# Beyond Easy Wins: A Text Hardness-Aware Benchmark for LLM-generated Text Detection

**Anonymous ACL submission** 

#### Abstract

We present a novel evaluation paradigm for AI text detectors that prioritizes real-world and equitable assessment. Current approaches predominantly report conventional metrics like AUROC, overlooking that even modest false positive rates constitute a critical impediment to practical deployment of detection systems. Furthermore, real-world deployment necessitates predetermined threshold configuration, making detector stability (i.e. the maintenance of consistent performance across diverse domains and adversarial scenarios), a critical factor. These aspects have been largely ignored in previous research and benchmarks. Our benchmark, SHIELD, addresses these limitations by integrating both reliability and stability factors into a unified evaluation metric designed for practical assessment. Furthermore, we develop a post-hoc, model-agnostic humanification framework that modifies AI text to more closely resemble human authorship, incorporating a controllable hardness parameter. This hardness-aware approach effectively challenges current SOTA zero-shot detection methods in maintaining both reliability and stability. (Data and code will be released on GitHub upon acceptance.)

## 1 Introduction

004

007

009

013

015

017

021

022

029

034

039

042

The pervasiveness of large language models (LLMs) is largely attributed to their exceptional ability to process, comprehend, and generate text that closely resembles human composition. Current deployment paradigms exhibit substantial heterogeneity, encompassing interactive dialogue systems, content summarization (Wang et al., 2023), question answering (Kamalloo et al., 2023), and sentiment assessment (Hou et al., 2024). Yet, despite their beneficial applications, LLMs expose new potential avenues for malicious exploitation. Such harmful practices include, but are not limited to, automated disinformation dissemination (Vykopal et al., 2024), academic plagiarism and cheating (Cotton et al., 2024; Wahle et al., 2022), and the fabrication of deceptive reviews (Chiang et al., 2023). Beyond deliberate misuse scenarios, the automated identification and filtration of LLMgenerated content from training corpora has become imperative for preserving the integrity of contemporary human-generated information in training datasets (Wu et al., 2025). This process facilitates the development of models with current knowledge and mitigating the risk of cascading hallucinations in LLMs (Rawte et al., 2023). 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

The subtlety of distinguishing recurring patterns in LLM-generated text renders human classification efforts scarcely better than chance (Uchendu et al., 2021; Clark et al., 2021; Dou et al., 2022). Consequently, research emphasis has shifted toward the development of automatic detection tools. Current detectors encounter several critical shortcomings that compromise their robustness and reliability. Most prominently, their inability to generalize to out-of-distribution cases leads to failures when analyzing texts generated by unseen models or characterized by unfamiliar stylistic nuances (Kuznetsov et al., 2024; Lai et al., 2024). Furthermore, detector efficacy is significantly diminished through minimal perturbations, text length modifications, or the application of adversarial techniques (Zhou et al., 2024; Huang et al., 2024) including paraphrasing (Hu et al., 2023), stylistic transformation, and intentional insertion of errors (Dugan et al., 2024).

In efforts to improve detector robustness, most existing studies predominantly report conventional metrics, such as accuracy (Kuznetsov et al., 2024), F1-score (Guo et al., 2024a), and AUROC (Yu et al., 2024b; Su et al., 2023; Mitchell et al., 2023; Bao et al., 2023) when assessing performance under diverse adversarial attacks. However, this evaluation paradigm manifests several critical limitations in **real-world assessment** of detectors. Primarily, even modest false positive rates (FPR) are fundamentally unacceptable in LLM text detection contexts. For instance, in academic integrity applications, where the objective is to ensure fairness by identifying instances of academic dishonesty, misclassification of legitimately authored student work introduces significant procedural inequity by penalizing students undeservedly. Consequently, some researchers have transitioned toward reporting true positive rates (TPR) at fixed FPR (e.g. 1%) (Hans et al., 2024; Yang et al., 2023). However, despite this shift, additional unresolved issues remain. When deployed in real-world applications, detection systems require configuration with a predetermined threshold, independent of the generative provenance of examined text. Within this practical framework, quantifying detector stability through analysis of threshold dynamics across diverse adversarial conditions becomes critically important, a dimension that previous studies have largely overlooked. Consequently, aforementioned metrics provide inadequate characterization of practical detector efficacy.

086

090

100

101

102

103

104

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

131

132

133

134

135

In this paper, we propose novel evaluation metrics that facilitate more equitable comparative assessment of detection methods by simultaneously accounting for FPR impact on performance and detector stability variation across diverse scenarios. We integrate these multidimensional considerations into a unified metric that comprehensively characterizes both the reliability and stability of detection systems under real-world implementation conditions. In addition, we present a post-hoc, model-agnostic framework designed to steer LLMgenerated texts toward more human-like word distributions across calibrated difficulty gradients. This humanification process spans multiple hardness levels and is implemented through three key strategies: a) Random meaning-preserving mutation, b) AI-flagged word swap, and c) Recursive humanization loop. These strategies specifically target vulnerabilities in contemporary zero-shot detection approaches (Mitchell et al., 2023; Bao et al., 2023), which predominantly operate by perturbing texts and measuring token statistical properties. By progressively diminishing these detection signals while maintaining semantic coherence, our framework provides increasingly sophisticated evaluation scenarios that advance detector robustness assessment and illuminate the limitations of current detection approaches. In essence, this paper evaluates state-of-the-art detection systems using our hardness-aware benchmark (incorporating both challenging samples and our fairness-oriented metrics), offering a broader and more real-world evaluation framework that pushes detection efforts **"beyond easy wins**"! The core contributions of this paper are the following: 136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

162

163

164

165

166

167

168

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

- We formulate a novel evaluation paradigm that integrates both detector performance and stability while specifically penalizing elevated FPRs, thus ensuring fair and rigorous comparative assessments.
- We develop a model-agnostic generation framework that produces LLM-generated texts with controlled difficulty gradients to systematically evaluate detectors' performance.
- We compiled the largest dataset to date, consisting of both human-written and LLMgenerated texts prior to adversarial manipulation, see Table 1.

## 2 Related work

### 2.1 LLM-generated text detection

Detection of LLM-generated text falls into three principal categories: watermarking (Kirchenbauer et al., 2023; Liu and Bu, 2024; Panaitescu-Liess et al., 2025), supervised techniques (Guo et al., 2024a,b; Abassy et al., 2024; Yu et al., 2024a), and zero-shot approaches (Hans et al., 2024; Ma and Wang, 2024; Yang et al., 2023; Bao et al., 2023). Watermarking embeds imperceptible signals during text generation. These approaches fail to protect unknowing third-party users and are susceptible to paraphrasing attacks (Pang et al., 2024). Supervised techniques train classifiers atop encoder-based backbones like RoBERTa (Liu et al., 2021) using annotated corpora. These approaches manifest considerable performance deterioration when applied to out-of-distribution contexts. Zero-shot methods operate without training requirements, exploiting LM generative mechanisms through statistical indicators including loglikelihood (Gehrmann et al., 2019), perplexity (Hans et al., 2024), token rank (Gehrmann et al., 2019; Su et al., 2023), and entropy (Lavergne et al., 2008). Many approaches require generating alternative text versions to detect statistical deviations (Mitchell et al., 2023; Yang et al., 2023), a computationally intensive process. This issue is mitigated through reducing required revisions, and

219

220

Benchmark name	Human samples	LLM samples	Num of Styles	Multiple LLMs	Hardness Levels?	Fair Metric?
BUST; (Cornelius et al., 2024)	3.2k	22k	3	1	×	×
DetectRL; (Wu et al., 2024)	11.2k	11.2k	4	1	×	×
HC3; (Guo et al., 2023)	59k	27k	1	×	×	×
MAGE; (Li et al., 2024)	154k	295k	7	1	×	×
M4GT-Bench; (Wang et al., 2024)	65k	88k	6	1	×	×
MGTBench; (He et al., 2024)	3k	18k	3	1	×	×
RAID; (Dugan et al., 2024)	15k	509k	8	1	×	×
SHIELD; (Ours)	87.5k	612.5k	7	1	1	1

Table 1: Comparative analysis of LLM-generated text detection benchmarks.

efficient sampling methods (Bao et al., 2023; Su et al., 2023).

#### 2.2 Benchmarks for LLM text detection

185

186

187

The literature presents multiple benchmarks for evaluating LLM-generated text detection, each 189 with varying characteristics in scale, diversity, and 190 evaluation methodology (Uchendu et al., 2021; Yu 191 et al., 2025; Pudasaini et al., 2025). RAID (Dugan et al., 2024) systematically examines robustness 193 across multiple decoding strategies. MAGE (Li et al., 2024) extends evaluation capabilities across 195 a broader spectrum of LLMs. DetectRL (Wu et al., 196 2024) focuses on vulnerability assessment through 197 implementation of adversarial attacks and pertur-198 bations. M4GT-Bench (Wang et al., 2024) con-199 tributes a multilingual evaluation framework, and HC3 (Guo et al., 2023) compiles one of the largest ChatGPT-centric datasets available. Despite these significant contributions, current benchmarks commonly lack samples with structured difficulty gra-204 dients and principled metrics that ensure fairness in practical comparisons. Our benchmark, SHIELD, represents the first benchmark to incorporate hu-207 manified samples with graduated hardness levels. Furthermore, SHIELD pioneers a fairness-aware evaluation methodology, thus filling critical gaps in 210 the current evaluation paradigm. Table 1 provides 211 a comparative analysis of our proposed benchmark 212 against existing benchmarks in English. The com-213 214 parison covers critical aspects including pre-attack dataset size, diversity of writing styles, utilization 215 of multiple LLMs, structured hardness levels, and 216 the presence of fairness-oriented evaluation metrics. 217 218

# **3** SHIELD benchmark: data creation, humanification, and metric design

This section introduces the methodology underlying our benchmark **SHIELD** (Scalable Hardness-Informed Evaluation of LLM Detectors).

#### 3.1 Data creation

SHIELD comprises seven diverse writing styles: semi-formal discourse from Medium posts, journalistic reporting from news sources, evaluative content from Amazon reviews, questionanswering text from Reddit's ELI5, scientific writing from arXiv abstracts, partisan-persuasion reporting from pink slime, and expository documentation from Wikipedia. Please refer to Appendix A.1 for additional data characteristics. To obtain the LLM-generated counterparts, we deployed seven models: Llama3.2-1b, Llama3.2-3b, Llama3.1-8b (Grattafiori et al., 2024), Mistral-7b (Jiang et al., 2023), Qwen-7b (Bai et al., 2023), Gemma2-2b, and Gemma2-9b (Mesnard et al., 2024) for rephrasing of human-written texts. Additional specifications regarding models and prompting are detailed in Appendix A.2. To guarantee human authorship, the dataset comprises exclusively pre-2021 data, predating the emergence of LLMs. The SHIELD dataset contains 87.5k human-written documents and 612.5k LLM-generated samples before the application of adversarial techniques or humanization processes. Complete statistical details are presented in Appendix A.3.

