Bypassing the Decoding: Detecting Copyright Infringement through LLM Internal States

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have revolutionized Natural Language Processing through advanced text generation capabilities. However, their use raises legal and ethical concerns, particularly related to copyright infringement. While traditional methods assess the entire generated output for potential violations, this study introduces a novel framework that detects copyright risks by analyzing LLMs' internal states before any text is generated. This proactive approach enhances efficiency by identifying issues early in the generation process. To implement this framework, we used a dataset of literary works to derive both the LLMs' internal states and reference materials. These were used to train a neural network classifier capable of detecting potential copyright con-017 cerns. Additionally, this method helps prevent the unintended release of copyrighted content, offering an extra layer of protection. We also integrated this framework into a Retrieval-Augmented Generation (RAG) system, using FAISS (Facebook AI Similarity Search) and SQLite to efficiently manage reference texts. These texts are sourced from a protected copyright database, improving the accuracy and reliability of our detection process. By comparing generated content to known copyrighted material, our system ensures better compliance with legal and ethical standards. Overall, our findings demonstrate the value of analyzing internal states for proactive copyright monitoring, providing a scalable and effective solution for responsible AI-driven text generation.

1 Introduction

042

Large Language Models have revolutionized text generation and dialogue systems in Natural Language Processing with their advanced capabilities (Zhang et al., 2023; Li et al., 2022). However, as these models generate content that may inadvertently reproduce protected material, they raise significant challenges related to copyright infringe-



Figure 1: The process of predicting copyright infringement risk involves proactively preventing potential violations by analyzing a large language model's hidden states before content generation. The model encodes the user query, processes it through layers, and decodes the output. By extracting information from the hidden states at intermediate layers, we gain a detailed semantic representation, which is then analyzed to anticipate potential copyright risks. This approach enhances both efficiency and accuracy in predicting infringement before any content is generated.

ment, making it essential to ensure compliance with legal and ethical standards across various applications (Peng et al., 2023; Xue et al., 2021).

Previous studies have highlighted the risks of LLM-generated copyright infringement (Xu et al., 2024), focusing primarily on case studies that analyze verbatim reproduction using techniques such as Longest Common Subsequence (LCS) (Karamolegkou et al., 2023) or custom similarity metrics (Mueller et al., 2024). To mitigate such risks, various methods have been proposed, including the SHIELD Defense Mechanism (Liu et al., 2024), prompt engineering, and Memfree Decoding (Chen et al., 2024). However, these approaches share two fundamental limitations: first, they all rely on decoding the full generated text before detection can occur, which introduces significant computational inefficiencies and delays. Second, de-

079

101

102

104

105

106

108

109

110

111

112

061

tecting risks only after the text has been generated exposes the possibility of disseminating inappropriate or infringing content, and once the text is released, the damage becomes irreversible. These challenges underscore the need for a more proactive and efficient approach that can detect copyright risks before full text generation.

To bridge this gap, our study introduces an innovative framework, the Internal State Analyzer for Copyright (ISAC), for assessing the risk of copyright infringement in LLM-generated text. ISAC utilizes internal states from the prefill phase of LLMs to evaluate potential infringement risks before any text is generated. Unlike traditional methods that require generating entire outputs, ISAC proactively detects potential copyright violations by analyzing the model's early-stage representations of input text, which encode the semantic and structural properties of the input. This approach offers a scalable, real-time, and precise risk assessment mechanism that strengthens copyright compliance in AI-generated content without needing to generate the full output. To enhance detection performance, ISAC is integrated into a Retrieval-Augmented Generation system. In this system, input text and reference counterparts are indexed using FAISS and stored in SQLite, enabling efficient retrieval of relevant texts during infringement risk evaluation. When a relevant reference is identified, it is concatenated with the models internal states to assess the likelihood of copyright infringement. This retrieval mechanism significantly improves the models ability to compare generated content against known copyrighted material, boosting both precision and real-time detection efficiency while ensuring compliance with legal and ethical standards.

In a certain series of experimental configurations, ISAC delivered impressive results, achieving various accuracy and F1 scores. Specifically, the accuracy ranged from 91.88% to 95.05%, and F1 scores varied between 0.9249 and 0.9468. In some configurations, ISAC even reached near-perfect detection rates. These results demonstrate ISAC's consistent ability to accurately identify potential copyright violations across multiple settings, maintaining high precision and recall. The findings emphasize ISAC's robustness in real-time, scalable risk detection for LLM-generated content, even without generating any text. For a detailed description of the experimental setup and results, please refer to Section 4.3. Our primary contributions are as follows:

• As illustrated in Figure 1, we propose a real-time framework "ISAC" for predicting copyright infringement in LLM-generated text by leveraging internal states extracted before any token is decoded, ensuring efficiency without relying on output generation. 113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

- ISAC is the first framework to proactively detect potential copyright violations by analyzing LLM internal states before content is generated. This approach ensures that neither users nor language models are exposed to any information, such as prompts, that could lead to infringement, thereby ensuring compliance with legal and ethical standards.
- We validate the ISACs effectiveness in largescale text generation scenarios and demonstrate its integration with a RAG system. This integration enables efficient and accurate text retrieval, making the approach suitable for industrial applications requiring real-time copyright compliance.

2 Related Work

2.1 Copyright Issues with LLMs

Scholars have emphasized the importance of protecting the intellectual property associated with the parameters of Large Language Models (Peng et al., 2023; Xue et al., 2021). This concern arises from the substantial investments in resources required for training LLMs, as well as the risk of unauthorized exploitation of these models, which can have significant economic and ethical implications (Zhang et al., 2018; He et al., 2022; Dale, 2021).

Copyright concerns are not limited to text; they span across various digital content creation formats, including scripts, images, videos (Moayeri et al., 2024; Kim et al., 2024), and code (Yu et al., 2023). This widespread impact underscores the urgency of addressing these complex issues (Lucchi, 2023).

2.2 Detecting Copyright Issues in Training Data

LLMs are capable of retaining and reproducing significant parts of their training datasets, which may include copyrighted materials and sensitive data (Karamolegkou et al., 2023; Carlini et al., 2019; Lee et al., 2023; Carlini et al., 2022; Kandpal et al., 2022). The potential for such memorization poses significant copyright infringement concerns, especially as these models scale up and face extraction attacks (Carlini et al., 2021; Ozdayi et al., 2023;



Figure 2: Overview of our Copyright Infringement Detection Framework: Our approach involves maintaining a database of copyright-protected materials to support the analysis of LLM hidden states. During inference, this database provides reference samples for potential violations, working in conjunction with the model's hidden states to predict whether the generated content poses a risk of copyright infringement. The pipeline is structured into three key stages: The left section focuses on the construction and extraction of data for Retrieval-Augmented Generation, a core component designed to enhance model performance and address copyright-related challenges. The right section illustrates the generation of training data, including the collection of internal states, labels, and reference embeddings, which are then used to train a Multi-Layer Perceptron as the final infringement risk detector. Lastly, the bottom section showcases real-world user interaction, where queries are submitted, and the system applies our framework to assess potential infringement risks effectively.

Chao et al., 2023; Ishihara, 2023).

