EvoBench: Towards Real-world LLM-Generated Text Detection Benchmarking for Evolving Large Language Models

Anonymous ACL submission

Abstract

With the widespread of Large Language Models (LLMs), there has been an increasing need to detect LLM-generated texts, prompting extensive research in this area. However, existing detection methods mainly evaluate on static benchmarks, which neglect the evolving nature of LLMs. Relying on existing static benchmarks could create a misleading sense of security, overestimating the real-world effectiveness of detection methods. To bridge this gap, we introduce EvoBench, a dynamic benchmark considering a new dimension of generalization across continuously evolving LLMs. EvoBench categorizes the evolving LLMs into (1) updates over time and (2) developments like finetuning and pruning, covering 7 LLM families and their 30 evolving versions. To measure the generalization across evolving LLMs, we introduce a new EMG (Evolving Model Generalization) metric. Our evaluation of 14 detection methods on EvoBench reveals that they all struggle to maintain generalization when confronted with evolving LLMs. To mitigate the generalization problems, we further propose improvement strategies, demonstrating EMG performance improvements up to 12.15%. Our research sheds light on critical challenges in real-world LLM-generated text detection and represents a significant step toward practical applications.¹

1 Introduction

002

006

007

011

017

027

031

Large Language Models (LLMs), such as Chat-GPT (OpenAI, 2022b), Claude (Anthropic, 2024), and LLaMA (Touvron et al., 2023a), have demonstrated remarkable capabilities in natural language understanding and task processing, leading to their widespread application (M Alshater, 2022; Yuan et al., 2022; Christian, 2023). However, concerns have been raised about the misuse of these models



Figure 1: Current benchmarks (Guo et al., 2023; Bao et al.; Wang et al., 2024b,a; Kwan et al., 2024) primarily focus on specific versions of LLMs and neglect vast ecosystem of evolving LLMs beneath the surface, including updates over time or developments through fine-tuning or pruning. The figure primarily includes LLMs before Jan 2025.

in areas like social media (Ahmed et al., 2021), education (Lee et al., 2023), and academic writing (Mitchell, 2022; Patrick Wood, 2023). For instance, LLMs can be used to manipulate the public by generating comments or to fabricate experimental data and statistical results in support of unverified hypotheses (Solaiman et al., 2019; Goldstein et al., 2023). The potential misuse of LLMs highlights the urgent need to detect LLM-generated text (Kaur et al., 2022; Chen and Shu, 2023).

The academic community has carried out extensive research to detect LLM-generated text effectively (Liu et al., 2019; Gehrmann et al., 2019; Su et al., 2023; Solaiman et al., 2019). Current methods can be categorized into supervised methods (Hu et al., 2023; Yu et al., 2024; Chen et al., 2024b) and zero-shot methods (Ippolito et al., 2020; Yang et al.; Mitchell et al., 2023; Su et al., 2023; Bao et al.). Supervised methods are typically trained with a binary classifier to distinguish be-

¹The Evobench is now available at:

https://anonymous.4open.science/r/EvoBench.



Figure 2: Detection accuracy (measured in AUROC) of Fast-DetectGPT when faced with the evolving LLMs. Figure (a) shows a clear decline in average detection performance as the LLM updates. Figure (b) illustrates the developments of LLMs, including fine-tuning and pruning.

tween texts generated by LLMs and those created by humans, while zero-shot methods primarily rely on statistical features gathered from pre-trained large language models. Although these methods show strong performance on existing benchmarks (Bao et al.; Hu et al., 2023; Chen et al., 2024b; Yu et al., 2024), they often fall short in realworld applications as current benchmarks neglect the evolving nature of LLMs (Wang et al., 2024b; Guo et al., 2023; Wang et al., 2024a; He et al., 2024; Macko et al., 2023). As illustrated in Figure 1, existing benchmarks include a limited versions of LLMs, much like observers seeing only the tip of an iceberg while overlooking the vast evolving LLMs hidden beneath the surface.

061

062

073

084

092

In real-world applications, LLMs are continuously evolving via regular updates, fine tuning (Touvron et al., 2023b; Zhang et al., 2023), or pruning (Sun et al.; Liang et al., 2021), all of which affect LLMs' output (Tao et al., 2024; Touvron et al., 2023a; Gunasekar et al., 2023), thereby impacting the detection performance. Figure 2 illustrates that the detection performance of current detection methods suffers up to 25% drop as LLMs evolve.² Therefore, relying on current static benchmarks for evaluation could create a misleading sense of security, leading the research community to overestimate the real-world effectiveness of detection methods. Therefore, it is crucial to enhance existing benchmarks to better capture the ongoing evolving LLMs.

In this paper, we propose EvoBench, a dynamic benchmark that extends traditional benchmarking to account for the evolving LLMs. EvoBench aims to provide a more accurate evaluation of detection methods by incorporating two dimensions of evolving LLMs: **updates** and **developments**, as shown in Figure 1. Updates refer to changes made by the LLM publishers, including LLM updates, *i.e.*, GPT-40 undergoes updates approximately every 3 months ³. On the other hand, developments refer to optimizations made by LLM developers for specific application scenarios, such as fine-tuning and pruning. EvoBench introduces a new dimension of generalization, enabling the evaluation of detection methods across evolving LLMs, covering 7 widely used LLMs and their 30 evolving versions.

097

100

101

103

104

105

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

Using EvoBench, we evaluate 14 widely used detection methods and find that all struggle to adapt to the evolving nature of LLMs. To intuitive quantify the generalization across evolving LLMs, we introduce the EMG (Evolving Model Generalization) metric. Results reveal that the performance of 14 current detection methods, including widely used Fast-DetectGPT (Bao et al.) and RADAR (Hu et al., 2023) detectors, significantly drop when faced with evolving versions, as reflected by low EMG values, while some achieve high AUROC (0.91) scores on individual models⁴. This phenomenon highlights the limitations of current detection methods in generalization capabilities.

To mitigate this problem, we further explore two strategies. For zero-shot methods like Fast-DetectGPT (Bao et al.), we propose to prune the scoring model to extract shared features for detection, which could lead to a 12.15% improvement in EMG performance across developments of the LLM. For supervised methods, improving the quality and distribution of training data - by contin-

²The widely used detection method, Fast-DetectGPT, declines up to 25% drop when detecting the Claude-3-5-haiku-20241022 compared to Claude-3-haiku-20240307.

³These versions are gpt-4o-2024-05-13, gpt-4o-2024-08-06, gpt-4o-2024-11-20, and chatgpt-4o-latest.

⁴Fast-DetectGPT achieve 0.91 AUROC when detecting Claude-Haiku-2024-03-07, but drop to 0.65 AUROC when detecting the version of 2024-10-22.



Figure 3: The construction pipeline and the dataset statistics of EvoBench. EvoBench includes three dimensions of generalization: dataset domains, generation strategies, and evolving LLMs. The evolving LLMs are further divided into two categories: updates and developments.

uously incorporating data from evolving LLMs could increase 0.24% EMG performance across updates of the GPT-4. However, ensuring effective generalization across both updates and developments simultaneously remains a significant challenge. This finding provides valuable insight into the complexities of adapting detection methods to the continuously evolving LLMs.

Overall, EvoBench serves as a complementary benchmark that focuses on the evolving nature of LLMs, offering a more accurate framework for evaluating the real-world applicability of detection methods. While it does not claim to be a one-sizefits-all solution, we envision EvoBench evolving alongside detection methods and other benchmarks, helping guide future developments in this area. In particular, we suggest that extracting shared features from evolving LLMs could enhance the generalization of detection methods. With EvoBench, the community can more effectively evaluate and refine detection methods, contributing to more robust and practical detection systems that mitigate LLM misuse in real-world applications.

2 EvoBench

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

152

153

154

155

156

158

159

160

2.1 Definition of Evolving LLMs

EvoBench primarily considers two evolving pathways for LLMs: (1) **updates**, which are released by publishers, and (2) **developments**, which involve optimization made by developers. These two evolving pathways present significant challenges to current detection methods.

Updates. We focus on two types of updates that
occur over time: model (Brown et al., 2020; Tao
et al., 2024) and version updates (Brown et al.,
2020). Model updates typically involve significant
changes to the model's architecture. For example,
the transition from GPT-4 to GPT-40 represents a

major architectural shift (Achiam et al., 2023). In contrast, version updates occur more frequently, with shorter cycles and less perception, such as the regular releases of new versions of GPT-4 every 2-3 months, such as updates in May, August, and November 2024. 167

168

169

170

171

172

173

174

175

176

177

178

179

181

182

183

184

185

186

187

188

189

190

191

192

193

194

196

197

198

199

200

201

202

203

204

Developments. In EvoBench, model developments include fine-tuning (Hu et al., 2022; Mangrulkar et al., 2022), pruning (Liang et al., 2021), often used for domain adaptation and optimizing efficiency for real-world applications. Unlike updates, these actions are nearly imperceptible to the detector (Bhattacharjee et al., 2023). For example, developers may fine-tune LLMs using private datasets, making the detection more challenging.

2.2 Dataset Collection

2.2.1 Data Collection Principle

EvoBench, incorporating the two evolving pathways of LLMs, extends existing benchmarks (Mitchell et al., 2023) widely used in current detection methods, including DetectGPT (Mitchell et al., 2023), Fast-DetectGPT (Bao et al.), and ImBD (Chen et al., 2024b) to assess three key dimensions of generalization: (1) evolving LLMs, (2) domain, and (3) generation paradigms. Figure 3 illustrates our data collection pipeline to cover these three key dimensions.

Generalization to LLM Evolution. To comprehensively assess evolving, we consider 7 widely used LLM families: GPT-40, GPT-4, Claude, Gemini, Qwen, LlaMA3, and LlaMA2, encompassing a total of 30 evolving versions. For updates, we focus on families of GPT-40, GPT-4, Claude, Gemini, Qwen, and LlaMA3, tracking how detection methods perform across updates of these models. For developments, we focus on the LlaMA2 family, incorporating fine-tuning and pruning techniques. Detailed specifications of all LLMs are provided in

255

252

256 257 258

260 261

259

262

267

265

266

269

270

271

272

273

274

276

277

278

279

280

281

283

284

285

287

289

290

292

the Appendix C.

205

216

217

218

221

225

226

227

229

230

234

235

241

Generalization to Domains. EvoBench in-206 cludes five datasets spanning diverse tasks: 207 XSum (Narayan et al., 2018) for (1)news articles, WritingPrompts (Fan et al., 2018) for (2) creative 209 story writing, and PubMed (Jin et al., 2019) for (3) biomedical research question answering. Additionally, we also include two datasets: SocialMe-212 213 dia (Kula and Gregor, 2024) for (4) datasets on social media and PeerRead (Kang et al., 2018) for 214 (5) peer-reviewed academic writing. 215

Generalization to Generation Paradigms. EvoBench incorporates three distinct generation paradigms: (1) continuation-based generation, includes XSum, Writing, and PeerRead; (2) question-answering generation, including *PubMed*; (3) paraphrased generation, including SocialMedia. By integrating these diverse paradigms, EvoBench offers a more comprehensive framework to assess the robustness of detection methods in various text generation paradigms.

2.2.2 Dataset Collection Process

We detail our systematic dataset construction process following the three-dimensional generalization framework. For each dataset, we first carefully selected 150 human-authored texts as reference samples. To collect AI-generated texts, we employed consistent prompting strategies across all evaluated models. Specifically, for continuationbased generation tasks (Xsum, Writing, PeerRead), we provided the prefix of human-written texts as context. For PubMed, we maintained the original question-answer format, while for SocialMedia, we implemented a paraphrasing approach where models were tasked with preserving semantic meaning while using different words. Details of prompts are shown in Appendix E.

2.3 **Dataset Statistics**

EvoBench consists of 7 LLM families, encom-243 passing a total of 30 evolving versions. For each evolving LLM, we collect samples across 5 distinct 245 datasets, with each dataset containing a balanced 247 distribution of 150 human-authored texts and 150 machine-generated texts. In total, EvoBench com-248 prises $30 \times 5 \times 150 = 22500$ machine-generated samples, resulting in a comprehensive evaluation benchmark of 22, 500 pairs of text samples. 251

2.4 **Evaluation Metrics**

To quantify the generalization ability of detection methods to evolving LLMs, we introduce the EMG Φ_E (Evolving Model Generalization) metric, inspired by the coefficient of variation. It evaluates the consistency and trend of detection performance across evolving LLMs.

For an LLM family with *m* evolving versions, the widely evaluated version is selected as the base model, and the remaining m-1 versions as evolving models. The performance (measured by AU-ROC Φ_A) change $\Delta \Phi_A^i$ between the *i*-th models and the base model is :

$$\Delta \Phi^i_A = \Phi^i_A - \Phi^{base}_A, \tag{1}$$

where Φ_A^i and Φ_A^{base} is the Φ_A of the *i*-th evolving model and *base* model, respectively.

The average performance change μ_E of $\Delta \Phi_A$, representing the overall performance change, is:

$$\mu_E = \frac{1}{m-1} \sum_{i=1}^{m-1} \left(\Delta \Phi_A^i \right).$$
 (2)

The volatility σ_E of $\Delta \Phi_A$ using the standard deviation is:

$$\sigma_E = \sqrt{\frac{1}{m-1} \sum_{i=1}^{m-1} \left(\Delta \Phi_A^i - \mu_E \right)^2}.$$
 (3)

The proposed Φ_E is defined as follows:

$$\Phi_E = \frac{\mu_E}{\sigma_E + \lambda} \times \gamma, \tag{4}$$

where λ is a regularization term to prevent fluctuations, and γ is a scaling factor. In this study, we set λ and γ as 1 and 100, respectively.

2.5 Detection Methods

We evaluate a range of supervised and zero-shot detection methods to assess their generalization ability against evolving LLMs. Additionally, we conducted repeated generation processes to ensure that output diversity does not significantly impact detection performance in Appendix B.

2.5.1 Existing Methods

For supervised detectors, we tested the GPT-2 detectors developed by OpenAI (Liu et al., 2019), the RADAR detector (Hu et al., 2023), Text Fluoroscopy (Yu et al., 2024), and ImBD (Imitate Before Detect) (Chen et al., 2024b). For zeroshot detectors, we included Likelihood (average

335

log probabilities) (Gehrmann et al., 2019), LRR (a hybrid method combining log probability and log-rank) (Su et al., 2023), LogRank (mean of the log ranks sorted in descending order of probabilities) (Solaiman et al., 2019), Entropy (average token-level entropy from the predictive distribution)(Ippolito et al., 2020), DNA-GPT (Yang et al.), DetectGPT (Mitchell et al., 2023), and its enhanced variants NPR (Su et al., 2023) and Fast-DetectGPT (Bao et al.).

2.5.2 Optimizing Strategies

294

295

307

310

312

313

314

317

318

319

321

322

325

326

327

330

331

334

Enhancing Supervised Detectors via optimizing Training Data. To enhance detection generalization across evolving models, we propose an incremental training approach: For a evolving LLM, we generate additional training data D_{evolving} and update the training set D_{old} :

$$D_{\text{new}} = D_{\text{old}} \cup D_{\text{evolving}}.$$
 (5)

Using the new dataset D_{new} , the detector f_{detector} is retrained to adapt to the evolving LLM.

Enhancing Zero-shot Detector via Extracting Shared Features. To enhance the detection generalization of the zero-shot method, we propose 315 to prune the scoring model $p_{\theta}^{\text{pruned}}$ to extract the shared features from the developments of LLMs. Based on Fast-DetectGPT, we define our detection metric $d(x, p_{\theta}^{\text{pruned}})$ based on conditional probability curvature:

$$d(x, p_{\theta}^{\text{pruned}}) = \frac{\log p_{\theta}^{\text{pruned}}(x|x) - \tilde{\mu}}{\tilde{\sigma}}, \quad (6)$$

where $\tilde{\mu}$ represents the expected log probability under the pruned model:

$$\tilde{\mu} = \mathbb{E}_{\tilde{x} \sim p_{\theta}^{\text{pruned}}(\tilde{x}|x)}[\log p_{\theta}^{\text{pruned}}(\tilde{x}|x)], \quad (7)$$

and $\tilde{\sigma}^2$ captures the variance of these log probabilities:

$$\tilde{\sigma}^2 = \mathbb{E}_{\tilde{x} \sim p_{\theta}^{\text{pruned}}(\tilde{x}|x)} \left[\left(\log p_{\theta}^{\text{pruned}}(\tilde{x}|x) - \tilde{\mu} \right)^2 \right].$$
(8)

Exprimental Results 3

3.1 Setup

Detection Methods. To mimic real-world scenarios, we set up a black-box environment in which the detector is assumed to be unaware of the source model of the text to be detected under evaluation. We tested 14 detection methods, including

5 supervised detectors and 9 zero-shot detectors. Among them, RADAR (Hu et al., 2023) and Fast-DetectGPT (Bao et al.) are well-known methods in their respective categories. Details of experimental settings can be found in the Appendix D and Appendix F.

3.2 Main Results

To begin with, we select two leading detection methods, Fast-DetectGPT and RADAR detector, and evaluate them on our EvoBench, the results are shown in Table 1. The evaluation results of other detection methods are shown in Appendix H.

First, a trend is that texts generated by more advanced LLMs are harder to detect in the same LLM family. For instance, comparing Claude-3-Haiku⁵ and Claude-3-Opus⁶, the latter demonstrates more advanced capabilities. When detecting these two LLMs, Fast-DetectGPT and RADAR showed decreased detection AUROC of 3.58% and 2.77%, respectively, as shown in Table 1. Similarly, comparing detecting SocialMedia datasets generated by GPT-4o-mini with GPT-4o, two detectors showed a decline in detection AUROC of 3.22% and 6.69%, respectively.

Our results also reveal an intriguing pattern: existing detection methods often over-optimize for specific LLM versions, rather than maintaining robust generalization capabilities across different evolving LLMs. Specifically, detectors with high AUROC on widely used LLM families often show performance degradation with newer versions of the same family. For example, Fast-DetectGPT performed excellently on the GPT-40-05-13 version, achieving an average detection accuracy of 0.8003 on EvoBench, but dropped to 0.7422 when facing the GPT-4o-latest version. Similar patterns were observed in other widely used LLM families, such as GPT-4 and Claude-Sonnet. In contrast, when evaluating less widely used LLMs like LlaMA3, detectors tend to demonstrate better generalization capabilities across evolving LLMs. For example, the detection performance of the RADAR remained stable across different versions of the LlaMA3 family, with AUROC scores across five datasets ranging from 0.7989 to 0.8348.

