# EvoStealer: Differential Evolution for Prompt Template Stealing Against Text-to-Image Synthesis

Anonymous ACL submission

### Abstract

Prompt trading has emerged as a significant intellectual property concern in recent years, where vendors entice users by showcasing sample images before selling prompt templates that can generate similar images. This work investigates a critical security vulnerability: at-007 tackers can steal prompt templates using only a limited number of sample images. To investigate this threat, we introduce PRISM, a prompt-stealing benchmark consisting of 50 templates and 450 images, organized into Easy and Hard difficulty levels. To identify the vul-013 nerabity of VLMs to prompt stealing, we propose EvoStealer, a novel template stealing method that operates without model fine-tuning by leveraging differential evolution algorithms. The system first initializes population sets using multimodal large language models (MLLMs) based on predefined patterns, then iteratively generates enhanced offspring through MLLMs. During evolution, EvoStealer identifies common features across offspring to derive generalized templates. Our comprehensive evaluation conducted across open-source (INTERNVL2-26B) and closed-source models (GPT-40 and GPT-40-MINI) demonstrates that EvoStealer's stolen templates can reproduce images highly similar to originals and effectively generalize to other subjects, significantly outperforming baseline methods with an average improvement of over 10%. Moreover, our cost analysis reveals that EvoStealer achieves template stealing with negligible computational expenses.

## 1 Introduction

034

042

Recent advancements in text-to-image generation (Liu et al., 2024a; Cao et al., 2024), particularly in multimodal large language models (MLLMs) (Liu et al., 2024a; Wang et al., 2024) and diffusion models (Ho et al., 2020; Sohl-Dickstein et al., 2015), have significantly improved image generation performance. However, crafting the perfect prompt to produce desired output images



Figure 1: Top: Illustrating the legitimate development of text-to-image prompt templates. Bottom: Depicting unauthorized extraction of proprietary prompt templates

043

044

045

047

051

057

060

061

062

063

064

065

results remains a meticulous process that requires significant expertise and time investment (Refer to Figure 1(top)). This challenge has catalyzed the emergence of prompt trading, a novel business model exemplified by platforms like PromptBase<sup>1</sup> and LaPrompt<sup>2</sup>. On these platforms, creators upload meticulously crafted prompt templates (viewable post-purchase) alongside multiple sample images (publicly visible). Customers attracted to these samples can purchase the template, then merely modify the subject specification to generate new images that preserve the original stylistic elements. In this context, the platform's copyright and security vulnerabilities raise significant concerns. If attackers reverse-engineer the proprietary templates by analyzing the visible samples, they could significantly compromise sellers' intellectual property rights and threaten the platform's business model (See Figure 1 (bottom)). We term this attack prompt template stealing.

Existing methods for prompt stealing attacks (Shen et al., 2024; Sha and Zhang, 2024; Naseh et al., 2024) focus on reconstructing individ-

<sup>&</sup>lt;sup>1</sup>https://promptbase.com/

<sup>&</sup>lt;sup>2</sup>https://laprompt.com/

ual prompts for each sampled image, rather than 066 recovering a general prompt template for the entire group of sampled images. As a result, the prompts 068 reconstructed by these methods are specific to each image and lack generalizability, which limits their applicability in practical scenarios, as illustrated in Figure 1. For example, in the case of the woman 072 image located in Figure 1, a stolen prompt might include the "golden sun" as a distinctive element. Nevertheless, a comparison with the other three images demonstrates that the "golden sun" is not a shared characteristic among them.

067

071

077

084

091

100

101

102

104

105

106

107

109

110

111

112 113

114

115

116

117

To fill this gap, we build a new and comprehensive benchmark named PRISM, comprising 50 prompt templates stratified across two difficulty levels (Easy and Hard) and spanning 9 distinct subjects, sourced from a specialized prompt trading platform. Utilizing DALL·E 3, we generated 450 images, with each group methodically partitioned into 5 in-domain and 4 out-of-domain images to systematically evaluate both model fitting capability and generalization performance.

Besides, we introduce EvoStealer, a novel template stealing methodology derived from the differential mutation algorithm in evolutionary computation. Our approach strategically leverages mutation and crossover operations within the search space to effectively mitigate overfitting and circumvent local optima, precisely aligning with template stealing objectives. We integrate large language models (LLMs) spanning both open-source and closedsource domains, specifically utilizing InternVL2-26B, GPT-4o, and GPT-4o-mini. By combining these models with a differential evolution algorithm, we generate prompt templates characterized by exceptional stability and robust generalization capabilities. Comprehensive experimental evaluations are conducted across easy and hard difficulty levels. The results demonstrate EvoStealer's remarkable performance: the methodology efficiently reproduces images highly similar to original templates while simultaneously exhibiting strong crosssubject generalizability. This enables large-scale image generation maintaining consistent stylistic characteristics.

Our main contributions are as follows:

(1) To the best of our knowledge, this is the first systematic study on prompt template stealing, revealing its severity as an emerging security threat and empirically demonstrating its significant risk to intellectual property protection;

(2) This study introduces PRISM, the first bench-

mark for prompt template stealing, and EvoStealer, a plug-and-play attack framework that requires no fine-tuning, significantly improving practicality and scalability;

118

119

120

121

122

123

124

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

164

165

(3) We conducted extensive experiments on both open-source models (INTERNVL2-26B) and closed-source models (GPT-40, GPT-40-MINI), with results validating the effectiveness of EvoStealer.

