

DivTrackee versus DynTracker: Promoting Diversity in Anti-Facial Recognition against Dynamic FR Strategy

Wenshu Fan*

University of Electronic Science and
Technology of China
fws@std.uestc.edu.cn

Minxing Zhang*

CISPA Helmholtz Center for
Information Security
minxing.zhang@cispa.de

Hongwei Li

University of Electronic Science and
Technology of China
hongweili@uestc.edu.cn

Wenbo Jiang†

University of Electronic Science and
Technology of China
wenbo_jiang@uestc.edu.cn

Hanxiao Chen

University of Electronic Science and
Technology of China
chenhanxiao.chx@gmail.com

Xiangyu Yue

The Chinese University of Hong Kong
xyyue@cuhk.edu.hk

Michael Backes

CISPA Helmholtz Center for
Information Security
backes@cispa.de

Xiao Zhang†

CISPA Helmholtz Center for
Information Security
xiao.zhang@cispa.de

Abstract

The widespread adoption of facial recognition (FR) models raises serious concerns about their potential misuse, motivating the development of anti-facial recognition (AFR) to protect user facial privacy. In this paper, we argue that the static FR strategy, predominantly adopted in prior literature for evaluating AFR efficacy, cannot faithfully characterize the actual capabilities of determined trackers who aim to track a specific target identity. In particular, we introduce DynTracker, a dynamic FR strategy where the model’s gallery database is iteratively updated with newly recognized target identity images. Surprisingly, such a simple approach renders all the existing AFR protections ineffective. To mitigate the privacy threats posed by DynTracker, we advocate for explicitly promoting diversity in the AFR-protected images. We hypothesize that the lack of diversity is the primary cause of the failure of existing AFR methods. Specifically, we develop DivTrackee, a novel method for crafting diverse AFR protections that builds upon a text-guided image generation framework and diversity-promoting adversarial losses. Through comprehensive experiments on various image benchmarks and feature extractors, we demonstrate DynTracker’s strength in breaking existing AFR methods and the superiority of DivTrackee in preventing user facial images from being identified by dynamic FR strategies. We believe our work can act as an important initial step towards developing more effective AFR methods for protecting user facial privacy against determined trackers.

1 Introduction

Facial recognition models, where the goal is to recognize the identity of individuals from digital images or videos, have been adopted in various real-world applications such as security checking and attendance management [3, 24, 26, 38, 39]. Although FR systems greatly enhance the convenience of our daily lives, their potential misuse raises serious concerns. In particular, after the release of privacy protection regulations such as GDPR [1] and CCPA [2],

more people have begun to realize that unauthorized FR models can significantly threaten their facial privacy. Therefore, there has been growing attention to developing AFR techniques to protect the facial privacy of users who post their images online [56]. Existing methods can be mainly divided into two categories, gallery-target AFR [9, 10, 21, 48] and query-target AFR [13, 14, 19, 29, 31, 32, 47, 52, 58, 60], depending on whether perturbations are crafted onto gallery or query images to avoid recognition of the target identity.

Witnessing the lack of a well-defined threat model in the existing literature on AFR, we start by formally defining the problem as a two-party security game between Tracker and Trackee (Section 3). We use *Tracker* to denote the malicious user who is determined to track the online images of a specific target identity using FR, while *Trackee* stands for the target identity who deploys AFR protections on their facial images before posting to dodge the tracking. Examining relevant literature on anti-facial recognition, we realize all the existing AFR methods, regardless of gallery-target or query-target, adopt a static FR strategy in their evaluation.

Contribution. We argue that such a static assumption largely restricts the Tracker’s behavior when using an FR model, which is insufficient to characterize the actual capabilities of determined Trackers (Section 4.1). In particular, we propose *DynTracker*, a dynamic FR strategy that iteratively updates the gallery database of Tracker’s FR model with newly identified Trackee images (Section 4.2). DynTracker gradually enriches the gallery database with more clues about Trackee’s facial features and the deployed AFR perturbation scheme, enabling it to uncover more of Trackee’s facial images. Based on a preliminary case study, we demonstrate that DynTracker can indeed achieve significantly higher tracking success rates than the static FR strategy, rendering all the existing AFR protection schemes almost completely ineffective (Section 4.3).

