Watermark Smoothing Attacks against Language Models

Anonymous ACL submission

Abstract

Watermarking is a key technique for detecting AI-generated text. In this work, we study its vulnerabilities and introduce the *Smoothing Attack*, a novel watermark removal method. By leveraging the relationship between the model's confidence and watermark detectability, our attack selectively smoothes the watermarked content, erasing watermark traces while preserving text quality. We validate our attack on open-source models ranging from 1.3B to 30B parameters on ten different watermarks, demonstrating its effectiveness. Our findings expose critical weaknesses in existing watermarking schemes and highlight the need for stronger defenses.

1 Introduction

017

Detecting whether a text is generated by language models is critical in domains like fraud detection, fake news identification, and plagiarism prevention. A common approach is watermarking, where subtle patterns are embedded in the generated text for later detection (Aaronson, 2023; Christ et al., 2023; Huang et al., 2023; Li et al., 2024). Watermarking has gained traction in both academia and industry (Dathathri et al., 2024) as a key safeguard for language model applications. While various watermarking techniques exist, they share a core principle: favoring certain tokens over others (detailed in Section 2).

In this work, we identify key scenarios where watermarks fail and introduce a novel watermark removal attack that exploits this weakness, revealing fundamental limitations in existing watermarking schemes.

Effectiveness of watermarks. We say a watermark is effective if (i) the watermarked text maintains high quality, comparable to those generated from the corresponding un-watermarked model, and (ii) the detector reliably identifies watermark traces, i.e., it can identify watermarked text without making a large error. We analytically and empirically show that these aspects are in tension: better text quality often implies lower watermark detectability, and vice versa. Moreover, both are connected through the model's confidence in generating output. We explain the high-level idea as follows (see more detail in Section 3). 041

042

043

044

045

047

050

051

054

060

061

062

063

065

066

067

069

070

071

072

074

075

076

077

079

Given a prefix, when the model is confident about the output token, watermarking has negligible impact on the output. In this case, the watermark trace is not obvious. Conversely, when the model is not confident, watermarking makes the model tend to select certain tokens (that are originally unlikely to get sampled) over others, making watermark trace more detectable while degrading the text quality.

Smoothing Attack. Leveraging this insight, we propose the *Smoothing Attack* for watermark removal. For each prefix, the attack first identifies if the output token contains the watermark trace, by estimating the target watermarked model's confidence in this output. If the confidence is low, then we replace the token with a freshly sampled one (see more detail in Section 4), removing watermark traces while maintaining text quality; otherwise, if the confidence is high, then we retain the watermarked model's output.

We evaluate our attack across ten diverse watermarking schemes and three different families of open-sourced models, OPT (Zhang et al., 2022) (from 1.3B to 30B parameters), Llama3-8B (Dubey et al., 2024) and Qwen2-1.5B (Chu et al., 2024). In certain cases, our attack completely removes the watermark (reducing watermark detection rates to zero) while preserving the text quality. Our attack can also outperform the state-of-the-art *Paraphrasing Attack*, which uses the strong GPT-3.5-turbo to paraphrase the watermarked text. Compared with *Paraphrasing Attack*, our attack is more costefficient, as it uses only much weaker reference models, e.g., OPT-125M (Zhang et al., 2022) when attacking OPT models from 1.3B to 30B parameters. These findings underscore critical weaknesses in existing watermarks and highlight the need for more robust defenses.

081

090

093

094

095

097

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

2 Preliminaries and Related Work

Given an auto-regressive language model (LM) Mwith vocabulary \mathcal{V} , the model outputs a probability distribution over tokens at each position t in a prompt by computing logits $l_t(v)$ and applying a softmax:

$$P_t(v) = \frac{\exp(l_t(v))}{\sum_{v' \in \mathcal{V}} \exp(l_t(v'))}.$$
 (1)

To sample the next token, common strategies include top-k sampling (Fan et al., 2018; Holtzman et al., 2018), selecting from the top k tokens by probability, and top-p (nucleus) sampling (Holtzman et al., 2019), selecting from the smallest set whose cumulative probability exceeds p.

Watermarking schemes subtly modify this sampling to embed patterns in the generated text (v_1, \ldots, v_T) . These patterns are later detected via a scoring function $d(v_1, \ldots, v_T)$; if the score exceeds a threshold τ , the text is deemed watermarked.

Below we briefly review representative watermarking approaches.

Green-red watermark (Kirchenbauer et al., 2023a). For each position t, a secret key and the current prefix deterministically partition the vocabulary \mathcal{V} into a green list \mathcal{G}_t (size $\gamma |\mathcal{V}|$) and a red list. The logits of green tokens are then boosted by a constant δ , giving the modified distribution

$$\widetilde{P}_t(v) = \frac{\exp(l_t(v) + \delta \cdot \mathbf{1}\{v \in \mathcal{G}_t\})}{\sum_{v' \in \mathcal{V}} \exp(l_t(v') + \delta \cdot \mathbf{1}\{v' \in \mathcal{G}_t\})}.$$
(2)

The detector checks whether green tokens appear more frequently than expected, computing the score

$$d(v_1,\ldots,v_T) = \frac{\sum_{t=1}^T \left(\mathbf{1}\{v_t \in \mathcal{G}_t\} - \gamma\right)}{\sqrt{T\gamma(1-\gamma)}}.$$
 (3)

119The score is high when green tokens are overrepre-120sented, indicating a watermark.

121Gumbel and Tournament watermarks (Aaron-122son, 2023; Dathathri et al., 2024). These meth-123ods introduce randomness via a secret key and se-124lection process without altering the distribution

 P_t , preserving average text quality. In Gumbel sampling, noise values $u_t(v) \in [0,1]$ are generated using a seed derived from recent tokens and a secret key; the token v_t^* is selected as $v_t^* =$ $\arg \max_v - \frac{\log u_t(v)}{P_t(v)}$. The detector uses the score $d(v_1, \ldots, v_T) = -\sum_t \log(1 - u_t(v_t))$.

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

165

166

167

168

169

170

171

172

173

Tournament sampling similarly uses m secret functions $g^{(1)}, ..., g^{(m)}$ to score each token based on a seed r_t . The token is chosen through m rounds of pairwise comparisons among 2^m sampled candidates. Detection relies on the average tournament score:

$$d(v_1, \dots, v_T) = \frac{1}{T} \sum_{t=1}^T \frac{1}{m} \sum_{l=1}^m g^{(l)}(v_t, r_t).$$
 (4)

If the score significantly exceeds 0.5, the text is predicted as watermarked.

Other related work. The Green-red list watermark (Kirchenbauer et al., 2023a) introduces tokenlevel bias by amplifying the probabilities of greenlisted tokens and is thus considered distortionary. Variants differ in list construction and detection methods (Kirchenbauer et al., 2023b; Lee et al., 2023; Liu et al., 2023; Wu et al.). In contrast, Gumbel and Tournament sampling (Kuditipudi et al., 2023; Aaronson, 2023; Dathathri et al., 2024) preserve the token distribution in expectation and are distortion-free. Other distortion-free schemes include those by Hu et al.; Christ et al. (2023); see Zhao et al. (2024a) for a broader review. We evaluate 10 representative watermarks from both categories to demonstrate the generality of our attack.

Watermark removal techniques typically disrupt token patterns using homoglyphs, emojis, or control characters (Pajola and Conti, 2021; Boucher et al., 2022; Goodside, 2023), but often degrade fluency. A more effective approach is *paraphrasing*, where a separate model rewrites the text (Kirchenbauer et al., 2023b; Krishna et al., 2023; Piet et al., 2023). These attacks depend on strong paraphrasers—e.g., using GPT-3.5-turbo to rewrite outputs from smaller models like LLaMA-7B (OpenAI, 2023; Touvron et al., 2023).

Recent work explores adaptive paraphrasers trained to evade specific watermark detectors (Diaa et al., 2024), often requiring detailed knowledge of the watermarking algorithm. Other methods combine paraphrasing with auxiliary models to select high-quality rewrites from large candidate pools (Jovanović et al., 2024; Zhang et al., 2024).

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

174In contrast, our smoothing attack is lightweight:175it leverages only top-K token probabilities from176the watermarked model and a small reference177model—without needing a strong paraphraser,178large-scale sampling, or exact knowledge of the179watermarking scheme.

3 On the Effectiveness of Watermarks

180

181

182

183

185

186

187

189

190

191

192

193

194

195

196

199

200

201

203

We investigate two key aspects of watermark effectiveness: *detectability* and *text quality*. Our key finding is that these aspects are inter-connected via the model's confidence in prediction and are inherent tension—improving detectability typically decreases text quality, and vice versa. This tradeoff arises directly from how watermarking algorithms exploit token-level decisions, revealing a fundamental vulnerability leveraged by our proposed attack. Full derivations and further analysis are provided in Appendix C.

3.1 Token-level detectability

Detectability. Watermark detectors aggregate token-wise signals. For the Green–red scheme, the contribution of position t is

$$S_t = \underbrace{\mathbb{E}_{v \sim \widetilde{P}_t} [1\{v \in \mathcal{G}_t\}]}_{\text{watermarked}} - \underbrace{\mathbb{E}_{v \sim P_t} [1\{v \in \mathcal{G}_t\}]}_{\text{original}},$$

where P_t and \tilde{P}_t are the original and logit-shifted distributions. The detector score in Eq. equation 3 is just the normalised sum of these S_t .

Link to model confidence. With logit shift δ , S_t can be expressed solely through $\mathbb{E}_{v \sim P_t}[1\{v \in \mathcal{G}_t\}]$ (full derivation in Appendix):

$$S_t = \frac{(e^{\delta} - 1) \left(1 - \mathbb{E}_{v \sim P_t} [1\{v \in \mathcal{G}_t\}] \right)}{1 + (e^{\delta} - 1) \mathbb{E}_{v \sim P_t} [1\{v \in \mathcal{G}_t\}]}.$$
 (5)

The variance of $1\{v \in \mathcal{G}_t\}$ under P_t equals 204 $\gamma(1-\gamma) \|P_t\|^2$; hence departures of the expectation from its mean γ grow with the confidence measure $||P_t||^2 = \sum_v P_t(v)^2$. Figure 1 plots S_t versus 207 $||P_t||^2$ for 400 prefixes (OPT-1.3B, $\gamma = 0.5$, $\delta =$ 1.0). High-confidence locations ($||P_t||^2 \approx 1$) yield small S_t , whereas low-confidence ones $(||P_t||^2 \approx$ 210 $|\mathcal{V}|^{-1}$) maximise S_t . The same inverse relation 211 holds for Gumbel and Tournament watermarks (see 212 213 Fig. 7 in Appendix C.5). The take-away is that when the model is confident about a token, that out-214 put token leaves little trace for watermark detection 215 (i.e., difficult to detect); uncertainty amplifies the 216 watermark signal (i.e., easy to detect). 217

3.2 Text quality

Ultimately, watermarking should minimize its negative impact on downstream *text quality*. While we empirically assess quality via perplexity and diversity in Section 5, we first quantify how watermarking affects token-level distributional *fidelity*, as changes at this level directly influence downstream quality.

To measure fidelity loss, we use the total-variation distance:

$$D_{TV}(P_t, \widetilde{P}_t) = \frac{1}{2} \sum_{v \in \mathcal{V}} |P_t(v) - \widetilde{P}_t(v)|.$$

This metric is prompt- and task-agnostic, precisely capturing how watermarking alters token probabilities. Importantly, D_{TV} provides an upper bound on any smooth token-level objective—including perplexity and diversity—thus directly linking distributional fidelity to measurable text quality.

We evaluate fidelity loss under a *fixed secret key*, which deterministically partitions tokens (e.g., into green/red lists) based on the prefix. This matches the practical setting, where detectors must know the exact watermark key used during text generation. Link to model confidence. Figure 2 (left) shows that distortion shrinks as confidence $||P_t||^2$ rises: sharply peaked distributions are barely perturbed, while flat ones incur substantial fidelity loss. Gumbel and Tournament watermarks exhibit the same pattern (Appendix Fig. 8). In addition, if we focus on the part where $||P_t||^2$ is small, we note there is very little difference between the D_{TV} measured between the original model and the watermarked model and the the D_{TV} measured between the original model and the watermarked reference model. Effectively, that means replacing low-confidence tokens with samples from a small reference model harms fidelity no more (and often less, according to our experiments) than the watermark itself. Conversely, at high-confidence positions the reference model can be worse than the watermarked model, which, in turn, underscores the watermark's muted effect per se.

3.3 Detectability-quality trade-off

Combining the findings of the previous two subsections, we now answer the question: are the tokens that are easiest to detect are also those that most distort the distribution?

The answer is yes. In Figure 2 (right), we see that tokens that boost detectability scores (S_t) are



Figure 1: Correlation between the token-level watermark signal S_t , the original model's expected green-token rate $\mathbb{E}_{v \sim P_t}[1v \in \mathcal{G}_t]$, and model confidence $||P_t||^2$, measured on OPT-1.3B with the Green-red watermark ($\gamma=0.5$, ; $\delta=1.0$). Prefixes are drawn from the Harry Potter Wikipedia article. Corresponding results for Gumbel and Tournament watermarks are provided in Fig. 7 (Appendix).



Figure 2: Left: Token-level distributional shift $D_{TV}(P_t, \tilde{P}_t)$ vs. model confidence $||P_t||^2$ evaluated on OPT-1.3B. The blue plot measures the distributional shift due to the watermark distortion for Green–red watermark (with $\gamma = 0.5$, $\delta = 1.0$). The red plot measures the distributional shift between the reference model (OPT-125M) and the original model. **Right:** $D_{TV}(P_t, \tilde{P}_t)$ vs. token-level detectability S_t for the Green-red watermark. Lower confidence leads to larger fidelity loss and stronger watermark detectability. Gumbel and Tournament show identical patterns (see Fig. 9).

exactly those with the largest fidelity loss (D_{TV}) . Hence token-level watermarks cannot achieve high detectability without sacrificing distributional fidelity—and, by extension, downstreaming text quality. This inherent trade-off motivates our smoothing attack to be introduced in Section 4: by targeting at low-confidence positions, we can successfully remove the watermark signal while preserving the overall text quality.

4 Smoothing Attack

266

273

274

275

276

277

278

281

Our attack aims to remove watermarks while preserving text quality. At each token position, we first estimate the model's confidence and then selectively smooth low-confidence tokens to weaken the watermark detection signal. Detailed justification and derivations are provided in Appendix C; here we present the core algorithm clearly.

