'text' field across datasets, except the *DAX ESG Media Dataset*, from which the 'content' field was extracted, and the *CLIMATE-FEVER*, where both the 'claim' and the 'evidence' fields within the 'evidence' array were extracted. The summary of datasets is shown in Table 6.

B ESG KEYWORDS

All keywords we used shown in Table 7 refer to [28].

C WORD CLOUDS OF KEYWORD SEARCH

After filtering the texts by keywords searching. The texts are categorized into Environmental (Env), Social (Soc), Governance (Gov), and Non-ESG groups. The word clouds generated from these texts shown in Figure 6 offer a visual representation of the predominant themes within each category. In the Environmental domain, the word cloud prominently features terms such as "GHG emission" and "climate change," highlighting the focus on environmental impact. Socially oriented texts are characterized by frequent mentions of "human rights," "product," and "customer," reflecting

the emphasis on societal concerns and stakeholder welfare. In the Governance category, words like "director," "financial statement," "management," and "shareholder" dominate, aligning with expectations for governance-related discourse. These visual insights from the word clouds roughly correspond with our anticipated highfrequency words for each ESG classification, underscoring the effectiveness of our keyword-based filtering approach. we got the preliminary results shown in Table 8.

Table 8: Summary of Processed Data

| Domain | Num. of Sentences | Avg. Num. of Wo | | Words |
|-------------|-------------------|-----------------|-------|-------|
| | | Q1 | Mean | Q3 |
| Environment | 2,143,453 | 19 | 30.43 | 36 |
| Social | 2,796,077 | 20 | 31.46 | 37 |
| Governance | 1,851,303 | 20 | 31.75 | 38 |
| Non-ESG | 11,392,832 | - | - | - |
| Total | 18,183,665 | - | - | - |

D DATA DISTRIBUTION

Pre-training Dataset. We performed a 90-10 train-eval split to create the training and evaluation datasets, as shown in Table 9.

Table 9: Pre-training Dataset Statistics

| Dataset | Num. of Sentences |
|---------|-------------------|
| Train | 5,257,347 |
| Valid | 584,150 |
| Total | 5,841,497 |

For the four-class dataset. We used a 70:15:15 splitting ratio to construct the train-dev-test sets. The training set consisted of 37,155 instances, with 10,144 'Soc', 9,799 'Non-ESG', 9,192 'Env', and 8,020 'Gov'. The validation and test set each contained 7,962 instances, with 2,174 'Soc', 2,100 'Non-ESG', 1,969 'Env', and 1,719 'Gov' for validation, and 2,174 'Soc', 2,100 'Non-ESG', 1,970 'Env', and 1,718 'Gov' for testing. Results are shown in Figure 7.



Figure 7: Four-class Label Distribution in Train, Val, Test Sets

¹⁰⁹⁷ ⁹https://huggingface.co/Abhijeet3922

¹⁰https://huggingface.co/datasets/TrajanovRisto/esg-sentiment

^{1099 &}lt;sup>11</sup>https://huggingface.co/datasets/ESGBERT/base_data

¹²https://huggingface.co/datasets/climatebert/environmental_claims

 ¹³https://www.kaggle.com/datasets/equintel/dax-esg-media-dataset
 ¹⁴https://www.sustainablefinance.uzh.ch/en/research/climate-fever.html

| Dataset Name | Content Format | Size |
|---------------------------------|---|---------------------------|
| ESG-Prospectus-Clarity-Category | <text, label=""></text,> | 2310 rows (546 kB) |
| Esg-sentiment | <text, environmental<br="">Negative,,Social Positive></text,> | 679 rows (80.1 kB) |
| ESGBERT base-data | <text></text> | 13,846,000 rows (2.33 GB) |
| Environmental_claims | <text, label=""></text,> | 2647 rows (272 kB) |
| DAX ESG Media | <company, content,="" data,<br="" datatype,="">domain, esg_topics, internal, symbol, title></company,> | 11455 rows (130.11 MB) |
| CLIMATE-FEVER | <claim_id, claim,="" claim_label,<br="">evidences></claim_id,> | 1,535 rows (3 MB) |

