# **Open-World Authorship Attribution**

## Anonymous ACL submission

#### Abstract

Recent years have witnessed rapid advance-001 ments in Large Language Models (LLMs). Nevertheless, it remains unclear whether state-004 of-the-art LLMs can infer the author of an anonymous research paper solely from the text, without any additional information. To inves-007 tigate this novel challenge, which we define as Open-World Authorship Attribution, we introduce a benchmark comprising thousands of research papers across various fields to quan-011 titatively assess model capabilities. Then, at the core of this paper, we tailor a two-stage 013 framework to tackle this problem: candidate selection and authorship decision. Specifically, 015 in the first stage, LLMs are prompted to generate multi-level key information, which are then 017 used to identify potential candidates through Internet searches. In the second stage, we introduce key perspectives to guide LLMs in deter-019 mining the most likely author from these candidates. Extensive experiments on our benchmark demonstrate the effectiveness of the proposed approach, achieving 60.7% and 44.3% accuracy in the two stages, respectively. We will release our benchmark and source codes to facilitate future research in this field.

## 1 Introduction

027

033

037

041

The advancement of Generative Artificial Intelligence (AI) and Large Language Models (LLMs) has revolutionized numerous fields due to their remarkable capabilities in Natural Language Processing (NLP) tasks (Hagos et al., 2024; Naveed et al., 2024; Cui et al., 2024; Zhang et al., 2024). Despite their widespread applications, their potential for authorship attribution—the task of identifying an author from anonymous text—remains largely unexplored. In this paper, we investigate an intriguing question: *Can state-of-the-art LLMs infer the author of an anonymous research paper without any additional information?* This problem is both practical and ambitious. On one hand, accurately attributing authorship in academic research is crucial for maintaining integrity, recognizing contributions, and detecting plagiarism or ghostwriting. On the other hand, directly applying modern LLMs to this task is challenging, as the relevant information is often dispersed across Internet-scale data, which makes it infeasible for these models to process efficiently. 042

043

044

047

048

054

057

061

062

063

064

065

066

067

068

069

070

071

074

075

076

077

078

In this paper, we define this challenging task as *Open-World Authorship Attribution*. Since no existing benchmark evaluates LLM performance in this area, we construct a dataset comprising thousands of academic papers from various research fields. Building on insights from this data, we propose a novel two-stage framework to address the task, including *Candidate Selection* and *Authorship Decision*.

Specifically, in the candidate-selection stage, we leverage LLMs to generate multi-level key representations of a target paper, which are then utilized to search the Internet for relevant authors and their publications at multiple levels of specificity. The retrieved authors along with authors in the citation list form our candidate pool. In the authorshipdecision stage, LLMs assess potential authorship in the candidate pool by evaluating the anonymous text against multiple guidelines. Finally, a holistic decision is made to determine the most probable author, serving as the final output.

We conduct experiments using multiple stateof-the-art LLMs, including both open-source and closed-source models.

Extensive evaluation validate the effectiveness and the superiority of the proposed solution. Specifically, our approach achieves 60.7% accuracy of candidate selection and 44.3% accuracy of authorship decision. The contribution of this work is summarized as below:

• We are the first to define, study, and benchmark the task of open-world authorship attri-

082	bution to the best of our knowledge.
083	• By leveraging impressive capacity of recent
084	LLMs, we devise a novel two-stage pipeline,
085	including candidate selection and authorship
086	decision, to tackle this challenge.
007	• Extensive evaluations showcase the potential
007	• Extensive evaluations showcase the potential
880	of modern LLWs and our proposed solution
089	for open-world authorship attribution. We will
090	release the dataset, prompts, and codes to sup-
091	port future research in this field.
092	2 Related Work

## 2.1 Large Language Models (LLMs)

097

100

102

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123 124

125

126

127

128

LLMs have demonstrated remarkable capability in solving various Natural Language Processing (NLP) tasks, such as mathematical reasoning, and text summarization (Xuanfan Ni, 2024; Desta Haileselassie Hagos, 2024). The unique characteristic of LLMs lies in utilizing a unified paradigm without additional training to address various tasks (Qin et al., 2024).

Language modelling: Language modelling as the core to current LLMs has developed from the traditional statistical methods like n-gram (Sharma et al., 2018a) models to Neural Network language models. The transformer language models with self-attention mechanisms further lay the foundation for the current rapid development of LLMs (Vaswani et al., 2017). The introduction of the revolutionized transformer helped the development of the GPT-1 transformer-decoder structure and Bert's transformer-encoder structure.

LLMs Tuning: Tuning techniques have evolved alongside the development of LLMs. Tuning consists of full-parameter and partial-parameter tuning. Due to computational constraints, research has focused on Parameter-Efficient Fine-Tuning (PEFT), including prompt tuning, Adapter-Tuning, and LoRA. In-context learning, a form of prompt learning, enables adaptation without parameter updates by providing example-based prompts.

Instruction tuning is also the current focus. The purpose is to transform NLP tasks with natural language instruction which improves the performance of LLMs in zero-shot learning. Chain-of-Thought (Wei et al., 2023) is another reasoning strategy to resolve the issue of low performance in arithmetic reasoning, normal inference and symbol inference.

### 2.2 AI-generated texts Detection with LLMs

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

The widespread accessibility of generative models has led to a proliferation of AI-generated texts across the internet. Several detection approaches have been developed to detect LLMgenerated works to address the issue of authenticity: (1) Training-based method adopt classifiers like Support Vector Machines (SVMs) or fine-tuned pre-trained language models like RoBERTa and T5(Yang et al., 2023; Tang et al., 2023). (2) Zeroshot Detection method directly uses the inherent properties embedded in LLMs (Yang et al., 2023). (3) Watermarking-based Detection like Inferencetime watermarking (Tang et al., 2023) embeds unique patterns into text during generation by manipulating decoding processes, while post-hoc watermarking retroactively modifies generated text using rule-based or neural techniques to ensure traceability (Tang et al., 2023).