Adversary's model access. We assume a practical scenario: the adversary queries the target watermarked model via an API, obtaining only the top-K token probabilities for a given prefix (with $K \ll |\mathcal{V}|$). The original, un-watermarked model is not available to the adversary.

Estimating model confidence. We estimate model confidence at position t by approximating the squared ℓ_2 -norm of the token probability vector P_t :

$$\widehat{c}_t = \sum_{v \in \mathcal{V}_{\text{top}}} P_t(v)^2 + \frac{\left(1 - \sum_{v \in \mathcal{V}_{\text{top}}} P_t(v)\right)^2}{|\mathcal{V}| - K},$$
29

285

286

287

288

289

290

291

292

294

295

296

297

299

300

301

302

303

304

305

306

307

308

309

310

311

where \mathcal{V}_{top} denotes the set of top-K tokens returned by the model. Here we assume uniform probabilities among unobserved tokens. - For distortion-free watermarks (Gumbel, Tournament), the observed probabilities $P_t(v)$ directly reflect the original distribution. - For distortionary watermarks (Greenred), the probabilities from the watermarked model slightly deviate from the original. However, as shown in Appendix C, this approximation still correctly ranks tokens by their confidence; hence, it remains effective without additional correction.

Normalizing confidence scores. To convert the confidence estimate \hat{c}_t into a normalized confidence score $c_t \in [0, 1]$, we first establish empirical bounds L (lower) and U (upper). Specifically, we query the watermarked model using N random prefixes (e.g., N = 200) and record their \hat{c}_t values. Given these bounds, we set

$$c_t = \frac{\widehat{c}_t - L}{U - L}.$$
312

Smoothing procedure. Using a smoothing param-313eter $\alpha > 0$, our attack proceeds as follows at each314token position. 1) For *distortion-free watermarks*,315with probability c_t^{α} , we retain the token originally316produced by the watermarked model; otherwise,317we resample from the observed top-K probabili-318ties. 2) For *distortionary watermarks*, we query a319

323

324

325

327

331

333

336

338

340

341

342

347

349

351

364

367

small reference model (e.g., smaller than the watermarked model) using the same prefix to obtain its top-K token probabilities to construct a mixture distribution:

$$c_t^{\alpha} \cdot P_{\rm wm} + (1 - c_t^{\alpha}) \cdot P_{\rm ref},$$

where $P_{\rm wm}$ and $P_{\rm ref}$ represent the watermarked and reference model distributions, respectively. We then sample from this mixture distribution. Intuitively, large α favors keeping tokens from the watermarked model, while small α smooths the watermark more aggressively by replacing tokens more frequently.

Finally, if the adversary is uncertain about whether the watermark is distortion-free or distortionary, they simply follow the procedure designed for distortionary watermarks (using the referencemodel mixture) and the attack's success is not affected, as demonstrated in Section 5.

Efficiency and assumptions. Our attack is computationally efficient: it requires a modest initial overhead (a few hundred queries) to establish confidence-score normalization bounds (L, U), followed by one query per token position during generation. Importantly, our attack requires no knowledge of the watermark's algorithm, secret key, or detector internals, ensuring practical applicability.

5 Experiments

Setup. We evaluate our attack on three open-source model families: Llama3 (8B parameters) (Dubey et al., 2024), OPT (1.3B to 30B parameters) (Zhang et al., 2022), and Qwen2 (1.5B parameters) (Chu et al., 2024). When attacking distortionary watermarks, we use smaller models as references: Llama3-1B, OPT-125M, and Qwen2-0.5B, respectively.

Following prior work (Kirchenbauer et al., 2023a; Pan et al., 2024), we conduct evaluations on the C4 dataset (Raffel et al., 2020), where watermark performance has been shown effective. We avoid datasets with inherently low entropy (e.g., code generation), since previous studies (Kirchenbauer et al., 2023b; Lee et al., 2023) have demonstrated that watermark effectiveness significantly reduces in such settings. For each text in the dataset, the first 30 tokens form the prompt, and models generate the subsequent 200 tokens. All reported results are averaged over 100 prompts. Experiments were conducted using RTX-Titan GPUs.

We evaluate our attack against 10 representative watermarking algorithms covering both distortionary and distortion-free methods: KGW (Kirchenbauer et al., 2023a), Unigram (Zhao et al., 2023), UPV (Liu et al., 2023), X-SIR (He et al., 2024), DIP (Wu et al.), SWEET (Lee et al., 2023), EWD (Lu et al., 2024), Unbiased (Hu et al.), SynthID (Tournament sampling) (Dathathri et al., 2024), and Gumbel sampling (Aaronson, 2023). Our implementations build on the Mark-LLM toolkit (Pan et al., 2024). Note that Gumbel and X-SIR evaluations are restricted to OPT models, since Gumbel sampling exceeds 100 GB GPU memory requirements for Llama/Qwen due to their large vocabulary sizes, and X-SIR's official implementation currently supports only OPT models.

368

369

370

371

372

373

374

375

376

377

378

379

381

384

385

386

388

390

391

392

393

394

395

396

397

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

We compare our Smoothing Attack with the state-of-the-art *Paraphrasing Attack* (Piet et al., 2023), which employs GPT-3.5-turbo to rewrite texts, as well as its enhanced variants: paraphrasing multiple times and using GPT-40 as a stronger paraphraser.

Performance metric. We measure attack effectiveness along two dimensions: *watermark removal* and *text quality preservation*. For watermark removal, we report the true positive rate (TPR) of watermark detection, under a fixed false positive rate (FPR) of less than 1%. A lower TPR indicates better watermark removal (TPR is 1% for un-watermarked texts and 100% for fully watermarked texts without any attack). To evaluate text quality, we follow established protocols (Kirchenbauer et al., 2023a; Pan et al., 2024; Kirchenbauer et al., 2023b) by reporting perplexity (lower is better) and diversity (higher is better).

Unless otherwise noted, we set the smoothing parameter $\alpha = 1.0$ and use the top-10 tokens from both the watermarked and reference models. Additional experimental details and parameter variations are provided in the appendix.

Performance in watermark removal. Our main results are summarized in Tables 1 and 2. The key finding is that the *Smoothing Attack* effectively removes watermarks across diverse models and watermarking algorithms, consistently outperforming the strong paraphrasing attacks. Specifically, our attack achieves very low watermark detection rates (TPR around 5%, occasionally reaching 0%), whereas paraphrasing attacks perform notably worse. For instance, on OPT-1.3B with the Unigram watermark (Table 1), paraphrasing leaves

Table 1: Performance of watermark removal attacks on OPT-1.3B, Llama3-8B, and Qwen-1.5B models. We report the watermark true positive rate (TPR, lower is better), perplexity (PPL, lower is better), and diversity (Div, higher is better). All TPR values are measured at a fixed false positive rate below 1%. Additional results for models from 1.3B to 30B parameters are presented in Appendix B.6, demonstrating consistent trends.

Watermark	Attack		OPT-1.3	BB		Llama3-	8B	(Qwen2-1	.5B
		TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div
Un-watermarked	-	1	11.39	8.22	1	3.47	6.82	1	12.26	8.10
Reference	-	1	19.57	7.69	1	4.40	6.52	1	16.02	8.06
	None	100	14.61	8.07	99	4.60	6.92	100	16.46	8.11
KGW (Kirchenbauer et al., 2023a)	Paraphrasing	3	14.82	9.56	2	5.35	8.0	2	10.45	9.42
	Smoothing	0	9.57	6.72	2	3.20	5.63	0	8.02	6.91
	None	100	14.99	7.29	99	4.61	6.56	100	15.41	7.37
Unigram (Zhao et al., 2023)	Paraphrasing	53	14.51	8.75	54	5.60	8.02	5	10.40	8.56
	Smoothing	5	9.44	6.73	24	3.10	5.44	1	7.77	6.71
	None	100	7.12	7.41	99	4.83	7.31	100	6.94	7.05
SynthID (Dathathri et al., 2024)	Paraphrasing	1	10.57	9.11	1	5.62	8.18	1	6.90	8.43
	Smoothing	0	10.40	8.64	0	3.40	6.86	0	10.21	8.04
	None	100	13.73	8.44	84	4.03	7.35	100	14.34	8.27
DIP (Wu et al.)	Paraphrasing	0	13.95	9.25	0	5.25	8.34	2	10.10	8.85
	Smoothing	6	9.34	6.84	6	3.17	5.67	11	7.62	6.92
	None	100	13.61	8.29	84	4.02	7.29	100	14.64	8.21
Unbiased (Hu et al.)	Paraphrasing	3	14.45	10.39	2	5.36	8.57	1	9.97	8.82
	Smoothing	27	9.19	6.84	5	3.17	5.75	5	7.68	6.94
	None	99	11.65	8.22	83	4.38	6.80	86	11.93	7.49
UPV (Liu et al., 2023)	Paraphrasing	34	13.73	9.92	2	5.43	8.00	2	9.03	8.58
	Smoothing	20	10.01	6.89	1	3.12	5.49	0	8.16	6.91
	None	100	15.23	7.92	100	4.56	6.71	100	16.31	7.85
EWD (Lu et al., 2024)	Paraphrasing	0	14.95	9.95	7	5.73	7.83	1	10.18	9.28
	Smoothing	0	9.93	6.78	3	3.13	5.38	0	7.82	6.85
	None	100	14.36	8.02	99	4.53	6.69	100	15.89	7.65
SWEET (Lee et al., 2023)	Paraphrasing	0	14.57	9.45	14	5.64	8.05	4	10.18	9.30
	Smoothing	0	9.59	6.72	4	3.09	5.40	0	7.85	6.92

53% of texts detectable, whereas our smoothing reduces detection to just 5%, despite using only a much smaller OPT-125M reference model.

419

420

421

422

423

424

425

426

427

428

429

430

431

432

Moreover, our attack is computationally inexpensive and practical: watermarking text with KGW on OPT-1.3B requires about 4.2 seconds, while our smoothing attack using OPT-125M takes just 6.5 seconds (on two TITAN RTX GPUs). This brings us an important message: *existing watermarking schemes are vulnerable to even resourcelimited adversaries, underscoring the significant real-world applicability of the smoothing attack and fundamental limitations of existing watermark defenses.*

Performance in text quality. Our smoothing at-433 tack effectively preserves text quality, maintain-434 ing low perplexity (PPL) and competitive diver-435 sity (Div) while significantly improving watermark 436 removal. For instance, on OPT-1.3B with the Un-437 438 igram watermark (Table 1), our attack achieves notably better perplexity (9.44 vs. 14.51) and com-439 parable diversity (6.73 vs. 8.75) relative to the para-440 phrasing attack, while dramatically reducing wa-441 termark detection (TPR of 5% vs. 50%). Further 442

Table 2: Performance of watermark removal attacks on OPT-1.3B with Gumbel (Aaronson, 2023) and X-SIR (He et al., 2024) watermarks (with FPR < 1%).

Watermark	Attack	TPR (%)	PPL	Div
Gumbal	None	98 12	2.96	4.35
Guilloei	Smoothing	9	14.21	8.30
X-SIR	None Paraphrasing Smoothing	94 34 9	13.99 14.13 9.47	7.96 8.80 6.75

details on these quality metrics are provided in Appendix B.3 (Figures 4 and 5).

The effectiveness of our smoothing method is also evident for Gumbel sampling (Table 2). Although Gumbel sampling itself yields artificially low perplexity by generating repetitive content, it substantially reduces text diversity and overall quality. Our smoothing attack, by comparison, slightly increases perplexity but significantly reduces undesirable repetition, thereby improving actual text readability and coherence (examples in Appendix B.5, Table 14).

While our attack may sometimes show slightly

454

455

443

444

Table 3: Effectiveness of watermark removal attacks on OPT-1.3B under Unigram and UPV schemes. Translation can reduce TPR but significantly degrades quality. Repeated paraphrasing using GPT-3.5 lowers TPR further, with higher cost. GPT-40 improves fluency (lower PPL) but still leaves detectable watermark traces. Smoothing achieves the best overall trade-off by reducing TPR while preserving fluency and diversity.

Watermark	Attack	TPR(%)	PPL	Div
	None	100	15.0	7.3
	Translate (ZH)	100	22.2	7.6
	Translate (FR)	0	23.0	7.3
Unigram	GPT-3.5 Paraphrase (1x)	53	14.5	8.8
Ulligrafii	GPT-3.5 Paraphrase (2x)	0	15.7	9.9
	GPT-3.5 Paraphrase (3x)	0	15.6	9.9
	GPT-40 Paraphrase (1x)	0	11.2	8.4
	Smoothing	5	9.4	6.7
	None	99	11.7	8.2
	Translate (ZH)	58	17.0	8.0
	Translate (FR)	83	17.2	7.7
	GPT-3.5 Paraphrase (1x)	34	13.7	9.9
UFV	GPT-3.5 Paraphrase (2x)	23	15.1	9.9
	GPT-3.5 Paraphrase (3x)	24	14.7	10.1
	GPT-40 Paraphrase (1x)	51	9.8	8.7
	Smoothing	20	10.0	6.9

lower diversity compared to paraphrasing, this primarily arises because our method samples only from the top-K most likely tokens instead of the entire vocabulary. Increasing the value of K (see Table 5) effectively restores diversity, though it requires additional probability information from the model. Interestingly, by restricting token selection to the top-K tokens, our smoothing attack consistently achieves *even lower perplexity than unwatermarked texts*. Thus, beyond effectively removing watermarks, our method also enhances overall text quality by preventing the selection of extremely unlikely tokens commonly encountered in standard sampling.

456

457

458

459

460

461

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

Comparison with translation and repeated paraphrasing attacks. We further contextualize our smoothing attack by evaluating two additional attack types—translation-based and repeated paraphrasing attacks—particularly targeting watermark schemes (Unigram and UPV) resilient to singlepass GPT-3.5 paraphrasing.

For translation-based attacks, texts are translated to an intermediate language (Chinese or French) and then back to English using Google Translate (Google, 2024). For repeated paraphrasing, we apply GPT-3.5 up to three times iteratively, simulating more aggressive rewriting. Additionally,

Table 4: Comparison with the adaptive paraphraser from Diaa et al. (2024) on Llama3-8B with the Unigram watermark. The adaptive paraphraser is fine-tuned on Llama3-3B. Our attack use the same base, non-fine-tuned model as our reference model.