Table 6: Summary of Collected ESG-Related Datasets

Domain Keywords Environmental adaptation, agricultural, air quality, biodiversity, biomass, climate, CO2, conservation, consumption, diversity, ecosystem, emissions, energy, environmental, flood, forest, fossil fuel, GHG, global warming, green, greenhouse, land use, methane, mitigation, nature, ozone, pollution, renewable, soil, solar, sustainability, water, recycling, clean energy, natural Social age, culture, race, accessibility, accident, accountability, awareness, charity, community, consumer protection, cyber security, data privacy, discrimination, diversity, education, employee benefit, empowerment, equality, ethics, fairness, gender, health, inclusion, mental well-being, parity, privacy, quality of life, religion, safety, social impact, volunteerism, welfare, wellbeing, workforce Governance audit, authority, bribery, compliance, corporate governance, corruption, crisis management, due diligence, ethics, framework, integrity, legal, lobby, oversight, policy, regulation, reporting, risk management, stakeholder engagement, transparency, whistleblower, board diversity, executive pay, shareholder rights, sustainable governance, corporate transparency, anti-corruption, business ethics



Figure 6: ESG Domain Word Clouds After Keywords Search

For the nine-class dataset. We applied an 81:09:10 splitting ratio. The training set had 17,419 instances, with each label ('Human Capital', 'Product Liability', 'Pollution and Waste', 'Business Ethics and Values', 'Corporate Governance', 'Community Relations', 'Non-ESG', 'Climate Change', 'Natural Capital'). The validation set contained 2,151 instances. Similarly, the test set had 1,936 instances.

These datasets were constructed using stratified sampling to ensure a balanced representation of each class in the train-dev-test splits. Lastly, we fine-tuned our pre-trained models on these two datasets to adapt them for the four-class and nine-class ESG text classification tasks. Results are shown in Figure 8.



Figure 8: Nine-class Label Distribution in Train, Val, Test Sets

E LLM LABELING PROMPTS DESIGN

We primarily utilize a combination of **few-shot learning** and Chain of Thought (**CoT**) in prompts design. Few-shot learning enables the model to learn from a limited quantity of text to align the acquired knowledge with our specific purpose. CoT [35] is a reasoning strategy that involves breaking down a problem into subproblems and connecting them in a specific logical order based on a chain structure. The purpose of using a few shots is to familiarize the model with the ESG classification strategy using a small sample. Using CoT is intended to enhance the model's reasoning process. Meanwhile, using CoT is intended to enhance the model's reasoning process in its responses, thereby improving its reasoning ability and enabling it to produce more data on the reasoning process for future SFT data construction.

For task 1: *Classify whether the text is high-quality ESG data: Yes or No.* The {Criteria} will be replaced by certain criteria, which are generated by GPT-4, and {Data} will be replaced by certain text we want to be classified.

```
You are a helpful assistant in data managing, and
      good at using high-quality data criteria for
      ESG content selection. To identify high-quality
      ESG data, we should consider the following criteria:
      {Criteria}
1325
      The following sentence is the data needed to define:
1326
      {Data}
1327
      Answer 'Yes' or 'No' first, then give an explanation.
1328
1329
      Let's think step by step.
1330
1331
1332
        For task 2: Classify whether the text is Env/Soc/Gov data: Yes or
      No. The Demonstrations are a few pairs of texts with their answer,
1333
```

Anon. Submission Id: 1

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388 1389

1390

1391

1392

| and {Type} can be environmental, social, and governance. | 1335 |
|--|------|
| You are an expert in ESG data classification | 1336 |
| | 1337 |
| especially {Type} ESG data classification. | 1338 |
| To identify {Type} ESG data, we should consider | 1339 |
| the following criteria: | 1340 |
| [Critoria] | 1341 |
| | 1342 |
| Answer 'Yes' or 'No' first, then give an explanation. | 1343 |
| Demonstrations: | 1344 |
| {Demonstrations} | 1345 |
| Toxt: (Data) | 1346 |
| Text. (Data) | 1347 |
| Answer: Let's think step by step. | 1348 |
| | 1349 |
| | 1350 |

For task 3: Assign texts to one of the following labels based on their domain. The {Classes} are corresponding to {Type}, which are defined previously. This work involves further subdividing the four-class data into nine classes.

