

STEALTHRL: REINFORCEMENT LEARNING PARAPHRASE ATTACKS FOR MULTI-DETECTOR EVASION OF AI-TEXT DETECTORS

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce *StealthRL*, a reinforcement-learning framework that paraphrases AI-generated text to evade detection while preserving meaning and fluency, enabling systematic red-teaming of AI-text detectors. The challenge is to reduce detector confidence at strict low-FPR operating points without collapsing semantic fidelity or overfitting to a single detector family. StealthRL fine-tunes a single Qwen3-4B paraphraser with LoRA Hu et al. (2021) and a multi-objective reward combining detector evasion, semantic similarity, and fluency. On the MAGE benchmark with three detectors (RoBERTa OpenAI, Fast-DetectGPT, Binoculars), StealthRL reduces mean TPR@1%FPR from 0.27 (no attack) to 0.04 while maintaining high semantic similarity (0.953 E5 cosine), outperforming simple paraphrasing and a detector-guided baseline. We release an anonymized code package in the supplementary material with a placeholder anonymous link for reproducibility.

1 INTRODUCTION

AI-generated text detectors are widely deployed in academic integrity, content moderation, and misinformation pipelines, yet their robustness to paraphrasing remains brittle. Recent work shows that detector performance can drop sharply under paraphrase attacks, motivating adversarial evaluation and red-teaming rather than static benchmarking. We focus on the adversarial aspect: how to generate high-fidelity paraphrases that reliably evade multiple detector families, while maintaining readability and semantic fidelity.

We present **StealthRL**, a reinforcement-learning paraphraser trained to reduce detector confidence under a strict low-FPR operating point. StealthRL is designed to (i) train against multiple detector families to promote transfer, (ii) preserve semantic content and fluency via auxiliary reward terms, and (iii) provide an evaluation harness that reports transfer and attack success at 1% FPR. We position StealthRL as a red-teaming tool that helps stress-test detector robustness and guide defensive improvements.

Contributions.

- **Multi-detector RL paraphrasing.** We train a single paraphrase policy against a detector ensemble using group-relative policy optimization with LoRA, enabling efficient adversarial fine-tuning.
- **Low-FPR evaluation protocol.** We report attack success across multiple detector families at TPR@1%FPR and release an evaluation pipeline that produces heatmaps and tradeoff curves.
- **Empirical results.** On MAGE, StealthRL reduces mean TPR@1%FPR to 0.04 (from 0.27 for no attack) while maintaining high semantic similarity.

Anonymized code release. An anonymized code release containing training and evaluation scripts is included in the supplementary material. A placeholder anonymous link is provided: <https://anonymous.4open.science/r/STEALTHRL>.

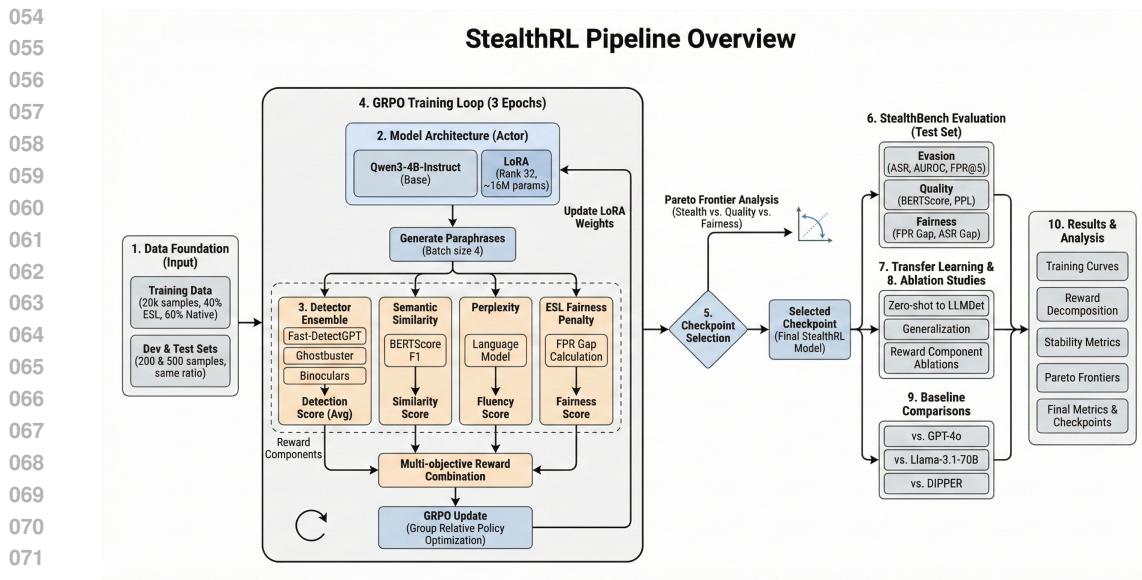


Figure 1: StealthRL training and evaluation pipeline. A single paraphrase policy is trained with detector-guided reward and quality constraints, then evaluated across multiple detector families with low-FPR metrics and tradeoff analysis.