Attack	TPR (%)	PPL	Div
None	67	4.6	6.6
Adaptive Paraphraser	0	6.5	7.4
Smoothing	0	3.1	6.6

we explore GPT-40, a stronger paraphraser with enhanced fluency and coherence.

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

Table 3 demonstrates that translation-based attacks often significantly degrade text quality (higher perplexity) while only inconsistently removing watermarks. Repeated paraphrasing improves watermark removal but substantially increases computational cost. Paraphrasing with GPT-40—a stronger model than GPT-3.5 turbo—enhances text quality (lower perplexity) but does not guarantee better watermark removal; detection rates remain inconsistent and can even worsen, as observed with the UPV watermark.

Overall, these baselines emphasize the core advantage of our smoothing attack: it achieves robust watermark removal while preserving or even enhancing text quality, without relying on costly LMs. Comparison with adaptive paraphrasers. We compare our smoothing attack with the adaptive paraphraser proposed by Diaa et al. (2024), which relies on white-box knowledge of watermark algorithms, including the secret key generation and embedding mechanisms. In contrast, our smoothing attack operates purely under black-box assumptions, requiring no detailed watermark knowledge. Despite this weaker assumption, Table 4 demonstrates our smoothing attack achieves comparable or superior watermark removal performance, while significantly preserving text quality (lower PPL). Thus, our method offers practical effectiveness without demanding unrealistic adversarial knowledge.

Ablation studies on K, α , and model size. We extensively study the sensitivity of our smoothing attack to critical parameters (K and α) and the target model's size.

Increasing K typically improves watermark removal (lower TPR) and diversity, at a slight cost of increased perplexity (Table 5). Notably, even with a minimal setting (K = 1), our attack remains effective, achieving a TPR of only 18%, significantly lower than GPT-3.5 paraphrasing (53%). Increas-

K**TPR** (%) PPL Div α 0.5 42 9.9 6.86 1.0 5 9.44 6.73 Fixed to 10 2.0 0 9.38 6.58 3.0 1 9.25 6.43 3.21 1 18 4.62 5 746 10 611 Fixed to 1 10 5 9.44 6.73 5 15 11.73 7.11

Table 5: Impact of K and α on Smoothing Attack per-

formance on OPT-1.3B with Unigram watermark.

Table 6: *Smoothing Attack* against Unigram watermarking on models of *different sizes*, with OPT-125M as the reference model.

Target model size	Setting	TPR (%)	PPL	Div
1.3B	Unwatermarked	0	12.95	8.67
	Watermarked	99	16.53	7.29
	Smoothing	6	10.37	6.83
2.7B	Unwatermarked	0	11.75	8.36
	Watermarked	100	14.31	7.41
	Smoothing	4	10.35	6.66
6.7B	Unwatermarked	0	10.20	8.45
	Watermarked	100	12.94	7.48
	Smoothing	6	10.54	6.68
13B	Unwatermarked	0	10.14	8.39
	Watermarked	100	12.44	7.39
	Smoothing	5	10.32	6.70
30B	Unwatermarked	0	8.46	8.44
	Watermarked	100	10.45	7.56
	Smoothing	7	10.15	6.75

ing the smoothing parameter α makes our attack more aggressive in replacing uncertain tokens, thus enhancing watermark removal and generally improving perplexity. However, higher α values can slightly reduce diversity. By adjusting α , adversaries can effectively balance between watermark removal and text diversity (see Table 5). Additional results across diverse watermarks and models are presented in Appendix B.4.

524

525

527

529

531

533

534

535

536

539

Finally, we evaluate the effect of varying the target model size within the OPT family (from 1.3B to 30B), using the much smaller OPT-125M as the reference model. Our results (Table 5) indicate minimal sensitivity to the target model's size, confirming the scalability of our attack. Comprehensive results are detailed in Appendix B.6.

Effect of reference model size. To evaluate the
influence of reference model size, we apply our
smoothing attack using OPT models ranging from
125M to 1.3B parameters as the reference model.
Table 7 clearly indicates that using larger reference
models improves text quality (lower PPL, higher
diversity), and can further reduce watermark de-

Table 7: Impact of size of reference model size on the performance of the smoothing attack. Larger models reduce perplexity; lower TPR; and maintain diversity.

Watermark	Attack / Model Size	TPR (%)	PPL	Div
	None (Watermarked)	100	14.6	8.1
KCW	Smoothing (OPT-125M)	0	9.6	6.7
KUW	Smoothing (OPT-350M)	0	8.5	7.0
	Smoothing (OPT-1.3B)	0	7.0	7.3
	None (Watermarked)	100	13.7	8.4
סוס	Smoothing (OPT-125M)	6	9.3	6.8
DIP	Smoothing (OPT-350M)	9	8.2	7.0
	Smoothing (OPT-1.3B)	6	6.9	7.2
	None (Watermarked)	99	11.7	8.2
	Smoothing (OPT-125M)	20	10.0	6.9
UPV	Smoothing (OPT-350M)	5	8.9	7.2
	Smoothing (OPT-1.3B)	4	7.4	7.3

Table 8: Impact of watermark-type knowledge on smoothing attack against the distortion-free SynthID watermark. Smoothing attack reduces true positive rate to zero whether or not it knows the watermark type.

Setting	TPR (%)	PPL	Div
Watermarked (no attack)	100	7.1	7.4
Smoothing (type known)	0	10.4	8.6
Smoothing (type unknown)	0	9.6	6.7

tectability. These results confirm that the expressiveness of the reference model positively impacts overall attack performance.

Impact of watermark type knowledge. Our smoothing attack requires knowledge of the watermark type—specifically, whether it is distortion-free or distortionary—only for deciding whether and how to re-sample tokens. To assess the importance of this assumption, we evaluate our attack on the distortion-free SynthID watermark both with and without access to this information (Table 8). The results show that our attack achieves identical performance (TPR of 0%) in both settings, demonstrating its robustness and practical applicability even without watermark-type knowledge.

6 Conclusion

8

We revealed limitations in existing watermarks for language models and examined their robustness against watermark removal attacks. We introduced *Smoothing Attack*, a novel method that leverages model confidence to selectively remove watermark traces while preserving text quality. Comprehensive evaluations demonstrated that *Smoothing Attack* can completely remove watermarks, outperforming the state-of-the-art attack and highlighting a critical gap in current watermarks, and calling for more robust solutions.

560

561

562

563

564

565

566

567

568

569

570

571

572

573

547

580

582

583

585

589

590

591

592

7 Limitations and Ethical Considerations

575 In conducting this study, we have carefully consid-576 ered several factors that could influence the general-577 izability and applicability of our findings. Here we 578 outline these considerations explicitly, along with 579 the steps taken to address potential limitations.

Access to confidence-related information. Our smoothing attack relies on estimating model confidence from top-*K* token probabilities, commonly provided by public APIs (e.g., OpenAI's API). Although limiting access to probability information could theoretically mitigate our attack, in practice, such restrictions are challenging to enforce and may conflict with broader ethical goals around AI transparency and interpretability (OECD, 2019; National Institute of Standards and Technology (NIST), 2023). Recognizing this tension, we highlight the need for watermarking methods robust to scenarios where confidence estimates remain partially accessible.

Dependence on reference models. For distor-594 tionary watermark removal, our attack uses a ref-595 erence model to generate alternative token candidates. We explicitly evaluated different reference 597 model sizes (Table 7), confirming strong performance even when using significantly smaller reference models. However, selecting an extremely mismatched or lower-quality reference model could impact both watermark removal and text quality. We recommend using reference models carefully matched to the domain or distribution of the target 604 model.

Dataset selection and evaluation scope. We evaluated our methods extensively on the C4 dataset due to its well-documented suitability for watermark evaluation (Kirchenbauer et al., 2023a; Pan et al., 2024). Prior research indicated that datasets with low-entropy generation (e.g., code) already significantly reduce watermark effective-612 ness (Kirchenbauer et al., 2023b; Lee et al., 2023). 613 Thus, while our findings clearly establish the effec-614 tiveness of our attack under typical high-entropy 615 generation conditions, results may differ in specialized, low-entropy contexts. 617

Future watermarking approaches. Our attack
exploits inherent structural vulnerabilities shared
by current token-level watermark schemes. We explicitly acknowledge that future watermarking algorithms could be designed specifically to counter

such confidence-based attacks. Recognizing this potential evolution, we strongly encourage further research into watermark robustness and the development of methods resilient to confidence-based adversaries. 623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

8 Ethical Considerations

Our work demonstrates that an adversary, under realistic assumptions, can successfully remove watermarks from texts without compromising text quality. Although robustness concerns regarding watermarking have been highlighted by prior studies, our research underscores that these risks may be even greater than previously assessed.

We conducted experiments on LLama (Dubey et al., 2024), OPT (Zhang et al., 2022), and Qwen (Chu et al., 2024), each of which has been released under their respective licenses, as detailed in their documentation. Our implementation is based on MarkLLM (under the Apache License), and all modifications we have introduced are clearly documented in the README file included with our submitted code. Furthermore, any external packages used in our evaluations have been explicitly presented in the code we submitted.

Artifact use in this research has been consistent with intended purposes. Our dataset derives from the publicly available C4 dataset, which, to the best of our knowledge, does not contain personally identifiable information or offensive content. Additionally, relevant statistics regarding the data utilized in our experiments have been comprehensively reported C4 dataset (Raffel et al., 2020).

We utilized ChatGPT to revise portions of the manuscript; however, all revisions were performed under direct human supervision, ensuring that the final text accurately reflects our intent and ethical standards.

We approach this research with a firm commitment to ethical standards and responsible disclosure. By openly illustrating vulnerabilities, recommending effective mitigation, and transparently sharing our methods and outcomes, our objective is to inform and assist the broader research community. Our goal is to facilitate advancements in watermarking techniques that effectively balance transparency, innovation, and security, aligning with emerging regulatory standards such as the EU AI Act and the U.S. AI safety policies (European Commission, 2021). Moreover, we explicitly discuss potential defensive strategies and evaluate their ef-

ficacy (see Appendix D and Table10), providing 673 actionable guidance for enhancing the resilience of watermarking techniques. Overall, by clearly outlining these ethical considerations and limitations, we believe our research contributes robust and actionable insights, responsibly addressing the ethical implications and boundaries inherent in our study.

References

674

675

679

684

687

688

690

699

703

704

706

710

711

712

715

717

718

719

720

721

722

723

725

- Scott Aaronson. 2023. Simons institute talk watermarking of large language modon https://simons.berkeley.edu/talks/ els. scott-aaronson-ut-austin-openai-2023-08-17.
- Nicholas Boucher, Ilia Shumailov, Ross Anderson, and Nicolas Papernot. 2022. Bad characters: Imperceptible nlp attacks. In 2022 IEEE Symposium on Security and Privacy (SP), pages 1987-2004. IEEE.
- Tom Brown, Benjamin Mann, Nick Ryder, Deepak Subbiah, Jack Kaplan, Prafulla Dhariwal, A. Neelakantan, Long Ouyang, and Dario Amodei. 2020. Language models are few-shot learners. Advances in Neural Information Processing Systems. Seminal paper introducing GPT-3, discussing the importance of probabilities in understanding model behavior.
- Miranda Christ and Sam Gunn. 2024. Pseudorandom error-correcting codes. In Annual International Cryptology Conference, pages 325–347. Springer.
- Miranda Christ, Sam Gunn, and Or Zamir. 2023. Undetectable watermarks for language models. arXiv preprint arXiv:2306.09194.
- Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, and 1 others. 2024. Qwen2-audio technical report. arXiv preprint arXiv:2407.10759.
- Aloni Cohen, Alexander Hoover, and Gabe Schoenbach. 2024. Watermarking language models for many adaptive users. In 2025 IEEE Symposium on Security and Privacy (SP), pages 84–84. IEEE Computer Society.
- Sumanth Dathathri, Abigail See, Sumedh Ghaisas, Po-Sen Huang, Rob McAdam, Johannes Welbl, Vandana Bachani, Alex Kaskasoli, Robert Stanforth, Tatiana Matejovicova, and 1 others. 2024. Scalable watermarking for identifying large language model outputs. Nature, 634(8035):818-823.
- Abdulrahman Diaa, Toluwani Aremu, and Nils Lukas. 2024. Optimizing adaptive attacks against content watermarks for language models. arXiv preprint arXiv:2410.02440.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. arXiv preprint arXiv:2407.21783.

The EU Artifi-European Commission. 2021. Available at: https:// cial Intelligence Act. artificialintelligenceact.eu/. Proposed regulation focusing on transparency and accountability in high-risk AI systems.