You are an expert in classifying {Type} ESG data into finer categories such as {Classes}.

```
For {Type} ESG data classification, consider:
{Criteria}
```

First assign a class label based on these categories, or 'other' if uncertain. Then give an explanation.

Demonstrations:

```
{Demonstrations}
```

```
Text: {Data}
```

Answer: Let's think step by step.

E.1 Annotation Examples

To illustrate the LLM's responses to each task, we use text "*The possible threat of natural disasters affecting the ability to trade.*" as an example, the responses are shown in Table 10.

F CPT PLMS DETAILS

The continue pre-training parameters for our models were set as follows: a batch size of 16 for both training and evaluation, a total of 25 epochs, a seed value of 42 for reproducibility, and a maximum sequence length of 512. The models were trained with a gradient accumulation step of 2, weight decay of 0.01, and 16-bit precision (fp16) to optimize memory usage and accelerate the training process. The model checkpoints were saved and evaluated at the end of each epoch, as specified by the evaluation_strategy and save_strategy parameters set to 'epoch'. We conducted experiments on BERT, DistilRoBERTa, and RoBERTa at three different stages: epoch 5, epoch 15, and epoch 25. The results indicated that the models performed best at epoch 25. We utilized torch.distributed.launch for single-machine multi-GPU training. The pretraining of BERT, DistilRoBERTa, and RoBERTa took approximately 27 hours, 23.5 hours, and 33 hours, respectively. The training loss, validation loss, and validation accuracy during the training process are visualized

1334

1277

1278

1279

1280 1281

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1451

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

| Item | Content |
|-----------------|--|
| Text | The possible threat of natural disasters affecting the ability to trade. |
| Task 1 Response | Yes. |
| | 1. Relevance: The sentence is directly related to environmental issues, specifically the potential impact of natural disasters on trade. |
| | 2. Accuracy: It is difficult to determine the accuracy of this statement without additional context or |
| | information about the specific location and type of natural disaster being referred to. However, in general, |
| | natural disasters can disrupt trade and supply chains, so the statement is plausible. |
| | 3. Source Credibility: Without knowing the source of this statement, it is difficult to assess its credibility. |
| | 4. Specificity: The statement provides a general idea of the potential threat posed by natural disasters to |
| | trade, but does not provide detailed insights or examples. |
| | 5. Objectivity: The statement appears to be objective and presents a balanced view of the potential impact |
| | of natural disasters on trade. |
| Task 2 Response | Answer: Yes, the possible threat of natural disasters affecting the ability to trade is an environmental |
| | factor. Natural disasters can have a significant impact on the environment, including infrastructure, |
| | supply chains, and the ability to conduct business operations. |
| Task 3 Response | Label: Climate Change. |
| - | Explanation: The text mentions the possible threat of natural disasters affecting the ability to trade. This |
| | indicates a concern related to climate change, as natural disasters such as hurricanes, floods, and wildfires |
| | are often associated with climate-related events. Climate change can disrupt trade by damaging |
| | infrastructure, disrupting supply chains, and impacting transportation systems. Therefore, the text falls |
| | under the category of Climate Change. |

Table 10: Example of Data Through Three Tasks

in Figure 9. The smooth decrease in the loss curve and the gentle rise in the accuracy curve suggest that the models were converging and the training process was proceeding as expected. Consequently, we selected the model with the smallest validation loss at epoch 25 as our final pretraining models: ESG-BERT, ESG-DistilRoBERTa, and ESG-RoBERTa.