162

163

165

166

168

170

173

174

175

176

177

178

To combat this, innovative strategies such as "copyright traps" have been introduced to detect copyrighted content in LLM training datasets (Shilov et al., 2024; Shi et al., 2023; Meeus et al., 2024). Studies have also investigated the likelihood of LLMs generating exact or near-verbatim copyrighted content and have quantified the legal risks associated with such reproductions (Carlini et al., 2021; Lee et al., 2021). Building on these efforts, our work explores how to quickly and accurately determine whether an LLM will generate copyrighted content.

2.3 Mitigating Copyright Issues in Model Serving

To mitigate copyright risks during model serving, recent studies have developed real-time intervention mechanisms. SHIELD (Liu et al., 2024) uses agent-based defenses and N-Gram models to dynamically verify copyright status, preventing copyrighted text generation while maintaining quality. MemFree Decoding (Chen et al., 2024) prevents verbatim copying during inference but fails to address non-literal copying, such as event or character overlaps, and may introduce hallucinations.

179

180

181

182

183

184

187

188

190

191

193

195

3 Internal State Judge: Detecting Copyright Infringement Before Decoding

3.1 Problem Formulation

The issue of copyright infringement in content generated by LLMs has attracted significant attention from both industry and academia. Existing approaches focus on detecting potential copyright violations only after the content has been generated.

252

253

254

257

258

259

260

261

262

263

264

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

287

This post-generation evaluation method presents several challenges, including high computational costs, delays in enforcement, and legal risks due to temporary exposure to infringing material.

196

197

198

199

201

202

204

207

210

211

213

214

215

216

217

218

219

221

224

227

229

236

In this paper, we present a framework (ISAC) designed to assess the risk of copyright infringement before an LLM generates any output. The inference process of an LLM for a given query can be divided into two phases: (1) **Prefill Phase**: The LLM processes the entire input query to create internal states. (2) **Decode Phase**: The LLM generates output based on these prefilled internal states.

This division of two phases leads us to the central question of our study: *Can the internal states produced during the prefill phase be used to predict the risk of copyright infringement before the decoding phase begins?*

To address this question, we argue that LLM's internal states of a query during the prefill phase capture critical contextual information linked to the likelihood of generating infringing content. We introduce an internal states judge designed to classify the copyright infringement risk of a query based on its internal states in this phase.

This method offers three key benefits:

- Efficiency: By evaluating internal states early in the prefill process, our approach can halt decoding if the internal states judge identifies potential risks, reducing unnecessary computational costs.
- **Proactive Copyright Compliance:** Our approach perform risk assessment occurs before content generation enables preventive actions rather than post-hoc interventions.
- Scalability: The internal states judge is designed to be adaptable across various LLM architectures and model sizes, facilitating broad deployment.

The following sections describe the design of the internal state judge, the methodology for training data collection, and the experimental evaluation of our approach.

3.2 Training An Internal States Judge

238Training Data Preparation.We construct a239dataset of triplets to train the classifier: (input, out-240put, reference).The input x is the query fed into241the LLM, the output y is the generated text, and the242reference t is the ground-truth continuation from243the source. Each generated output is assigned a risk

label based on its similarity to the reference text using the Rouge-L score:

$$\mathcal{H}^{\text{train}} = \mathcal{T}(j, \text{Rouge-L}(t, y)) \tag{1}$$

where the threshold-based function \mathcal{T} determines risk labels, and *j* represents the partitioning criterion:

$$\mathcal{T}(j, \text{Rouge-L}) = \begin{cases} 0, \text{ if } \mathcal{P}_2 \leq \text{Rouge-L} \leq 1\\ 1, \text{ if } 0 \leq \text{Rouge-L} \leq \mathcal{P}_1\\ \text{undefined, otherwise} \end{cases}$$
(2)

where \mathcal{P}_1 and \mathcal{P}_2 are predefined thresholds used to classify an output as either high or low risk.

Our dataset is structured as pairs of internal states and their associated risk labels: $\mathcal{D}_{\theta} = \{\langle S_{x_i}^{\text{train}}, \mathcal{H}_i^{\text{train}} \rangle\}_{i=1}^N$.

Internal States of Query in Prefill Phase of LLMs. A crucial step in ISAC is the extraction of internal states during the prefill phase of LLMs. In this phase, the model processes the entire input sequence to compute intermediate representations (such as keys and values) before generating any output tokens. This stage involves highly parallelized matrix-matrix operations, allowing the model to efficiently encode the semantic and structural properties of the input.

During forward propagation, the input text x from the dataset triplet is fed into the LLM, and we extract the internal states S from a specific layer in the prefill phase. These internal states are computed through multiple layers of non-linear transformations, activations, and information flow, formally represented as:

$$\mathcal{S}_{l} = f\left(\mathcal{W}_{l} \cdot \mathcal{S}_{l-1} + \mathcal{B}_{l}\right), \quad l = 1, 2, \dots, L \quad (3)$$

where S_l represents the internal states at layer l, W_l and B_l are the learnable weights and biases of the *l*-th layer, and *f* is the activation function. At each layer, the model refines its understanding of the input query *x*, progressively building increasingly sophisticated representations of syntax, context, and meaning (Devlin et al., 2019; Radford and Narasimhan, 2018). These internal states encode both token-level details and broader semantic relationships, providing a rich representation of the inputs meaning (Clark et al., 2019).

In our experiments, we extract internal states from the final encoder layer during the prefill phase and compute their mean across all tokens. This

334

335

336

338

340

341

342

343

344

345

346

347

348

349

350

351

354

355

provides a concise yet informative representation of the inputs semantics, effectively capturing both 289 local and contextual information. We hypothesize that these representations contain early indicators of potential copyright violations based on input queries. By analyzing these internal states before the decoding stage, we aim to proactively identify 294 and mitigate potential risks (Zellers et al., 2020).

296

297

298

310

311

312

313

314

315

317

319

320

324

328

Training Objectives of Internal States Judge. The objective of training the internal states judge is to create a classifier that predicts the likelihood of copyright infringement based on the internal states of the model. This classifier learns to assess the Rouge-L similarity score, distinguishing between high-risk and low-risk outputs. It is implemented using an MLP model:

$$\mathcal{M} = \operatorname{down}(\operatorname{up}(\mathcal{S}) \times \operatorname{SiLU}(\operatorname{gate}(\mathcal{S}))) \quad (4)$$

where SiLU serves as the activation function, and the linear layers down, up, and gate handle projection and gating mechanisms. This model enables efficient real-time risk prediction without requiring full output decoding.

Enhancing Internal States Judge with 3.3 **Retrieved References**

Leveraging References to Enhance Internal **States Judge.** Relying solely on input text may lack sufficient context for detecting copyright infringement. To improve detection, ISAC incorporates external references using RAG technology (Lewis et al., 2021), enhancing the model's ability to assess potential risks.

Formally, given an input query x, we first extract its internal states S_x from the prefill phase of the LLM, then retrieve a set of relevant reference texts $T = \{t_1, t_2, \dots, t_m\}$ from an external knowledge base. The retrieved references are encoded into an aggregated representation S_T , which is then concatenated with S_x to form the final combined representation. An MLP classifier is then applied to predict the infringement probability:

 $p = \sigma \left(\mathcal{M} \left(\text{concat} \left(f_{\theta}(x), h_{\phi}(\mathcal{G}(x)) \right) \right) \right),$

where f_{θ} represents the transformation function of the LLMs prefill phase, G is the retrieval function that selects references most relevant to x, h_{ϕ} encodes the retrieved references, \mathcal{M} denotes the 332

MLP model, and σ represents the sigmoid activation function that outputs the probability of copyright infringement.