726

727

728

729

730

731

732

733

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

757

759

760

762

763

764

765

766

767

768

769

770

771

772

774

775

776

777

778

779

- Jaiden Fairoze, Sanjam Garg, Somesh Jha, Saeed Mahloujifar, Mohammad Mahmoody, and Mingyuan Wang. 2023. Publicly detectable watermarking for language models. Cryptology ePrint Archive.
- 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.
- Jiayi Fu, Xuandong Zhao, Ruihan Yang, Yuansen Zhang, Jiangjie Chen, and Yanghua Xiao. 2024. Gumbelsoft: Diversified language model watermarking via the gumbelmax-trick. arXiv preprint arXiv:2402.12948.
- Surendra Ghentiyala and Venkatesan Guruswami. 2024. New constructions of pseudorandom codes. Cryptology ePrint Archive.
- Noah Golowich and Ankur Moitra. Edit distance robust watermarks via indexing pseudorandom codes. In The Thirty-eighth Annual Conference on Neural Information Processing Systems.
- Riley Goodside. 2023. There are adversarial attacks for that proposal as well - in particular, generating with emojis after words and then removing them before submitting defeats it. Twitter. URL: https://twitter.com/goodside/status/ 1610682909647671306.
- Google. 2024. Google Translate. https://translate. google.com. Accessed: 2024-05-01.
- Zhiwei He, Binglin Zhou, Hongkun Hao, Aiwei Liu, Xing Wang, Zhaopeng Tu, Zhuosheng Zhang, and Rui Wang. 2024. Can watermarks survive translation? on the cross-lingual consistency of text watermark for large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4115–4129, Bangkok, Thailand. Association for Computational Linguistics.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In International Conference on Learning Representations.
- Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, David Golub, and Yejin Choi. 2018. Learning to write with cooperative discriminators. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1638-1649.
- Zhengmian Hu, Lichang Chen, Xidong Wu, Yihan Wu, Hongyang Zhang, and Heng Huang. Unbiased watermark for large language models. In The Twelfth

781	International Conference on Learning Representa-	Aiwei Liu, Leyi Pan, Xuming Hu, Shiao Meng, and	835
782	tions.	Lijie Wen. 2024. A semantic invariant robust wa-	836
		termark for large language models. In The Twelfth	837
783	Baihe Huang, Banghua Zhu, Hanlin Zhu, Jason D Lee,	International Conference on Learning Representa-	838
784	Jiantao Jiao, and Michael I Jordan, 2023. Towards	tions	839
785	ontimal statistical watermarking arXiv preprint		000
786	arViv:2312.07930	Viiion Lu, Aiwai Liu, Dionzhi Vu, Jingjing Li, and Igwin	040
100	u/Alv.2512.07950.	I Jian Lu, Aiwei Liu, Dianzin Tu, Jingjing Li, and Ifwin	840
707	Mingile Hug. Soi Ashish Someyoiyle, Veyyyei Liong	King. 2024. An entropy-based text watermarking	841
/8/	Mingjia Huo, Sai Asilish Solhayajula, Touwei Liang,	detection method. arXiv preprint arXiv:2403.13485.	842
788	Ruisi Zhang, Farinaz Koushantar, and Pengtao		
789	Xie. Token-specific watermarking with enhanced de-	National Institute of Standards and Technology (NIST).	843
790	tectability and semantic coherence for large language	2023. Ai risk management framework. Technical	844
791	models. In Forty-first International Conference on	report, U.S. Department of Commerce. Framework	845
792	Machine Learning.	to improve AI system trustworthiness and manage	846
	·	risks emphasizing transparency	847
793	Nikola Jovanović, Robin Staab, and Martin Vechev.	nsks, emphasizing transparency.	041
794	2024. Watermark stealing in large language mod-	OECD 2010 Occid ai minainlas Available et lettras	0.40
795	els arXiv preprint arXiv:2402 19361	OECD. 2019. Oecd al principies. Available al. https:	848
100	c is. <i>arxiv</i> preprint arxiv.2702.19501.	//oecd.al/en/dashboards/al-principles/.	849
706	John Kirchenhauer Jonas Geining Vuvin Wen	Guidelines for ethical, trustworthy AI. Transparency	850
790	Jonathan Katz Jan Miars and Tom Coldstein 2022a	and accountability are key principles.	851
797	Johannan Katz, fair whers, and foin Goldstein. 2025a.		
798	A watermark for large language models. In Proceed-	OpenAI. 2023. Openai api documentation. Available	852
799	ings of the 40th International Conference on Machine	at: https://platform.openai.com/docs/. De-	853
800	Learning, volume 202 of Proceedings of Machine	veloper documentation highlighting the use of ton-K	854
801	Learning Research, pages 17061–17084. PMLR.	sampling beam search and probability outputs	955
	· · · ·	sampling, beam search, and probability outputs.	000
802	John Kirchenbauer, Jonas Geiping, Yuxin Wen, Manli	D. O. M. 2022 Cot 4 to being the state	0.50
803	Shu, Khalid Saifullah, Kezhi Kong, Kasun Fer-	R OpenAI. 2023. Gpt-4 technical report. arXiv	856
804	nando Aniruddha Saha Micah Goldblum and Tom	2303.08774. View in Article, 2(5).	857
005	Goldstein 2023b On the reliability of water		
000	mortes for large language models arViv mornint	Luca Pajola and Mauro Conti. 2021. Fall of giants:	858
806	marks for large language models. <i>arxiv preprim</i>	How popular text-based mlaas fall against a simple	859
807	arXiv:2300.04034.	evasion attack. In 2021 IEEE European Symposium	860
		on Security and Privacy (EuroS&P), pages 198–211.	861
808	Kalpesh Krishna, Yixiao Song, Marzena Karpinska,	IFFF	862
809	John Wieting, and Mohit Iyyer. 2023. Paraphras-		002
810	ing evades detectors of ai-generated text, but re-	Lavi Dan Aiwai Liu Zhiwai Ha Zitian Cao Yuandana	000
811	trieval is an effective defense. arXiv preprint	Zhao X ¹ ¹ ¹ ¹ ² L. D ¹ ¹ ² Zhao Oh 1 ¹ ² and L ¹ X and	863
812	arXiv:2303.13408.	Zhao, Yijian Lu, Binglin Zhou, Shuliang Liu, Xum-	864
		ing Hu, Lijie Wen, and I others. 2024. Marklim:	865
813	Rohith Kuditipudi, John Thickstun, Tatsunori	An open-source toolkit for llm watermarking. arXiv	866
814	Hashimoto and Percy Liang 2023 Robust	preprint arXiv:2405.10051.	867
915	distortion free watermarks for language models		
010	arVin proprint arVin 2207 15502	Julien Piet, Chawin Sitawarin, Vivian Fang, Norman	868
810	arxiv preprint arxiv:2507.15595.	Mu, and David Wagner, 2023. Mark my words: An-	869
		alyzing and evaluating language model watermarks	870
817	Taenyun Lee, Seokhee Hong, Jaewoo Ahn, Ilgee Hong,	arViv proprint arViv:2312 00273	071
818	Hwaran Lee, Sangdoo Yun, Jamin Shin, and Gunhee	urxiv preprini urxiv.2512.00275.	0/1
819	Kim. 2023. Who wrote this code? watermarking for	Culle D. C. I. N Channel A. Lee D. Lee a. K. (Lee)	070
820	code generation. arXiv preprint arXiv:2305.15060.	Colin Raffel, Noam Snazeer, Adam Roberts, Katherine	872
		Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	873
821	Xiang Li, Feng Ruan, Huiyuan Wang, Qi Long, and Wei-	Wei Li, and Peter J Liu. 2020. Exploring the lim-	874
822	iie J Su. 2024. A statistical framework of watermarks	its of transfer learning with a unified text-to-text	875
823	for large language models: Pivot detection efficiency	transformer. Journal of machine learning research,	876
82/	and optimal rules arXiv preprint arXiv:2404.01245	21(140):1–67.	877
024			
90E	Yiang Lieg Li Ari Haltzman Danial Eriad Dargy Liega	Hugo Touvron Louis Martin Kevin Stone Peter Al	272
620 000	Alang Lisa Li, Ali Holizinali, Daliel Fileu, Felcy Liang,	hart Amied Almahairi Veemina Dahaai Nikalay	070
020	Jason Eisner, Taisunori Hasnimoto, Luke Zettle-	Dashlukay County Dates Destinat Diagonal Chart	0/9
827	moyer, and Mike Lewis. 2022. Contrastive decoding:	Dasinykov, Soumya Batra, Prajjwal Bhargava, Shruti	088
828	Open-ended text generation as optimization. arXiv	Bnosale, and 1 others. 2023. Llama 2: Open foun-	881
829	preprint arXiv:2210.15097.	dation and fine-tuned chat models. arXiv preprint	882
		arXiv:2307.09288.	883
830	Aiwei Liu, Leyi Pan, Xuming Hu, Shuang Li, Lijie		
831	Wen, Irwin King, and S Yu Philip. 2023. An unforge-	Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Di-	884
832	able publicly verifiable watermark for large language	nan, Kyunghyun Cho, and Jason Weston. Neural text	885
833	models. In The Twelfth International Conference on	generation with unlikelihood training In Interna-	886
834	Learning Representations	tional Conference on Learning Representations	887
007	Learning Representations.	normi conjerence on Learning Representations.	007

987

939

John Wieting, Kevin Gimpel, Graham Neubig, and Taylor Berg-Kirkpatrick. 2022. Paraphrastic representations at scale. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 379–388.

891

892

897

901

902

903

904

907

908

909

910

911

912

913

914

917

918

919

921

922

923

924

925

926

927

928

931

932

934

935

938

- Yihan Wu, Zhengmian Hu, Junfeng Guo, Hongyang Zhang, and Heng Huang. A resilient and accessible distribution-preserving watermark for large language models. In *Forty-first International Conference on Machine Learning*.
- Hanlin Zhang, Benjamin L Edelman, Danilo Francati, Daniele Venturi, Giuseppe Ateniese, and Boaz Barak.
 2024. Watermarks in the sand: Impossibility of strong watermarking for language models. In *Fortyfirst International Conference on Machine Learning*.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, and 1 others. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- Xuandong Zhao, Prabhanjan Ananth, Lei Li, and Yu-Xiang Wang. 2023. Provable robust watermarking for ai-generated text. *arXiv preprint arXiv:2306.17439*.
- Xuandong Zhao, Sam Gunn, Miranda Christ, Jaiden Fairoze, Andres Fabrega, Nicholas Carlini, Sanjam Garg, Sanghyun Hong, Milad Nasr, Florian Tramer, and 1 others. 2024a. Sok: Watermarking for aigenerated content. *arXiv preprint arXiv:2411.18479*.
- Xuandong Zhao, Lei Li, and Yu-Xiang Wang. 2024b. Permute-and-flip: An optimally robust and watermarkable decoder for llms. *arXiv preprint arXiv:2402.05864*.
- Tong Zhou, Xuandong Zhao, Xiaolin Xu, and Shaolei Ren. 2024. Bileve: Securing text provenance in large language models against spoofing with bi-level signature. *arXiv preprint arXiv:2406.01946*.

A More on Related Work

Variations of Green-red list watermark. Different variations of Green-red list watermark, e.g., see (Kirchenbauer et al., 2023b; Lee et al., 2023; Liu et al., 2023; Wu et al.; Huo et al.; Zhou et al., 2024; Lu et al., 2024; Liu et al., 2024; He et al., 2024; Zhao et al., 2023; Kirchenbauer et al., 2023a), mainly differ in the assignment of the green lists and the detection process. In particular, the assignment of \mathcal{G}_t could depend on the prefix, e.g., the preceding *h* tokens in the generated text. When h = 0, we say the assignment is context-independent and is referred to as the *Unigram* watermark (Zhao et al., 2023); when h = 1, the assignment depends on the previous token and is referred to as the *KGW* watermark (Kirchenbauer et al., 2023a)

Scalable Tournament sampling. As shown in their paper, the original tournament process in (Dathathri et al., 2024) can be costly to implement, as there are $O(2^m)$ times of sampling and pair-wise comparison of tokens. Instead, they obtain a modified distribution for tokens. With $\tilde{P}_t^{(0)} = P_t$, they iteratively compute $\tilde{P}_t^{(l)}(v) = \left(1 + g^{(l)}(v, r_t) - \sum_{v' \in \mathcal{V}} \left(g^{(l)}(v', r_t) \cdot P_t^{(l-1)}(v')\right)\right) \cdot \tilde{P}_t^{(l-1)}(v)$, for $l = 1, \ldots, m$, and

then sample the token from $\widetilde{P}^{(m)}.$

Distortion-free watermark. There are also other distortion-free watermarks, which aim to preserve the original model's token distribution and avoid detectable shifts in probabilities of output tokens, e.g., see Hu et al.; Zhao et al. (2024b); Fu et al. (2024); Christ et al. (2023); Fairoze et al. (2023); Christ and Gunn (2024); Cohen et al. (2024); Ghentiyala and Guruswami (2024); Golowich and Moitra; Dathathri et al. (2024); Wu et al..

Comparison with paraphrasing attacks. When attacking OPT models (from 1.3B to 30B parameters), our attack only leverages the OPT-125M as the reference model when attacking distortionary watermarks such as the Unigram watermark. When attacking distortion-free watermarks, our attack sometimes resamples from the target watermarked model. In either case, the resource used in our attack is significantly smaller than the state-of-the-art paraphrasing attack, which uses the much larger GPT-3.5-turbo. Despite using fewer resources, our approach achieves higher watermark removal rates and comparable text quality. This highlights that even resource-limited adversaries can thwart watermarks, underscoring the need for stronger watermark defenses.

B More on Experiments

B.1 Implementation

We evaluate the smoothing attack on eight different watermarking algorithms, including KGW (Kirchenbauer et al., 2023a), Unigram (Zhao et al., 2023), SWEET (Lee et al., 2023), UPV (Liu et al., 2023), EWD (Lu et al., 2024), X-SIR (He et al., 2024), DIP (Wu et al.), Unbiased (Hu et al.), SynthID (Dathathri et al., 2024) and Gumbel (Aaronson, 2023). We use the implementations and default configurations provided by Mark-LLM (Pan et al., 2024). For completeness, we

Notation	Meaning	Definition Location
M	Auto-regressive language model (LM), which generates text sequentially based on a given prompt.	Section 2
\widetilde{M}	Watermarked model, a variant of M that embeds watermarks into generated text.	Section 3.1
\mathcal{V}	Vocabulary of the LM, the set of all possible tokens that can be generated.	Section 2
t	Token position in the generated sequence, indicating the index of a specific token.	Section 2
$l_t(v)$	Logit assigned by the model to token v at position t before applying softmax.	Eq. equation 1
$P_t(v)$	Probability of token v at position t after applying the softmax function.	Eq. equation 1
$\widetilde{P}_t(v)$	Modified probability distribution in the watermarked model after logit ma- nipulation.	Eq. equation 2
(v_1,\ldots,v_T)	Sequence of tokens forming the output text from the language model.	Section 2
$d(v_1,\ldots,v_T)$	Detection score function used to determine whether a text is watermarked.	Section 2
au	Threshold value for watermark detection; if $d(v_1, \ldots, v_T) > \tau$, the text is classified as watermarked.	Section 2
${\cal G}_t$	Green list, a subset of vocabulary containing tokens whose logits are in- creased in green-red list watermarking.	Section 2
γ	Fraction of the vocabulary included in the green list \mathcal{G}_t , determining the probability of token selection.	Section 2
δ	Logit increase applied to tokens in the green list \mathcal{G}_t , influencing token selection probabilities.	Eq. equation 2
T	Length of the generated sequence, i.e., the total number of tokens in the output text.	Section 2
$u_t(v)$	Randomly sampled value from $[0, 1]$ for token v in Gumbel sampling water-marking.	Section 2
v_t^*	Token selected using Gumbel sampling watermarking by maximizing a transformed probability.	Section 2
S_t	Contribution of the token at position t to the overall watermark detection score.	Eq. equation ??
$\mathbb{E}_{v \sim P_t}[1\{v \in \mathcal{G}_t\}]$	Expected probability mass assigned to green tokens at position t from probability distribution P_t .	Eq. equation ??
$\ P_t\ ^2$	\mathcal{L}_2 norm of the probability vector, measuring model confidence at position t. A higher value means greater confidence.	Section 3.1
$D_{TV}(P_t, \widetilde{P}_t)$	Total variation distance between original and watermarked probability distributions, measuring distortion.	Section 3.1
$D_{TV}(P_t, P_t^{\text{ref}})$	Total variation distance between the original model and a reference model's probability distributions.	Section 3.1
K	Number of most probable tokens that the adversary has access to from the watermarked model.	Section 4
$\mathcal{V}_{ ext{Top-K}}$	Set of top- K most probable tokens observed by the adversary.	Section 4
eta	Scaling factor used to estimate $ P_t ^2$ from watermarked probabilities in Green-red list watermarking.	Section 4
c	Normalized confidence score in $[0, 1]$ based on estimated \mathcal{L}_2 norm.	Section 4
U,L	Upper and lower bounds for normalizing \mathcal{L}_2 norms into the confidence score c .	Section 4
α	Exponential factor controlling the aggressiveness of the smoothing attack. A larger α favors keeping watermarked tokens, while a smaller α favors replacement.	Section 4
P_t^{ref}	Token probability distribution from a much smaller, un-watermarked reference model.	Section 4

Table 9: Table of notation definitions and their locations.

provide details of the algorithms below.