#### 3.2 Hardness-aware humanification

The core hypothesis of our approach is to replace words that strongly indicate LLM authorship with words indicative of human authorship. Initially, we quantify each word's impact on authorship inference by the following scoring function:

$$MI_i = \sum_{x} P(x|w_i) log(\frac{P(x|w_i)}{P(x)}) \qquad (1)$$

where x represents authorship, and the mutual information (MI) quantifies the extent to which observing the word  $w_i$  shifts our probabilistic belief regarding the text's authorship. This scoring is performed by the ranker module illustrated in Figure 1(a). Subsequently, we partition the vocabulary into two subsets: AI-associated A, and humanassociated  $\mathbb{H}$  vocabularies based on their usage frequencies  $f_i$ . This separation reflects whether a word predominantly contributes to the distribution of AI- or human-written texts. For humanification
of text, we implement three strategies: a) Random meaning-preserving mutation (RMM), b)
AI-flagged word swap (AWS), and c) Recursive
humanization loop (RHL). Please see Appendix
B for a sample text from each strategy.

### 3.2.1 Random meaning-preserving mutation

272

275

276

277

278

281

284

287

292

296

298

302

304

307

309

This approach simulates scenarios wherein malicious users substitute random words to circumvent detection or when  $\mathbb{A}$  and  $\mathbb{H}$  are inaccessible. Figure 1.b illustrates this process. Let  $\mathcal{D} = \{w_1, ..., w_n\}$ be an AI-generated text consisting of a sequence of words. We define  $\mathcal{S} \subset \mathcal{D}$  as the set of non-stopwords,  $S = \{w_i \in D | w_i \notin \text{StopWords}\}$ . Next, we randomly sample a subset  $\mathcal{M} \subset \mathcal{S}$  such that  $|\mathcal{M}| = p.|S|$ , with p representing the sampling ratio. Then, we construct a masked version  $\mathcal{D}^{mask} =$  $\mathbf{Mask}(\mathcal{D}, \mathcal{M})$  by replacing each word  $w_i \in \mathcal{M}$  in  $\mathcal{D}$  with the special token <mask>.  $\mathcal{D}^{mask}$  is fed into a masking language model  $\mathbf{f_{MLM}}$  which outputs a ranked list of predictions  $\mathbf{f}_{\mathbf{MLM}}^{(i)}(\mathcal{D}^{mask})$  at each masked position i. For each i, the first candidate  $\hat{w}_i$  with minimum rank is selected such that it differs from original word in  $\mathcal{D}, \hat{w}_i \neq w_i$ . Finally, the edited text  $\mathcal{D}^{edit}$  is produced by replacing each  $w_i \in \mathcal{M}$  with the corresponding  $\hat{w}_i$ :

$$\mathcal{D}^{edit} = \mathbf{Replace}(\mathcal{D}, \{(w_i, \hat{w}_i)\}_{w_i \in \mathcal{M}}) \quad (2)$$

#### 3.2.2 AI-flagged word swap

···

The second strategy, depicted in Figure leverages  $\mathbb{A}$  and  $\mathbb{H}$  to substitute AI-1(c),indicative words with human-characteristic alternatives. Without loss of generality, let  $\mathcal{S}' = \mathbf{Sort}(\mathcal{S}, \mathbf{MI}_{\mathbb{A}}) = [w_{(1)}, w_{(2)}, \dots, w_{(|\mathcal{S}|)}]$ such that  $\mathbf{MI}_{\mathbb{A}}(w_{(1)}) \geq \mathbf{MI}_{\mathbb{A}}(w_{(2)}) \geq \dots \geq \mathbf{MI}_{\mathbb{A}}(w_{(|\mathcal{S}|)}).$ Here, S' is the set S reordered in descending order based on MI scores with respect to the A. To construct set  $\mathcal{M}$ , we extract the top p% of words from  $\mathcal{S}'$ ,

$$\mathcal{M} = \{ w_{(i)} \in \mathcal{S}' | 1 \le i \le p.|\mathcal{S}| \}$$
(3)

The parameter p acts as a tunable knob that controls the hardness level of the humanified text. The words in  $\mathcal{M}$  undergo masking in the initial text.  $\mathbf{f}_{\mathbf{MLM}}$  then predicts candidate replacements for each masked position. For each masked position i, we select the word with the highest score in  $\mathbb{H}$ ,

311 
$$\hat{w}_i = \underset{w \in \mathbf{f}_{\mathbf{MLM}}^{(i)}(\mathcal{D}^{\mathbf{mask}})}{\arg \max \mathbf{MI}_{\mathbb{H}}(w)}$$
 (4)



Figure 1: Humanification strategies based on the ranked vocabularies  $\mathbb{A}$  and  $\mathbb{H}$  produced by the Ranker in (a). (b) Random meaning-preserving mutation (RMM), (c) AI-flagged word swap (AWS), (d) Recursive humanization loop (RHL).

The edited text  $\mathcal{D}^{edit}$  is derived similarly as Equation 2.

### 3.2.3 Recursive humanization loop

Figure 1(d) shows the third strategy which extends second strategy (AWS) through implementation of a recursive refinement. Let  $\mathcal{D}^{(0)}$  be the original AI text. We define the recursive editing process for Rrounds. At each round  $r \in \{1, 2, ..., R\}$ , the set of non-stop words are formed,

$$\mathcal{S}^{(r-1)} = \{ w_i \in \mathcal{D}^{(r-1)} | w_i \notin \text{StopWords} \}$$
(5)

Subsequently, words are ordered according to  $MI_{\mathbb{A}}$  scores,

$$\mathcal{S}^{\prime(r-1)} = \mathbf{Sort}(\mathcal{S}^{(r-1)}, \mathrm{MI}_{\mathbb{A}})$$
 324

312

313

314

315

316

317

318

319

321

322

$$=[w_{(1)}^{(r-1)}, w_{(2)}^{(r-1)}, ..., w_{(|\mathcal{S}|)}^{(r-1)}] \qquad (6)$$

389

390

391

392

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

369

370

371

326 327

011

328

329

332

333

334

335

336

337

341

342

343

345

346

352

356

357

358

360

364

365

368

A fixed proportion  $p=p_o$  (set at 10%) of words with maximal scores are selected for masking,

$$\mathcal{M}^{(r)} = \{ w_{(i)}^{(r-1)} \in \mathcal{S}'^{(r-1)} | 1 \le i \le p_o. |\mathcal{S}| \}$$
(7)

 $\mathbf{f}_{\mathbf{MLM}}$  predicts the masked words in  $D^{mask(r)}$ , which is obtained by  $\mathbf{Mask}(D^{(r-1)}, \mathcal{M}^{(r)})$ . For each masked position *i*, we choose the word with the highest score in  $\mathbb{H}$ ,

$$\hat{w}_{i}^{(r)} = \underset{w \in \mathbf{f}_{\mathbf{MLM}}^{(i)}(\mathcal{D}^{\mathbf{mask}, (\mathbf{r})})}{\arg \max} \mathbf{MI}_{\mathbb{H}}(w) \qquad (8)$$

Then the edited document for round r is:

$$\mathcal{D}^{(r)} = \operatorname{\mathbf{Replace}}(\mathcal{D}^{(r-1)}, \{(w_i, \hat{w}_i^{(r)})\}_{w_i \in \mathcal{M}^{(r)}})$$
(9)

The final humanified text is  $\mathcal{D}^{edit} = \mathcal{D}^{(R)}$ . The parameter R acts as a controllable knob to adjust the hardness level of the resulting text, with larger values yielding increasingly human-like phrasing.

#### **3.3** Fairness-oriented evaluation metric

#### **3.3.1** Reliability and performance metric

The AUROC metric can be interpreted as the probability that a randomly selected positive instance receives a higher score than a randomly selected negative one. AUROC neglects consideration of practical operational threshold regions. It uniformly weights the entire ROC curve, including high-FPR regions that are impractical for real-world deployment of AI-text detection systems. It also fails to capture performance instability, that is, significant changes in TPR or FPR due to small threshold adjustments. To more precisely evaluate the reliability of AI-text detection systems, we introduce weighted-AUROC (W-AUROC), defined as the expectation of TPR over a non-uniform probability distribution p(t) across FPR,

$$W-AUROC = \mathbb{E}_{t \sim p(t)}[TPR(t)]$$
(10)

where the weighting function is given by,  $p(t) = \frac{1}{Z} exp(-kt)$ , with decay parameter k > 0 and normalization constant  $\mathcal{Z} = \frac{1 - exp(-k)}{k}$ . To determine the decay parameter k, we set the exponential weighting function exp(-k.FPR) to decay to 50% of its initial value at FPR=0.05. This decision is inspired by prior works in AI-text detection that routinely report TPR at a fixed FPR of less than 5% as a key performance indicator, reflecting its status as a standard deployment-level operating point. This constraint yields  $k=20 \ln 2$ .

#### **3.3.2** Stability under FPR deviation (SFD)

To assess stability across different scenarios, we compute the standard deviation of FPR at decision thresholds determined by Youden's J statistic (Youden, 1950) on ROC curve. Our selection of FPR as the target variable stems from two considerations: 1) the threshold determination in each detection system is intrinsically dependent on scoring function values with ranges varying between methods, 2) this approach enables direct penalization of significant FPR fluctuations for stability assessment, as substantial FPR variability constitutes unacceptable performance in AI-text detection frameworks. For each detection system and across Mevaluation scenarios (e.g., generative model, attack types, or writing styles), we extract the threshold  $t_i^*$  that maximizes J(t) = TPR(t) - FPR(t)for each scenario *i*, and record the corresponding  $FPR_i^* = FPR_i(t_i^*),$ 

$$t_i^* = \arg\max_i \left[ \text{TPR}_i(t) - \text{FPR}_i(t) \right]$$
(11)

The standard deviation of these optimal FPRs across all M scenarios is calculated and denoted as  $\sigma_{\text{FPR}}$ . Finally, we define the stability metric as,

$$SFD = exp(-\lambda.\sigma_{FPR})$$
(12)

where  $\lambda > 0$  is a tunable hyperparameter controlling the sensitivity to instability. Lower standard deviation in FPR results in higher stability scores, with perfect stability ( $\sigma_{\rm FPR}=0$ ) yielding a maximum value of 1. To determine a principled value for the decay parameter  $\lambda$ , we calibrate it so that a moderate but practically noticeable instability in FPRs corresponds to a mid-range stability score. Therefore, we set it such that the stability score reduces to 0.5 when the standard deviation  $\sigma_{\rm FPR}$ reaches 0.1. This yields  $\lambda=10 \ln 2$ .

We adopt a multiplicative formulation to combine W-AUROC (performance) and SFD (stability) into a single unified reliability-stability score (URSS),

URSS = 
$$\left(\frac{1}{M}\sum_{i=1}^{M} \text{W-AUROC}_{i}\right)$$
. SFD (13)

Multiplication enforces a **non-compensatory relationship**: a high W-AUROC cannot mask poor stability, and conversely, robust stability does not necessarily indicate high discrimination capability. This reflects real-world deployment requirements, where even a highly accurate detector is unusable if unstable, and vice versa.