## 2.3 Authorship Identification

Several studies have already researched the authorship identification capabilities of LLMs, highlighting the importance of authorship attribution in forensic investigations, cybersecurity, and tackling misinformation (Huang et al., 2024).

Traditionally, authorship attribution and verification focus on analyzing writing styles to measure similarities and make authorial decisions. Early methods employed natural language processing (NLP) techniques, such as n-grams (Sharma et al., 2018a), part-of-speech (POS) tags (Sundararajan and Woodard, 2018), and Linguistic Inquiry and Word Count (LIWC) (Uchendu et al., 2020). These handcrafted features are designed to quantify stylistic patterns, including vocabulary richness, syntactic complexity, and semantic focus, for effective analysis (Huang et al., 2024).

More recently, with advancements in deep learning, text embeddings have become a prominent tool in authorship attribution. Text embeddings represent textual data as vectorized numerical representations, enabling models to encode both semantic and stylistic nuances. (Kumarage and Liu, 2023) emphasizes the potential of leveraging large pretrained language models (LLMs) like BERT and GPT to generate embeddings that capture deeper stylistic and contextual patterns, thus applying them to authorship attribution tasks. Another significant change in authorship attribution is the integration of contrastive learning techniques into



Figure 1: Violin plot illustrating the distribution of information entropy among 300 authors. An information entropy value of 1.5219 indicates that, among the five collected articles for a given author, two articles share the same topic group, another two belong to a different topic group, and the remaining article falls into a separate group

embedding-based methods(Patel et al., 2023).

Some methods developed the prompt pipeline for authorship identification, leveraging the inherent stylistic and linguistic extraction capabilities of LLMs (Huang et al., 2024; Wen et al., 2024). The results demonstrate the ability of LLMs to capture nuanced stylistic features without explicit feature engineering. However, limitations of the study are noticeable, such as dependency on pre-collected candidate authors which hinders its application in large-scale candidate pools. In most cases, the number of candidate authors is fewer than 50, making the approach impractical for real-world applications. Their focus on stylometric feature analysis and prioritizing explainability in the authorship decision-making limits its efficiency.

## 3 Methods

179

181

182

185

187

188

189

190

191

192

193

194

195

196

198

199

200

201

204

210

In this section, we elaborate on the proposed benchmark and two-stage approach for open-world authorship attribution. Sec. 3.1 introduces the data sources and the construction of the benchmark. Secs. 3.2 and 3.3 describe the main pipelines of the two stages: candidate selection and authorship decision, respectively, in the proposed solution. In the first stage, candidate papers are retrieved from the Internet using LLM-generated keywords. In the second stage, LLMs determine the most likely author from these candidates. Fig. 3 provides an overview of the entire streamline.

## 3.1 Data Curation

Considering there is no off-the-shelf benchmark for the task of open-world authorship attribution, we



Figure 2: This heatmap shows the extent to which authors choose the same topics across their publications. Each cell represents the co-occurrence strength between topics for the same author, with darker shades indicating a higher likelihood of an author selecting the same topic in their papers. Each number in axises indicates different topic groups 5

construct a dataset in this work. Specifically, To ensure diversity, we select papers from CVPR 2024, spanning 30 subfields in computer vision. Moreover, to guarantee sufficient online reference materials for candidate selection, we filter out authors with fewer than five first-author papers. This results in a dataset comprising 300 authors and 1,500 papers. More details are provided in Sec. 4.1.

211

212

213

214

215

216

217

218

219

221

222

223

224

225

227

228

229

230

231

233

234

## 3.2 Candidate Selection

The core challenge of open-world authorship attribution lies in handling Internet-scale data, which significantly exceeds the processing capabilities of LLMs. Therefore, identifying common patterns among papers by the same author is crucial for narrowing down potential candidates from such a vast data source.

During our investigation and collection of the dataset, we observe that many authors' published papers demonstrate a correlation in research topics. To analyze the relationship between authors and research topics, we leverage Information Entropy and perform a statistical analysis. Specifically, for each author, we evaluate the randomness and diversity of the involved research fields via:

$$H(T) = -\sum_{i=1}^{N} p(T_i) \log_2 p(T_i), \qquad (1)$$



Figure 3: The whole process of the proposed automatic end-to-end authorship identification consists of 2 stages. The first stage will contribute to the collection of relevant candidates. The second stage performs 2 rounds of inference decision based on different metrics with different input information from candidate pool. The final decision is based on the 2 inferences.

where  $p(T_i)$  represents the probability of an author's papers belonging to topic  $T_i$ . Higher entropy suggests greater diversity in research areas, while lower entropy indicates topic consistency, aiding in author identification.

In the violin plot shown in Fig. 1, we observe that only a small subset of authors exhibit high entropy and diversity across research fields. In Fig. 2, we also visualize the extent to which authors in a given field also conduct research in other fields. Strong diagonal activations suggest a high likelihood of topic overlap within an author's publications.

Based on the correlation of topics between authors' different papers above, we can rely on this to search for the author. Therefore, we need some keywords which can summarize the topics from anonymous text and be utilized for searching the relevant articles. This is where LLMs can help in our candidate selection stage - keyword generation. These candidate articles will be used further in the next stage of decision-making.

**Keywords Generation.** The keywords generated need to accurately capture the contents of the anonymous input for the effective searching of the relevant articles. Therefore, we decided to generate the relevant keywords in hierarchies to describe the anonymous content. The different levels should range from general (level 1) to specific (level 5). In this way, we can search from the most specific to the most general to get the relevant articles as our candidates prepare for the next stage of decisionmaking.

**Few-shot Prompting.** If we ask LLMs to generate the keywords directly in hierarchies, the output may be in different formats and the quality of generation is not guaranteed by the simple instruction of "Generate 5-level keywords". Few-shot prompting is utilized for LLMs to demonstrate how the response should look. We will manually create an example 6 and use it as an example within the prompt. This will maximize the possibility of effective keyword generation which follows our proposed method.

