# *Is your Reviewer 2 Human?* Tracing a Novel Problem of LLM Authorship in Academic Peer Reviews

Anonymous ACL submission

#### Abstract

The integrity of the peer review process is crucial for maintaining scientific rigor and trust in academic publishing. This process relies on (human) domain experts critically evaluating the merit of the submitted manuscripts. However, the peer review system faces growing strain from increased submissions and limited reviewer availability, prompting lazy reviewing practices in which reviewers use large language models (LLMs) to generate reviews, raising concerns about the quality, reliability, 011 and accountability of those evaluations. Pre-012 vious work has focused on estimating the pro-014 portion of AI-generated peer reviews or developing AI-generated text detectors. However, existing detectors are not resistant to adversarial attacks and often require domain or modelspecific retraining. To address these challenges, 019 we propose a framework for peer review watermarking. Our method includes a Query-Aware Response Generation module that selectively embeds subtle yet detectable signals while preserving scientific terminology, based on the user's submission of a research paper, along with a watermarking detection mechanism that enables editors to reliably verify the authenticity of reviews. Extensive experiments on ICLR and NeurIPS data demonstrate that our method outperforms various AI text detectors under adversarial attacks. We hope that this work will facilitate the further development of watermarking and responsible use of LLM systems. We make our code and dataset public<sup>1</sup>.

# 1 Introduction

034

The emergence of frontier large language models (LLMs), such as Claude, Gemini, GPT-4 (Achiam et al., 2023), LLaMa, etc. has revolutionized natural language generation. The sophisticated human-like fluency and coherence exhibited by texts pro-

<sup>1</sup>https://anonymous.4open.science/r/ PeerWatermarking-51DD/ duced by these models present considerable challenges in discerning whether such content is humangenerated or machine-generated, even for domain experts (Shahid et al., 2022). Peer review remains a foundational practice in academia, serving as a critical mechanism through which expert scrutiny ensures the integrity and credibility of scholarly outputs prior to publication (Alberts et al., 2008). Nevertheless, the escalating volume of manuscript submissions (Bornmann and Mutz, 2015; McCook, 2006) has increasingly burdened the peer review system, amplifying concerns regarding the system's sustainability and efficacy (Arns, 2014). 041

042

043

044

045

047

049

052

053

054

057

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

Scientific peer review fundamentally depends on expert reviewers to provide insightful, critical, and constructive evaluations of submitted manuscripts or proposals (Shah, 2022). According to the Association for Computational Linguistics (ACL)  $policy^2$ , Artificial Intelligence (AI) tools may assist with paraphrasing and proofreading tasks, particularly benefiting non-native English speakers; however, reviewers are required to independently generate substantive review content. Recent research (Liang et al., 2024) examining peer reviews from AIrelated conferences identified that approximately 6.5% to 16.9% of review text may have been substantially modified using LLMs. The study highlights a significant increase in LLM usage, particularly ChatGPT, immediately preceding review deadlines, with higher reliance detected among reviewers not engaging in author rebuttals at prominent venues such as ICLR and NeurIPS. Additionally, increased ChatGPT usage was associated with diminished self-reported reviewer confidence. A relevant research (Ye et al., 2024) demonstrated that manipulating just 5% of reviews could disrupt rankings sufficiently to displace approximately 12% of papers from the top 30%. Further, this

<sup>&</sup>lt;sup>2</sup>https://2023.aclweb.org/blog/review-acl23/#faq-can-iuse-ai-writing-assistants-to-write-my-review

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

157

158

159

160

161

162

163

164

165

166

167

169

170

171

172

173

study revealed intrinsic limitations of LLM-based 079 reviews, including potential biases such as favor-080 ing incomplete submissions over fully developed 081 manuscripts and preferentially rating submissions by prominent authors in single-blind review scenarios. Moreover, authors can intentionally embed covert content within manuscripts to deliberately manipulate LLM-generated reviews, leading to artificially inflated assessments misaligned with human evaluations. These findings collectively suggest that current LLMs are insufficiently robust for deployment as primary reviewers due to inherent vulnerabilities and susceptibility to manipulation. Consequently, rigorous safeguards and enhanced evaluation frameworks must be implemented to ensure review fairness and accuracy before broader adoption of LLM-based review processes.

In this paper, we propose a novel framework for watermarking LLM-generated peer reviews. Our approach consists of several key components. First, we introduce a Query-Aware Response Generation module, which selectively applies watermarking when user uploads a research paper and there is a risk of peer review misuse. Then, our Watermark Injection Mechanism embeds subtle yet detectable signals in peer reviews while preserving scientific terminology. Finally, we implement Watermark Detection, which allows editors and conference chairs to verify the authenticity of peer reviews. Across ICLR and NeurIPS data, our watermarking framework achieves significantly higher detection accuracy than existing methods, maintaining performance even in adversarial settings. Our work aims to safeguard the peer review process against misuse of generative language models, thereby reinforcing ethical norms in scholarly communication and contributing to a trustworthy research ecosystem.

098

100

101

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

121

122

123

124 125

126

127

129

Our contributions are summarized as follows:-

- We introduce the novel task of watermarking AI-generated peer reviews to ensure authenticity and traceability.
- We propose a new lightweight framework that (i) employs a gating mechanism to watermark only potentially unsafe peer-review generation requests and (ii) introduces a simple yet effective watermarking strategy that markedly improves the detection of machine-generated peer reviews.
  - Our watermarking technique outperforms existing AI-based text detectors, even under adversarial conditions.

# 2 Related Work

Zero-shot text detection identifies AI-generated text without requiring training on specific data (Mitchell et al., 2023). Solaiman et al. (2019) detect AI-generated text by measuring its average log probability under the generative model. Detect-GPT (Mitchell et al., 2023) leverages the tendency of AI-generated text to reside in negative curvature regions of the model's log probability function for detection. Fast-DetectGPT (Bao et al., 2023a) enhances efficiency by applying conditional probability curvature instead of raw probability. Guo et al. (2023) developed the OpenAI text classifier by training it on a large dataset comprising millions of texts. However, heavy dependence on training data makes many of these models susceptible to adversarial attacks (Wolff, 2020).

Watermarking AI-generated text, introduced by Wiggers (2022), embeds an imperceptible pattern to verify authorship, similar to encryption. Watermarks can be embedded without requiring modifications to the underlying language model, allowing standard models to generate watermarked text (Kirchenbauer et al., 2023). Rather than focusing on individual detection, Liang et al. (2024) proposed a method that estimates the proportion of AI-generated text within a large corpus using maximum likelihood estimation of probability distributions.

As far as we know, this is the first work to address AI-generated peer review detection through watermarking. Unlike existing AI text detectors, which are vulnerable to adversarial attacks, our approach embeds traceable markers directly into generated content, improving the detection of AIgenerated reviews. Additionally, current AI text detection models require task-specific training for each conference and dataset, making large-scale deployment challenging. In contrast, our watermarking method provides a scalable solution that eliminates the need for continuous retraining, enhancing both reliability and adaptability across diverse academic settings.