Figure 2: Radar charts of comparing detectors in different writing styles across all generative models.

#### 4 Experiments and discussion

416

417

418

419

420

421

499

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452 453

454

455

456

457

In this section, we organize our experiments and discussion around three core questions: 1) Why is AUROC insufficient for evaluating detectors in real-world settings? 2) How effective is our proposed humanification strategies in degrading SOTA zero-shot detectors? and 3) How robust are these detectors across different levels of humanification hardness? We evaluate six SOTA zero-shot detectors, Binoculars (Hans et al., 2024), Fast-DetectGPT (Bao et al., 2023), LRR (Su et al., 2023), Log-Likelihood, Log-Rank, and Rank (Gehrmann et al., 2019), along with a supervised model, Radar (Hu et al., 2023), to examine the comparative vulnerability of zero-shot methods. Further information on the detectors used can be found in Appendix C.

#### 4.1 Why is AUROC NOT enough?

For each writing style, we compare detectors across all generative models using the original LLMgenerated texts (without humanification). The resulting radar charts are presented in Figure 2 ( Please refer to Appendix D for radar charts of different LLMs across all writing styles). While traditional AUROC can exaggerate the superiority of a detector, our proposed framework reveals an illuminating truth. AUROC may obscure cases where detectors perform equivalently in practice. In contrast, URSS effectively exposes equivalences by examining performance parity within operationally low FPR regions and assessing stability. For instance, in the Reddit chart, although Radar exhibits a substantially higher AUROC than Rank, both achieve identical URSS, highlighting their practical equivalence. Conversely, there are cases where detectors achieve similar AUROC scores, yet one outperforms the other in low-FPR operational regions while also exhibiting greater threshold stability. Such multidimensional superiority is entirely masked when relying solely on the conventional AUROC metric. Representative cases include Rank vs. Log-Rank and Log-Rank vs. LogLikelihood in the Medium chart, and Binoculars vs. Fast-DetectGPT across News, Wikipedia, Medium, and Review charts. Despite significant disparities in either W-AUROC or SFD, identical AUROC across detectors can cause a misleading impression of equivalence. This potentially results in suboptimal choices for real-world deployment. Such patterns are seen in Fast-DetectGPT vs. Rank and Log-Likelihood vs. Log-Rank on arXiv, Log-Rank vs. Log-Likelihood in Reddit, and Rank vs. Log-Likelihood in Wikipedia. URSS suggests that meaningful performance equivalency between detectors can only be established through comprehensive evaluation that simultaneously considers both operational region sensitivity and crossscenario stability. For example, in the Wikipedia chart, URSS observes that Log-Likelihood has slightly lower W-AUROC but higher SFD than Log-Rank, and assigns them equal scores, reflecting fairness when each method excels in one dimension and the trade-off is not substantial.

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

A significant observation from our experiments is that methods exhibiting severe deficiencies in any critical performance dimension are heavily penalized, acknowledging that such limitations undermine real-world deployment potential. This evaluation principle is exemplified by Binoculars which, while achieving the highest AUROC and second-highest W-AUROC scores, was ultimately positioned last in the Reddit writing style analysis based on URSS metric due to its exceptionally poor stability.

#### 4.2 Effectiveness of humanification strategies

Table 2 shows the impact of different strategies on detector performance across all writing styles, using the largest LLM model from each family evaluated in this study. Corresponding results for smaller models are available in Appendix D. To enhance readability, all metrics are reported as percentages. The baseline consists of original LLMgenerated texts (specifically paraphrased texts without humanification). Zero-shot detectors suffer

$LLM \rightarrow$		Gemn	na2 9b			Llama3.1 8b				Mistral 7b				Owen 7b		
Detector $\downarrow$ Metric (%) $\rightarrow$	AUC	W-A	SFD	URSS	AUC	W-A	SFD	URSS	AUC	W-A	SFD	URSS	AUC	W-A	SFD	URSS
Paraphrase (baseline)																
Binoculars Fast-DetectGPT Log-Likelihood Log-Rank LRR Radar Rank	93.8 93.1 77.6 77.6 76.2 76.7 75.9	69.9 74.7 28.7 30.1 34.0 34.2 35.6	47.0 60.7 41.2 47.3 58.4 17.6 49.3	32.9 45.3 11.8 14.3 19.9 6.0 17.5	95.4 94.6 82.2 82.5 81.5 79.8 76.3	74.0 80.3 36.1 38.6 45.1 42.3 36.8	55.2 60.5 41.8 44.5 55.5 14.8 57.1	40.9 48.6 15.1 17.2 25.0 6.3 21.0	90.6 89.2 64.8 63.3 58.5 71.0 62.0	59.4 62.1 16.6 15.7 13.5 25.5 16.9	46.2 62.1 22.8 20.8 16.9 15.9 16.3	27.4 38.6 3.8 3.3 2.3 4.0 2.8	85.5 91.7 66.0 65.6 62.5 83.7 59.3	35.2 76.2 18.0 18.1 17.8 50.7 18.1	42.5 81.2 42.6 37.6 27.2 21.1 16.0	14.9 61.8 7.7 6.8 4.8 10.7 2.9
Random meaning-preserving mutation (RMM)																
Binoculars Fast-DetectGPT Log-Likelihood Log-Rank LRR Radar Radar Rank	77.9 77.0 34.4 36.9 50.4 89.0 48.4	35.4 37.8 4.7 5.4 9.6 52.4 7.0	35.1 35.8 5.9 7.3 11.8 34.0 9.5	12.4 13.5 0.3 0.4 1.1 17.8 0.7	83.6 82.2 41.2 44.2 58.5 88.1 50.1 <b>AI-fla</b>	44.7 47.8 8.4 9.6 16.4 52.3 8.4 gged w	43.8 46.3 6.2 6.9 10.3 27.8 9.2 ord swa	19.6 22.1 0.5 0.7 1.7 14.5 0.8 <b>p (AWS</b>	73.9 71.7 26.1 27.6 39.4 83.7 39.7	28.9 27.2 2.2 2.2 3.4 41.9 2.8	53.0 45.8 4.0 4.5 8.9 24.8 6.7	15.3 12.5 0.1 0.1 0.3 10.4 0.2	77.7 83.5 31.3 32.9 43.1 89.4 38.8	30.4 50.4 4.8 4.8 5.4 60.7 3.3	46.2 52.8 8.1 7.6 5.7 29.9 6.0	14.126.60.40.40.318.20.2
Binoculars Fast-DetectGPT Log-Likelihood Log-Rank LRR Radar Radar Rank	81.6 83.1 21.5 22.7 32.6 85.9 46.1	45.4 45.9 2.5 2.7 4.0 47.6 4.9	54.4 31.0 4.3 4.3 7.8 30.0 10.3	$\begin{array}{c} 24.7 \\ 14.2 \\ 0.1 \\ 0.1 \\ 0.3 \\ 14.3 \\ 0.5 \end{array}$	81.9 83.3 25.0 26.4 36.1 82.3 46.5	$\begin{array}{r} 46.1 \\ 49.6 \\ 4.7 \\ 5.1 \\ 6.8 \\ 40.3 \\ 5.5 \end{array}$	37.2 36.3 5.3 5.3 4.9 26.1 10.5	$17.2 \\18.0 \\0.3 \\0.3 \\0.3 \\10.5 \\0.6$	78.6 79.5 16.1 16.7 24.7 80.4 39.5	41.2 37.7 1.3 1.2 1.4 39.0 2.4	57.4 37.8 3.7 4.0 4.3 19.0 11.9	$\begin{array}{c} 23.6 \\ 14.3 \\ 0.0 \\ 0.0 \\ 0.1 \\ 7.4 \\ 0.3 \end{array}$	73.2 78.5 18.8 18.8 22.3 86.0 31.3	31.3 42.5 2.7 2.4 1.7 52.3 2.0	42.6 43.1 18.4 13.6 3.7 26.5 5.7	$13.3 \\18.3 \\0.5 \\0.3 \\0.1 \\13.9 \\0.1$
				R	ecursive	huma	nization	loop (R	HL)							
Binoculars Fast-DetectGPT Log-Likelihood Log-Rank LRR Radar Rank	75.8 76.5 17.5 20.2 36.4 85.4 45.1	$32.2 \\ 30.9 \\ 0.6 \\ 0.8 \\ 3.3 \\ 43.9 \\ 4.4$	59.3 50.2 3.5 3.7 11.6 34.5 8.8	$19.1 \\ 15.5 \\ 0.0 \\ 0.0 \\ 0.4 \\ 15.1 \\ 0.4$	$78.0 \\78.2 \\21.6 \\24.5 \\40.6 \\81.8 \\45.1$	36.2 37.3 1.8 2.4 6.4 37.9 4.5	40.0 42.4 3.8 4.2 5.8 27.0 9.0	$14.5 \\ 15.8 \\ 0.1 \\ 0.1 \\ 0.4 \\ 10.2 \\ 0.4$	75.0 74.8 12.4 14.4 29.4 79.9 39.0	$\begin{array}{c} 32.3 \\ 27.1 \\ 0.5 \\ 0.5 \\ 1.1 \\ 33.9 \\ 1.7 \end{array}$	63.4 50.2 3.3 3.2 10.9 27.1 7.2	$20.5 \\ 13.6 \\ 0.0 \\ 0.0 \\ 0.1 \\ 9.2 \\ 0.1$	72.6 77.4 17.2 18.3 27.6 86.6 32.6	26.6 36.7 1.9 1.8 1.7 51.2 1.9	44.9 61.0 8.3 6.7 3.6 30.9 4.5	11.9 22.4 0.2 0.1 0.1 15.8 0.1

Table 2: Performance of detectors under paraphrasing (baseline), RMM, AWS, and RHL strategies.

526

527

marked average URSS degradation of 41%, 62%, 97%, 96%, 92%, and 95% in **RMM**; 27%, 66%, 98%, 98%, 98%, and 95% in AWS; and 38%, 65%, 99%, 99%, 97%, and 97% in **RHL**, corresponding to Binoculars, Fast-DetectGPT, Log-Likelihood, Log-Rank, LRR, and Rank, respectively. This demonstrates that without devising robustness enhancements, zero-shot detectors fail to withstand word-level humanification, resulting in severe performance loss. Interestingly, even though AWS and RHL introduce more complex humanification than RMM, our experiments reveal that multiple detectors demonstrate similarly compromised URSS performance under the simpler RMM. For instance, Log-Likelihood's URSS drops to 0.1 under RMM for Mistral 7b, nearly equivalent to its 0.0 under AWS and RHL. This observation indicates that zero-shot detectors depend on fragile token-level statistical signatures that collapse under even modest perturbations. This vulnerability becomes particularly concerning considering that RMM more closely approximates natural user editing behaviors, rendering these detection systems unreliable in practical deployment scenarios. This limitation exacerbates the previously documented challenge of generative model dependency (Wu et al., 2024), as further evidenced in Table 2, which demonstrates significant variation in URSS metrics across different LLMs when subjected to identical strategies.