**Candidates Search & Collect.** Different levels of the keywords will be sent to the search engine to collect the potential candidates. The search engine we choose is Scholar Inbox which has a semantic section to input the different level keywords. In

237

240

241

242

243

245

246

247

248

252

253

259

283

260

261

262

263

264



Figure 4: This figure presents the cumulative density plot (CDF) of citation similarity as it increases. The varying rate of cumulative percentage growth indicates that articles within the same topic tend to exhibit higher citation similarity.

each search result of the level keywords, we decided to collect the first 20 papers with their respective titles, authors and abstracts. If the searched article provides an arxiv link, we will try to retrieve the introduction and citation to help the stage 2 decision-making.

**Appending Self-citation.** Most authors tend to cite their previous works, especially when focused on a specific research area. This characteristic increases the likelihood of successfully including the true author in the candidate pool, which serves as the foundation for Stage 2. To leverage this, we incorporate self-citation into our candidate selection process, and their metadata is retrieved through web scraping to enhance the effectiveness of our decision-making.

#### 3.3 Authorship Decision

284

289

294

296

297

300

301

303

307

310

311

313

LLMs also show their great capability in analyzing large-scale data. Therefore In this stage, after creating our potential candidate pool, we need to enable LLMs to decide the most possible author. However, direct instruction like "Please decide the most possible author from the candidate pool" is not detailed enough for LLMs to accurately decide the most possible author. Therefore, we establish different metrics to guide LLMs in identifying the most probable author step by step. Chain-of-Thought (CoT) (Wei et al., 2023) prompting is a technique designed to enhance the reasoning capabilities of large language models (LLMs) by guid-

#### Algorithm 1 Open World Authorship Attribution

- 1: **Input:** Anonymous article x
- 2: **Prompts:**  $metrics = \{style, citation\}$
- 3: Output: Final attributed author 4: Step 1: Keyword Generation
  - Step 1: Keyword Generation# Extract representative keywords using LLMs:
- 5: keywords  $\leftarrow LLMs(x, Prompts, Example)$
- 6: Step 2: Candidates Collection
- 7: for each level in keywords do
- 8: # Use web scraping to collect potential authors and articles:
- 9: CandidatePool += WebScraping(level)
- 10: end for
- 11: Step 3: Iterative Filtering by Metrics
- 12: for i = 1 to |metrics| do
- 13: # Use LLMs to rank authors based on current metric:
- 14:  $TopAuthors \leftarrow LLMs(x, CPool, metrics_i)$
- # Append with the top-ranked authors: 15:  $CandidatePool \leftarrow TonAuthors$
- 15: CandidatePool ← TopAuthors
  16: end for
- 17: Step 4: Final Attribution
- #determine the most likely author from *TopAuthors*:
- 18:  $FinalAuthor \leftarrow LLMs(x, TopAuthors)$
- 19: **Return:** *FinalAuthor*

ing them to generate intermediate reasoning steps before arriving at a final answer. This approach mirrors the human thought process, and by simulating it, LLMs can achieve more accurate analysis and decision-making outcomes. 314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

331

332

333

334

335

336

337

338

340

341

342

343

For each metric, we ask the LLMs to list the top 5 authors that most match the metrics. In the last step, we input all the decision results from its decision and let LLMs decide the most possible one holistically. In this way, we can identify the most possible author. LLMs will perform 3 rounds of most possible author inference based on contents, writing styles, and citation similarities. The full prompting can be referred to in Appendix 6.

### 3.3.1 Contents

As mentioned in the above section, due to the similarities of authors chosen topic in their published paper, content is the important metric to determine the actual author (Halvani and Graner, 2021; Potha and Stamatatos, 2019). the specific metrics within the contents is preferred topics and Domain-Specific term used in the writing. We will input author names, titles, abstracts and introduction (not every paper have) from the candidates pool for LLMs inference.

## 3.3.2 Writing Style

Relying on stylometry is the traditional way of authorship attribution. Evolving from the human skills in identifying (Argamon et al., 2009), computational methods gradually become the main trend

433

434

435

436

437

438

439

440

441

442

443

444

395

in the analysis of authors' unique linguistic features in stylometry methods(Lagutina et al., 2019; Neal et al., 2017). Machine learning methods with LLMs further advance the computational methods with their powerful ability to extract features (Boenninghoff et al., 2019; Kojima et al., 2022).

345

346

351

353

357

365

370

374

376 377

386

390

391

394

In our proposed method, we also utilize writing style as an important factor in identifying the author from our collected candidate pool. Different people have different habits or underlying characteristics in writing. Some typical metrics are repetition Patterns in words and phrases (Sharma et al., 2018b), sentence complexity, paragraph structure, sentence length and variation.

To demonstrate the importance of writing style in differentiating authors, we select three paragraphs from three articles—two written by the same author and one by a different author(Appendix 7). To eliminate the influence of topic variation, all three articles cover the same subject. We then prompt GPT-40 to analyze the texts and determine which two paragraphs are authored by the same individual. GPT-40 demonstrates its ability to make this distinction based on factors such as writing tone, repetition patterns, sentence complexity, and paragraph structure.

In this metric, we have the same input as contentbased inference which is title, author, abstract and introduction.

#### 3.3.3 Citation Similarity

As mentioned above, the topics for an author's published papers are largely overlapped. This may also indicate that the literature review - citations of an author tend to be similar. In other words, the author has a preference for some citations and they may prefer to reuse these citations in their other works. Based on this assumption we establish the third metric which is the citation similarities. In this stage, we only utilize the author names, articles titles with the corresponding citations as the input for the LLMs.