This suggests that the essence of generalization lies in the overfitting problem. Current detectors are often optimized to achieve superior performance on existing benchmarks, which leads to

⁵claude-3-haiku-20240307

⁶claude-3-opus-20240229

LIMs	Version Time/	Fast-DetectGPT								RA	DAR		
LLIVIS	Version Name	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
					١	Updates							
	2024-05-13	0.90	0.97	0.73	0.58	0.83	0.80	0.99	0.83	0.83	0.60	0.68	0.79
	2024-08-06	0.87	0.94	0.70	0.59	0.77	0.77	0.99	0.86	0.82	0.62	0.69	0.80
-	2024-08-00	(-2.75%)	(-2.41%)	(-3.54%)	(+ 1.76%)	(-5.89%)	(-2.57%)	(-0.11%)	(+3.27%)	(-0.53%)	(+2.12%)	(+ 0.75%)	(+1.10%)
GPI-40	2024-11-20	0.72	0.90	0.70	0.57	0.78	0.73	0.99	0.72	(10000)	0.62	0.66	0.75
		(-17.0%)	0.91	0.70	0.58	(-4.23%) 0.82	(-0.34%) 0.74	(-0.87%)	0.75	(- 4.08%)	(+2.02%)	(-1.47%)	0.77
	Latest	(-20.0%)	(-6.13%)	(-3.36%)	(+0.64%)	(-0.14%)	(-5.81%)	(- 0.79%)	(-7.49%)	(-2.83%)	(+4.46%)	(-0.53%)	(-1.44%)
GPT-40-mini	2024-07-18	0.91	0.97	0.73	0.61	0.80	0.80	1.00	0.87	0.82	0.66	0.64	0.80
011-40-111111	2024=07=18	(+ 0.91%)	(+0.12%)	(-0.42%)	(+3.02%)	(-2.14%)	(+0.30%)	(+0.33%)	(+3.82%)	(-0.68%)	(+ 6.69%)	(-3.96%)	(+1.24%)
	2023-06-13	0.93	0.98	0.76	0.49	0.75	0.78	0.99	0.87	0.89	0.66	0.75	0.83
	2023-11-06	0.92	0.93	0.73	0.55	0.83	0.79	0.99	0.79	0.80	0.66	0.63	0.77
GPT-4		(- 0.85 %)	(-4.76%)	(-2.94%)	(+5.43%)	(+7.71%)	(+0.92%)	(-0.01%)	(-8.49%)	(-8.81%)	(+0.24%)	(-12.3%)	(-5.88%)
011-4	2024-01-25	(-3.22%)	(-3.41%)	(-5.75%)	$(\pm 4.23\%)$	$(\pm 7.34\%)$	(-0.16%)	(-0.23%)	(-7.23%)	(-5.09%)	(-0.46%)	(-8.68%)	(-4.34%)
		0.83	0.91	0.72	0.58	0.79	0.77	0.99	0.79	0.82	0.60	0.60	0.76
	2024-04-09	(-9.30%)	(-6.56%)	(-4.39%)	(+ 8.27%)	(+4.08%)	(-1.58%)	(- 0.36%)	(-8.45%)	(-7.07%)	(-5.92%)	(-14.8%)	(-7.32%)
	2024-02-29	0.95	0.98	0.81	0.58	0.93	0.85	0.97	0.86	0.84	0.74	0.61	0.80
	2024 06 20	0.97	0.99	0.77	0.66	0.90	0.86	0.99	0.81	0.82	0.72	0.55	0.78
Claude-Sonnet	2024-00-20	(+2.33%)	(+0.80%)	(-4.32%)	(+ 7.69%)	(-2.90%)	(+ 0.72%)	(+1.61%)	(-5.21%)	(-2.44%)	(-1.05%)	(-5.57%)	(-2.53%)
	2024-10-22	0.90	0.93	0.68	0.51	0.78	0.76	0.95	0.72	0.74	0.67	0.55	0.73
		(-4.43%)	(-5.39%)	(-13.9%)	(-7.80%)	(-15.3%)	(-9.37%)	(-1.52%)	(-13.4%)	(-9.85%)	(-0./4%)	(-0.38%)	(-1.59%)
a 1 a 1	2024-03-07	1.00	1.00	0.86	0.75	0.94	0.91	1.00	0.93	0.84	0.77	0.67	0.84
Claude-Haiku	2024-10-22	(15.5%)	(756%)	(21.7%)	(35.9%)	(47.4%)	0.05	(0.07%)	0.80	(10.0%)	(6.38%)	(114.7%)	(2.04%)
		0.97	0.96	0.82	0.72	0.89	0.87	0.99	0.87	0.82	0.77	0.62	0.81
Claude-Opus	2024-02-29	(- 2.49%)	(-3.89%)	(-3.80%)	(-3.27%)	(-4.46%)	(-3.58%)	(- 0.33%)	(-6.01%)	(-2.23%)	(+ 0.29%)	(-5.57%)	(-2.77%)
	Qwen1.5-7B	0.92	0.99	0.66	0.49	0.91	0.80	0.93	0.90	0.72	0.51	0.48	0.71
	Owen2 7P	0.99	1.00	0.74	0.46	0.89	0.82	0.89	0.87	0.82	0.52	0.54	0.73
Qwen	Qweii2*/B	(+6.86%)	(+ 0.29%)	(+7.35%)	(-2.75%)	(-1.37%)	(+2.08%)	(-4.13%)	(-3.42%)	(+ 9.46%)	(+1.35%)	(+6.14%)	(+1.88%)
	Qwen2.5-7B	0.99	1.00	0.79	0.47	0.90	0.83	0.90	0.91	0.79	0.52	0.53	0.73
	*****	(+0.51%)	(+0.10%)	(+12.3%)	(-2.11%)	(-0.78%)	(+3.22%)	(-2.00%)	(+0.00%)	(+0.85%)	(+1.31%)	(+4.95%)	(+2.23%)
	LlaMA-3.1-8B	0.99	0.98	0.85	0.69	0.97	0.90	1.00	0.89	0.78	0.72	0.63	0.80
	LlaMA-3.1-70B	$(\pm 0.31\%)$	(+1.50%)	(-3.72%)	(-1.73%)	(-0.96%)	(-0.92%)	(-0.19%)	(-1.32%)	(-1.57%)	(+1.29%)	(-1.22%)	(-0.60%)
LlaMA3		0.96	0.99	0.92	0.86	0.96	0.94	0.96	0.94	0.82	0.75	0.71	0.83
	LIaMA-3.2-1B	(-3.46%)	(+1.09%)	(+ 7.29%)	(+17.5%)	(-0.30%)	(+4.43%)	(-3.51%)	(+4.39%)	(+3.74%)	(+ 2.80%)	(+7.52%)	(+ 2.99%)
	LlaMA-3.2-3B	0.99	0.99	0.91	0.77	0.96	0.93	0.99	0.92	0.82	0.77	0.66	0.83
		(-0.13%)	(+0.31%)	(+6.56%)	(+8.32%)	(-0.30%)	(+2.95%)	(-0.44%)	(+2.71%)	(+3.70%)	(+4.23%)	(+2.80%)	(+2.60%)
	LlaMA-3.3-70B	$(\pm 0.32\%)$	(+1.61%)	(-5.53%)	(+0.15%)	(-6.38%)	(-1.97%)	(-0.32%)	(-0.74%)	(+ 0.94 %)	(+0.75)	(+2.29%)	(+0.59%)
		((()=====)	(,	Day	alonmante	((((1.00.000)	(()))	() //	(1.1.1.1)
	LI-MA 2 7D shot he	0.07	0.00	0.01	0.00	0.07	0.05	0.69	0.(2	0.40	0.64	0.64	0.(2
	LIaMA-2-/B-chat-hf	0.97	0.99	0.91	0.90	0.97	0.95	0.68	0.63	0.49	0.64	0.64	0.62
Fine-tuning	Vicuna-1.5-7B	(+2.36%)	(+0.52%)	(-10.4%)	(-1.58%)	(-2.54%)	(-2.34%)	(+29.3%)	(+29.7%)	(+31.5%)	(+ 0.94 %)	(-0.25%)	(+18.2%)
	Wenneth 7D	0.92	0.98	0.74	0.85	0.84	0.87	0.75	0.76	0.70	0.69	0.54	0.69
	wizardmatn-/B	(-5.17%)	(-1.00%)	(-17.0%)	(-4.73%)	(-12.8%)	(-8.14%)	(+ 6.68%)	(+12.5%)	(+20.5%)	(+4.61%)	(-9.33%)	(+7.01%)
	Sheared-LlaMA1.3B	0.65	0.89	0.48	0.86	0.71	0.72	0.77	0.66	0.59	0.73	0.54	0.66
		(-31.7%)	(-10.1%)	(-42.0%)	(-3.29%)	(-26.2%)	(-22.7%)	(+8.98%)	(+ 2.96%)	(+ 9.45%)	(+ 8.59%)	(-10.2%)	(+3.96%)
Pruning	Sheared-LlaMA1.3B-pruned	(-56.0%)	(-18.2%)	(-55.3%)	(-14.3%)	(-51.0%)	(-38.9%)	(+24.0%)	(+15.7%)	(+20.8%)	(+22.7%)	$(\pm 1.59\%)$	(+16.9%)
	Channel LI-MAD 7D	0.45	0.87	0.37	0.75	0.52	0.59	0.84	0.76	0.59	0.78	0.55	0.70
	Sneared-LIaMA2./B-pruned	(-52.4%)	(-12.3%)	(-53.9%)	(-14.2%)	(-44.9%)	(-35.5%)	(+15.7%)	(+12.8%)	(+ 9.59%)	(+14.3%)	(-9.25%)	(+8.66%)
	Sheared-LlaMA2.7B	0.70	0.92	0.50	0.82	0.65	0.72	0.75	0.68	0.58	0.67	0.47	0.63
		(-26.7%)	(-7.50%)	(-40.0%)	(-7.72%)	(-31.7%)	(-22.7%)	(+ 6.88%)	(+4.81%)	(+ 8.83%)	(+3.47%)	(-17.2%)	(+1.35%)

Table 1: The detection performance (measured in AUROC) of two leading detection methods, Fast-DetectGPT and RADAR, on EvoBench. For each LLM family, we set the LLMs that are widely evaluated in text generation detection benchmarks as anchor points. 'Latest' refers to the LLM currently used in the web version of GPT-40.

overfitting to a specific LLM version rather than being truly effective. This also reveals the vulnerability of current detection methods across evolving LLMs. From a practical perspective, a key priority for future research lies in developing more adaptable detection frameworks that can simultaneously preserve detection performance while demonstrating robust generalization across dataset domains and evolving LLMs.

3.2.1 Detection Generalization across Evolving LLMs

To visually demonstrate the generalization ability of detection methods across evolving LLMs, we evaluate the 14 detection methods using the EMG metric, with results shown in Table 2.

Table 2 demonstrates that no method can

maintain stable performance across all evolving LLMs. Specifically, detectors like ImBD and Text-Fluoroscopy perform well on complex models like GPT-40 and GPT-4. Meanwhile, detectors like entropy and RADAR exhibit better generalizations of newer model families, including Qwen, LlaMA3, and the developments of LlaMA2. This variation in performance across different detectors suggests that improving generalization might require adopting strategies that combine the strengths of different approaches, which may help improve the adaptability of the detector and reduce sensitivity to specific LLM evolution trends. 400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

3.2.2 Performance of Optimizing Strategies

To improve the generalization of current detection methods across evolving LLMs, we have pre-

384

395 396

397

Methods			Upd	ates			Develop	Avσ	
methods	GPT-40	GPT-4	Claude-Sonnet	Claude-haiku	Qwen	LlaMA3	Fine-tuning	Pruning	
			Supe	rvised Detectors					
RoBERTa-base	-2.252	-4.593	-8.325	-15.95	+5.862	+1.457	-2.378	+6.871	-2.413
RoBERTa-large	-1.878	-5.577	-10.06	-14.63	+4.646	+3.448	-0.725	+12.54	-1.530
RADAR	-0.537	-5.775	-4.934	-2.394	+2.053	+1.372	+11.96	+7.302	+1.131
Text-Fluoroscopy	-0.493	-0.655	+0.902	-0.123	-5.241	+1.895	-11.03	-15.77	-3.815
Imitate Before Detect	Imitate Before Detect $+0.440 + 2.397 - 1.572$				+0.023	-0.054	+8.942	-9.842	-0.513
Likelihood	-0.873	-0.047	-3.810	-8.280	-1.174	+0.439	-6.991	-32.99	-6.717
Rank	-2.127	+1.071	-0.178	-12.09	+4.249	+1.023	-9.661	-7.109	-3.103
LogRank	-1.316	-0.106	-4.272	-9.443	-0.608	+0.837	-6.851	-30.14	-6.488
Entropy	-3.202	+0.416	+0.602	-3.411	+4.138	+2.145	+5.150	+25.12	+3.869
LRR	-2.511	-0.075	-5.864	-14.26	+2.629	+1.911	-6.964	-17.45	-5.325
NPR	-2.846	+7.436	-2.237	-15.26	+7.132	+0.866	-8.063	-31.18	-5.520
DNA-GPT	-3.805	+1.249	-4.621	-12.26	-1.332	+1.337	-6.355	-36.11	-7.738
DetectGPT	-4.095	+10.04 +2.258		-11.09	+2.357 -0.33		-5.442	-20.18	-3.312
Fast-DetectGPT	tectGPT -3.558 -0.271 -4.115			-13.16	+2.632	+1.094	-5.095	-27.94	-6.302

Table 2: EMG performance of 14 detection methods on EvoBench. Red indicates a negative EMG, signifying a decrease in AUROC when facing evolving LLMs, while a larger EMG value reflects a better generalization of the detection method.



Figure 4: EMG performance of two optimizing strategies compared with their corresponding original methods. The left panel illustrates the improvement for zero-shot detection using the pruned scoring model, while the right panel shows the optimization for supervised detection using the newer dataset.

liminarily explored two possible approaches, with results presented in Figure 4. For the zero-shot detection method, we compared the EMG performance of the Fast-DetectGPT (Bao et al.) detector with different scoring models, including GPT-J, LlaMA2, and the pruned version of LlaMA2. The results show that our proposed method outperforms by 12.152% in EMG when dealing with developments, including fine-tuning and pruning. However, for fine-tuned developments, using the pruned LlaMA2 sacrifices some performance compared to directly using LlaMA2.

416

417

418

419

420

421

499

423

424

425

426

427

For the supervised detection method, we chose to improve Text-Fluoroscopy (Yu et al., 2024). The original training dataset was generated by GPT-3.5turbo. We enhanced it by regenerating it with the initial version of the corresponding LLM family and retrained the detector. The details are shown in the Appendix D. We observed improvements in performance for GPT-4 and Claude-Sonnet but a decline for GPT-40. This suggests that while dataset enhancement can improve performance in some LLMs, it may not always yield better results when faced with different evolving LLMs.

433

434

435

436

437

438

439

440

441

4 Related Work

4.1 Evolving LLMs

Updates of Pre-trained LLMs. The advent of442LLMs, such as ChatGPT (OpenAI, 2022a) and443GPT-4 (OpenAI, 2023), has marked a paradigm444shift in text generation. However, these pre-trained445models are not static (Tao et al., 2024; Touvron446et al., 2023a; Gunasekar et al., 2023; Biderman447et al., 2023); rather, they undergo continuous evo-448

lution, with models like ChatGPT frequently updating their parameters to enhance performance and
adapt to new application scenarios (Zheng et al.,
2024; Touvron et al., 2023a).

Developments on Pre-trained LLMs. The rise of 453 open-source LLMs has facilitated their widespread 454 usage in diverse domains. Developers can flex-455 ibly tailor LLMs to specific tasks through fine-456 tuning (Touvron et al., 2023a,b; Zhang et al., 2023), 457 pruning (Sun et al.; Liang et al., 2021; Xia et al.), 458 and quantization (Du et al., 2024; Chen et al., 459 2024a; Dettmers et al., 2023) techniques that mod-460 ify the model's parameters or structure, ultimately 461 influencing output characteristics. In this pro-462 cess, developers could use domain-specific datasets 463 to enhance the model's understanding of special-464 ized fields and their contextual knowledge (Kerner, 465 2024). Furthermore, depending on practical ap-466 plication scenarios, users may select models with 467 varying parameter scales to strike an optimal bal-468 ance between performance, speed, and resource 469 consumption (Kim et al., 2024; Zhao et al., 2023). 470

4.2 LLM-generated Text Detection

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

Supervised Methods. Supervised methods (Liu et al., 2019; Chen et al., 2024b) are usually trained to differentiate between texts generated by LLMs and texts created by humans. For example, RADAR (Hu et al., 2023) introduces the idea of adversarial learning to train a detector that can resist paraphrase attacks. Text Fluoroscopy (Yu et al., 2024) extracts discriminative features from the intermediate layers of the language model and utilizes them to train a binary classifier, which enhances the generalization of supervised detectors across texts from different semantic domains. However, as LLMs continue to evolve, frequent adaptation mechanisms-including updates, finetuning, pruning, and other optimization strategies-introduce changes in their generated text, making it difficult to guarantee the effectiveness of existing detection methods on newer model versions. This creates a significant gap between academic detection models and real-world applications, ultimately limiting the practical utility of these methods.

494Zero-shot Methods. Existing zero-shot methods495primarily rely on statistical features extracted us-496ing pre-trained large language models (Bao et al.;497Mitchell et al., 2023). These features include like-498lihood (Gehrmann et al., 2019), probability curva-499ture (Mitchell et al., 2023), divergence between

multiple completions of a truncated passage (Yang et al.), and conditional probability curvature (Bao et al.). Zero-shot detection methods are immune to domain-specific degradation, demonstrating superior generalization in detection tasks (Gehrmann et al., 2019; Mitchell et al., 2023). However, the effectiveness of zero-shot detection heavily depends on the alignment between the pre-trained LLM used for detection and the generation model that produced the text to be detected (Bao et al.; Mitchell et al., 2023). While current methods achieve state-of-the-art results on existing benchmarks (Guo et al., 2023; Bao et al.; Wang et al., 2024b,a; Kwan et al., 2024), the evolving nature of LLMs introduces significant shifts in this alignment. As models continue to evolve, the pre-trained LLMs used for detection may become increasingly misaligned with the updated generation models, leading to a decline in detection performance. This explains why many methods, despite ranking highly on leaderboards, often fail when deployed in real-world scenarios.

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

5 Conclusion

In this paper, we introduce EvoBench, a novel benchmark for evaluating the generalization of LLM-generated text detection methods across evolving LLMs. EvoBench defines two key dimensions of LLM evolution: (1) updates made by LLM publishers over time and (2) developments carried out by developers, ensuring a comprehensive understanding of how evolving LLMs impact detection performance. To quantify this evolving dimension generalization, we propose the EMG (Evolving Model Generalization) metric.

We evaluate 14 widely used detection methods using EvoBench, revealing their vulnerabilities when facing evolving LLMs. In response, to improve the generalization of zero-shot methods across developing LLMs, we propose two strategies: for zero-shot methods, we suggest pruning the scoring model to extract shared features across LLM developments; for supervised methods, we recommend augmenting training data with data from LLM updates. Our benchmark represents a key step towards bringing detection methods into real-world scenarios for evolving LLMs. We also envision continuously updating this benchmark to cover a broader range of evolving LLMs, enabling more comprehensive evaluations across various domains and evolving LLMs.

550 Limitations

566

567

569

570

572

573

574

575

582

583

584

589

591

592

593

595

596

597

598

599

Although we have thoroughly considered the three key dimensions of generalization toward real-world 552 scenarios, certain limitations remain that warrant 553 further exploration in future research. First, it is 554 important to emphasize that EvoBench is not in-555 tended to replace existing benchmarks for evaluating the generalization of LLM-generated text 557 detection methods but rather to complement them. Therefore, EvoBench primarily focuses on three key dimensions in real-world scenarios: domains, 560 generation strategies, and evolving models. We 561 have not covered other aspects of generalization, such as language. A potential direction for future research is to extend EvoBench to incorporate additional dimensions of generalization. 565

> Second, robustness has not been addressed in this study, as our focus was mainly on detection generalization. However, EvoBench could serve as a foundation for inspiring current robustness research, as robustness also involves tasks such as rewriting text using different LLMs.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Alim Al Ayub Ahmed, Ayman Aljabouh, Praveen Kumar Donepudi, and Myung Suh Choi. 2021. Detecting fake news using machine learning: A systematic literature review. *arXiv preprint arXiv:2102.04458*.
 - AI Anthropic. 2024. Claude 3.5 sonnet model card addendum. *Claude-3.5 Model Card*, 3(6).
- Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. Fast-detectgpt: Efficient zeroshot detection of machine-generated text via conditional probability curvature. In *The Twelfth International Conference on Learning Representations*.
- Amrita Bhattacharjee, Tharindu Kumarage, Raha Moraffah, and Huan Liu. 2023. ConDA: Contrastive domain adaptation for AI-generated text detection. In Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 598–610, Nusa Dua, Bali. Association for Computational Linguistics.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit,

USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR. 601

602

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

- Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Canyu Chen and Kai Shu. 2023. Combating misinformation in the age of llms: Opportunities and challenges. *arXiv preprint arXiv: 2311.05656*.
- Hong Chen, Chengtao Lv, Liang Ding, Haotong Qin, Xiabin Zhou, Yifu Ding, Xuebo Liu, Min Zhang, Jinyang Guo, Xianglong Liu, and Dacheng Tao. 2024a. DB-LLM: Accurate dual-binarization for efficient LLMs. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 8719–8730, Bangkok, Thailand. Association for Computational Linguistics.
- Jiaqi Chen, Xiaoye Zhu, Tianyang Liu, Ying Chen, Xinhui Chen, Yiwen Yuan, Chak Tou Leong, Zuchao Li, Tang Long, Lei Zhang, et al. 2024b. Imitate before detect: Aligning machine stylistic preference for machine-revised text detection. *arXiv preprint arXiv:2412.10432*.
- Jon Christian. 2023. Cnet secretly used ai on articles that didn't disclose that fact, staff say. *Futurusm*, *January*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115.
- DaYou Du, Yijia Zhang, Shijie Cao, Jiaqi Guo, Ting Cao, Xiaowen Chu, and Ningyi Xu. 2024. BitDistiller: Unleashing the potential of sub-4-bit LLMs via self-distillation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 102–116, Bangkok, Thailand. Association for Computational Linguistics.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 889–898.
- Sebastian Gehrmann, Hendrik Strobelt, and Alexander M Rush. 2019. Gltr: Statistical detection and visualization of generated text. In *Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics (ACL).