It is noteworthy that Radar exhibits a reversed trend: across all humanification strategies, word replacements improve its performance in both W-AUROC and SFD, leading to higher URSS. This may be attributed to Radar's design, which emphasizes robustness to paraphrasing. We leave further investigation into whether this behavior generalizes to other supervised methods for future work.

529

530

531

532

533

534

535

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

#### 4.3 Robustness under hardness levels

Figure 3 presents the trend of all evaluation metrics for each humanification strategy across their respective control parameters (i.e., p ranging from 10% to 100% for RMM and AWS, and R from 5 to 40 for RHL) for zero-shot detectors. Dashed lines indicate the baseline, corresponding to paraphrased text without any humanification. Notably, both performance metric (W-AUROC) and the stability metric (SFD) consistently fall below the baseline across all knob values, confirming the effectiveness of the proposed strategies. RHL achieves significant degradation in both W-AUROC and SFD starting from a low  $R \approx 15$ , maintaining this effect across the full range. AWS follows a similar trend at a moderate  $p \approx 60\%$ , while RMM requires a higher  $p \approx 80\%$  to reach comparable impact. This can be attributed to the fact that AWS and RHL deliberately replace high-entropy AI words with highentropy human words, whereas RMM introduces more random substitutions, resulting in less tar-



Figure 3: The performance of zero-shot detectors in RMM, AWS, and RHL strategies based on their hardness levels (p%, p%, and R, for RMM, AWS, and RHL, respectively).

geted degradation. A surprising observation is that some detectors exhibit increased stability as the hardness knob surpasses a certain threshold. This 560 phenomenon is particularly pronounced in Binocu-561 lars and Fast-DetectGPT. A possible explanation is that as the text becomes more human-like, the discriminative signal diminishes (evidenced by lower W-AUROC), leading to more homogeneity across LLMs and writing styles. Consequently, detec-566 tors struggle to differentiate content and converge to consistent decision thresholds, resulting in reduced  $\sigma_{\rm FPR}$ . While these systems exhibit improved consistency under intense adversarial conditions, their substantially degraded discriminative performance ultimately results in significantly penalized overall URSS scores. 573

An additional noteworthy observation derived 574 from Figure 3 is the pronounced convergence of 575 baseline trajectories in AUROC plots, which subsequently differentiates into distinctly separated lines across W-AUROC, SFD, and URSS plots. This transformation from convergent to divergent patterns across different evaluation metrics provides 581 further empirical validation for the conclusions presented in Section 4.1, why AUROC is not enough, and illustrates how AUROC may offer misleadingly "easy wins" to certain detectors that may fail in real-world applications. 585

## 5 Conclusion

In this work, we presented SHIELD, a comprehensive benchmark designed to advance the fair evaluation of AI text detectors under realistic deployment scenarios. Our results showed that conventional metrics like AUROC, despite their ubiquity, can significantly overestimate detector efficacy by neglecting critical operational constraints: the necessity of maintaining low FPR in practical AI text detection applications and the stability of detectors across varying threshold under deployment conditions where both writing style and generative model characteristics are typically unknown. SHIELD addresses these limitations through the introduction of USRR, a metric that integrates both performance measurement in low FPR regions and stability assessment. Complementarily, SHIELD introduces a scalable humanification framework for generating humanified texts across graded difficulty levels, facilitating robust stress-testing of detection systems. Our experimental findings reveal significant vulnerabilities in zero-shot detection methodologies, which exhibit performance degradation of approximately 80% on average when subjected to even the most rudimentary word manipulation scenarios evaluated in this study, perturbations that closely approximate natural user editing behaviors without requiring specialized knowledge of AI- or humanauthored text characteristics.

586

587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

## 615 Limitations

Several constraints merit acknowledgment within 616 our experimental framework. Primarily, our in-617 vestigation was confined to monolingual English 618 text analysis, precluding examination of multilin-619 gual detection scenarios that represent increasingly important deployment contexts given the global nature of AI text generation. This linguistic lim-622 itation potentially restricts the generalizability of our findings to diverse language environments. Additionally, the inherently dynamic nature of LLM 625 development presents a significant temporal con-626 straint; as generative architectures evolve through version updates and architectural innovations, their statistical signatures undergo corresponding transformations, necessitating periodic collection of representative text samples to maintain benchmark cur-631 rency and relevance. Furthermore, budgetary constraints confined our experimental protocol to opensource models, resulting in the exclusion of closedsource systems including ChatGPT and Claude. 635 The inclusion of these commercial platforms would enhance the comprehensiveness of our evaluation 637 framework, particularly considering their extensive deployment and potentially distinctive generative characteristics that may present unique challenges to detection methodologies. 641

## Ethical considerations

This investigation aims to evaluate the robustness of contemporary AI-text detection methods against adversarial manipulations. The proliferation of LLM-generated content and its potential for malicious applications necessitates robust detection 647 mechanisms to serve as effective countermeasures against synthetic text deception. Vulnerabilities in these detection frameworks could precipitate significant complications in computational forensics and 651 information verification processes. Consequently, this research endeavors to provide detection system engineers with rigorous adversarial testing frame-654 works for comprehensive validation of their algorithmic approaches against sophisticated evasion techniques. We explicitly stipulate that the methodologies and results documented in this study are intended exclusively for detection system improve-659 ment and validation, and not for circumvention of existing detection systems. 661

### References

Mervat Abassy, Kareem Elozeiri, Alexander Aziz, Minh Ngoc Ta, Raj Vardhan Tomar, Bimarsha Adhikari, Saad El Dine Ahmed, Yuxia Wang, Osama Mohammed Afzal, Zhuohan Xie, Jonibek Mansurov, Ekaterina Artemova, Vladislav Mikhailov, Rui Xing, Jiahui Geng, Hasan Iqbal, Zain Muhammad Mujahid, Tarek Mahmoud, Akim Tsvigun, and 5 others. 2024. LLM-DetectAIve: a tool for fine-grained machine-generated text detection. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 336–343, Miami, Florida, USA. Association for Computational Linguistics. 662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. 2023. Fast-detectgpt: Efficient zero-shot detection of machine-generated text via conditional probability curvature. *arXiv preprint arXiv:2310.05130*.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv:2004.05150*.
- Hung-Yun Chiang, Yi-Syuan Chen, Yun-Zhu Song, Hong-Han Shuai, and Jason S. Chang. 2023. Shilling black-box review-based recommender systems through fake review generation. In *Proceedings* of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23, page 286–297, New York, NY, USA. Association for Computing Machinery.
- Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 7282–7296, Online. Association for Computational Linguistics.
- Colin B Clement, Matthew Bierbaum, Kevin P O'Keeffe, and Alexander A Alemi. 2019. On the use of arxiv as a dataset. *arXiv preprint arXiv:1905.00075*.
- Joseph Cornelius, Oscar Lithgow-Serrano, Sandra Mitrovic, Ljiljana Dolamic, and Fabio Rinaldi. 2024. BUST: Benchmark for the evaluation of detectors of LLM-generated text. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8029–8057, Mexico City, Mexico. Association for Computational Linguistics.

- 718 719 721
- 725 726 729 731 732 733 736 739 740
- 741 742

- 744 745 747 748 749 750 751
- 755 756
- 758 759 761

- 767
- 768 770

771

- 773
- 774 775

- Debby R. E. Cotton, Peter A. Cotton, and J. Reuben Shipway and. 2024. Chatting and cheating: Ensuring academic integrity in the era of chatgpt. Innovations in Education and Teaching International, 61(2):228-239.
- Yao Dou, Maxwell Forbes, Rik Koncel-Kedziorski, Noah A. Smith, and Yejin Choi. 2022. Is GPT-3 text indistinguishable from human text? scarecrow: A framework for scrutinizing machine text. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7250–7274, Dublin, Ireland. Association for Computational Linguistics.
- Liam Dugan, Alyssa Hwang, Filip Trhlík, Andrew Zhu, Josh Magnus Ludan, Hainiu Xu, Daphne Ippolito, and Chris Callison-Burch. 2024. RAID: A shared benchmark for robust evaluation of machinegenerated text detectors. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 12463-12492, Bangkok, Thailand. Association for Computational Linguistics.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: Long form question answering. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 3558-3567, Florence, Italy. Association for Computational Linguistics.
- Sebastian Gehrmann, Hendrik Strobelt, and Alexander Rush. 2019. GLTR: Statistical detection and visualization of generated text. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pages 111–116, Florence, Italy. Association for Computational Linguistics.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. arXiv preprint arXiv:2407.21783.
- Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. arXiv preprint arXiv:2301.07597.
- Hanxi Guo, Siyuan Cheng, Xiaolong Jin, Zhuo Zhang, Kaiyuan Zhang, Guanhong Tao, Guangyu Shen, and Xiangyu Zhang. 2024a. Biscope: Ai-generated text detection by checking memorization of preceding tokens. In Advances in Neural Information Processing Systems, volume 37, pages 104065–104090. Curran Associates, Inc.
- Xun Guo, Shan Zhang, Yongxin He, Ting Zhang, Wanquan Feng, Haibin Huang, and Chongyang Ma. 2024b. Detective: Detecting ai-generated text via multi-level contrastive learning. In Advances in Neural Information Processing Systems, volume 37, pages 88320-88347. Curran Associates, Inc.

Abhimanyu Hans, Avi Schwarzschild, Valeriia Cherepanova, Hamid Kazemi, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2024. Spotting llms with binoculars: zero-shot detection of machine-generated text. In Proceedings of the 41st International Conference on Machine Learning, ICML'24. JMLR.org.

776

783

784

785

786

787

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

- Xinlei He, Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. 2024. Mgtbench: Benchmarking machine-generated text detection. In Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security, CCS '24, page 2251–2265, New York, NY, USA. Association for Computing Machinery.
- Benjamin D. Horne and Maurício Gruppi. 2024. Nelaps: A dataset of pink slime news articles for the study of local news ecosystems. Proceedings of the International AAAI Conference on Web and Social Media, 18(1):1958-1966.
- Guiyang Hou, Yongliang Shen, and Weiming Lu. 2024. Progressive tuning: Towards generic sentiment abilities for large language models. In Findings of the Association for Computational Linguistics: ACL 2024, pages 14392–14402, Bangkok, Thailand. Association for Computational Linguistics.
- Xiaomeng Hu, Pin-Yu Chen, and Tsung-Yi Ho. 2023. Radar: Robust ai-text detection via adversarial learning. In Advances in Neural Information Processing Systems, volume 36, pages 15077-15095. Curran Associates, Inc.
- Guanhua Huang, Yuchen Zhang, Zhe Li, Yongjian You, Mingze Wang, and Zhouwang Yang. 2024. Are AIgenerated text detectors robust to adversarial perturbations? In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6005-6024, Bangkok, Thailand. Association for Computational Linguistics.
- 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.
- Ehsan Kamalloo, Nouha Dziri, Charles Clarke, and Davood Rafiei. 2023. Evaluating open-domain question answering in the era of large language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5591–5606, Toronto, Canada. Association for Computational Linguistics.
- John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages 17061-17084. PMLR.