#### 2 **Related Work**

We discuss two lines of related work: the text-toimage prompt stealing attacks and the evolutionary algorithms in LLMs.

### 2.1 Text-to-Image Prompt Stealing Attack

Prompt stealing attacks, or prompt extraction attacks, aim to infer the input from a model's output. A successful attack infringes on intellectual property and poses significant risks to prompt trading platforms in the era of LLMs. However, such attacks are more challenging in text-to-image generation due to the greater uncertainty in image generation compared to text. CLIP Interrogator employs CLIP (Radford et al., 2021) to extract the subject and then selects phrases matching the target image from predefined sets (Udo and Koshinaka, 2023). Shen et al. (2024) fine-tunes 2 models to extract image subjects and modifiers separately, combining them for the attack. Building on this, Naseh et al. (2024) employ GPT-4V to iteratively optimize the prompt, resulting in higher quality. Unlike these works, EvoStealer targets the extraction of a generalizable prompt template, offering greater practical value compared to stealing individual prompts.

### 2.2 Evolutionary Algorithms in LLMs

Recent researches combining evolutionary algorithms with LLMs have demonstrated strong and stable performance across various tasks (Yang et al., 2023; Liu et al., 2023). Some studies leverage the rich domain knowledge and powerful text analysis capabilities of LLMs to accelerate the search process in evolutionary algorithms, particularly in tasks involving complex reasoning (Meyerson et al., 2024; Liu et al., 2024b; Lange et al., 2024; Brahmachary et al., 2024) and interpretability (Chiquier et al., 2025). Conversely, some studies capitalize on the stability of evolutionary algorithms to utilize LLMs for generating higher-quality prompt words (Xu et al., 2022; Prasad et al., 2022; Guo

et al., 2023; Fernando et al., 2023). In this pa-166 per, we use evolutionary algorithms to progres-167 sively generate style descriptors that closely resem-168 ble multiple example images, thereby achieving 169 prompt template stealing.

#### 3 **Data Consturction**

171

172

173

174

176

178

179

180

181

182

183

184

185

186

187

188

191

192

194

195

196

198

199

201

206

210

211

In this section, we introduce the threat model of prompt template stealing, providing a detailed description of the attacker's existing conditions, constraints, and objectives. We then detail our methodology for developing PRISM, a comprehensive benchmark designed to realistically simulate this attack scenario. The specifics are presented below.

## 3.1 Threat Model

The attack scenario is grounded in real-world applications. Attackers have access to two pieces of information from the prompt trading platform: 9 sample images and the generative model (e.g., DALL-E  $3^{3}$  or Midjourney<sup>4</sup>). While attackers can interact with the model via an API, they are not privy to its internal parameters. Their objectives are twofold: first, to generate images that closely resemble, or even replicate, the sample images by using the same subject with the stolen prompt template; and second, to alter the subject within the template and generate images that retain the same style as the sample images.

### 3.2 Benchmark Construction

Currently, no specialized benchmark exists for prompt template stealing research. To address this gap, we introduce PRISM, a novel benchmark comprising 50 freely available prompt templates sourced from PromptBase and LaPrompt. These templates are divided into two equal groups of 25 templates each, categorized as "Easy" and "Hard" based on complexity. Each group encompasses 9 distinct subject categories. We utilize DALL·E 3 as our generation model, combining each prompt template with the 9 subjects to produce 450 unique images. To ensure quality control, we implemented a comprehensive manual review process focusing on two key criteria: subject-prompt alignment and stylistic consistency across template-generated images. For each group, we designated the first 5 generated images as in-domain data to assess similarity between original and stolen prompt. The

remaining 4 images serve as out-of-domain data, enabling evaluation of prompt template generaliza-213 tion capabilities across diverse subjects. For com-214 prehensive details on the benchmark construction 215 methodology, please refer to Appendix A. 216

212

217

218

219

221

222

223

224

225

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

#### **EvoStealer** 4

In this section, we introduce the three main steps of EvoStealer: Image Element Extraction, Differential Evolution, and Fitness Function. The details are presented below.

## 4.1 Image Element Extraction

High-quality prompts for text-to-image generation typically consist of a subject and several modifiers (Liu and Chilton, 2022; Oppenlaender, 2024). The subject defines the object or scene depicted in the image, such as "a woman with a flower crown" or a more intricate description like "Woodland creatures gather around a shimmering pond, surrounded by trees and glowing flowers, creating a peaceful scene". The modifiers specify the style of the image, including aspects such as artistic style and resolution. While multimodal models can accurately identify simple subjects, they often misinterpret complex subjects, mistakenly treating parts of the subject as style modifiers. For instance, in the case of the complex subject "peaceful scene", the model may misinterpret "peaceful" as a style modifier, contaminating the intended description.