764

765

766

768

769

Josh A Goldstein, Girish Sastry, Micah Musser, Renee DiResta, Matthew Gentzel, and Katerina Sedova. 2023. Generative language models and automated influence operations: Emerging threats and potential mitigations. *arXiv preprint arXiv:2301.04246*.

658

662

668

671

673

674

675

678

688

690

700

701

703

704

707

708

709

710

711

712

- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*.
- 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*.
- 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*, pages 2251–2265.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Xiaomeng Hu, Pin-Yu Chen, and Tsung-Yi Ho. 2023. Radar: Robust ai-text detection via adversarial learning. Advances in Neural Information Processing Systems, 36:15077–15095.
- Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. 2020. Automatic detection of generated text is easiest when humans are fooled. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1808–1822.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. 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 2567–2577.
- Dongyeop Kang, Waleed Ammar, Bhavana Dalvi, Madeleine van Zuylen, Sebastian Kohlmeier, Eduard Hovy, and Roy Schwartz. 2018. A dataset of peer reviews (PeerRead): Collection, insights and NLP applications. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1647–1661, New Orleans, Louisiana. Association for Computational Linguistics.
- Davinder Kaur, Suleyman Uslu, Kaley J Rittichier, and Arjan Durresi. 2022. Trustworthy artificial intelligence: a review. ACM Computing Surveys (CSUR), 55(2):1–38.

- Tobias Kerner. 2024. Domain-specific pretraining of language models: A comparative study in the medical field. *arXiv preprint arXiv:2407.14076*.
- Sanghoon Kim, Dahyun Kim, Chanjun Park, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, Changbae Ahn, Seonghoon Yang, Sukyung Lee, Hyunbyung Park, Gyoungjin Gim, Mikyoung Cha, Hwalsuk Lee, and Sunghun Kim. 2024. SOLAR 10.7B: Scaling large language models with simple yet effective depth upscaling. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 6: Industry Track), pages 23–35, Mexico City, Mexico. Association for Computational Linguistics.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Sebastian Kula and Michal Gregor. 2024. Multilingual models for check-worthy social media posts detection. *arXiv preprint arXiv:2408.06737*.
- Wai-Chung Kwan, Xingshan Zeng, Yufei Wang, Yusen Sun, Liangyou Li, Yuxin Jiang, Lifeng Shang, Qun Liu, and Kam-Fai Wong. 2024. M4LE: A multiability multi-range multi-task multi-domain longcontext evaluation benchmark for large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15568–15592, Bangkok, Thailand. Association for Computational Linguistics.
- Jooyoung Lee, Thai Le, Jinghui Chen, and Dongwon Lee. 2023. Do language models plagiarize? In *Proceedings of the ACM Web Conference 2023*, pages 3637–3647.
- Tailin Liang, John Glossner, Lei Wang, Shaobo Shi, and Xiaotong Zhang. 2021. Pruning and quantization for deep neural network acceleration: A survey. *Neurocomputing*, 461:370–403.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Muneer M Alshater. 2022. Exploring the role of artificial intelligence in enhancing academic performance: A case study of chatgpt. *Available at SSRN*.
- Dominik Macko, Robert Moro, Adaku Uchendu, Jason Lucas, Michiharu Yamashita, Matúš Pikuliak, Ivan Srba, Thai Le, Dongwon Lee, Jakub Simko, et al. 2023. Multitude: Large-scale multilingual machinegenerated text detection benchmark. In *Proceedings* of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 9960–9987.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin

876

877

878

879

880

Bossan. 2022. Peft: State-of-the-art parameterefficient fine-tuning methods. https://github. com/huggingface/peft.

770

774

775

778

789

790

791

793

807

810

812

813

814

815

816

817

818

819

820

822

823

- Alex Mitchell. 2022. Professor catches student cheating with chatgptl: 'i feel abject terror'. Retrieved from https://nypost.com/2022/12/26/ students-using-chatgpt-to-cheat-professor\ -warns/. Accessed on February 17, 2023.
- Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. 2023. Detectgpt: Zero-shot machine-generated text detection using probability curvature. In *International Conference on Machine Learning*, pages 24950–24962. PMLR.
- Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807.
- OpenAI. 2022a. ChatGPT. https://chat.openai. com/.
- OpenAI. 2022b. Introducing chatgpt.
 - OpenAI. 2023. GPT-4 Technical Report. arXiv preprint arXiv:2303.08774.
- Mary Louise Kelly Patrick Wood. 2023. 'everybody is cheating': Why this teacher has adopted an open chatgpt policy. Retrieved from https: //www.npr.org/2023/01/26/1151499213/ chatgpt-ai-education-cheating-classroom\ -wharton-school. Accessed on January 26, 2023.
 - Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. 2019. Release strategies and the social impacts of language models. *arXiv preprint arXiv:1908.09203*.
 - 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.
 - Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. In *The Twelfth International Conference on Learning Representations*.
 - Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. 2024. A survey on self-evolution of large language models. *arXiv preprint arXiv*:2404.14387.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal

Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/ mesh-transformer-jax.
- 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. 2024a. 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.
- Yuxia Wang, Jonibek Mansurov, Petar Ivanov, Jinyan Su, Artem Shelmanov, Akim Tsvigun, Chenxi Whitehouse, Osama Mohammed Afzal, Tarek Mahmoud, Toru Sasaki, et al. 2024b. M4: Multi-generator, multi-domain, and multi-lingual black-box machinegenerated text detection. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1369–1407.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language model pre-training via structured pruning. In *The Twelfth International Conference on Learning Representations*.
- Xianjun Yang, Wei Cheng, Yue Wu, Linda Ruth Petzold, William Yang Wang, and Haifeng Chen. Dna-gpt: Divergent n-gram analysis for training-free detection of gpt-generated text. In *The Twelfth International Conference on Learning Representations*.
- Xiao Yu, Kejiang Chen, Qi Yang, Weiming Zhang, and Nenghai Yu. 2024. Text fluoroscopy: Detecting llmgenerated text through intrinsic features. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15838–15846.
- Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. 2022. Wordcraft: story writing with large language models. In 27th International Conference on Intelligent User Interfaces, pages 841–852.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. arXiv preprint arXiv:2303.18223.

884

885

887

890

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient finetuning of 100+ language models. arXiv preprint arXiv:2403.13372.

Α **Examples of Texts Generated by Different LLMs**

In this section, we provided examples of responses from different updates of the same LLMs to the same question, offering a visual representation of how changes in versions lead to variations in the model's output in Table 3 and Table 4.

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

927

B Additional Experiment

Additionally, to eliminate the influence of output diversity on detection performance, we first generated 150 texts from each model for every dataset, resulting in a total of 300 texts. Additionally, we conducted repeated experiments with several models to ensure that output diversity does not significantly impact changes in detection difficulty. By doing so, we ensure that the variations in detection difficulty primarily stem from model evolution itself, rather than biases introduced by the diversity of the generated texts. The results of this experimental process are presented in Table 5, where we highlight the potential range of biases caused by the diversity of generated texts. This helps further to understand the relationship between model evolution and detection performance.

С **Details of Evolving LLMs**

The details of evolving LLMs are shown in Table 6.

D **Details of Optimizing Strategies**

In this section, we provide the details of two op-919 timizing strategies. For the zero-shot detection 920 method, we chose to optimize the Fast-DetectGPT 921 detector. We use the pruned model 'princeton-922 nlp/Sheared-LLaMA-2.7B-Pruned' to replace the 923 scoring model 'meta-llama/Llama-2-7b-hf' model 924 as the scoring model. For the supervised detector, 925 we chose to optimize the Text Fluoroscopy. Follow-926 ing Text Fluoroscopy, we used the first 200 entries of the open-source Human-ChatGPT Comparison 928 Corpus (HC3) (Guo et al., 2023) dataset collected 929 by previous researchers as a training set. The ratio 930 for splitting the training and validation is 8:1. We 931 use a specific LLM of the LLM families to regener-932 ate machine text in the training dataset and then use 933 the new dataset and the original dataset to retrain 934 the detector. Specifically, we regenerated this part 935 of the training set using GPT-4-0613 and tested 936 the results on the GPT-4 family. For the GPT-40 937

Does a	irway surgery lower serum lipid levels in obstructive sleep apnea patients?
Human Answer	Surgery may improve blood lipid levels in patients with OSA, especially in patients with preoperative dyslipidemia, potentially yielding a major benefit in metabolism and cardiovascular sequelae. Prospective studies should examine this
GPT-40 2024-05-13	Airway surgery for obstructive sleep apnea (OSA) has been shown to reduce serum lipid levels in some patients. The improvement is likely due to better oxygenation and reduced inflammation following successful surgery,
GPT-40 2024-08-06	Airway surgery can potentially lower serum lipid levels in obstructive sleep apnea (OSA) patients. Surgical interventions that alleviate OSA symptoms may improve metabolic profiles, including lipid levels, due to reduced hypoxia and
GPT-40 2024-11-20	Yes, airway surgery in obstructive sleep apnea (OSA) patients may lower serum lipid levels. Studies suggest that by improving oxygenation, reducing inflammation, and addressing the metabolic impacts of OSA, surgeries like
GPT-40 Latest	Yes, airway surgery for obstructive sleep apnea (OSA) can lower serum lipid levels. By improving airway patency and reducing hypoxia, surgery may enhance lipid metabolism and decrease systemic inflammation, which are linked
GPT-40-mini 2024-07-18	Airway surgery in obstructive sleep apnea patients may lead to improvements in overall health and potentially lower serum lipid levels. However, the extent of lipid level changes can vary among individuals and

Table 3: An example of content generated by different versions of the GPT-40 and GPT-40-mini

Table 4: An example of content generated by different versions of the Claude-Haiku and Calude-Opus

Nec	ek pain treatment with acupuncture: does the number of needles matter?
Human Answer	For this pathology, the number of needles, 5 or 11, seems not to be an important variable in determining the therapeutic effect when the time of stimulation is the same in the two groups.
Claude-Haiku 2024-03-07	Acupuncture for neck pain can be effective, but the number of needles used may not be the most critical factor. The placement and selection of acupuncture points tailored to the individual's condition are more important
Claude-Haiku 2024-10-22	Research suggests that the number of acupuncture needles may not significantly impact neck pain treatment. Effectiveness depends more on precise needle placement, practitioner skill, and individual patient response. Further studies are needed to definitively determine
Claude-Opus 2024-02-29	The number of needles used in acupuncture for neck pain treatment may impact its effectiveness. However, factors such as needle placement, depth of insertion, and stimulation technique are likely more important than the specific number

family, we used GPT-4o-2024-05-13, and for the Claude-Haiku family, we used Claude-3-sonnet-20240229. we use gte-Qwen1.5-7B-instruct⁷ as the encoder and the classifier consists of three fully connected layers with Tanh function. The dimensions of the intermediate layers in the classifier are 1024 and 512, respectively. The batch size is set to 16, and Adam (Kingma and Ba, 2014) optimizer is employed with an initial learning rate of 3e - 3.

944

945

946

947

948

949

950

⁷https://huggingface.co/Alibaba-NLP/gte-Qwen1.5-7B-instruct

E Data Collection.

In this section, we describe the process of data collection. Most of these datasets were generated

	Fast-DetectGPT	RADAR	Text-Fluoroscopy	Imitate Before Detect	RoBERTa-base	RoBERTa-large	Likelihood	Rank	LogRank	Entropy	DetectGPT	LRR	NPR	DNA-GPT
GPT-3.5-turbo-2024-01-25	0.8301	0.8280	0.8904	0.8992	0.6621	0.6723	0.7701	0.6432	0.7663	0.4163	0.4910	0.7098	0.6072	0.8160
	0.8420	0.8319	0.8917	0.9011	0.6629	0.6696	0.7822	0.6511	0.7787	0.4118	0.5008	0.7120	0.6276	0.8266
	(+1.19%)	(+ 0.39 %)	(+0.13%)	(+ 0.19%)	(+0.08%)	(- 0.26 %)	(+1.21%)	(+0.78%)	(+1.24%)	(-0.44%)	(+0.98%)	(+0.22%)	(+2.05%)	(+ 1.06 %)
GPT-4o-mini-2024-07-18	0.8053	0.7954	0.8800	0.8831	0.5539	0.5799	0.7238	0.6474	0.7133	0.4617	0.5374	0.6341	0.6045	0.7255
	0.8033	0.7976	0.8777	0.8886	0.5566	0.5827	0.7264	0.6494	0.7154	0.4543	0.5417	0.6356	0.6101	0.7291
	(-0.20%)	(+0.22%)	(-0.23%)	(+ 0.55 %)	(+0.27%)	(+0.28%)	(+0.26%)	(+0.20%)	(+0.21%)	(-0.74%)	(+0.43%)	(+0.15%)	(+0.56%)	(+ 0.36 %)
GPT-4-turbo-2024-04-09	0.7644	0.7459	0.8535	0.8534	0.5053	0.5146	0.6959	0.6315	0.6856	0.4526	0.5229	0.6182	0.5982	0.7114
	0.7661	0.7595	0.8504	0.8559	0.5034	0.5099	0.6953	0.6278	0.6848	0.4533	0.5208	0.6109	0.5925	0.6952
	(+ 0.17%)	(+1.36%)	(-0.31%)	(+ 0.25 %)	(-0.18%)	(- 0.47%)	(-0.06%)	(-0.37%)	(-0.08%)	(+ 0.07%)	(-0.20%)	(-0.72%)	(-0.57%)	(-1.63%)

Table 5: We selected several LLMs and produced LLM-generated text many times, and then conducted inspections to evaluate the impact of diversity on the detector.

following the dataset construction methodology used in DetectGPT and Fast-detectGPT, which are widely recognized in the literature on AI-generated text detection. This setting minimizes the influence of factors, other than model evolution, on the detection performance.

951

952

954

955

956

957

961

962

963

965

968

970

971

972

973

974

975

977

978

981

982

983 984

985

988

To assess the domain generalization capability of detection methods, we included the following datasets: Writing, PubMed, Community, and Peer-Read. We also incorporated three different generation paradigms: continuation, question-answering, and paraphrasing. Specifically, the prompts used for generating each dataset are as follows:

- Xsum: "Please write an article with about 150 words starting exactly with: prefix>"
- Writing: "Please write an article with about 150 words starting exactly with: cprefix>"
- **PubMed**: "Please answer the question in about 50 words: <question>"
- SocialMedia: "Generate text similar to the input social media text but using different words and sentence composition: <Full text>"
- **PeerRead**: "Please write a peer review with about 150 words starting exactly with: <prefix>"

F Experimental Setup Details

In this section, we provide a detailed description of the experimental setup. Specifically, for zeroshot, we follow Fast-DetectGPT (Bao et al.) and use GPT-J-6B (Wang and Komatsuzaki, 2021) and GPT-Neo-2.7B (Black et al., 2021) as the scoring models for zero-shot methods. For supervised detectors, we used the pre-trained detectors provided by the authors. All experiments are conducted on a workstation equipped with 4 NVIDIA A100 GPUs.

For the valuation metric, we measure the detection performance using two metrics: AUROC(the area under the receiver operating characteristic) and EMG (evolving model generalization). A higher AUROC value indicates better detection quality, while a higher EMG metric indicates better generalization ability across evolving LLMs. 989

990

991

992

993

994

995

996

997

998

999

1000

1001

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

G The explanation of EMG

In this section, we provide an explanation of EMG. A positive EMG value indicates that the detection performance improves as the LLM evolves, while negative values reflect a decline in generalization. The magnitude of the EMG value also reflects the degree of improvement or decline in generalization; smaller values indicate worse stability. EMG reflects the direction and degree of fluctuations. However, it does not directly indicate the actual detection performance regarding AUROC.

For example, as shown in Table 2, in the GPT-40 family, the EMG values for RADAR and Text-Fluoroscopy are -0.537% and -0.493%, respectively. However, as shown in Table 1, the AUROC values for RADAR range from 0.8416 to 0.8694, while those for Text-Fluoroscopy fluctuate between 0.75426 to 0.79618. Although the EMG values are similar, the AUROC differs significantly, suggesting that both methods have comparable generalization abilities in the face of evolving LLMs, but their actual detection performance varies. Therefore, we recommend evaluating both AUROC and EMG together.

H Results of Other 12 Detection Methods

In this section, we provide the detailed results of the other detection methods in Table 7, 8, 9, 10, 11, 12, 13. The results are presented with four decimal places in the tables.

Evolving LLMs	Version Time	Source
GPT-40	2024-05-13	gpt-4o-2024-05-13
GPT-40	2024-08-06	gpt-4o-2024-08-06
GPT-40	2024-11-20	gpt-4o-2024-11-20
GPT-40	Latest	chatgpt-40-latest
GPT-4o-mini	2024-07-18	gpt-4o-mini-2024-07-18
GPT-4	2023-06-13	gpt-4-0613
GPT-4	2023-11-06	gpt-4-1106-preview
GPT-4	2024-01-25	gpt-4-0125-preview
GPT-4	2024-04-09	gpt-4-turbo-2024-04-09
Claude-Sonnet	2024-02-29	claude-3-sonnet-20240229
Claude-Sonnet	2024-06-20	claude-3-5-sonnet-20240620
Claude-Sonnet	2024-10-22	claude-3-5-sonnet-20241022
Claude-Haiku	2024-03-07	claude-3-haiku-20240307
Claude-Haiku	2024-10-22	claude-3-5-haiku-20241022
Claude-Opus	2024-02-29	claude-3-opus-20240229
Qwen	Qwen1.5-7B	Qwen/Qwen1.5-7B-Chat
Qwen	Qwen2-7B	Qwen2-7B-Instruct
Qwen	Qwen2.5-7B	Qwen/Qwen2.5-7B-Instruct
LlaMa3	Llama-3.1-8B	meta-llama/Meta-Llama-3.1-8B-Instruct
LlaMa3	Llama-3.1-70B	meta-llama/Meta-Llama-3.1-70B-Instruct
LlaMa3	Llama-3.2-1B	meta-llama/Meta-Llama-3.2-1B-Instruct
LlaMa3	Llama-3.2-3B	meta-llama/Meta-Llama-3.2-3B-Instruct
LlaMa3	Llama-3.3-70B	meta-llama/Meta-Llama-3.3-70B-Instruct
Fine-tuning	LlaMA-2-7B-chat-hf	meta-llama/Llama-2-7b-chat-hf
Fine-tuning	Vicuna-1.5-7B	lmsys/vicuna-7b-v1.5
Fine-tuning	Wizardmath-7B	WizardLMTeam/WizardMath-7B-V1.0
Pruning	Sheared-LlaMA1.3B	princeton-nlp/Sheared-LLaMA-1.3B
Pruning	Sheared-LlaMA1.3B-pruned	princeton-nlp/Sheared-LLaMA-1.3B-Pruned
Pruning	Sheared-LlaMA2.7B-pruned	princeton-nlp/Sheared-LLaMA-2.7B
Pruning	Sheared-LlaMA2.7B	princeton-nlp/Sheared-LLaMA-2.7B-Pruned

Table 6: Details of evolving LLMs.