# 3 Methodology

Figure 1 illustrates the framework, which con-<br/>sists of two key components: (a) Query-Aware174Response Generation, where a user uploads or sub-<br/>mits a research paper along with a query related to<br/>it, which is classified by the Query Type Identifier.175If identified as an unsafe query (indicating poten-178



Figure 1: Overview of the Proposed Watermarking Framework. (a) Watermark Generation (b) Watermark Detection; Here red indicates red token, green indicates green token, blue indicates blue token

tial peer review misuse), it is processed through the LLM generator with watermarking. Otherwise, queries proceed through normal Generation. (b) Watermark Detection, where an editor or chair submits a research paper and its corresponding review for verification.

# 3.1 Document Type Identifier

180

183

186

190

192

193

197

198

199

206

210

To determine whether a document is a research paper, we implemented a simple rule-based approach. Documents with fewer than 500 words were filtered out as unlikely to be research papers. We detected key section headers such as Abstract, Introduction, Methodology, Results, Conclusion, and References using regular expressions and header position analysis. A document was classified as a research paper if at least three core sections were present, along with a reference section and in-text citations (e.g., (Author, Year) or [1]). This method provided a lightweight and efficient first-pass classification before passing the query to the Query Type Identifier, which then determines whether the query is safe or unsafe.

# 3.2 Query Type Identifier

The Query Type Identifier determines whether a query is classified as Safe (S) or Unsafe (UN). A query is considered Unsafe (UN) if it requests a peer review in a way that allows the reviewer to directly submit the generated content as a peer review. Any query that does not fall into the Unsafe category is classified as Safe (S), including those that seek explanations, summaries, or clarifications related to the paper's content. We employ a few-<br/>shot prompting approach to classify user queries.211We discuss this in detail in Section F.213

214

215

216

217

218

219

220

221

222

224

225

226

# 3.3 Watermark Injection

#### Algorithm 1 Watermark Injection

- **Require:** Vocabulary V, Paper P , Watermark Strength  $\delta$ , Green List Fraction  $\gamma$
- 1: Compute seed  $S_T$  using paper information and secret keys  $K_{secret}$  and p.
- 2: Generate green list G and red list R from V based on  $\gamma$
- 3: Extract blue list B as domain-specific terms from P
- 4: for each generation step  $t = 1, \ldots, T$  do
- 5: Obtain logits  $l_w^{(t)}$  from LLM

6: Adjust logits: 
$$l_w^{(t)} = l_w^{(t)} + \delta \cdot \mathbb{K}[w \in G \cup B]$$

$$p_{w}^{(t)} = \frac{e^{l_{w}^{(t)}}}{\sum\limits_{w' \in V} e^{l_{w'}^{(t)}}}$$

8: Sample next token  $w^{(t)}$  from adjusted distribution 9: end for

In this section, we introduce our watermarking injection technique, which ensures that LLMgenerated text is subtly embedded with verifiable signals without significantly altering fluency or coherence. In this work, we utilize the soft watermarking technique that introduces probabilistic biases in token selection to subtly mark text generated by large language models (LLMs). Our approach is inspired by prior watermarking method (Kirchenbauer et al., 2023). However, this approach face challenges in peer review text generation due to the inherent trade-off between watermark strength and text quality. The random selection of green and red tokens introduces high variability, weakening the watermark signal. Our novel approach incorporates 'blue tokens'—key technical terms from the research paper itself—enhancing watermark robustness by grounding signals in semantically meaningful content. Algorithm 1 outlines the complete watermark injection process.

# 3.3.1 Paper Seed Generation

227

228

237

239

241

242

243

244

245

246

247

257

258

261

262

263

270

272

Our seed generation mechanism ensures unique, deterministic, and secure watermarking by leveraging context-aware encoding, cryptographic hashing, and a secret key. The input text T, which can be a paper title, abstract, or any small portion of text from any section of the paper. To strengthen security, a secret key  $K_{\text{secret}}$  is concatenated with T, ensuring that different users generate distinct seeds:

$$I = T \| K_{\text{secret}} \tag{1}$$

 $K_{\text{secret}}$  is a fixed, confidential key (text) known only to trusted parties (e.g., conference chairs or editors). It is not generated at runtime. Instead, it is concatenated with the paper text to produce a deterministic seed for watermarking. This ensures that watermarking is reproducible and verifiable only by those with access to the key. The encoded representation is then hashed using SHA-256 for collision resistance:

$$H(I) = SHA-256(E(I)) \tag{2}$$

Here, E(I) denotes the full encoding of the input text I, computed by applying the character-level encoding function f(c, n) to each character c in I. Each character is mapped to its alphabetical index and applying a shift cipher based on the text length:

$$f(c,n) = ((\operatorname{ord}(c) - \operatorname{ord}('A') + n) \mod 26) + 1$$
 (3)

where c is the character and n is the total number of characters in T. This ensures that encoding remains text-dependent, enhancing uniqueness. Finally, the hash is mapped to a bounded numeric space using modular reduction:

$$S_T = H(I) \mod pk \tag{4}$$

Even if an attacker identifies the specific paper text used for watermarking, they would still require two secret keys,  $K_{\text{secret}}$  and pk, to decode the green list and verify the watermark. These keys ensure that only authorized individuals, such as the editor or program chair, can perform detection. To maintain the integrity and security of the watermarking system,  $K_{\text{secret}}$  and pk keys must be kept strictly confidential and accessible only to authorized personnel.

273

274

275

276

278

279

281

282

283

284

285

286

287

289

290

291

292

293

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

#### 3.3.2 Green-Red Token Partitioning

Given a vocabulary set V, we define a subset of tokens, G, termed as the "green list" which are favored during text generation. The remaining tokens form the "red list" R. Instead of a probability-based split, we use a deterministic *random permutation* seeded by the paper seed generator. A fraction  $\gamma$  of tokens is selected as green, ensuring consistency across runs. Let |V| denote the vocabulary size :-

 $|G| = \gamma |V|, \quad |R| = (1 - \gamma)|V|.$  (5)

#### 3.3.3 Blue Token Selection

We define a subset of tokens, denoted as the *blue list B*, which consists of important technical terms extracted from a given research paper. These blue tokens represent domain-specific terminology that is crucial for maintaining the technical accuracy and coherence of the generated text.

Unlike the green list G, which is deterministically selected based on a fixed fraction  $\gamma$  of the vocabulary, the blue list is explicitly derived from the research content, ensuring a stronger alignment with the subject matter of the paper. To construct B, we utilize a language model (LLM) to extract key technical terms from the paper by prompting it to identify domain-relevant terminology. We found that average number of extracted terms per paper is approximately 43.83. We discuss this in detail in Appendix C.

# 3.3.4 Logit Adjustment Mechanism

During inference, given a token sequence  $w_0, w_1, \ldots, w_{t-1}$ , the language model produces a logit vector  $l_w^{(t)}$  representing the probability distribution over V. We modify these logits using a biasing function:

$$l_w^{(t)} = l_w^{(t)} + \delta \cdot \mathscr{W}[w \in G \cup B], \tag{6}$$

where  $\delta$  is a tunable parameter controlling the watermarking strength, and  $\mathbb{M}[w \in G \cup B]$  is an indicator function returning 1 if w is in the green list G or blue list B, and 0 otherwise. The resulting logits are then passed through the softmax function to obtain the final token probabilities:

$$p_w^{(t)} = \frac{e^{l_w^{(t)}}}{\sum\limits_{w' \in V} e^{l_{w'}^{(t)}}}.$$
(7)

This ensures that tokens in G and B are more likely to be sampled while discouraging tokens from R, thereby reinforcing both the structured watermarking and the preservation of domain-specific terminology.