We conduct a citation similarity analysis on our collected dataset, calculating the similarity between the citations of test articles and titles authored by the same or different authors. The plotted cumulative distribution function (CDF) Figure 4 illustrates the distribution of citation similarity under the same or different authors. Approximately 40% of the articles exhibit similar citation match probabilities between the same author and different authors. However, the remaining 60% show a clear distinction in citation similarity. In general, articles written by the same author tend to have higher citation similarity values.

## 4 Experiments

### 4.1 Experimental Setup

**Models.** We choose nowadays popular LLMs to conduct the test. We download the open source Meta-Llama-3.1-8B-Instruct (Grattafiori et al., 2024) for the initial test with both abstract-only and abstract-plus-introduction as input information input. We also use the recent GPT-4o-mini (OpenAI et al., 2024) to conduct the test. GPT-4o-mini was accessed and tested via API requests.

**Dataset.** Due to a lack of academic paper datasets available online, We self-collected our dataset for our testing. Our dataset includes 300 authors, with 5 papers selected for each author. To ensure relevance and keep the writing style of the author, we only selected the paper where the author is listed as the first author or second author of the paper. For every author, the first paper we collected is ensured to be the most recent published from the author (simultaneously ensuring the author is the first author), which can minimize the possible bias that the paper is included as the pretraining data of the popular LLMs. To facilitate the extraction of relevant information such as authors, titles, abstracts, introductions or citations, we collect the paper link in the sample. Additionally, we assume that the authors' papers are published on the arxiv.org website. Hence, all the paper links are arxiv links. For every test, we use every author's first paper as the anonymous text input. The rest 4 papers are used for other analysis.

**Implementations.** Meta-Llama-3.1-8B-Instruct model was downloaded and deployed on 4-10 RTX 4090 GPUs, with the max new token set to 2000 to guarantee complete output. GPT-4o-mini was utilized through API request. The maximum input tokens for GPT-4o-mini is about 200K. During the searching process, if we collect the same article as our anonymous test input, we will ignore it as the assumption of our method is that the test paper should not appear on the website.

**Evaluations.** Our evaluations are divided into 2 parts, the first part is to examine the effectiveness of searching and collecting the true author from the crawling based on the generated keywords from LLMs. The second part is to examine the ability of the LLMs to identify the true author from the

	STAGE 1 (%)	STAGE 2 (%)						
MODEL		CON	TENT	WRITIN	NG STYLE	Сітатіо	N SIMILARITY	FINAL ACCURACY
		TOP 1	TOP 5	Top 1	Top 5	TOP 1	Top 5	(%)
LLAMA-3.1-8B	60.7	19.0	47.0	21.3	21.3	52.0	57.0	31.3
GPT-40-mini	56.3	27.3	52.3	36.3	61.3	49.3	66.0	41.3

Table 1: Test Results of Two Stages: Candidate Selection and Authorship Decision.

Methods	MODEL	WRITING STYLE		CITATION	FINAL SCORE (%)
		Тор1	Тор5	TOP1 TOP5	
LIP	Llama-3.1-8B	10.0	18.3	-	8.7
(HUANG ET AL., 2024)	GPT-40-mini	6.7	17.7	-	9.3
AIDBENCH	LLAMA-3.1-8B	9.7	18.0	-	10.0
(WEN ET AL., 2024)	GPT-40-mini	11.3	21.7	-	11.3
OUR METHOD	LLAMA-3.1-8B	8.7	16.0	12.3 20.0	11.6
OUR METHOD	GPT-40-mini	18.0	29.7	15.0 26.7	17.7

Table 2: Final Score Benchmark with Baselines Comparison.

candidates pool. If the first stage fails to collect the true author, we add the same authors' other papers from our dataset to the candidate pool and allow the LLMs to reattempt the stage 2 test.

Additionally, we conducted an overall score evaluation, as our proposed method is an automated, end-to-end process. This test measures the percentage of cases where the correct author is successfully selected and identified throughout the entire pipeline.

Multi-Conversation Handling. The input token limits of OpenAI API requests are 128,000 to 200,000 tokens. When the input token number exceeds the token limitations due to additional information such as citation and introduction, or a large number of papers collected in the candidate pool, we need to split the input into 2-4 bathes to guarantee input size is compatible with the token limits. We assume that API requests for LLMs do not support conversation memory. In this case, we need to handle the situation of multi-round conversation. We prompt the LLMs to decide on the topranked author from mini-batches and extract the corresponding information of these top-ranked authors from candidate pools again. This information is then further provided to the LLMs in subsequent rounds until the top 5 authors are identified for each



Figure 5: The final Score comparison across different models using our methods.

metric. In this way, we treat every API request as a new conversation and the complete information will be provided to ensure the decision is correct. 472

473

474

475

476

477

478

479

480

481

482

483

### 4.2 Results

First, the experiment results of Stage 1 and Stage 2 are summarized in Table 1. In Stage 1, we utilize our hierarchical levels to guide the models in emulating the keyword generation process. Each level serves as a basis for retrieving candidate authors, which are then used in Stage 2 for decision-making. Additionally, we incorporate citation information from the anonymous articles into our candidate

468

469

470

471

445

INPUT		FINAL			
	CONTENT	WRITING STYLE	CITATION SIMILARITY	ACCURACY (%)	
Abstract	$\checkmark$	$\checkmark$		26.0	
Abstract+Intro+Citation	$\checkmark$		$\checkmark$	39.7	
Abstract+Intro+Citation	$\checkmark$	$\checkmark$	$\checkmark$	41.3	
Abstract+Intro+Citation		$\checkmark$	$\checkmark$	44.3	

Table 3: Ablation experiment by using different input and different prompts. The experiment was conducted using GPT-4o-mini.

Table 4: Evaluating the accuracy of searching the correct author's other published articles in our stage 1. The keywords are generated by Llama-3.1-8B.