LLMs	Version Time/			Fast-De	tectGPT					RA	DAR		
LEMO	Version Name	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
					τ	Jpdates							
	2024-05-13	0.8971	0.9688	0.7346	0.5756	0.8256	0.8003	0.9945	0.8281	0.8276	0.5976	0.6781	0.7852
	2024-08-06	0.8696	0.9447	0.6992	0.5932	0.7667	0.7747	0.9934	0.8608	0.8223	0.6188	0.6856	0.7962
CIDT 1	2024-00-00	(-2.75%)	(-2.41%)	(-3.54%)	(+1.76%)	(-5.89%)	(-2.57%)	(-0.11%)	(+3.27%)	(-0.53%)	(+2.12%)	(+0.75%)	(+1.10%)
GP1-40	2024-11-20	0.7202	0.8998	0.6969	0.5747	0.7831	0.7349	0.9858	0.7175	0.7868	0.6178	0.6634	0.7543
		0.6963	0.9075	0 7010	0.5820	0 8242	0 7422	0.9866	0.7532	0 7993	(+2.02%) 0.6422	0.6728	0 7708
	Latest	(-20.0%)	(-6.13%)	(-3.36%)	(+ 0.64 %)	(-0.14%)	(-5.81%)	(- 0.79%)	(-7.49%)	(-2.83%)	(+4.46%)	(-0.53%)	(-1.44%)
CPT 40 mini	2024 07 18	0.9062	0.9700	0.7304	0.6058	0.8042	0.8033	0.9978	0.8663	0.8208	0.6645	0.6385	0.7976
011-40-11111	2024-07-18	(+ 0.91%)	(+ 0.12%)	(-0.42%)	(+3.02%)	(-2.14%)	(+ 0.30 %)	(+ 0.33%)	(+3.82%)	(- 0.68%)	(+ 6.69%)	(-3.96%)	(+1.24%)
	2023-06-13	0.9262	0.9776	0.7611	0.4937	0.7511	0.7819	0.9944	0.8725	0.8918	0.6562	0.7490	0.8328
	2023-11-06	0.9177	0.9300	0.7317	0.5480	0.8282	0.7911	0.9943	0.7876	0.8037	0.6586	0.6259	0.7740
GPT-4		(-0.85%)	(-4.76%)	(-2.94%)	(+5.43%)	(+7.71%)	(+0.92%)	(-0.01%)	(-8.49%)	(-8.81%)	(+0.24%)	(-12.3%)	(-5.88%)
0114	2024-01-25	(-3.22%)	(-3.41%)	(-5.75%)	(+4.23%)	(+7.34%)	(-0.16%)	(-0.23%)	(-7.23%)	(-5.09%)	(-0.46%)	(-8.68%)	(-4.34%)
	****	0.8332	0.9120	0.7172	0.5764	0.7919	0.7661	0.9908	0.7880	0.8211	0.5970	0.6008	0.7595
	2024-04-09	(- 9.30%)	(-6.56%)	(-4.39%)	(+ 8.27 %)	(+4.08%)	(-1.58%)	(- 0.36%)	(-8.45%)	(-7.07%)	(-5.92%)	(-14.8%)	(-7.32%)
	2024-02-29	0.9459	0.9804	0.8150	0.5836	0.9320	0.8514	0.9700	0.8582	0.8403	0.7355	0.6106	0.8029
	2024 06 20	0.9692	0.9884	0.7718	0.6605	0.9030	0.8586	0.9861	0.8061	0.8159	0.7250	0.5549	0.7776
Claude-Sonnet	2024-00-20	(+2.33%)	(+0.80%)	(-4.32%)	(+ 7.69%)	(-2.90%)	(+ 0.72%)	(+1.61%)	(-5.21%)	(-2.44%)	(-1.05%)	(-5.57%)	(-2.53%)
	2024-10-22	0.9016	0.9265	0.6760	0.5056	0.7789	0.7577	0.9548	0.7238	0.7418	0.6681	0.5468	0.7271
		(-4.43%)	(-5.39%)	(-13.9%)	(-7.80%)	(-15.3%)	(-9.37%)	(-1.52%)	(-13.4%)	(-9.85%)	(-6.74%)	(-6.38%)	(-7.59%)
	2024-03-07	0.9952	0.9993	0.8614	0.7486	0.9389	0.9087	0.9972	0.9328	0.8388	0.7669	0.6745	0.8420
Claude-Haiku	2024-10-22	0.8396	0.9237	0.6439	0.3892	0.4641	0.6521	0.9965	0.8578	0.7293	0.7031	0.8217	0.8217
		(-15.5%)	(-7.50%)	(-21.7%)	(-35.9%)	(-47.4%)	(-25.6%)	(-0.07%)	(-7.50%)	(-10.9%)	(-0.38%)	(+14.7%)	(-2.04%)
Claude-Opus	2024-02-29	(-2.49%)	(-3.89%)	(-3.80%)	(-3.27%)	(-4.46%)	(-3.58%)	(-0.33%)	(-6.01%)	(-2.23%)	$(\pm 0.29\%)$	(-5.57%)	(-2.77%)
	Owen1 5-7B	0.9243	0.9945	0.6641	0.4897	0.9076	0.7960	0.9269	0.9034	0.7222	0.5074	0.4761	0.7072
	Qweii1.5=7B	0.9245	0.9974	0.7376	0.4622	0.8939	0.8168	0.9209	0.8692	0.8168	0.5209	0.5375	0.7260
Qwen	Qwen2-7B	(+6.86%)	(+ 0.29%)	(+7.35%)	(-2.75%)	(-1.37%)	(+2.08%)	(-4.13%)	(-3.42%)	(+9.46%)	(+1.35%)	(+6.14%)	(+1.88%)
	Owen2.5.7P	0.9894	0.9955	0.7878	0.4686	0.8998	0.8282	0.9009	0.9100	0.7907	0.5205	0.5256	0.7295
	Qwell2.5=7B	(+6.51%)	(+0.10%)	(+12.3%)	(-2.11%)	(-0.78%)	(+3.22%)	(-2.60%)	(+0.66%)	(+6.85%)	(+1.31%)	(+4.95%)	(+2.23%)
	LlaMA-3.1-8B	0.9926	0.9839	0.8476	0.6866	0.9680	0.8957	0.9987	0.8939	0.7784	0.7232	0.6304	0.8049
	LlaMA-3.1-70B	0.9957	0.9989	0.8104	0.6693	0.9584	0.8865	0.9968	0.8807	0.7627	0.7361	0.6182	0.7989
11-1442		(+ 0.31%)	(+1.50%)	(-3.72%)	(-1.73%)	(- 0.96%)	(-0.92%)	(- 0.19%)	(-1.32%)	(-1.57%)	(+1.29%)	(-1.22%)	(-0.60%)
LIAMAS	LlaMA-3.2-1B	0.9580	0.9948	(17.20%)	0.8617	0.9650	0.9400	0.9636	0.9378	0.8158	0.7512	0.7056	0.8348
		0.9913	(+1.09%)	(+1.29%) 0.9132	(+17.5%) 0.7698	0.9650	0.9253	(-3.51%) 0.9943	(+4.39%) 0.9210	0.8154	(+2.80%)	0.6584	0.8309
	LlaMA-3.2-3B	(-0.13%)	(+0.31%)	(+6.56%)	(+8.32%)	(-0.30%)	(+2.95%)	(-0.44%)	(+2.71%)	(+3.70%)	(+4.23%)	(+2.80%)	(+2.60%)
	LI0MA 2 2 70P	0.9958	1.0000	0.7923	0.6881	0.9042	0.8761	0.9955	0.8865	0.7878	0.7308	0.6533	0.8108
	LiawiA=5.5=70B	(+0.32%)	(+1.61%)	(-5.53%)	(+0.15%)	(-6.38%)	(-1.97%)	(-0.32%)	(-0.74%)	(+ 0.94 %)	(+ 0.76%)	(+2.29%)	(+0.59%)
					Dev	elopments							
	LlaMA-2-7B-chat-hf	0.9722	0.9904	0.9055	0.8963	0.9693	0.9467	0.6808	0.6330	0.4911	0.6396	0.6380	0.6165
	Vicupa-1 5-7B	0.9958	0.9956	0.8008	0.8805	0.9439	0.9233	0.9745	0.9305	0.8064	0.6490	0.6355	0.7992
Fine-tuning	vicula-1.5-715	(+ 2.36%)	(+ 0.52%)	(-10.4%)	(-1.58%)	(-2.54%)	(-2.34%)	(+29.3%)	(+ 29.7%)	(+31.5%)	(+ 0.94%)	(-0.25%)	(+18.2%)
-	Wizardmath-7B	0.9205	0.9804	0.7354	0.8490	0.8412	0.8653	0.7476	0.7582	0.6969	0.6857	0.5447	0.6866
		(- 5.17%)	(-1.00%)	(-17.0%) 0.4810	(-4.73%)	(-12.8%) 0.7504	(- 8.14%)	(+ 0.08%)	(+12.5%)	(+20.5%)	(+4.01%) 0.5787	(-9.33%) 0.4464	(+7.01%)
	Llama-2-7B-hf	(-30.4%)	(-8.81%)	(-42.4%)	(-11.8%)	(-20.9%)	(-22.9%)	(+23.6%)	(+17.2%)	(+20.5%)	(-6.09%)	(-19.1%)	(+7.23%)
		0.6548	0.8888	0.4846	0.8634	0.7072	0.7198	0.7706	0.6626	0.5856	0.7255	0.5360	0.6561
	Sneared-LlaMA1.3B	(-31.7%)	(-10.1%)	(-42.0%)	(- 3.29%)	(-26.2%)	(-22.7%)	(+ 8.98%)	(+ 2.96%)	(+9.45%)	(+8.59%)	(-10.2%)	(+3.96%)
	Sheared-LlaMA1 3B-pruned	0.4118	0.8077	0.3524	0.7530	0.4593	0.5568	0.9213	0.7903	0.6991	0.8669	0.6539	0.7863
Pruning		(-56.0%)	(-18.2%)	(-55.3%)	(-14.3%)	(-51.0%)	(-38.9%)	(+24.0%)	(+15.7%)	(+20.8%)	(+22.7%)	(+1.59%)	(+16.9%)
	Sheared-LlaMA2.7B-pruned	0.4480	0.8672	0.3664	0.7543	0.5196	0.5911	0.8383	0.7618	0.5870	0.7831	0.5455	0.7031
2	~	(-52.4%) 0.7049	(-12.5%) 0.9154	(- 53.9%) 0 5047	(-14.2%) 0.8191	(-44.9%) 0.6520	(-35.5%) 0.7192	(+15.7%) 0.7496	(+12.8%) 0.6811	(+ 9.39%) 0 5794	(+14.3%) 0.6743	(-9.25%) 0.4655	(+ 8.00%)
	Sheared-LlaMA2.7B	(-26.7%)	(-7.50%)	(-40.0%)	(-7.72%)	(-31.7%)	(-22.7%)	(+6.88%)	(+4.81%)	(+8.83%)	(+3.47%)	(-17.2%)	(+1.35%)
		S	1	1	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		S	100 m 100 m	1	100 C 100 C 100 C	1 S 1 S 1 S 1	

Table 7: The detection performance (measured in AUROC) of Fast-DetectGPT and RADAR on EvoBench.

