ATTACK AS DEFENSE: RUN-TIME BACKDOOR IMPLANTATION FOR IMAGE CONTENT PROTECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

As generative models achieve great success, tampering and modifying the sensitive image contents (i.e., human faces, artist signatures, commercial logos, etc.) have induced a significant threat with social impact. The backdoor attack is a method that implants vulnerabilities in a target model, which can be activated through a trigger. In this work, we innovatively prevent the abuse of image content modification by implanting the backdoor into image-editing models. Once the protected sensitive content on an image is modified by an editing model, the backdoor will be triggered, making the editing fail. Unlike traditional backdoor attacks that use data poisoning, to enable protection on individual images and eliminate the need for model training, we developed the first framework for run-time backdoor implantation, which is both time- and resource- efficient. We generate imperceptible perturbations on the images to inject the backdoor and define the protected area as the only backdoor trigger. Editing other unprotected insensitive areas will not trigger the backdoor, which minimizes the negative impact on legal image modifications. Evaluations with state-of-the-art image editing models show that our protective method can increase the CLIP-FID of generated images from 12.72 to 39.91, or reduce the SSIM from 0.503 to 0.167 when subjected to malicious editing. At the same time, our method exhibits minimal impact on benign editing, which demonstrates the efficacy of our proposed framework. The proposed run-time backdoor can also achieve effective protection on the latest diffusion models.

031

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

032 033 034

035

1 INTRODUCTION

In recent years, advances in generative models have been remarkable, demonstrating exceptional performance in image-editing tasks. Numerous editing models and their variant applications, such 037 as LaMa(Suvorov et al. (2022); Yu et al. (2023)) and Stable Diffusion(Rombach et al. (2022b); Lugmayr et al. (2022)), have been developed, expanding the capabilities of image modification. Image inpainting stands out as a typical task of image editing, which aims to repaint a specific area 040 in the image. The image inpainting receives an image-mask pair and uses the mask as positional 041 guidance to edit the image. Text-guided editing/inpainting receives an image-mask-text triple as 042 input and repaints the mask area of image following the text instruction. As these image inpaint-043 ing models are widely used, they can be maliciously employed for copyright logo erasuring, face 044 replacement, background replacement, and other abusing and tampering of image content, which brings new challenges to the field of artificial intelligence security.

Backdoor attack(Liu et al. (2018)) is a well-known threat on deep neural networks, which has received more attention with the rise of generative models such as diffusion models and large language models Ramesh et al. (2021); Achiam et al. (2023). Attackers can implant a backdoor into deep learning models by poisoning training data or manipulating model weights Chen et al. (2023); Chou et al. (2023; 2024); Zhai et al. (2023) in the training stage, then activate the backdoor in the inference stage through a predefined trigger in the input. In this work, we innovatively leverage the backdoor of image-editing models as a defense to protect sensitive image content from being tampered. In other words, when an editing model is utilized to modify a protected region on the image, the model backdoor will be activated and the modification could fail.



Figure 1: Paradigm comparison of traditional backdoor framework (top row) and the proposed Runtime implant framework (bottom row). The traditional approach requires obtaining a compromised
model via Trojan training prior to deployment, with the backdoor being activated during the inference stage. In contrast, our runtime implant framework bypasses the need for prior poisoning,
enabling the backdoor to be activated solely during inference. Conventional backdoor typically relies on explicit trigger, and our method leverages region-aware trigger, which is imperceptible and
can be activated with editing location.

082

083 However, the unique characteristics of backdoor implantation present key challenges for this attackas-defense process: as shown in the top row of Figure 1, **0** it is difficult for photographers or 084 creators to finetune or manipulate an image editing model, as attackers did before, to implant the 085 backdoor. Even if it is possible, it is still hard to force others to use the released backdoor-ed model for image editing. **2** The backdoor trigger on the image should be invisible, which can avoid dis-087 turbing the original image content created by the artists. Since the image editing operation could 088 be diverse (e.g., various mask shapes and locations), backdoor activation must remain robust when 089 the protected region is under modification. Additionally, editing other unprotected insensitive re-090 gions should be allowed to minimize the negative impact on legal modifications. To address these 091 challenges, we propose our novel run-time backdoor implantation framework capable of implant-092 ing backdoors into image editing models without requiring model training. Figure 1 illustrates the overview of the attack framework, where the top row indicates the traditional backdoor implanta-094 tion and the triggering process while the bottom row demonstrates ours. Specifically, for **0**, instead of fine-tuning the subject model, we generate protective noises that can induce inherent backdoors 095 within the editing model when applied to the input images. For **2**, we establish the protected area 096 on the image as the trigger for the backdoor, thus avoiding the risk of introducing visually perceptible disturbances, caused by explicit triggers, in the image content. To address $\boldsymbol{\Theta}$, we propose three 098 collaborative optimization objectives to achieve region-aware backdoor activation. The implant loss is designed to activate the backdoor target when the edited area aligns with the protected region. 100 The incomplete-trigger loss aims to activate the backdoor robustly when the edit area partially over-101 laps with the protected region. Additionally, the hide loss seeks to mitigate the impact of editing 102 operations on other unprotected regions. 103

In summary, our contributions are as follows:

105

107

• We are the first to propose a run-time backdoor implantation framework for image models, enabling the injection of backdoors through protective noise during the inference stage, without requiring model retraining.

- We introduce the use of a protected area as a trigger mechanism to ensure minimal perceptible interference in the image, making the backdoor activation less detectable.
- We design a region-aware backdoor activation mechanism with three collaborative optimization objectives that allow the protective noise to reliably activate the backdoor across various image editing operations and remain minimal impact on legal image modifications.
- We evaluate the proposed framework on LaMa using inpainting datasets across six distinct scenarios. In the task of comprehensive inpainting, our approach reduces structural similarity index(SSIM) from 0.503 to 0.167, or improves CLIP-based fréchet inception distance(CLIP-FID) from 12.72 to 39.91. Furthermore, extended experiments demonstrate that our method effectively implants backdoors in diffusion models.

2 RELATED WORKS

108

109

110 111

112

113

114 115

116

117

118

119

120 121 122

123 124

125 **Image editing** involves the modification or enhancement of images using various models collec-126 tively known as image editing models. These models encompass tasks such as image inpainting Yu 127 et al. (2023); Lugmayr et al. (2022); Yang et al. (2023); Corneanu et al. (2024), style transfer Kwon & 128 Ye (2022); Zhang et al. (2023); Deng et al. (2022), and text-guided editing Tao et al. (2023); Nichol 129 et al. (2021); Wang et al. (2023). Image inpainting is a key subtask focused on filling missing or 130 corrupted regions of an image, and can be classified into two types: mask-guided inpainting and 131 text-guided inpainting. In mask-guided inpainting Suvorov et al. (2022); Rombach et al. (2022b), models use an image-mask pair as input, where the mask defines the region to be reconstructed. The 132 inpainting process relies on the surrounding unmasked regions to fill the masked area in a contextu-133 ally coherent manner. In contrast, text-guided inpainting Manukyan et al. (2023); Ni et al. (2023); 134 Xie et al. (2023) uses an image-mask-text triplet, where the missing region is filled based on both 135 the surrounding image context and the textual input. Given the widespread use of inpainting-based 136 image editing models, these functions can be exploited for tasks such as removing copyright logos, 137 replacing faces, or altering backgrounds-raising significant concerns about image manipulation 138 and misuse, which pose new challenges to AI security. Zhang et al. (2024) and Hu et al. (2023) are 139 designed to ensure the authenticity of digital images by embedding imperceptible protection signals 140 that can detect tampering and unauthorized modifications, but we want to prevent the tampering 141 process itself. Although Salman et al. (2023) has performed adversarial learning on images to make 142 them resistant to manipulation by diffusion models, it is not designed in the backdoor setting, where the protection will only be activated on specific conditions (i.e., triggers). Hence, it is not able to 143 differentiate between malicious and benign editing operations, and could overreact on legal modi-144 fications. In light of this, we aim to strengthen protection against the misuse of fine grained image 145 content (e.g. watermark, human face) while allowing for benign editing in authorized changeable 146 areas (e.g., removing garbage in the background). 147