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

945

946

947

891

892

834

843 844 845

842

- 847 848 849 850 851 852
- 853 854 855 856 857 858
- 860 861

862 863 864

870

871

- 878
- 879 880
- 881
- 883 884
- 885 886 887
- 8
- 88
- 889

Kristian Kuznetsov, Eduard Tulchinskii, Laida Kushnareva, German Magai, Serguei Barannikov, Sergey Nikolenko, and Irina Piontkovskaya. 2024.
Robust AI-generated text detection by restricted embeddings. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 17036–17055, Miami, Florida, USA. Association for Computational Linguistics.

- Zhixin Lai, Xuesheng Zhang, and Suiyao Chen. 2024. Adaptive ensembles of fine-tuned transformers for Ilm-generated text detection. In 2024 International Joint Conference on Neural Networks (IJCNN), pages 1–7.
- Thomas Lavergne, Tanguy Urvoy, and François Yvon. 2008. Detecting fake content with relative entropy scoring. In *Proceedings of the 2008 International Conference on Uncovering Plagiarism, Authorship and Social Software Misuse - Volume 377*, PAN'08, page 27–31, Aachen, DEU. CEUR-WS.org.
- Yafu Li, Qintong Li, Leyang Cui, Wei Bi, Zhilin Wang, Longyue Wang, Linyi Yang, Shuming Shi, and Yue Zhang. 2024. MAGE: Machine-generated text detection in the wild. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 36–53, Bangkok, Thailand. Association for Computational Linguistics.
- Yepeng Liu and Yuheng Bu. 2024. Adaptive text watermark for large language models. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.
- Zhuang Liu, Wayne Lin, Ya Shi, and Jun Zhao. 2021. A robustly optimized bert pre-training approach with post-training. In *Chinese Computational Linguistics*, pages 471–484, Cham. Springer International Publishing.
- Shixuan Ma and Quan Wang. 2024. Zero-shot detection of LLM-generated text using token cohesiveness.
  In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 17538–17553, Miami, Florida, USA. Association for Computational Linguistics.
- Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, and 1 others. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.
- Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. Detectgpt: zero-shot machine-generated text detection using probability curvature. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.
- Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *Proceedings*

of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 188–197, Hong Kong, China. Association for Computational Linguistics.

- Michael-Andrei Panaitescu-Liess, Zora Che, Bang An, Yuancheng Xu, Pankayaraj Pathmanathan, Souradip Chakraborty, Sicheng Zhu, Tom Goldstein, and Furong Huang. 2025. Can watermarking large language models prevent copyrighted text generation and hide training data? *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(23):25002– 25009.
- Qi Pang, Shengyuan Hu, Wenting Zheng, and Virginia Smith. 2024. No free lunch in llm watermarking: Trade-offs in watermarking design choices. *arXiv preprint arXiv:2402.16187*.
- Shushanta Pudasaini, Luis Miralles, David Lillis, and Marisa Llorens Salvador. 2025. Benchmarking AI text detection: Assessing detectors against new datasets, evasion tactics, and enhanced LLMs. In *Proceedings of the 1stWorkshop on GenAI Content Detection (GenAIDetect)*, pages 68–77, Abu Dhabi, UAE. International Conference on Computational Linguistics.
- Vipula Rawte, Swagata Chakraborty, Agnibh Pathak, Anubhav Sarkar, S.M Towhidul Islam Tonmoy, Aman Chadha, Amit Sheth, and Amitava Das. 2023. The troubling emergence of hallucination in large language models - an extensive definition, quantification, and prescriptive remediations. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pages 2541–2573, Singapore. Association for Computational Linguistics.
- Jinyan Su, Terry Zhuo, Di Wang, and Preslav Nakov. 2023. DetectLLM: Leveraging log rank information for zero-shot detection of machine-generated text. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12395–12412, Singapore. Association for Computational Linguistics.
- Adaku Uchendu, Zeyu Ma, Thai Le, Rui Zhang, and Dongwon Lee. 2021. TURINGBENCH: A benchmark environment for Turing test in the age of neural text generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2001–2016, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ivan Vykopal, Matúš Pikuliak, Ivan Srba, Robert Moro, Dominik Macko, and Maria Bielikova. 2024. Disinformation capabilities of large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 14830–14847, Bangkok, Thailand. Association for Computational Linguistics.
- Jan Philip Wahle, Terry Ruas, Frederic Kirstein, and Bela Gipp. 2022. How large language models are

transforming machine-paraphrase plagiarism. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 952–963, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

952

957

959

960

961

962

963

965

967

969 970

971

972

973

974 975

976

977

978

981

982 983

987

993

995

996

997 998

1000

1001

1002

1003

1004

- Yiming Wang, Zhuosheng Zhang, and Rui Wang. 2023. Element-aware summarization with large language models: Expert-aligned evaluation and chain-ofthought method. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8640–8665, Toronto, Canada. Association for Computational Linguistics.
- Yuxia Wang, Jonibek Mansurov, Petar Ivanov, Jinyan Su, Artem Shelmanov, Akim Tsvigun, Osama Mohammed Afzal, Tarek Mahmoud, Giovanni Puccetti, Thomas Arnold, Alham Aji, Nizar Habash, Iryna Gurevych, and Preslav Nakov. 2024. M4GTbench: Evaluation benchmark for black-box machinegenerated text detection. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3964– 3992, Bangkok, Thailand. Association for Computational Linguistics.
- Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Lidia Sam Chao, and Derek Fai Wong. 2025. A survey on llm-generated text detection: Necessity, methods, and future directions. *Computational Linguistics*, 51(1):275–338.
- Junchao Wu, Runzhe Zhan, Derek F. Wong, Shu Yang, Xinyi Yang, Yulin Yuan, and Lidia S. Chao. 2024.
  Detectrl: Benchmarking llm-generated text detection in real-world scenarios. In Advances in Neural Information Processing Systems, volume 37, pages 100369–100401. Curran Associates, Inc.
- Xianjun Yang, Wei Cheng, Yue Wu, Linda Petzold, William Yang Wang, and Haifeng Chen. 2023. Dnagpt: Divergent n-gram analysis for training-free detection of gpt-generated text. *arXiv preprint arXiv*:2305.17359.
- WJ Youden. 1950. Index for rating diagnostic tests. *Cancer*, 3(1):32–35.
- Sungduk Yu, Man Luo, Avinash Madusu, Vasudev Lal, and Phillip Howard. 2025. Is your paper being reviewed by an llm? a new benchmark dataset and approach for detecting ai text in peer review. *arXiv preprint arXiv*:2502.19614.
- Xiao Yu, Kejiang Chen, Qi Yang, Weiming Zhang, and Nenghai Yu. 2024a. Text fluoroscopy: Detecting LLM-generated text through intrinsic features. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 15838–15846, Miami, Florida, USA. Association for Computational Linguistics.
- Xiao Yu, Yuang Qi, Kejiang Chen, Guoqiang Chen, Xi Yang, Pengyuan Zhu, Xiuwei Shang, Weiming Zhang, and Nenghai Yu. 2024b. Dpic: Decoupling

prompt and intrinsic characteristics for llm generated1005text detection. In Advances in Neural Information1006Processing Systems, volume 37, pages 16194–16212.1007Curran Associates, Inc.1008

Ying Zhou, Ben He, and Le Sun. 2024. Humanizing<br/>machine-generated content: Evading AI-text detec-<br/>tion through adversarial attack. In Proceedings of<br/>the 2024 Joint International Conference on Compu-<br/>tational Linguistics, Language Resources and Eval-<br/>uation (LREC-COLING 2024), pages 8427–8437,<br/>Torino, Italia. ELRA and ICCL.1009<br/>1010

#### A Dataset

1016

1018

1019

1021

1022

1023

1025

1048

1050

1051

1052

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

#### A.1 Data domains

Medium: We utilized the Medium Articles dataset curated by Fabio Chiusano, hosted on Hugging Face<sup>1</sup>. This dataset consists of more than 190k English-language articles from Medium.com, each containing textual and metadata fields, including title, body, URL, authorship, timestamp, and tags. We randomly selected a subset of 12.5k articles with word counts between 400 and 2000 that were published before 2021.

News: We curated a collection of 200k news arti-1027 cles from reputable news agency websites, includ-1028 1029 ing ABC News, Al Jazeera, American Press, Associated Press News, CBS News, CNN, NBC News, 1030 Reuters, and The Guardian. From this corpus, we 1031 randomly sampled 12.5k articles with lengths ranging from 250 to 2000 words. The selected articles 1033 span diverse topical domains such as politics, sci-1034 1035 ence, social issues, religion, technology, sports, and culture. All articles were published prior to the advent of modern LLMs. 1037

Reviews: We utilized the Amazon Reviews dataset 1038 collected by (Ni et al., 2019), which comprises 1039 1040 over 233 million customer reviews spanning from 1996 to October 2018. This large-scale dataset includes rich metadata, such as review text, star 1042 ratings, helpfulness scores, and product attributes 1043 (e.g., category, brand, price, and image features). 1044 For our benchmark, we randomly sampled 12.5k 1045 reviews from various product categories, retaining 1046 only those with 30 or more words. 1047

**Reddit:** To build the Reddit component of our dataset, we extracted question–answer pairs from the ELI5 subreddit, following a collection strategy similar to (Fan et al., 2019). We restricted the dataset to answers with lengths between 400 and 2000 words, all posted before 2021. A final sample of 12.5k answers was randomly selected for inclusion in our benchmark.

**arXiv:** We utilized the arXiv dataset introduced by (Clement et al., 2019), which contains over 1.5 million preprint articles spanning disciplines such as physics, mathematics, and computer science. Each article includes metadata such as title, abstract, authors, categories, and citation data. To construct our dataset, we randomly sampled 12.5k abstracts with word counts ranging from 150 to 500, limited to publications before 2021.

<sup>1</sup>https://huggingface.co/datasets/fabiochiu/ medium-articles **Pink slime:** We employed the NELA-PS dataset 1065 introduced by (Horne and Gruppi, 2024), which 1066 encompasses 7.9 million articles from 1093 local 1067 news sources, commonly referred to as "pink slime" 1068 journalism, spanning from March 2021 to January 1069 2024. These outlets generate content that mimics 1070 legitimate local journalism in structure and pre-1071 sentation while frequently advancing partisan nar-1072 ratives. From this corpus, we sampled 12.5k ar-1073 ticles published in 2021, with document lengths 1074 constrained to 250-2000 words. 1075

1076

1077

1078

1079

1080

1082

1083

1084

1085

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

**Wikipedia** We utilized the Plain Text Wikipedia 2020-11 dataset accessible through Kaggle<sup>2</sup>. This corpus comprises a comprehensive Wikipedia dump of 23 GB, containing articles spanning diverse topics and domains. From this dataset, we randomly sampled 12.5k articles, each with a word count between 400 and 2000, for inclusion in our benchmark.

### A.2 Utilized LLMs and input prompts

### A.2.1 Models

We employed seven distinct open-source large language models and their instruction-tuned variants sourced from the Hugging Face repository: Llama3.2-1b, Llama3.2-3b, Llama3.1-8b, Mistral-7b, Qwen-7b, Gemma2-2b, and Gemma2-9b.