#### 3.4 Watermark Detection

320

321

322

323

324

325

327

329

331

333

334 335

336

337

339

Given a research paper, the proposed algorithm generates a deterministic seed to ensure consistency between encoding and detection. Since each research paper is unique, the generated seed remains identical to that used during watermark insertion. As a result, the same random token list (formerly the green list) is reconstructed. Similarly, the blue token list, consisting of technical terms extracted from the paper, is also reproduced, as these terms remain unchanged. Consequently, the marked token set, i.e., the union of random and blue tokens, remains identical, enabling accurate watermark detection. We discuss the algorithm in detail in Algorithm 2.

#### Algorithm 2 Watermark Detection

Require: Peer Review Text R, Paper Tokens P

- **Ensure:** Marked Token Fraction  $f_m$ , Z-Score z
- 1: Tokenize the peer review R using tokenizer  $\mathcal{T}$
- 2: Generate a deterministic seed  $S_T$  from the paper using the seed generator
- 3: Partition vocabulary V into random tokens G and red tokens  $R_{\rm red}$  using  $S_T$
- 4: Extract blue tokens from the paper:  $O = P \cap R_{red}$
- 5: Compute marked tokens:  $M = G \cup O$
- 6: Extract bigrams B from R and initialize marked token hit count  $M_c = 0$
- 7: for each bigram  $(x, y) \in B$  do
- 8: Increment  $M_c$  if  $y \in M$
- 9: end for
- 10: Compute marked token fraction:

$$f_m = \frac{M_c}{|B|}$$

11: Compute expected marked token fraction:

$$\mathbb{E}[f_m] = \frac{|M|}{|V|}$$

12: Compute z-score:

$$z = \frac{M_c - |B|\mathbb{E}[f_m]}{\sqrt{|B|\mathbb{E}[f_m](1 - \mathbb{E}[f_m])}}$$

13: return  $f_m, z$ 

We extract bigrams (k=2) as part of our detec-

tion pipeline, building on the k-gram watermarking framework introduced by (Kirchenbauer et al., 2023) and extended theoretically by (Zhao et al., 2023). These works justify using k-grams, where k can be tuned based on task requirements. Bigrams offer a practical balance between local contextual awareness and statistical reliability. Specifically, they reduce token-level noise (e.g., repetition) and improve robustness to paraphrasing.

We computed z score which is a standard test statistic (Zhao et al., 2023; Kirchenbauer et al., 2023), used to distinguish between watermarked and non-watermarked text, with theoretical guarantees on false positive and false negative rates (see (Zhao et al., 2023), Theorems 3.3–3.5). Under this model, the expected count of green tokens is  $\gamma T$ , with variance  $T \cdot \gamma (1 - \gamma)$ , yielding the normalized score:

$$z = \frac{|G| - \gamma T}{\sqrt{T \cdot \gamma (1 - \gamma)}} \tag{8}$$

Here y indicates whether the review is watermarked.

#### 4 **Experiments**

#### 4.1 Implementation Details

We used 1,090 papers from ICLR and NeurIPS (year: 2022) for our experiments from (Kumar et al., 2024). The average number of reviews per paper is 3.88, and the average token length of the reviews is 566.42. For generation, we used the Llama-3.1-8B-Instruct<sup>3</sup> in our experiments. We discuss the implementation details in Appendix A.

#### 4.2 Main Result

Mod	lel	w/o	Low P	High P	Token
Baseline Models	Radar	48.02	16.16	4.24	14.14
	LLM-Det	34.24	33.38	32.72	19.30
	Fast Detect	60.36	13.09	3.44	43.24
	Deep Fake	66.00	57.03	35.44	63.78
	TF-Model	88.06	68.58	66.10	18.70
	RR-Model	78.38	63.51	61.60	64.12
	SynthID Text	83.65	77.27	70.21	72.34
Our Model	$\delta = 3.0$	91.45	85.29	76.42	77.81
	$\delta = 4.0$	95.20	88.14	79.56	80.32
	$\delta = 5.0$	98.31	92.79	84.36	83.87

Table 1: F1 Score Performance Comparison Under Different Attack Scenarios (values in %). Here P  $\rightarrow$  Paraphrasing; Token  $\rightarrow$  Token attack; w/o  $\rightarrow$  without any attack

<sup>3</sup>https://huggingface.co/meta-llama/Llama-3. 1-8B-Instruct 352 353 354

341

342

343

345

346

347

348

350

351

355

356

357

358

360 361

362 363 364

365 366 367

368 369 370

371

We evaluate our method against multiple AIgenerated text detectors, including RADAR (Hu et al., 2023), DEEP-FAKE (Li et al., 2023) and Fast-Detect GPT (Bao et al., 2023b). Additionally, we evaluated against specialized AI-generated text detectors for peer review, such as TF-Model (which leverages term frequency of AI-generated tokens) and RR-Model (a regeneration-based method) (Kumar et al., 2024). For *Watermarking based* methods, we have included two baselines: WLLM (Kirchenbauer et al., 2023) and SynthID Text (Dathathri et al., 2024).

372

373

374

378

381

390

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418 419

420

421

499

423

We used both AI-generated reviews and human reviews for this experiment. During the attack phase, we targeted only the AI-generated reviews, as they are the ones intended to evade detection. As shown in the Table 1, existing AI detectors exhibit extreme sensitivity to adversarial attacks, with Fast Detect suffering a 94.30% drop  $(60.36\% \rightarrow 3.44\%)$ and Radar declining by 91.17% (48.02%  $\rightarrow$  4.24%) under high paraphrasing. Similarly, TF-Model's F1 score decreases by 78.76% (88.06%  $\rightarrow$  18.70%) under token attack, highlighting the brittleness of non-watermarked approaches. In contrast, our proposed watermarking method retains a performance of 84.36% under high paraphrasing and 83.87% under token attack ( $\delta$ =5.0), outperforming all baselines by a substantial margin. Even with lower  $\delta$ values, our model demonstrates resilience, with  $\delta$ =3.0 yielding 76.42% and  $\delta$ =4.0 yielding 79.56% under high paraphrasing, indicating consistent adversarial robustness. These findings emphasize that existing AI text detectors alone are insufficient for detecting AI-generated text under adversarial conditions. Our watermarking approach provides a promising solution for improving the resilience of AI generated peer review detection, even in challenging settings.

We also compare our approach with WLLM (Kirchenbauer et al., 2023), which relies solely on randomly selected green tokens for watermarking. In contrast, our method incorporates domainspecific tokens in addition to green tokens. Our results demonstrate that integrating domain-specific tokens significantly enhances watermark detectability, highlighting the importance of semantically meaningful token selection. We discussed this in detail in Section 4.4. Additionally, we compared our method with SynthID Text by injecting watermarks using its speculative sampling technique and detecting them based on the weighted mean detection score (ranging from 0 to 1), optimized



Figure 2: Effect of Watermarking on varying  $\gamma$  on with only green token (without Term) and with green token and blue token (With Term);  $\delta = 2.0$ .

on a validation set. As shown in Table 1, **SynthID** achieves F1 scores of 83.65 (no attack), 77.27 (low-perplexity paraphrasing), 70.21 (high-perplexity paraphrasing), and 72.34 (token substitution). Our method, evaluated under multiple threshold settings ( $\delta$ ), consistently obtains higher F1 scores across all conditions, demonstrating strong robustness and reliability under adversarial peer review scenarios.