Finally, the predicted probability p is compared with a predefined threshold τ to make the final infringement risk decision:

$$\mathcal{H}^{\text{predict}} = \begin{cases} 1, & \text{if } p \ge \tau \\ 0, & \text{otherwise} \end{cases}$$
(6)

where τ is a tunable threshold that determines the sensitivity of infringement detection. By integrating external references into the internal state analysis and applying a threshold-based decision rule, this enhanced approach significantly improves the models predictive capabilities, reducing both false positives and false negatives.

Retrieving References from Indexed Documents. To facilitate Retrieval-Augmented Generation, as shown in Figure 3, we construct a RAG-Enhanced Reference Database that efficiently stores and retrieves references for infringement detection. This database is designed to manage copyright materials effectively, ensuring quick access to relevant references and supporting robust content analysis and decision-making. The construction details of the RAG-based database are provided in Appendix B.



Figure 3: Process of constructing a vector database for the RAG system and handling user queries.

4 Experiments

In this section, we evaluate the effectiveness of the internal states judge in identifying literal copying

(5)

357

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

409

410

in text continuations. Specifically, we address the following research questions (RQ):

- **RQ 1:** How well does our method detect literal copying across various LLMs, such as the Llama and Mistral series, and how does model size influence performance?
- **RQ 2:** Can our method accurately identify nonliteral copying, such as paraphrased content, and how does its performance compare to that of literal copying detection?
- **RQ 3:** What factors affect the performance of our method, including the role of the RAG system, the choice of LLM internal state layers, and the strategies used for dataset division?

To investigate these questions, we conduct experiments using a structured dataset that includes both literal and non-literal copying tasks. For literal copying, we evaluate the risk of copyright infringement in text continuations by using excerpts from well-known fiction books. For non-literal copying, we focus on identifying event and character copying within paraphrased content. We test our method on LLMs from the Llama and Mistral series, ranging from 7B to 70B parameters, and compare it with baseline approaches. Our findings show that our method is both effective and accurate in detecting literal and non-literal copying, while also revealing the challenges involved in identifying paraphrased content.

4.1 Dataset

367

371

373

374

378

379

390

395

400

401

402

403

404

405

406

407

408

We used the COPYBENCH dataset (Chen et al., 2024) to evaluate LLM infringement risks on fiction texts (Meeus et al., 2024; Chang et al., 2023; Shi et al., 2023).

4.2 Model Selection

We used LLMs from the Llama (Touvron et al., 2023) and Mistral (Jiang et al., 2023) series to generate text continuations and extract internal states, ensuring accurate dataset classification. To capture true continuations, we extracted reference embeddings using BERT (Devlin et al., 2019), which effectively captured the semantic content for training.

4.3 Detecting Literal Copying through LLM Internal States

In this section, we empirically evaluate the effectiveness of our method for detecting literal copying across different LLMs, including Llama and Mistral, as well as a range of model sizes from 7B to 70B parameters. To assess model performance, we use standard metrics such as Accuracy and F1score, described in appendix C, providing insights into the models' precision and effectiveness in detecting infringement risks. Our approach involves extracting internal states from the last layer of the model during the pre-filling phase, which are then used to train a classifier for predicting copyright risk.

Baselines. In our experiment, we established a baseline model using LLMs to assess potential infringement in content generation tasks. It includes two configurations: "Input Only" (LLM-w/oRAG), where decisions are made based solely on the input text, and "Input with RAG system" (LLM-w/RAG), where both the input text and reference materials are considered. Similar to our proposed method, the baseline evaluates potential infringement without generating the next text segment. The task is to identify whether the continuation text contains elements that may raise infringement concerns. Predicted outcomes are compared to ground truth labels, which are derived from the dataset and based on Rouge-L scores. Details of the baseline prompt settings are provided in Table 9.

Results and Analysis. The results are based on three dataset splits, determined by Rouge-L scores: 10%, 20%, and 30%. Each split classifies the dataset into high-scoring (infringing) and low-scoring (non-infringing) samples. We assess the model's ability to distinguish between these groups and examine how incorporating reference embeddingsretrieved from a databaseenhances performance across various levels of textual similarity.

We also compare our method to the "LLM as Judge" approach. As shown in Table 1, we analyze the performance differences across dataset splits and model configurations, demonstrating the practical advantages of our approach.

Several key insights emerge from the analysis. First, our method significantly improves efficiency. The pre-trained MLP-based binary classifier provides faster inference and better accuracy compared to the "LLM as Judge" method, which relies on direct LLM predictions. This indicates that our approach is not only more efficient but also more precise in identifying potential copyright infringement. Second, using original reference text retrieved from the database during training enhances accuracy, outperforming models that rely solely on LLM-extracted internal states. This highlights the importance of external reference ma-

Table 1: The results on the literal dataset evaluate the performance of various models and methods. We compare four approaches: LLM-w/oRAG and LLM-w/RAG, which represent the "LLM as Judge (Without RAG system)" and "LLM as Judge (With RAG system)" methods. In these approaches, we use the LLM directly to detect potential copyright infringement in the input texteither based solely on the input (LLM-w/oRAG) or using both the input and the RAG system (LLM-w/RAG). Additionally, we evaluate the Internal States Judge (IS) methods: IS-w/oRAG and IS-w/RAG, which represent the "Internal States Judge (Without RAG system)" and "Internal States Judge (With RAG system)" methods. We report accuracy (ACC) and F1 scores for dataset divisions at 10%, 20%, and 30%.