ImageYears <th< th=""><th>LLMs</th><th>Version Time/</th><th></th><th></th><th>Text-Flu</th><th>ioroscopy</th><th></th><th></th><th></th><th></th><th>Imitate Be</th><th>fore Detect</th><th></th><th></th></th<>	LLMs	Version Time/			Text-Flu	ioroscopy					Imitate Be	fore Detect		
Update GP7-40 202-40-13 0.0994 0.9870 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9871 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9881 0.9880 0.9880 0.9880 <	LLWIS	Version Name	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
3024-66-13 0.0994 0.0821 0.8670 0.0874 0.0864 0.0964 0.9855 0.7129 0.0861 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0256 0.0257 0.7366 0.0287 0.7366 0.0287 0.7366 0.0287 0.7366 0.0287 0.0288 0.0						1	Updates							
GPT-40 0.998 0.9877 0.9823 0.7072 0.816 0.8975 0.9724 0.9036 0.7254 0.7156 0.9386 0.7254 0.9376 0.9386 0.7254 0.9376 0.9386 0.7254 0.9376 0.9386 0.9376 0.9386 <td></td> <td>2024-05-13</td> <td>0.9994</td> <td>0.9823</td> <td>0.8510</td> <td>0.6870</td> <td>0.8124</td> <td>0.8664</td> <td>0.9964</td> <td>0.9835</td> <td>0.7189</td> <td>0.6983</td> <td>0.9126</td> <td>0.8619</td>		2024-05-13	0.9994	0.9823	0.8510	0.6870	0.8124	0.8664	0.9964	0.9835	0.7189	0.6983	0.9126	0.8619
CPT-Lo C+0.01% C+0.01% <th< td=""><td></td><td>2024-08-06</td><td>0.9995</td><td>0.9867</td><td>0.8423</td><td>0.7072</td><td>0.8116</td><td>0.8695</td><td>0.9974</td><td>0.9836</td><td>0.7254</td><td>0.7416</td><td>0.9056</td><td>0.8707</td></th<>		2024-08-06	0.9995	0.9867	0.8423	0.7072	0.8116	0.8695	0.9974	0.9836	0.7254	0.7416	0.9056	0.8707
GP1-40 202-11-20 0.997 0.0581 0.1300 0.0581 0.0380 0.2889 0.2899 0.289		2024-00-00	(+0.01%)	(+ 0.44%)	(-0.87%)	(+2.02%)	(-0.08%)	(+ 0.30 %)	(+ 0.10%)	(+0.01%)	(+ 0.65%)	(+4.33%)	(-0.70%)	(+ 0.88%)
Latest 10099 0.033 0.0360 0.0399 0.0389 0.0399 0.0580 0.0389 0.0386 0.0389 0.0386 0.0389 0.0386 0.0389 0.0386 0.0389 0.0386 0.0386 0.0386 0.0386 0.0386 0.0386 0.0386 0.0386 <td>GPI-40</td> <td>2024-11-20</td> <td>0.9997</td> <td>0.9584</td> <td>0.8130</td> <td>0.6877</td> <td>0.7496</td> <td>0.8417</td> <td>0.9802</td> <td>0.9480</td> <td>0.6889</td> <td>0.7358</td> <td>0.8998</td> <td>0.8505</td>	GPI-40	2024-11-20	0.9997	0.9584	0.8130	0.6877	0.7496	0.8417	0.9802	0.9480	0.6889	0.7358	0.8998	0.8505
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			(+0.05%)	(-2.39%)	(-3.80%)	(+0.07%)	(-0.28%)	(-2.47%)	(-1.02%)	(-3.55%)	(-3.00%)	(+3.75%)	(-1.28%)	(-1.14%)
GPT-do.mini 2024-07.18 L0000 0.0589 0.0387 0.0589 0.0589 0.0587 0.0517 0.0586 GPT-do.mini 2023-06-13 0.9997 0.0850 0.2245 (-11.058) (+0.06%) (+0.26%)		Latest	(-0.00%)	(-5.90%)	(-4.90%)	(+4.30%)	(+1.69%)	(-0.96%)	(-1.26%)	(-2.69%)	(-3.42%)	(+4.33%)	(-0.05%)	(-0.62%)
GPT-4o-min 2024-06-13 0.9997 0.9950 0.9550 0.9728 (-2.23%) (-2.36%) (-9.94%) (-0.02%) (-2.36%) (-9.94%) (-0.02%) (-2.36%) (-0.23%) (-2.37%) <th< td=""><td></td><td></td><td>1.0000</td><td>0.9894</td><td>0.8438</td><td>0.7185</td><td>0.8369</td><td>0.8777</td><td>0.9974</td><td>0.9897</td><td>0.7555</td><td>0.7887</td><td>0.9117</td><td>0.8886</td></th<>			1.0000	0.9894	0.8438	0.7185	0.8369	0.8777	0.9974	0.9897	0.7555	0.7887	0.9117	0.8886
2023-06-13 0.9997 0.9580 0.0738 0.0730 0.7280 0.0577 0.7094 0.0739 0.08459 0714 0.0255 0.1722 0.2580 0.7380 0.0850 0.0971 0.0784 0.0784 0.0784 0.0784 0.0784 0.0784 0.0784 0.0845 2024-01-50 0.9998 0.9750 0.0877 0.1885 0.0978 0.1885 0.0978 0.0783 0.0864 0.9998 0.7214 0.0804 0.9988 0.716 0.8991 0.716 <td>GPT-40-mini</td> <td>2024-07-18</td> <td>(+0.06%)</td> <td>(+0.71%)</td> <td>(-0.72%)</td> <td>(+3.15%)</td> <td>(+2.45%)</td> <td>(+1.13%)</td> <td>(+0.10%)</td> <td>(+0.62%)</td> <td>(+3.66%)</td> <td>(+9.04%)</td> <td>(-0.09%)</td> <td>(+2.67%)</td>	GPT-40-mini	2024-07-18	(+0.06%)	(+0.71%)	(-0.72%)	(+3.15%)	(+2.45%)	(+1.13%)	(+0.10%)	(+0.62%)	(+3.66%)	(+ 9.04%)	(-0.09%)	(+2.67%)
GPT-I 2023-11-06 0.9995 0.9728 0.7288 0.7288 0.7288 0.8333 0.9996 0.7397 0.7848 0.9707 0.7848 0.9707 0.7848 0.9707 0.7848 0.9707 0.7848 0.9707 0.7845 0.7335 0.7313 0.73		2023-06-13	0.9997	0.9850	0.8748	0.6930	0.7289	0.8563	0.9951	0.9877	0.7094	0.6753	0.8619	0.8459
$ \begin{array}{c} {\rm GPT-4} & (-0.02^{+}) & (-1.28^{+}) & (-5.18^{+}) & (+3.88^{+}) & (+0.78^{+}) & (-0.40^{+}) & (-1.98^{+}) & (+3.18^{+}) & (+1.10^{+}) & (+4.28^{+}) & (+4.48^{+}) \\ (-0.01^{+}) & (-1.01^{+}) & (-1.01^{+}) & (-2.67^{+}) & (-2.18^{+}$		2023-11-06	0.9995	0.9722	0.8230	0.7298	0.7368	0.8523	0.9970	0.9684	0.7907	0.7854	0.9071	0.8897
OF H 2024-01-25 (0.9998 0.0999 (0.0099 (0.0997) (0.8097) (0.8097) (0.8085) (0.9086) (0.9086) (0.9086) (0.9086) (0.1998) (-1.2185)	CDT 4	2020 11 00	(-0.02%)	(-1.28%)	(-5.18%)	(+3.68%)	(+ 0.79%)	(-0.40%)	(+0.19%)	(-1.93%)	(+8.13%)	(+11.0%)	(+4.52%)	(+4.38%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	GF I-4	2024-01-25	0.9998	0.9750	0.8079	0.7099	0.7397	0.8465	0.9956	0.9708	0.7433	0.7214	0.8936	0.8649
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			0 9998	0.9779	0.8171	0.6901	0 7671	0.8504	(+0.05%)	0.9615	07136	(+4.01%) 0.7186	0.8930	(+1.91%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		2024-04-09	(+0.01%)	(-0.71%)	(-5.77%)	(-0.29%)	(+3.82%)	(-0.59%)	(-0.21%)	(-2.62%)	(+0.42%)	(+4.33%)	(+3.11%)	(+1.01%)
Claude-Somet 0.0930 0.9759 0.9759 0.9818 0.9807 0.9837 0.9837 0.9817 0.9817 0.9022 Claude-Somet 0.2024 0.202 0.9906 0.9916 0.9554 0.7356 (-1.355) (2024-02-29	0.0013	0.9810	0.8494	0.7582	0 7972	0.8754	0.9804	0.9850	0.7732	0.8864	0.9080	0.9066
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		2024=02=29	0.9915	0.9810	0.8494	0.8188	0.7972	0.87.54	0.9804	0.9850	0.8092	0.9127	0.9057	0.9000
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Claude-Sonnet	2024-06-20	(-0.23%)	(-0.51%)	(+2.27%)	(+6.06%)	(+4.96%)	(+2.51%)	(+1.53%)	(+0.33%)	(+3.60%)	(+2.63%)	(-0.23%)	(+1.57%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		2024 10 22	0.9906	0.9421	0.8096	0.7843	0.8167	0.8687	0.9916	0.9564	0.7138	0.7866	0.8437	0.8584
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		2024-10-22	(-0.07%)	(-3.89%)	(-3.98%)	(+ 2.61%)	(+1.95%)	(-0.68%)	(+ 1.12%)	(-2.86%)	(-5.94%)	(- 9.98%)	(-6.43%)	(-4.82%)
		2024-03-07	0.9998	0.9856	0.8482	0.8182	0.7997	0.8903	0.9996	0.9981	0.8428	0.9192	0.9331	0.9386
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Claude-Haiku	2024 10 22	0.9999	0.9650	0.8436	0.7944	0.8591	0.8924	0.9981	0.9796	0.7884	0.8375	0.7420	0.8691
Chaude-Opus 2024-02-29 0.9968 0.9801 0.9803 0.0800 0.9116 0.8900 0.9176 0.9803 0.0800 0.9116 0.8900 0.9176 0.9176 0.9581 0.7974 0.8143 (-1.4385)		2024-10-22	(+0.01%)	(-2.06%)	(-0.46%)	(-2.38%)	(+5.94%)	(+0.21%)	(-0.15%)	(-1.85%)	(-5.44%)	(-8.17%)	(-19.1%)	(-6.94%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Claude-Opus	2024-02-29	0.9968	0.9801	0.8403	0.7974	0.8140	0.8857	0.9899	0.9803	0.8080	0.9116	0.8950	0.9170
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $			(-0.30%)	(-0.55%)	(-0.79%)	(-2.08%)	(+1.43%)	(-0.46%)	(-0.97%)	(-1.78%)	(-3.48%)	(-0.76%)	(-3.81%)	(-2.16%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		Qwen1.5-7B	0.9017	0.9581	0.7955	0.5083	0.5593	0.7446	0.9493	0.9822	0.7517	0.4690	0.7928	0.7890
Qwen C (-2.04%) (-2.04%) (+1.55%) (+1.55%) (+2.07%) (+0.01%) (-4.97%) (+4.08%) (-1.44%) (-1.43%) (-1.44%) (-1.43%) (-1.44%) (-1.43%) (-1.24%) (-1.2	0	Owen2-7B	0.6711	0.9378	0.7617	0.4819	0.5746	0.6854	0.9765	0.9883	0.7018	0.4779	0.7784	0.7846
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Qwen		(-23.0%)	(-2.03%)	(-3.38%)	(-2.64%)	(+1.53%)	(-5.92%)	(+2.72%)	(+0.61%)	(-4.99%)	(+ 0.89%)	(-1.44%)	(-0.44%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		Qwen2.5-7B	(12.8%)	(124%)	(747%)	(301%)	(1138%)	(463%)	(183%)	(1034%)	(10.7394)	(1078%)	(127%)	(10.7939)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		X1 X4 A 1 AD	0.0007	0.0705	0.0100	0.7554	(11.50 %)	(4.05 %)	(1100 %)	(10.0470)	0.0(20	(10.7072)	0.0040	0.0000
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		LIaMA-3.1-8B	0.9997	0.9795	0.8108	0.7554	0.6524	0.8396	0.9988	0.9896	0.8630	0.8855	0.9040	0.9282
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		LlaMA-3.1-70B	$(\pm 0.02\%)$	(+1.13%)	(-1.30%)	$(\pm 1.46\%)$	(-3.45%)	(-0.45%)	(-0.07%)	$(\pm 0.61\%)$	(-10.3%)	$(\pm 2.48\%)$	(-1.05%)	(-168%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	LlaMA3		0.9838	0.9844	0.8541	0.8428	0.6758	0.8682	0.9915	0.9963	0.8981	0.9476	0.9146	0.9496
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		LlaMA-3.2-1B	(-1.59%)	(+ 0.49%)	(+4.33%)	(+8.74%)	(+2.34%)	(+2.86%)	(-0.73%)	(+ 0.67%)	(+3.51%)	(+6.21%)	(+1.06%)	(+2.14%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		LIOMA 2.2.2P	0.9994	0.9779	0.8383	0.8077	0.6673	0.8581	0.9992	0.9952	0.8818	0.9579	0.9192	0.9507
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		LiawiA-5.2-5D	(-0.03%)	(-0.16%)	(+2.75%)	(+5.23%)	(+1.49%)	(+1.86%)	(+ 0.04%)	(+ 0.56%)	(+1.88%)	(+7.24%)	(+1.52%)	(+2.25%)
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		LlaMA-3.3-70B	1.0000	0.9883	0.8240	0.7978	0.7587	0.8738	0.9969	0.9954	0.7246	0.9142	0.8629	0.8988
Developments Bit in the system of the system			(+0.03%)	(+0.88%)	(+1.32%)	(+4.24%)	(+10.6%)	(+3.42%)	(-0.19%)	(+0.58%)	(-13.8%)	(+2.87%)	(-4.11%)	(-2.94%)
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $						Dev	elopments							
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		LlaMA-2-7B-chat-hf	0.8320	0.8940	0.7480	0.8427	0.7311	0.8096	0.7310	0.7203	0.5381	0.8248	0.9436	0.7516
$ \begin{array}{c} \mbox{Fine-tuning} \\ \mbox{Fine-tuning} \\ \mbox{Wizardmath-7B} \\ \mbox{Wizardmath-7B} \\ \mbox{Uizardmath-7B} \\ Uizardmath-7$		Vicuna-1.5-7B	0.9138	0.9150	0.6650	0.6332	0.6375	0.7529	0.9974	0.9952	0.7742	0.7408	0.8704	0.8756
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Fine-tuning		(+8.18%)	(+ 2.10%)	(-8.30%)	(-20.9%)	(- 9.36%)	(-5.67%)	(+26.6%)	(+27.4%)	(+23.6%)	(-8.40%)	(-7.32%)	(+12.4%)
$ \begin{array}{c} \mbox{Hama-2-7B-hf} & (-12.6\%) & (-11.6\%) & (-11.6\%) & (-12.6\%) & (-1$		Wizardmath-7B	(-28.9%)	(-11.2%)	(-17.0%)	(-187%)	(-12.6%)	(-17.7%)	(+15.1%)	$(\pm 20.5\%)$	$(\pm 17.3\%)$	(-7.26%)	(-15.5%)	$(\pm 6.05\%)$
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $			0.4764	0.6419	0.4751	0.7159	0.5464	0.5711	0.9748	0.9700	0.8896	0.7722	0.7270	0.8667
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		Llama-2-7B-hf	(-35.5%)	(-25.2%)	(-27.2%)	(-12.6%)	(-18.4%)	(-23.8%)	(+24.3%)	(+24.9%)	(+35.1%)	(-5.26%)	(-21.6%)	(+11.5%)
Pruning Sheared-LlaMA1.3B-pruned (-30.6%) (-22.6%) (-4.16%) (-20.2%) (-20.6%) (+1.68%) (+2.34%) (+2.10%) (-5.48%) (-21.9%) (-4.25%) Pruning Sheared-LlaMA1.3B-pruned (-30.5%) (-30.6%) (-9.416%) (-20.6%) (-4.16%) (-10.6%) (+2.34%) (+2.10%) (-5.48%) (-21.9%) (-4.25%) Sheared-LlaMA1.3B-pruned (-30.5%) (-9.49%) (-4.56%) (-8.71%) (-3.46%) (-5.48%) (-11.1%) (-17.4%) (-3.33%) (-17.4%) (-3.41%) (-2.15%) (-7.72%) Sheared-LlaMA2.7B 0.5864 0.6410 0.5716 0.5846 0.5841 0.742 0.748 0.5405 0.6300 0.6713 Sheared-LlaMA2.7B 0.5864 0.6440 0.5300 0.6716 0.5846 0.5414 0.7472 0.7484 0.5405 0.6300 0.6713 Sheared-LlaMA2.7B 0.5864 0.6515 0.6637 0.5266 0.5959 0.7030 0.6804 0.4982 0.6782 0.6321		Chaorad LloMA 1 2D	0.5258	0.6373	0.5212	0.8011	0.5287	0.6028	0.7478	0.7437	0.5591	0.7700	0.7246	0.7090
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		Sneared-LIaMA1.3B	(-30.6%)	(-25.6%)	(-22.6%)	(-4.16%)	(-20.2%)	(-20.6%)	(+ 1.68%)	(+2.34%)	(+2.10%)	(-5.48%)	(-21.9%)	(-4.25%)
Prunng (+0.39%) (+4.98%) (-9.40%) (-4.56%) (-8.71%) (-3.46%) (-5.48%) (-11.4%) (-17.4%) (-3.83%) (-17.4%) Sheared-LlaMA2.7B-pruned 0.4902 0.6440 0.5300 0.6716 0.5841 0.7472 0.7484 0.5405 0.7056 0.6300 0.6743 Sheared-LlaMA2.7B (-34.5%) (-21.8%) (-17.1%) (-14.6%) (-22.5%) (+1.62%) (+2.81%) (+0.24%) (-11.9%) (-7.72%) Sheared-LlaMA2.7B 0.5864 0.6454 0.5215 0.6637 0.5266 0.5959 0.7030 0.6804 0.492 0.6789 0.6021 0.6324 Sheared-LlaMA2.7B (-24.5%) (-24.8%) (-22.6%) (-17.9%) (-16.8%) (-21.3%) (-2.30%) (-3.99%) (-3.49%) (-11.9%) (-3.41%) (-11.9%) (-11.9%) (-11.9%) (-11.9%) (-11.9%) (-17.2%)	. ·	Sheared-LlaMA1.3B-pruned	0.8359	0.9438	0.6540	0.7971	0.6440	0.7750	0.6762	0.5737	0.4265	0.6502	0.5600	0.5773
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Pruning	in manual pranoa	(+ 0.39%)	(+4.98%)	(- 9.40%)	(-4.56%)	(-8.71%)	(-3.46%)	(-5.48%)	(-14.6%)	(-11.1%)	(-17.4%)	(-38.3%)	(-17.4%)
$ \begin{array}{c} (-34.1\%) & (-2.30\%) & (-21.4\%) & (-11.4\%) & (-14.0\%) & (-24.5\%) & (+1.02\%) & (+1.02\%) & (+1.02\%) & (+1.02\%) & (+1.02\%) & (-11.9\%) & (-11.9\%) & (-11.9\%) \\ \hline \\ Sheared-LiaMA2.7B & \begin{array}{c} 0.5864 & 0.6454 & 0.5215 & 0.6637 & 0.5626 & 0.5959 & 0.7030 & 0.6804 & 0.4982 & 0.6789 & 0.6021 & 0.6328 & (-21.3\%) & (-2.80\%) & (-3.99\%) & (-3.99\%) & (-3.99\%) & (-3.99\%) & (-3.91\%) & (-11.9\%) \\ \hline \\ \end{array} $		Sheared-LlaMA2.7B-pruned	0.4902	0.6440	0.5300	0.6/16	0.5846	0.5841	0.7472	0./484	0.5405	0.7056	0.6300	0.6743
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			0 5864	0 6454	0 5215	0.6637	0.5626	0 5959	(+1.02%) 0.7030	0 6804	0 4982	0.6789	0.6021	0.6325
		Sheared-LlaMA2.7B	(-24.5%)	(-24.8%)	(-22.6%)	(-17.9%)	(-16.8%)	(-21.3%)	(-2.80%)	(-3.99%)	(-3.99%)	(-14.5%)	(-34.1%)	(-11.9%)

Table 8: The detection performance (measured in AUROC) of Text-Fluoroscopy and Imitate Before Detect on EvoBench.

LLMs	Version Time/			RoBEI	RTa-base					RoBEF	Ta-large		
LLING	Version Name	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
					τ	Jpdates							
	2024-05-13	0.6696	0.4668	0.4680	0.4946	0.5584	0.5315	0.6124	0.3163	0.5102	0.5804	0.6031	0.5245
	2024-08-06	0.5975	0.4108	0.4644	0.5487	0.5425	0.5128	0.6478	0.3191	0.5079	0.5897	0.5646	0.5258
	2024-08-00	(-7.21%)	(-5.60%)	(-0.36%)	(+5.41%)	(-1.59%)	(-1.87%)	(+3.54%)	(+0.28%)	(-0.23%)	(+ 0.93%)	(-3.85%)	(+0.13%)
GPT-40	2024-11-20	0.5858	0.3992	0.4331	0.4644	0.5061	0.4777	0.4226	0.2562	0.4536	0.5414	0.5402	0.4428
		(-8.38%) 0.5417	(-0.70%) 0.3060	(- 3.49%) 0.4386	(-3.02%)	(-5.25%) 0.5269	(- 5.38%) 0.4859	(-18.9%) 0.4151	(-0.01%) 0.2583	(- 5.00%) 0.4500	(-3.90%)	(-0.29%) 0.5880	(- 8.17%)
	Latest	(-12.7%)	(-6.99%)	(-2.94%)	(+3.10%)	(-3.15%)	(-4.55%)	(-19.7%)	(-5.80%)	(-5.03%)	(+3.46%)	(-1.42%)	(-5.70%)
CPT 40 mini	2024 07 18	0.7376	0.4991	0.4542	0.5608	0.5314	0.5566	0.7445	0.4498	0.5256	0.6047	0.5887	0.5827
GF 1-40-11111	2024-07-18	(+6.80%)	(+3.23%)	(-1.38%)	(+6.62%)	(-2.70%)	(+2.51%)	(+13.2%)	(+13.3%)	(+1.54%)	(+ 2.43 %)	(-1.44%)	(+5.82%)
	2023-06-13	0.7349	0.5198	0.5286	0.5664	0.5205	0.5740	0.7368	0.3972	0.5917	0.6340	0.5932	0.5906
	2023-11-06	0.6978	0.4235	0.4483	0.5488	0.5584	0.5354	0.6471	0.3444	0.5056	0.5924	0.6205	0.5420
GPT-4		(-3.71%)	(-9.63%)	(-8.03%)	(-1.76%)	(+3.79%)	(-3.87%)	(- 8.97%)	(-5.28%)	(-8.61%)	(-4.16%)	(+2.73%)	(-4.86%)
0114	2024-01-25	(-6.21%)	(-2.50%)	(-7.41%)	(-212%)	$(\pm 2.79\%)$	(-3.09%)	(-7.34%)	(-3.83%)	(-6.96%)	(-4.39%)	$(\pm 2.05\%)$	(-4.09%)
		0.5949	0.3947	0.4817	0.5029	0.5430	0.5034	0.6030	0.2916	0.5179	0.5550	0.5820	0.5099
	2024-04-09	(-14.0%)	(-12.5%)	(-4.69%)	(-6.35%)	(+2.25%)	(-7.06%)	(-13.3%)	(-10.5%)	(-7.38%)	(-7.90%)	(-1.12%)	(-8.07%)
	2024-02-29	0.7459	0.5915	0.5312	0.6558	0.5397	0.6128	0.6757	0.4837	0.5679	0.7196	0.6077	0.6109
	2024 06 20	0.7550	0.5054	0.4776	0.5935	0.5355	0.5734	0.6328	0.3337	0.5254	0.6633	0.5861	0.5483
Claude-Sonnet	2024-06-20	(+0.91%)	(-8.61%)	(-5.36%)	(-6.23%)	(-0.42%)	(-3.94%)	(-4.29%)	(-15.0%)	(-4.25%)	(-5.63%)	(-2.16%)	(-6.27%)
	2024-10-22	0.6080	0.3702	0.3741	0.5283	0.5082	0.4778	0.5696	0.2660	0.4148	0.5533	0.5147	0.4637
		(-13.7%)	(-22.1%)	(-15.7%)	(-12.7%)	(-3.15%)	(-13.5%)	(-10.6%)	(-21.7%)	(-15.3%)	(-16.6%)	(-9.30%)	(-14.7%)
	2024-03-07	0.9427	0.8251	0.5400	0.6579	0.6523	0.7236	0.9162	0.6806	0.6079	0.7053	0.6716	0.7163
Claude-Haiku	2024-10-22	0.6408	0.3626	0.3556	0.5740	0.5146	0.4895	0.6422	0.2771	0.4116	0.6527	0.5997	0.5167
		(-30.1%)	(-46.2%)	(-18.4%)	(-8.39%)	(-13.7%)	(-23.4%)	(-27.4%)	(-40.3%)	(-19.6%)	(-5.26%)	(-7.19%)	(-19.9%)
Claude-Opus	2024-02-29	(-10.9%)	(-19.4%)	(-5.25%)	(-7.54%)	(-9.58%)	(-10.5%)	(-14.8%)	(-18.4%)	(-8.99%)	(-3.42%)	(-7.62%)	(-10.6%)
	Owen1 5-7B	0.8357	0.8064	0.6017	0.4882	0 5330	0.6532	0.8891	0.7393	0.6413	0.5203	0.5213	0.6623
	Qweiri.5-7B	0.9770	0.9259	0.6218	0.4980	0.5404	0.7126	0.9624	0.8716	0.6588	0.4953	0.5440	0.7064
Qwen	Qwen2-7B	(+14.1%)	(+11.9%)	(+2.01%)	(+ 0.98%)	(+ 0.65 %)	(+5.94%)	(+7.33%)	(+13.2%)	(+1.75%)	(-2.50%)	(+2.27%)	(+4.42%)
	Owen2.5.7P	0.9776	0.8910	0.6249	0.5027	0.5592	0.7111	0.9645	0.8508	0.6841	0.5017	0.5552	0.7113
	Qwell2.5=7B	(+14.1%)	(+8.46%)	(+2.32%)	(+1.45%)	(+2.53%)	(+5.79%)	(+7.54%)	(+11.1%)	(+4.28%)	(-1.86%)	(+3.39%)	(+4.90%)
	LlaMA-3.1-8B	0.9590	0.8822	0.5254	0.5528	0.6138	0.7066	0.9344	0.7625	0.5748	0.5896	0.6216	0.6966
	LlaMA-3.1-70B	0.9676	0.8311	0.4566	0.5980	0.5769	0.6860	0.9080	0.7065	0.5002	0.6352	0.5824	0.6665
LI-MA2		(+ 0.86%)	(-5.11%)	(-6.88%)	(+4.52%)	(-3.69%)	(-2.06%)	(-2.64%)	(- 5.60%)	(-7.46%)	(+4.56%)	(-3.92%)	(-3.01%)
LIawiA5	LlaMA-3.2-1B	(1303%)	(1842%)	0.0148	(1848%)	(1188%)	(1633%)	(1602%)	(177%)	(115.2%)	(1131%)	0.0772	(1115%)
		0.9748	0.9216	0.5327	0.6107	0.6326	0.7345	0.9736	0.8459	0.6316	0.6768	0.6772	(+11.3 %) 0.7610
	LlaMA-3.2-3B	(+1.58%)	(+3.94%)	(+ 0.73%)	(+5.79%)	(+ 1.88%)	(+2.78%)	(+3.92%)	(+8.34%)	(+5.68%)	(+8.72%)	(+5.56%)	(+6.44%)
	L1aMA-3 3-70B	0.9534	0.8542	0.4533	0.6084	0.6123	0.6963	0.9315	0.7169	0.5106	0.6604	0.6438	0.6926
	Example 515 Fob	(- 0.56%)	(-2.80%)	(-7.21%)	(+5.56%)	(-0.15%)	(-1.03%)	(-0.29%)	(-4.56%)	(-6.42%)	(+7.08%)	(+2.22%)	(-0.39%)
					Dev	elopments							
	LlaMA-2-7B-chat-hf	0.9314	0.8601	0.6764	0.6858	0.6604	0.7628	0.9249	0.8129	0.7219	0.7028	0.6582	0.7641
	Vicuna-1.5-7B	0.9606	0.9720	0.6742	0.5445	0.5766	0.7456	0.9647	0.9424	0.6928	0.5659	0.6096	0.7551
Fine-tuning		(+ 2.92%)	(+11.1%)	(-0.22%)	(-14.1%)	(-8.38%)	(-1.72%)	(+3.98%)	(+12.9%)	(-2.91%)	(-13.6%)	(-4.86%)	(-0.91%)
	Wizardmath-7B	0.8680	0.8914	0.6902	(125%)	(12.2%)	0.7322	0.8845	0.8839	(1, 1, 7, 7, 7, 7, 7, 7, 7, 7, 7, 7, 7, 7, 7,	(1176%)	0.5649	0.7587
		0.8712	0.8686	(+1.30%)	0 7207	0.6179	0.7375	0.9268	0.9083	0.6849	0.7586	0.6668	0 7891
	Llama-2-7B-hf	(-6.02%)	(+ 0.85%)	(-6.71%)	(+3.49%)	(-4.25%)	(-2.53%)	(+0.19%)	(+9.54%)	(-3.70%)	(+5.58%)	(+ 0.86%)	(+2.49%)
	Choosed LloMA 1 2D	0.9268	0.9089	0.6990	0.8616	0.7411	0.8275	0.9864	0.9868	0.7582	0.9155	0.7953	0.8884
	Sneared-LiawA1.3B	(-0.46%)	(+4.88%)	(+2.26%)	(+17.5%)	(+8.07%)	(+6.47%)	(+6.15%)	(+17.3%)	(+3.63%)	(+21.2%)	(+13.7%)	(+12.4%)
. ·	Sheared-LlaMA1.3B-pruned	0.9966	0.9924	0.8346	0.9131	0.7516	0.8977	1.0000	1.0000	0.9345	0.9827	0.8686	0.9572
Pruning	.	(+6.52%)	(+13.2%)	(+15.8%)	(+22.7%)	(+9.12%)	(+13.4%)	(+7.51%)	(+18.7%)	(+21.2%)	(+ 27.9%)	(+21.0%)	(+19.3%)
	Sheared-LlaMA2.7B-pruned	(-477%)	(+3.76%)	$(\pm 0.54\%)$	(1898%)	(-7 69%)	(+0.16%)	$(\pm 4.44\%)$	(+14.7%)	(-0.51%)	$(\pm 14.4\%)$	(-3.60%-)	(±5 89%)
		0.9767	0.9649	0.7646	0.8799	0.6624	0.8497	0.9996	0.9999	0.8647	0.9480	0.7589	0.9142
	Sheared-LlaMA2.7B	(+4.53%)	(+10.4%)	(+8.82%)	(+19.4%)	(+0.20%)	(+8.69%)	(+7.47%)	(+18.7%)	(+14.2%)	(+24.5%)	(+10.0%)	(+15.0%)
		· · · · ·			· · · · · ·			· · · · · ·	· · · · · · · · · · · · · · · · · · ·				

Table 9: The detection performance (measured in AUROC) of RoBERTa-base and RoBERTa-large on EvoBench.