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

Additionally, we found that the Query Type Identifier achieves an accuracy of 95.5% on the test set. Further, we studied the effect of varying the underlying base language model on detection accuracy. As shown in Figure 4, our approach maintains robust performance across LLMs of different sizes and architectures. For extended results and further discussion, please refer to Appendix H.

# **4.3** Effect of varying $\gamma$ on Detection Accuracy

Figure 2 shows that at low green token fractions  $(\gamma = 0.1 \text{ to } \gamma = 0.3)$ , detectability remains weak due to an insufficient statistical signal. When too few green tokens are available, the sampling algorithm operates largely unconstrained, following the model's natural probability distribution with minimal watermarking influence. As a result, the watermark imprint is inconsistent, leading to higher variance in detection scores. However, at  $\gamma = 0.3$ , detectability peaks, indicating an optimal balance where the watermarking method biases the sampling process enough to be recognized while still allowing diverse token choices. Beyond  $\gamma = 0.3$ , an interesting shift occurs. As  $\gamma$  increases, the green token fraction introduces greater randomness into the sampling process, allowing the model more flexibility in token selection. At  $\gamma = 0.4$ , this increased entropy makes the watermark signal less distinct, leading to a temporary decline in

detectability. Interestingly, at  $\gamma = 0.5$ , detectabil-460 ity recovers, possibly due to an optimal trade-off 461 which watermarking constraints are still strong 462 enough for recognition while allowing sufficient 463 linguistic variation to stabilize detection. Beyond 464 this point, performance declines again as higher 465 green token fractions ( $\gamma > 0.6$ ) further increase 466 randomness, making the text appear more natural 467 and reducing watermark signal strength. At very 468 high  $\gamma$  values (e.g.,  $\gamma = 0.9$ ), nearly all tokens in 469 the sampling space are green, making the sampling 470 distribution indistinguishable from unwatermarked 471 text, effectively neutralizing detectability. 472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

493

494

495

496

497

498

499

501

505

509

# 4.4 Effect of Domain-Specific Token Selection on Watermark Detectability

The results in Figure 2 demonstrate that incorporating important domain-specific tokens (blue tokens) alongside green tokens significantly improves watermark detectability across all thresholds compared to using only green tokens. This also highlights that our approach improves upon WLLM. The improvement is particularly notable at lower thresholds ( $\gamma = 0.1$  to  $\gamma = 0.3$ ), with performance gains exceeding 10% at  $\gamma = 0.1$  and  $\gamma = 0.2$ . This suggests that while random green token selection introduces high variability, leading to a weaker watermark signal, integrating important technical terms from the research paper enhances detection robustness by grounding the watermark in semantically meaningful and contextually significant words. Interestingly, at higher thresholds  $(\gamma > 0.6)$ , the performance difference reduces, likely due to the increased randomness in token selection making the watermark less distinguishable. These findings underscore the effectiveness of domain-aware token selection in improving watermark detectability while maintaining text fluency.

# **4.5** Effect of Watermarking Strength ( $\delta$ ) on Detection Accuracy

The graph demonstrates a positive correlation between watermarking strength ( $\delta$ ) and detection accuracy. As  $\delta$  increases from 2 to 6, the accuracy of watermark detection improves from 86.51% to 99.54%. This trend shows that increasing the watermarking bias enhances the distinguishability of AI-generated text. The primary reason for this improvement is that a higher  $\delta$  more strongly biases the model's token selection toward a predefined set of "green list" and "blue list" tokens, making it easier to detect the watermark statistically.



Figure 3: Effect Of Watermarking Strength ( $\delta$ ) on Perplexity and Accuracy;  $\delta = 2.0$ 

This controlled alteration in token probabilities increases the reliability of detection algorithms, as deviations from a natural distribution become more pronounced.

# **4.6** Effect of Watermarking Strength $(\delta)$ on Perplexity

Perplexity is a fundamental metric used to evaluate the confidence of a language model in its predictions. Lower perplexity values indicate that the model assigns higher probabilities to its predicted tokens, signifying more fluent and coherent text generation. Conversely, higher perplexity suggests greater uncertainty, implying that the text deviates from the model's natural distribution. In watermarking studies, minimizing the impact on perplexity is crucial to ensure that the watermarked text remains natural and human-like (Kirchenbauer et al., 2023).

From Figure 3, we observe a consistent increase in perplexity as the watermarking strength  $\delta$  increases from 2 to 6. Specifically, perplexity rises from 4.27 at  $\delta = 2$  to 6.30 at  $\delta = 6$ . At the same time, accuracy improves from 86.51% at  $\delta = 2$  to 99.54% at  $\delta = 6$ . This behavior occurs because watermarking forces the model to prefer certain tokens ("green list" or "blue list"), which may not always align with the most natural token choices. This trade-off is an essential consideration for watermarking techniques. While higher  $\delta$  ensures more robust watermark detection, excessive perplexity increases can negatively impact readability and coherence.

# 4.7 Impact of Watermarking on Downstream Decision Prediction

To further evaluate the practical quality of watermarked text, we performed a downstream task-

based analysis using aspect-based decision predic-546 547 tion. Specifically, we used the DeepASPeer (Kumar et al., 2022) which predicts paper acceptance 548 decisions from peer reviews using sentiment information for aspects such as novelty, substance, and soundness. Our results show that the model 551 achieved an accuracy of 75.6% on unwatermarked 552 reviews and 74.23% on watermarked reviews. This marginal drop of 1.74% suggests that watermarking has minimal impact on aspect-based sentiment structure and does not significantly degrade the 556 informativeness required for downstream decision-557 making tasks. This automatic evaluation further 558 complements our qualitative findings, indicating that watermarked reviews retain their functional quality for scholarly applications.

# **5** Robustness Analysis

562

563

569

570

571

574

575

582

584

586

588

590

591

594

Since reviewers may deliberately alter watermarked text to evade detection, we evaluate the robustness of our method.

# 5.1 GPT Paraphrasing

Given GPT's effectiveness in high-fidelity paraphrasing (Hassanipour et al., 2024), we employed it in two distinct settings: low-degree and highdegree paraphrasing. A detailed discussion of the paraphrasing procedure is provided in Appendix E. Experimental results reveal that GPT-based paraphrasing attacks substantially compromise the detection performance of existing models. Under high-degree paraphrasing, Radar and FastDetect-GPT perform particularly poorly, with F1 scores declining sharply to 4.24% and 3.44%, respectively. Even the strongest baseline, TF-Model, experiences a performance drop from 88.06% (no attack) to 66.10%, highlighting the susceptibility of current detectors to paraphrastic transformations. In contrast, our model exhibits significantly greater robustness, achieving F1 scores of 92.79% and 84.36% under low- and high-degree paraphrasing, respectively, at a watermarking strength of  $\delta = 5.0$ . As  $\delta$  increases, the model maintains higher detection accuracy, even under aggressive paraphrasing conditions.