| | | | Division (10%) | | Division (20%) | | Division (30%) | |
|-----------------|------------|---------------|----------------|---------------|----------------|--------------|----------------|--------------|
| LLMs | Method | Time (s) | ACC (%) | F1 (%) | ACC (%) | F1 (%) | ACC (%) | F1 (%) |
| | | | Llam | a | | | | |
| Llama-3.1-8B | LLM-w/oRAG | 0.4914 | 57.11 | 49.81 | 53.41 | 48.74 | 50.45 | 45.29 |
| | IS-w/oRAG | 0.0564 | 91.53 | 92.96 | 78.05 | 79.25 | 73.73 | 77.36 |
| | LLM-w/RAG | 0.7012 | 63.14 | 62.58 | 58.57 | 58.66 | 54.59 | 54.43 |
| | IS-w/RAG | 0.0592 | 92.37 | 93.71 | 83.26 | 82.67 | 77.11 | 78.62 |
| Llama-2-13b | LLM-w/oRAG | 0.5412 | 61.41 | 50.93 | 60.02 | 50.74 | 57.94 | 48.73 |
| | IS-w/oRAG | 0.0642 | 91.75 | 93.37 | 82.46 | 81.47 | 78.83 | 76.44 |
| | LLM-w/RAG | 0.8109 | 63.21 | 61.54 | 60.40 | 60.71 | 58.62 | 55.63 |
| | IS-w/RAG | 0.0696 | 93.23 | 94.18 | 86.52 | 85.57 | 80.03 | 79.15 |
| Llama-3.1-70B | LLM-w/oRAG | 1.1492 | 62.23 | 51.48 | 57.89 | 50.22 | 58.23 | 49.68 |
| | IS-w/oRAG | 0.1274 | 100.00 | 100.00 | 94.55 | 94.63 | 91.88 | 92.49 |
| | LLM-w/RAG | 1.4335 | 65.43 | 64.12 | 62.67 | 60.78 | 60.89 | 62.45 |
| | IS-w/RAG | 0.1389 | 100.00 | 100.00 | 95.05 | 94.68 | 94.48 | 94.64 |
| Mistral | | | | | | | | |
| Mistral-7B-v0.1 | LLM-w/oRAG | 0.5238 | 52.90 | 51.85 | 49.01 | 48.85 | 50.67 | 51.12 |
| | IS-w/oRAG | 0.0623 | 97.96 | 98.00 | 79.58 | 82.97 | 70.75 | 76.24 |
| | LLM-w/RAG | 0.6876 | 58.52 | 54.36 | 55.49 | 52.44 | 51.87 | 53.23 |
| | IS-w/RAG | 0.0677 | 98.98 | 98.99 | 83.25 | 85.59 | 78.01 | 82.35 |
| Mistral-7B-v0.3 | LLM-w/oRAG | 0.5324 | 53.78 | 51.56 | 53.23 | 52.90 | 51.67 | 40.52 |
| | IS-w/oRAG | 0.0597 | 91.75 | 92.59 | 83.52 | 84.21 | 79.46 | 83.04 |
| | LLM-w/RAG | 0.6343 | 57.45 | 55.60 | 53.29 | 54.78 | 53.13 | 48.75 |
| | IS-w/RAG | 0.0614 | 93.76 | 95.30 | 87.27 | 86.24 | 84.86 | 87.39 |

terial, which offers richer context and enables the model to more accurately detect potential copyright violations. Additionally, we observe that the performance of different LLMs varies. Larger Llama models are more sensitive to infringement, suggesting that their increased size allows them to better capture subtle text similarities. In contrast, Llama and Mistral models show different capabilities in capturing textual nuances, which affects their effectiveness in this task. Finally, the dataset division strategy plays a key role. Larger Rouge score differences between high- and low-scoring samples make it easier for the model to differentiate between them. This emphasizes the importance of carefully selecting dataset splits, as they have a significant impact on the model's ability to accurately identify infringement risks.

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477 Variability in FN & FP Rates, but Stable Over478 all Accuracy & F1. To further analyze model
479 performance, we selected four representative con480 figurations and generated confusion matrix plots, as
481 shown in Figure 4. These configurations combine
482 two factors: the model (Llama-3.1-8B or Llama-

3.1-70B) and whether a reference is included, with the Rouge-L 30% split strategy applied.

Its important to note that the figures shown here represent a single instance from repeated experiments. Since the training and test sets are randomly split, some variability in the False Negative (FN) and False Positive (FP) rates is expected. However, despite this variability, we found that the overall prediction accuracy and F1 score remain consistently stable across different runs. This suggests that, while there are fluctuations in specific error types, the model's overall performance is reliable and robust.

Time Efficiency Comparison. We conducted experiments to compare the time efficiency of infringement prediction methods, and the results show that the proposed methods using internal states (IS-w/oRAG and IS-w/RAG) are significantly faster than the traditional basic method. In the basic method, each input text is processed sequentially by the LLM to generate the next segment, which is then compared with the reference text to assess potential infringement. The majority

483

495 496

494

497 498

499

500

501

502

503

504



Figure 4: Confusion matrix plots showing the effect of model size and RAG system on prediction performance, with Llama-3.1-8B and Llama-3.1-70B models, both with and without reference information, using a Rouge-L 30% threshold for dataset splitting.

of the time in this approach is spent on text generation, while the comparison step takes up very little time. As a result, the basic method is much slower, as indicated by its higher time values compared to the internal states-based methods. These methods streamline the process, eliminating the need for text generation and leading to faster, more efficient predictions. The detailed results of this comparison are shown in Table 2.

506

507

510

511

512

513

514

Table 2: This table shows the average time efficiency comparison (in seconds) for infringement prediction based on a single data point, testing three methods: predicting infringement using internal states without (IS-w/oRAG) and with (IS-w/RAG) RAG system, and the basic method of generating continuation text and comparing it with reference text.

| Method Model | Basic | IS-w/oRAG | IS-w/RAG |
|-----------------|--------|-----------|----------|
| Llama-3.1-8B | 0.4319 | 0.0564 | 0.0592 |
| Llama-2-13b | 0.6584 | 0.0642 | 0.0696 |
| Llama-3.1-70B | 1.6796 | 0.1274 | 0.1389 |
| Mistral-7B-v0.1 | 0.3571 | 0.0623 | 0.0677 |
| Mistral-7B-v0.3 | 0.3463 | 0.0597 | 0.0614 |

5 Conclusion and Future Work

This study presents a new framework "ISAC" for detecting potential copyright infringement in text generated by LLMs by analyzing their internal states. Unlike traditional methods that require decoding the generated output, our approach uses internal states to enable real-time detection, improving efficiency. Experiments with models like Llama and Mistral show that larger models achieve higher classification accuracy due to more detailed internal representations. By integrating RAG with FAISS for vector search and SQLite for structured storage, ISAC enhances retrieval and prediction reliability. This method strikes a balance between computational efficiency and legal compliance. 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

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

Future work will focus on expanding the framework to address non-literal copyright concerns, such as conceptual similarity and paraphrasing, and refining the classifier to improve robustness across different model sizes. We will also explore ways to enhance the interpretability of internal states to build trust and improve explainability. Collaborations with legal experts will be essential to ensure alignment with evolving copyright laws. Additionally, we plan to create an LLM agent that actively prevents copyright infringement by crossreferencing generated text against a curated corpus of licensed or public-domain material. This agent will help ensure compliance with copyright guidelines in real-time, providing a practical solution for applications focused on legal compliance.

Limitations

Despite its advantages, ISAC has some limitations. Detection accuracy in smaller models requires improvement, as these models often have less nuanced internal representations, which can affect reliability. Moreover, this study focuses mainly on assessing the ability of LLM internal states to identify copyright infringement, but more precise criteria for determining infringement are needed for practical applications. In particular, clearer standards are required to address complex cases like conceptual similarity or paraphrasing.

Ethics Statement

We all comply with the ACL Ethics Policy¹ during our study. All datasets used contain anonymized consumer data, ensuring strict privacy protections.