KEYWORD INSTRUCTION	Stage-1 Acc. (%)
NO INSTRUCTONS	7.3
FROM GENERAL TO SPECIFIC	12.0
WITH LEVELS	15.7
Self-Citation	55.3
WITH LEVELS + SELF-CITATION	60.7

pool to enhance selection accuracy. Using keywords generated by Llama-3.1-8B, we achieved an accuracy of 59.3%. Keywords emulated and generated by GPT-4o-mini are slightly less effective, yielding an accuracy of 56.3%.

In Stage 2, we collect information about the candidate authors and input this data into LLMs for reference-based evaluation using three key metrics: content similarity, writing style, and citation similarity. Across all three metrics, Top-5 accuracy consistently exceeded 50%. Among these metrics, writing style outperform content similarity in distinguishing authors. However, citation similarity achieves the highest accuracy, with Top-1 accuracy reaching 52% for Llama-3.1-8B and 49.3% for GPT-40-mini. Finally, by integrating these three metrics, our final decision accuracies are 31.3% for Llama-3.1-8B and 41.3% for GPT-40-mini.

We also experiment to evaluate the overall accuracy by combining Stage 1 and Stage 2 (Table 2). The final score indicates the probability of correctly identifying the author from searching to authorship-decision, Our method achieves the best results using GPT-40-mini (17.6%), outperforming the baseline models LIP (Huang et al., 2024) and AID-Bench (Wen et al., 2024) in both Llama and GPT models.

To validate the effectiveness of our proposed hierarchical keyword levels for candidate collection, we conduct ablation experiments on keyword analysis in Table 4. The results demonstrate that using different prompts leads to varying search performance. When applying our proposed levels, we achieve an accuracy rate of 15.7%. Furthermore, when combined with self-citation, the overall accuracy increases significantly to 60.7%. 511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

We also investigate the impact of incorporating the introduction and citation in author attribution in Table 3, finding that it significantly improves the LLMs' ability to identify the correct author, with accuracy increasing from 26% to approximately 40%. When prompting the model to perform inference based on content, writing style, and citation similarity, the results are slightly lower than the accuracy achieved using only writing style and citation similarity (44.3%).

Finally, in Figure 5, we achieve the highest overall score of 18.3% using the latest GPT-40 model.

## 5 Conclusion

In this paper, we introduce the first benchmark and dedicated solution for *Open-World Authorship Attribution*. Leveraging recent advancements in LLMs, we propose a two-stage pipeline: candidate selection and authorship decision. In the first stage, multi-levels keywords extracted from the target paper are used to search the Internet. The retrieved results, combined with citation lists, form a pool of potential candidates. In the second stage, LLMs infer authorship based on writing style and citation similarity from these candidates. Extensive experiments demonstrate the effectiveness and superiority of our approach over multiple potential baseline methods.

510

484

## 6 Limitations

547

576

578

581

582

583

584

585

586

587

588

589

590

591

592 593

594

595

596

Table & Figure Features. Another distinguishing feature in authorship attribution is the unique preferences authors showcase in structuring and de-550 signing their tables and figures. These characteris-551 tics manifest in various ways, including the choice 552 of color palettes, where some authors consistently favor specific hues or grayscale representations. 554 Differences also emerge in plotting styles, such 555 as the use of bar charts, scatter plots, heatmaps, or line graphs, along with variations in grid us-557 558 age, axis formatting, and legend placement. Labeling and annotation preferences also contribute to stylistic distinctions, as authors may differ in font choices, caption positioning, and the inclusion of callout markers. Additionally, the structuring of 562 tables varies, with some researchers favoring de-563 tailed grid layouts while others opt for minimalistic 564 designs with selective use of horizontal and vertical lines. Another notable characteristic is the numbering and referencing approach, with some authors 567 preferring "Figure 1" while others use "Fig. 1," 568 along with variations in how they cross-reference 569 visual elements within the text. In future work, we 570 aim to systematically analyze and quantify these 571 stylistic preferences, leveraging feature extraction techniques and deep learning models to explore how visual elements can enhance authorship attri-574 bution accuracy.

Large input. Since our method follows an openworld authorship attribution approach with an endto-end pipeline, it requires collecting a substantial amount of information as input for the LLMs. This often exceeds the maximum token limit of many models, which results in extra strategies to handle multi-turn conversations, as these models do not have built-in memory functions.

## References

- Shlomo Argamon, Moshe Koppel, James W. Pennebaker, and Jonathan Schler. 2009. Automatically profiling the author of an anonymous text. *Communications of the ACM*, 52(2):119–123.
- Benedikt T. Boenninghoff, Steffen Hessler, Dorothea Kolossa, and Robert M. Nickel. 2019. Explainable authorship verification in social media via attentionbased similarity learning. In *Proceedings of the 2019 IEEE International Conference on Big Data (Big Data)*, pages 36–45, Los Angeles, CA, USA. IEEE.
- Jing Cui, Yishi Xu, Zhewei Huang, Shuchang Zhou, Jianbin Jiao, and Junge Zhang. 2024. Recent ad-

vances in attack and defense approaches of large language models. *Preprint*, arXiv:2409.03274. 597

598

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

- Danda B. Rawat Desta Haileselassie Hagos, Rick Battle. 2024. Recent advances in generative ai and large language models: Current status, challenges, and perspectives. *arXiv preprint arXiv:2407.14962*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seo-

hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-670 ney Meers, Xavier Martinet, Xiaodong Wang, Xi-671 aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-674 schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing 677 Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, 679 Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei 680 Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew 684 Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, 691 Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, 694 Changkyu Kim, Chao Zhou, Chester Hu, Ching-695 Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, 701 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-703 ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, 704 Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat 705 Ozgenel, Francesco Caggioni, Frank Kanayet, Frank 706 Seide, Gabriela Medina Florez, Gabriella Schwarz, 707 Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanaz-710 eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun 711 Habeeb, Harrison Rudolph, Helen Suk, Henry As-712 pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim 713 Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, 714 Irina-Elena Veliche, Itai Gat, Jake Weissman, James 715 Geboski, James Kohli, Janice Lam, Japhet Asher, 716 Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy 717 Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe 718 719 Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, 720 721 Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-722 delwal, Katayoun Zand, Kathy Matosich, Kaushik

Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models. Preprint, arXiv:2407.21783.