# 5.2 Token Attack

We also performed a token attack (adjective) (Kumar et al., 2024). The Adjective Attack targets frequently occurring adjectives in AI-generated text and replaces them with their less frequent synonyms while preserving the overall meaning. The results show that baseline models struggle to maintain performance under this attack. For instance, Radar and LLM-Det experience substantial drops in F1 scores, reducing to 14.14% and 19.30%, respectively. Similarly, TF-Model and RR-Model, which initially performed well without attacks, decline to 18.70% and 64.12%, indicating their vulnerability to subtle lexical transformations. In contrast, our model remains highly robust, achieving 83.87% F1 at  $\delta = 5.0$ , demonstrating its ability to detect AI-generated text even when common adjectives are perturbed. These findings underscore the susceptibility of existing detectors to lexical style attacks and the effectiveness of our method under such perturbations. 595

596

597

598

599

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

# 6 Human Analysis

We conducted a qualitative analysis of 80 peer reviews generated under different watermarking intensities ( $\delta = 3.0, 4.0, 5.0$ ) to assess their Coherence, Consistency, and Fluency. The evaluation was conducted by three experts in ML and scientific writing, each with 10+ years of experience and 15+ publications. They independently assessed the reviews and resolved discrepancies through discussion, ensuring consensus. We found that  $\delta$ -3.0 was the most readable,  $\delta$ -4.0 introduces some rewording but remains logically coherent and effective, and  $\delta$ -5.0 resulted in overly complex phrasing that could hinder comprehension. Additionally, the blue list contributed to an increased density of technical terms in  $\delta$ -5.0, making the reviews more complex but not necessarily more informative. We discuss this in detail with examples in Appendix B.

# 7 Conclusion and Future Work

In this work, we introduced a novel watermarking framework for detecting LLM-generated peer reviews. Through extensive evaluation on ICLR and NeurIPS data, our method demonstrated consistently higher detection accuracy than existing baselines, especially under various adversarial attacks. While watermarking is still an emerging technique, we believe our framework offers a scalable and low-overhead approach to enhancing the reliability of peer review by enabling traceable detection of AI-generated peer reviews supporting editors and chairs in preserving trust within scholarly communication.

In future, we aim to extend detection to hybrid AI-human-generated reviews.

# Limitations

Our method of generating paper seed is sensitive to paper text. If a paper text is highly modified, the green token selection could change unpredictably, 647 making wrong detection. A more robust hashing mechanism (e.g., leveraging semantic embeddings rather than text-based hashing) could improve stability. Our method is tailored for reviews that are entirely AI-generated. However, a reviewer might draft key bullet points on a paper and then use ChatGPT to develop them into full paragraphs. We 654 655 recommend investigating this aspect in future research. While our current evaluation focuses on ML conferences to ensure experimental rigor and comparability, we agree that extending the evaluation to other domains (e.g., journals or interdisciplinary venues) would provide valuable generalization insights. Also, watermark effectiveness can be affected by the model's familiarity with domainspecific content. If an LLM fails to appropriately incorporate key technical terms, it may underutilize watermarked tokens, potentially weakening the signal or resulting in false negatives.

Also, our proposed generative watermarking framework, like other watermarking approaches, does not provide a complete solution for detecting AI-generated text; rather, it serves as a complement to other detection methods. In particular, applying such watermarks requires cooperation among the entities deploying LLM-based peer review systems. We discuss more about the practical deployment of this framework in detail in Appendix G. Detecting AI-generated text from entities that choose not to use watermarking requires alternative strategies, such as post hoc analysis. Additionally, the growing prevalence of open-source models poses a significant challenge, as their decentralized deployment makes watermark enforcement difficult (Dathathri et al., 2024).

# Ethics Statement

667

671

673

675

679

For this study, we used an open-source dataset. We do not take a stance on whether using AI tools for peer reviews is inherently positive or negative, nor do we claim definitive evidence that reviewers are relying on ChatGPT for drafting. The primary goal of this system is to aid editors/chair in detecting potentially AI-generated reviews, and it is designed solely for internal editorial use, not for authors or reviewers.

Although this watermarking method is designed

to mitigate the misuse of AI in peer reviews, it also 694 introduces potential risks. For instance, if the water-695 marking mechanism of a specific LLM were to be 696 publicly exposed, a malicious actor could exploit 697 it to generate unethical content embedded with the 698 model's watermark. To prevent such misuse, we 699 strongly recommend safeguarding the integrity of 700 the system by keeping key components such as the 701 hash function keys used for green and red list par-702 titioning confidential and restricted to authorized 703 users. 704

# References

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

706

707

708

709

710

711

712

713

714

715

716

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

- Bruce Alberts, Brooks Hanson, and Katrina L Kelner. 2008. Reviewing peer review.
- Martijn Arns. 2014. Open access is tiring out peer reviewers. *Nature*, 515(7528):467–467.
- Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. 2023a. Fast-detectgpt: Efficient zero-shot detection of machine-generated text via conditional probability curvature. *CoRR*, abs/2310.05130.
- Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. 2023b. Fast-detectgpt: Efficient zero-shot detection of machine-generated text via conditional probability curvature. *arXiv preprint arXiv:2310.05130*.
- Lutz Bornmann and Rüdiger Mutz. 2015. Growth rates of modern science: A bibliometric analysis based on the number of publications and cited references. *Journal of the association for information science and technology*, 66(11):2215–2222.
- Sumanth Dathathri, Abigail See, Sumedh Ghaisas, Po-Sen Huang, Rob McAdam, Johannes Welbl, Vandana Bachani, Alex Kaskasoli, Robert Stanforth, Tatiana Matejovicova, Jamie Hayes, Nidhi Vyas, Majd Al Merey, Jonah Brown-Cohen, Rudy Bunel, Borja Balle, A. Taylan Cemgil, Zahra Ahmed, Kitty Stacpoole, Ilia Shumailov, Ciprian Baetu, Sven Gowal, Demis Hassabis, and Pushmeet Kohli. 2024. Scalable watermarking for identifying large language model outputs. *Nat.*, 634(8035):818–823.
- Alexander R Fabbri, Wojciech Kryściński, Bryan Mc-Cann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2020. Summeval: Re-evaluating summarization evaluation. *arXiv preprint arXiv:2007.12626*.
- Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng

abs/2301.07597. lodging complaint after complaint about the process at top-tier journals. what's wrong with peer review? 805 Soheil Hassanipour, Sandeep Nayak, Ali Bozorgi, *The scientist*, 20(2):26–35. 806 Mohammad-Hossein Keivanlou, Tirth Dave, Abdulhadi Alotaibi, Farahnaz Joukar, Parinaz Mellatdoust, Mary L McHugh. 2012. Interrater reliability: the kappa 807 Arash Bakhshi, Dona Kuriyakose, et al. 2024. The statistic. Biochemia medica, 22(3):276–282. 808 ability of chatgpt in paraphrasing texts and reducing plagiarism: a descriptive analysis. JMIR Medical Eric Mitchell, Yoonho Lee, Alexander Khazatsky, 809 *Education*, 10(1):e53308. Christopher D. Manning, and Chelsea Finn. 2023. 810 Detectgpt: Zero-shot machine-generated text detec-811 Xiaomeng Hu, Pin-Yu Chen, and Tsung-Yi Ho. 2023. tion using probability curvature. In International 812 RADAR: robust ai-text detection via adversarial Conference on Machine Learning, ICML 2023, 23-813 learning. In Advances in Neural Information Pro-29 July 2023, Honolulu, Hawaii, USA, volume 202 814 cessing Systems 36: Annual Conference on Neural of Proceedings of Machine Learning Research, pages 815 Information Processing Systems 2023, NeurIPS 2023, 24950-24962. PMLR. 816 New Orleans, LA, USA, December 10 - 16, 2023. National Institute of Standards and Technology (NIST). 817 ISO/IEC JTC 1/SC 42. 2023. Standards for artificial 2023. Ai 100-4: Reducing risks posed by synthetic 818 intelligence: Trustworthiness. https://www.iso. content. Technical report, NIST. 819 org/committee/6794475.html. Accessed: 2024-Nihar B Shah. 2022. Challenges, experiments, and com-820 putational solutions in peer review. Communications 821 John Kirchenbauer, Jonas Geiping, Yuxin Wen, of the ACM, 65(6):76-87. 822 Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In Interna-Wajiha Shahid, Yiran Li, Dakota Staples, Gulshan Amin 823 tional Conference on Machine Learning, ICML 2023, Gilkar, Saqib Hakak, and Ali A. Ghorbani. 2022. Are 824 23-29 July 2023, Honolulu, Hawaii, USA, volume you a cyborg, bot or human? - A survey on detecting 825 202 of Proceedings of Machine Learning Research, fake news spreaders. IEEE Access, 10:27069-27083. 826 pages 17061-17084. PMLR. Irene Solaiman, Miles Brundage, Jack Clark, Amanda 827 Sandeep Kumar, Hardik Arora, Tirthankar Ghosal, and Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, 828 Asif Ekbal. 2022. Deepaspeer: towards an aspectand Jasmine Wang. 2019. Release strategies and 829 level sentiment controllable framework for decision the social impacts of language models. CoRR, 830 prediction from academic peer reviews. In Proceedabs/1908.09203. 831 ings of the 22nd ACM/IEEE Joint Conference on Digital Libraries, pages 1–11. Kyle Wiggers. 2022. Openai's attempts to watermark ai 832 text hit limits. TechCrunch, December, 10. Sandeep Kumar, Mohit Sahu, Vardhan Gacche, 833 Tirthankar Ghosal, and Asif Ekbal. 2024. 'quis custodiet ipsos custodes?' who will watch the watch-Max Wolff. 2020. Attacking neural text detectors. 834 men? on detecting AI-generated peer-reviews. In CoRR, abs/2002.11768. 835 Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages Rui Ye, Xianghe Pang, Jingyi Chai, Jiaao Chen, Zhenfei 836 22663-22679, Miami, Florida, USA. Association Yin, Zhen Xiang, Xiaowen Dong, Jing Shao, and 837 for Computational Linguistics. Siheng Chen. 2024. Are we there yet? revealing the 838 risks of utilizing large language models in scholarly 839 Yafu Li, Qintong Li, Leyang Cui, Wei Bi, Longyue peer review. arXiv preprint arXiv:2412.01708. 840 Wang, Linyi Yang, Shuming Shi, and Yue Zhang. 2023. Deepfake text detection in the wild. arXiv Xuandong Zhao, Prabhanjan Ananth, Lei Li, and 841 *preprint arXiv:2305.13242.* Yu-Xiang Wang. 2023. Provable robust water-842 marking for ai-generated text. arXiv preprint 843 Weixin Liang, Zachary Izzo, Yaohui Zhang, Haley Lepp, arXiv:2306.17439. 844 Hancheng Cao, Xuandong Zhao, Lingjiao Chen, Haotian Ye, Sheng Liu, Zhi Huang, Daniel A. McFarland, **Details on Implementation** Α 845 and James Y. Zou. 2024. Monitoring ai-modified content at scale: A case study on the impact of chatgpt on Our goal is to simulate a real-world scenario in 846 AI conference peer reviews. CoRR, abs/2403.07183. which a reviewer might use a LLM to generate a 847 Aiwei Liu, Leyi Pan, Yijian Lu, Jingjing Li, Xuming peer review by providing the manuscript as input. 848 Hu, Xi Zhang, Lijie Wen, Irwin King, Hui Xiong, Therefore, in our framework, if the input consists 849 and Philip Yu. 2024. A survey of text watermarking of a research paper and an unsafe query, the output 850 in the era of large language models. ACM Computing is the corresponding peer review. Surveys, 57(2):1-36. 851

Alison McCook. 2006. Is peer review broken? submis-

sions are up, reviewers are overtaxed, and authors are

802

803

746

747

749

751

752

755

756

757

760

762

763

773

774

775

778

779 780

781

785

791

792

793

795

796

797

801

05-17.

Wu. 2023. How close is chatgpt to human experts?

comparison corpus, evaluation, and detection. CoRR,

857

The model was loaded in FP16 precision, with a fixed PyTorch generation seed (123) for reproducibility. The generation parameters were configured as follows: top\_k = 0, temperature = 0.7, and beam size = 1. We use the below generation prompt for our experiments :-

**System**: You are a Research Scientist. Your task is to thoroughly and critically read the paper and write a peer review of it. **User:** Instructions 1. Read the paper critically and only write a peer review. Do not include any other content.

2. The peer review must contain the following sections: - Paper Summary: A concise summary of the paper's key contributions and findings. - Strengths: Highlight the notable strengths of the paper. - Weaknesses: Identify any limitations or areas of concern. - Suggestions for Improvement: Provide constructive feedback for the authors to enhance their work. - Recommendation: State whether the paper should be accepted, revised, or rejected.

Paper: {paper\_content}

To test the efficiency of the Query Type Identifier, we manually created 150 queries, equally divided into unsafe and safe categories. We divided this into 50% for validation and 50% test. We used the same model for this task as we did for generation, i.e., Llama-3.1-8B-Instruct. The watermark classifier was trained using a fully connected neural network with two hidden layers (16 and 8 neurons, both with ReLU activation) and an output layer of size 2 for binary classification. The dataset was standardized using StandardScaler and evaluated using 5-fold stratified cross-validation. Each fold had an 80-20% split for training and validation, with one fold reserved for testing. The model was optimized using the Adam optimizer with a learning rate of 0.001 and weight decay of 1e-4, and trained using cross-entropy loss. Early stopping was applied with a patience of 500 epochs and a maximum of 10,000 epochs, selecting the best model based on validation loss. Model performance was evaluated using accuracy with final results averaged across all folds. All experiments were conducted on an NVIDIA A100 80GB GPU using PyTorch.

# **B** Detailed Human Evaluation

Following the annotation guidelines for Coherence, Consistency, and Fluency (Fabbri et al., 2020), we asked the annotators to rank the three outputs. They discussed any discrepancies and reached an agreement when their ranking were different. The annotators were paid 20 USD per hour. We found that  $\delta = 3.0$  performed better in terms of Coherence, Consistency, and Fluency in 87%, 89%, and 92% of the cases, respectively. Similarly, for  $\delta = 4.0$ , we found that it performed better than  $\delta = 5.0$  in 77%, 79%, and 82% of the cases for Coherence, Consistency, and Fluency, respectively. To quantify agreement, we computed the inter-annotator agreement, achieving a Cohen's Kappa coefficient (McHugh, 2012) of 0.8629, which indicates strong agreement among annotators. Based on their comments we discuss the below observation:-

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

# **B.1** Linguistic Fluency and Readability

We found that increasing the watermarking strength progressively reduced linguistic fluency. Reviews generated with  $\delta$ -3.0 exhibited natural and wellstructured sentences, while  $\delta$ -4.0 introduced slight verbosity and rewording. However,  $\delta$ -5.0 resulted in excessive sentence expansion, leading to unnatural phrasing and reduced readability.