LLMs	Version Time/	Likelihood					Rank						
LLING	Version Name	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
					τ	Jpdates							
	2024-05-13	0.7820	0.8560	0.8033	0.5452	0.7356	0 7444	0.6596	0.7265	0 5969	0 5359	0.6840	0.6406
	2021 05 15	0.6980	0.8023	0.7717	0.5578	0.6823	0.7024	0.6751	0.7358	0.5920	0.5422	0.6643	0.6419
	2024-08-06	(-8.40%)	(-5.37%)	(-3.16%)	(+1.26%)	(-5.33%)	(-4.20%)	(+1.55%)	(+ 0.93%)	(- 0.49%)	(+ 0.63%)	(-1.97%)	(+0.13%)
GPT-40	2024-11-20	0.8764	0.8361	0.7752	0.5889	0.6910	0.7535	0.4497	0.7078	0.6022	0.5432	0.6640	0.5934
	2024-11-20	(+ 9.44%)	(-1.99%)	(-2.81%)	(+4.37%)	(-4.46%)	(+0.91%)	(-20.9%)	(-1.87%)	(+0.53%)	(+ 0.73%)	(-2.00%)	(-4.72%)
	Latest	0.8708	0.8846	0.7745	0.5498	0.7183	0.7596	0.4356	0.6846	0.5960	0.5584	0.6768	0.5903
		(+8.88%)	(+2.80%)	0.7716	(+0.40%)	(-1./3%)	(+1.52%)	0.7024	0.7420	(-0.09%)	(+2.25%)	0.6635	(-5.05%)
GPT-4o-mini	2024-07-18	(-5.16%)	(-0.47%)	(-3.17%)	(+3.94%)	(-4.15%)	(-1.80%)	(+4.28%)	(+1.64%)	(+0.12%)	(+0.41%)	(-2.05%)	(+0.88%)
	2022.06.13	0.7840	0.8302	0.8132	0.4306	0.6313	0.7015	0.6750	0.7207	0.6013	0.4415	0.6378	0.6154
	2023-00-15	0.7546	0.7389	0.7907	0.4666	0.7287	0.6959	0.6718	0.6640	0.6010	0.4851	0.6710	0.6186
	2023-11-06	(-2.94%)	(-10.0%)	(-2.25%)	(+2.70%)	(+9.74%)	(-0.56%)	(- 0.41%)	(-5.67%)	(-0.03%)	(+4.36%)	(+3.32%)	(+ 0.31%)
GPT-4	2024 01 25	0.6970	0.8465	0.7933	0.4902	0.7316	0.7117	0.6644	0.7181	0.5956	0.5064	0.6767	0.6322
	2024-01-25	(-8.70%)	(+0.73%)	(-1.99%)	(+ 5.06%)	(+10.0%)	(+1.03%)	(-1.15%)	(-0.26%)	(-0.57%)	(+ 6.49%)	(+3.89%)	(+1.68%)
	2024-04-09	0.6850	0.7897	0.7978	0.5264	0.6777	0.6953	0.6795	0.7041	0.6011	0.4868	0.6676	0.6278
		(-9.90%)	(-4.95%)	(-1.54%)	(+8.68%)	(+4.64%)	(-0.61%)	(+0.36%)	(-1.66%)	(-0.02%)	(+4.53%)	(+2.98%)	(+1.24%)
	2024-02-29	0.8796	0.9386	0.8372	0.4578	0.8621	0.7951	0.7069	0.7629	0.6062	0.4769	0.7132	0.6532
Claude Connect	2024-06-20	0.9123	0.9338	0.7980	0.5643	0.7927	0.8002	0.6962	0.7419	0.6055	0.5523	0.6988	0.6589
Claude-Sonnet		(+3.27%) 0.8453	(-0.48%) 0.8024	(-3.92%) 0.7476	(+10.6%)	(- 6.94%)	(+0.52%)	(-1.0/%) 0.6076	(-2.10%) 0.6054	(-0.07%) 0.6016	(+7.54%)	(-1.44%) 0.6044	(+0.57%)
	2024-10-22	(-3.43%)	(-13.6%)	(-8.96%)	(+0.06%)	(-16.4%)	(-8.48%)	(-0.93%)	(-6.75%)	(-0.46%)	(+5.36%)	(-1.88%)	(-0.93%)
	2024 03 07	0.0611	0.0821	0.8612	0.6122	0.8746	0.9595	0.7240	0.7049	0.6112	0.5925	0.7165	0.6979
Claude Haiku	2024-03-07	0.9611	0.9821	0.8013	0.0133	0.8746	0.8585	0.7340	0.7948	0.6115	0.5825	0.7165	0.6878
Claude-Haiku	2024-10-22	(-17.2%)	(-14.3%)	(-16.5%)	(-30.5%)	(+5.45%)	(-14.6%)	(-24.8%)	(-22.2%)	(-1.73%)	(-15.0%)	(-56.1%)	(-24.0%)
Cl	2024 02 20	0.9288	0.9581	0.8415	0.6053	0.8140	0.8295	0.6964	0.7513	0.6074	0.5458	0.7006	0.6603
Claude-Opus	2024-02-29	(-3.23%)	(-2.40%)	(-1.98%)	(-0.80%)	(-6.06%)	(-2.89%)	(-3.76%)	(-4.35%)	(-0.39%)	(-3.67%)	(-1.59%)	(-2.75%)
	Qwen1.5-7B	0.9613	0.9724	0.7250	0.5598	0.8303	0.8098	0.5473	0.7721	0.4532	0.4838	0.6862	0.5885
	Owen2-7B	0.9897	0.9881	0.6697	0.4810	0.8378	0.7933	0.7598	0.8230	0.4206	0.4603	0.6856	0.6299
Qwen	Qwell2-7B	(+2.84%)	(+1.57%)	(-5.53%)	(-7.88%)	(+ 0.75%)	(-1.65%)	(+21.2%)	(+ 5.09%)	(-3.26%)	(-2.35%)	(- 0.06%)	(+4.13%)
	Qwen2.5-7B	0.9834	0.9905	0.7348	0.4825	0.8221	0.8027	0.7359	0.8209	0.4551	0.4652	0.6843	0.6323
	-	(+2.21%)	(+1.81%)	(+0.98%)	(-7.73%)	(-0.82%)	(-0./1%)	(+18.8%)	(+4.88%)	(+0.19%)	(-1.86%)	(-0.19%)	(+4.38%)
	LlaMA-3.1-8B	0.9568	0.9577	0.8251	0.5852	0.9287	0.8507	0.7410	0.7805	0.6072	0.5506	0.7351	0.6829
	LlaMA-3.1-70B	0.9377	0.9754	0.7928	0.5636	0.9168	0.8373	0.7304	0.7894	0.5918	0.5403	0.7212	0.6746
LIaMA3		(-1.91%) 0.9406	(+1.//%)	(-3.23%) 0.8842	(-2.16%) 0.7317	(-1.19%)	(-1.34%) 0.8030	(-1.00%) 0.7616	(+ 0.89%)	(-1.54%)	(-1.03%)	(-1.39%)	(- 0.83%)
Liuini	LlaMA-3.2-1B	(-1.62%)	(+2.62%)	(+5.91%)	(+14.6%)	(-0.43%)	(+4.23%)	(+2.06%)	(+3.67%)	$(\pm 1.20\%)$	(+12.1%)	(-0.02%)	(+3.81%)
		0.9766	0.9729	0.8718	0.6604	0.9244	0.8812	0.7503	0.8090	0.6142	0.6173	0.7349	0.7051
	LIaMA-3.2-3B	(+1.98%)	(+1.52%)	(+4.67%)	(+7.52%)	(-0.43%)	(+3.05%)	(+ 0.93%)	(+2.85%)	(+ 0.70%)	(+ 6.67%)	(-0.02%)	(+2.23%)
	LlaMA-3.3-70B	0.9324	0.9741	0.7772	0.5779	0.7860	0.8095	0.7344	0.7873	0.5992	0.5592	0.6827	0.6726
		(-2.44%)	(+1.64%)	(-4.79%)	(-0.73%)	(-14.2%)	(-4.12%)	(- 0.66%)	(+0.68%)	(-0.80%)	(+ 0.86 %)	(-5.24%)	(-1.03%)
					Dev	elopments							
	LlaMA-2-7B-chat-hf	0.9604	0.9640	0.8568	0.8764	0.9433	0.9202	0.7120	0.7754	0.5840	0.6566	0.7868	0.7030
	Vicuna-1.5-7B	0.9961	0.9945	0.7624	0.7932	0.8948	0.8882	0.5429	0.7583	0.4262	0.4776	0.6820	0.5774
Fine-tuning		(+3.57%)	(+3.05%)	(-9.44%)	(-8.32%)	(-4.85%)	(-3.20%)	(-16.9%)	(-1.71%)	(-15.7%)	(-17.9%)	(-10.4%)	(-12.5%)
	Wizardmath-7B	0.836/	0.9352	0.6/85	0.7574	0.8253	0.8066	0.6458	(278%)	(15.0%)	0.6588	(11.2%)	(7.28%)
		0.6500	0.8027	0.5198	0 7768	0.7620	0 7023	0.6650	07413	0.4683	0 7340	0.7706	0.6758
	Llama-2-7B-hf	(-31.0%)	(-16.1%)	(-33.7%)	(-9.96%)	(-18.1%)	(-21.7%)	(-4.70%)	(-3.41%)	(-11.5%)	(+7.74%)	(-1.62%)	(-2.71%)
	Chaorad LioMA 1 2D	0.5791	0.7824	0.4899	0.7466	0.6532	0.6502	0.6724	0.7359	0.5050	0.7945	0.7018	0.6819
	Sheared-LiawiA1.3B	(-38.1%)	(-18.1%)	(-36.6%)	(-12.9%)	(-29.0%)	(-26.9%)	(- 3.96%)	(-3.95%)	(-7.90%)	(+13.7%)	(-8.50%)	(-2.10%)
Deres	Sheared-LlaMA1.3B-pruned	0.2053	0.6235	0.3066	0.5647	0.4081	0.4216	0.5174	0.7134	0.4448	0.6632	0.5871	0.5852
Pruning		(-75.5%)	(-34.0%)	(-55.0%)	(-31.1%)	(-53.5%)	(-49.8%)	(-19.4%)	(-6.20%)	(-13.9%)	(+ 0.66%)	(- 19.9%)	(-11.7%)
	Sheared-LlaMA2.7B-pruned	(-34.4%)	(-14.7%)	(-38.0%-)	(-14.2%)	(_28 3 %)	(-25.9%)	(-7.08%)	(-2.78%)	(-6 56%)	(±393%)	(-11.9%-)	(_4 89%)
		0.3425	0.6982	0.3254	0.6110	0.4911	0.4936	0.5199	0.7284	0.4564	0.6521	0.6174	0.5948
	Sheared-LlaMA2.7B	(-61.7%)	(-26.5%)	(-53.1%)	(-26.5%)	(-45.2%)	(-42.6%)	(-19.2%)	(-4.70%)	(-12.7%)	(-0.45%)	(-16.9%)	(-10.8%)

Table 10: The detection performance (measured in AUROC) of Likelihood and Rank on EvoBench.

LLMs	Version Time/			Log	Rank					Ent	ropy		
DLing	Version Name	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
					τ	Jpdates							
	2024-05-13	0 7804	0.8312	0 7888	0.5328	0 7234	0.7313	0.4602	0 3304	0.3091	0.5161	0 5264	0.4284
	2021 05 15	0.6969	0.7743	0.7607	0.5484	0.6670	0.6895	0.5289	0.3723	0.3269	0.4856	0.5467	0.4521
	2024-08-06	(-8.35%)	(-5.69%)	(-2.81%)	(+1.56%)	(-5.64%)	(-4.19%)	(+6.87%)	(+ 4.19%)	(+1.78%)	(-3.05%)	(+2.03%)	(+2.36%)
GPT-40	2024-11-20	0.8408	0.8010	0.7668	0.5725	0.6758	0.7314	0.1507	0.2264	0.3105	0.4423	0.5325	0.3325
	2024-11-20	(+ 6.04%)	(-3.02%)	(-2.20%)	(+ 3.97%)	(-4.76%)	(+0.01%)	(-30.9%)	(-10.4%)	(+0.14%)	(-7.38%)	(+ 0.61%)	(-9.60%)
	Latest	0.8299	0.8481	0.7605	0.5412	0.6972	0.7354	0.1458	0.1671	0.3216	0.5112	0.5505	0.3392
		(+4.95%)	(+1.69%)	(-2.83%)	(+0.84%)	(-2.62%)	(+0.41%)	(-31.4%)	(-16.3%)	(+1.25%)	(-0.49%)	(+2.41%)	(-8.92%)
GPT-4o-mini	2024-07-18	(-3.66%)	(-0.93%)	(-2.58%)	(+3.92%)	(-4.69%)	(-1.59%)	(+8.06%)	(+0.5555)	(+4.01%)	(-3.36%)	(+3.72%)	(+2.59%)
	2022.06.12	0.7925	0.8120	0.8048	0.4220	0.6275	0.6022	0.4676	0.2722	0.2194	0.5742	0.5799	0.4625
	2023-00-13	0.7823	0.8139	0.8048	0.4529	0.0273	0.6923	0.4070	0.3732	0.3124	0.5726	0.5788	0.4023
	2023-11-06	(-3.06%)	(-10.7%)	(-1.94%)	(+3.13%)	(+9.25%)	(-0.68%)	(+8.65%)	(+9.14%)	(-0.60%)	(-0.17%)	(-2.32%)	(+2.94%)
GPT-4	2024.01.25	0.7042	0.8260	0.7808	0.4832	0.7229	0.7034	0.5677	0.3238	0.2975	0.5511	0.5346	0.4549
	2024-01-25	(-7.83%)	(+1.21%)	(-2.40%)	(+5.03%)	(+9.54%)	(+1.11%)	(+10.0%)	(-4.94%)	(-2.09%)	(-2.32%)	(-4.42%)	(-0.75%)
	2024-04-09	0.6912	0.7563	0.7851	0.5158	0.6756	0.6848	0.5367	0.3414	0.3041	0.5134	0.5709	0.4533
	2021 01 05	(-9.13%)	(-5.76%)	(-1.97%)	(+ 8.29%)	(+4.81%)	(-0.75%)	(+ 6.91%)	(-3.18%)	(-1.43%)	(-6.09%)	(- 0.79%)	(-0.92%)
	2024-02-29	0.8854	0.9281	0.8322	0.4366	0.8678	0.7900	0.3882	0.2235	0.2986	0.5934	0.5103	0.4028
	2024-06-20	0.9047	0.9156	0.7927	0.5397	0.7974	0.7900	0.3372	0.2327	0.3312	0.5531	0.5516	0.4012
Claude-Sonnet	2021 00 20	(+1.93%)	(-1.25%)	(-3.95%)	(+10.3%)	(-7.04%)	(-0.00%)	(-5.10%)	(+ 0.92%)	(+ 3.26%)	(-4.03%)	(+4.13%)	(-0.16%)
	2024-10-22	0.8404	0.7813	0.7426	0.4459	0.6936	0.7008	0.3324	0.3329	0.3414	0.5526	0.5236	0.4166
		(-4.50%)	(-14.0%)	(-8.90%)	(+0.95%)	(-17.4%)	(-8.93%)	(-5.58%)	(+10.9%)	(+4.28%)	(-4.08%)	(+1.33%)	(+1.38%)
	2024-03-07	0.9675	0.9781	0.8583	0.5935	0.8824	0.8560	0.4216	0.2080	0.3073	0.5430	0.5176	0.3995
Claude-Haiku	2024-10-22	0.7687	0.7928	0.6859	0.3033	0.8854	0.6872	0.3788	0.2764	0.4005	0.6444	0.0585	0.3517
		(-19.8%)	(-18.5%)	(-17.2%)	(-29.0%)	(+0.30%)	(-16.8%)	(-4.28%)	(+6.84%)	(+9.32%)	(+10.1%)	(-45.9%)	(-4.78%)
Claude-Opus	2024-02-29	(-3.80%)	(-312%)	(-252%)	(-0.85%)	(-6.10%)	(-3.30%)	(-547%)	(-3.84%)	(± 0.3080)	(-0.97%)	(-0.5123)	(-2.14%)
	0.1670	0.0(21	0.0722	0.7105	0.5502	0.0270	0.0000	(3.47 %)	0.0005	(10.10%)	0.4227	(0.00 %)	0.2255
	Qwen1.5-/B	0.9631	0.9723	0.7135	0.5583	0.8379	0.8090	0.2406	0.2085	0.3407	0.4227	0.4648	0.3355
Owen	Qwen2-7B	$(\pm 3.00\%)$	$(\pm 1.63\%)$	(-3.88%)	(-7.32%)	(±1.13%)	(-1.09%)	$(\pm 4.98\%)$	(-4.20%)	$(\pm 12.8\%)$	$(\pm 7.95\%)$	(-1.16%)	$(\pm 4.09\%)$
X		0.9892	0.9905	0.7334	0.4858	0.8395	0.8077	0.3210	0.1523	0.4165	0.5007	0.4962	0.3773
	Qwen2.5-7B	(+2.61%)	(+1.82%)	(+1.99%)	(-7.25%)	(+ 0.16%)	(-0.13%)	(+8.04%)	(-5.62%)	(+7.58%)	(+7.80%)	(+3.14%)	(+4.19%)
	L1aMA-3 1-8B	0.9625	0.9513	0.8099	0.5759	0.9393	0.8478	0.4501	0.2765	0 3204	0 5331	0.4246	0.4009
		0.9498	0.9726	0.7823	0.5460	0.9267	0.8355	0.5097	0.2724	0.3192	0.5419	0.4308	0.4148
	LlaMA-3.1-70B	(-1.27%)	(+2.13%)	(-2.76%)	(- 2.99%)	(-1.26%)	(-1.23%)	(+5.96%)	(-0.41%)	(-0.12%)	(+ 0.88%)	(+0.62%)	(+1.39%)
LlaMA3	LIMA 2.2.1P	0.9423	0.9842	0.8808	0.7349	0.9418	0.8968	0.4843	0.2871	0.3330	0.5038	0.4826	0.4182
	LIaWA-3.2-1D	(-2.02%)	(+3.29%)	(+ 7.09%)	(+15.9%)	(+0.25%)	(+4.90%)	(+3.42%)	(+1.06%)	(+1.26%)	(-2.93%)	(+5.80%)	(+1.72%)
	LlaMA-3.2-3B	0.9798	0.9734	0.8629	0.6417	0.9418	0.8799	0.4505	0.2719	0.3169	0.5083	0.4826	0.4060
		(+1.73%)	(+2.21%)	(+5.30%)	(+6.58%)	(+0.25%)	(+3.21%)	(+0.04%)	(- 0.46%)	(-0.35%)	(-2.48%)	(+5.80%)	(+0.51%)
	LlaMA-3.3-70B	0.9514	(1104%)	(330%)	0.56/6	(13.7%)	(342%)	0.5398	(0.2738)	(10.5270)	0.5397	0.5801	0.4521
		(-1.11 //)	(+1.94 %)	(-3.39 %)	(-0.85 %)	(-13.7 %)	(-3.4270)	(+0.97 %)	(-0.27 %)	(+0.00 %)	(+0.00 %)	(+13.370)	(+3.11%)
					Dev	elopments							
	LlaMA-2-7B-chat-hf	0.9659	0.9636	0.8550	0.8860	0.9536	0.9248	0.2688	0.2043	0.3181	0.4155	0.3612	0.3136
	Vicuna-1.5-7B	0.9968	0.9937	0.7628	0.7992	0.9140	0.8933	0.3321	0.1489	0.3581	0.4209	0.4604	0.3441
Fine-tuning		(+3.09%)	(+3.01%) 0.0282	(-9.22%) 0.6720	(-8.68%) 0.7846	(-3.96%) 0.8257	(-3.15%)	(+0.33%)	(-5.54%) 0.2360	(+ 4.00%)	(+ 0.54%)	(+9.92%) 0.2002	(+3.05%)
	Wizardmath-7B	(-11.7%)	(-2.53%)	(-18.3%)	(-10.1%)	(-12.7%)	(-11.1%)	(+13.6%)	(+3.17%)	(+12.2%)	(+4.52%)	(+3.81%)	(+7.48%)
		0.6918	0.8354	0.5242	0.8093	0.7749	0.7271	0.4626	0.3657	0.4919	0.3611	0.4684	0.4299
	Llama-2-7B-hf	(-27.4%)	(-12.8%)	(-33.0%)	(-7.67%)	(-17.8%)	(-19.7%)	(+19.3%)	(+16.1%)	(+17.3%)	(-5.44%)	(+10.7%)	(+11.6%)
	Choosed LioMA1 2D	0.6329	0.8232	0.5013	0.7989	0.6832	0.6879	0.6036	0.4374	0.5483	0.5286	0.5848	0.5405
	Sneared-LiawA1.3B	(-33.3%)	(-14.0%)	(-35.3%)	(-8.71%)	(-27.0%)	(-23.6%)	(+33.4%)	(+23.3%)	(+23.0%)	(+11.3%)	(+22.3%)	(+22.7%)
. ·	Sheared-LlaMA1.3B-pruned	0.2612	0.6781	0.3235	0.6188	0.4304	0.4624	0.8336	0.5806	0.6668	0.6672	0.6402	0.6777
Pruning	pranou	(-70.4%)	(-28.5%)	(-53.1%)	(-26.7%)	(-52.3%)	(-46.2%)	(+56.4%)	(+ 37.6%)	(+34.8%)	(+25.1%)	(+ 27.9%)	(+36.4%)
	Sheared-LlaMA2.7B-pruned	0.6664	0.8491	0.5012	0.7650	0.6758	0.6915	0.5406	0.3740	0.5650	0.4914	0.5100	0.4962
	Sheared-LlaMA2.7B-pruned	(-29.9%) 0.4036	(-11.4%) 0.7400	(-35.5%) 0 3/29	(-12.1%) 0.6612	(-27.7%)	(-23.3%) 0.5218	(+27.1%) 0.7128	(+10.9%) 0.4800	(+24.0%) 0.6750	(+7.59%)	(+14.8%) 0.5094	(+18.2%)
	Sheared-LlaMA2.7B	(-56.2%)	(-22.3%)	(-51.1%)	(-22.4%)	(-44.3%)	(-39.3%)	(+44.5%)	(+27.5%)	(+35.6%)	(+19.0%)	(+23.7%)	(+30.1%)
		(001270)	((244 (0)	((0)	((0,0,0)	((), (0, (0))	(1410.00)	(1001070)	(0,0000)	(1.4011.70)	(10012.00)