988

989

990 991 KGW (Kirchenbauer et al., 2023a): The green set G_t at each position t is selected based on the previous h tokens and a secret key known to the service provider. The hyperparameters are set as follows: $\gamma = 0.5$, $\delta = 2.0$, and h = 1.

992 993

994

995

• Unigram (Kirchenbauer et al., 2023a): The

green set \mathcal{G}_t is fixed for each token t and each prefix, depending solely on the secret key known to the service provider. No dynamic updates are performed based on previous tokens. The parameters are: $\gamma = 0.5$, $\delta = 2.0$.

996

997

998

1000

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1028

1029

1030

1031

1032

1034

1035

1036

1037

1038

1039

1040

- SWEET (Lee et al., 2023): A shift is applied only when the entropy of the probability distribution at position t is high, improving text quality, particularly for code generation tasks. The parameters are set as: γ = 0.5, δ = 2.0, the entropy threshold is 0.9, and h = 1.0.
- UPV (Liu et al., 2023): The green token selection process is similar to the previous approaches. However, this method requires training two additional models: a generator network to separate red and green tokens and a detector network for classification based on the input text. The watermarks are introduced using $\gamma = 0.5$, $\delta = 2.0$, and h = 1.0. The detector produces a binary prediction rather than a continuous score like a z-score.
- EWD (Lu et al., 2024): Watermark introduction follows a similar process as the previous methods. The hyperparameters are $\gamma = 0.5$, $\delta = 2.0$, and h = 1.0. During detection, tokens are assigned different weights based on their entropy, with higher entropy tokens receiving greater weight to improve detectability in low-entropy scenarios.
- X-SIR (He et al., 2024): Instead of operating at the token level, the red-green partition is applied at the level of semantic clusters, grouping similar words together and adding bias at the group level. This improves robustness against Cross-lingual Watermark Removal Attacks (CWRA).
- DIP (Wu et al.): Similar to Kirchenbauer et al. (2023), this method selects green tokens but uses a distribution-preserving reweight function to adjust token probabilities. This increases the probability of green tokens while maintaining the overall distribution. The reweighting is controlled by the parameter α. The hyperparameters are set as γ = 0.5, h = 5.
- 1041Implementation of the paraphrasing attack.1042We include the strongest baseline that paraphrases1043the given text based on the GPT-3.5-turbo (Piet

et al., 2023), denoted as P-GPT3.5 using the prompt: "Please rewrite the following text:". As shown in (Kirchenbauer et al., 2023b), GPT-3.5turbo is more powerful in removing the watermarks compared to Dipper model (Krishna et al., 2023).

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

Text quality metric. We use Llama3-8B, Qwen2-7B, and OPT-2.7B to evaluate the perplexity of the text generated from Llama3, Qwen2, and OPT models. We also report the log diversity of the text (Welleck et al.; Kirchenbauer et al., 2023b; Li et al., 2022), following the definition in (Kirchenbauer et al., 2023b) considering the 2-gram, 3-gram, and 4-gram repetition in the generated text. A higher diversity score represents a more diverse text.

B.2 Performance of the smoothing attack

Figure 3 shows three scatterplots of TPR vs. PPL for text generated under different watermarking and attack settings. Each point is colored by the watermarking method and corresponds to one of three models (OPT-1.3B, Llama3-8B, Qwen2-1.5B). Overall, the smoothing attack yields substantially lower TPR relative to the watermarked setting, demonstrating its performance at watermark removal. Notably, smoothing's TPR is on par with that of the paraphrasing attack, which uses a more powerful model (GPT-3.5-turbo). In terms of perplexity (PPL), smoothing also generates text that is competitive with (and sometimes lower than) both the watermarked text and the paraphrased text, indicating that it preserves text quality while removing the watermark.

B.3 Text Quality Evaluation

Figure 4 and Figure 5 present boxplots of the perplexity (PPL) and diversity of text generated from different sources using the OPT-1.3B model. We observe that the smoothing attack generally yields text with lower PPL than the watermarked model, except in cases involving the Gumbel watermark. This suggests that, according to the PPL metric, the smoothing attack can generate high-quality text. In terms of diversity, the constrained selection process—where sampling is restricted to the top-K candidates from both the reference and target models—results in lower diversity for the smoothing attack. These findings are consistent with the average PPL results reported in Table 1 in the main paper.

In addition, we compute the P-SP score (Wieting



Figure 3: Each subfigure shows how the true positive rate (TPR) varies with perplexity (PPL) for a specific attack. No attack (a) corresponds to watermarked text without modifications, paraphrasing (b) uses GPT-3.5-turbo to rewrite the text, and smoothing (c) randomly replaces some tokens to remove the watermark. Colors indicate the particular watermarking method and each point corresponds to one of three models (OPT-1.3B, Llama3-8B, Qwen2-1.5B).



Figure 4: Text Quality Comparison – Perplexity (OPT-1.3B). Box plots of perplexity for text generated from different sources, with perplexity computed using the OPT-2.7B model. Our smoothing attack produces text with quality comparable to, and in some cases better than, that of the watermarked model.



Figure 5: Text Quality Comparison – Diversity (OPT-1.3B). Box plots of text diversity for outputs generated from different sources. Our smoothing attack produces text with diversity comparable to, and in some cases lower than, that of the watermarked model due to its constrained selection process.

et al., 2022), which quantifies the similarity between pairs of texts in the embedding space, with higher scores indicating greater similarity. Specifically, we calculate P-SP scores for text generated from different sources and visualize the results in the heatmap shown in Figure 6. We observe that, aside from the paraphrasing case, texts from different sources generally exhibit low similarity. For instance, text generated by the watermarked model has a P-SP score of 53.6 on Unigram, whereas the similarity between the watermarked text and its paraphrased version reaches 82.3. Our smoothing

1094

1095

1096

1099

1100

1101

1102

1103

1104

attack produces a P-SP score (measuring similarity between text from the smoothing attack and unwatermarked text) comparable to that of the watermarked text (measuring similarity between watermarked text and unwatermarked text). The generally low P-SP scores across different sources reflect the natural variability in generated responses, as multiple reasonable outputs can exist for the same prompt. Therefore, P-SP metrics may not be a reliable measure for assessing text quality degradation due to watermarking or smoothing.

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

		KG	W				Unig	gram				Synt	thID		
Unwatermarked	100 5	55.9	54.9	45.9	55.0	100	53.6	54.0	43.9	55.4	100	56.3	54.0	46.8	63.0
Watermarked	55.9	100	53.7	44.6	53.3	53.6	100	52.0	82.3	54.2	56.3	100	54.1	85.0	56.3
Reference	54.9 5	53.7	100	43.3	55.9	54.0	52.0	100	42.3	56.1	54.0	54.1	100	44.3	54.0
Paraphrasing				100	43.4	43.9	82.3	42.3	100	43.2	46.8	85.0	44.3	100	
Smoothing	55.0 5	53.3	55.9	43.4	100	55.4	54.2	56.1	43.2	100	63.0	56.3	54.0	47.5	100
		D	IP				Unbi	iased				XS	IR		
Unwatermarked	100 5	57.0	54.0	47.2	54.5	100	56.6	54.0	46.8	54.5	100	54.5	54.0	44.8	54.8
Watermarked	57.0	100	54.8	83.7	55.2	56.6	100	54.0	84.1	53.5	54.5	100	53.4	83.7	53.3
Reference	54.0 5	54.8	100	44.6	56.4	54.0	54.0	100	44.2	55.7	54.0	53.4	100	43.4	56.5
Paraphrasing	47.2 8	83.7		100	44.1	46.8	84.1	44.2	100	42.5	44.8	83.7	43.4	100	
Smoothing	54.5 5	55.2	56.4	44.1	100	54.5	53.5	55.7	42.5	100	54.8	53.3	56.5	42.1	100
	Unwatermarked	Watermarked	Reference	Paraphrasing	Smoothing	Unwatermarked	Watermarked	Reference	Paraphrasing	Smoothing	Unwatermarked	Watermarked	Reference	Paraphrasing	Smoothing

Figure 6: Text Quality Comparison - P-SP (OPT-1.3B). Heatmap comparing the similarity of text generated by different models in the sentence embedding space. Text from the watermarked model has a low similarity score compared to unwatermarked text, reflecting the inherent variability in generated responses. However, the paraphrased text (Paraphrasing vs. watermarked) exhibits a high similarity score, suggesting that the P-SP metric is more suitable for evaluating paraphrasing rather than assessing text quality degradation due to watermarking or smoothing.

Table 10: Effect of K on Smoothing Attack Performance (OPT-1.3B). Evaluation of the smoothing attack	's
effectiveness against different watermarking algorithms on the OPT-1.3B model, varying the number of top-	K
tokens accessible to the attacker.	

Κ		KGW		τ	Jnigram		SynthID			DIP			Unbiased		
	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div
1	9%	3.22	4.54	18%	3.21	4.62	0%	10.45	8.47	1%	3.36	4.57	7%	3.36	4.56
3	0.0%	5.76	5.71	8.0%	5.9	5.68	0.0%	10.5	8.31	4.0%	5.58	5.66	14.0%	5.59	5.68
5	2.0%	7.27	6.17	10.0%	7.46	6.11	0.0%	10.35	8.71	3.0%	6.97	6.23	19.0%	7.11	6.29
7	1.0%	8.14	6.46	5.0%	8.48	6.55	0.0%	10.42	8.63	7.0%	7.97	6.46	29.0%	8.06	6.47
10	0.0%	9.57	6.72	5.0%	9.44	6.73	0.0%	10.4	8.64	6.0%	9.34	6.84	27.0%	9.19	6.84
	XSIR UPV														
K		XSIR			UPV			Gumbel			EWD		S	WEET	
K	TPR	XSIR PPL	Div	TPR	UPV PPL	Div	TPR	Gumbel PPL	Div	TPR	EWD PPL	Div	TPR	WEET PPL	Div
К 1	TPR 14%	XSIR PPL 3.31	Div 4.5	TPR 22%	UPV PPL 3.62	Div 4.63	TPR 0%	Gumbel PPL 20.8	Div 8.2	TPR 1%	EWD PPL 3.31	Div 4.49	TPR 2%	WEET PPL 3.41	Div 4.57
К 1 3	TPR 14% 17.0%	XSIR PPL 3.31 5.69	Div 4.5 5.52	TPR 22% 14.0%	UPV PPL 3.62 6.22	Div 4.63 5.87	TPR 0% 2.0%	Gumbel PPL 20.8 21.72	Div 8.2 8.47	TPR 1% 0.0%	EWD PPL 3.31 5.78	Div 4.49 5.71	TPR 2% 0.0%	WEET PPL 3.41 5.64	Div 4.57 5.75
K 1 3 5	TPR 14% 17.0% 8.0%	XSIR PPL 3.31 5.69 6.8	Div 4.5 5.52 6.04	TPR 22% 14.0% 16.0%	UPV PPL 3.62 6.22 7.68	Div 4.63 5.87 6.3	TPR 0% 2.0% 8.0%	Gumbel PPL 20.8 21.72 20.3	Div 8.2 8.47 8.23	TPR 1% 0.0% 0.0%	EWD PPL 3.31 5.78 7.32	Div 4.49 5.71 6.18	TPR 2% 0.0% 0.0%	WEET PPL 3.41 5.64 7.15	Div 4.57 5.75 6.23
K 1 3 5 7	TPR 14% 17.0% 8.0% 10.0%	XSIR PPL 3.31 5.69 6.8 8.26	Div 4.5 5.52 6.04 6.48	TPR 22% 14.0% 16.0% 7.0%	UPV PPL 3.62 6.22 7.68 8.75	Div 4.63 5.87 6.3 6.65	TPR 0% 2.0% 8.0% 9.0%	Gumbel PPL 20.8 21.72 20.3 21.15	Div 8.2 8.47 8.23 8.15	TPR 1% 0.0% 0.0% 0.0%	EWD PPL 3.31 5.78 7.32 8.65	Div 4.49 5.71 6.18 6.52	TPR 2% 0.0% 0.0% 0.0%	WEET PPL 3.41 5.64 7.15 8.47	Div 4.57 5.75 6.23 6.45

B.4 Effect of K and α

Table 10 and Table 12 show the performance of 1117 smoothing attacks against different watermarking algorithms under varying values of K. In a smoothing attack, the adversary has access only to the top-K tokens and their probabilities from both the

Table 11: Effect of α on Smoothing Attack Performance (OPT-1.3B). Evaluation of the smoothing attack's effectiveness against different watermarking algorithms on the OPT-1.3B model, varying the parameter α . A larger α indicates that the attack relies more on the reference model's output, while a smaller α means the attack is more influenced by the watermarked text.