FINLLAMA TRAINING DETAILS G

Datasets. This fine-tuning was conducted on a specialized instruction-1430 tuning dataset on financial domain delineated in FinGPT [33]. Fur-1431 thermore, we enhanced the dataset by incorporating the ESG_Chat 1432 dataset, which consists of dialogues between humans and Large 1433 1434 Language Models (LLMs) focusing on methodologies to improve ESG scores. These conversations are structured to provide step-by-1435 step guidance, with the LLM responses specifically tailored to offer 1436 1437 structured, actionable advice. The characteristics of these datasets are detailed in Table 11. 1438

Hyperparameters. Each stage was meticulously conducted through-1439 1440 out the training regimen over 3 epochs to ensure the models' robust 1441 assimilation of the task-specific nuances. A consistent set of hyperparameters characterized the training to maintain uniformity 1442 across the models. Specifically, the batch size per device was set 1443 to 4, coupled with a gradient accumulation strategy involving four 1444 steps. This setup facilitated optimal resource utilization and stable 1445 training dynamics. The learning rate scheduler employed was of 1446 1447 the cosine type, which aided in gradual learning rate adjustments, contributing to smoother convergence. For monitoring and model 1448 checkpointing, logging intervals were established at every 10 steps, 1449

and model states were preserved at every 100 steps, ensuring detailed progress tracking and the ability to revert to the most effective model state. The learning rate was judiciously chosen as 5×10^{-5} , balancing rapid adaptation and the preservation of pre-learned representations.

The training progression for both models was visually documented through loss curves, providing insightful glimpses into the models' learning trajectories. Notably, a significant loss reduction was observed after the initial epoch for both models, indicative of their swift adaptation to the training objectives. For ESGLlama, the training culminated with the loss stabilizing around 0.4, shown in Figure 10a, suggesting effective learning. Conversely, FinLlama exhibited a distinct two-phase training dynamic; the initial pretraining phase concluded with a loss of around 2.4, shown in Figure 10b, which, upon undergoing the subsequent Supervised Fine-Tuning (SFT) phase, settled at approximately 1.15 shown in Figure 10c. This delineation in training phases for FinLlama underscores the layered approach to model refinement, first broadening its financial domain comprehension, followed by targeted instruction-based fine-tuning to hone its capabilities for specific financial tasks. These models will be tested on our labeled ESG classification data.

¹https://huggingface.co/datasets/FinGPT/fingpt-sentiment-train ²https://huggingface.co/datasets/FinGPT/fingpt-finred ³https://huggingface.co/datasets/FinGPT/fingpt-headline

⁴https://huggingface.co/datasets/FinGPT/fingpt-ner

⁵https://huggingface.co/datasets/FinGPT/fingpt-fiqa_qa

⁶https://huggingface.co/datasets/FinGPT/fingpt-fineval

⁷https://huggingface.co/datasets/zadhart/ESG_Chat



Figure 9: Continue Pre-training Log Loss and Accuracy across epochs

| | Table | 11: | Instruction | Financial | Dataset | Overview |
|--|-------|-----|-------------|-----------|---------|----------|
|--|-------|-----|-------------|-----------|---------|----------|

| Datasets | Train Rows | Test Rows | Description |
|-------------------------------------|------------|-----------|---------------------------------|
| fingpt-sentiment-train ¹ | 76.8K | N/A | Sentiment Analysis Training |
| | | | Instructions |
| fingpt-finred ² | 27.6K | 5.11K | Financial Relation Extraction |
| | | | Instructions |
| fingpt-headline ³ | 82.2K | 20.5K | Financial Headline Analysis |
| | | | Instructions |
| fingpt-ner ⁴ | 511 | 98 | Financial Named-Entity |
| | | | Recognition Instructions |
| fingpt-fiqa_qa ⁵ | 17.1K | N/A | Financial Q&A Instructions |
| fingpt-fineval ⁶ | 1.06K | 265 | Chinese Multiple-Choice |
| | | | Questions Instructions |
| ESG_Chat ⁷ | 914 | N/A | Chat History about Improve |
| | | | ESG Score step-by-step |