Backdoor attack(Liu et al. (2018); Hayase & Oh (2022); Qi et al. (2023)) is a classic threat to 148 neural networks. Attackers usually implant backdoor into the model by data poisoning(Chen et al. 149 (2017); Liu et al. (2020); Zhai et al. (2023)) during the training or fine-tuning phase, and then 150 activate the backdoor by inputting samples containing trigger in the inference phase. In contrast to 151 adversarial attacks(Cheng et al. (2024a;b)), which primarily exploit model vulnerabilities during 152 inference, backdoor attacks are characterized by a predefined target orientation and require model 153 training. Previous works(Chou et al. (2023); Sun et al. (2024); Chou et al. (2024); Li et al. (2024)) 154 proposed to backdoor diffusion models by adding a special pattern as a trigger and train the model to 155 make incorrect outputs when this trigger is encountered. However, since this modification is explicit 156 to the model parameters, it makes the backdoor easily detectable. Yet, some backdoor detection 157 methods(An et al. (2024; 2023); Hao et al. (2024)) for diffusion models have been proposed to 158 detect the data distribution of models, especially on models released by third parties. Therefore, 159 we propose a run-time backdoor framework to implant backdoor without manipulating the model, making it much more stealthy. We regard the position of the editing mask as the trigger. When 160 protected area of the image is attempted to be edited, the backdoor is activated, leading to failed and 161 unrealistic outputs.



Figure 2: Optimization target of our run-time backdoor. We use three different edit regions as input to guide the optimization of protective noise. In the first row, the entire trigger region is employed to optimize the implant loss $\mathcal{L}_{implant}$. The second row utilizes an expanded trigger region to address incomplete activation loss $\mathcal{L}_{incomplete}$. The hide loss \mathcal{L}_{hide} in third row applies editing to regions without trigger to minimize interference with benign modifications, thereby preserving the image's editability on non-trigger inputs.

3 Method

175

176

177

178

179

180 181

182 183

185

187 188

189

The overview of our proposed run-time backdoor framework is shown in Figure 1. In Section 3.1, we introduce the threat model and the scenario at first. Then, we describe the run-time backdoor implant framework in Section 3.2 and region-aware backdoor activation mechanism in Section 3.3.

3.1 THREAT MODEL

Image inpainting has gained widespread adoption as researchers increasingly release pre-trained 190 image editing models to the public. In this work, we conceptualize three key parties involved in 191 the process: The *developer* of inpainting models, responsible for training these models and making 192 both the code and model weights available on platforms such as *hugging face*, allowing any user 193 to download and utilize them. The user, who seeks to apply publicly available inpainting models 194 and may download images shared on public platforms. The defender, representing image content 195 protectors and copyright holders, such as artists and photographers, who aim to prevent unauthorized 196 editing or misuse of their images. 197

In our proposed implant framework, the *developer* solely provides pre-trained inpainting models and does not engage in image editing or in the insertion of backdoors. The *user* may download images posted by content owners and employ the publicly available inpainting models for image manipulation. However, the user group may include malicious actors who exploit sensitive content for malicious purposes, such as illegal profiteering or infringement.

The *defender*, in turn, employs a run-time backdoor implantation mechanism in inpainting models, introducing protective noise perturbations into their images. A specific location within the image is designated as a trigger, producing a modified image that is visually indistinguishable from the original. The defender (the image owner) can then publicly post the protected image. If a malicious user attempts to edit the protected content using the inpainting model, the model will produce a distorted output, thus safeguarding the image from unauthorized manipulation.

209 210

3.2 RUN-TIME BACKDOOR IMPLANTATION

In this section, we formalize our *run-time backdoor* method. We describe the defense scenario using the image inpainting task, and the method can be easily transferred to applications of inpainting models such as Stable Diffuison(Rombach et al. (2022b)).

We first summarize the training paradigm of state-of-the-art image editing models, which can be described as follows: given an image editing model \mathcal{IM} , the inpainting function of \mathcal{IM} is trained

with loss function:

$$\mathcal{L}_{IM} = \alpha \cdot \mathcal{L}_{recon} + \beta \cdot \mathcal{L}_{adv} + \gamma \cdot \mathcal{L}_{perc}, \tag{1}$$

(2)

In this formulation, α , β , and γ serve as hyperparameters that balance the contributions of various components of the loss. Specifically, \mathcal{L}_{adv} corresponds to the adversarial loss associated with the generative adversarial training process. The reconstruction loss, \mathcal{L}_{recon} , assesses the fidelity of the inpainted regions at the pixel level, while the perception loss, \mathcal{L}_{perc} , captures the differences in high-level features between the generated and target images, ensuring perceptual consistency. Specifically, the basic losses \mathcal{L}_{recon} and \mathcal{L}_{perc} for inpainting function can be jointly formalized as:

224

228

where the x and $\mathcal{G}(x)$ denote the input image and the generated result. Generally, traditional backdoor method requires to trojan the model \mathcal{IM} by training the inpainting loss $\mathcal{L}_{inpaint}$ with an extra objective:

 $\mathcal{L}_{inpaint} = \mathbb{E}_{x,\mathcal{G}} \left[\left\| x - \mathcal{G}(x) \right\|^2 \right],$

$$\mathcal{L}_{trojan} = \begin{cases} \mathbb{E}_{x \sim \mathcal{S}, \mathcal{G}} \left[\|x - \mathcal{G}(x)\|^2 \right], & \text{if } S = Clean \\ \mathbb{E}_{x \sim \mathcal{S}, y, \mathcal{G}} \left[\|y - \mathcal{G}(\mathcal{T}(x))\|^2 \right]. & \text{if } S = Poisonous \end{cases}$$
(3)

Here, y represents the output target associated with the backdoor mechanism. The function $\mathcal{T}(x)$ denotes the image input containing the trigger, which is designed to activate the backdoor target. The training set distribution, denoted by S, consists of both clean data and poisoned data, reflecting a mixture of benign and manipulated samples used during model training.

However, the above paradigm of traditional backdoor is not suitable for our attack-as-defense sce-237 nario. Training a trojan model is both time- and resource-expensive for defenders. Our run-time 238 backdoor is completely based on released models. Given an inpainting model \mathcal{IM} , it receives an 239 image x and a mask m as input, and returns an image r as the repainted result. This step can be 240 written as $r = \mathcal{IM}(x, m)$. Our objective is to train a set of protective perturbations, denoted as \mathcal{P} . 241 These perturbations are applied to the original image such that $\mathcal{P}(x) = x \otimes \mathcal{P}$. We use perturba-242 tions of the same shape as the image to achieve protection. The resulting perturbed image, $\mathcal{P}(x)$, is 243 designed to resist malicious modification attempts while remaining visually indistinguishable from 244 the original image to the human observer. 245

246 We achieve this by optimizing the two fundamental objectives $\mathcal{L}_{implant}$ and \mathcal{L}_{hide} :

$$\mathcal{L}_{implant} = \mathbb{E}_{\phi, \mathcal{P}(x), m} \left[\left\| \phi_x - \mathcal{I}\mathcal{M}(\mathcal{P}(x), \mathcal{T}(m)) \right\|^2 \right], \tag{4}$$

259

260

247

$$\mathcal{L}_{hide} = \mathbb{E}_{\mathcal{P}(x),m} \left[\left\| \mathcal{I}\mathcal{M}(x,m_0) - \mathcal{I}\mathcal{M}(\mathcal{P}(x),m_0) \right\|^2 \right],$$
(5)

We designate the region of the protected content within the image as the trigger. The backdoor target ϕ_x is generated specifically for each input sample. Detailed implementation of backdoor target can be found in Appendix B. The mask $\mathcal{T}(m)$ here is the entire protected region and we force the model to returns a distorted image similar to target ϕ_x . The implant loss, denoted as $\mathcal{L}_{implant}$, is designed to optimize the perturbation such that $\mathcal{P}(x)$ effectively activates the backdoor in the presence of a trigger m. The hide loss, \mathcal{L}_{hide} , ensures operations outside the trigger region produce results similar to those from the original image x. The mask m_0 used in \mathcal{L}_{hide} is specifically defined to exclude any overlap with the protected region, thereby ensuring benign editing behavior in those regions.