Llama3.1-8b: The Llama 3.1-8b-Instruct model is an instruction-tuned variant of Meta's 8Bparameter LLM from the Llama 3.1 series, optimized for dialogue and assistant-like tasks. It is fine-tuned using supervised fine-tuning (SFT) and reinforcement learning with human feedback (RLHF), enhancing its ability to align with human preferences and generate helpful, and safe responses. The model supports up to 128k token context lengths and employs Grouped-Query Attention (GQA), rotary positional embeddings (RoPE), and SwiGLU activations to improve scalability and inference efficiency. Trained on over 15T tokens of publicly available data, it supports multilingual capabilities across several major languages and demonstrates strong performance in reasoning, text, and code generation tasks.

Llama3.2-1b and Llama3.2-3b: The Llama 3.2-1b-Instruct and Llama 3.2-3b-Instruct models are instruction-tuned variants of Meta's Llama 3.2 series, comprising 1.23 billion and 3.21 billion pa-

<sup>&</sup>lt;sup>2</sup>https://www.kaggle.com/datasets/ltcmdrdata/ plain-text-wikipedia-202011

rameters respectively. Both models are optimized 1113 for multilingual dialogue applications, including 1114 tasks like agentic retrieval and summarization, and 1115 have demonstrated superior performance compared 1116 to many open-source and proprietary chat mod-1117 els on standard industry benchmarks. They em-1118 ploy an auto-regressive transformer architecture 1119 enhanced with GQA for improved inference scala-1120 bility and support a context length of up to 128k to-1121 kens. Trained on up to 9 trillion tokens of publicly 1122 available online data, these models support mul-1123 tiple languages such as English, German, French, 1124 Italian, Portuguese, Hindi, Spanish, and Thai. 1125

Mistral-7b: We utilized Mistral-7b-Instruct-v0.3 1126 model, that is an instruction-tuned variant of Mis-1127 tral AI's 7.25b-parameter language model, de-1128 signed to excel in a wide range of natural lan-1129 guage processing tasks. This version introduces 1130 several enhancements over its predecessors, includ-1131 ing an expanded vocabulary of 32,768 tokens and 1132 support for the v3 tokenizer, which improves its 1133 ability to handle complex text inputs. The v0.3 1134 instruction-tuned variant has been specifically opti-1135 mized through reinforcement learning from RLHF 1136 techniques to better follow user instructions and 1137 generate more helpful responses. Despite its rela-1138 tively compact size compared to models with hun-1139 dreds of billions of parameters, Mistral-7b-Instruct-1140 v0.3 demonstrates competitive performance across 1141 various benchmarks, including reasoning, coding, 1142 and language understanding tasks. 1143

Qwen-7b: The DeepSeek-R1-Distill-Qwen-7b 1144 model is a 7.62b-parameter instruction-tuned lan-1145 guage model developed by DeepSeek AI, based 1146 on the Qwen2.5 architecture. It was fine-tuned 1147 using reinforcement learning and SFT techniques, 1148 leveraging reasoning data generated by the larger 1149 DeepSeek-R1 model. This training approach en-1150 hances the model's capabilities in reasoning, math-1151 ematics, and coding tasks. The model supports a 1152 context length of up to 128k tokens, enabling it to 1153 handle extensive inputs effectively. 1154

Gemma2-2b and Gemma2-9b: The Gemma fam-1155 ily, developed by Google, includes the Gemma2-2b-1156 instruct and Gemma2-9b-instruct models, which 1157 are instruction-tuned variants of the base Gemma 1158 models. These models are part of Google's effort to 1159 1160 provide lightweight, high-performing open models, drawing from the same research and technology 1161 used to create the Gemini models. Both models are 1162 decoder-only large language models, available in 1163 English, and are designed to be versatile for various 1164

applications, including text generation, conversational AI, and summarization.

#### A.2.2 Prompts

We employed LLMs in their chat-based configura-1168 tions to generate AI texts across multiple writing 1169 styles. Each data domain received tailored prompts 1170 to elicit appropriate responses. For the "pink slime" 1171 data, we supplemented prompts with definitional 1172 context to ensure semantic alignment. To preserve 1173 comparability in length, the LLMs were instructed 1174 to produce outputs approximately equal in word 1175 count to the corresponding human-written text (de-1176 noted as <N>). Moreover, in the Amazon Reviews, 1177 Reddit, Pink slime, and Wikipedia datasets, we 1178 incorporated the respective titles, product name, 1179 post title, news headline, or document title, into 1180 the input prompt to enhance coherence and topical 1181 relevance. Below, we detail the specific prompts 1182 used to condition LLM outputs across the various 1183 writing styles examined in this work. 1184

Medium:

[{'role': 'system', 'content': 'You are a blog writer in Medium website. You paraphrase the Medium article I give you in about <N> words as if you are the original author, maintaining the same ideas and tone while using your own words.'}, {'role': 'user', 'content': 'The article is <HUMAN TEXT>'},]

News:

[{'role': 'system', 'content': 'You journalist are а working for а reputable news agency. You paraphrase the news article I give you in about <N> words as if you are the original writer, maintaining the same ideas and tone while using your own words.'}, {'role': 'user', 'content': 'The article is <HUMAN TEXT>"},]

1186

1185

1167

1188

1190	<pre>with title <title> in about <n> words as if you are the original review writer, maintaining the same ideas and tone while using your own words.'}, {'role': 'user', 'content': 'The article is <human text="">"},]</human></n></title></pre>	<pre>with title <title> in about <n> words as if you are the original article writer, maintaining the same ideas and tone while using your own words.'}, {'role': 'user', 'content': 'The article is <human text="">"},]</human></n></title></pre>
1191	Reddit:	A.3 Dataset statistics
1192	<pre>[{'role': 'system', 'content': 'You are a Reddit user. You paraphrase the Reddit post with title <title> I give you in about <n> words as if you are the original Reddit post writer, maintaining the same ideas and tone while using your own words.'}, {'role': 'user', 'content': 'The article is <human text="">"},]</human></n></title></pre>	Table 3 reports the number of text samples included in our benchmark for each LLM evaluated. The "Human" column denotes the number of collected human-written texts. The "Paraphrased" column corresponds to the original LLM-generated outputs without humanification. The "Humanified" column indicates the number of paraphrased samples that were modified using the humanification strategies specified in the corresponding rows. Additionally, Figure 4 illustrates the word count distributions
1193	arXiv:	of human-written texts and LLM-generated para-
1104	<pre>[{'role': 'system', 'content': 'You are a scientific paper writer. You paraphrase the abstract of a scientific paper I give you in about <n> words as if you are the original author, maintaining the same ideas and tone while using your own words.'}, {'role': 'user', 'content': 'The article is <human text="">"},]</human></n></pre>	<ul> <li>B Text samples</li> <li>For masked word prediction within our MLM framework, we employed the Longformer-base-4096 architecture developed by the Allen Institute for AI (Beltagy et al., 2020). In this section, we provide exemplar texts demonstrating the three humanification strategies employed in our method-</li> </ul>
1195	<pre>Pink slime:     [{'role': 'system', 'content': 'Pink     slime journalism is a practice in     which American news outlets, or fake     partisan operations masquerading as     such, publish poor-quality news reports     which appear to be local news. You are     a pink slime journalist. You paraphrase     the pink slime article I give you about     a subject with title <title> in about     <n> words as if you are the original     article writer, maintaining the same     ideas and tone while using your own</n></title></pre>	<ul> <li>ology. The words enclosed in brackets and high-lighted in red represent words predicted by the MLM architecture, which subsequently replaced their corresponding precedent words in the paraphrased text. Text highlighted in blue represents the baseline AI-paraphrased content prior to humanification processing.</li> <li><b>B.1 Sample humanified text from RMM</b> strategy</li> <li>The White House is escalating its efforts to persuade Congress to approve limited [military] strikes [action] against Syria, as President Obama faces a formidable challenge in convincing lawmakers to</li> </ul>
1196 1197	<pre>words.'}, {'role': 'user', 'content': 'The article is <human text="">"},] Wikipedia:</human></pre>	back a new military campaign in the Middle East. The president has been personally engaging with skeptical lawmakers over the weekend, delivering [making] a tailored pitch to Democrats and Repub- licans who remain undecided or open [close] to

[{'role': 'system', 'content': 'You

are a Wikipedia writer. You paraphrase

the article I give you about a subject

[{'role': 'system', 'content': 'You

are an Amazon customer. You paraphrase

the review I give you about a product

Table 3: Data statistics	s for each	LLM	utilized.
--------------------------	------------	-----	-----------