α	KGW			Unigram			SynthID			DIP			Unbiased		
	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div
0.5	11.0%	10.03	7.02	42.0%	9.9	6.86	2.0%	9.33	7.9	29.0%	9.27	7.11	63.0%	8.92	7.09
1.0	0.0%	9.57	6.72	5.0%	9.44	6.73	0.0%	10.4	8.64	6.0%	9.34	6.84	27.0%	9.19	6.84
2.0	0.0%	9.35	6.65	0.0%	9.38	6.58	0.0%	11.16	8.26	1.0%	9.03	6.71	9.0%	8.89	6.59
3.0	0.0%	9.45	6.46	1.0%	9.25	6.43	0.0%	11.33	8.61	0.0%	9.32	6.82	1.0%	9.05	6.65
α	X-SIR UPV				Gumbel			EWD			SWEET				
	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div
0.5	28.0%	9.47	6.94	42.0%	10.01	7.14	80.0%	13.73	7.54	0.0%	9.76	7.01	6.0%	9.66	7.13
1.0	9.0%	9.47	6.75	20.0%	10.01	6.89	9.0%	19.25	8.3	0.0%	9.93	6.78	0.0%	9.59	6.72
2.0	6.0%	9.45	6.46	4.0%	9.28	6.59	0.0%	25.39	9.04	0.0%	9.63	6.58	0.0%	9.29	6.45
20	0.007	0.10	6 11	1.007	0.05	6 57	0.007	25 77	0.5	0.007	0.42	6 60	0.007	0.22	652

Table 12: Effect of K on Smoothing Attack Performance (Llama3-8B). Evaluation of the smoothing attack's effectiveness against different watermarking algorithms on the Llama3-8B model, varying the number of top-K tokens accessible to the attacker.

Κ		KGW		ا	Unigran	ı		SynthID)	DIP			
	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	
1	6%	2.37	4.67	19%	2.41	4.67	0%	3.6	6.86	2%	2.53	4.84	
3	1%	2.81	5.17	27%	2.8	5.2	0%	3.42	6.87	4%	2.91	5.47	
5	3%	2.99	5.36	24%	2.92	5.31	0%	3.41	6.89	1%	2.97	5.55	
7	2%	3.14	5.55	23%	3.03	5.43	0%	3.41	6.86	4%	3.1	5.78	
10	2%	3.2	5.63	24%	3.1	5.44	0%	3.4	6.86	6%	3.17	5.67	
	Unbiased												
K	τ	Unbiase	d		UPV			EWD			SWEET		
K	TPR	Unbiase PPL	d Div	TPR	UPV PPL	Div	TPR	EWD PPL	Div	TPR	SWEET PPL	Div	
К 1	TPR 1%	Unbiase PPL 2.5	d Div 4.8	TPR 1%	UPV PPL 2.48	Div 4.76	TPR 3%	EWD PPL 2.43	Div 4.68	TPR 3%	$\frac{\text{SWEET}}{\frac{\text{PPL}}{2.41}}$	Div 4.72	
K 1 3	TPR 1% 4%	Unbiase PPL 2.5 2.9	d Div 4.8 5.44	TPR 1% 1%	UPV PPL 2.48 2.97	Div 4.76 5.37	TPR 3% 4%	EWD PPL 2.43 2.94	Div 4.68 5.33	TPR 3% 3%	SWEET PPL 2.41 2.91	Div 4.72 5.27	
K 1 3 5	TPR 1% 4% 2%	Unbiase PPL 2.5 2.9 2.95	d Div 4.8 5.44 5.53	TPR 1% 1% 0%	UPV PPL 2.48 2.97 3.02	Div 4.76 5.37 5.55	TPR 3% 4% 3%	EWD PPL 2.43 2.94 3.06	Div 4.68 5.33 5.48	TPR 3% 3% 4%	SWEET PPL 2.41 2.91 3.01	Div 4.72 5.27 5.43	
K 1 3 5 7	TPR 1% 4% 2% 7%	Unbiase PPL 2.5 2.9 2.95 3.14	d Div 4.8 5.44 5.53 5.72	TPR 1% 1% 0% 1%	UPV PPL 2.48 2.97 3.02 3.1	Div 4.76 5.37 5.55 5.54	TPR 3% 4% 3% 6%	EWD PPL 2.43 2.94 3.06 3.09	Div 4.68 5.33 5.48 5.43	TPR 3% 3% 4% 5%	SWEET PPL 2.41 2.91 3.01 3.01	Div 4.72 5.27 5.43 5.37	

Table 13: Effect of α on Smoothing Attack Performance (Llama3-8B). Evaluation of the smoothing attack's effectiveness against different watermarking algorithms on the Llama3-8B model, varying the parameter α . A larger α indicates greater reliance on the reference model's output, while a smaller α means the attack text is more influenced by the watermarked model.

α		KGW		۱	Unigran	1		SynthIE)	DIP			
	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	
0.5	13%	3.45	5.92	62%	3.4	5.77	0%	3.78	6.88	35%	3.34	6.19	
1.0	2%	3.2	5.63	24%	3.1	5.44	0%	3.4	6.86	6%	3.17	5.67	
2.0	0%	3.05	5.28	12%	2.93	5.21	0%	3.49	6.87	3%	2.99	5.23	
3.0	0%	2.93	5.17	12%	2.99	5.26	0%	3.52	6.83	1%	2.96	5.16	
α	τ	Jnbiase	d	UPV				EWD		SWEET			
	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	TPR	PPL	Div	
0.5	26%	3.37	6.09	10%	3.47	6.08	28%	3.44	5.84	44%	3.38	5.9	
1.0	5%	3.17	5.75	1%	3.12	5.49	3%	3.13	5.38	4%	3.09	5.4	
2.0	3%	2.98	5.28	0%	2.96	5.2	0%	3.0	5.36	1%	3.06	5.38	
3.0	3%	2.96	5.21	0%	2.99	5.2	0%	2.89	5.18	0%	2.93	5.22	

reference and target models. Even with K = 1, 1122 the attack can drastically reduce the true positive 1123 rate (TPR) from 99% (without any attack) to an 1124 extremely low value, sometimes reaching 0%. This 1125 indicates that even with minimal access to both 1126 models, the smoothing attack can effectively re-1127 move watermarks. Furthermore, we observe that 1128 increasing K leads to more diverse text generation, 1129 as discussed in the main paper. This is because a 1130 higher K provides the attack with a larger selection 1131 of candidate tokens, allowing for greater variation 1132 in the generated text. This observation remains con-1133 sistent across both the OPT-1.3B and Llama3-8B 1134 models. 1135

Table 11 and Table 13 analyze the performance 1136 of smoothing attacks against different watermark-1137 ing algorithms under varying values of α . In this 1138 attack, the weight assigned to the top-K tokens 1139 from the watermarked model is given by c^{α} , while 1140 the weight for the top-K tokens from the reference 1141 model is $1 - c^{\alpha}$, where c is a confidence score be-1142 tween 0 and 1. A larger α shifts the token selection 1143 preference toward the reference model, making the 1144 generated text more aligned with it. Conversely, a 1145 smaller α biases the attack toward the watermarked 1146 model, producing text that more closely resembles 1147 the watermarked output. As α increases, the true 1148 positive rate (TPR) decreases, leading to a higher 1149 watermark removal rate-an effect consistently ob-1150 served across all watermarking methods for both 1151 the OPT-1.3B and Llama3-8B models. 1152

> In terms of text quality, when α is lower, the generated text is more influenced by the watermarked model, which generally exhibits higher quality than the reference model. Consequently, decreasing α can improve text quality. This provides an adversary with a way to adjust α to balance watermark removal and text quality preservation.

B.5 Text example

1153

1154

1155

1156

1157

1158

1159

1160

Table 14 presents text generated by the Gumbel 1161 sampling algorithm and the smoothing attack. We 1162 observe that, although the perplexity of the water-1163 marked text is significantly lower than that of the 1164 1165 text from the smoothing attack, this is primarily due to repetition in the generated text. This be-1166 havior may stem from the deterministic nature of 1167 Gumbel sampling, which can lead to less diverse 1168 outputs. 1169

B.6 Impact of model size

Table 15 presents the performance of the smooth-
ing attack across different watermarking algorithms
and varying sizes of OPT models. Perplexity
(PPL) is computed with respect to the OPT-30B
model, while the reference model remains consis-
tent across all settings—the OPT-125M.

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

For unwatermarked models, the True Positive Rate (TPR) is consistently 0%. In contrast, watermarked models achieve near-perfect TPR. However, the smoothing attack significantly reduces TPR across all model sizes, with its impact increasing as the model size grows—for instance, TPR drops to 0% for the KGW watermark in the 30B model.

Watermarked models exhibit a notable increase in perplexity, indicating that watermarking impacts text fluency. The smoothing attack reduces perplexity, bringing it closer to unwatermarked levels, suggesting a partial recovery of fluency. Regarding diversity, the unwatermarked text demonstrates the highest variation, while watermarking constrains generation patterns, resulting in a noticeable drop in diversity. The smoothing attack further reduces diversity, primarily because tokens are sampled only from the top-K tokens of both the watermarked and reference models, limiting the range of possible candidates.

C Analysis

C.1 Contribution Depends on the Confidence Score of the Unwatermarked Model

We first demonstrate that the contribution of each token to the detection score is influenced by the confidence score of the unwatermarked model, as measured by its probability distribution.

C.1.1 Case Study: Green-Red List Watermark

Suppose that l_t is the logit vector for predicting the *t*-th token from the unwatermarked model, and \mathcal{G}_t is the green list used by the watermarked model at position *t*, with size $\gamma |\mathcal{V}|$. Given the watermark shift δ , the probabilities assigned by the unwatermarked and watermarked models are expressed as:

$$P_t(v) = \frac{\exp(l_t(v))}{\sum_{v' \in \mathcal{V}} \exp(l_t(v'))}.$$
121:

$$\widetilde{P}_t(v) = \frac{\exp(l_t(v) + \delta \cdot \mathbf{1}_{\{v \in \mathcal{G}_t\}})}{\sum_{v' \in \mathcal{V}} \exp(l_t(v') + \delta \cdot \mathbf{1}_{\{v' \in \mathcal{G}_t\}})}.$$
1214

Table 14: Text generated by watermarked model and smoothing attack (OPT-1.3B against Gumbel sampling).

Watermarked	Smoothing Attack
Cluster comprises IBM's Opteron-based eServer 325	Cluster comprises IBM's Opteron-based eServer 325
server and systems management software and storage	server and systems management software and storage
devices that can run Linux and Windows operating	devices that can run Linux and Windows operating
systems.	systems.
The data center will be built on top of existing IBM	IBM will start selling customers a prototype of the
Power servers. The company will offer a variety of	cluster by July, according to Jim Bessen, the executive
services, including cloud-based services, as well as a	vice president of middleware architecture and services
"plug and play" environment that will allow users to	at IBM.
transition to other hardware.	"The cluster is just the tip of the iceberg," he said.
The data center will be built on top of existing IBM	The cluster firm computing is encased in a step-by-step
Power servers. The company will offer a variety of	process starting with activities such as how to set up
services, including cloud-based services, as well as a	clinks that back up disks, store data, create back-up
"plug and play" environment that will allow users to	schemes for data, and at what resolutions to save the
transition to other hardware. (Image: IBM)	data.
The data center will be built on top of existing IBM	"We are not done yet," he said. As new kinds of
Power servers. The company will offer a variety of	programming software and other technologies come out,
services, including cloud-based services, as well as a	users need to act on it.
"plug and play" environment that will allow users to	While IBM expects internal clusters with the capability
transition to other hardware.	to run multiple operating systems in fuse over the next
The data center will be built on top of existing IBM	year, this capability will be available only to Enterprise
Power servers. The company will offer a variety of	Software Group (ESG) customers.

ESG will not sell its cluster technology to anyone else, Bessen said.

1215 Rewriting
$$\widetilde{P}_t(v)$$
, we observe:

1216
$$\widetilde{P}_t(v) = P_t(v) \times \frac{\exp(\delta \mathbf{1}_{\{v \in \mathcal{G}_t\}})}{\sum_{v' \in \mathcal{V}} P_t(v') \exp(\delta \mathbf{1}_{\{v' \in \mathcal{G}_t\}})}.$$

transition to other hardware.

Define the normalization factor:

218
$$Z_{\delta} = \frac{\sum_{v' \in \mathcal{V}} \exp(l_t(v') + \delta \mathbf{1}_{\{v' \in \mathcal{G}_t\}})}{\sum_{v' \in \mathcal{V}} \exp(l_t(v'))} \quad (6)$$

$$= \sum_{v' \in \mathcal{V}} P_t(v') \exp(\delta \mathbf{1}_{\{v' \in \mathcal{G}_t\}}).$$
(7)

services, including cloud-based services, as well as a

"plug and play" environment that will allow users to

Then:

1217

1219

1220

1221

1222

1223

1224

1226

1227

1228

1229

$$\widetilde{P}_t(v) = \begin{cases} \frac{e^{\delta}}{Z_{\delta}} P_t(v), & v \in \mathcal{G}_t, \\ \frac{1}{Z_{\delta}} P_t(v), & v \notin \mathcal{G}_t. \end{cases}$$

The expected fraction of tokens belonging to the green list under the unwatermarked model is given by:

1225
$$\mathbb{E}_{v \sim P_t}[\mathbf{1}(v \in \mathcal{G}_t)] = \sum_{v \in \mathcal{G}_t} P_t(v) = P_{\mathcal{G}_t},$$

where P_{G_t} represents the probability mass assigned to green tokens in the unwatermarked model.

> Similarly, the expected fraction of green tokens in the watermarked model is:

$$\mathbb{E}_{v \sim \widetilde{P}_t}[\mathbf{1}(v \in \mathcal{G}_t)] = \sum_{v \in \mathcal{G}_t} \widetilde{P}_t(v) = \frac{e^{\delta}}{Z_{\delta}} P_{\mathcal{G}_t}.$$
 (8) 123

Since $Z_{\delta} = (e^{\delta} - 1)P_{\mathcal{G}_t} + 1$, the difference in 1231 green token probabilities (i.e., the detection contri-1232 bution at token position t) is: 1233

$$S_t = \mathbb{E}_{v \sim \widetilde{P}_t} [\mathbf{1}(v \in \mathcal{G}_t)] - \mathbb{E}_{v \sim P_t} [\mathbf{1}(v \in \mathcal{G}_t)]$$
(9) 1234

$$=\frac{-(e^{\delta}-1)P_{\mathcal{G}_t}+(e^{\delta}-1)}{(e^{\delta}-1)+\frac{1}{P_{\mathcal{G}_t}}}.$$
 (10)

In other words, the token-level detection contribution S_t is a function of the probability mass $P_{\mathcal{G}_t}$ assigned to green tokens by the unwatermarked model.