To address this issue, we define an image element extraction pattern: <Subject, Modifiers, Supplements>. The subject describes the object or scene, while the modifiers are categorized into four types: Artistic Style, Visual Composition and Structure, Aesthetic and Emotional Atmosphere, and Medium and Material. For further details, please refer to Appendix B. We have imposed the aforementioned constraints on the modifiers to ensure that the model describes the image solely from the relevant perspectives within the four predefined categories. This restriction contributes to the stability and controllability of modifier extraction. However, such a constraint may limit the model's ability to fully capture the diversity of style features. To mitigate this limitation, we incorporate supplements as a compensatory measure. Supplements encompass descriptions outside the four categories and can include individual words, phrases, or even sentences, such as "radiating lines suggesting motion" or "subtle transitions between colors".

<sup>&</sup>lt;sup>3</sup>https://openai.com/index/dall-e-3/

<sup>&</sup>lt;sup>4</sup>https://www.midjourney.com/home



Figure 2: The key steps of EvoStealer in differential evolution, including the identification of differences and commonalities, mutation, mutation addition, and crossover operations.

### 4.2 Differential Evolution

261

262

266

267

271

273

274

275

279

Firstly, we introduce the theory of differential evolution. The process of generating offspring through the differential evolution algorithm is represented using numerical vectors. Initially, each vector in the population is sequentially selected as the base vector, denoted as  $\alpha$ . Then, three individuals,  $x_1, x_2, x_3$ , are randomly chosen from the population to perform the mutation operation. Specifically, the difference between  $x_2, x_3$  is calculated, and this difference undergoes mutation. The mutated difference is then added to  $x_1$  to produce a new vector, denoted as  $\beta$ . The mutation operation is mathematically expressed as:  $\beta = x_1 + F(x_2 - x_3)$ , where F represents the mutation factor, which controls the magnitude of the mutation. Finally, a crossover operation is applied to the vectors  $\alpha$  and  $\beta$  to generate the offspring.

Figure 2 illustrates the differential evolution process implemented in EvoStealer. In Step 1, we differentiate between modifiers and supplements due to their distinct characteristics, particularly in terms of controllability and unpredictability. For modifiers, we focus on identifying the differences between the two sets, while for supplements, we concentrate on their common components. This approach is grounded in the understanding that the uncontrollability of supplements introduces unique features specific to individual images. Additionally, supplements typically contain more tokens than modifiers, which results in a greater influence on the visual representation of the image and, consequently, on the generalization ability of the template. In Step 2, we randomly select an image from the in-domain dataset to influence the mutation process. This strategy serves two purposes: first, it helps filter out modifiers that do not align with the image (e.g., in the case of a surrealistic style image, modifiers such as "cartoon style" are excluded); second, the image, serving as a mutation variable, introduces additional contextual information. As mentioned earlier, the initial version of EvoStealer directly derives image element extraction, which results in an over-reliance on the quality of this extraction. By incorporating the image in Step 2, we enable the population to gain valuable information that may otherwise be overlooked, thus mitigating the drawback of over-dependence on image element extraction. In Step 3, no modifications are made, and the two components are simply combined to generate the mutated description. In Step 4, in contrast to the direct crossover used in genetic

287

288

289

290

291

292

293

294

295

296

298

299

301

302

303

304

305

307

308

309

310

311

Method	DINO	<b>CLIP</b> <sub>img</sub>	$\mathbf{CLIP}_{txt}$	SigLIp <sub>img</sub>	SigLIp <sub>txt</sub>	Average	Human Evaluation
Easy Benchmark							
BLIP 2 (Li et al., 2023)	62.07	79.38	48.35	82.32	52.69	64.96	3.42
CLIP Interrogator	69.93	82.76	54.14	85.86	62.59	70.86	4.02
PromptStealer (Shen et al., 2024)	63.73	77.90	49.21	82.73	61.93	67.10	3.78
EvoStealer (INTERNVL2-26B) EvoStealer (GPT-40-MINI) EvoStealer (GPT-40)	74.68 73.87 <b>75.83</b>	84.46 84.79 <b>85.30</b>	68.94 72.12 <b>74.41</b> Hard Benchr	87.88 88.38 <b>89.14</b> nark	<b>74.93</b> 71.80 72.75	78.18 78.19 <b>79.49</b>	4.32 4.30 4.52
BLIP 2 (Li et al., 2023)	61.16	76.67	46.04	80.51	50.74	63.02	3.24
CLIP Interrogator	66.45	78.26	54.62	82.45	60.78	68.51	3.66
PromptStealer (Shen et al., 2024)	60.01	75.58	47.10	79.20	59.71	64.32	3.48
EvoStealer (INTERNVL2-26B)	70.16	80.63	63.02	84.66	68.14	73.32	4.17
EvoStealer (GPT-40-MINI)	<b>71.05</b>	81.02	67.64	84.88	69.00	74.72	4.12
EvoStealer (GPT-40)	69.24	<b>81.34</b>	<b>70.61</b>	<b>85.28</b>	<b>69.27</b>	<b>75.15</b>	4.24

Table 1: The overall evaluation results for the in-domain data, with the bolded values indicating the best scores.

algorithms, EvoStealer first identifies the common parts between the two individuals. When generating the new offspring, the common parts are fully inherited, while only the differing parts undergo the crossover operation. This design approach strikes a balance between the exploration and exploitation of the algorithm, facilitating effective exploration while ensuring that generalization constraints are preserved.