Figure 10: Training loss analysis during each stage of fine-tuning

H ESG CLASSIFICATION SFT DATASET

The dataset we used for supervised fine-tuning is constructed from ESG SFT data in Section 3.3. The ESG classification SFT data was sampled and reconstructed from ESG SFT data by only selecting classification data and straight-forward simplifying the result by retaining the text label without any additional explanations. There are two main classification tasks contained in this dataset: fourclass classification and nine-class classification. Finally, we obtained approximately 24k ESG Classification SFT Data. An example of the

1626

1627

1628

1629

1630

1631

1669

1670

1671

1672

1673

1674

1675

1676

1677

1678

1679

1680

1681

1682

1683

1684

1693 1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739

1740

ESG classification SFT dataset regarding these two tasks is shown in Table 12. Using this dataset, we aim to enhance the baseline's ESG text classification capability. This is because the baseline's pre-training data may include financial text data that is partially related to ESG. We intend to modify the baseline for this task and evaluate its performance during the SFT training phase.

| 532 | |
|-----|--|
| 633 | Format: [{"instruction": "", "input": "", "output": ""]}] |
| 34 | Four-class Classification: |
| | instruction: Classify the following text into one of the four ESG |
| | categories, choose an answer from {Env/Soc/Gov/Non-ESG}. |
| | input: We maintain a health and safety management system aligned |
| | to ISO legal requirements in Australia and New Zealand. |
| | |
| | Nine-class Classification: |
| | FSG categories choose an answer from {Climate Change/Natu- |
| | ral Capital/Pollution and Waste/Human Capital/Product Liability/- |
| | Community Relations/Corporate Governance/Business Ethics and |
| | Values/Non-ESG}. |
| | input: Grievance mechanisms forms an important part of our stake- |
| | holder engagement process, and our human rights policy states that |
| | if we have caused or contributed to adverse human rights impacts |
| | output: Human Capital |
| | Table 12: An example of ESG classification SFT dataset. |
| | 1 |
| | |
| | I CLASSIFICATION PROMPTS |
| | |
| | System Prompt: "You are an expert in classifying ESG |
| | data You will start your response with 'label.' " |
| | Hear Brownet, "Classify the following text into an of |
| | User Prompt: classify the following text into one of |
| | the four ESG categories, choose an answer from |
| | {Categories} |
| | Demonstrations: |
| | {Demonstrations} |
| | Text: {Text} |
| | Labol. So the answer is" |
| | Label. 50, the diswer is |
| | |
| | |

For Four-Class classification task, we should specify the {Categories} by:

{Env/Soc/Gov/Non-ESG}

For Nine-Class classification task, we should specify the {Categories} by:

{Climate Change/Natural Capital/Pollution and Waste/ Human Capital/Product Liability/Community Relations/ Corporate Governance/Business Ethics&Values/Non-ESG}

To employ a chain-of-thought (CoT) setting, we need to slightly modify the system prompt and add let's think step by step at the end of the user prompt:

| System Prompt: "You are an expert in classifying ESG | 1685 | | | | | |
|--|------|--|--|--|--|--|
| data. You will response in this format: | 1686 | | | | | |
| 'Label:xxx. Explanation:xxx'. | | | | | | |
| Your responses should be precise and concise " | 1688 | | | | | |
| Har Brent " | 1689 | | | | | |
| User Prompt: | 1690 | | | | | |
| Label: Let's think step by step. So, the answer is" | 1691 | | | | | |
| | 1692 | | | | | |

J ADDITIONAL LLM CLASSIFICATION ANALYSIS

For Four-class classification. In evaluating our models, ESGLlama and FinLlama, within our experimental framework, we employed Precision, Recall, F1 Score, and Accuracy as our performance metrics. Initially, let us delve into the precision aspect, which serves to illustrate the models' exactness in classification tasks. Through the analysis of precision scores and the accompanying graphical representations shown in Figure 4, it becomes evident that both ESGLlama and FinLlama surpass the baseline model across most experimental configurations. Furthermore, even the baseline model, when subjected to Supervised Fine-Tuning (SFT) using our constructed ESG classification dataset, demonstrates enhanced performance compared to its original state. Interestingly, the Freeze fine-tuning approach generally outperforms the LoRA method, except in zero-shot settings. This observation could be attributed to the Freeze technique requiring a broader range of parameters for fine-tuning, thereby facilitating a deeper understanding of downstream tasks. In contrast, LoRA's superior performance in zero-shot scenarios might hint at a slight overfitting issue; external demonstration examples, not included in the training set, could potentially disrupt the model's inference processes. The Freeze approach, in this context, better preserves the model's generalization capabilities and intrinsic reasoning faculties.