723

724

725

726

727

730

731

732

733

734

735

736

737

738

740

741

742

743

744

745

747

748

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

781

782

- Desta Haileselassie Hagos, Rick Battle, and Danda B. Rawat. 2024. Recent advances in generative ai and large language models: Current status, challenges, and perspectives. *Preprint*, arXiv:2407.14962.
- Oren Halvani and Lukas Graner. 2021. Posnoise: An effective countermeasure against topic biases in authorship analysis. In *Proceedings of the 16th Inter-*

national Conference on Availability, Reliability and
Security, pages 1–12.

788

790

793

794

795

796

797

805

810

811

812

813

814

815

816

817

818

819

820 821

822

823

824

826

827

830

831

832 833

834

837

840

841

- Baixiang Huang, Canyu Chen, and Kai Shu. 2024. Can large language models identify authorship? *Preprint*, arXiv:2403.08213.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in Neural Information Processing Systems*, 35:22199– 22213.
- Tharindu Kumarage and Huan Liu. 2023. Neural authorship attribution: Stylometric analysis on large language models. *arXiv preprint arXiv:2308.07305*. Computation and Language (cs.CL); Artificial Intelligence (cs.AI).
  - Ksenia Lagutina, Nadezhda Lagutina, Elena Boychuk, Inna Vorontsova, Elena Shliakhtina, Olga Belyaeva, Ilya Paramonov, and P. G. Demidov. 2019. A survey on stylometric text features. In 2019 25th Conference of Open Innovations Association (FRUCT), pages 184–195. IEEE.
- Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. 2024. A comprehensive overview of large language models. *Preprint*, arXiv:2307.06435.
- Tempestt Neal, Kalaivani Sundararajan, Aneez Fatima, Yiming Yan, Yingfei Xiang, and Damon Woodard. 2017. Surveying stylometry techniques and applications. ACM Computing Surveys (CSuR), 50(6):1–36.
- OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine

Choi, Christine McLeavey, Christopher Hesse, Clau-842 dia Fischer, Clemens Winter, Coley Czarnecki, Colin 843 Jarvis, Colin Wei, Constantin Koumouzelis, Dane 844 Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, 845 David Carr, David Farhi, David Mely, David Robin-846 son, David Sasaki, Denny Jin, Dev Valladares, Dim-847 itris Tsipras, Doug Li, Duc Phong Nguyen, Duncan 848 Findlay, Edede Oiwoh, Edmund Wong, Ehsan As-849 dar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, 850 Eric Kramer, Eric Peterson, Eric Sigler, Eric Wal-851 lace, Eugene Brevdo, Evan Mays, Farzad Khorasani, 852 Felipe Petroski Such, Filippo Raso, Francis Zhang, 853 Fred von Lohmann, Freddie Sulit, Gabriel Goh, 854 Gene Oden, Geoff Salmon, Giulio Starace, Greg 855 Brockman, Hadi Salman, Haiming Bao, Haitang 856 Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, 857 Heather Whitney, Heewoo Jun, Hendrik Kirchner, 858 Henrique Ponde de Oliveira Pinto, Hongyu Ren, 859 Huiwen Chang, Hyung Won Chung, Ian Kivlichan, 860 Ian O'Connell, Ian O'Connell, Ian Osband, Ian Sil-861 ber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya 862 Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, 863 Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, 866 Jason Kwon, Jason Phang, Jason Teplitz, Jason 867 Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-868 avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui 869 Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, 870 Joaquin Quinonero Candela, Joe Beutler, Joe Lan-871 ders, Joel Parish, Johannes Heidecke, John Schul-872 man, Jonathan Lachman, Jonathan McKay, Jonathan 873 Uesato, Jonathan Ward, Jong Wook Kim, Joost 874 Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, 875 Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, 876 Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai 877 Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin 878 Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, 879 Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, 880 Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle 881 Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau-882 ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia 883 Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil-884 ian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kon-886 draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, 887 Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine 888 Boyd, Madeleine Thompson, Marat Dukhan, Mark 889 Chen, Mark Gray, Mark Hudnall, Marvin Zhang, 890 Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, 891 Max Johnson, Maya Shetty, Mayank Gupta, Meghan 892 Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao 893 Zhong, Mia Glaese, Mianna Chen, Michael Jan-894 ner, Michael Lampe, Michael Petrov, Michael Wu, 895 Michele Wang, Michelle Fradin, Michelle Pokrass, 896 Miguel Castro, Miguel Oom Temudo de Castro, 897 Mikhail Pavlov, Miles Brundage, Miles Wang, Mi-898 nal Khan, Mira Murati, Mo Bavarian, Molly Lin, 899 Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na-900 talie Cone, Natalie Staudacher, Natalie Summers, 901 Natan LaFontaine, Neil Chowdhury, Nick Ryder, 902 Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, 903 Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel 904 Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, 905

Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, 906 Olivier Godement, Owen Campbell-Moore, Patrick 907 Chao, Paul McMillan, Pavel Belov, Peng Su, Pe-908 ter Bak, Peter Bakkum, Peter Deng, Peter Dolan, 910 Peter Hoeschele, Peter Welinder, Phil Tillet, Philip 911 Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming 912 Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul 913 Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, 914 Reza Zamani, Ricky Wang, Rob Donnelly, Rob 915 Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-916 dani, Romain Huet, Rory Carmichael, Rowan Zellers, 917 Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan 918 Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, 919 Sam Toizer, Samuel Miserendino, Sandhini Agar-921 wal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu 923 Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shi-924 rong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stew-925 art Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, 927 Tejal Patwardhan, Thomas Cunninghman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, 931 Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, 933 Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne 934 Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, 937 Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and 938 Yury Malkov. 2024. Gpt-4o system card. Preprint, 940 arXiv:2410.21276.