## 4.3 Fitness Function

313 314

315

316

317

319

320

322

323

325

326

327

331

332

333

334

336

338

340

341

342

343

The fitness function is employed to assess the quality of offspring, with those exhibiting higher fitness scores being retained for progression to the next iteration. While the fitness function does not directly influence the offspring generation, it guides the search direction throughout the evolution process. Our fitness function incorporates both the semantic similarity of the text and the style similarity of the image. Specifically, for each offspring (i.e., a prompt template), we sequentially replace the subject within the template and calculate its semantic similarity with the ground truth. Additionally, we randomly select a subject and use the target model (DALL·E 3) to generate the corresponding image, subsequently calculating the similarity between this generated image and the corresponding image from the in-domain dataset. The mathematical formulation is as follows:

$$\begin{split} F &= \frac{1}{n} \sum_{i=1}^{n} \left( \lambda \left( \frac{\mathbf{T}_{\text{off}}(i) \cdot \mathbf{I}_{\text{gt}}(i)}{\|\mathbf{T}_{\text{off}}(i)\| \|\mathbf{I}_{\text{gt}}(i)\|} \right) \right) \\ &+ (1-\lambda) \left( \frac{\mathbf{I}_{\text{off}} \cdot \mathbf{I}_{\text{gt}}}{\|\mathbf{I}_{\text{off}}\| \|\mathbf{I}_{\text{gt}}\|} \right) \end{split}$$

Where off and gt denote the offspring and ground truth, respectively, and **T** and **I** represent

the text and image embeddings. The parameter  $\lambda$  serves as a balance factor to weight the two similarity measures.

344

345

347

348

350

351

352

353

355

356

357

358

360

361

362

363

364

365

366

367

369

370

371

372

373

## **5** Experiments

We employ PRISM to evaluate the vulnerability of image generation models to prompt template stealing. Following recent works (Shen et al., 2024; Naseh et al., 2024; Huang et al., 2024), we employ subject similarity, style similarity, and semantic similarity metrics to evaluate the performance of image generation models against prompt template stealing (Section 5.3). These metrics demonstrate higher agreement with human annotations than pre-vious approaches. Additionally, we conduct human evaluation to measure the quality of prompt steal-ing.

## 5.1 Baselines

Our baselines encompass models for both caption generation (BLIP-2) and prompt stealing attack (CLIP Interrogator and PromptStealer).

- **BLIP-2**: BLIP-2 (Li et al., 2023) is a multimodal model that aligns text with images using a lightweight Querying Transformer to connect a frozen image encoder with LLMs. In this study, we employ the BLIP-2-opt-2.7b model to generate image descriptions.
- **CLIP Interrogator**: CLIP Interrogator <sup>5</sup> uses CLIP to generate image descriptions, incorporating prompts from preset categories such as artists, flavors, and mediums. It encodes

<sup>&</sup>lt;sup>5</sup>https://github.com/pharmapsychotic/clipinterrogator/tree/main

Method	DINO	<b>CLIP</b> <sub>img</sub>	<b>CLIP</b> <sub>txt</sub>	SigLIp <sub>img</sub>	SigLIp <sub>txt</sub>	Average	Human Evaluation
Easy Benchmark							
CLIP Interrogator	64.02	78.72	53.95	82.98	63.73	68.68	3.80
PromptStealer (Shen et al., 2024)	60.53	75.53	51.37	81.19	61.16	65.96	3.64
EvoStealer (INTERNVL2-26B)	72.93	83.13	68.63	<b>87.24</b>	<b>74.54</b>	77.29	4.36
EvoStealer (GPT-40-MINI)	74.53	83.60	71.87	85.28	73.30	77.71	4.47
EvoStealer (GPT-40)	<b>75.14</b>	<b>83.91</b>	<b>74.18</b>	85.75	73.53	<b>79.10</b>	4.60
Hard Benchmark							
CLIP Interrogator	62.23	69.90	51.66	75.19	58.51	63.50	3.74
PromptStealer (Shen et al., 2024)	58.53	70.42	45.29	74.38	55.07	60.74	3.46
EvoStealer (INTERNVL2-26B)	68.92	78.96	61.29	83.37	67.87	72.08	4.32
EvoStealer (GPT-40-MINI)	67.76	79.55	66.91	84.13	68.84	73.44	4.35
EvoStealer (GPT-40)	<b>67.00</b>	<b>80.50</b>	<b>69.27</b>	<b>84.55</b>	<b>69.79</b>	<b>74.22</b>	4.48

Table 2: The overall evaluation results for the out-of-domain data, with the bolded values indicating the best scores.

both the image and text with the CLIP model, calculates their similarity, and generates the most matching description.

• **PromptStealer**: PromptStealer (Shen et al., 2024) consists of two modules: the Subject Generator, fine-tuned on BLIP to extract image subjects, and the Modifier Detector, a multi-class classifier that selects style modifiers based on similarity to predefined categories. The final prompt is generated by concatenating the subject and selected modifiers.