¹https://www.aclweb.org/portal/content/ acl-code-ethics

References

562

569

576

577

584

586

592

594

596

597

598

601

606

610

611

612

613

614

615

- Amos Azaria and Tom Mitchell. 2023. The internal state of an llm knows when it's lying. *Preprint*, arXiv:2304.13734.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2022. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*.
- Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks. In 28th USENIX security symposium (USENIX security 19), pages 267–284.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting training data from large language models. In 30th USENIX Security Symposium (USENIX Security 21), pages 2633–2650.
- Kent Chang, Mackenzie Cramer, Sandeep Soni, and David Bamman. 2023. Speak, memory: An archaeology of books known to ChatGPT/GPT-4. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 7312–7327, Singapore. Association for Computational Linguistics.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. 2023. Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*.
- Tong Chen, Akari Asai, Niloofar Mireshghallah, Sewon Min, James Grimmelmann, Yejin Choi, Hannaneh Hajishirzi, Luke Zettlemoyer, and Pang Wei Koh. 2024. Copybench: Measuring literal and non-literal reproduction of copyright-protected text in language model generation. *Preprint*, arXiv:2407.07087.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of BERT's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Robert Dale. 2021. Gpt-3: Whats it good for? *Natural Language Engineering*, 27(1):113–118.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library. 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

658

659

660

661

662

663

664

665

666

667

668

669

- Xuanli He, Qiongkai Xu, Lingjuan Lyu, Fangzhao Wu, and Chenguang Wang. 2022. Protecting intellectual property of language generation apis with lexical watermark. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10758–10766.
- Shotaro Ishihara. 2023. Training data extraction from pre-trained language models: A survey. *arXiv* preprint arXiv:2305.16157.
- Ziwei Ji, Delong Chen, Etsuko Ishii, Samuel Cahyawijaya, Yejin Bang, Bryan Wilie, and Pascale Fung. 2024. Llm internal states reveal hallucination risk faced with a query. *arXiv preprint arXiv:2407.03282*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. In *International Conference on Machine Learning*, pages 10697–10707. PMLR.
- Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. 2023. Copyright violations and large language models. *arXiv preprint arXiv:2310.13771*.
- Minseon Kim, Hyomin Lee, Boqing Gong, Huishuai Zhang, and Sung Ju Hwang. 2024. Automatic jailbreaking of the text-to-image generative ai systems. *arXiv preprint arXiv:2405.16567*.
- Jooyoung Lee, Thai Le, Jinghui Chen, and Dongwon Lee. 2023. Do language models plagiarize? In *Proceedings of the ACM Web Conference 2023*, pages 3637–3647.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2021. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. Retrieval-augmented generation for knowledgeintensive nlp tasks. *Preprint*, arXiv:2005.11401.
- Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2022. Pretrained language models for text generation: A survey. *Preprint*, arXiv:2201.05273.

727

728

Xiaoze Liu, Ting Sun, Tianyang Xu, Feijie Wu, Cunxiang Wang, Xiaoqian Wang, and Jing Gao. 2024. Shield: Evaluation and defense strategies for copyright compliance in llm text generation. *Preprint*, arXiv:2406.12975.

671

672

673

675

684

685

687

700

702

704

705

706

710

711

713

714

715

716

717

718

719

721

722

724

725

726

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
 Roberta: A robustly optimized bert pretraining approach. *Preprint*, arXiv:1907.11692.
- Nicola Lucchi. 2023. Chatgpt: a case study on copyright challenges for generative artificial intelligence systems. *European Journal of Risk Regulation*, pages 1–23.
- Matthieu Meeus, Igor Shilov, Manuel Faysse, and Yves-Alexandre de Montjoye. 2024. Copyright traps for large language models. *arXiv preprint arXiv:2402.09363*.
- Mazda Moayeri, Samyadeep Basu, Sriram Balasubramanian, Priyatham Kattakinda, Atoosa Chengini, Robert Brauneis, and Soheil Feizi. 2024. Rethinking artistic copyright infringements in the era of text-to-image generative models. *arXiv preprint arXiv:2404.08030*.
- Felix B Mueller, Rebekka Görge, Anna K Bernzen, Janna C Pirk, and Maximilian Poretschkin. 2024. Llms and memorization: On quality and specificity of copyright compliance. arXiv preprint arXiv:2405.18492.
- Mustafa Safa Ozdayi, Charith Peris, Jack FitzGerald, Christophe Dupuy, Jimit Majmudar, Haidar Khan, Rahil Parikh, and Rahul Gupta. 2023. Controlling the extraction of memorized data from large language models via prompt-tuning. *arXiv preprint arXiv:2305.11759*.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. *Preprint*, arXiv:1912.01703.
- Wenjun Peng, Jingwei Yi, Fangzhao Wu, Shangxi Wu, Bin Zhu, Lingjuan Lyu, Binxing Jiao, Tong Xu, Guangzhong Sun, and Xing Xie. 2023. Are you copying my model? protecting the copyright of large language models for eaas via backdoor watermark. *arXiv preprint arXiv:2305.10036*.
- Alec Radford and Karthik Narasimhan. 2018. Improving language understanding by generative pretraining.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. Detecting pretraining

data from large language models. *arXiv preprint arXiv:2310.16789*.

- Igor Shilov, Matthieu Meeus, and Yves-Alexandre de Montjoye. 2024. Mosaic memory: Fuzzy duplication in copyright traps for large language models. *arXiv preprint arXiv:2405.15523*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface's transformers: State-of-the-art natural language processing. *Preprint*, arXiv:1910.03771.
- Jialiang Xu, Shenglan Li, Zhaozhuo Xu, and Denghui Zhang. 2024. Do llms know to respect copyright notice? *arXiv preprint arXiv:2411.01136*.
- Mingfu Xue, Yushu Zhang, Jian Wang, and Weiqiang Liu. 2021. Intellectual property protection for deep learning models: Taxonomy, methods, attacks, and evaluations. *IEEE Transactions on Artificial Intelli*gence, 3(6):908–923.
- Zhiyuan Yu, Yuhao Wu, Ning Zhang, Chenguang Wang, Yevgeniy Vorobeychik, and Chaowei Xiao. 2023. Codeipprompt: Intellectual property infringement assessment of code language models. In *International Conference on Machine Learning*, pages 40373–40389. PMLR.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2020. Defending against neural fake news. *Preprint*, arXiv:1905.12616.
- Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. 2023. A survey of controllable text generation using transformer-based pre-trained language models. *Preprint*, arXiv:2201.05337.
- Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph Stoecklin, Heqing Huang, and Ian Molloy. 2018. Protecting intellectual property of deep neural networks with watermarking. In *Proceedings of the* 2018 on Asia conference on computer and communications security, pages 159–172.

A Implementation Details

The input dimension of our classifier is defined by the number of features in the training dataset, ensuring that the model can properly process the input

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

829

830

791

781

782

786

787

790

GPUs.

B

792 793

79 79

79

797

799

80

80

804 805

806 807

80

80

810 811

812

813

81

814 815

816 817

8

818 819

8

822 823

824 825

826

827 828 where C is the set of centroids and C_k is the set of embeddings in cluster k. During retrieval, a query

data. The hidden dimension is fixed at 256, a value

that aligns with the design of our models and sup-

ports effective learning. We train our classifier with

the following settings and hyper-parameters: the

epoch is 250, the batch size is 4, the learning rate is

1e-3, and the AdamW optimizer has a linear sched-

uler. We conduct all the experiments using Pytorch

(Paszke et al., 2019) and HuggingFace library(Wolf

et al., 2020) on 4 NVIDIA A100-SXM4-80GB

Data Preparation. To establish a comprehensive

retrieval system, we use datasets representing both

infringement and non-infringement cases. Each

dataset consists of input-reference text pairs (x, t),

where the input text x acts as a query, and the refer-

ence text t provides contextual information, mean-

ing the surrounding content in a specific context, such as the following text in a classic work. The

entire dataset is stored as a structured collection:

 $\mathcal{D} = \{(x_i, t_i)\}_{i=1}^N$, where N is the total number of

pairs in the dataset. By merging multiple datasets

into a unified pool, we ensure broad coverage of

potential scenarios, forming a strong foundation

semantic relationships between input and refer-

ence texts, we encode each text into a dense vector

representation using a pre-trained Sentence Transformer \mathcal{E} (all-roberta-large-v1) (Liu et al., 2019):

 $v_x = \mathcal{E}(x), \quad v_t = \mathcal{E}(t), \text{ where } v_x, v_t \in \mathbb{R}^d \text{ are }$

the dense embeddings of the input query and the

reference text, respectively, and d is the embedding

dimension. To enhance efficiency, we implement

batch encoding with GPU acceleration, ensuring

scalable processing of large datasets while main-

Indexing with FAISS & Document Storage in

SQLite. For efficient nearest-neighbor retrieval,

we use FAISS (Douze et al., 2024) with the Index-IVFFlat method, which clusters the vector space to accelerate query execution. Given a set of indexed

reference embeddings $\{v_{t_i}\}_{i=1}^N$, FAISS partitions

them into K clusters, with each vector assigned to

 $C = \{\mu_k\}_{k=1}^K, \quad \mu_k = \frac{1}{|C_k|} \sum_{v \in C_*} v,$

To capture the

for benchmarking and future improvements.

Dense Representation Encoding.

taining retrieval accuracy.

its nearest cluster center:

RAG System Construction

embedding v_x is assigned to the closest centroid μ_k , and the nearest neighbors are searched within that cluster: $\hat{t} = \operatorname{argmin}_{t_i \in C_k} ||v_x - v_{t_i}||_2$. This reduces search complexity from O(N) to O(N/K), ensuring fast retrieval even for large datasets.

Additionally, we use SQLite for structured text storage, where each document entry (including original input and reference texts) is indexed with its corresponding embedding. This allows efficient retrieval of both vector embeddings and textual data based on semantic similarity and exact text matches: $\mathcal{T} = \{(x_i, t_i, v_{t_i})\}_{i=1}^N$.

Retrieval Accuracy Since our input and reference pairs are stored in the external knowledge base as structured pairs, our retrieval method achieves a 100% accuracy rate in search matching within the current dataset:

$$\underset{t_i}{\operatorname{argmax}}\operatorname{Sim}(v_x, v_{t_i}) = t_j, \text{ where } (x, t_j) \in \mathcal{D}.$$

Here, $Sim(\cdot, \cdot)$ denotes the similarity function (e.g., cosine similarity), ensuring that the retrieved reference always corresponds to the correct pair in our dataset. By integrating dense vector retrieval with structured text storage, ISAC provides efficient and accurate reference retrieval, forming a crucial component of our infringement detection system.

C Metric Details

ACC & F1. For the classification task where the predictions are discrete, we use F1 score and Accuracy as the metrics to assess the performance of the predicted categories.

In classification tasks, accuracy and F1 score are two important metrics used to evaluate the performance of a model. Accuracy represents the proportion of correctly classified instances among the total number of instances, providing a general measure of how often the model makes the right prediction. It is calculated as:

$$\mathcal{A} = \frac{\mathcal{T}_p + \mathcal{T}_n}{\mathcal{N}_{\text{total}}} \tag{7}$$

where \mathcal{T}_p and \mathcal{T}_n represent true positives and true negatives, respectively, and \mathcal{N}_{total} is the total number of samples. Accuracy is simple and intuitive but may be unreliable with imbalanced datasets, where one class dominates the others. A model predicting only the majority class can achieve high accuracy but fail to detect minority instances.

The F1 score provides a more balanced evaluation by considering both precision and recall. Precision (\mathcal{P}) is the fraction of correctly predicted positive observations out of all positive predictions, while recall (\mathcal{R}) is the fraction of true positives out of all actual positive samples. The F1 score is defined as:

869

870

871

874

875

876

878

900

901

902

904

906

$$\mathcal{F}_1 = 2 \times \frac{\mathcal{P} \times \mathcal{R}}{\mathcal{P} + \mathcal{R}} \tag{8}$$

The F1 score is particularly useful in imbalanced datasets, balancing false positives and false negatives to provide a comprehensive view of performance. While accuracy works well for balanced data, the F1 score is more informative for assessing real-world classification problems.

ROUGE. ROUGE (Recall-Oriented Understudy for Gisting Evaluation) is a set of metrics commonly used to evaluate the quality of automatic 885 text summarization and natural language generation systems by comparing the overlap between generated text and reference text. ROUGE includes several variations: Rouge-N evaluates the overlap of N-grams, Rouge-L focuses on the longest common subsequence (LCS), and Rouge-S uses skipbigram matching. Rouge-L specifically measures the sequence similarity between generated text and 893 reference text by identifying the longest common subsequence. It captures both content and sequential structure. The Rouge-L score comprises Precision, Recall, and F-score, representing different perspectives of text similarity, where Recall emphasizes content coverage and Precision reflects exact matching accuracy. In our experiments, we calculate Rouge-1 and Rouge-L scores using the rouge_score library, and we utilize the Rouge-L score as a key metric for classifying and evaluat-903 ing the quality of datasets based on the sequential similarity of text pairs. 905

D Dataset

We provide the data source of copyrighted ma-907 terial in Table 6. For the literal copying task, 908 which assesses copyright risks in text continua-909 tions, the dataset includes excerpts from 16 fiction 910 911 titles in BookMIA (Shi et al., 2023), likely part of ChatGPT's training data (Chang et al., 2023). 912 To increase diversity, we supplemented these with 913 works by J.K. Rowling. For the non-literal copying 914 task, focused on event and character copying, we 915

used CliffsNotes study guides paired with humanwritten summaries. To ensure all texts remain under copyright, we excluded non-fiction and pre-1923 books.

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

Ε **Prompt Design**

In designing the baseline for our experiment on detecting text infringement risks through internal states, we adopted the "LLM as Judge" approach. This method leverages LLMs to evaluate potential infringement risks in text generation tasks. To ensure robust and accurate assessment, we carefully crafted evaluation prompts tailored to capture nuanced scenarios of potential infringement, as shown in Table 9. This design allows for a systematic comparison between traditional heuristic-based methods and our proposed internal state detection framework.

F **Ablation Studies**

Effect of Internal States Layers F.1

Unlike previous studies emphasizing the importance of later layers in LLMs for tasks like hallucination detection (Ji et al., 2024), our experiments on copyright detection show a different trend based on model size. For smaller models like Llama-3.1-8B, layer selection doesn't significantly affect the prediction of potential copyright infringement. However, for larger models such as Llama-3.1-70B, deeper layers significantly improve performance, especially in accuracy and F1 score.



Figure 5: Impact of layer selection on copyright infringement risk prediction: A comparative analysis across different layers in Llama models with 8B and 70B parameters. For smaller models (Llama-3.1-8B), the prediction performance is relatively consistent across layers, with minimal variation in accuracy and F1 score. For larger models (Llama-3.1-70B), deeper layers significantly enhance performance, capturing more nuanced semantic features and improving the prediction of potential copyright infringement in text continuation tasks.