> Ajay Patel, Delip Rao, Ansh Kothary, Kathleen McKeown, and Chris Callison-Burch. 2023. Learning interpretable style embeddings via prompting llms. *Preprint*, arXiv:2305.12696.

941 942

943

944

945

947

950

951

955

957

958

959

960

961

962 963

- Nektaria Potha and Efstathios Stamatatos. 2019. Improving author verification based on topic modeling. *Journal of the Association for Information Science and Technology*, 70(10):1074–1088.
- Libo Qin, Qiguang Chen, Xiachong Feng, Yang Wu, Yongheng Zhang, Yinghui Li, Min Li, Wanxiang Che, and Philip S. Yu. 2024. Large language models meet nlp: A survey. *Preprint*, arXiv:2405.12819.
- Abhay Sharma, Ananya Nandan, and Reetika Ralhan. 2018a. An investigation of supervised learning methods for authorship attribution in short hinglish texts using char & word n-grams. *arXiv preprint*, arXiv:1812.10281. Available at https://arxiv. org/abs/1812.10281.
- Abhay Sharma, Ananya Nandan, and Reetika Ralhan. 2018b. An investigation of supervised learning methods for authorship attribution in short hinglish texts using char & word n-grams. *arXiv preprint*, abs/1812.10281.

Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. 2023. A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2306.11695*. 964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1001

1002

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

- Kalaivani Sundararajan and Damon Woodard. 2018. What represents "style" in authorship attribution? In Proceedings of the 27th International Conference on Computational Linguistics (COLING), pages 2814– 2822.
- Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. 2023. The science of detecting llm-generated texts. *arXiv preprint arXiv:2303.07205*. Accessed: January 7, 2025.
- Adaku Uchendu, Thai Le, Kai Shu, and Dongwon Lee. 2020. Authorship attribution for neural text generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing* (*EMNLP*), pages 8384–8395.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.
- Zichen Wen, Dadi Guo, and Huishuai Zhang. 2024. Aidbench: A benchmark for evaluating the authorship identification capability of large language models. *arXiv preprint arXiv:2411.13226*.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. 2024. Sheared llama: Accelerating language model pre-training via structured pruning. In *Proceedings of the International Conference on Learning Representations (ICLR)*. Presented at ICLR 2024.
- Mengzhou Xia, Zexuan Zhong, and Danqi Chen. 2022. Structured pruning learns compact and accurate models. *arXiv preprint*, arXiv:2204.00408. Accepted to ACL 2022. The code and models are available at https://doi.org/10.48550/arXiv.2204.00408.
- Piji Li Xuanfan Ni. 2024. A systematic evaluation of large language models for natural language generation tasks. *arXiv preprint arXiv:2405.10251*.
- Xianjun Yang, Liangming Pan, Xuandong Zhao, Haifeng Chen, Linda Petzold, William Yang Wang, and Wei Cheng. 2023. A survey on detection of llmsgenerated content. *arXiv preprint arXiv:2310.15654*. Accessed: January 7, 2025.
- Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. 2024. Mmllms: Recent advances in multimodal large language models. *Preprint*, arXiv:2401.13601.

#### Appendix A

#### **Ethical Discussion** A.1

The application of this method may lead to unintended consequences, such as identifying authors during the Open Review stage of anonymous submissions, which would constitute an inappropriate and unethical use of our approach. Misuse of this technique in peer review processes could compromise the integrity of double-blind evaluation systems, introducing bias in scholarly assessments.

To mitigate such risks, we strongly advocate for the ethical application of our method in domains where it can serve as a tool for transparency, accountability, and integrity. These include plagiarism detection, where it can help identify unauthorized reproduction of content; authenticity verification, which ensures the legitimacy of texts to detect spam and fraudulent writing; and forensic linguistic analysis, where attribution techniques contribute to research integrity.

Furthermore, we emphasize the importance of responsible deployment, encouraging institutions, publishers, and AI practitioners to implement strict ethical guidelines when leveraging authorship attribution technologies. By doing so, we can ensure that such methods are used only in contexts that promote fairness, trust, and the credibility of research and publishing.

#### **Cost Discussion** A.2

During testing, the primary cost is associated with API access to OpenAI models. GPT-4o-mini and GPT-40 are priced at 0.15 and 2.50 USD per million input tokens, respectively. Due to the large input size generated by our candidate pools, each test round involving 300 authors incurs an estimated cost of 5 USD when using GPT-4o-mini and 50 USD per round when using GPT-4o.

#### **Effectiveness Analysis of Multi-Level Keywords Retrieval** A.3



Figure 6: Our proposed Multi-Levels keywords serve as the foundation for searching and collecting candidate authors. These keywords are used as example prompts for LLMs to generate emulated results. The Keywords Level Analysis evaluates the effectiveness of each multi-levels in the search process, comparing the results against those obtained without any guiding instructions.

1018

1020

1021

1022

1023

1024

1027

1028

1029

1030

1031

1032

1033

1035

1037

Groups	Categories	Groups	Categories
1	3D Vision and Reconstruction	16	Natural Language Processing (NLP)
2	Medical Imaging and Diagnostics	17	Scene Understanding and Parsing
3	Image Processing and Enhancement	18	Tracking and Re-Identification
4	Video Understanding and Generation	19	Federated and Distributed Learning
5	Vision-Language Models	20	AI Ethics and Explainability
6	Generative Models and Techniques	21	Physics and Scientific Topics
7	Object Detection and Recognition	22	Data Representation and Augmentation
8	Semantic and Instance Segmentation	23	Audio and Speech Processing
9	Adversarial Techniques and Robustness	24	Emotion and Human-Centric Applications
10	Optimization and Efficiency	25	Novel Applications and Emerging Topics
11	Robotics and Navigation	26	Large Language Models (LLM)
12	Graph Neural Networks and Hyperbolic Models	27	Neural Architecture Optimization
13	Multimodal and Hybrid Models	28	Other
14	Scientific Modeling and Mathematics	29	Deep Learning and Foundational Models
15	Self-Supervised and Semi-Supervised Learning		

## Table 5: Categories & Groups in Heatmap

Table 6: Full Instruction & Prompting used in decision stage based on 3 different metrics: Content, Writing Style and Citation Similarities.