The incorporation of Chain of Thought (CoT) prompts leads to a performance decline in zero-shot and one-shot settings, except for the Iterated Chain of Learning (ICL) tasks. This decline could stem from the absence of stepwise reasoning chains in our training data, coupled with the inherent incompatibility of the CoT methodology with classification tasks—CoT primarily suits logic-based problem-solving. Nevertheless, the addition of demonstrations in ICL tasks enriches the model's learning of classification logic through increased sample exposure, culminating in the most favorable outcomes under ICL CoT configurations.

Further examination of performance metrics, as detailed in the corresponding table shown in Table 13, reveals that the LoRA method, applied directly to the baseline on our ESG classification dataset, achieves the highest precision (0.6928), recall (0.5557), F1 score (0.5488), and accuracy (0.5697) in zero-shot tasks. This outcome not only underscores the constructed dataset's validity but also establishes a benchmark for subsequent comparisons. Furthermore, the bold formatting in the table highlights the highest precision scores across six method settings for each model, underscoring the best-performing configurations. The underlined values denote the top performance metrics across all models and settings,

establishing a benchmark for comparison. The star symbol (*) identifies the best baseline result for the LoRA and Freeze fine-tuning methods, serving as a reference point for assessing the fine-tuned models' enhancements. The directional arrows $(\uparrow\downarrow)$ provide a visual cue for performance fluctuations in comparison to the baseline, elucidating the impact of our fine-tuning strategies on model preci-sion. Against this backdrop, both ESGLlama and FinLlama exhibit a decline, albeit still outperforming the baseline, especially in ICL settings. Notably, FinLlama achieves superior precision over ES-GLlama with the addition of CoT, underscoring the nuanced impact of our training methodologies on model performance. In summary, the table elucidates the nuanced interplay between fine-tuning methodologies, the inclusion of CoT prompts, and the iterative learning approach on model precision. The discernible improve-ment in precision with ESGLlama and FinLlama, particularly in ICL settings, reaffirms the efficacy of our fine-tuning strategies in embedding ESG-specific knowledge into large language models.

For Nine-class classification, the analysis of performance metrics, particularly precision, elucidates a notable trend: as the complexity and diversity of classification categories increase, the task inherently becomes more challenging, as evidenced by the overall diminished performance compared to the four-class scenario. This trend underscores the escalated difficulty in distinguishing among a greater number of classes.

The precision score visualization (Figure 5) demonstrates that both ESGLlama and FinLlama significantly outperform the baseline model across most methodological settings. This superiority highlights our fine-tuned models' enhanced understanding and classification capability in the context of ESG-related texts. FinLlama

demonstrates superior proficiency in iterative contrastive learning (ICL), particularly in scenarios with increased sample availability, indicating a profound comprehension of financial texts and their nuances. The analysis further reveals a pronounced decrement in performance for the one-shot learning setting across more granular classification tasks. Providing only one example per class introduces considerable bias and may confound the model's judgment due to the high variance associated with minimal data. Conversely, enriching the model with a broader set of examples significantly ameliorates performance, aligning with the expected benefits of expanded training data. This intricate classification landscape observes a notable divergence in the efficacy of the LoRA and Freeze fine-tuning methods. Interestingly, The LoRA approach exhibits superior performance in the one-shot setting compared to Freeze, suggesting that LoRA's parameter adaptation might be more conducive to effectively assimilating sparse class-specific information.