Table 11: The detection performance (measured in AUROC) of LogRank and Entropy on EvoBench.

LLMs	Version Time/			Dete	ctGPT			LRR					
DLing	Version Name	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
					τ	Jpdates							
	2024-05-13	0.5182	0.8413	0.4455	0.4921	0 4643	0 5523	0.7032	0 7057	0.6719	0.4871	0.6213	0.6378
	2021 05 15	0.4785	0.7437	0.4558	0.5058	0.4450	0.5258	0.6655	0.6443	0.6601	0.5065	0.5744	0.6102
	2024-08-06	(-3.97%)	(-9.76%)	(+1.03%)	(+1.37%)	(-1.93%)	(-2.65%)	(-3.77%)	(-6.14%)	(-1.18%)	(+1.94%)	(-4.69%)	(-2.77%)
GPT-40	2024 11 20	0.2796	0.7283	0.4636	0.5469	0.4668	0.4970	0.6201	0.6486	0.6593	0.5000	0.5696	0.5995
	2024-11-20	(-23.8%)	(-11.3%)	(+1.81%)	(+5.48%)	(+0.25%)	(-5.52%)	(-8.31%)	(-5.71%)	(-1.26%)	(+ 1.29%)	(-5.17%)	(-3.83%)
	Latest	0.2796	0.6685	0.4704	0.5108	0.4539	0.4766	0.6139	0.6611	0.6480	0.5103	0.5879	0.6042
		(-23.8%)	(-17.2%)	(+2.49%)	(+1.8/%)	(-1.04%)	(-7.50%)	(-8.93%)	(-4.40%)	(-2.39%)	(+2.32%)	(-3.34%)	(-3.30%)
GPT-4o-mini	2024-07-18	(-1.92%)	(-7.01%)	(+0.19%)	(+ 0.87%)	(+2.59%)	(-1.06%)	(+2.07%)	(-1.39%)	(-1.50%)	(+3.65%)	(-3.95%)	(-0.22%)
	2022.06.12	0.2592	0.6976	0.2562	0.2699	0.2566	0.4255	0.7280	0.6061	0.6811	0.4211	0.5677	0.6208
	2025-00-15	0.5332	0.0870	0.4392	0.4385	0.3500	0.5333	0.7280	0.5723	0.6838	0.4747	0.5077	0.6121
	2023-11-06	(+18.8%)	(+5.89%)	(+8.29%)	(+6.97%)	(+13.8%)	(+10.7%)	(-2.60%)	(-12.3%)	(+0.27%)	(+4.36%)	(+5.98%)	(-0.87%)
GPT-4	2024 01 25	0.5225	0.7256	0.4187	0.4772	0.4818	0.5252	0.6974	0.7293	0.6573	0.4677	0.6338	0.6371
	2024-01-25	(+16.4%)	(+3.80%)	(+6.24%)	(+10.8%)	(+12.5%)	(+ 9.97%)	(-3.06%)	(+3.32%)	(-2.38%)	(+ 3.66%)	(+6.61%)	(+1.63%)
	2024-04-09	0.5240	0.7389	0.4031	0.4856	0.4526	0.5208	0.6788	0.6120	0.6684	0.4918	0.6037	0.6109
	2021 01 05	(+16.5%)	(+5.13%)	(+4.68%)	(+11.6%)	(+ 9.60%)	(+9.53%)	(-4.92%)	(-8.41%)	(-1.27%)	(+ 6.07%)	(+3.60%)	(-0.99%)
	2024-02-29	0.6532	0.8303	0.4666	0.3497	0.6006	0.5801	0.8454	0.8629	0.7312	0.3942	0.7993	0.7266
	2024-06-20	0.6551	0.8853	0.4738	0.4700	0.5718	0.6112	0.8236	0.8143	0.6860	0.4499	0.7341	0.7016
Claude-Sonnet	2021 00 20	(+ 0.19%)	(+ 5.50%)	(+ 0.72%)	(+12.0%)	(-2.88%)	(+3.11%)	(-2.18%)	(-4.86%)	(-4.52%)	(+5.57%)	(-6.52%)	(-2.50%)
	2024-10-22	0.5885	0.8112	0.5158	0.4866	0.5704	0.5945	0.7560	0.6939	0.6557	0.4221	0.6230	0.6301
		(-0.47%)	(-1.91%)	(+4.92%)	(+13.0%)	(-3.02%)	(+1.44%)	(-8.94%)	(-10.9%)	(-1.55%)	(+2.79%)	(-17.0%)	(-9.05%)
	2024-03-07	0.5728	0.7993	0.4760	0.4571	0.5511	0.5713	0.9439	0.9340	0.7639	0.5228	0.8174	0.7964
Claude-Haiku	2024-10-22	0.2622	0.4994	0.4924	0.3608	0.0210	0.3272	0.6454	0.5772	0.6065	0.3321	0.4288	0.5180
		(-31.0%)	(-29.9%)	(+1.64%)	(-9.63%)	(-53.0%)	(-24.4%)	(-29.8%)	(-35.0%)	(-15.7%)	(-19.0%)	(-38.8%)	(-27.8%)
Claude-Opus	2024-02-29	(-3.28%)	(-251%)	$(\pm 2.85\%)$	(± 0.4575)	(± 0.5582)	(-0.44%)	(-6.40%)	(-5.84%)	(-3.17%)	(-0.64%)	(-4.39%)	(-4.09%)
	0.1670	0.4416	0.7.107	(12.0070)	0.5055	0.((0)	0.5470	0.0200	0.0500	0 (102	0.52(0	0.7000	0.050(
	Qwen1.5-7B	0.4416	0.7407	0.3880	0.5055	0.6606	0.5473	0.9200	0.9508	0.6103	0.5260	0.7808	0.7576
Owen	Qwen2-7B	(1161%)	(1106%)	0.4380	(10.12%)	(372%)	(1373%)	0.9800	(1342%)	(1333%)	(264%)	(1131%)	(1242%)
Q		0.5188	0.6895	0.4255	0.5146	0.6406	0.5578	0.9826	0.9819	0.6673	0.4999	(+1.51%) 0.7989	0.7861
	Qwen2.5-7B	(+7.72%)	(-5.12%)	(+3.75%)	(+ 0.91%)	(-2.00%)	(+1.05%)	(+6.26%)	(+3.11%)	(+5.70%)	(-2.61%)	(+1.81%)	(+2.85%)
	LIOMA 2.1.8P	0.5006	0.7026	0.4803	0.5480	0.6603	0.5064	0.0435	0.0057	0.6842	0.5301	0.0031	0.7033
	LiawA-5.1-6D	0.5000	0.7920	0.4429	0.4987	0.0005	0.5904	0.9455	0.9037	0.6669	0.4895	0.8832	0.7935
	LlaMA-3.1-70B	(+9.27%)	(+4.34%)	(-3.74%)	(-4.93%)	(+1.62%)	(+1.31%)	(-1.22%)	(+3.59%)	(-1.73%)	(-4.06%)	(-1.99%)	(-1.08%)
LlaMA3	11-MA 2.2.1D	0.4201	0.7794	0.4736	0.6325	0.6274	0.5866	0.9443	0.9827	0.7754	0.6932	0.9235	0.8638
	LIAMA-3.2-IB	(-8.05%)	(-1.32%)	(-0.67%)	(+8.45%)	(-3.29%)	(-0.98%)	(+0.08%)	(+ 7.70%)	(+9.12%)	(+16.3%)	(+2.04%)	(+7.05%)
	L1aMA-3 2-3B	0.5017	0.8217	0.4506	0.5150	0.6274	0.5833	0.9725	0.9612	0.7428	0.5544	0.9235	0.8309
	Entrin 1 512 515	(+ 0.11%)	(+ 2.91%)	(-2.97%)	(-3.30%)	(-3.29%)	(-1.31%)	(+ 2.90%)	(+5.55%)	(+5.86%)	(+2.43%)	(+2.04%)	(+3.76%)
	LlaMA-3.3-70B	0.5910	0.8195	0.447/8	0.5073	0.5968	0.5925	0.9510	0.9366	0.7000	0.5253	0.7636	0.7753
		(+9.04%)	(+2.09%)	(-3.25%)	(-4.07%)	(-0.35%)	(-0.39%)	(+0.75%)	(+3.09%)	(+1.58%)	(-0.48%)	(-13.9%)	(-1.80%)
					Dev	elopments							
	LlaMA-2-7B-chat-hf	0.4919	0.7052	0.5017	0.7118	0.6650	0.6151	0.9579	0.9436	0.7791	0.8559	0.9300	0.8933
	Vicuna-1.5-7B	0.4796	0.6116	0.3794	0.6358	0.5986	0.5410	0.9863	0.9815	0.7012	0.7583	0.8787	0.8612
Fine-tuning		(-1.23%)	(-9.36%)	(-12.2%)	(-7.60%)	(-6.64%)	(-7.41%)	(+2.84%)	(+3.79%)	(-7.79%)	(-9.76%)	(-5.13%)	(-3.21%)
	Wizardmath-7B	0.5770	0.60/6	0.4606	0.6652	0.5814	0.5784	0.8240	0.9186	0.60/6	0.8084	0.7438	0.7805
		(+0.51%)	0.4501	(-4.11%)	(-4.00%) 0.6717	(-8.30%)	(-3.08%)	(-13.3%)	0.8604	0.5288	(-4.75%)	(-10.0%)	(-11.2%)
	Llama-2-7B-hf	(-1.03%)	(-25.5%)	(+1.80%)	(-4.01%)	(-11.2%)	(-8.00%)	(-19.5%)	(-7.42%)	(-25.0%)	(-1.41%)	(-18.5%)	(-14.3%)
		0.4112	0.3640	0.4510	0.5996	0.3168	0.4285	0.7537	0.8866	0.5291	0.8554	0.7000	0.7450
	Sheared-LlaMA1.3B	(-8.07%)	(-34.1%)	(-5.07%)	(-11.2%)	(-34.8%)	(-18.6%)	(-20.4%)	(-5.70%)	(-25.0%)	(-0.05%)	(-23.0%)	(-14.8%)
	Sheared I aMA1 3B propod	0.0297	0.0538	0.2284	0.3150	0.1406	0.1535	0.5272	0.7888	0.4092	0.7304	0.5318	0.5975
Pruning	Shearen-Liawirt i 3D-plutteu	(-46.2%)	(-65.1%)	(-27.3%)	(-39.6%)	(-52.4%)	(-46.1%)	(-43.0%)	(-15.4%)	(-36.9%)	(-12.5%)	(-39.8%)	(-29.5%)
	Sheared-LlaMA2.7B-pruned	0.5573	0.6327	0.4669	0.6359	0.4899	0.5565	0.9317	0.9638	0.5611	0.7985	0.6615	0.7833
	in the prime of the second sec	(+6.54%)	(-7.25%)	(-3.48%)	(-7.59%)	(-17.5%)	(-5.86%)	(- 2.62%)	(+ 2.02%)	(-21.8%)	(-5.74%)	(-26.8%)	(-11.0%)
	Sheared-LlaMA2.7B	0.4398	0.4220	0.4388	0.4282	0.2548	0.3967	0.7510	0.8904	0.5631	0.7493	0.5490	0.7006
		(-3.41%)	(-40.3%)	(-0.29%)	(-20.3%)	(-41.0%)	(-21.0%)	(-20.0%)	(-3.34%)	(-41.0%)	(-10.0%)	(-30.1%)	(-19.2%)

Table 12: The detection performance (measured in AUROC) of DetectGPT and LRR on EvoBench.

