



# Edit Away and My Face Will not Stay: Personal Biometric Defense against Malicious Generative Editing

Hanhui Wang<sup>1\*</sup>, Yihua Zhang<sup>2\*</sup>, Ruizheng Bai<sup>3</sup>, Yue Zhao<sup>1</sup>, Sijia Liu<sup>2†</sup>, Zhengzhong Tu<sup>3†</sup>

<sup>1</sup>University of Southern California, <sup>2</sup>Michigan State University, <sup>3</sup>Texas A&M University

\* Equal Contribution, † Corresponding Authors



CVPR Nashville JUNE 11-15, 2025

Paper

Code

## Motivation

*Can we design adversarial perturbations that cause edited images to lose their biometric information, making the edited image biometrically unrecognizable and thereby causing the edit to fail?*

## Contributions

- We present a novel perspective for protecting **personal images** from malicious editing by focusing on **making biometric features unrecognizable** after edits.
- We conduct critical analyses on quantitative evaluation metrics commonly used in image editing, exposing their vulnerabilities and the potential for manipulation to achieve deceptive results.
- We introduce **FaceLock**, which incorporates facial recognition models and feature embedding penalties to effectively protect against diffusion-based image editing.
- Extensive experiments demonstrate that **FaceLock** effectively alters human facial features against various editing prompts, achieving superior defense performance compared to baselines.



Figure 1. Illustration of the two requirements of image editing: prompt fidelity and image integrity.

## Related Works

[1] Salman Hadi et al. "Raising the Cost of Malicious AI-Powered Image Editing."

[2] Ruoxi Chen et al. "EditShield: Protecting Unauthorized Image Editing by Instruction-guided Diffusion Models." ECCV 2024.

## What defines a successful image edit?

- Prompt fidelity:** Edits should accurately reflect the instructions provided in the prompt. (Fig. 1 (b) and (d))
- Image integrity:** Other elements in the image should remain intact after editing. (Fig. 1 (b) and (c))

## FaceLock: Adversarial Biometrics Erasure

- FaceLock: Making edited images biometrically unrecognizable rather than blocking edits outright via perturbation optimization on facial disruption and feature embedding disparity.
- We formulate an optimization objective that jointly enforces biometric disruption and high-level feature disparity, defined as follows:

$$\delta = \arg \max_{\|\delta\|_\infty \leq \epsilon} f_{\text{FR}}(\mathcal{D}(\mathcal{E}(\mathbf{x} + \delta)), \mathbf{x}) + \lambda f_{\text{FE}}(\mathcal{D}(\mathcal{E}(\mathbf{x} + \delta)), \mathbf{x}),$$

FR: the facial recognition loss, FE: feature embedding loss between the input images.  $\mathcal{E}/\mathcal{D}$ : Encoder/Decoder.

## Pitfalls on Existing Evaluation Metrics

- CLIP-based scores overemphasize the presence of elements from the editing instructions, which often leads to prioritizing over-editing. (**Fig. 2**)
- SSIM and PSNR over-rely on differences between the edited image and the undefended source, potentially leading to a false sense of successful defense. (**Fig. 3**)
- We use LPIPS scores as a more robust alternative to SSIM and PSNR for evaluating high-level similarity between edited images.
- We propose to use the facial recognition (FR) similarity score to assess the preservation of biometric identity between the source and edited images.



CLIP-S=N/A CLIP-S=0.091 CLIP-S=0.103 CLIP-S=0.118  
(a) Source Image (b) Edited I (c) Edited II (d) Edited III

Figure 2. CLIP score (CLIP-S) of different editing results. The CLIP score provides a contradictory ranking (III > II > I) compared to the visual quality (I > II > III).



SSIM=N/A PSNR=N/A SSIM=0.869 PSNR=16.44 SSIM=0.746 PSNR=11.60  
(a) Source Image (b) No Defense (c) Defense I (d) Defense II

Figure 3. SSIM and PSNR scores of different defenses. Defense II (d) receives better scores than Defense I (c) due to greater pixel differences from the unprotected edit (b), despite being less effective.