Previous research (Azaria and Mitchell, 2023) emphasized the effectiveness of the final layer for 946 hallucination detection, but our analysis indicates that for copyright risk prediction, deeper layers are more essential in larger models. As shown in Figure 5, deeper layers in larger models are better 950 at capturing textual similarities to existing liter-951 ary works, which is crucial for identifying potential infringement. In contrast, for smaller models, early and intermediate layers perform similarly to the final layer, suggesting that while semantic and 955 contextual information is spread across all layers, deeper layers in larger models are more effective 957 in detecting the finer details needed for accurate predictions.

> One possible explanation for this is that copyright detection requires identifying both local and global semantic patterns, which are essential for spotting similarities and potential plagiarism. In smaller models, these patterns are well-represented across various layers, whereas larger models excel in capturing the more subtle textual similarities through their deeper layers. Unlike hallucination detection, which focuses on long-range dependencies and uncertainty captured in later layers, copyright detection benefits from the ability of larger models to focus on detailed patterns across deeper layers.

F.2 Effect of Model Size

961

962

963

964

965

968

969

971

972

973

974

975

976

978

979

982

983

985

987

991

995

This section investigates how model size influences the efficacy of LLM's internal states in classifier training, comparing Llama models with 1B, 3B, 8B, 13B, and 70B parameters. Experimental results demonstrate that smaller Llama models generate internal states that yield lower F1 scores and accuracy in classification tasks compared to larger models, regardless of whether the input data is presented in isolation or supplemented with reference information provided by RAG system. As shown in Figure 6, the performance of the models improves significantly with increasing size, highlighting the importance of model scale in enhancing classification accuracy and F1 scores.

As shown in Figure 7, Larger models not only outperform smaller ones in producing higherquality internal states for classification but also excel in text generation tasks. They exhibit a stronger ability to comprehend context, maintain coherence, and produce semantically rich text. These capabilities lead to more accurate continuations that closely align with the input text, facilitating the generation



Figure 6: Impact of model size on behavior prediction performance: a comparative analysis of classification accuracy and F1 scores across Llama models with 1B to 70B parameters

of datasets that better represent the original data. Consequently, this improves the precision of subsequent dataset categorization processes. 996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1021

1023

1024

1025

1026

1027

1028

1029

1031

To address the behavioral variations arising from differences in internal state quality and data generation strategies across models of varying sizes, it is essential to design separate, model-specific databases. These databases should capture the unique characteristics of the internal states and outputs generated by each model size. For smaller models, stricter control over Rouge-based segmentation thresholds may be necessary to achieve clearer distinctions between potentially infringing and non-infringing data. Such measures are particularly important because smaller models tend to produce less semantically rich internal states, potentially diminishing classification accuracy.

By refining the dataset segmentation strategyparticularly for smaller models the accuracy of infringement risk predictions can be significantly improved. This ensures that even resourceconstrained models are well-prepared for robust downstream classification tasks, enabling reliable performance across diverse use cases.

F.3 Effect of Generation Prompts

In this section, we discuss the impact of varying prompt design strategies used as input to the LLM on the prediction accuracy of the trained model during the dataset construction process. Building on the prompt configurations from prior work (Chen et al., 2024), we modify them as the sole variable in our experiments. Table 3 presents the results of these experiments, highlighting how different prompt formulations influence the overall performance. The prompt design is presented in Table 7 for clarity and reference.



Figure 7: Distribution of upper and lower 30% Rouge-L scores for LLMs of different sizes based on continuation outputs. Larger models tend to generate continuation outputs with a higher risk of copyright infringement, as they are more likely to produce content with a high similarity to reference texts.

Table 3: The table illustrates how prompt selection affects text generation by comparing F1 scores and accuracy across different prompts used in preparing the training dataset for the Llama-3.1-70B model. It evaluates two methods: IS-w/oRAG (Internal States Judge without the RAG system) and IS-w/RAG (Internal States Judge with the RAG system).

| | | Division (10%) | | Division (20%) | | Division (30%) | |
|---------|-----------|----------------|--------|----------------|--------|----------------|--------|
| Prompt | Method | ACC (%) | F1 (%) | ACC (%) | F1 (%) | ACC (%) | F1 (%) |
| Prompt1 | IS-w/oRAG | 97.01 | 96.00 | 88.79 | 87.43 | 85.24 | 88.07 |
| | IS-w/RAG | 97.34 | 95.13 | 90.57 | 89.34 | 87.29 | 89.94 |
| Prompt2 | IS-w/oRAG | 85.71 | 89.50 | 75.12 | 79.52 | 67.55 | 75.25 |
| | IS-w/RAG | 91.73 | 93.17 | 89.27 | 89.42 | 73.84 | 75.06 |
| Prompt3 | IS-w/oRAG | 91.41 | 93.33 | 74.51 | 79.22 | 62.54 | 71.29 |
| | IS-w/RAG | 98.44 | 98.73 | 87.75 | 88.29 | 70.03 | 75.53 |

As shown in this table, the design corresponding to Prompt 2 exhibits relatively lower performance compared to the designs associated with Prompt 1 and Prompt 3. Both the IS-w/oRAG and IS-w/RAG methods yield weaker results under this configuration, with ACC and F1 scores declining as the dataset division percentage increases. In conclusion, variations in each prompt used for data generation have a noticeable impact on the prediction accuracy of models trained with the resulting datasets. Therefore, when predicting infringement risks, multiple models utilizing datasets generated with different prompt designs can be employed. By applying this approach, it becomes possible to identify and prioritize data associated with higher infringement risk, enhancing the effectiveness of the risk detection process.

1032

1033

1034

1035

1036

1038

1039

1040

1041

1044

1045

1046

1047

1048

1050

F.4 Effect of Internal States Extraction Methods

1051In our experiments, we examined the impact of dif-1052ferent internal state extraction methods at a given1053layer for copyright detection, specifically compar-1054ing the effectiveness of using the average internal1055state across all tokens versus extracting only the1056internal state of the last token. Our results indicate1057that, for a fixed layer, computing the mean internal

state across all tokens provides significantly higher prediction accuracy than relying solely on the internal state of the last token, as shown in Table 4. 1058

1059

1061

1063

1064

1065

1066

1067

1069

1070

1072

1074

1075

1077

1078

1080

1081

1082

1084

When taking the average internal state, the representation is aggregated across all token embeddings within the selected layer. This method ensures that the extracted feature captures a comprehensive understanding of the entire sequence, incorporating both local token-level details and global contextual relationships. As a result, this approach is particularly effective for copyright detection, where recognizing semantic and structural similarities across a text is crucial.

Conversely, extracting the last token's internal state from the same layer restricts the representation to a single token position, potentially losing valuable contextual information present in the earlier tokens. While this method is commonly used in classification tasks, our analysis shows that, in copyright risk prediction, it leads to a weaker overall representation, as the key signals indicating similarity to existing works may be distributed throughout the sequence rather than concentrated in the final token.