2 RELATED WORK

Detector families include curvature-based methods (DetectGPT and Fast-DetectGPT) and paired-LM detectors such as Binoculars. Mitchell et al. (2023); Bao et al. (2024); Hans et al. (2024) Recent evasion work focuses on paraphrase attacks and RL-based humanization, including Adversarial Paraphrasing and AuthorMist, while character-level attacks such as homoglyph substitution demonstrate strong but often less readable perturbations. Cheng et al. (2025); David & Gervais (2025); Creo & Pudasaini (2025) StealthRL builds on these directions by training a single paraphrase policy against a detector ensemble and evaluating transfer at strict low-FPR operating points.

3 THREAT MODEL

We assume **black-box access to detector scores**, i.e., the attacker can query detector confidence but does not require gradients. In practice, detectors are often open-source or deployed with a confidence API, so both black-box scoring and open-source replication are realistic. We evaluate transfer to a held-out detector family to test robustness beyond the training ensemble.

4 METHOD

Given AI-generated text x , we learn a paraphrase policy $\pi_\theta(y | x)$ that produces y with low detector confidence while preserving meaning and fluency. We optimize a composite reward:

$$R(x, y) = \alpha R_{\text{det}}(y) + \beta R_{\text{sem}}(x, y) + \gamma R_{\text{ppl}}(y), \quad (1)$$

where R_{det} is the ensemble detector score (lower AI probability), R_{sem} is E5 embedding cosine similarity, and R_{ppl} is a fluency proxy derived from a frozen LM. We train with group-relative policy optimization (GRPO) and LoRA adapters Hu et al. (2021) on Qwen/Qwen3-4B-Instruct-2507 (rank 32, $\alpha = 32$, dropout 0.05). Training uses group size 4, learning rate 2.8×10^{-4} , batch size 16, and two epochs. The detector ensemble used for training is RoBERTa OpenAI + Fast-DetectGPT with weights 0.6/0.4.

Method	TPR@1%FPR ↓ (R)	F (Fast)	B (Binoc)	Mean ↓	ASR ↑	AUROC ↓
M0 No Attack	0.22	0.28	0.32	0.27	0.73	0.71
M1 Simple Para	0.14	0.06	0.02	0.07	0.93	0.54
M2 StealthRL-v1	0.08	0.02	0.02	0.04	0.96	0.43
M3 Adv. Para (guided)	0.14	0.06	0.04	0.08	0.92	0.56
M4 AuthorMist	0.02	0.00	0.00	0.01	0.99	0.24
M5 Homoglyph	0.02	0.00	0.00	0.01	0.99	0.41

Table 1: Main results on MAGE (TPR@1%FPR, ASR, and AUROC). Lower TPR/AUROC is better for the attacker; higher ASR is better. R/F/B denote RoBERTa/Fast-DetectGPT/Binoculars.

5 EXPERIMENTAL SETUP

Dataset. We evaluate on MAGE (Machine-generated Text Detection in the Wild) using the test split with 50 human and 50 AI samples, filtered to a 100–500 token window. Li et al. (2024)

Detectors. Evaluation uses three detector families: RoBERTa OpenAI classifier, Fast-DetectGPT (curvature-based), and Binoculars (paired-LM). Mitchell et al. (2023); Bao et al. (2024); Hans et al. (2024)

Methods. We compare six methods: M0 (no attack), M1 (simple paraphrase), M2 (StealthRL-v1 checkpoint), M3 (detector-guided selection inspired by Adversarial Paraphrasing), M4 (AuthorMist), and M5 (homoglyph substitution). Cheng et al. (2025); David & Gervais (2025); Creo & Pudasaini (2025)

Metrics. We calibrate a 1% FPR threshold per detector on human texts and report TPR@1%FPR (lower is better for the attacker), attack success rate (ASR = 1 - TPR@1%FPR), and AUROC. We also report semantic similarity (E5 cosine) to quantify fidelity. All results use one candidate per sample ($n=1$) to match the current evaluation run.

6 RESULTS

Table 1 summarizes detection evasion at 1% FPR. StealthRL reduces mean TPR@1%FPR to 0.04 (vs. 0.27 for no attack) and improves mean AUROC from 0.71 to 0.43. Per-detector TPR@1%FPR drops from 0.22/0.28/0.32 (M0) to 0.08/0.02/0.02 (M2) on RoBERTa/Fast-DetectGPT/Binoculars, indicating strong transfer to the held-out detector family. StealthRL maintains high semantic similarity (0.953 E5 cosine), comparable to simple paraphrasing (0.959), while substantially outperforming the detector-guided baseline (M3).