## 5.2 Experimental Settings

374

375

378

385

Due to the inherent difficulties in subject identifi-386 cation and replacement within BLIP-2-generated prompts, its evaluation is limited to in-domain data only. For both CLIP Interrogator and Prompt-Stealer methods, we first extract subjects and mod-390 ifiers from 5 in-domain samples and concatenate them to create prompts. We then randomly select a prompt and systematically replace its subject with subjects from the out-of-domain group. For PromptStealer, we maintain a threshold value of 0.6. In EvoStealer's implementation, we extract prompt templates from in-domain data and perform sequential subject substitutions, using 9 different 398 subjects to generate the final prompts. Both the population size and generation count are set to 5, 400 with the temperature parameter set to 0 to ensure 401 402 consistent results. We employ SigLIP (Zhai et al., 2023) for fitness score calculations and set  $\lambda$  to 0.5. 403 All image generation is performed using DALL·E 404 3 with a resolution of 1024×1024 and standard 405 quality settings. 406

## 5.3 Evaluation Metric

We adopt the evaluation framework proposed by Huang et al. (2024) and employ the following metrics to assess the performance of EvoStealer and baseline methods: 407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

- **Subject Similarity**: To evaluate the similarity between subjects in paired images, we utilize the self-supervised model DINO (Oquab et al., 2023), as subject comparison is a crucial aspect of image similarity assessment.
- Style Similarity: To measure style consistency, we employ CLIP and SigLIP to extract style features from images generated using stolen prompts and compare them with the original images.
- Semantic Similarity: To assess prompt similarity, we compute the cosine similarity between embeddings of the stolen and target prompts, generated using CLIP and SigLIP.
- Human Evaluation: We recruit 3 external evaluators to rate the similarity between generated and target images on a scale of 1-5, where higher scores indicate greater similarity. For each group, we randomly sample 2 images from both in-domain and out-of-domain categories and calculate average scores. For out-of-domain samples, the evaluation focuses exclusively on style similarity.

## 5.4 Main Results

Tables 1 and 2 present comparative performance evaluations between EvoStealer and baseline methods using both in-domain and out-of-domain data. The results demonstrate that EvoStealer consistently outperforms baseline approaches across all
evaluation metrics.

Performance on in-domain data. EvoStealer 442 443 outperforms other methods in both the Easy and Hard categories. For example, EvoStealer (GPT-444 40) leads the second-best method, CLIP Interroga-445 tor, by 8.63% and 6.64% on the two datasets, 446 demonstrating its ability to generate more accurate 447 prompt templates and better stealing performance. 448 Notably, EvoStealer excels in textual semantic com-449 parison, as its prompts are significantly more ef-450 fective. CLIP Interrogator and PromptStealer rely 451 on simple concatenation of [subject] and modifiers, 452 limiting variability. Additionally, CLIP's length 453 restriction hampers modifier extraction. In contrast, 454 EvoStealer generates diverse templates iteratively, 455 avoiding these limitations. This aligns with the 456 findings of Naseh et al. (2024). 457

**Performance on out-of-domain data.** As shown 458 in Table 2, EvoStealer outperforms other meth-459 ods, especially on out-of-domain data, where it 460 demonstrates a larger advantage compared to in-461 domain data. For instance, EvoStealer (GPT-40) 462 leads by more than 10% across all data types, indi-463 cating better generalization of stolen templates to 464 different subjects. Furthermore, as seen in Table 1, 465 466 EvoStealer's performance on out-of-domain data 467 shows minimal degradation, while CLIP Interrogator and PromptStealer experience average degrada-468 tions of 3.60% and 2.36%, respectively. This is due to EvoStealer's effective template stealing by ex-470 471 tracting common features across multiple images.

**Comparison of performance across different** 472 models. As shown in Tables 1 and 2, GPT-40 473 outperforms the other models, followed by GPT-474 40-mini, with InternVL2-26B performing slightly 475 worse. However, the performance differences 476 among these models are minimal. This is pri-477 marily due to EvoStealer's reliance on the models' 478 text and image analysis capabilities, indicating that 479 EvoStealer is highly compatible and not dependent 480 on a specific multimodal model. 481

#### 6 Analysis

482

In this section, we analyze the effects of
EvoStealer's components, the iteration number, and
the experimental costs.

Method	InDom.	OutDom.	Average
Ours	77.32	76.67	77.00
w/o. supp.	73.56	74.85	74.21
w/o. img.	75.89	75.57	75.73

Table 3: Results of the ablation study: Impact of omitting supplements in the extraction pattern (w/o supp.) and excluding image similarity in the fitness function (w/o img.), with the model employed being GPT-4.

#### 6.1 Ablation Study

We remove the supplements from the extracted templates and the image similarity evaluation from the fitness function to examine their impact on EvoStealer. The results are shown in Table 3. As observed, removing either module results in decreased performance, with a more significant drop when supplements are removed-specifically, an average similarity reduction of 2.79%. This is because supplements provide additional details, such as image features and style information. As noted in Section 4.2, supplements are longer than individual modifiers, so their removal has a more pronounced effect on visual performance. A comparison of the performance before and after removing supplements is provided in Appendix D. Removing the image similarity evaluation from the fitness function causes a performance decrease of 1.27%, suggesting that including the comparison between the generated and target images in the fitness function helps guide the evolutionary process and accelerate convergence.



Figure 3: The convergence curve of EvoStealer, with the left half showing changes in fitness score and the right half depicting performance changes of the optimal prompt template for in-domain and out-of-domain data.