3.3 REGION-AWARE BACKDOOR ACTIVATION MECHANISM

261 The loss optimization procedure is illustrated in Figure 2. Through the optimization of the implant 262 loss $\mathcal{L}_{implant}$, the core backdoor mechanism is implanted at run-time via the introduction of protec-263 tive noise. However, in certain cases, malicious users do not modify the entire sensitive area, and 264 malicious manipulation may cover only a portion of the trigger region. In these instances, the input 265 trigger is incomplete, we still anticipate the successful activation of the backdoor mechanism. Based on this, we expand the mask region used in training perturbations as a potential protection region. 266 We design a mask expansion strategy $\mathcal{E}(\cdot)$, which generates a new mask $\mathcal{E}(m)$ to expand the masked 267 region in the image. We conduct our incomplete activation loss $\mathcal{L}_{incomplete}$ as: 268

$$\mathcal{L}_{incomplete} = \mathbb{E}_{\phi, \mathcal{P}(x), m} \left[\| \phi_x - \mathcal{IM}(\mathcal{P}(x), \mathcal{E}(m)) \|^2 \right], \tag{6}$$

$$\mathcal{E}(m) = conv2d(m, size_{kernel}, size_{padding}),\tag{7}$$

271 given that the mask image is binary, the value in the masked (white) region is 1, while the rest is 0. Thus, the operation $\mathcal{E}(\cdot)$ expands the white region through convolution. By optimizing Equation (6), the backdoor can be activated even when trigger is incomplete. However, a corresponding challenge 274 is emerged: the backdoor may also be activated when editing regions outside the trigger region (e.g., 275 performing benign erasure in the background), as shown in the yellow box in Figure 4. To reduce 276 the impact of editing on non-trigger regions of the image, we propose to jointly optimize \mathcal{L}_{hide} :

$$\mathcal{L}_{hide} = \mathbb{E}_{\mathcal{P}(x),m} \left[\left\| \mathcal{I}\mathcal{M}(x, \mathcal{E}'(m)) - \mathcal{I}\mathcal{M}(\mathcal{P}(x), \mathcal{E}'(m)) \right\|^2 \right],$$
(8)

$$\mathcal{E}'(m) = max(0, min(1, \mathcal{E}(m) - \mathcal{T}(m))), \tag{9}$$

281 By optimizing Equation (8), the trigger region in image can be effectively distinguished by model, preventing the backdoor from being mis-activated when editing does not involve the protected area. 282 We optimizing above objectives in parallel. The total loss function is: 283

$$\mathcal{L}_{total} = \mathcal{L}_{implant} + \mathcal{L}_{incomplete} + \mathcal{L}_{hide}.$$
 (10)

4 **EVALUATION**

In this section, we evaluate our Run-time Backdoor on image inpainting task at first. Then, we present ablation studies of loss functions and perturbation bounds. We discuss transferability of our framework in the end.

292 4.1 EXPERIMENT SETUP 293

294 Model & Scenario Selection. In our evaluation, we utilized two state-of-the-art image inpainting 295 models: LaMa(Suvorov et al. (2022)) and MAT(Li et al. (2022)). These models represent con-296 volutional neural network architecture and transformer architecture, respectively. Additionally, we 297 included Latent Diffusion(Rombach et al. (2022b)) as a comparison for diffusion-based inpaint-298 ing paradigm. To assess the performance of resistance to editing, we employed subsets from four 299 datasets with different scenarios: Places2(Zhou et al. (2017)), CelebA-HQ(Karras (2017)), Watermark, and Car Brands Images. Each subset consisting of 50 randomly selected images. CelebA-HQ 300 is an enhanced version of the CelebA dataset, consisting of high-resolution images of celebrities, 301 we apply the 256×256 image sizes. Places2 is a comprehensive dataset comprising over 400 scene 302 categories, we apply the 512×512 image sizes. Watermark dataset including 30 image-mask pairs. 303 We collected released watermark-mask pairs from Kaggle and supplemented it with another water-304 mark dataset with manually annotated masks. The image sizes include 768×1024 and 768×768 . 305 Car Brands Images consists of a collection of car images with logos. 306

Implant Settings. We define the region in the image corresponding to the mask to serve as trigger. 307 For the trigger region of CelebA-HQ(Karras (2017)), and Places2(Zhou et al. (2017)) datasets, we 308 apply standard 64×64 centering mask. We apply manually annotated mask on Watermark dataset 309 and Car Brands Images, the masks locate at the regions of watermark and car logo respectively. 310 Subsequently, we randomly retain half of the expanded mask region generated by Equation (7) to 311 act as the incomplete trigger. We use the incomplete trigger to subtract the part that intersects with 312 the trigger mask to get a mask without trigger. We empirically set the weight of the \mathcal{L}_{hide} terms to 313 2, and the number of training iterations for the protective noise to 20. We adopt pure color image or 314 inverted color image as the target, the detail will be provided in detail in the Appendix B.

315 Metric Settings. Given the lack of prior research on the run-time backdoor paradigm in image in-316 painting tasks, there is no established baseline for direct comparison. To address this, we employ 317 several metrics to assess the resistance of implanted image to malicious editing. Specifically, we 318 use Contrastive Language-Image Pretraining FID (CLIP-FID)(Radford et al. (2021)) and Learned 319 Perceptual Image Patch Similarity (LPIPS) (Zhang et al. (2018)) to quantify the semantic plausibil-320 ity of the edited image, where higher values indicate greater distortion. The Structural Similarity 321 Index Measure (SSIM)(Wang et al. (2004)) is employed to evaluate the degradation in fidelity of local structure and texture, where lower values signify greater distortion. Additionally, we use Peak 322 Signal-to-Noise Ratio (PSNR) to assess the loss of the refinement and restoration quality in the im-323 age, with lower scores indicating higher distortion. CLIP-FID is calculated on the entire image to

272 273

284 285

287 288

289

290

291



Figure 3: Examples illustrating the qualitative resistance of implanted imge to editing. The redcircled area in the figure highlights the inpainting result.