# **B.1.1** Example ( $\delta$ -3.0 vs. $\delta$ -4.0 vs. $\delta$ -5.0)

- δ-3.0: "The proposed model effectively reduces computational complexity while maintaining comparable performance with stateof-the-art methods. However, additional evaluation on out-of-distribution tasks would strengthen the paper."
- δ-4.0: "The proposed model provides an effective approach to reducing computational complexity while ensuring that performance remains competitive with state-of-theart methodologies. Further assessment on out-of-distribution tasks could help verify its robustness."
- $\delta$ -5.0: "The proposed model, as introduced 923 by the authors, offers a compelling approach 924 to addressing computational complexity while 925 ensuring that performance levels remain com-926 petitive with current state-of-the-art method-927 ologies. Nevertheless, to comprehensively val-928 idate the robustness of the approach, further 929 evaluation on out-of-distribution tasks should 930

860

861

871

876

878

879

882

1006

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

977

978

979

980

981

982

983

931 932

933

935

936

937

939

941

943

944

948

951

952

955

956

957

959

960

961

962

963

965

966

967

970

971

972

973

974

975

976

be conducted to provide a more complete assessment."

We observed that  $\delta$ -3.0 maintained conciseness,  $\delta$ -4.0 introduced slightly more complex phrasing without significant readability loss, and  $\delta$ -5.0 contained excessive verbosity, making the review harder to read.

# 938 B.2 Logical Coherence and Idea Flow

Logical coherence was largely preserved in  $\delta$ -3.0 and  $\delta$ -4.0, but  $\delta$ -5.0 introduced redundancy that disrupted idea flow. Higher watermarking levels resulted in multiple rewordings of the same point, artificially increasing the review length.

- **B.2.1** Example ( $\delta$ -3.0 vs.  $\delta$ -4.0 vs.  $\delta$ -5.0)
  - δ-3.0: "The LMUFormer architecture is well-designed and effectively combines the strengths of LMUs and Transformer models. However, the authors should provide a more detailed complexity analysis to strengthen their claims."
  - δ-4.0: "The LMUFormer architecture successfully integrates the advantages of LMUs and Transformer models while maintaining computational efficiency. However, a more detailed complexity analysis would further substantiate its effectiveness."
- δ-5.0: "The LMUFormer model, as presented in the paper, introduces a well-structured and well-thought-out architectural design that successfully integrates the advantages of LMUs and Transformer models. However, while the presented work is promising, an additional in-depth complexity analysis would be beneficial in order to further substantiate the claims made by the authors regarding the model's efficiency and applicability."

We found that  $\delta$ -3.0 was direct and logically structured,  $\delta$ -4.0 introduced slight elaboration while maintaining coherence, and  $\delta$ -5.0 resulted in unnecessary repetition, disrupting logical progression.

# **B.3 Redundancy and Verbosity**

We observed that  $\delta$ -5.0 significantly increased redundancy, whereas  $\delta$ -4.0 introduced only minor rewording.  $\delta$ -3.0 remained the most precise and concise.

# **B.3.1** Example ( $\delta$ -3.0 vs. $\delta$ -4.0 vs. $\delta$ -5.0)

- δ-3.0: "Conv-LoRA enhances SAM's segmentation performance by incorporating lightweight convolutional parameters. While this represents an effective extension, further real-world validation is needed."
- δ-4.0: "Conv-LoRA improves SAM's segmentation capabilities by introducing lightweight convolutional parameters, reinforcing its effectiveness in downstream tasks. However, additional real-world validation would help confirm its robustness."
- δ-5.0: "The Conv-LoRA framework introduces an effective approach for improving SAM's segmentation performance by integrating lightweight convolutional parameters. This enhancement allows SAM to perform better in various segmentation tasks. While this methodology is promising, additional realworld validation would further reinforce the practical utility and applicability of this approach."

We found that  $\delta$ -3.0 was the most precise,  $\delta$ -4.0 introduced slight elaboration without unnecessary repetition, and  $\delta$ -5.0 contained inflated and redundant phrasing.

# B.4 Technical Terminology and the Blue List Effect

We observed that  $\delta$ -5.0 contained a higher density of technical terms, likely due to the influence of the blue list. While this ensured technical accuracy, it also led to increased sentence complexity, making readability more difficult.

# **B.4.1** Example ( $\delta$ -3.0 vs. $\delta$ -4.0 vs. $\delta$ -5.0)

- δ-3.0: "The proposed fine-tuning approach effectively adapts the model to domain-specific segmentation tasks, ensuring efficient performance without significantly increasing parameter count."
- δ-4.0: "The fine-tuning strategy optimizes the model for domain-specific segmentation tasks, maintaining efficiency while minimizing parameter growth."
- δ-5.0: "The fine-tuning methodology proposed by the authors strategically integrates parameter-efficient training techniques
   1020

1023within the optimization framework to en-1024hance domain-specific segmentation tasks1025while maintaining computational efficiency1026and preserving model scalability."

# 1027 C Details about Blue Token Selection

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1042

1043

1044

1046

1048

1049

1050

1051

1052

1053

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1068

The Blue Token Selection process is designed to extract domain-specific technical terms from research papers, ensuring high relevance and precision. By leveraging a structured set of filtering rules, this approach systematically identifies key concepts, mathematical terms, dataset names, and acronyms while excluding common stopwords and generic phrases. Table 2 contains the detailed prompt we used for our experiment.

# D Details on output Reviews

Our generated peer reviews average 546 tokens, which is sufficient for reliable watermark detection. Prior studies (Liu et al., 2024) have demonstrated that watermarking techniques are effective on texts of moderate length. For example, watermarked texts exceeding 600 tokens have been shown to be generally robust against various attacks, including paraphrasing and rewriting. Moreover, our experimental results confirm that the proposed watermarking method performs reliably at the typical length of our generated reviews. We show an output with and without watermarking in Table 3.

#### E GPT Paraphrasing

Table 4 presents the prompt used for GPT-based paraphrasing. We used GPT-40 model for paraphrasing.

The decoded output is the watermarked text when compared with our models, ensuring watermark retention analysis, and the non-watermarked text for AI text detectors, allowing evaluation of AI-generated content detection.