Case Study: Tournament Sampling C.1.2 Watermark

In the Tournament Sampling watermark, when gen-1242 erating the *t*-th token, the algorithm assigns scores 1243 to each token using m independent watermarking 1244 functions $g^{(1)}, ..., g^{(m)}$. These scores depend on a 1245 random seed generated based on the recent context 1246 and a secret watermarking key. The token selection 1247 follows a multi-round elimination process, where 1248 2^m tokens are first sampled from $P_t(\cdot)$, then com-1249 pete in m rounds to determine the final output. 1250

1235

1236

1237

1238

1239

1240

Table 15: Impact of Model Size on the Smoothing Attack (OPT). Performance of the smoothing attack across different watermarking algorithms and various sizes of OPT models. The perplexity (PPL) is computed with respect to the OPT-30B model, while the reference model is consistently the OPT-125M. The table reports True Positive Rate (TPR), Perplexity (PPL), and Diversity (Div.) for unwatermarked, watermarked, and smoothed settings.

Size	Setting	1	KGW		Unigram			S	SynthID			DIP			Unbiased		
DILU	Setting	TPR	PPL	Div.	TPR	PPL	Div.	TPR	PPL	Div.	TPR	PPL	Div.	TPR	PPL	Div.	
1.3B	Unwatermarked	0.0%	12.95	8.67	0.0%	12.95	8.67	0.0%	12.95	8.67	0.0%	12.95	8.67	0.0%	12.95	8.67	
	Watermarked	100.0%	15.94	8.09	99.0%	16.53	7.29	100.0%	7.7	7.41	100.0%	15.16	8.44	99.0%	15.14	8.29	
	Smoothing	4.0%	10.48	6.72	6.0%	10.37	6.83	1.0%	11.37	8.67	6.0%	10.03	7.03	4.0%	9.94	6.79	
2.7B	Unwatermarked	0.0%	11.75	8.36	0.0%	11.75	8.36	0.0%	11.75	8.36	0.0%	11.75	8.36	0.0%	11.75	8.36	
	Watermarked	100.0%	13.94	7.88	100.0%	14.31	7.41	99.0%	6.86	7.55	97.0%	13.86	8.61	97.0%	13.6	8.69	
	Smoothing	4.0%	10.35	6.77	4.0%	10.35	6.66	6.0%	9.84	8.0	13.0%	9.87	6.84	6.0%	9.85	6.88	
6.7B	Unwatermarked	0.0%	10.2	8.45	0.0%	10.2	8.45	0.0%	10.2	8.45	0.0%	10.2	8.45	0.0%	10.2	8.45	
	Watermarked	100.0%	13.16	8.06	100.0%	12.94	7.48	98.0%	6.21	7.48	98.0%	11.8	8.48	97.0%	11.79	8.59	
	Smoothing	4.0%	10.07	6.92	6.0%	10.54	6.68	3.0%	8.98	8.31	8.0%	9.78	6.86	8.0%	9.68	6.74	
13B	Unwatermarked	0.0%	10.14	8.39	0.0%	10.14	8.39	0.0%	10.14	8.39	0.0%	10.14	8.39	0.0%	10.14	8.39	
	Watermarked	100.0%	12.88	8.56	100.0%	12.44	7.39	100.0%	5.88	7.8	96.0%	11.67	9.34	93.0%	11.42	8.77	
	Smoothing	2.0%	10.24	6.82	5.0%	10.32	6.7	8.0%	8.07	7.8	8.0%	9.6	6.88	7.0%	9.37	6.77	
30B	Unwatermarked	0.0%	8.46	8.44	0.0%	8.46	8.44	0.0%	8.46	8.44	0.0%	8.46	8.44	0.0%	8.46	8.44	
	Watermarked	100.0%	10.23	8.34	100.0%	10.45	7.56	100.0%	5.27	7.72	94.0%	9.43	8.78	97.0%	9.89	9.08	
	Smoothing	0.0%	9.5	6.8	7.0%	10.15	6.75	5.0%	6.96	8.04	4.0%	9.34	6.89	4.0%	9.36	6.88	
Size	Setting	X-SIR			UPV			Gumbel			EWD			SWEET			
	6	TPR	PPL	Div.	TPR	PPL	Div.	TPR	PPL	Div.	TPR	PPL	Div.	TPR	PPL	Div.	
1.3B	Unwatermarked	1.0%	12.95	8.67	0.0%	12.95	8.67	0.0%	12.95	8.67	0.0%	12.95	8.67	0.0%	12.95	8.67	
	Watermarked	94.0%	15.42	7.96	99.0%	12.79	8.22	98.0%	3.15	4.35	100.0%	16.88	7.92	100.0%	15.99	8.02	
	Smoothing	13.0%	10.3	6.72	20.0%	10.78	6.89	9.0%	20.94	8.30	1.0%	10.71	6.75	1.0%	10.54	6.81	
2.7B	Unwatermarked	3.0%	11.75	8.36	0.0%	11.75	8.36	0.0%	11.75	8.36	0.0%	11.75	8.36	0.0%	11.75	8.36	
	Watermarked	91.0%	14.07	8.25	99.0%	12.30	8.01	99.0%	2.96	4.38	100.0%	14.88	7.98	100.0%	14.07	8.32	
	Smoothing	10.0%	10.34	6.77	18.0%	10.56	6.90	10.0%	19.46	8.41	1.0%	10.43	6.86	3.0%	10.49	6.86	
6.7B	Unwatermarked	0.0%	10.2	8.45	0.0%	10.20	8.45	0.0%	10.20	8.45	0.0%	10.20	8.45	0.0%	10.20	8.45	
	Watermarked	91.0%	13.04	8.19	97.0%	10.92	7.75	100.0%	2.97	4.49	100.0%	13.42	8.69	100.0%	13.05	8.41	
	Smoothing	9.0%	10.01	6.7	8.0%	10.60	7.05	9.0%	14.85	8.62	0.0%	10.60	6.79	1.0%	10.07	6.89	
13B	Unwatermarked	0.0%	10.14	8.39	0.0%	10.14	8.39	0.0%	10.14	8.39	0.0%	10.14	8.39	0.0%	10.14	8.39	
	Watermarked	88.0%	12.29	8.05	99.0%	10.59	7.91	98.0%	2.96	4.63	100.0%	13.09	8.74	100.0%	12.32	8.35	
	Smoothing	11.0%	9.84	6.79	12.0%	10.84	6.88	12.0%	15.06	8.27	0.0%	10.16	6.73	2.0%	10.15	6.74	
30B	Unwatermarked	0.0%	8.46	8.44	0.0%	8.46	8.44	0.0%	8.46	8.44	0.0%	8.46	8.44	0.0%	8.46	8.44	
	Watermarked	91.0%	10.43	8.43	97.0%	8.59	8.13	97.0%	2.89	4.79	100.0%	10.75	8.54	100.0%	9.98	8.25	
	Smoothing	16.0%	9.65	6.74	17.0%	10.06	7.11	9.0%	11.92	8.39	2.0%	10.02	6.99	2.0%	9.55	6.84	

Despite the complex sampling mechanism, the probability of each token in the modified distribution \tilde{P}_t is adjusted by a factor dependent on its assigned g value. Specifically, for any token v:

$$\widetilde{P}_{t}(v) = \begin{cases} P_{t}(v) \cdot (1 - P_{\mathcal{G}_{t}}) & \text{if } g(v) = 0, \\ P_{t}(v) \cdot (2 - P_{\mathcal{G}_{t}}) & \text{if } g(v) = 1. \end{cases}$$
(11)

During watermark detection, the detector computes the average g value across all tournament layers, i.e., $\frac{1}{m} \sum_{l=1}^{m} g^{(l)}(v)$, as the watermark score for the token.

Single Tournament Layer (m = 1). Consider the simplest case where m = 1, meaning only one tournament round is used. Let \mathcal{G}_t denote the set of tokens where $g^{(1)}(v) = 1$. The probability modification simplifies to:

$$\widetilde{P}_t(v) = \begin{cases} P_t(v) \cdot (1 - P_{\mathcal{G}_t}) & \text{if } v \notin \mathcal{G}_t, \\ P_t(v) \cdot (2 - P_{\mathcal{G}_t}) & \text{if } v \in \mathcal{G}_t. \end{cases}$$
(12)

The expected g value for tokens sampled from \widetilde{P}_t is $(2 - P_{\mathcal{G}_t}) \cdot P_{\mathcal{G}_t}$, while the expectation under P_t is simply $P_{\mathcal{G}_t}$. Thus, the detection contribution S_t is:

$$S_t = (1 - P_{\mathcal{G}_t}) \cdot P_{\mathcal{G}_t}.$$
(13)

This mirrors the Green-Red List watermark, showing that the detection contribution per token is fundamentally tied to $P_{\mathcal{G}_t}$.

C.2 Low Model Confidence Leads to Large Variance in the Watermark Score for Unwatermarked Text

Thus far, we have established that the contribution1277of each token to the detection score is correlated1278

1281 1282

1283

1285

1286

1289 1290

1291

1292

1293

129

1295

1297

1298

1300

1301

1302

1303

1304

1306

1307

1309

with the expected watermark score under the unwatermarked model. We now analyze what affects the watermark score of the unwatermarked model.

Let $P_t = (p_1, p_2, ..., p_d)$ be the probability vector from the unwatermarked model at token position t, where $p_i \in [0, 1]$ and $\sum_{i=1}^d p_i = 1$. Typically, $d = |\mathcal{V}|$ is large. We randomly select a subset $\mathcal{G}_t \subset \{1, ..., d\}$ of indices of size $\gamma |\mathcal{V}|$. Define the random variable:

$$P_{\mathcal{G}_t} = \sum_{i \in \mathcal{G}_t} p_i$$

We analyze how $P_{\mathcal{G}_t}$ is distributed over all possible assignments of \mathcal{G}_t . Define the indicator variable X_i as follows:

$$X_i = \begin{cases} 1, & \text{if } i \in \mathcal{G}_t, \\ 0, & \text{otherwise.} \end{cases}$$

Since each token is independently assigned to \mathcal{G}_t with probability γ , we have:

$$\mathbb{E}[X_i] = \gamma$$
, and $\operatorname{Var}(X_i) = \gamma(1 - \gamma)$.

For different token indices $i \neq j$, the covariance between their assignments is:

$$Cov(X_i, X_j) = \mathbb{E}[X_i X_j] - \mathbb{E}[X_i] \mathbb{E}[X_j].$$

For Poisson sampling (i.e., assigning each token to \mathcal{G}_t independently with probability γ), the covariance is zero. However, under a fixed-size sampling setup (i.e., selecting exactly $\gamma |\mathcal{V}|$ tokens), we have:

$$\operatorname{Cov}(X_i, X_j) = \frac{\gamma |\mathcal{V}|}{d} \cdot \frac{\gamma |\mathcal{V}| - 1}{d - 1} - \gamma^2 = -\frac{\gamma (1 - \gamma)}{|\mathcal{V}| - 1}$$

Expressing $P_{\mathcal{G}_t}$ in terms of X_i , we obtain:

$$P_{\mathcal{G}_t} \;=\; \sum_{i=1}^d X_i \, p_i$$

Expectation and Variance of $P_{\mathcal{G}_t}$. The expectation is:

1308
$$\mathbb{E}[P_{\mathcal{G}_t}] = \sum_{i=1}^d \mathbb{E}[X_i] p_i = \gamma \sum_{i=1}^d p_i = \gamma$$

The variance is:

$$\operatorname{Var}(P_{\mathcal{G}_t}) = \sum_{i=1}^d p_i^2 \operatorname{Var}(X_i) + \sum_{i \neq j} p_i p_j \operatorname{Cov}(X_i, X_j).$$
 1310

Substituting
$$\operatorname{Var}(X_i) = \gamma(1 - \gamma)$$
 and 1311
 $\operatorname{Cov}(X_i, X_j) = -\frac{\gamma(1 - \gamma)}{|\mathcal{V}| - 1}$: 1312

$$\operatorname{Var}(P_{\mathcal{G}_t}) = \gamma(1-\gamma) \sum_{i=1}^d p_i^2 - \frac{\gamma(1-\gamma)}{|\mathcal{V}| - 1} \sum_{i \neq j} p_i p_j.$$
 1313

For the first term,

$$\gamma(1-\gamma)\sum_{i=1}^{d}p_i^2 = \gamma(1-\gamma)\sigma^2,$$
1315

1314

1318

1320

1325

1327

1328

where $\sigma^2 = \sum_{i=1}^d p_i^2$ represents the squared ℓ_2 1316 norm of the probability vector. 1317

For the second term, using the identity:

$$\sum_{i \neq j} p_i p_j = \left(\sum_{i=1}^d p_i\right)^2 - \sum_{i=1}^d p_i^2 = 1 - \sigma^2,$$
 1319

and we obtain:

$$\frac{\gamma(1-\gamma)}{|\mathcal{V}|-1} \sum_{i \neq j} p_i p_j = \frac{\gamma(1-\gamma)}{|\mathcal{V}|-1} (1-\sigma^2).$$
 1321

For large $|\mathcal{V}|$, the correction term $\frac{\gamma(1-\gamma)}{|\mathcal{V}|-1}(1-\sigma^2)$ 1322 becomes negligible, and we approximate: 1323

$$\operatorname{Var}(P_{\mathcal{G}_t}) \approx \gamma (1 - \gamma) \sigma^2.$$
 1324

Interpretation. This analysis shows that $P_{\mathcal{G}_t}$ depends on the probability mass distribution.