				LaMa	a		MAT					
Datasets	Editing Region	Input	CLIP-FID↑	LPIPS↑	SSIM↓	PSNR↓	CLIP-FID↑	LPIPS↑	SSIM↓	PSNR↓		
		Imp.	39.91	0.083	0.167	8.98	23.92	0.083	0.294	9.56		
	Trigger	Ben.	12.72	0.017	0.503	18.24	2.12	0.008	0.631	22.57		
		Diff.	+27.18	+0.066	-0.336	-9.26	+21.80	+0.075	-0.337	-13.01		
CalabA		Imp.	23.39	0.047	0.277	13.31	9.46	0.039	0.398	15.05		
LO	Incmp.	Ben.	3.25	0.001	0.556	19.93	1.59	0.006	0.653	22.78		
-nų		Diff.	+20.14	+0.046	-0.279	-6.62	+7.86	+0.033	-0.255	-7.74		
		Imp.	0.845	0.004	0.496	22.55	0.61	0.003	0.585	23.98		
	W/o.	Ben.	0.610	0.003	0.556	24.26	0.47	0.002	0.653	25.89		
		Diff.	+0.235	+0.001	-0.059	-1.71	+0.14	+0.001	-0.132	-1.91		
		Imp.	11.40	0.075	0.099	8.46	29.37	0.073	0.219	10.04		
	Trigger	Ben.	2.88	0.029	0.375	15.81	24.73	0.083	0.238	8.08		
		Diff.	+8.52	+0.046	-0.276	-7.35	4.63	-0.010	-0.019	1.95		
		Imp.	5.06	0.039	0.202	12.52	2.50	0.029	0.342	15.82		
Places2	Incmp.	Ben.	1.94	0.019	0.401	16.55	1.19	0.023	0.449	15.81		
	_	Diff.	+3.12	+0.02	-0.199	-4.03	+1.31	+0.007	-0.107	+0.01		
	W/o.	Imp.	0.63	0.008	0.408	19.74	0.22	0.003	0.612	24.00		
		Ben.	0.49	0.005	0.401	20.98	0.11	0.002	0.449	25.51		
		Diff.	+0.14	+0.003	+0.007	-1.24	+0.11	+0.001	-0.151	-1.51		
		Imp.	22.62	0.141	0.091	6.14	17.93	0.064	0.296	9.81		
	Trigger	Ben.	8.65	0.047	0.362	13.26	14.16	0.063	0.312	8.51		
		Diff.	+13.97	+0.094	-0.271	-7.12	+3.77	+0.001	-0.017	+1.29		
Car		Imp.	9.39	0.069	0.168	11.45	4.34	0.026	0.324	14.27		
Drugala	Incmp.	Ben.	3.34	0.024	0.421	15.58	3.25	0.022	0.434	14.56		
Brands	-	Diff.	+6.05	+0.045	-0.253	-4.13	+1.09	+0.004	-0.110	-0.29		
		Imp.	0.594	0.008	0.340	18.83	0.11	0.003	0.599	23.38		
	W/o.	Ben.	0.259	0.004	0.421	20.74	0.08	0.002	0.434	24.59		
		Diff.	+0.335	+0.004	-0.081	-1.91	+0.03	+0.001	-0.165	-1.21		
	Trigger	Imp.	35.60	0.108	0.173	7.74	31.95	0.069	0.268	11.29		
		Ben.	26.44	0.057	0.459	15.41	27.14	0.061	0.421	16.07		
		Diff.	+9.16	+0.050	-0.286	-7.67	+4.80	+0.008	-0.153	-4.77		
		Imp.	18.15	0.061	0.221	10.02	17.21	0.037	0.309	14.14		
Watermark	Incmp.	Ben.	12.44	0.031	0.523	16.21	13.31	0.033	0.476	16.64		
	_	Diff.	+5.71	+0.030	-0.302	-6.19	+3.89	+0.004	-0.167	-2.49		
		Imp.	0.11	0.004	0.454	26.62	0.16	0.004	0.522	27.45		
	W/o.	Ben.	0.05	0.001	0.523	32.26	0.08	0.002	0.476	32.14		
		Diff.	+0.06	+0.003	-0.069	-5.64	+0.08	+0.002	-0.046	-4.69		

measure the global distortion, and the remaining metrics are calculated on the edited region to reflect the local distortion.

4.2 **Resistance performance to editing**

We conducted a comparative analysis of the editing resistance of run-time backdoor implant method in the LaMa(Suvorov et al. (2022)) and MAT(Li et al. (2022)) models across four scenarios. We report the average distortion levels for each dataset, utilizing four inpainting iterations for each input sample to derive the metrics presented in Table 1. The metrics evaluate the discrepancies in editing results between benign images and those embedded with perturbations across three different modification regions. Table 1 outlines these findings, with the first column indicating the datasets.



Figure 4: Example of qualitative ablation study on loss functions.

	Loss	(LIP-FID↑			LPIPS↑			SSIM↓			PSNR↓	
Dataset	\mathcal{L}_a \mathcal{L}_i \mathcal{L}_h	Trigger	Incmp.	W/o.↓	Trigger	Incmp.	W/o.↓	Trigger	Incmp.	W/o.↑	Trigger	Incmp.	W/o.↑
	Benign	12.72	3.26	0.612	0.017	0.010	0.003	0.503	0.556	0.556	18.24	19.93	24.26
CelebA	· 🗸	36.50	9.071	0.867	0.098	0.021	0.005	0.175	0.368	0.405	8.61	17.57	21.82
-HQ	\checkmark \checkmark	36.01	22.53	4.48	0.098	0.052	0.012	0.158	0.265	0.324	8.69	12.84	19.33
	\checkmark \checkmark \checkmark	39.91	23.39	0.845	0.083	0.047	0.004	0.167	0.277	0.496	8.98	13.31	22.55
	Benign	2.87	1.94	0.043	0.029	0.019	0.005	0.375	0.401	0.401	15.81	16.55	20.98
Diana 2	· √ _	10.28	2.84	0.592	0.083	0.028	0.007	0.072	0.271	0.364	8.09	15.22	19.22
r lace2	\checkmark \checkmark	10.84	5.54	1.59	0.084	0.048	0.013	0.074	0.175	0.273	8.28	11.27	16.63
	\checkmark \checkmark \checkmark	11.40	5.06	0.625	0.075	0.039	0.007	0.099	0.202	0.408	8.46	12.52	19.74

Table 2: Ablation study on loss functions.

The second column delineates the operational regions: "Trigger" refers to the editing of the trigger region, "Incmp." indicates the intersection of the editing area with the trigger, and "W/o." signifies the editing of non-trigger regions. The third column distinguishes between the input images, with "Ben." representing benign images and "Imp." indicating images post-perturbation. The subsequent columns present the resistance performance to editing of the run-time backdoor against various image editing models, evaluated when the perturbation bound is set at 6/255.

As illustrated in Table 1, the run-time implantation demonstrates robust resistance performance to 410 editing across various scenarios when applied to editing operations within protected regions. Run-411 time implantation causes significant damage to the structural similarity metrics and semantic con-412 sistency of image editing results. In comprehensive scenarios with Places2(Zhou et al. (2017)), our 413 approach reduces the structural similarity metrics(SSIM) by 0.336. In the context of facial editing 414 using the CelebA-HQ(Karras (2017)) dataset, global coherence (CLIP-FID) exhibited a minimum 415 performance damage of 21.80, with the local coherence in the editing region (LPIPS) by 0.066. In 416 cases where the editing on the region without trigger (W/o.), the differences between editing results 417 of images with perturbations(Imp.) and those of original images(Ben.) are minimal. Moreover, the metrics in the rows of "W/o." indicating that the perturbation has a negligible compromise the visual 418 quality or realism of the edited images. 419

As shown in Figure 3, we present some qualitative results on Watermark dataset to analyze the resistance performance of our run-time backdoor to malicious editing. The red-circled area in the figure highlights the inpainting result. The original image exhibits minimal resistance to watermark removal, whereas the implanted image demonstrates significantly enhanced protection for the watermark region. More qualitative results can be found in Appendix C. These results highlight the effectiveness of our run-time backdoor framework in maintaining resistance across different editing region, ensuring minimal compromise of the protected regions under subtle perturbation conditions.

427

429

378

379 380

381 382

384

386 387

388

389 390

391 392

393

396 397

399 400

401 402 403

428 4.3 ABLATION STUDIES

430 **Loss function ablation**. Table 2 presents the ablation study on different loss functions, where $\mathcal{L}_a, \mathcal{L}_i$ 431 and \mathcal{L}_h represent implant loss $\mathcal{L}_{implant}$, incomplete activation loss $\mathcal{L}_{incomplete}$ and hide loss \mathcal{L}_{hide} , respectively. As shown in the first column of Table 2, the three loss functions are progressively

458

459

460

461

462



Table 3: Ablation study on perturbation bound.