## 6.2 Effect of Number of Iterations

We select 10 groups of easy and 10 groups of hard cases to examine EvoStealer's convergence (we use GPT-40 as the analysis model), with results shown in Figure 3. The left section of the figure

512

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505



(b) Out-of-domain Examples

Figure 4: The attack results of EvoStealer compared to three baseline methods on both easy and hard examples. (a)-(d) represent EvoStealer, CLIP-Interrogator, PromptStealer, and BLIP2, respectively.

displays changes in the fitness score as evolution progresses, while the right section shows changes in the scores of the optimal templates for both indomain and out-of-domain data. We observe that as evolution progresses, both the optimal and average fitness scores gradually increase, indicating that EvoStealer generates offspring with higher adaptability. The performance of the prompt templates steadily improves for both in-domain and out-ofdomain data. Two examples are provided in Appendix E

#### 6.3 Cost Analysis

513

514

515

516

517

518

519

520

521

522

To assess the practicality of EvoStealer, we ana-525 526 lyzed the cost of stealing a prompt template. The primary overhead of EvoStealer consists of three 527 components: population initialization, differential 528 evolution (including the fitness function), and image synthesis. A detailed cost estimation process is 530 provided in Appendix F. The results indicate that EvoStealer requires 144 API calls, generates 34 532 images (including 9 final synthesized images), and consumes approximately 119.1k tokens, amounting to a total cost of \$1.70. While this is lower than the platform's pricing range of \$3-9, the cost advantage is not substantial. However, as demonstrated in the ablation study in Section 6.1, costs can be 539 further reduced by using open-source models or omitting image similarity calculations in the fit-540 ness function, enabling near-zero-cost stealing. Al-541 though this cost-reduced version performs slightly worse than the full EvoStealer model, it still signif-543

icantly outperforms alternative approaches.

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

565

566

567

568

569

570

571

572

573

## 7 Case Study

To clearly demonstrate EvoStealer's advantages over baseline methods, we select an easy and a hard example for case study, with the results shown in Figure 4. The results show that, on in-domain data, EvoStealer generates images that closely match the style of the original images, with all four synthesized images maintaining stylistic consistency. In contrast, the four images generated by the other baseline methods exhibit significant style variation. On out-of-domain data, EvoStealer maintains the same style as in-domain images, successfully achieving subject generalization. In contrast, the other baseline methods fail to generalize. Additionally, we analyze three distinct failure cases (see Appendix G for details).

## 8 Conclusion

This paper investigates prompt template stealing—whether attackers can extract generalizable templates that maintain stylistic consistency using minimal sample images. To explore this scenario, we provide PRISM, a two-tier benchmark consisting of 50 templates and 450 images, organized into Easy and Hard difficulty levels. We also introduce EvoStealer, a template stealing method that combines differential evolution algorithms with MLLMs, enabling template stealing without the need for fine-tuning. Extensive experiments and analysis validate its effectiveness and practicality.

## 9 Limitations

574

575

582

583

584

585

587

590

591

592

593

596

598

604

607

611

612

613

615

616

617

618

The current implementation of EvoStealer and benchmark presents several methodological limitations:

- 1. EvoStealer's MLLM-based design offers simplified implementation without fine-tuning requirements and maintains robust performance across open datasets. However, this approach inherently limits the system's maximum performance to the capabilities of the underlying MLLMs.
  - Resource constraints restricted our benchmark to DALL·E-3 generated images, excluding other prominent models like Midjourney and Stable Diffusion. Nevertheless, the current benchmark adequately evaluates stealing method performance, with planned expansion to additional models in future work.
  - The benchmark's single-subject design facilitates comparative analysis but does not address multi-subject templates in real-world applications—a limitation to be addressed in subsequent research.

### 10 Ethical Considerations

EvoStealer's ability to extract prompt templates from minimal image examples enables attackers to generate multiple stylistically similar images through minor template modifications, posing significant risks to creators' intellectual property. This research highlights this security vulnerability, as understanding such threat models is essential for developing effective countermeasures.

While watermarking offers some protection, its implementation on trading platforms presents practical challenges. Watermarks can obscure image details, potentially deterring buyers or leading to customer dissatisfaction when purchased prompts fail to meet expectations. Our findings suggest that limiting the number of displayed images to 2-4 examples provides a simple yet effective defensive strategy.

Future research should prioritize developing robust protection mechanisms to safeguard both creators' rights and the integrity of the AI-generated content marketplace.

#### References

Shuvayan Brahmachary, Subodh M Joshi, Aniruddha Panda, Kaushik Koneripalli, Arun Kumar Sagotra, Harshil Patel, Ankush Sharma, Ameya D Jagtap, and Kaushic Kalyanaraman. 2024. Large language model-based evolutionary optimizer: Reasoning with elitism. *arXiv preprint arXiv:2403.02054*. 619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

- Pu Cao, Feng Zhou, Qing Song, and Lu Yang. 2024. Controllable generation with text-to-image diffusion models: A survey. *arXiv preprint arXiv:2403.04279*.
- Mia Chiquier, Utkarsh Mall, and Carl Vondrick. 2025. Evolving interpretable visual classifiers with large language models. In *European Conference on Computer Vision*, pages 183–201. Springer.
- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. 2023. Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv* preprint arXiv:2309.16797.
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. 2023. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. *arXiv preprint arXiv:2309.08532*.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840– 6851.
- Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. 2024. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818.
- Robert Lange, Yingtao Tian, and Yujin Tang. 2024. Large language models as evolution strategies. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, pages 579–582.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Fei Liu, Xi Lin, Zhenkun Wang, Shunyu Yao, Xialiang Tong, Mingxuan Yuan, and Qingfu Zhang. 2023. Large language model for multi-objective evolutionary optimization. *arXiv preprint arXiv:2310.12541*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024a. Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Shengcai Liu, Caishun Chen, Xinghua Qu, Ke Tang, and Yew-Soon Ong. 2024b. Large language models as evolutionary optimizers. In 2024 IEEE Congress on Evolutionary Computation (CEC), pages 1–8. IEEE.