Ablation (implemented; results pending). We implemented a *guidance-transfer* ablation that mimics the “guidance vs. deploy” setting from Adversarial Paraphrasing: candidate selection is guided by (i) RoBERTa, (ii) Fast-DetectGPT, or (iii) the ensemble mean, and evaluated across all detectors. This isolates how guidance detector choice affects transfer. Results are pending for the final model checkpoints.

7 LIMITATIONS AND SAFETY

Our evaluation currently focuses on three detector families and one primary benchmark; additional detectors (e.g., watermark-based) and datasets can broaden coverage. Adversarial paraphrasing is dual-use: we present StealthRL as a red-teaming tool to expose detector fragility and inform defensive calibration, and we release code anonymized with the expectation of responsible use.

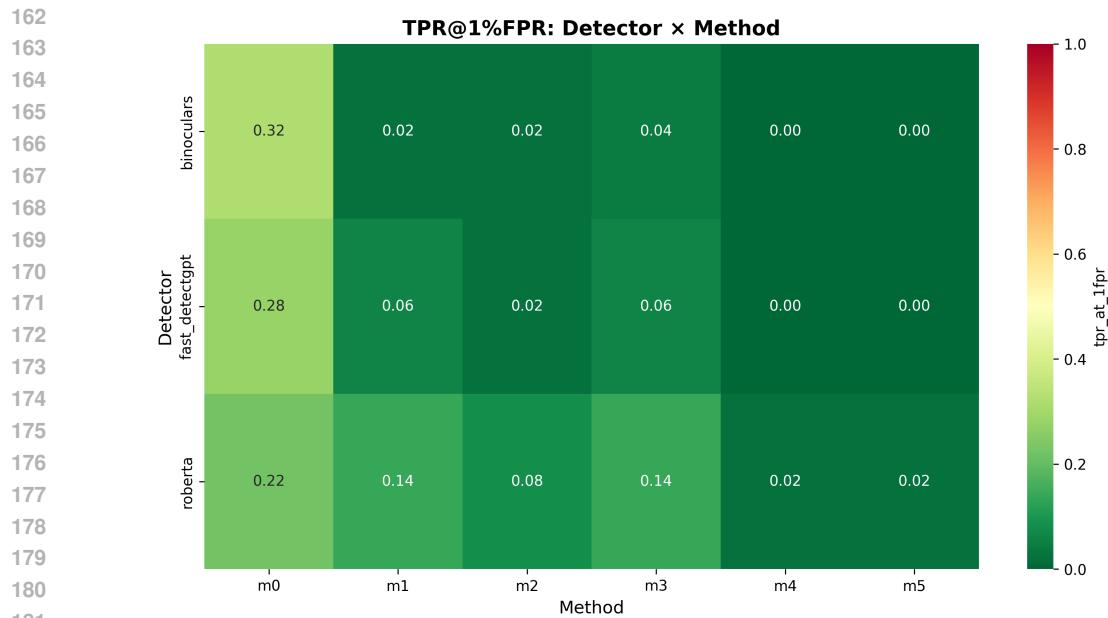


Figure 2: Detector \times method heatmap of TPR@1%FPR. StealthRL shows consistently low TPR across detectors, indicating transfer beyond the training ensemble.