	Perturbation	0	LIP-FID↑			LPIPS↑			SSIM↓			PSNR↓	
Dataset	Bound	Trigger	Incmp.	W/o.↓	Trigger	Incmp.	W/o.↓	Trigger	Incmp.	W/o.↑	Trigger	Incmp.	W/o.↑
	Benign	12.72	3.26	0.612	0.017	0.010	0.003	0.503	0.556	0.556	18.24	19.93	24.26
CelebA	$\ell_{\infty} = 2/255$	38.08	12.70	0.604	0.048	0.018	0.003	0.234	0.441	0.702	11.12	16.80	24.18
-HQ	$\ell_{\infty} = 3/255$	38.08	18.80	0.642	0.064	0.031	0.003	0.223	0.377	0.660	9.97	14.81	23.70
	$\ell_{\infty} = 6/255$	39.91	23.39	0.845	0.083	0.047	0.004	0.167	0.277	0.496	8.98	13.31	22.55
	$\ell_{\infty} = \frac{13}{255}$	38.81	22.10	0.996	0.098	0.054	0.006	0.107	0.194	0.346	9.15	12.86	20.94
	Benign	2.87	1.94	0.485	0.029	0.019	0.004	0.375	0.401	0.401	15.81	16.55	20.98
Dlage?	$\ell_{\infty} = 2/255$	7.26	2.75	0.594	0.050	0.025	0.004	0.215	0.329	0.547	10.51	14.80	20.69
F lace2	$\ell_{\infty} = 3/255$	9.00	4.74	0.608	0.060	0.032	0.005	0.167	0.278	0.512	9.61	13.70	20.46
	$\ell_{\infty} = 6/255$	11.40	5.06	0.625	0.075	0.039	0.007	0.099	0.202	0.408	8.46	12.52	19.74
	$\ell_{\infty} = \frac{13}{255}$	10.61	4.81	0.792	0.080	0.044	0.008	0.044	0.143	0.321	8.44	12.04	18.63

incorporated into the optimization process. Figure 4 presents the qualitative results of the ablation study on the three loss functions. The first column represents the input masks, with the first one corresponding precisely to the trigger region. The second column shows the inpainting result of the original image, while the subsequent three columns display the inpainting results of images with noise implanted using different loss functions.

The results in Table 2 and Figure 4 indicate that when an incomplete trigger (Incmp.) is introduced, 463 the activation of the backdoor is significantly enhanced after optimizing \mathcal{L}_i . For instance, in the 464 comprehensive scene Places2 Zhou et al. (2017), the semantic rationality distortion index CLIP-FID 465 improves from 2.84 to 5.54 (compared to 1.94 for the original image) after optimizing \mathcal{L}_i , while 466 LPIPS increases from 0.028 to 0.048 (compared to 0.019 for the original image). Additionally, 467 without optimizing hide loss \mathcal{L}_h , the editing results (W/o.) in regions without trigger exhibit signif-468 icant distortion. The yellow box in the Figure 4 illustrates that when optimizing only for implant 469 loss and incomplete loss, editing to non-trigger regions may still mis-activate some backdoor tar-470 gets. Consequently, it is essential to optimize the \mathcal{L}_h , to mitigate false activations of backdoors and minimize interference with benign edits. In terms of structural rationality, the SSIM index is highest 471 when \mathcal{L}_h is optimized in both scenarios. Remarkably, in the Places2 scenario of Table 1, the SSIM 472 of the image's editing result after perturbation implantation even surpasses that of the original image 473 (0.408 vs. 0.401). 474

475 **Perturbation bound ablation**. Table 3 provides a quantitative assessment of the impact of varying 476 perturbation bounds of protective noise on resistance of editing performance. Specifically, perturba-477 tion bounds were set at 1/255, 3/255, and 6/255, respectively. As the perturbation bound increases, the distortion in the resulting image edits becomes more pronounced, with more evident alterations in 478 the visual appearance of the image. When the perturbation bound increases from 6/255 to 13/255, the 479 distortion in the image editing results remains comparable is close, with the latter not consistently 480 showing improvement. This suggests that increasing the perturbation bound does not necessarily 481 lead to better performance and may even plateau or diminish in effectiveness. At a perturbation 482 bound of 13/255, noticeable pixel-level changes begin to emerge. 483

484 Moreover, the "W/o." column of Table 3 indicates reducing the perturbation bound can mitigate 485 the influence of the editing operation in unprotected areas, though this reduction leads to weaker protection in more sensitive regions. For instance, under a perturbation bound of 1/255, the fidelity of

 Benign
 Mask
 Benign Result
 LaMa Result
 MAT Result
 LDM Result
 Target

 Ca pluncidaist cire
 Image: Ca pluncidaist cir

Figure 6: Qualitative example of the resistance performance of editing on different models.

the edited images in the "W/o." condition appears superior; however, qualitative analysis in Figure 5 indicates that this bound renders the backdoor mechanism largely ineffective. At a perturbation bound of 3/255, edits involving regions containing the trigger can successfully activate the backdoor mechanism, although activation becomes unreliable when the trigger is incomplete.

Therefore, it is advisable to increase the noise intensity cautiously, ensuring that it does not introduce visible pixel changes. Empirical observations suggest that a perturbation bound of 3/255 strikes an appropriate balance, offering sufficient protection without perceptible distortion. Detailed selection criteria can be found in Appendix E

507 508

486

497 498

4.4 TRANSFERABILITY

509 We compare the inpainting performance of different models in Figure 6, demonstrating the effec-510 tiveness of the run-time backdoor injection framework across various architectures. The latent dif-511 fusion model (LDM) is tested with input images of size 256×256 due to computational limitations. 512 The analysis shows that LaMa(Suvorov et al. (2022)), a convolutional model, produces outputs 513 closest to the target when using a solid color as the target, followed by MAT(Li et al. (2022)), a 514 transformer-based model, which maintains better global structure due to its attention mechanism. 515 LDM(Rombach et al. (2022b)), due to its denoising process, generates less accurate results, often 516 producing unwanted blue pixels. perturbation bounds for LaMa and MAT were set at 6/255, while 517 LDM used a higher amplitude of 13/255. Furthermore, our implantation framework can keep the 518 backdoor effective even after data augmentation operations such as flipping, cropping, or rescaling. 519 Further details on these parameters are discussed in the Appendix D.

5 CONCLUSION

523 We propose a resource-efficient run-time implant framework designed to expose pre-existing back-524 doors in models while minimizing both time and space consumption. This framework only requires 525 learning protective noise for images during the inference phase and activates the backdoor through a 526 trigger based on the editing region. It has been thoroughly evaluated across various image inpainting 527 models, demonstrating its effectiveness in distorting editing operations when the protected region is involved. The results show that our implanted protective noise can significantly degrade restora-528 tion performance, increasing the CLIP-FID score from 12.72 to 39.90, or reducing the SSIM of the 529 generated content from 0.503 to 0.167 on average. 530