- 674 675
- 678
- 686
- 687
- 691
- 692
- 695
- 700 701
- 702 704 705

709

- 710 711 712 713
- 714
- 715 716
- 717
- 718
- 719 720 721

- 722 723
- 725 726

- Vivian Liu and Lydia B Chilton. 2022. Design guidelines for prompt engineering text-to-image generative models. In Proceedings of the 2022 CHI conference on human factors in computing systems, pages 1–23.
- Elliot Meyerson, Mark J Nelson, Herbie Bradley, Adam Gaier, Arash Moradi, Amy K Hoover, and Joel Lehman. 2024. Language model crossover: Variation through few-shot prompting. ACM Transactions on Evolutionary Learning, 4(4):1-40.
- Ali Naseh, Katherine Thai, Mohit Iyyer, and Amir Houmansadr. 2024. Iteratively prompting multimodal llms to reproduce natural and ai-generated images. arXiv preprint arXiv:2404.13784.
- Jonas Oppenlaender. 2024. A taxonomy of prompt modifiers for text-to-image generation. Behaviour & Information Technology, 43(15):3763–3776.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. 2023. Dinov2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193.
- Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. 2022. Grips: Gradient-free, edit-based instruction search for prompting large language models. arXiv preprint arXiv:2203.07281.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748-8763. PMLR.
- Zeyang Sha and Yang Zhang. 2024. Prompt stealing attacks against large language models. arXiv preprint arXiv:2402.12959.
- Xinyue Shen, Yiting Qu, Michael Backes, and Yang Zhang. 2024. Prompt stealing attacks against {Textto-Image} generation models. In 33rd USENIX Security Symposium (USENIX Security 24), pages 5823-5840.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In International conference on machine learning, pages 2256–2265. PMLR.
- Honori Udo and Takafumi Koshinaka. 2023. Image captioners sometimes tell more than images they see. arXiv preprint arXiv:2305.02932.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191.

Hanwei Xu, Yujun Chen, Yulun Du, Nan Shao, Yanggang Wang, Haiyu Li, and Zhilin Yang. 2022. Gps: Genetic prompt search for efficient few-shot learning. arXiv preprint arXiv:2210.17041.

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

774

775

776

- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2023. Large language models as optimizers. arXiv preprint arXiv:2309.03409.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 11975-11986.

#### Data Collection Α

Our data collection and preprocessing pipeline consists of several systematic steps. Initially, we collect 100 free templates (excluding specific images) from PromptBase and LaPrompt, subsequently generating 900 corresponding images using DALL-E 3. Through manual curation, we eliminate templates that prove challenging to reproduce, such as those containing specific artistic style descriptors (e.g., Arkhip Kuindzhi). We then perform deduplication to remove images with highly similar styles and subjects, thereby preventing evaluation bias from data redundancy. The subsequent quality control process encompasses two primary aspects: subject alignment verification and style consistency assessment. When anomalous data is identified, we regenerate images using DALL-E 3 until they meet our quality criteria. Finally, we categorize our dataset into Easy and Hard classifications. The Hard category is characterized by: uncommon modifiers, abstract subject descriptions and rich image details. The token distribution of the complete dataset is illustrated in Figure 5, while Table 4 presents a detailed breakdown of token statistics for both Easy and Hard categories.

#### **Extraction Pattern Detail** B

Describing the style of an image requires including different perspectives. The style description of EvoStealer includes four categories: Artistic Style, Visual Composition and Structure, Aesthetic and Emotional Atmosphere, and Medium and Material.

- Artistic Style: Include Genre, Era or Historical Style, Cultural and Technological Style.
- Visual Composition and Structure: Include Composition and Layout, Form and Structure, Scale, Movement, Perspective, Pattern and Ornamentation and Detail Level.

	Ε	asy	Hard		
Number	Subject	Modifier	Subject	Modifier	
Min.	1	23	1	16	
Max.	27	154	24	107	
Avg.	8.03	65.00	3.60	43.64	

Table 4: Token Statistics for Easy and Hard Benchmarks.

• Aesthetic and Emotional Atmosphere: Include Tone and Atmosphere, Emotional Atmosphere, Lighting and Shadow Effects,

777

778

779

782

785

786

790

791

794

797

798

799

 Medium and Material: Include Medium, Material, Technique, Texture, Surface, Color Palette, Brushwork, Line Quality, Strokes, Layering, Transparency, Opacity and Resolution.