## Experiment Results Highlights

Table 1. Quantitative evaluation on prompt fidelity (CLIP-S, PSNR, SSIM, LPIPS) and image integrity (CLIP-I, FR). Arrows indicate whether a higher or lower value is preferred for a successful defense.

Method	Prompt Fidelity				Image Integrity	
	CLIP-S ↓	PSNR ↓	SSIM ↓	LPIPS ↑	CLIP-I ↓	FR ↓
No Defense	0.118±0.037	-	-	-	0.808±0.074	0.833±0.111
PhotoGuard Encoder attack	<b>0.108</b> ±0.030	<b>15.44</b> ±2.01	0.612±0.056	0.403±0.071	0.670±0.118	0.590±0.264
EditShield	0.110±0.026	17.74±2.20	0.593±0.072	0.382±0.071	0.677±0.096	0.641±0.231
Untargeted Encoder attack	0.116±0.023	16.74±2.27	<b>0.589</b> ±0.084	0.371±0.094	0.653±0.090	0.563±0.236
CW L2 attack	0.115±0.031	19.64±2.46	0.701±0.060	0.247±0.062	0.733±0.089	0.725±0.173
VAE attack	0.114±0.034	19.40±1.70	0.715±0.039	0.251±0.060	0.786±0.061	0.846±0.097
FACELOCK (ours)	0.114±0.024	17.11±2.36	<b>0.589</b> ±0.079	<b>0.436</b> ±0.065	<b>0.648</b> ±0.089	<b>0.315</b> ±0.109

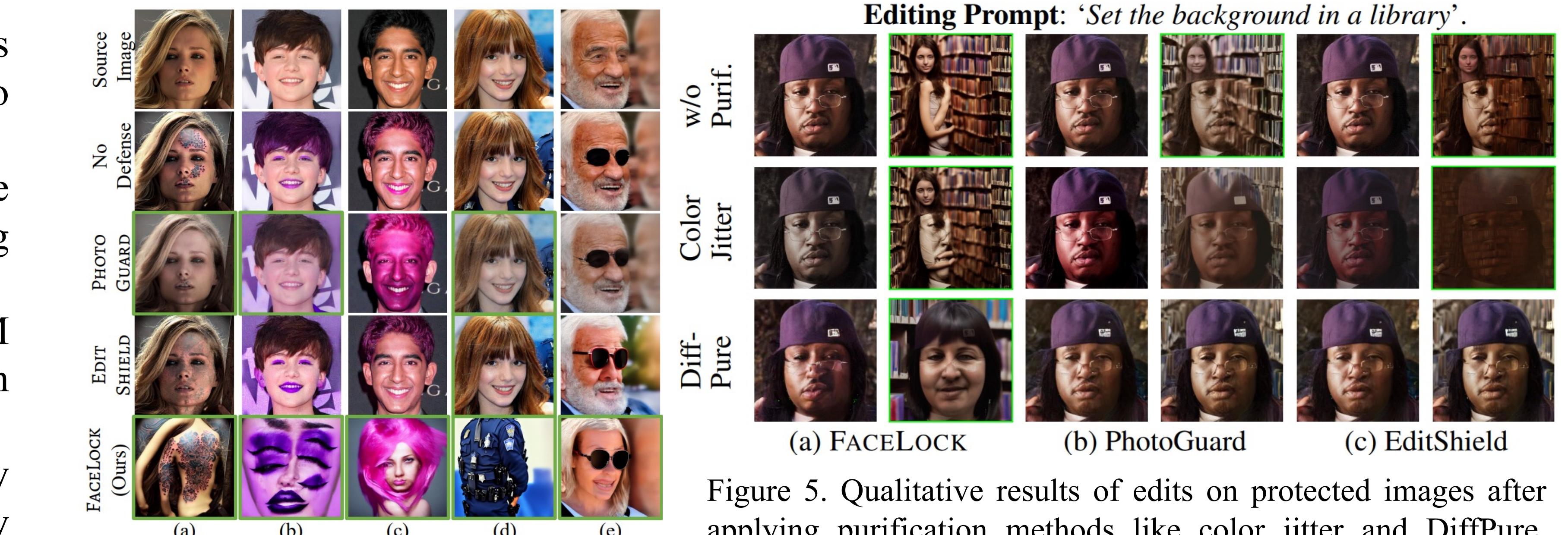


Figure 4. Qualitative results of defensive methods on various editing prompts.