mmValueVa	LLMs	Version Time/ Version Name	NPR						DNA-GPT					
Update 2024-06-13 0.8971 0.9896 0.236 0.8971 0.9896 0.2361 0.276 0.576 0.2371 0.2761 0.5771 0.9991 0.8981 0.2281 0.2755 0.5176 0.5771 0.9991 0.8981 0.0282 0.0581 0.5781 0.5781 0.5781 0.5781 0.5781 0.5781 0.5785 0.5181 0.5789 0.5283 0.1783 0.6681 0.9752 0.5181 0.5789 0.5787 0.5181 0.5781 0.1781 0.5881 0.1781 0.5881 0.1781 0.5881 0.1781 0.5851 0.1781 0.5851 0.1485 0.5851 0.1485 0.5851 0.1485 0.5851 0.1485 0.5851 0.1485 0.5851 0.1485 0.5851 <th< td=""><td>Xsum</td><td>Writing</td><td>PubMed</td><td>SocialMedia</td><td>PeerRead</td><td>Avg.</td><td>Xsum</td><td>Writing</td><td>PubMed</td><td>SocialMedia</td><td>PeerRead</td><td>Avg.</td></th<>			Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.	Xsum	Writing	PubMed	SocialMedia	PeerRead	Avg.
3024-06.1 0.987 0.988 0.7346 0.9796 0.9745 0.9281 0.9766 0.9781 0.9876 0.9785 0.9784 0.9878 0.9875 0.9785 0.9786 0.9785 0.9786 0.9785 0.9786 0.9786 0.9785 0.9786 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9785 0.9795 0.5855 0.9795 0.5856 0.9795 0.5856 0.9795 0.5856 0.9795 0.5856 0.9795 0.5856 0.9795 0.5856 0.9795 0.5856 0.6855 0.9795 0.5856 0.5856 0.9795 0.5856 0.5856 0.9795 0.5866 0.5856 0.9795 0.5866 0.5856 0.9795 0.5866 0.5856 0.9796 0.5866 0.5856 0.9796 0.5866 0.5856 0.9796 0.5866 0.5856 0.9776 0.5876	Updates													
0297-00 03920 0.067 0.072 0.0747 0.0747 0.0747 0.0747 0.0747 0.0748 0.0223 0.0188 0.0235 (-1116) 0291-1100 (-12756) (-12166) (-1266) (-12676) (-126		2024-05-13	0.8971	0.9688	0.7346	0.5756	0.8256	0.8003	0 9945	0.8281	0.8276	0 5976	0.6781	0.7852
100-1600 (-2,57) (-2,178) (-2,178) (-0,178) (-2,178) (-1,178) (-2,178) (-1,178) (-2,178) (-1,178) (-2,178) (-1,178)	GPT-40	2024-08-06	0.8696	0.9447	0.6992	0.5932	0.7667	0.7747	0.9934	0.8608	0.8223	0.6188	0.6856	0.7962
GPT-lo 2024-11-20 C17.80 C0.999 C.0979 C.1378 C.0378 C.0			(-2.75%)	(-2.41%)	(-3.54%)	(+1.76%)	(-5.89%)	(-2.57%)	(-0.11%)	(+3.27%)	(-0.53%)	(+2.12%)	(+ 0.75%)	(+1.10%)
Lines: 06651 0095 0.7076 (-3.27%) (-0.09%) (-4.25%) (-6.55%) (-1.14%) (-0.14%) (-2.02%) (-1.25%) (-2.07%) (-2.0		2024 11 20	0.7202	0.8998	0.6969	0.5747	0.7831	0.7349	0.9858	0.7175	0.7868	0.6178	0.6634	0.7543
Lassi 0.0905 0.0915 0.7310 0.7332 </td <td>2024-11-20</td> <td>(-17.6%)</td> <td>(-6.90%)</td> <td>(-3.77%)</td> <td>(-0.09%)</td> <td>(-4.25%)</td> <td>(-6.54%)</td> <td>(-0.87%)</td> <td>(-11.0%)</td> <td>(-4.08%)</td> <td>(+2.02%)</td> <td>(-1.47%)</td> <td>(-3.09%)</td>		2024-11-20	(-17.6%)	(-6.90%)	(-3.77%)	(- 0.09%)	(-4.25%)	(-6.54%)	(- 0.87%)	(-11.0%)	(-4.08%)	(+2.02%)	(-1.47%)	(-3.09%)
Characterization Caluery		Latest	0.6963	0.9075	0.7010	0.5820	0.8242	0.7422	0.9866	0.7532	0.7993	0.6422	0.6728	0.7708
GPT-4e-mini 202407.18 0.0000 <th< td=""><td></td><td></td><td>(-20.0%)</td><td>(-0.13%)</td><td>(-3.30%)</td><td>(+0.64%)</td><td>(-0.14%)</td><td>(-5.81%)</td><td>(-0.79%)</td><td>(-7.49%)</td><td>(-2.83%)</td><td>(+4.40%)</td><td>(-0.53%)</td><td>(-1.44%)</td></th<>			(-20.0%)	(-0.13%)	(-3.30%)	(+0.64%)	(-0.14%)	(-5.81%)	(-0.79%)	(-7.49%)	(-2.83%)	(+4.40%)	(-0.53%)	(-1.44%)
GPT-4 2023-06-13 2024-11-06 (-0.5%) 0.9950 (-0.5%) 0.9171 (-0.5%) 0.9371 (-0.5%) 0.7519 (-0.5%) 0.9311 (-0.5%) 0.7525 (-0.5%) 0.9319 (-0.5%) 0.0535 (-0.5%) 0	GPT-4o-mini	2024-07-18	(+0.9002)	(+0.12%)	(-0.42%)	(+3.02%)	(-2.14%)	(+0.30%)	$(\pm 0.33\%)$	(+3.82%)	(-0.68%)	(+ 6.69 %)	(-3.96%)	$(\pm 1.24\%)$
Balesons	GPT-4	2022.06.12	0.0262	0.0776	0.7611	0.4027	0.7511	0.7910	0.0044	0.8725	0.9019	0.6562	0.7400	0.9229
GPT-4 222:11-06 (-0.85%) (-7.1%) (-9.2%) (-9.2%) (-8.1%) (-8.1%) (-1.2.4%) (-1.2.4%) (-3.7%) 2024-01-25 (-3.3%) (-3.2%) (-3.1%) (-5.3%) (-7.3%) (-5.0%) (-4.0%) (-3.7%) (-3.7%) (-5.0%) (-4.0%) (-8.4%) (-3.7%) (-5.0%) (-4.4%) (-3.7%) (-3.7%) (-5.0%) (-4.4%) (-3.7%) (-3.7%) (-5.0%) (-4.1%) (-1.3%) (-7.3%) (-5.0%) (-4.4%) (-3.7%) (-3.7%) (-3.7%) (-3.4%) (-7.3%)		2025-00-15	0.9202	0.9300	0.7317	0.5480	0.8282	0.7911	0.9944	0.7876	0.8037	0.6586	0.6259	0.7740
GPT-4 2024-01-25 0.3940 0.9315 0.0706 0.7380 0.7380 0.9210 0.9200 0.9300 0.6360 0.6365 0.6250 0.7385 0.6456 0.6250 0.7385 0.6300 0.6380 0.7380 0.9300 0.6300 0.7380 0.9310 0.6300 0.7380 0.9310 0.6300 0.7380 0.9310 0.6300 0.7380 0.9310 0.6300 0.7380 0.9310 0.6300 0.7380 0.73		2023-11-06	(-0.85%)	(-4.76%)	(-2.94%)	(+5.43%)	(+7.71%)	(+ 0.92%)	(-0.01%)	(-8.49%)	(-8.81%)	(+0.24%)	(-12.3%)	(-5.88%)
1000-05 (-3.22%) (-3.41%) (-4.23%) (-7.23%) (-5.03%) (-4.04%) (-8.04%) (-3.78) 2024-04-09 (-9.30%) (-6.57%) (-4.23%) (+2.37%) (-1.88%) (-0.65%) (-0.45%) (-7.37%) (-5.05%) (-1.48%) (-7.37%) (-7.37%) (-7.37%) (-1.64%) (-1.88%) (-7.37%) (-		2024 01 25	0.8940	0.9435	0.7036	0.5360	0.8245	0.7803	0.9921	0.8002	0.8409	0.6516	0.6622	0.7894
b b c		2024-01-25	(-3.22%)	(-3.41%)	(-5.75%)	(+4.23%)	(+7.34%)	(-0.16%)	(-0.23%)	(-7.23%)	(-5.09%)	(- 0.46%)	(-8.68%)	(-4.34%)
1000000000000000000000000000000000000		2024-04-09	0.8332	0.9120	0.7172	0.5764	0.7919	0.7661	0.9908	0.7880	0.8211	0.5970	0.6008	0.7595
Barbon Barbon Output Outpu Outpu Outpu <td>(-9.30%)</td> <td>(-6.56%)</td> <td>(-4.39%)</td> <td>(+8.27%)</td> <td>(+4.08%)</td> <td>(-1.58%)</td> <td>(-0.36%)</td> <td>(-8.45%)</td> <td>(-7.07%)</td> <td>(-5.92%)</td> <td>(-14.8%)</td> <td>(-7.32%)</td>			(-9.30%)	(-6.56%)	(-4.39%)	(+8.27%)	(+4.08%)	(-1.58%)	(- 0.36%)	(-8.45%)	(-7.07%)	(-5.92%)	(-14.8%)	(-7.32%)
Barborner Barborner <t< td=""><td rowspan="5">Claude-Sonnet</td><td>2024-02-29</td><td>0.9459</td><td>0.9804</td><td>0.8150</td><td>0.5836</td><td>0.9320</td><td>0.8514</td><td>0.9700</td><td>0.8582</td><td>0.8403</td><td>0.7355</td><td>0.6106</td><td>0.8029</td></t<>	Claude-Sonnet	2024-02-29	0.9459	0.9804	0.8150	0.5836	0.9320	0.8514	0.9700	0.8582	0.8403	0.7355	0.6106	0.8029
Claude-Some (+2.33%) (+4.98%) (+4.25%) (-1.76%) (-2.90%) (+1.61%) (-5.21%) (-1.61%) (-5.25%)		2024-06-20	0.9692	0.9884	0.7718	0.6605	0.9030	0.8586	0.9861	0.8061	0.8159	0.7250	0.5549	0.7776
2024-10-22 0.9016 0.9265 0.6769 0.9789 0.9789 0.9787 0.9548 0.7218 0.6748 0.6438 0.6		2024-00-20	(+2.33%)	(+ 0.80%)	(-4.32%)	(+ 7.69%)	(-2.90%)	(+ 0.72%)	(+ 1.61%)	(-5.21%)	(-2.44%)	(-1.05%)	(-5.57%)	(-2.53%)
(-4.37) (-5.39%) (-7.30%) (-7.30%) (-1.54%)		2024-10-22	0.9016	0.9265	0.6760	0.5056	0.7789	0.7577	0.9548	0.7238	0.7418	0.6681	0.5468	0.7271
Image: state			(-4.43%)	(-5.39%)	(-13.9%)	(-7.80%)	(-15.3%)	(-9.37%)	(-1.52%)	(-13.4%)	(-9.85%)	(-6.74%)	(-6.38%)	(-7.59%)
Chande-Haiku 2024-10-22 0.8396 0.9237 0.6439 0.3820 0.4744 0.6521 0.9955 0.8788 0.7293 0.7031 0.8217 Claude-Opus 2024-02-29 0.9703 0.9644 0.8234 0.7159 0.8474 0.8729 0.9939 0.8727 0.8165 0.7698 0.6188 0.8134 Claude-Opus 2024-02-29 0.9243 0.9945 0.8374 0.7150 0.8497 0.8969 0.9269 0.9034 0.7222 0.5074 0.4761 0.7726 Qwen1.5-7B 0.9292 0.9974 0.7376 0.4622 0.8989 0.8856 0.8856 0.8856 0.8856 0.8570 0.5755 0.77260 Qwen2.5-7B 0.9894 0.9955 0.7878 0.4666 0.9808 0.8356 0.8857 0.8987 0.9987 0.9387 0.8168 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256 0.5256	Claude-Haiku	2024-03-07	0.9952	0.9993	0.8614	0.7486	0.9389	0.9087	0.9972	0.9328	0.8388	0.7669	0.6745	0.8420
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		2024-10-22	0.8396	0.9237	0.6439	0.3892	0.4641	0.6521	0.9965	0.8578	0.7293	0.7031	0.8217	0.8217
Claude-Opus 09/00 09/00 08/24 0.8/14 0.8/19 0.89/19 0.99/19 0.8/17 0.818 0.6188 0.8188 0.8183 Qwen1.5-7B 0929 09309 0110 07907 05205 05205 07295 07285 07285 07289 07361 0520 <td< td=""><td></td><td>(-15.5%)</td><td>(-7.56%)</td><td>(-21.7%)</td><td>(-35.9%)</td><td>(-47.4%)</td><td>(-25.6%)</td><td>(-0.07%)</td><td>(-7.50%)</td><td>(-10.9%)</td><td>(-6.38%)</td><td>(+14.7%)</td><td>(-2.04%)</td></td<>			(-15.5%)	(-7.56%)	(-21.7%)	(-35.9%)	(-47.4%)	(-25.6%)	(-0.07%)	(-7.50%)	(-10.9%)	(-6.38%)	(+14.7%)	(-2.04%)
Provent (=2,49%) (=3,30%) (=3,30%) (=3,30%) (=4,35%) (=2,35%) (=2,35%) (=2,35%) (=2,35%) (=2,35%) (=2,35%) (=3,35%)	Claude-Opus	2024-02-29	0.9703	0.9604	0.8234	0.7159	0.8943	0.8729	0.9939	0.8727	0.8165	0.7698	0.6188	0.8143
Qwen1.5-7B 0.9243 0.0945 0.6661 0.4897 0.0976 0.7960 0.9269 0.9034 0.7222 0.5074 0.4071 0.7020 Qwen2.5-7B 0.9929 0.977 0.7357 0.4225 0.8168 0.8860 0.6929 0.5175 0.7260 Qwen2.5-7B 0.9984 0.9955 0.7878 0.4686 0.8987 0.8989 0.8778 0.468 0.8927 0.9987 0.8395 0.7784 (-1.35%) (+1.31%) (+1.45%) (+2.23%) LlaMA.3.1-8B 0.9957 0.9889 0.8176 0.6866 0.9897 0.8987 0.7861 0.7784 0.7222 0.6314 0.7899 LlaMA.3.1-70B 0.9957 0.9789 0.8101 0.6650 0.9400 0.9087 0.8937 0.7784 0.7225 0.7784 0.7280 0.7784 0.7825 0.7784 0.7825 0.7784 0.7825 0.7784 0.7825 0.7784 0.7825 0.7835 0.7863 0.8984 0.8986 0.8987 0.8987			(-2.49%)	(-3.89%)	(-3.80%)	(-3.21%)	(-4.40%)	(-3.58%)	(-0.33%)	(-0.01%)	(-2.23%)	(+0.29%)	(-5.57%)	(-2.77%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Qwen	Qwen1.5-7B	0.9243	0.9945	0.6641	0.4897	0.9076	0.7960	0.9269	0.9034	0.7222	0.5074	0.4761	0.7072
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		Qwen2-7B	0.9929	0.9974	0.7376	0.4622	0.8939	0.8168	0.8856	0.8692	0.8168	0.5209	0.5375	0.7260
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			(+0.80%) 0.0804	(+0.29%)	(+7.35%) 0.7878	(-2.75%) 0.4686	(-1.37%)	(+2.08%) 0.8282	(-4.13%)	(- 3.42%)	(+ 9.40%)	(+1.35%)	(+ 0.14%)	(+1.88%)
$ \begin{tabular}{ c c c c c c c c c c c c c c c c c c c$		Qwen2.5-7B	(+6.51%)	$(\pm 0.10\%)$	(+12.3%)	(-2.11%)	(-0.78%)	(+3.22%)	(-2.60%)	$(\pm 0.66\%)$	(+6.85%)	(+1.31%)	(+4.95%)	(+2.23%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	LlaMA3		(10027,0)	0.0000	(1121070)	0.000	(01/07/07)	(100270)	0.0007	(100070)	0.7704	((101%)	(1 100 10)	(12000)
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		LIAMA-3.1-8B	0.9926	0.9839	0.8476	0.6602	0.9680	0.8957	0.9987	0.8939	0.7784	0.7252	0.6304	0.8049
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		LlaMA-3.1-70B	$(\pm 0.31\%)$	$(\pm 1.50\%)$	(-372%)	(-173%)	(-0.9584)	(-0.92%)	(-0.19%)	(-1.32%)	(-1.57%)	$(\pm 1.29\%)$	(-1.22%)	(-0.60%)
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			0.9580	0.9948	0.9205	0.8617	0.9650	0.9400	0.9636	0.9378	0.8158	0.7512	0.7056	0.8348
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		LlaMA-3.2-1B LlaMA-3.2-3B	(-3.46%)	(+ 1.09%)	(+7.29%)	(+17.5%)	(-0.30%)	(+4.43%)	(-3.51%)	(+4.39%)	(+3.74%)	(+2.80%)	(+7.52%)	(+2.99%)
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $			0.9913	0.9870	0.9132	0.7698	0.9650	0.9253	0.9943	0.9210	0.8154	0.7655	0.6584	0.8309
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			(-0.13%)	(+0.31%)	(+6.56%)	(+8.32%)	(-0.30%)	(+2.95%)	(-0.44%)	(+ 2.71%)	(+3.70%)	(+4.23%)	(+2.80%)	(+ 2.60%)
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		LlaMA-3.3-70B	0.9958	1.0000	0.7923	0.6881	0.9042	0.8761	0.9955	0.8865	0.7878	0.7308	0.6533	0.8108
$ \begin{tabular}{ c c c c c c c } \hline $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $			(+0.32%)	(+1.61%)	(-5.53%)	(+0.15%)	(-6.38%)	(-1.97%)	(-0.32%)	(-0.74%)	(+0.94%)	(+0.76%)	(+2.29%)	(+0.59%)
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Developments													
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Fine-tuning	LlaMA-2-7B-chat-hf	0.9722	0.9904	0.9055	0.8963	0.9693	0.9467	0.6808	0.6330	0.4911	0.6396	0.6380	0.6165
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		Vicuna-1.5-7B	0.9958	0.9956	0.8008	0.8805	0.9439	0.9233	0.9745	0.9305	0.8064	0.6490	0.6355	0.7992
			(+2.36%)	(+ 0.52%)	(-10.4%)	(-1.58%)	(-2.54%)	(-2.34%)	(+29.3%)	(+29.7%)	(+31.5%)	(+ 0.94%)	(-0.25%)	(+18.2%)
$ \begin{array}{c} \text{Hummand} D & (-5.17\%) & (-1.00\%) & (-17.0\%) & (-4.73\%) & (-12.8\%) & (-8.14\%) & (+6.68\%) & (+12.5\%) & (+20.5\%) & (+4.61\%) & (-9.33\%) & (+7.01\%) \\ \text{Llama-2-7B-hf} & (-5.17\%) & (-0.0\%) & (-4.11\%) & (-5.5\%) & (-12.8\%) & (-6.64\%) & (-12.5\%) & (+20.5\%) & (+4.61\%) & (-9.33\%) & (+7.01\%) \\ \text{Llama-2-7B-hf} & (-0.0\%) & (-3.7.0\%) & (-44.1\%) & (-15.6\%) & (-31.1\%) & (-33.7\%) & (-4.65\%) & (+9.80\%) & (+0.46\%) & (+10.7\%) & (+13.3\%) & (+5.94\%) \\ \text{Sheared-LlaMA1.3B} & (-5.64\%) & (-4.11\%) & (-12.6\%) & (-3.11\%) & (-2.27\%) & (+8.88\%) & (+2.96\%) & (+9.45\%) & (+5.95\%) & (-10.2\%) & (+3.95\%) \\ \text{Sheared-LlaMA1.3B-pruned} & (-10.1\%) & (-10.1\%) & (-42.0\%) & (-2.27\%) & (-8.88\%) & (+2.27\%) & (+8.98\%) & (+2.0\%) & (+9.45\%) & (+10.2\%) & (+15.9\%) & (+10.2\%) & (+16.9\%) \\ \text{Sheared-LlaMA2.7B-pruned} & (-4.4\%) & (-8.53\%) & (-14.3\%) & (-14.3\%) & (-52.5\%) & (-14.3\%) & (-52.5\%) & (-14.3\%) & (-51.0\%) & (-35.7\%) & (+15.7\%) & (+2.8\%) & (+2.2.7\%) & (+14.3\%) & (+16.9\%) \\ \text{Sheared-LlaMA2.7B} & (-0.12.3\%) & (-53.9\%) & (-14.2\%) & (-44.9\%) & (-55.7\%) & (+15.7\%) & (+2.8\%) & (+3.9\%) & (+3.47\%) & (-7.2\%) & (+1.3\%) \\ \text{Sheared-LlaMA2.7B} & (-7.50\%) & (-40.0\%) & (-7.2\%) & (-31.7\%) & (-22.7\%) & (+6.8\%) & (+4.81\%) & (+8.83\%) & (+3.47\%) & (-17.2\%) & (+1.3\%) \\ \end{array}$		Wizardmath-7B Llama-2-7B-hf	0.9205	0.9804	0.7354	0.8490	0.8412	0.8653	0.7476	0.7582	0.6969	0.6857	0.5447	0.6866
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $			(-5.17%)	(-1.00%)	(-17.0%)	(-4.73%)	(-12.8%)	(-8.14%)	(+6.68%)	(+12.5%)	(+ 20.5%)	(+ 4.61%)	(-9.33%)	(+ 7.01%)
$ \begin{array}{c} -44.9\% \\ \text{Pruning} \end{array} \begin{array}{c} -44.9\% \\ \text{Sheared-LlaMA1.3B} \\ \text{Sheared-LlaMA1.3B} \\ \text{Sheared-LlaMA1.3B} \end{array} \begin{array}{c} -6.3.1\% \\ -0.554 \\ -0.571 \\ -0.56\% \\ -31.7\% \\ -10.1\% \\ -10.1\% \\ -5.6\% \\ -5.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2.4\% \\ -5.2\% \\ -5.2\% \\ -5.2\% \\ -5.2\% \\ -7.2\% $			0.5632	0.6202	0.4642	0.7395	0.6582	0.6091	0.6343	0.7310	0.4957	0.7474	0.7711	0.6759
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $			(-40.9%)	(-37.0%)	(-44.1%)	(-15.6%)	(-31.1%)	(-33.7%)	(-4.05%)	(+9.80%)	(+0.40%)	(+10.7%)	(+13.3%)	(+5.94%)
$ \begin{array}{c} \text{Pruning} \\ \text{Pruning} \\ \begin{array}{c} \text{Sheared-LlaMA1.3B-pruned} \\ \text{Sheared-LlaMA2.7B-pruned} \\ \text{Sheared-LlaMA2.7B} \\ \begin{array}{c} \text{Out} (10, 0, 10) \\ 0, 118 \\ 0, 007 \\ 0, 128 \\ 0, 100 \\$	Pruning	Sheared-LlaMA1.3B	(-31.7%)	(-10.1%)	(-42.0%)	(-3.29%)	(-26.2%)	(-22.7%)	(+8.98%)	(+2.96%)	(+9.45%)	(+8.59%)	(-10.2%)	(+3.96%)
Pruning Sheared-LlaMA1.3B-pruned		Sheared-LlaMA1.3B-pruned	0.4118	0.8077	0.3524	0.7530	0.4593	0.5568	0.9213	0.7903	0.6991	0.8669	0.6539	0.7863
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $			(-56.0%)	(-18.2%)	(-55.3%)	(-14.3%)	(-51.0%)	(-38.9%)	(+24.0%)	(+15.7%)	(+20.8%)	(+22.7%)	(+1.59%)	(+16.9%)
$\frac{(-52.4\%)}{(-26.7\%)} \left(\begin{array}{c} (-12.3\%) \\ (-53.9\%) \\ (-12.3\%) \\ (-53.9\%) \\ (-14.2\%) \\ ($		Channed LinMAA 7D 1	0.4480	0.8672	0.3664	0.7543	0.5196	0.5911	0.8383	0.7618	0.5870	0.7831	0.5455	0.7031
$ Sheared-LlaMA2.7B \begin{array}{cccccccccccccccccccccccccccccccccccc$		Silcared-LiawiA2./B-pruned	(-52.4%)	(-12.3%)	(-53.9%)	(-14.2%)	(-44.9%)	(-35.5%)	(+15.7%)	(+12.8%)	(+ 9.59%)	(+14.3%)	(-9.25%)	(+8.66%)
(-26.7%) (-7.50%) (-40.0%) (-7.72%) (-31.7%) (-22.7%) (+6.88%) (+4.81%) (+8.83%) (+3.47%) (-17.2%) (+1.35%)		Sheared-LlaMA2 7B	0.7049	0.9154	0.5047	0.8191	0.6520	0.7192	0.7496	0.6811	0.5794	0.6743	0.4655	0.6300
			(-26.7%)	(-7.50%)	(-40.0%)	(-7.72%)	(-31.7%)	(-22.7%)	(+6.88%)	(+4.81%)	(+8.83%)	(+3.47%)	(-17.2%)	(+1.35%)

Table 13: The detection performance (measured in AUROC) of NPR and DNA-GPT on EvoBench.