These findings highlight that, even when working with the same layer, the choice of how internal states are extracted plays a crucial role in model

Table 4: This table explores the effectiveness of different internal state extraction methods under the Llama-3.1-70B model. The results show that, at a fixed layer, averaging the internal states across all tokens significantly outperforms using only the last token's internal state, as the averaging method better captures contextual information, making it more suitable for copyright detection.

| | Division (10%) | | Division (20%) | | Division (30%) | |
|-------------------|----------------|--------|----------------|--------|----------------|--------|
| Methods | ACC (%) | F1 (%) | ACC (%) | F1 (%) | ACC (%) | F1 (%) |
| Last Token-w/oRAG | 68.57 | 75.56 | 66.83 | 74.33 | 62.99 | 72.46 |
| Last Layer-w/oRAG | 100.00 | 100.00 | 94.55 | 94.63 | 93.18 | 93.62 |
| Last Token-w/RAG | 88.57 | 89.09 | 88.61 | 88.78 | 83.77 | 85.47 |
| Last Layer-w/RAG | 100.00 | 100.00 | 95.05 | 94.68 | 94.48 | 94.64 |

performance. Averaging across all tokens allows for a more robust and contextually rich representation, making it a preferable choice for copyright infringement detection. Future studies could further explore whether weighting token contributions or applying attention-based pooling strategies can further refine the effectiveness of internal state-based detection methods.

F.5 Non-literal Copying Detection

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1099

1100

1101

1102

1103

1104

1105

1106 1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

In this section, we examine infringement detection for non-literal paraphrasing (Chen et al., 2024). We measure the overlap between generated and reference texts at the character and event levels to assess potential infringement. This approach is similar to the literal copying task, but in the non-literal case, the continuation is based on paraphrasing instead of direct copying. As shown in Table 5, we evaluate prediction accuracy across three prompt types, detailed in Table 8.

Despite the smaller dataset, the results show that detecting infringement in paraphrased texts is more challenging for large language models than in literal data. This leads to lower prediction accuracy in non-literal paraphrasing, as paraphrased texts are harder to compare directly with the reference text due to structural, vocabulary, and expression differences. This complexity reduces the model's ability to generalize, resulting in lower classification performance. Even with additional reference information by using RAG system, the model struggles to capture the intricate features required for accurate prediction.

| | Prompt 1 | | Prompt 2 | | Prompt 3 | |
|------------|----------|--------|----------|--------|----------|--------|
| Method | ACC (%) | F1 (%) | ACC (%) | F1 (%) | ACC (%) | F1 (%) |
| IS-w/oRAG | 53.33 | 57.89 | 46.67 | 54.72 | 51.11 | 62.30 |
| IS-w/RAG-C | 63.33 | 70.27 | 56.67 | 41.67 | 56.67 | 31.58 |
| IS-w/RAG-E | 55.56 | 65.60 | 52.22 | 58.93 | 55.56 | 64.29 |

Table 5: The experiment utilizes non-literal data for both training and testing, with the training set divided based on the upper and lower 30% of Rouge scores. "C" denotes character-related copying while "E" represents event-related copying. Additionally, test results are extracted from the internal states of Llama-3.1-70B.

Table 6: List of Book Titles and Authors for literal task

| Title | Author | |
|--|--------------------|--|
| 1984 | George Orwell | |
| A Game of Thrones | George R.R. Martin | |
| Casino Royale | Ian Fleming | |
| Dune | Frank Herbert | |
| Fahrenheit 451 | Ray Bradbury | |
| Fifty Shades of Grey | E.L. James | |
| Five on a Treasure Island | Enid Blyton | |
| Harry Potter and the Sorcerer's Stone | J.K. Rowling | |
| Hitchhiker's Guide to the Galaxy | Douglas Adams | |
| Lord of the Flies | William Golding | |
| The Da Vinci Code | Dan Brown | |
| The Hunger Games | Suzanne Collins | |
| The Silmarillion | J.R.R. Tolkien | |
| Their Eyes Were Watching God | Zora Neale Hurston | |
| Things Fall Apart | Chinua Achebe | |
| To Kill a Mockingbird | Harper Lee | |
| Harry Potter and the Philosopher's Stone | J.K. Rowling | |
| Harry Potter and the Chamber of Secrets | J.K. Rowling | |
| Harry Potter and the Prisoner of Azkaban | J.K. Rowling | |
| The Hobbit | J.R.R. Tolkien | |

Table 7: Three prompt templates for generating passage completion to evaluate literal copying. (Chen et al., 2024)

Prompt 2

You will be given several passage examples from famous literary works. Each → example includes a prefix and its completion. After these examples, you → will be provided with the final prefix, which you need to complete in → the style of the famous work. Your completion should be at least 100 → words long and formatted as follows: "Completion: {your completion to → the prefix}" Generate only the completion for the last given prefix → without copying the examples, the given prefixes, or adding extra → commentary. Beginning: {demonstration-input}

```
Completion: {demonstration input}
Beginning: {input}
Completion:
```

Prompt 3

Table 8: Three prompt templates for generating creative story to evaluate non-literal copying. (Chen et al., 2024)

Prompt 1

Prompt 2



Prefix: {input}
Completion:

Prompt 3

You will be provided with a series of creative writing examples. Each example \rightarrow includes a story beginning and its completion. After these \rightarrow demonstrations, you will be given the last beginning, which you are \rightarrow asked to complete. You should compose an original story of about 500 \rightarrow words that is distinct from any existing copyrighted novels. An original \rightarrow story should have unique characters, settings, plots, and dialogues \rightarrow that are not taken from existing works. Avoid using specific phrases, \rightarrow sequences of events, or character traits that directly mirror those \rightarrow found in copyrighted material. Your output should be in the format: " \rightarrow Completion: your completion to the beginning". Generate only the \rightarrow completion to the last given beginning without copying the examples, the \rightarrow given beginnings, or adding extra commentary.

```
Given Prefix: {demonstration-input}
Provided Completion: {demonstration-output}
Given Prefix: {input}
Completion:
```

Table 9: Baseline Evaluation Prompts for LLM-as-Judge Framework

Version A: Input-Only Assessment

```
{
     "role": "user"
     "content": f"""
          Please only output 'label: 0' or 'label: 1', nothing else. Do not
               \hookrightarrow provide intermediate steps or explanations; simply output the
               \hookrightarrow generated label.
          You are a language model that determines whether the next generated
               \hookrightarrow sentence based on the given input text might cause potential
               ↔ copyright infringement.
          If you believe the next sentence could significantly resemble or repeat
          → existing literary works, indicating possible infringement,
→ return "label: 0" (possible infringement).
If you believe the next sentence will not show high similarity to any

→ literary work, return "label: 1" (no infringement).

          Here is the input to evaluate:
Input text: '{input_text}'
          Based on the input text provided, determine if generating the next
               \hookrightarrow sentence might result in infringement.
}
```

Version B: Reference-Aware Assessment

```
{
     "role": "user"
     "content": f"""
         Please only output 'label: 0' or 'label: 1', nothing else. Do not
             \hookrightarrow provide intermediate steps or explanations; simply output the
              \hookrightarrow generated label.
         You are a language model that evaluates whether the output text
              \hookrightarrow potentially involves infringement based on the relationship
              \hookrightarrow between the reference and output sequences.
         If you believe the next sentence could significantly resemble or repeat \hookrightarrow existing literary works, indicating possible infringement,
              ↔ return "label: 0" (possible infringement).
         If you believe the next sentence will not show high similarity to any
              \hookrightarrow literary work, return "label: 1" (no infringement).
         Here are the inputs to evaluate:
         Input text: '{input_text}'
         Reference: '{reference}' (This is a segment from the original source
             \hookrightarrow text corresponding to the input text.)
         Based on the comparison between the reference and output above,
             \hookrightarrow determine if generating the next sentence might result in

→ infringement.

}
```