High-Uncertainty Case (Uniform Distribution): If $p_i = \frac{1}{|\mathcal{V}|}$ for all *i*, then

$$\sigma^2 = \sum_{i=1}^{|\mathcal{V}|} \frac{1}{|\mathcal{V}|^2} = \frac{1}{|\mathcal{V}|}.$$
 1329

For large $|\mathcal{V}|$, σ^2 is small, meaning that the distribution of $P_{\mathcal{G}_t}$ concentrates tightly around γ with small variance. This corresponds to a scenario where the model has high uncertainty, spreading probability mass nearly uniformly over all tokens. 1334

Low-Uncertainty Case (Dominant Tokens): In1335practice, language models often assign high proba-
bility mass to a small number of dominant tokens.1336Suppose $p_j \ge 0.8$ for some token j, then:1338

$$\sigma^2 \ge p_j^2 = 0.64$$

In this case, σ^2 is much larger than $1/|\mathcal{V}|$ (which is on the order of 10^{-5} for large models). Consequently, $P_{\mathcal{G}_t}$ exhibits a bimodal distribution: it is either close to 0 or close to 1, depending on whether the dominant tokens are in \mathcal{G}_t . The probability of $P_{\mathcal{G}_t} \approx \gamma$ is nearly zero.

Thus, when the model is confident in its predictions (low uncertainty), the variance of $P_{\mathcal{G}_t}$ is large, leading to a higher variance in the watermark score. Conversely, when the model is uncertain, the watermark score is more stable and centered around γ .

Connection to Watermark Detection. Since the contribution to the detection score S_t depends on $P_{\mathcal{G}_t}$ (Eq. equation ??), its variance is governed by $Var(P_{\mathcal{G}_t})$. This means that tokens generated with high confidence contribute more variability to the detection score, whereas tokens generated under uncertainty contribute less variability.

C.3 Estimating the Confidence Score of the Unwatermarked Model Using the Watermarked Model

1362Our goal is to estimate the squared ℓ_2 norm of the1363probability distribution $||P_t||^2$, which serves as a1364confidence measure for the unwatermarked model,1365using only access to the watermarked model \widetilde{P}_t .1366This estimation is critical for adaptive attacks and1367for understanding how watermarking affects text1368quality.

Setup. We consider the Green-Red List watermarking scheme, where the probability distribution \widetilde{P}_t is obtained by modifying P_t as:

1372
$$\widetilde{P}_t(v) = \frac{e^{\delta \mathbf{1}_{\{v \in \mathcal{G}_t\}}}}{Z_{\delta}} P_t(v),$$

where the normalization factor Z_{δ} is defined as:

$$Z_{\delta} = (1 - P_{\mathcal{G}_t}) + e^{\delta} P_{\mathcal{G}_t}.$$

1375 We aim to construct an estimator \hat{U} for the con-1376 fidence measure:

1377
$$||P_t||^2 = \sum_{v \in \mathcal{V}} P_t(v)^2.$$

Expected Squared Norm of the Watermarked 1378 **Model.** Since each probability mass in P_t is 1379 scaled by either e^{δ}/Z_{δ} (if in \mathcal{G}_t) or $1/Z_{\delta}$ (if not 1380 in \mathcal{G}_t), we have: 1381

$$\mathbb{E}[\tilde{P}_t(v)^2] = (1-\gamma)\frac{1}{Z_{\delta}^2}P_t(v)^2 + \gamma \frac{e^{2\delta}}{Z_{\delta}^2}P_t(v)^2.$$
 1382

Summing over all tokens in \mathcal{V} , we obtain:

$$\mathbb{E}[\|\widetilde{P}_t\|^2] = \frac{(1-\gamma) + \gamma e^{2\delta}}{Z_{\delta}^2} \|P_t\|^2.$$
 1384

1383

1388

1394

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

Unbiased Estimator.Rearranging the above ex-1385pression, we define an unbiased estimator:1386

$$\widehat{U} = \frac{Z_{\delta}^2}{(1-\gamma) + \gamma e^{2\delta}} \|\widetilde{P}_t\|^2.$$
1387

Taking expectation, we confirm:

$$\mathbb{E}[\widehat{U}] = \|P_t\|^2.$$
1389

Practical Approximation.Since Z_{δ} depends on1390 $P_{\mathcal{G}_t}$, which is unknown to an adversary, we approximate it using γ :1391

$$Z_{\delta} \approx (1 - \gamma) + \gamma e^{\delta}.$$
 1393

Thus, the practical estimator becomes:

$$\widetilde{U} = \frac{[(1-\gamma)+\gamma e^{\delta}]^2}{(1-\gamma)+\gamma e^{2\delta}} \|\widetilde{P}_t\|^2.$$
1395

This provides a computationally efficient way to estimate $||P_t||^2$ using only \tilde{P}_t , making it useful for designing attacks.

C.4 Estimating the ℓ_2 Norm Using Top-KProbabilities

While we have established the connection between the squared ℓ_2 norm $||P_t||^2$ of the probability distribution and its contribution to the watermark detection score, direct access to this quantity is often unavailable, even for the watermarked model. In this section, we show how to estimate $||P_t||^2$ using only limited access to the model's top-K probabilities.

Suppose we only have access to the top-K probabilities:

$$p_1 \ge p_2 \ge \dots \ge p_K, \tag{1411}$$

1339

1340

1341

1342

1343

1344

1345

1346

1347

1349

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1369

1370

1371

1373

where the remaining probabilities $p_{K+1}, \ldots, p_{|\mathcal{V}|}$ are unknown. Define the remaining probability mass of the tail as:

1415
$$R = 1 - \sum_{i=1}^{K} p_i$$

1417
$$||P_t||^2 = \sum_{i=1}^{|\mathcal{V}|} p_i^2,$$

1418 given only p_1, \ldots, p_K and R.

We bound $||P_t||^2$ by considering two extreme ways in which the unknown tail probabilities could be distributed:

Our goal is to estimate the squared ℓ_2 norm:

- 1. Uniform Tail: The remaining probability mass *R* is evenly distributed across the unknown tokens, minimizing the sum of squares.
- 2. Concentrated Tail: The entire probability mass R is assigned to a single token, maximizing the sum of squares.

Uniform Tail (Lower Bound) If the tail probability mass R is *uniformly* spread among the remaining $|\mathcal{V}| - K$ tokens, then each unknown probability is $\frac{R}{|\mathcal{V}|-K}$. The squared sum of the tail probabilities is then:

$$\sum_{i=K+1}^{|\mathcal{V}|} p_i^2 = (|\mathcal{V}| - K) \left(\frac{R}{|\mathcal{V}| - K}\right)^2 = \frac{R^2}{|\mathcal{V}| - K}$$

Since distributing the mass uniformly minimizes the squared sum (due to convexity), this scenario provides a lower bound for $||P_t||^2$:

$$\|P_t\|^2 \geq \sum_{i=1}^K p_i^2 + rac{R^2}{|\mathcal{V}| - K}.$$

Concentrated Tail (Upper Bound) At the other extreme, if the entire remaining probability mass R is assigned to a single token, then the squared sum of the tail probabilities is simply:

1442
$$\sum_{i=K+1}^{|\mathcal{V}|} p_i^2 = R^2$$

Since concentrating all probability mass in one entry maximizes the sum of squares, this provides an upper bound for $||P_t||^2$:

$$|P_t||^2 \le \sum_{i=1}^K p_i^2 + R^2.$$
 1446

Combining both bounds, we obtain:

$$\sum_{i=1}^{K} p_i^2 + \frac{R^2}{|\mathcal{V}| - K} \leq ||P_t||^2 \leq \sum_{i=1}^{K} p_i^2 + R^2,$$
 1448

where
$$R = 1 - \sum_{i=1}^{K} p_i$$
. 1449

Practical Approximation.A commonly used1450practical heuristic is to assume that the remaining1451probability mass R follows a uniform distribution1452across the unknown probabilities.1453sumption, we approximate:1454

$$||P_t||^2 \approx \sum_{i=1}^{K} p_i^2 + \frac{R^2}{|\mathcal{V}| - K}.$$
 1455

This estimate tends to be slightly lower than the true value, since in reality, the tail probabilities are rarely perfectly uniform—some tokens may have slightly higher probabilities than others. However, in the case of language modeling, probability distributions often exhibit a "long tail" where the remaining probability mass is spread across many small values. In such cases, the uniform assumption serves as a reasonable first-order approximation.

C.5 Additional Numerical Analysis

Generalization to other watermarking solutions. For Gumbel sampling, we define the token-level contribution to watermark detection as $S_t = -\log(1 - U_{v^*}) - \mathbb{E}_{v \sim P_t}[-\log(1 - U_v)]$, where v^* is the token selected by the watermarked model. Note that the choice of v^* is deterministic after the secret key held by the LM provider and the prefix content are fixed. For Tournament sampling, we define the token-level contribution as $S_t = \mathbb{E}_{v \sim \tilde{P}_t} \left[\frac{1}{m} \sum_{l=1}^m g^{(l)}(v, r) \right] - \mathbb{E}_{v \sim P_t} \left[\frac{1}{m} \sum_{l=1}^m g^{(l)}(v, r) \right]$, where \tilde{P}_t is the modified probability distribution.

For these two watermarks, we still observe the same correlation between S_t and $||P_t||^2$ as we have for Green-list watermarks, as shown in Figure 7. Namely, the token-level contribution S_t to the watermark detectability is negatively correlated to the model's confidence at position t.



Figure 7: The correlation between S_t (watermark contribution score) and $||P_t||^2$ (model confidence) evaluated on model OPT-1.3B with the Gumbel and Tournament sampling (with *m* tournaments) watermarks, using the same setup as in Figure 1. Each sample corresponds to a specific prefix and secret key. $||P_t||^2$ is computed from the original un-watermarked model. The overall observation is similar to what we have for the *Green-red list* watermarking: S_t decreases as $||P_t||^2$ increases.



Figure 8: The correlation between $D_{TV}(P_t, \tilde{P}_t)$, i.e., the negative impact on text quality due to watermarks (in color blue), and $||P_t||^2$, measured on OPT-1.3B with the Green-red list and Gumbel and Tournament sampling watermarks. We also plot $D_{TV}(P_t, P_t^{\text{ref}})$, which measures the negative impact on text quality if we use tokens sampled from the reference model OPT-125M (in color red).

Impact of watermarking on text quality We also plot $D_{TV}(P_t, P_t^{\text{ref}})$, which measures the negative impact on text quality if we alternatively sample from the reference model OPT-125M (in color red). We note that when the model is not confident in its output, i.e., when $||P_t||^2$ is small, sampling from the reference model's token distribution, i.e., P_t^{ref} , does not hurt the text quality. In particular, under the Green-red list watermarking scheme, $D_{TV}(P_t, P_t^{ref})$ is comparable to $D_{TV}(P_t, \widetilde{P}_t)$ when $||P_t||^2$ is small (observe that the red points generally overlap with the blue ones). For Gumbel and Tournament sampling, $D_{TV}(P_t, P_t^{\text{ref}})$ is even smaller than $D_{TV}(P_t, \widetilde{P}_t)$ when $||P_t||^2$ is small (observe that the red points are generally below the blue ones). Conversely, when the model is confident in its output, i.e., when $||P_t||^2$ is large, replacing the watermarked

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1498

1499

1500

1501

1502

model with a reference model may hurt the text quality (observe that the red points are above the blue ones).

1503

1504

1505

Trade-off between detectability and text qual-1506 ity In Figure 9, we plot the correlation between 1507 $D_{TV}(P_t, P_t)$ and S_t , empirically measured on 1508 OPT-1.3B model using the same setup as the above 1509 simulations. When the watermark has little im-1510 pact on text quality (i.e., smaller total variation distance), the watermark is also less detectable (i.e., 1512 smaller S_t). Conversely, tokens that contribute 1513 more to watermark detection also lead to more 1514 notable text quality degradation. This finding, in 1515 turn, reveals the crucial limitation of existing wa-1516 termarking schemes: high watermark detectability 1517 and high text quality cannot be achieved at the same 1518 time, since the very same set of tokens causes qual-1519 ity degradation while contributing to watermark 1520



Figure 9: The correlation between $D_{TV}(P_t, \tilde{P}_t)$, i.e., the negative impact of watermarking on the text quality, and S_t , i.e., the token-level contribution to watermark detectability. We measure this on OPT-1.3B. For all three watermarking schemes, $D_{TV}(P_t, \tilde{P}_t)$ increases as S_t increases.

detectability simultaneously.

1521

1522

D Possible Defenses to Smoothing Attack

Our attack exploits the correlation between a to-1523 ken's contribution to the watermark detection score 1524 and the confidence level of the unwatermarked model in predicting that token. One possible defense against this attack is to restrict access to 1527 1528 confidence-related information, such as returning only the most probable token without revealing its probability. Note that, if the probability of the 1530 most likely token is available, our attack remains effective. However, such a defense is challeng-1532 ing to enforce in practice. Many existing LLM services provide top-K probabilities (e.g., Ope-1534 nAI's API returns probabilities for the top 20 to-1535 kens), which is already sufficient to approximate 1536 model confidence and execute our attack. Moreover, service providers often release these prob-1538 abilities to enhance transparency and build trust 1539 by providing insights into the model's reasoning, 1540 addressing concerns about the opacity of AI systems (European Commission, 2021; OECD, 2019). 1542 Access to probability distributions is also essential for debugging and evaluating model perfor-1544 mance, as it allows developers to identify biases, 1545 diagnose overconfidence, and improve reliability (National Institute of Standards and Technology 1547 (NIST), 2023). Probabilities support explainable 1548 AI (XAI) by revealing model uncertainty, enabling users to interpret predictions and explore alterna-1551 tive suggestions (Brown et al., 2020). From an ethical standpoint, making probability distributions available facilitates bias auditing and aligns with 1553 broader efforts to promote fairness and accountability in AI (OECD, 2019). Additionally, probability 1555

information empowers developers and end users 1556 by enabling advanced decision-making strategies, 1557 such as re-ranking, rejection sampling, and beam 1558 search (OpenAI, 2023). Furthermore, it helps mit-1559 igate risks associated with model overconfidence 1560 and hallucinations, which is particularly crucial in 1561 high-stakes domains such as healthcare and law 1562 (National Institute of Standards and Technology 1563 (NIST), 2023). Given the practical difficulties in restricting access to confidence-related information, 1565 our findings suggest that existing watermarking 1566 techniques may be vulnerable when model confi-1567 dence can be estimated. This highlights the need 1568 for developing watermarking schemes that remain effective even in scenarios where adversaries have 1570 partial access to confidence estimates. Future re-1571 search should explore watermarking methods that 1572 explicitly account for the model's confidence and 1573 ensure robustness against adversarial attacks that 1574 exploit confidence information. 1575