Dataset	Medium				Nows			Amazon revi	owe	Reddit			
Strategy	Human	Paraphrased	Humanified	Human	Paranhrased	Humanified	Human	Paraphrased	Humanified	Human	Paraphrased	Humanified	
Baranhrasing	2000	2000	Trainanned	2000	3000	Tramanned	2000	3000	Trantannea	2000	2000	Tumanned	
rarapinasing	3000	5000		3000	5000		3000	3000		3000	3000		
RMM @ p=10%	500	500	500	500	500	500	500	500	500	500	500	500	
RMM @ p=20%	500	500	500	500	500	500	500	500	500	500	500	500	
RMM @ p=40%	500	500	500	500	500	500	500	500	500	500	500	500	
RMM @ p=60%	500	500	500	500	500	500	500	500	500	500	500	500	
RMM @ p=80%	500	500	500	500	500	500	500	500	500	500	500	500	
RMM @ p=100%	500	500	500	500	500	500	500	500	500	500	500	500	
AWS @ $p{=}10\%$	500	500	500	500	500	500	500	500	500	500	500	500	
AWS @ p=20%	500	500	500	500	500	500	500	500	500	500	500	500	
AWS @ p=40%	500	500	500	500	500	500	500	500	500	500	500	500	
AWS @ p=60%	500	500	500	500	500	500	500	500	500	500	500	500	
AWS @ p=80%	500	500	500	500	500	500	500	500	500	500	500	500	
AWS @ p=100%	500	500	500	500	500	500	500	500	500	500	500	500	
RHL @ $R{=}5\%$	500	500	500	500	500	500	500	500	500	500	500	500	
RHL @ R=15%	500	500	500	500	500	500	500	500	500	500	500	500	
RHL @ R=20%	500	500	500	500	500	500	500	500	500	500	500	500	
RHL @ R=25%	500	500	500	500	500	500	500	500	500	500	500	500	
RHL @ R=30%	500	500	500	500	500	500	500	500	500	500	500	500	
RHL @ R=35%	500	500	500	500	500	500	500	500	500	500	500	500	
RHL @ R=40%	500	500	500	500	500	500	500	500	500	500	500	500	
Dataset $\rightarrow$		arXiv			Pink slime			Wikipedia	1				
				1									
Strategy $\downarrow$	Human	Paraphrased	Humanified	Human	Paraphrased	Humanified	Human	Paraphrased	Humanified				
Strategy ↓ Paraphrasing	Human 3000	Paraphrased 3000	Humanified	Human 3000	Paraphrased 3000	Humanified	Human 3000	Paraphrased 3000	Humanified				
Strategy $\downarrow$ ParaphrasingRMM @ $p=10\%$	Human 3000 500	Paraphrased 3000 500	Humanified  500	Human 3000 500	Paraphrased 3000 500	Humanified — 500	Human 3000 500	Paraphrased 3000 500	Humanified — 50				
$\begin{tabular}{c} Strategy \downarrow \\ \hline Paraphrasing \\ \hline RMM @ p{=}10\% \\ RMM @ p{=}20\% \end{tabular}$	Human 3000 500 500	Paraphrased 3000 500 500	Humanified — 500 500	Human 3000 500 500	Paraphrased 3000 500 500	Humanified — 500 500	Human 3000 500 500	Paraphrased 3000 500 500	Humanified — 50 500				
Strategy ↓           Paraphrasing           RMM @ p=10%           RMM @ p=20%           RMM @ p=40%	Human 3000 500 500 500	Paraphrased 3000 500 500 500	Humanified — 500 500 500	Human 3000 500 500 500	Paraphrased 3000 500 500 500	Humanified — 500 500 500	Human 3000 500 500 500	Paraphrased 3000 500 500 500	Humanified 				
Strategy ↓           Paraphrasing           RMM @ p=10%           RMM @ p=20%           RMM @ p=40%           RMM @ p=60%	Human 3000 500 500 500 500	Paraphrased 3000 500 500 500 500	Humanified 	Human 3000 500 500 500 500	Paraphrased 3000 500 500 500 500	Humanified 	Human 3000 500 500 500 500	Paraphrased 3000 500 500 500 500 500	Humanified 				
Strategy ↓           Paraphrasing           RMM @ p=10%           RMM @ p=20%           RMM @ p=40%           RMM @ p=60%           RMM @ p=80%	Human 3000 500 500 500 500 500	Paraphrased 3000 500 500 500 500 500	Humanified 	Human 3000 500 500 500 500 500	Paraphrased 3000 500 500 500 500 500	Humanified 	Human 3000 500 500 500 500 500	Paraphrased 3000 500 500 500 500 500	Humanified 				
$\begin{tabular}{lllllllllllllllllllllllllllllllllll$	Human 3000 500 500 500 500 500 500	Paraphrased 3000 500 500 500 500 500 500	Humanified 	Human           3000           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500	Humanified 	Human 3000 500 500 500 500 500 500	Paraphrased 3000 500 500 500 500 500 500 500	Humanified 				
$\begin{tabular}{lllllllllllllllllllllllllllllllllll$	Human 3000 500 500 500 500 500 500 500	Paraphrased 3000 500 500 500 500 500 500 500	Humanified 	Human           3000           500           500           500           500           500           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500 500	Humanified 	Human 3000 500 500 500 500 500 500 500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
$\begin{tabular}{lllllllllllllllllllllllllllllllllll$	Human           3000           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500           500           500           500           500           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500 500 500	Humanified — 500 500 500 500 500 500 500	Human           3000           500           500           500           500           500           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
Strategy ↓           Paraphrasing           RMM @ p=10%           RMM @ p=20%           RMM @ p=40%           RMM @ p=60%           RMM @ p=100%           AWS @ p=20%           AWS @ p=20%           AWS @ p=40%	Human           3000           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 500 500 500 500 500 500 500	Human           3000           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
Strategy ↓           Paraphrasing           RMM @ $p=10\%$ RMM @ $p=20\%$ RMM @ $p=60\%$ RMM @ $p=60\%$ RMM @ $p=100\%$ AWS @ $p=10\%$ AWS @ $p=20\%$ AWS @ $p=20\%$ AWS @ $p=0\%$ AWS @ $p=0\%$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 500 500 500 500 500 500 500	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 500 500 500 500 500 500 500	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 50 500 500 500 50 50 50 500 500 500 500 500				
$\begin{tabular}{lllllllllllllllllllllllllllllllllll$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
$\begin{array}{c} {\rm Strategy} \downarrow \\ \hline {\rm Paraphrasing} \\ {\rm RMM} @ p = 10\% \\ {\rm RMM} @ p = 20\% \\ {\rm RMM} @ p = 40\% \\ {\rm RMM} @ p = 60\% \\ {\rm RMM} @ p = 80\% \\ {\rm RMM} @ p = 80\% \\ {\rm AWS} @ p = 10\% \\ {\rm AWS} @ p = 10\% \\ {\rm AWS} @ p = 40\% \\ {\rm AWS} @ p = 40\% \\ {\rm AWS} @ p = 80\% \\ {\rm AWS} @ p = 100\% \end{array}$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
$\begin{array}{c} {\rm Strategy} \downarrow \\ \hline {\rm Paraphrasing} \\ {\rm RMM} @ p = 10\% \\ {\rm RMM} @ p = 20\% \\ {\rm RMM} @ p = 40\% \\ {\rm RMM} @ p = 60\% \\ {\rm RMM} @ p = 80\% \\ {\rm RMM} @ p = 80\% \\ {\rm AWS} @ p = 20\% \\ {\rm AWS} @ p = 20\% \\ {\rm AWS} @ p = 40\% \\ {\rm AWS} @ p = 60\% \\ {\rm AWS} @ p = 80\% \\ {\rm AWS} @ p = 80\% \\ {\rm AWS} @ p = 80\% \\ {\rm AWS} @ p = 100\% \\ {\rm RHL} @ R = 5\% \end{array}$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 500 500 500 500 500 500 500	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 50 500 500 500 500 500 500 5				
$\begin{tabular}{lllllllllllllllllllllllllllllllllll$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 500 500 500 500 500 500 500	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 500 500 500 500 500 500 500	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 50 500 500 500 500 500 500 5				
Strategy ↓           Paraphrasing           RMM @ $p=10\%$ RMM @ $p=20\%$ RMM @ $p=60\%$ RMM @ $p=60\%$ RMM @ $p=100\%$ AWS @ $p=10\%$ AWS @ $p=20\%$ AWS @ $p=20\%$ AWS @ $p=60\%$ AWS @ $p=80\%$ AWS @ $p=100\%$ RHL @ $R=15\%$ RHL @ $R=15\%$ RHL @ $R=20\%$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
Strategy ↓           Paraphrasing           RMM @ $p=10\%$ RMM @ $p=20\%$ RMM @ $p=60\%$ RMM @ $p=60\%$ RMM @ $p=100\%$ AWS @ $p=20\%$ AWS @ $p=60\%$ AWS @ $p=50\%$ AWS @ $p=50\%$ RHL @ $R=5\%$ RHL @ $R=25\%$ RHL @ $R=20\%$ RHL @ $R=25\%$ RHL @ $R=25\%$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
$\begin{array}{c} \textbf{Strategy} \downarrow \\ \hline Paraphrasing \\ \hline RMM @ p=10\% \\ RMM @ p=20\% \\ RMM @ p=40\% \\ RMM @ p=60\% \\ RMM @ p=80\% \\ RMM @ p=100\% \\ \hline AWS @ p=20\% \\ AWS @ p=20\% \\ AWS @ p=40\% \\ AWS @ p=60\% \\ AWS @ p=80\% \\ AWS @ p=100\% \\ \hline RHL @ R=5\% \\ RHL @ R=15\% \\ RHL @ R=25\% \\ RHL @ R=25\% \\ RHL @ R=25\% \\ RHL @ R=30\% \\ \end{array}$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 500 500 500 500 500 500 500	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 				
$\begin{array}{l} \textbf{Strategy} \downarrow \\ \hline Paraphrasing \\ \hline RMM @ p=10\% \\ RMM @ p=20\% \\ RMM @ p=40\% \\ RMM @ p=60\% \\ RMM @ p=80\% \\ RMM @ p=100\% \\ \hline AWS @ p=20\% \\ AWS @ p=20\% \\ AWS @ p=40\% \\ AWS @ p=60\% \\ AWS @ p=60\% \\ AWS @ p=100\% \\ \hline RHL @ R=5\% \\ RHL @ R=15\% \\ RHL @ R=25\% \\ RHL @ R=25\% \\ RHL @ R=35\% \\ \end{array}$	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified 	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified	Human           3000           500	Paraphrased 3000 500 500 500 500 500 500 50	Humanified — 50 500 500 500 500 500 500 5				

reconsidering their opposition to military action.

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

According to White House officials, Obama's argument [case] centers on the dual imperative of both moral responsibility and national security interest. The president believes that the United States has a critical obligation to respond to the devastating chemical weapons attack [attacks] in Syria, which has left countless civilians, including children, dead or injured. This perspective is underscored by a series of disturbing videos, obtained by ABC News, which were shown to lawmakers in a classified briefing [session] last week.

These graphic images, which depict the harrowing aftermath of the chemical attack [attacks], are being used by the administration to make a powerful [compelling] case to Congress and the American public. Secretary of State John Kerry, who has been leading the charge to build international support [consensus] for military action, referenced the videos in a speech in Paris on Saturday, emphasizing the atrocities committed by the Syrian [Assad] regime against its own people [citizens].

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

Kerry's impassioned plea, which highlighted the tragic fate of innocent civilians, including children, is a stark reminder [illustration] of the human cost of inaction. As he noted, the use of chemical weapons in the middle of the night, when people [children] should have been sleeping safely in their beds, is an unconscionable act that demands [warrants] a response from the international community [body].

The administration is aware that it faces a tough sell in convincing Congress to approve military action, with a recent ABC News survey indicating deep opposition among lawmakers. However, officials remain hopeful that they can build a coalition of 60 [enough] Democrats and Republicans to overcome the threat of a filibuster and secure approval for the military strike.

In a bid to build support [consensus], Vice President Biden is hosting a dinner for over a dozen Republican senators on Sunday night, with the guest

list including several key lawmakers who have ex-1280 pressed reservations about military action. The administration is also counting on the support [en-1282 dorsement] of retired General David Petraeus, a 1283 respected figure in military circles, who has publicly urged lawmakers to back the military strike. 1285

1281

1286

1287

1289

1291

1292

1293

1294

1295

1297

1298

1299

1300

1301

1302

1303

1305

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1321

1322

1324 1325

1326

1327

1328

1329

As the vote on the military strike looms in the Senate as early as Wednesday, the White House is stepping [speeding] up its efforts to persuade lawmakers to support [back] the president's plan. Obama will make his case to the American public in a series of interviews with network television anchors [stations] on Monday, followed by a televised address on Tuesday. The outcome of this high-stakes debate will have far-reaching [widereaching] implications for the United States and its role in the Middle East, and the White House is leaving [left] no stone unturned in its bid to secure approval for military action.