531

520 521

- 532
- 53
- 534
- 535
- 536
- 537
- 538

540 REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 report. *arXiv preprint arXiv:2303.08774*, 2023.
- Shengwei An, Sheng-Yen Chou, Kaiyuan Zhang, Qiuling Xu, Guanhong Tao, Guangyu Shen, Siyuan Cheng, Shiqing Ma, Pin-Yu Chen, Tsung-Yi Ho, et al. How to remove backdoors in diffusion models? In *NeurIPS 2023 Workshop on Backdoors in Deep Learning-The Good, the Bad, and the Ugly*, 2023.
- Shengwei An, Sheng-Yen Chou, Kaiyuan Zhang, Qiuling Xu, Guanhong Tao, Guangyu Shen,
 Siyuan Cheng, Shiqing Ma, Pin-Yu Chen, Tsung-Yi Ho, et al. Elijah: Eliminating backdoors
 injected in diffusion models via distribution shift. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 10847–10855, 2024.
- Weixin Chen, Dawn Song, and Bo Li. Trojdiff: Trojan attacks on diffusion models with diverse tar gets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 pp. 4035–4044, 2023.
- Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. In *arXiv preprint arXiv:1712.05526*, 2017.
- Zhiyuan Cheng, Hongjun Choi, Shiwei Feng, James Chenhao Liang, Guanhong Tao, Dongfang Liu,
 Michael Zuzak, and Xiangyu Zhang. Fusion is not enough: Single modal attacks on fusion models
 for 3d object detection. In *The Twelfth International Conference on Learning Representations*,
 2024a.
- ⁵⁶⁴ Zhiyuan Cheng, Zhaoyi Liu, Tengda Guo, Shiwei Feng, Dongfang Liu, Mingjie Tang, and Xiangyu
 ⁵⁶⁵ Zhang. Badpart: Unified black-box adversarial patch attacks against pixel-wise regression tasks. In *The Forty-first International Conference on Machine Learning*, 2024b.
- Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. How to backdoor diffusion models? In *Proceed-ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4015–4024, 2023.
- Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. Villandiffusion: A unified backdoor attack frame work for diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ciprian Corneanu, Raghudeep Gadde, and Aleix M Martinez. Latentpaint: Image inpainting in latent
 space with diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 4334–4343, 2024.
- 577 Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11326–11336, 2022.
- Jiang Hao, Xiao Jin, Hu Xiaoguang, and Chen Tianyou. Diff-cleanse: Identifying and mitigating
 backdoor attacks in diffusion models. In *arXiv preprint arXiv:2407.21316*, 2024.
- Jonathan Hayase and Sewoong Oh. Few-shot backdoor attacks via neural tangent kernels. In *arXiv preprint arXiv:2210.05929*, 2022.
- Xiaoxiao Hu, Qichao Ying, Zhenxing Qian, Sheng Li, and Xinpeng Zhang. Draw: Defending
 camera-shooted raw against image manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22434–22444, 2023.
- Tero Karras. Progressive growing of gans for improved quality, stability, and variation. *arXiv* preprint arXiv:1710.10196, 2017.
- Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style transfer with a single text condition.
 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 18062–18071, 2022.

- Sen Li, Junchi Ma, and Minhao Cheng. Invisible backdoor attacks on diffusion models. In *arXiv* preprint arXiv:2406.00816, 2024.
- Wenbo Li, Zhe Lin, Kun Zhou, Lu Qi, Yi Wang, and Jiaya Jia. Mat: Mask-aware transformer for
 large hole image inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10758–10768, 2022.
- Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In 25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-221, 2018. The Internet Society, 2018.
- Yunfei Liu, Xingjun Ma, James Bailey, and Feng Lu. Reflection backdoor: A natural backdoor attack on deep neural networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16*, pp. 182–199. Springer, 2020.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool.
 Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11461–11471, 2022.
- Hayk Manukyan, Andranik Sargsyan, Barsegh Atanyan, Zhangyang Wang, Shant Navasardyan, and
 Humphrey Shi. Hd-painter: high-resolution and prompt-faithful text-guided image inpainting
 with diffusion models. *arXiv preprint arXiv:2312.14091*, 2023.
- Minheng Ni, Xiaoming Li, and Wangmeng Zuo. Nuwa-lip: language-guided image inpainting with defect-free vqgan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14183–14192, 2023.
- Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- Kiangyu Qi, Tinghao Xie, Yiming Li, Saeed Mahloujifar, and Prateek Mittal. Revisiting the assumption of latent separability for backdoor defenses. In *The eleventh international conference on learning representations*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,
 and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pp. 8821–8831. Pmlr, 2021.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Highresolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer- ence on computer vision and pattern recognition*, pp. 10684–10695, 2022a.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer- ence on computer vision and pattern recognition*, pp. 10684–10695, 2022b.
- Hadi Salman, Alaa Khaddaj, Guillaume Leclerc, Andrew Ilyas, and Aleksander Madry. Raising the
 cost of malicious ai-powered image editing. *arXiv preprint arXiv:2302.06588*, 2023.
- Wenli Sun, Xinyang Jiang, Dongsheng Li, and Cairong Zhao. Diffphysba: Diffusion-based physical backdoor attack against person re-identification in real-world. *arXiv preprint arXiv:2405.19990*, 2024.
- Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha,
 Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky.
 Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 2149–2159, 2022.

675

682

683

684

696 697

699 700

- Ming Tao, Bing-Kun Bao, Hao Tang, Fei Wu, Longhui Wei, and Qi Tian. De-net: Dynamic text-guided image editing adversarial networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 9971–9979, 2023.
- Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricut, et al. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 18359–18369, 2023.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:
 from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22428–22437, 2023.
- Shiyuan Yang, Xiaodong Chen, and Jing Liao. Uni-paint: A unified framework for multimodal image inpainting with pretrained diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 3190–3199, 2023.
- Tao Yu, Runseng Feng, Ruoyu Feng, Jinming Liu, Xin Jin, Wenjun Zeng, and Zhibo Chen. Inpaint anything: Segment anything meets image inpainting. *arXiv preprint arXiv:2304.06790*, 2023.
- Shengfang Zhai, Yinpeng Dong, Qingni Shen, Shi Pu, Yuejian Fang, and Hang Su. Text-to-image diffusion models can be easily backdoored through multimodal data poisoning. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 1577–1587, 2023.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.
- Xuanyu Zhang, Runyi Li, Jiwen Yu, Youmin Xu, Weiqi Li, and Jian Zhang. Editguard: Versatile image watermarking for tamper localization and copyright protection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11964–11974, 2024.
- Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Chang sheng Xu. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10146–10156, 2023.
 - Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

703 704 705

706

708 709

710 711

712

713

714 715

728 729 730

731 732 733

734

735

736 737 738

739 740

741

742 743 744

748

Appendix

Attack as Defense: Run-time Backdoor Implantation for Image Content Protection

A DETAILED IMPLEMENTATION OF IMAGE INPAINTING

In image inpainting models, an image-mask pair $\langle x, m \rangle$ is typically provided as input. In the mask m, the masked area is represented as white, with a corresponding value of 1, while the non-masked area is represented as black, with a corresponding value of 0.

$$m(i,j) = \begin{cases} 1, & if \quad (i,j) \quad is \quad in \quad the \quad mask \quad region \quad (white) \\ 0, & if \quad (i,j) \quad is \quad outside \quad the \quad mask \quad region \quad (black) \end{cases}$$
(11)

This binary mask is used to delineate the regions for inpainting, where the white areas (value 1) guide the model to predict and fill the missing content, and the black areas (value 0) preserve the original image content. The term r is used to denote the final outcome image, while $\mathcal{G}(x)$ refers to the content predicted by the inpainting model. The relationship between r and $\mathcal{G}(x)$ can be expressed as follows: $r = m \odot \mathcal{C}(r) + (1 - m) \odot r$ (12)

$$r = m \odot \mathcal{G}(x) + (1 - m) \odot x, \tag{12}$$

where the ⊙ denotes the element-wise multiplication. It is evident that in mask-guided editing tasks,
the output result is not directly equivalent to the predicted result from the model. Instead, the final
output is constructed by splicing the model's prediction for the masked regions with the original
unmasked parts of the image. Consequently, even if the model's predicted result closely aligns with
a preset backdoor target, the final inpainting outcome will only reflect the target within the masked
region, while the unmasked region remains unchanged, as the content in the non-mask area is frozen.



Figure 7: The final result of inpainting is the element-by-element multiplication of the model prediction and the frozen part of the image.