# F Query Type Identifier

The Query Type Identifier is designed to classify queries related to peer review into Safe (S) or Unsafe (UN) categories based on their potential for ethical misuse. This classification system helps ensure that AI-generated content is not directly used in peer review submissions, thereby maintaining the integrity of the review process. Queries explicitly requesting a full peer review that could be submitted as-is are marked as Unsafe (UN), while those seeking explanations, clarifications, or sum-1069maries are classified as Safe (S). Table 5 presents1070the prompt used for the classification task.1071

1072

# **G** Implementation Strategy

Our watermarking mechanism integrates at the de-1073 coding stage of text generation and thus does not 1074 require direct access to proprietary model internals 1075 or explicit model identification. Editors or chairs 1076 would not need to detect which specific LLM re-1077 viewers use; instead, they can mandate a standard 1078 watermarking plugin provided as a lightweight li-1079 brary to ensure watermark insertion regardless of 1080 the LLM's origin. Policymakers have recognized 1081 watermarking as essential for ensuring content au-1082 thenticity, with standards bodies such as NIST ex-1083 plicitly recommending watermark integration for 1084 synthetic content provenance (of Standards and 1085 , NIST). Moreover, international standards communities (e.g., (ISO/IEC JTC 1/SC 42, 2023)) are 1087 actively developing watermarking methodologies, 1088 facilitating their inclusion into broader AI gover-1089 nance frameworks. Thus, chairs can convincingly 1090 advocate for widespread adoption by referencing these emerging guidelines and incentivizing com-1092 pliance through established governmental and in-1093 stitutional policies. 1094

#### H Effect of Varying Base LLMs

To rigorously assess the generalizability of our wa-1096 termark detection method, we evaluated our frame-1097 work in different settings using gemma-2-2b-it 1098 (2B parameters), Llama-3.1-8B-Instruct (8B 1099 parameters), and Qwen-14B (14B parameters). 1100 These models span an order of magnitude in size 1101 and differ notably in their tokenization schemes, 1102 decoding strategies, and inductive biases. Despite 1103 these substantial variations, our watermark detec-1104 tion approach consistently achieves high accuracy, 1105 yielding F1 scores of 96.9%, 98.3%, and 97.8%, 1106 respectively (see Figure 4). The minimal variation 1107 of only 1.4 percentage points underscores the wa-1108 termark's resilience to differences in the underlying 1109 language model. The method's consistent perfor-1110 mance across multiple LLMs demonstrates that our 1111 framework is model-agnostic and readily transfer-1112 able, making it a practical tool for watermarking in 1113 the peer review domain. 1114

Role	Content	
System User	You are a highly advanced AI specializing in scientific text processing. Your task is to extract important technical terms from a given research paper. These terms will be used for further analysis	
Instructions	1. Extract the following types of terms:	
	• <b>Technical Concepts</b> (e.g., "self-attention", "hyperparameter tuning", "zero-shot learning")	
	• Mathematical & Statistical Terms (e.g., "gradient descent", "log- likelihood estimation", "Bayes theorem")	
	• Machine Learning/Dataset Names (e.g., "ResNet", "BERT", "ImageNet", "MNIST")	
	• Key Nouns & Phrases Related to the Paper's Topic (e.g., "archi- tecture design", "model convergence", "loss function")	
	• Acronyms of Important Models & Techniques (e.g., "LSTM", "CNN", "SVM", "GAN")	
	• Scientific Terminology (e.g., "thermodynamic equilibrium", "quan- tum entanglement", "protein folding")	
	2. Do NOT include:	
	• Common Stopwords (e.g., "and", "or", "the", "but", "therefore")	
	• General Academic Phrases (e.g., "this paper presents", "in conclu- sion", "as shown in Figure")	
	• Adverbs or Common Verbs (e.g., "significantly", "appears", "seems", "performs")	
	• Generic Words Unrelated to the Paper's Topic (e.g., "data", "study", "results", "important", "analysis")	
	3. Output Format:	
	• Provide the extracted terms in a single, comma-separated string without duplicates.	
Input	Paper: {paper_content}	

Table 2: Prompt for Extracting Technical Terms from Research Papers

Aspect	Peer Review Without Watermark	Peer Review With Watermark
Paper Summary	Proposes ModernTCN, a pure convo- lutional model for time series anal- ysis. Modernizes TCN to handle long- and short-term tasks efficiently, outperforming Transformer models across five key tasks.	Introduces ModernTCN with cross- time and cross-variable dependency modeling, inspired by computer vi- sion. Shows state-of-the-art per- formance with efficient architecture across five tasks.
Strengths	<ul> <li>Identifies limitations in traditional TCNs.</li> <li>Demonstrates task generality and efficiency.</li> <li>Outperforms Transformerbased models.</li> <li>Offers comprehensive experimental results.</li> </ul>	<ul> <li>Addresses prior convolutional limitations with modern design.</li> <li>Leverages decoupling structure for dependency modeling.</li> <li>Includes ablation studies and ERF analysis.</li> <li>Strong empirical results across tasks.</li> </ul>
Weaknesses	<ul> <li>Assumes data stationarity.</li> <li>Lacks detail on TCN modifications.</li> <li>No comparison with other convolution models.</li> <li>No computational complexity analysis.</li> </ul>	<ul> <li>Assumes multivariate input; no univariate discussion.</li> <li>Omits complexity comparison with prior work.</li> <li>No application examples in specific domains.</li> </ul>
Suggestions for Improvement	<ul> <li>Detail time series modifications to TCN.</li> <li>Compare with convolutional baselines.</li> <li>Include complexity analysis.</li> <li>Discuss potential real-world applications.</li> </ul>	<ul> <li>Clarify assumptions around data types.</li> <li>Add complexity comparison with other models.</li> <li>Present domain-specific use cases.</li> </ul>

Table 3: Comparison of Peer Reviews Generated Without and With Watermark

Prompt Type	Content
GPT-40 Para- phrasing Prompt	<ul><li>System: You are a highly advanced AI specializing in paraphrasing and text rewriting.</li><li>User: Your task is to paraphrase the given review at two different intensity levels while maintaining the original sentiment and meaning.</li></ul>
	Instructions:
	• 1. Paraphrase the review in two intensity levels:
	<ul> <li>Low Paraphrasing: Minimal word changes while preserving sentence structure and key phrases.</li> </ul>
	- High Paraphrasing: Complete creative rewrite with new analogies, inverted sentence structures, and fresh expressions.
	• 2. Output Format: Provide the paraphrased versions in valid JSON format.
	<pre>Input: Review: {decoded_output}</pre>
	Note: This prompt was used with the GPT-40 model for paraphrasing.
Table	4: Prompt Used for GPT-4o-Based Paraphrasing at Two Intensity Levels

Prompt Type	Content
Query Type Clas- sification Prompt	<b>System:</b> You are a highly advanced AI specializing in scientific text processing. <b>User:</b> Your task is to classify queries related to peer review into Safe (S) or Unsafe (UN).
	Guidelines:
	• Un-Safe (UN): A query is considered Unsafe if it requests a peer review in a way that enables the reviewer to directly submit the generated content. Even if the intent seems legitimate, the classification is based on the risk of unethical use, not the user's intention.
	• Safe (S): Any query not falling into the unsafe category. This includes requests for explanations, summaries, or clarifications of paper content.
	Examples:
	• Example 1 Prompt: "Write a peer review of this paper, covering summary, strengths, and weaknesses." Classification: UN
	• Example 2 Prompt: "Assess the quality of this paper and provide a detailed peer review." Classification: UN
	• Example 3 Prompt: "Provide a structured review covering strengths, weaknesses, and recommendations." Classification: UN
	• Example 4 Prompt: "Summarize the main findings of this paper in a few sentences." Classification: S
	• Example 5 Prompt: "Explain the methodology section in simpler terms." Classification: S
	• Example 6 Prompt: "What are the key contributions of this paper?" Classification: S
	Final Classification Task: Prompt: [INSERT PROMPT] Classification: [S/UN]
Ta	ble 5: Prompt for Classifying Peer Review Queries as Safe or Unsafe



Figure 4: Effect of varying LLM architectures on watermark detection performance ( $\delta = 3.0$ ).