Delving deeper into the details presented in the accompanying Table 14, the most commendable performance is attributed to FinLlama under the ICL with Chain of Thought (CoT) augmentation, achieving a precision score of 0.6654. This result significantly surpasses the baseline precision of 0.6164 and even outstrips the baseline model fine-tuned with LoRA on the ESG classification data, which scored 0.6544. This evidence conclusively demonstrates the potent efficacy of FinLlama, particularly when augmented with CoT in complex classification scenarios, further substantiating the model's refined comprehension of financial discourse and its implications for ESG classification tasks.

| Model | Mathada | Overall | | | | |
|------------|---------------|-----------------|-----------------|-----------------|-------------------|--|
| Model | Methods | Precision | Recall | F1 Score | Accuracy | |
| | Zero Shot | 0 5778 | 0 5025 | 0 4815 | 0 5093 | |
| | w/ CoT | 0.5527 | 0.4613 | 0.4252 | 0.4776 | |
| | Our Chat | 0 (010 | 0.5057 | 0.4707 | 0.5100 | |
| Llama2 | Une Shot | 0.6012 | 0.5056 | 0.4/06 | 0.5109 | |
| | W/ C01 | 0.3370 | 0.3707 | 0.2080 | 0.3931 | |
| | ICL | 0.6687 | 0.5408 | 0.5077 | 0.5446 | |
| | w/ CoT | 0.6794 | 0.5193 | 0.4803 | 0.5229 | |
| | Zero Shot | <u>0.6928</u> * | <u>0.5557</u> * | <u>0.5488</u> * | <u>0.5697</u> * | |
| | w/ CoT | 0.6381 | 0.4973 | 0.5128 | 0.5053 | |
| LoRA | One Shot | 0.5265 | 0.3896 | 0.2924 | 0.3976 | |
| | w/ CoT | 0.5646 | 0.3291 | 0.2442 | 0.3360 | |
| | ICL | 0.6148 | 0 5157 | 0 4821 | 0 5232 | |
| | w/ CoT | 0.6213 | 0.3971 | 0.3247 | 0.4019 | |
| | Zana Chat | 0 5741 | 0 5000 | 0 4797 | 0 50(9 | |
| | Lero Shot | 0.5741 | 0.5000 | 0.4/8/ | 0.5068 | |
| | W/ C01 | 0.3400 | 0.4015 | 0.4270 | 0.4775 | |
| Freeze | One Shot | 0.6085 | 0.5113 | 0.4761 | 0.5168 | |
| | w/ CoT | 0.6168 | 0.3932 | 0.2873 | 0.4073 | |
| | ICL | 0.6611 | 0.5382 | 0.5036 | 0.5422 | |
| | w/ CoT | 0.6749 | 0.5181 | 0.4767 | 0.5216 | |
| | Zero Shot | 0.5770 | 0.4997 | 0.4768 | 0.5054 | |
| | w/ CoT | 0.5502 | 0.4594 | 0.4205 | 0.4753 | |
| FSGLlama | One Shot | 0.6106 | 0 5373 | 0 5140 | 0 5389 | |
| LUCLIAIIIA | w/ CoT | 0.6064 | 0.3984 | 0.3128 | 0.4147 | |
| | ICI | 0 6 7 2 8 | 0 5500 | 0 5202 | 0 55 4 9 | |
| | ICL w/ CoT | 0.0730 | 0.3500 | 0.3203 | 0.3340↓ 0.4935 | |
| | W/ C01 | 0.07 10 | 0.4002 | 0.4525 | 0.4755 | |
| | Zero Shot | 0.5766 | 0.4961 | 0.4745 | 0.5024 | |
| | w/ C01 | 0.5665 | 0.4669 | 0.4297 | 0.4828 | |
| FinLlama | One Shot | 0.6139 | 0.5375 | 0.5139 | 0.5394 | |
| | w/ CoT | 0.5724 | 0.3856 | 0.3011 | 0.4017 | |
| | ICL | 0.6698 | 0.5497↓ | 0.5174↓ | 0.5535↓ | |
| | w/ CoT | 0.6797↓ | 0.4917 | 0.4365 | 0.4971 | |