B COMPARISON OF DIFFERENT BACKDOOR TARGETS

In this work, we utilize pure color images as backdoor targets. The procedure for generating these targets is outlined as follows. Specifically, we compute the average color of original image x:

$$\mu(x) = \frac{1}{N} \sum_{i=1}^{N} x_i,$$
(13)

where the N is the total number of pixels in image x and x_i represents the RGB value of each pixel. Then we compute the difference between the $\mu(x)$ and the three primary colors—red, green, and blue:

$$\phi_x = \arg \max_{c \in \{c_r, c_g, c_b\}} \|\mu(x) - c\|,$$
(14)

where the c_r, c_g and c_b represent the RGB value of red, green and blue. The color with the largest difference is then selected to generate a pure color image ϕ_x , which serves as the backdoor target. As shown in Figure 8, the red-circled area in the figure highlights the inpainting result. The first row (Benign) presents the editing outcome on the original image, while the second row (Imp.Anti) and third row (Imp.RGB) display the results of editing on perturbed images. In the second row, the backdoor target is the inverted image, whereas in the third row, the target is the pure color image. Backdoor results generated by using pure color images as targets.



Figure 8: Comparison of different backdoor target.

C MORE QUALITATIVE RESULTS OF OUR RUN-TIME BACKDOOR

We present the performance of the run-time backdoor on the CelebA-HQ(Karras (2017)) dataset. The inpainting model is LaMa(Suvorov et al. (2022)). In this experiment, we utilize the default 64×64 centering mask as the designated trigger area. The backdoor target is defined as the inverted color of the image, with the first image serving as an illustrative example.



Figure 9: Qualitative example of the anti-editing performance of the runtime backdoor on the CelebA-HQ dataset.

D QUALITATIVE ANALYSIS OF OUR RUN-TIME BACKDOOR IMPLANTION FRAMEWORK ON DIFFERENT MODELS.

We compare the inpainting performance across different models in Figure 10, highlighting the effectiveness of our run-time backdoor implantation framework in multiple model architectures. Due to computational resource constraints, the input image size for the diffusion model(Rombach et al. (2022b)) is set to 256×256 , and the model employed is latent diffusion (LDM).



Figure 10: Qualitative example of the anti-editing performance on different models.

The qualitative analysis of the three models reveals that, when using a solid color image as the target, the output produced by LaMa(Suvorov et al. (2022)) is closest to the target, followed by MAT(Li et al. (2022)), with LDM yielding the result furthest from the target. This discrepancy arises from the underlying architecture of the models. LaMa, being a convolutional network, is more vulnerable to attacks on local features. In contrast, MAT, as a transformer-based architecture, leverages the attention mechanism, which enables better reconstruction of global structures, leading to more coherent inpainting results. On the other hand, the denoising process in LDM inherently reduces noise, causing some perturbations to fail, and in this case, predominantly generating blue pixels in the inpainting results.

Additionally, the perturbation bound used for LaMa and MAT is set to $^{6}/_{255}$, whereas LDM operates with a perturbation bound of $^{13}/_{255}$.

E QUALITATIVE ANALYSIS ON PROTECTIVE NOISE BOUND.

The analysis presented in Table 3 indicates that increasing the perturbation bound of the protective noise does not invariably enhance the distortion in the editing results. On the contrary, excessively large perturbation bounds introduce perceptible pixel changes in the image. A comparative evaluation between the original image and those with perturbation bounds of 6/255 and 13/255, as shown in the Figure 11, demonstrates that the perturbation at a bound of 6/255 remains within acceptable limits, whereas a bound of 13/255 introduces significant and visually noticeable noise artifacts.

848 849

850

810

826 827

828

829

830

831

832

833

834

835

836

837

838 839 840

841

F DISCUSSION OF COMPUTING RESOURCE REQUIREMENTS

⁸⁵¹ Due to the limitations of computational resources, all our experiments were conducted on a single ⁸⁵² NVIDIA A100-SXM4-40GB GPU. The proposed run-time backdoor method is both efficient and ⁸⁵³ resource-saving. Our framework only introduces perturbations of size $C \times H \times W$ as learnable ⁸⁵⁴ parameters, where C is the number of image channels, $H \times W$ is the image size, and the specific ⁸⁵⁵ memory consumption depends on the inference cost of the target model. Importantly, this approach ⁸⁵⁶ does not require retraining or modification of the model itself, as we simply retain the original ⁸⁵⁷ model's forward pass.

Considering the case of implanting noise into each sample during inference, using the LaMa(Suvorov et al. (2022)) model as an example, the smallest image in our dataset is 256×256 , requiring approximately 3068MB of memory for optimization. For the largest image in the dataset, 768×1024 , the optimization process requires up to 24.96GB of memory. On average, embedding a backdoor into a single image takes approximately 14 seconds. In principle, our method is modelagnostic and can be applied to any model, enabling run-time backdoor implantation with minimal computational overhead.



Figure 11: Qualitative examples of perturbation bounds.

G ROBUSTNESS OF GLOBAL IMAGE TRANSFORMATIONS

We add an experiment to verify the robustness of the proposed method to global image transformations. In our experiments, we evaluated the robustness of the proposed attack under several common image transformations. The results in Figure12 demonstrate that the solution remains effective even after the application of standard image enhancements. Specifically, we tested brightness adjustment (increased by a factor of 1.5), scaling (doubling the image size), blurring (with a blur radius of 3), a 90-degree counterclockwise rotation, and horizontal flipping. In all cases, the attack continued to succeed, indicating that the proposed solution is robust to variations in resolution, rotation, and image reflections. In the case of blurring, our attack method is destructive to the image in some degree, as blur itself degrades the image quality to an unusable level, indicating that such transformations have a significant impact on image usability regardless of whether or not it is attacked. These findings underscore the resilience of the attack to typical image augmentations and suggest that the solution is robust to common preprocessing operations to some extent.

908 909

896 897

898

899

900

901

902

903

904

905

906

907

910 911 912

H PROTECTION PERSISTENCE ACROSS MULTIPLE EDITING CYCLES

We use LaMa model to show the results of four rounds continuous editing, to evaluate the protection persistence. As shown in Figure 13, the image inpainting process freezes the area outside the mask and restores the content of the masked area based on the pixels of these non-masked areas. Therefore, the surrounding disturbances are rarely changed in this process, which can resist multiple rounds of editing.



I MORE EVALUATIONS ON DIFFUSION MODELS

967 968

We found Salman et al. (2023) called image immunization, which is close to our threat model (introduced in Section 2). Their method has two paradigms: encoding attack and diffusion attack.
The encoder attack aims to attack the encoder of the diffusion model, causing the diffusion model to receive inferior image embeddings. The diffusion attack aims to attack the entire diffusion process,



Figure 14: Evaluation of cross-model transferability. "MAT.Imp." and "LaMa.Imp." present the implanted image based on MAT and LaMa, "Infer" presents the editing results (inference).

obtaining better resistance at a greater computational cost. We provide comparison results for both paradigms. Specifically, We use the same model as them, the latent diffusion model Rombach et al. (2022a), and the comparison dataset is CelebA-HQ. As shown in Table 4, our method is still the best compared to the two paradigms with additional baseline.

Table 4: Resistance performance on latent diffusion model with an additional baseline. Whether
 attacking only the encoder of diffusion model or the entire diffusion process, our method can effectively resist image editing and achieve SOTA performance.