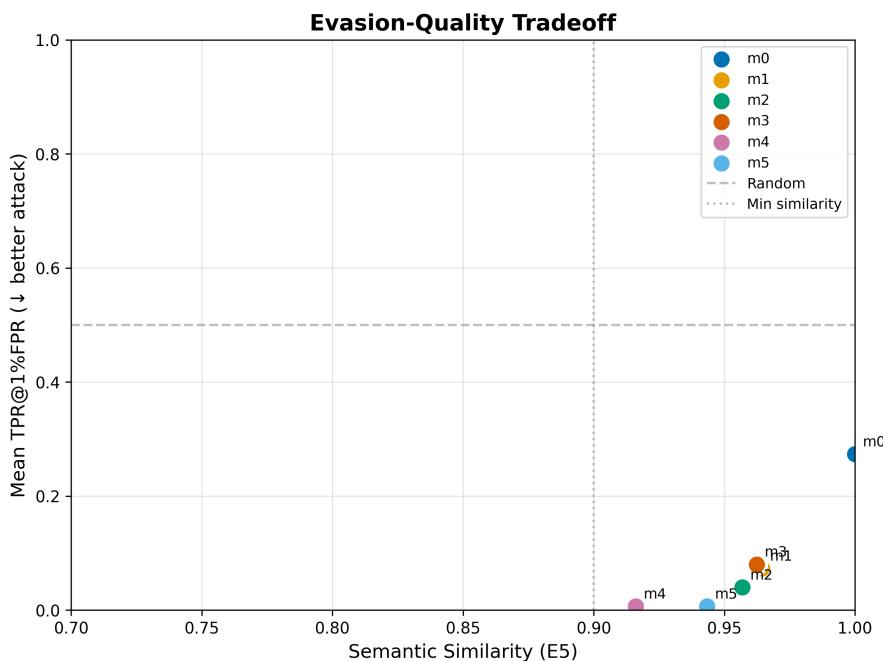


Figure 3: Evasion–quality tradeoff (mean TPR@1%FPR vs. semantic similarity). StealthRL achieves strong evasion while preserving meaning; homoglyph attacks improve evasion but degrade fluency and readability.

216 **8 LLM USAGE DISCLOSURE**
 217

218 Large language models were used to assist with code scaffolding and with editing text for clarity.
 219 All modeling decisions, experiments, evaluations, and interpretations were designed, executed, and
 220 validated by the authors.

222 **9 CONCLUSION**
 223

224 StealthRL provides a practical adversarial RL framework for generating high-fidelity paraphrases
 225 that evade multiple detector families under strict low-FPR constraints. Our results show strong
 226 transfer to a held-out detector while preserving semantic fidelity, and our released evaluation pipeline
 227 supports reproducible red-teaming of AI-text detectors.

229 **REFERENCES**
 230

231 Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. Fast-detectgpt: Efficient
 232 zero-shot detection of machine-generated text via conditional probability curvature, 2024. URL
 233 <https://arxiv.org/abs/2310.05130>.

234 Yize Cheng, Vinu Sankar Sadasivan, Mehrdad Saberi, Shoumik Saha, and Soheil Feizi. Adver-
 235 sarial paraphrasing: A universal attack for humanizing ai-generated text, 2025. URL <https://arxiv.org/abs/2506.07001>.

238 Aldan Creo and Shushanta Pudasaini. Silverspeak: Evading ai-generated text detectors using homo-
 239 glyphs, 2025. URL <https://arxiv.org/abs/2406.11239>.

240 Isaac David and Arthur Gervais. Authormist: Evading ai text detectors with reinforcement learning,
 241 2025. URL <https://arxiv.org/abs/2503.08716>.

243 Abhimanyu Hans, Avi Schwarzschild, Valeria Cherepanova, Hamid Kazemi, Aniruddha Saha,
 244 Micah Goldblum, Jonas Geiping, and Tom Goldstein. Spotting llms with binoculars: Zero-shot
 245 detection of machine-generated text, 2024. URL <https://arxiv.org/abs/2401.12070>.

246 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 247 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.

250 Yafu Li, Qintong Li, Leyang Cui, Wei Bi, Zhilin Wang, Longyue Wang, Linyi Yang, Shuming Shi,
 251 and Yue Zhang. Mage: Machine-generated text detection in the wild, 2024. URL <https://arxiv.org/abs/2305.13242>.

253 Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. De-
 254 tectgpt: Zero-shot machine-generated text detection using probability curvature, 2023. URL
 255 <https://arxiv.org/abs/2301.11305>.

257 **A QUALITATIVE EXAMPLES**
 258

259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289

Original (M0)	StealthRL (M2)
<p>291 During cardio the heart increases its workload and all 292 the body's other systems adjust to help support that 293 endeavor. The blood vessels dilate, the muscles do 294 their best to help pump blood back to the heart, and 295 the lungs work harder to take in oxygen and remove 296 waste gases like carbon dioxide.</p>	<p>In cardio, the heart ramps up its workload, prompting the body's other systems to adapt and assist. Blood vessels widen, muscles strive to push blood back to the heart, and lungs intensify their efforts to absorb oxygen and expel waste gases such as carbon dioxide.</p>
<p>297 The engine has to endure the torque of powering two 298 axles and a drive shaft generally the transfer casing 299 connects to the drive shaft with a ujoint and same for 300 power distribution. The electricity goes into the bat- 301 tery first then is sent to the alternator where it gen- 302 erates voltage that powers all other electrical compo- 303 nents like the lights, radio.</p>	<p>The engine must handle the torque from two axles and a drive shaft, typically linked via a universal joint in the transfer casing for power distribution. Electricity first charges the battery, then flows to the alternator, which converts it into voltage to power essential electrical systems such as lights and the radio.</p>

303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Table 2: Representative paraphrases from the evaluation run (MAGE test split).