Table 13: Four-class evaluation results compare with baseline and our fine-tuned LLMs. Blod shows the best results in six method settings according to each model, and underline illustrates the best performance in each column. Star (*) is the best baseline result for two fine-tuning methods (LoRA and Freeze). Arrow (↑↓) signifies performance compared with Star (*).

| Model | Methode | Overall | | | | | |
|----------|-----------|-----------------|---------|----------|----------|--|--|
| mouer | methous | Precision | Recall | F1 Score | Accuracy | | |
| | Zero Shot | 0.5875 | 0.4404 | 0.4454 | 0.4886 | | |
| | w/ CoT | 0.5826 | 0.4106 | 0.4171 | 0.4654 | | |
| Llama2 | One Shot | 0.5049 | 0.4322 | 0.3877 | 0.4737 | | |
| | w/ CoT | 0.4314 | 0.3556 | 0.2895 | 0.3838 | | |
| | ICL | 0.6108 | 0.4029 | 0.4017 | 0.4411 | | |
| | w/ CoT | 0.6164 | 0.4624 | 0.4932 | 0.5057 | | |
| | Zero Shot | 0.5681 | 0.4901 | 0.4759 | 0.5294* | | |
| | w/ CoT | 0.5180 | 0.4112 | 0.3895 | 0.4473 | | |
| LoRA | One Shot | 0.6256 | 0.5347* | 0.4795* | 0.5186 | | |
| | w/ CoT | 0.5751 | 0.3915 | 0.3450 | 0.3972 | | |
| | ICL | 0.6242 | 0.1946 | 0.1340 | 0.2257 | | |
| | w/ CoT | 0.6544* | 0.1834 | 0.1465 | 0.2123 | | |
| | Zero Shot | 0.5911 | 0.4458 | 0.4488 | 0.4974 | | |
| | w/ CoT | 0.5799 | 0.4122 | 0.4161 | 0.4664 | | |
| Freeze | One Shot | 0.5258 | 0.4445 | 0.4148 | 0.4866 | | |
| | w/ CoT | 0.4922 | 0.4005 | 0.3353 | 0.4323 | | |
| | ICL | 0.6285 | 0.4189 | 0.4265 | 0.4649 | | |
| | w/ CoT | 0.5719 | 0.2432 | 0.2337 | 0.2862 | | |
| | Zero Shot | 0.5866 | 0.4271 | 0.4340↓ | 0.4778 | | |
| | w/ CoT | 0.5914 | 0.4190 | 0.4258 | 0.4726 | | |
| ESGLlama | One Shot | 0.5138 | 0.4446↓ | 0.4136 | 0.4855↓ | | |
| | w/ CoT | 0.4785 | 0.4031 | 0.3373 | 0.4318 | | |
| | ICL | 0.6201↓ | 0.4143 | 0.4235 | 0.4576 | | |
| | w/ CoT | 0.5773 | 0.2533 | 0.2470 | 0.2965 | | |
| | Zero Shot | 0.5608 | 0.4293 | 0.4301↓ | 0.4830↓ | | |
| | w/ CoT | 0.5750 | 0.4123 | 0.4164 | 0.4664 | | |
| FinLlama | One Shot | 0.5219 | 0.4376↓ | 0.4069 | 0.4757 | | |
| | w/ CoT | 0.4886 | 0.4062 | 0.3399 | 0.4349 | | |
| | ICL | 0.6168 | 0.4127 | 0.4163 | 0.4638 | | |
| | w/ CoT | <u>0.6654</u> ↑ | 0.2504 | 0.2478 | 0.2908 | | |

Table 14: Nine-class evaluation results compare with baseline and our fine-tuned LLMs. Bold shows the best results in six method settings according to each model, and <u>underline</u> illustrates the best performance in each column. Star (*) is the best baseline result for two fine-tuning methods (LoRA and Freeze). Arrow (↑↓) signifies performance compared with Star (*).