Figure 5: Token frequency distribution of the dataset

### **C** Human Evaluation

We implement a rigorous human evaluation protocol using a blinded manual scoring approach. Each evaluator is presented with the original benchmark images alongside extracted results from all methods, comprising two in-domain and two out-ofdomain images per set. To maintain objectivity, evaluators are blinded to the generation methods and conduct their assessments independently, without inter-evaluator communication. The evaluation criteria are differentiated by image category:

- For in-domain data: Evaluators assess both subject matter and stylistic similarity to measure template reproduction fidelity
- For out-of-domain data: Evaluation focuses

exclusively on stylistic similarity to assess template generalization capability

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

Images are rated using a 5-point Likert scale, with higher scores indicating greater similarity. Final results are reported as mean scores across all evaluators. The detailed scoring criteria are presented below.

- 1. **Completely Different:** The generated image exhibits no discernible similarities to the original, presenting entirely distinct content and stylistic elements.
- 2. **Barely Similar:** While minimal thematic or elemental commonalities may exist between the original and generated images, they demonstrate significant divergence in both content and stylistic execution.
- 3. **Somewhat Similar:** The generated image maintains recognizable correspondence to the original's content or subject matter, although notable stylistic variations are present.
- 4. **Closely Similar:** The generated image demonstrates substantial fidelity to the original's content and subject matter, with only minor compositional variations.
- 5. Very Similar: The generated image achieves near-identical reproduction, maintaining high fidelity to the original's content, style, and intricate details.

### **D** Ablation Comparison

Our extraction template incorporates controllable Subjects and Modifiers, complemented by a flexible Supplements module designed to address potential gaps in subject and modifier extraction. Figure 6 demonstrates the impact of the Supplements module on EvoStealer's effectiveness.

The first case study illustrates how the Supplements module enhances feature detection. While analyzing images individually may cause oversight of shared characteristics—such as the presence of petals in 'a floating umbrella covered in flowers'—the Supplements module successfully captures these overlooked elements in Subject, thereby improving extraction accuracy. In the second case, the module demonstrates its ability to detect visual attributes that are overlooked by predefined modifier categories, such as 'dark yellow tone' within the 'Visual Composition and Structure' categories.



Figure 6: Three examples are used to demonstrate the impact of removing supplements. "w/o. supp." represents the removal of supplements extracted from the pattern.

The third case exemplifies the module's capacity to identify fundamental aesthetic properties like symmetry, which fall outside established modifier categories. These examples highlight how the Supplements module's flexibility enables the detection of additional key features, ultimately enhancing the quality of image generation.

#### Е **Evolution Progress**

847

848

852

857

861

Figure 7 presents the iterative results of EvoStealer on in-domain data across two distinct styles. The figure demonstrates that with each iteration, the generated images progressively converge toward the ground truth style. This progression indicates that EvoStealer successfully refines the quality of style descriptors throughout its iterative process, resulting in images that increasingly approximate the target style. The visual comparison clearly illustrates the algorithm's capacity to incrementally improve stylistic fidelity through successive refinements.

#### F **Cost Estimate**

The execution process of EvoStealer comprises three main stages: population initialization, differential evolution (including the fitness function), and image synthesis. We assess the cost from three perspectives: API call frequency, token consumption, and image generation. While API calls and image generation can be accurately and directly measured, token consumption is estimated. Given the instability of the model's output, only the input

862

863

864

865

866

867

868

869

870

871

872

873

874

875



Figure 8: Three failure cases in EvoStealer.

portion is estimated. For this analysis, we evaluate the cost of stealing a prompt template using EvoStealer, based on GPT-40.

877

878

879

882

During the population initialization phase, EvoStealer performs two key operations: image element extraction (which generates <subject, modifiers, supplements> triples) and initial template synthesis. On average, this requires 10 calls to GPT-40, consuming 1.6k tokens, with an estimated cost of approximately \$0.04. In the differential evolution phase, EvoStealer performs operations such as difference and commonality identification, mutation, mutation addition, and crossover. Additionally, for each offspring, template synthesis and image generation are required for both creation and evaluation. On average, this phase involves 125 API calls, consumes 117.5k tokens, and generates 25 images, resulting in a total cost of approximately \$1.30. In the image synthesis phase, only the optimal template is used to generate 9 images. This requires 9 API calls and 9 image generations, totaling \$0.36. Thus, the overall cost amounts to approximately \$1.70.

## G Failure Cases

900

901

902

903 904

905

907 908

909

910

911

912

913

914

915

916

917

918

919 920

921

922

923

924

925

926

929

In this section, we will examine several typical failure cases. These failures stem either from the complexity of the images themselves and vague descriptions, or from the inherent limitations of the current EvoStealer method. Figure 8 illustrates representative examples.

A primary limitation is the system's inadequate interpretation of specific artistic styles. Analysis of PromptBase and LaPrompt platforms reveals that many prompt templates incorporate stylistic modifiers, such as "Arshile Gorky style," "Disney style," and "Renaissance style." However, the system struggles to accurately identify and replicate the distinctive characteristics of individual artists' techniques or historical artistic movements, resulting in significant stylistic disparities between generated and source images.

A second limitation concerns text recognition capabilities. The current EvoStealer implementation lacks explicit protocols for extracting textual elements from images. Despite MLLMs' inherent text recognition capabilities, this functionality remains underutilized in the present version—a limitation scheduled for address in future iterations.

The third limitation involves comprehensive detail preservation. When processing images with complex color palettes and rich content, EvoStealer may fail to capture fine-grained features, leading to degraded quality in the resultant prompt templates.