	Method	CLIP-FID	LPIPS	SSIM	PSNR
	Benign	18.28	0.047	0.316	13.94
	Encoder Attack	27.85	0.069	0.181	12.20
En	coder Backdoor(ours)	28.94	0.075	0.156	10.98
	Diffusion Attack	35.56	0.096	0.123	7.28
Rur	ntime Backdoor(ours)	40.55	0.101	0.097	7.20

And in Table 5, Our method can significantly affect the latent diffusion model with a perturbation of $\frac{6}{255}$. There are two points worth noting here. First, the diffusion model itself has a denoising function, which will weaken the perturbation to a certain extent. Second, when the perturbation bound increases from $\frac{6}{255}$ to $\frac{13}{255}$, the CLIP-FID in the image editing results remains comparable is close, with the latter not consistently showing improvement (this is consistent with what we described in Section 4.3). To get a better protection, we still recommend appropriately increasing the perturbation bound for the diffusion model.

Perturbation Bound CLIP-FID LPIPS SSIM **PSNR** 13.94 Benign 18.28 0.047 0.316 $\ell_{\infty} = 3/255$ 19.53 0.048 0.262 13.11 $\ell_{\infty} = 6/255$ 40.55 7.20 0.101 0.097 $\ell_{\infty} = \frac{13}{255}$ 36.26 0.128 0.110 5.29

Table 5: Ablation of perturbation bound in latent diffusion model.

1012 1013

1008

1009

1010

1011

982

983 984 985

986

987

988

1014

1015

1017

1016 J TRANSFER ANALYSIS

To assess the cross-model transferability of protective perturbations, we trained images with perturbations based on MAT and then edited them using LaMa. As a result, the protective perturbations still retained a certain resistance effect on LaMa, even though the datasets used for MAT and LaMa training were different. The reverse is not success. As shown in Figure 14. We believe that this is because the vision transformer architecture of MAT has a stronger ability to capture image features than the convolutional network architecture of LaMa, and therefore has a certain degree of transferability. This is also reflected in previous work on similar protective perturbations.

1025 We also discussed two training methods to improve the transferability of different models. One method is parallel backdoor implantation, as shown in Algorithm F. Taking the optimization of two

1026 1027 1028 1029 Algorithm 1 Parallel Backdoor Implantation for Muti Models 1030 1: Input: Image x, Sensitive region m, Inpainting models $\mathcal{IM}_1, \mathcal{IM}_2$, perturbation limit l, num-1031 ber of iterations N1032 2: Output: Optimized perturbation P 3: Initialize perturbation P_0 as a zero matrix of the same size as x 1033 4: Compute target image $\mathcal{T} = f(x, m)$ 1034 5: for iteration t = 1 to N do 1035 Perturb the image: $x' = x + \mathcal{P}_{t-1}$ 6: 1036 7: Obtain inpainting results: 1037 $y_1 = \mathcal{IM}_1(x', m)$ 8: 9: $y_2 = \mathcal{IM}_2(x', m)$ 1039 10: Compute similarity with target: 1040 11: $S_1 = \operatorname{Sim}(y_1, \mathcal{T})$ 1041 12: $S_2 = \operatorname{Sim}(y_2, \mathcal{T})$ 1042 13: Compute total similarity $S = S_1 + S_2$ 1043 14: Compute gradient of the similarity w.r.t. perturbation P: 15: $\nabla \mathcal{P}_t = \nabla (S_1 + S_2)$ 1044 16: Update perturbation: 1045 17: $\mathcal{P}_t = \mathcal{P}_{t-1} - \eta \nabla \mathcal{P}_t$ 1046 18: Project perturbation to stay within the limit: 1047 19: $\mathcal{P}_t = \operatorname{clip}(\mathcal{P}_t, -l, l)$ 1048 20: end for 1049 21: **Return:** Optimized perturbation P_N 1050 1051 1052 1053 1054 1055 1056 Algorithm 2 Sequential Backdoor Implantation for Multi Models 1057 1058 1: Input: Image x, Sensitive region m, Models $\mathcal{IM}_1, \mathcal{IM}_2$, limit l, iterations N_1 , fine-tune N_F 2: **Output:** Optimized perturbation P 3: Initialize perturbation P_0 as zero 4: Compute target $\mathcal{T} = f(x, m)$ 1061 5: Step 1: Train on \mathcal{IM}_1 1062 6: **for** t = 1 to N_1 **do** 1063 $x' = x + \mathcal{P}_{t-1}$ 7: 1064 $y_1 = \mathcal{IM}_1(x', m)$ 8: 1065 9: $S_1 = \operatorname{Sim}(y_1, \mathcal{T})$ $\nabla \mathcal{P}_t = \nabla S_1$ 10: 1067 $\mathcal{P}_t = \operatorname{clip}(\mathcal{P}_{t-1} - \eta \nabla \mathcal{P}_t, -l, l)$ 11: 1068 12: end for 13: Step 2: Fine-tune on \mathcal{IM}_2 1069 1070 14: for t = 1 to N_F do $x' = x + \mathcal{P}_{t-1}$ 1071 15: $y_2 = \mathcal{IM}_2(x', m)$ 16: 1072 17: $S_2 = \operatorname{Sim}(y_2, \mathcal{T})$ 1073 $\nabla \mathcal{P}_t = \nabla S_2$ 18: 1074 19: $\mathcal{P}_t = \operatorname{clip}(\mathcal{P}_{t-1} - \eta \nabla \mathcal{P}_t, -l, l)$ 1075 20: end for 21: Return: Optimized perturbation P 1077 1078 1079



Figure 15: Evaluation on Text-Guided Editing. "Imp.Res" denotes the editing result of implanted image. "Benign" denotes the editing results of original image

models at the same time as an example, we add perturbations to the image and input them to the two models at the same time, and optimize the generated results in parallel towards the Target. This method usually requires large computing resources, so we provide a second method - sequential backdoor implantation, as shown in Algorithm F, which first implants on one model and then fine-tuning the perturbed image on the other model.

1097 1098 1099

1100

1090 1091 1092

1093

1094

1095

1096

1080

K EVALUATION ON TEXT-GUIDED EDITING

1101 In text-guided editing, text prompt will be used to generate content to fill the mask area. Pixels 1102 outside the mask area are only used to smooth the content of text generation. Therefore, when using 1103 this method to resist the method of text guidance, performance is limited. We use Stable Diffusion 1104 for the editor of text guidance to evaluate our methods. The results of the qualitative analysis are 1105 shown in the Figure 15. The prompt we use is "Tail of a dog, high resolution, sitting on a bench in the park." Although it is difficult for us to intercept the production of text content, protective 1106 disturbances can still bring certain distortion to the generated images, such as making pixel colors 1107 close to target. 1108

L DISCUSSION ON LIMITATIONS

1110 1111

Our proposed method introduces the first run-time backdoor implantion framework that eliminates 1112 the need for model retraining. As a result, we prioritize simple and effective backdoor targets over 1113 more complex ones. This choice arises from the observation that the similarity between the gener-1114 ated outputs and the intended target is influenced not only by the introduced perturbations but also 1115 by the inherent generative capacity of the model. In our attack-as-defense scenario, the generative 1116 ability of the model is not manipulated by the backdoor implanter(*defender*). Moreover, our method exhibits good robustness to data augmentation, while the blurring operation may weaken the effect 1117 to a certain extent. This does not impact protection of the image. Please see the Appendix G for 1118 details. 1119

Furthermore, our approach targets open-source white-box models, which implies that the computational resources required depend heavily on the forward computation process of the model, making it difficult to quantify. Nevertheless, the method itself only introduces parameters proportional to the size of the image (we have already analyzed computational resources in Appendix F). Future work could explore strategies such as reinforcement learning or greedy search to implement attackas-defense in black-box models, which would make the computational resource requirements more quantifiable.

Finally, there is a consensus that it is impossible to completely protect against malicious image editing. The method we provide can be used to resist model-driven editing, aiming to raise the threshold
for malicious image manipulation based on AI models, making it difficult for non-professional technicians to save time and effort in infringing image copyrights. However, resisting traditional editing
methods such as Photoshop may require more careful design, which also brings more room for
exploration of image content protection in the future.