Metrics	Instruction as Input
Content	I will provide the information of the anonymous article's title, abstract, dataset, introduc- tion or extra information, please remember them. Then, Please choose the top 5 possible articles' author(s) among all the candidates with their corresponding information. In this time, Decide the author(s) based on the Content like topics covered. You need to evaluate based on metrics focused on contents includes:(a) Preferred Topics: Common themes or subjects frequently addressed by the author. (b)Domain-Specific Terms: Use of jargon or technical language tied to the author's expertise. Now i will start to give you the list of candidates for you to decide!
Writing Style	Please choose the top 5 most possible articles' author(s) among all the candidates with their corresponding information. In this time Decide the author based on the writing style. Metrics for evaluation include: (a) Writing Tone: Formal, casual, emotional, or neutral tone in the text; (b) Repetition Patterns: Tendency to repeat certain ideas, phrases, or structures; (c) Complexity: Use of compound or complex sentences, and overall readability level; (d) Paragraph Structure: Length and organization of paragraphs; (e) Vocabulary Usage: Word choices, diversity, and domain-specific terms; (f) Punctuation Patterns: Frequency and style of punctuation usage; (g) Sentence Length and Variation: Average length and variability of sentences; (h) Personal Pronouns and Voice: Usage of pronouns and active/passive voice; (i) Lexical Density: Ratio of content words to function words; (j) Rhythm and Flow: Natural sentence progression and rhythm.
Citation Similarity	Please choose the top 5 most possible articles' author(s) among all the candidates with their corresponding information. In this time Decide the author based on the citations. Different papers with the same author tend to share similarities in references. Therefore, please refer to References and Sources: Citation patterns, including the types of resources cited (e.g., scholarly papers, blogs).

#### Examples

#### GPT-40 Output Analysis



Figure 7: The GPT-40 analysis was conducted on three selected texts. To minimize the influence of topic variation on the LLM's ability to determine authorship, all three texts were chosen to focus on the same topics: model pruning and large language models (LLMs). The first and third text are from the same author (Xia et al., 2022, 2024). Therefore the analysis shown LLMs' ability to distinguish and identify authors based on writing style.

Table 7: Full Example of Different Writing Styles. Example 1 (Xia et al., 2022), Example 2(Sun et al., 2023), Example 3(Xia et al., 2024).

Class	Abstract
Example 1	The growing size of neural language models has led to increased attention in model compression. The two predominant approaches are pruning, which gradually removes weights from a pre-trained model, and distillation, which trains a smaller compact model to match a larger one. Pruning methods can significantly reduce the model size but hardly achieve large speedups as distillation. However, distillation methods require large amounts of unlabeled data and are expensive to train. In this work, we propose a task-specific structured pruning method CoFi (Coarse- and Fine- grained Pruning), which delivers highly parallelizable subnetworks and matches the distillation methods in both accuracy and latency, without resorting to any unlabeled data. Our key insight is to jointly prune coarse-grained (e.g., layers) and fine-grained (e.g., heads and hidden units) modules, which controls the pruning decision of each parameter with masks of different granularity. We also devise a layerwise distillation strategy to transfer knowledge from unpruned to pruned models during optimization. Our experiments on GLUE and SQuAD datasets show that CoFi yields models with over 10x speedups with a small accuracy drop, showing its effectiveness and efficiency compared to previous pruning and distillation approaches.
Example 2	As their size increases, Large Language Models (LLMs) are natural candidates for network pruning methods: approaches that drop a subset of network weights while striving to preserve performance. Existing methods, however, require either retraining, which is rarely affordable for billion-scale LLMs, or solving a weight reconstruction problem reliant on second-order information, which may also be computationally expensive. In this paper, we introduce a novel, straightforward yet effective pruning method, termed Wanda (Pruning by Weights and Activations), designed to induce sparsity in pretrained LLMs. Motivated by the recent observation of emergent large magnitude features in LLMs, our approach prunes weights with the smallest magnitudes multiplied by the corresponding input activations, on a per-output basis. Notably, Wanda requires no retraining or weight update, and the pruned LLM can be used as is. We conduct a thorough evaluation of our method Wanda on LLaMA and LLaMA-2 across various language benchmarks. Wanda significantly outperforms the established baseline of magnitude pruning and performs competitively against recent methods involving intensive weight update.
Example 3	The popularity of LLaMA (Touvron et al., 2023a;b) and other recently emerged moderate-sized large language models (LLMs) highlights the potential of building smaller yet powerful LLMs. Regardless, the cost of training such models from scratch on trillions of tokens remains high. In this work, we study structured pruning as an effective means to develop smaller LLMs from pre-trained, larger models. Our approach employs two key techniques: (1) targeted structured pruning, which prunes a larger model to a specified target shape by removing layers, heads, and intermediate and hidden dimensions in an end-to-end manner, and (2) dynamic batch loading, which dynamically updates the composition of sampled data in each training batch based on varying losses across different domains. We demonstrate the efficacy of our approach by presenting the Sheared-LLaMA series, pruning the LLaMA2-7B model down to 1.3B and 2.7B parameters. Sheared-LLaMA models outperform state-of-the-art open-source models of equivalent sizes, such as Pythia, INCITE, and OpenLLaMA models, on a wide range of downstream and instruction tuning evaluations, while requiring only 3% of compute compared to training such models from scratch. This work provides compelling evidence that leveraging existing LLMs with structured pruning is a far more cost-effective approach for building smaller LLMs.