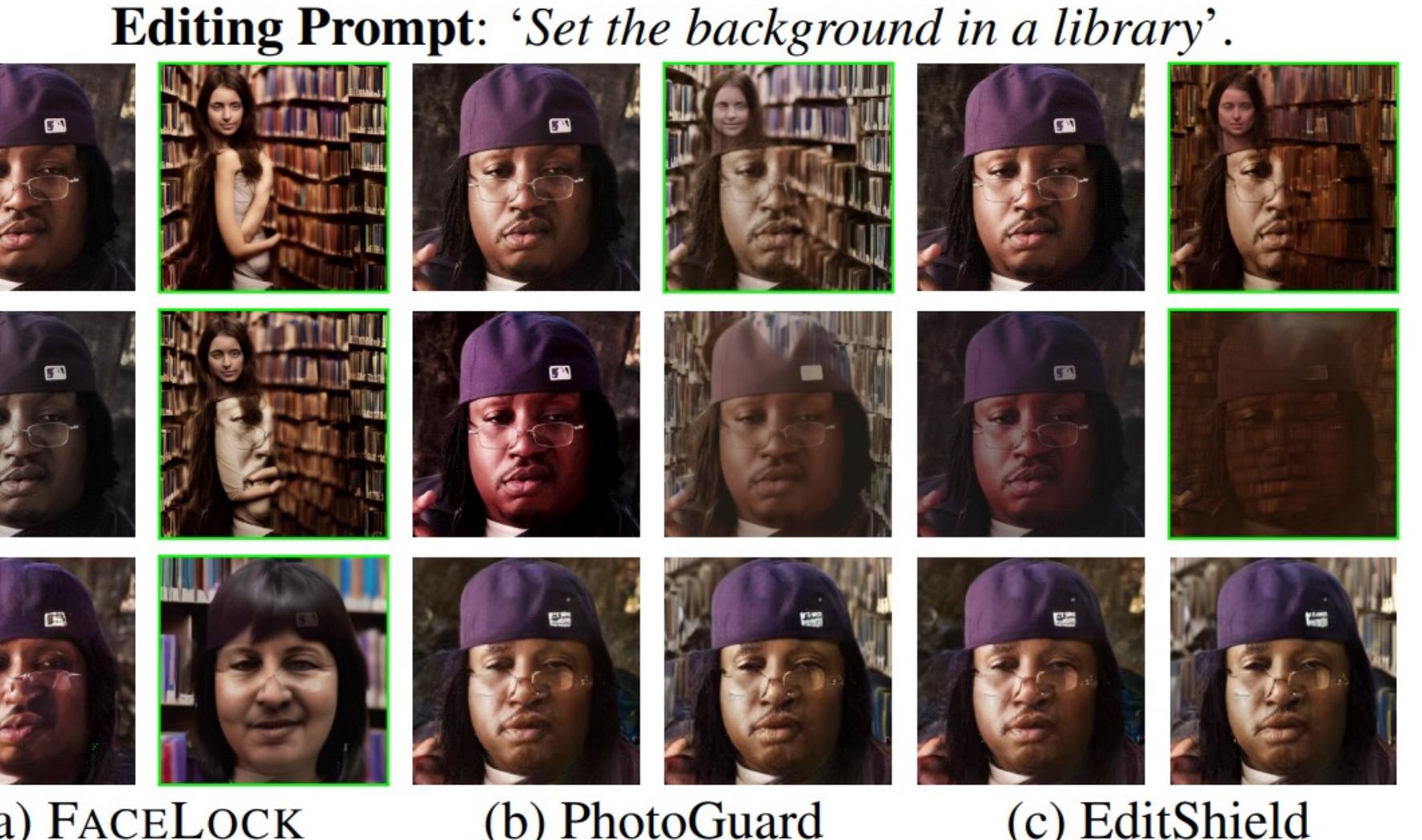


Figure 5. Qualitative results of edits on protected images after applying purification methods like color jitter and DiffPure. Compared to other methods, FaceLock more effectively prevents identity recovery after purification.