#### Sample humanified text from AWS **B.2** strategy

Pleased with the affordable [good] price and adorable [good] design[ [look], I've been using this cardigan for work over the past few months. However [but], I've noticed minor [some] issues [things] where my arms have caused small fabric balls to form, which can be a bit bothersome but not entirely unexpected [bad] for a budget-friendly [work-type] item. Sizing has been an issue [problem] for me, as I often [have] find [that] cardigans in my size still gape open. If you're larger-chested [big-chested] like me, I'd suggest sizing up for a better fit. Despite [Besides] these concerns [things], I might [have] still purchase [get] more due [because] to the low cost, but a higher-quality [highercut] garment might [do] offer [get] better longevity [looks].

#### **B.3** Sample humanified text from RHL strategy

I'm feeling [really] so hurt [upset] and confused after our last conversation. I said I was sorry and wanted to make amends [up], but he's now acting like I'm the one who's done him wrong. I don't see how I deserve [want] to be treated [loved] this way, especially [just] when I've been trying [able] to be understanding [kind] and accommodating [loving]. Our relationship was pretty [so] intense [good] from the start, and we were monogamous pretty [very] quickly. However [unfortunately], he refused [wanted] to acknowledge [believe] that we were in a real relationship for a few months, which 1330 was confusing and frustrating [annoying]. He had 1331 high expectations [opinion] of me, and I felt like 1332 I had to constantly prove [explain] myself to him. 1333 I've never [only] cheated on anyone before, and 1334 I thought I was able to handle [explain] myself 1335 around [in] people with good intentions [reason]. 1336 One of the things that really bothered me was how 1337 he would belittle [tell] me about my male friends. 1338 He would get [be] upset if I spent time with them 1339 or made plans with them, and he would even get 1340 [be] angry if I didn't include them in everything. 1341 I felt like he was conditioning me into isolation, 1342 and it was suffocating [annoying]. There were a lot 1343 [couple] of fights in our relationship, and many of 1344 them were petty [stupid]. Like the time he thought 1345 I smelled [was] of cologne [stupid] and swore [that] 1346 I had just slept with someone [her]. Or the time he 1347 came over and claimed [said] there were tire marks 1348 [nuts] in his parking [sweet] spot, even though [if] 1349 it was just a dirt [back] road. It was like he was 1350 looking for any excuse to get [be] upset. I also [really] felt like he was playing games with me. Like 1352 the time I told him I was lazy [sick] and didn't want 1353 to drive [go] to my friends' house, and he got upset 1354 because I didn't spend the night [time] with him. 1355 Or the time he broke up with me and then said we 1356 weren't in a relationship, even though [if] we had 1357 been hanging out and talking. It was like he was 1358 trying [able] to manipulate [trick] me into feeling 1359 [being] guilty [stupid] or responsible [sorry] for 1360 his own emotions [faults]. But the thing that really 1361 hurt [upset] me was when he introduced me to his 1362 ex-girlfriend [girl-friends] as his friend, and then 1363 acted like I was the one who had done him wrong. 1364 He said I was "fucking with his heart and emotions 1365 [mind]," and that I needed to earn [lose] his trust. 1366 It was like he was trying [out] to make me feel 1367 [look] like I was the problem [one], even though 1368 [if] I had done nothing wrong. I feel [felt] like 1369 I'm losing [wasting] my mind, to be honest [fair]. 1370 I'm trying [attempting] to be understanding [kind] 1371 and accommodating [sweet], but it feels [was] like 1372 he's not giving me any space or trust. I'm starting 1373 [beginning] to wonder [think] if I'm just not good 1374 enough for him, or if he's just not willing [able] 1375 to work through our issues [shit] together. Do I 1376 deserve [have] to be treated [left] this way? 1377

## 1378 C Detectors

**Binoculars** is a zero-shot, model-agnostic method 1379 for detecting machine-generated text that requires 1380 no training data. It operates by comparing the per-1381 plexity of a text as evaluated by two language mod-1382 els: an "observer" model and a "performer" model. The observer computes the perplexity of the text 1384 1385 directly, while the performer generates next-token predictions, which are then evaluated by the ob-1386 server to compute cross-perplexity. The ratio of 1387 perplexity to cross-perplexity serves as a strong 1388 indicator of whether the text is human- or machine-1389 generated.

1391Fast-DetectGPT is a zero-shot method, building1392upon the principles of DetectGPT. It introduces1393the concept of conditional probability curvature to1394distinguish between human- and AI-authored con-1395tent. The method operates by sampling alternative1396word choices for a given text and evaluating the1397conditional probabilities using a language model.1398By analyzing the curvature of these probabilities,1399Fast-DetectGPT identifies AI text.

LRR (Log-Likelihood Log-Rank Ratio) is a zero-1400 shot approach. It combines two statistical mea-1401 sures: the log-likelihood, which assesses the abso-1402 lute confidence of a language model in predicting 1403 a sequence, and the log-rank, which evaluates the 1404 relative ranking of the predicted tokens. By com-1405 puting the ratio of these two measures, LRR cap-1406 1407 tures nuanced differences between human-written and LLM-generated text. 1408

Log-Likelihood calculates the log-probability of 1409 each token in a text sequence using a language 1410 model, assessing how predictable each word is 1411 within its context. In this framework, human-1412 written text typically exhibits a mix of high- and 1413 low-probability tokens, reflecting natural linguistic 1414 variability. In contrast, LLM-generated text often 1415 contains a higher proportion of high-probability 1416 tokens, indicating more predictable word choices. 1417 **Rank** evaluates the predictability of each token 1418 in a text by determining its rank within the lan-1419 guage model's probability distribution. Tokens that 1420 consistently appear among the top-ranked predic-1421 tions indicate higher predictability, a characteristic 1422 often associated with LLM-generated text. Ana-1423 lyzing the distribution of these token ranks assists 1424 in distinguishing between human-authored and AI-1425 generated content. 1426

Log-Rank enhances the performance of Rankmethod by applying a logarithmic transformation

to the rank of each token within a language model's predicted probability distribution.

1429

1430

**RADAR** is a framework designed to enhance the 1431 detection of AI-generated text, particularly against 1432 paraphrased content that often evades traditional 1433 detectors. It employs an adversarial training ap-1434 proach involving two components: a paraphraser 1435 and a detector. The paraphraser aims to rewrite AI-1436 generated text to resemble human-authored content, 1437 thereby challenging the detector's ability to identify 1438 machine-generated text. Conversely, the detector is 1439 trained to distinguish between human-written and 1440 AI-generated texts. 1441

## **D** Additional results

Figure 5 shows the radar charts of comparing detec-1443tors for each LLM across all writing styles. Table 41444shows the impact of different strategies on detector1445performance across all writing styles for smaller1446models from each family evaluated in this study.1447



Figure 4: Word count distribution of human-written and LLM-generated texts aggregated across all LLMs. "All" represents all samples across both writing styles and LLMs.



Figure 5: Radar charts of comparing detectors in different generative models across all writing styles.

$LLM \rightarrow$	Gemma2 2b					Llam	a3.2 1b		Llama3.2 3b				
Detector $\downarrow$ Metric (%) $\rightarrow$	AUC	W-A	SFD	URSS	AUC	W-A	SFD	URSS	AUC	W-A	SFD	URSS	
			P	araphra	se (base	eline)							
Binoculars	97.3	84.0	57.6	48.4	85.8	33.1	47.3	15.7	93.6	63.4	64.0	40.6	
Fast-DetectGPT	96.8	87.6	60.6	53.0	86.4	64.7	58.2	37.7	93.3	75.4	56.2	42.4	
Log-Likelihood	86.6	40.1	45.2	18.1	75.9	33.9	51.8	17.6	78.5	30.8	42.4	13.0	
Log-Rank	88.1	46.1	63.0	29.0	76.7	36.9	58.9	21.7	79.1	33.3	41.3	13.7	
LRR	89.5	61.7	54.3	33.5	77.0	42.9	58.2	24.9	79.1	41.0	54.7	22.5	
Radar	83.1	49.6	16.1	8.0	74.2	36.6	14.7	5.4	78.7	38.9	17.3	6.7	
Rank	79.4	42.7	58.8	25.1	69.1	29.2	62.8	18.3	73.6	32.7	56.8	18.6	
Random meaning-preserving mutation (RMM)													
Binoculars	89.6	58.6	40.3	23.6	80.2	29.2	50.9	14.8	82.8	39.8	50.1	19.9	
Fast-DetectGPT	88.9	63.2	50.8	32.1	80.1	46.9	57.8	27.1	81.8	45.8	49.2	22.5	
Log-Likelihood	48.0	10.5	8.5	0.9	42.5	11.7	11.8	1.4	38.2	6.8	6.0	0.4	
Log-Rank	52.2	12.9	9.3	1.2	45.7	13.3	18.4	2.5	41.0	7.8	7.6	0.6	
LRR	67.6	25.5	14.5	3.7	59.6	20.0	13.2	2.6	55.3	13.8	10.7	1.5	
Radar	91.7	63.1	49.3	31.1	86.2	50.9	28.9	14.7	88.8	53.1	30.9	16.4	
Rank	55.7	11.5	13.2	1.5	50.2	8.5	9.6	0.8	48.7	7.4	9.7	0.7	
			AI-fl	agged wo	o <mark>rd swa</mark>	p (AWS	5)						
Binoculars	88.4	59.5	49.7	29.6	76.0	28.1	32.2	9.1	81.7	43.5	36.5	15.9	
Fast-DetectGPT	89.0	62.3	45.8	28.5	78.2	43.2	49.7	21.4	83.1	47.9	43.0	20.6	
Log-Likelihood	29.5	5.9	6.0	0.3	25.7	7.0	13.1	0.9	22.2	3.4	4.8	0.2	
Log-Rank	31.9	6.6	5.5	0.4	26.9	7.3	16.2	1.2	23.5	3.6	4.2	0.2	
LRR	45.1	10.9	7.2	0.8	34.5	7.9	5.2	0.4	33.5	5.2	4.1	0.2	
Radar	89.0	56.6	35.5	20.1	79.8	38.1	24.4	9.3	83.5	43.1	25.7	11.1	
Rank	50.5	7.0	12.5	0.9	42.4	4.9	7.0	0.3	44.6	4.6	10.3	0.5	
Recursive humanization loop (RHL)													
Binoculars	86.8	51.1	52.5	26.8	73.2	23.5	51.9	12.2	78.3	35.1	40.9	14.4	
Fast-DetectGPT	86.8	55.1	50.6	27.9	75.0	34.3	53.5	18.4	78.9	37.7	45.2	17.0	
Log-Likelihood	28.3	3.5	4.3	0.2	22.9	4.6	12.5	0.6	20.0	1.7	3.3	0.1	
Log-Rank	31.9	4.4	4.7	0.2	25.3	5.1	12.5	0.6	22.7	2.1	3.5	0.1	
LRR	49.8	10.1	11.4	1.2	38.9	7.2	6.8	0.5	38.4	5.1	6.9	0.3	
Radar	89.0	54.3	41.5	22.5	80.4	37.5	27.2	10.2	83.1	40.5	26.3	10.7	
Rank	50.4	6.7	13.8	0.9	41.4	3.8	6.6	0.2	43.8	4.0	7.7	0.3	

Table 4: Performance of detectors under paraphrasing (baseline), RMM, AWS, and RHL strategies.