Moreover, we analyze the underlying reasons behind the catastrophic failure of existing AFR methods against our DynTracker, where we hypothesize the main cause is the lack of diversity in the AFR-protected images (Section 5.1). For generative-based AFR, in particular, a fixed auxiliary image of a different identity from

*These authors contributed equally to this work.

†Corresponding authors.

Trackee is typically used to guide the generation process, leading to similar patterns shared among AFR-protected Trackee images. Therefore, once a few protected Trackee images are recognized by the FR model and included in the gallery database, such similarities will significantly decrease the protection success rates of AFR. To address the privacy risks posed by DynTracker, we develop *DivTrackee*, a novel text-guided generative-based AFR method that builds upon diversity-promoting modules and adversarial losses (Section 5.2). In particular, DivTrackee explicitly promotes more diverse generations of AFR perturbations by randomly selecting an auxiliary image from a pool of candidate images and penalizing the similarities of newly generated images to previously produced AFR-protected Trackee images stored in a queue.

We conduct extensive experiments across facial image benchmarks and feature extractors to validate DynTracker’s strength in identifying Trackee images and DivTrackee’s protection efficacy against dynamic FR strategies (Section 6). We also vary the Tracker’s initial knowledge about the Trackee and consider facial verification models to account for possible variations in Tracker’s tracking behaviors. Our results confirm that all the existing AFR methods are highly vulnerable to dynamic FR strategies. For instance, a state-of-the-art generative-based AFR scheme, Clip2Protect [47], achieves only less than 20% protection success rates with respect to Trackee images against our DynTracker. In sharp contrast, our DivTrackee significantly lowers DynTracker’s tracking success rates, consistently achieving more than a 40% increase in protection success while preserving high visual quality of Trackee’s AFR-protected images. Our work highlights the importance of adopting dynamic FR strategies, such as our DynTracker, for more rigorous evaluations of AFR methods and also reveals the effectiveness of promoting more diverse AFR protections as potential countermeasures. We hope our work can serve as an important first step towards developing more reliable AFR protection schemes against determined trackers.

To summarize, our key contributions are as follows:

- By characterizing the limitations of static FR strategies, we propose a dynamic FR strategy *DynTracker*, which is easy to implement and more powerful.
- To adaptively mitigate the privacy breaches induced by DynTracker, we propose *DivTrackee*, which crafts diverse AFR protections using text-guided image generation and diversity-promoting adversarial losses.
- Extensive experiments show that DynTracker effectively breaks existing AFR methods across various facial image benchmarks and feature extractors, and DivTrackee demonstrates superior performance in preventing user facial images from being identified by dynamic FR strategies.

2 Background and Related Work

In this section, we introduce mathematical definitions, necessary preliminaries, and the most relevant literature on facial recognition models and anti-facial recognition technology.

2.1 Facial Recognition

Facial recognition models are designed to identify people by matching their facial characteristics in an unknown image with a database of known faces [30], typically consisting of three major components:

a pre-trained facial feature extractor, a gallery database of known faces, and a query matching step based on some similarity metric.

Feature Extractor. Let $\mathcal{X} \subseteq \mathbb{R}^n$ be the input space of facial images and $\mathcal{Z} \subseteq \mathbb{R}^d$ be some low-dimensional latent feature space. A feature extractor in a modern FR model is usually a deep neural network $f : \mathcal{X} \rightarrow \mathcal{Z}$ trained to extract distinctive facial features of different identities from their digital images. Since human faces often contain large variations between identities and even within a single identity, training a high-quality feature extractor from scratch is expensive in terms of both data collection and computational costs. To encourage collaboration and promote facial recognition applications, various well-established facial feature extractors are publicly available online [8, 12, 33, 46, 59]. Such open-source efforts provide practitioners and researchers with the feasibility and ease of deploying facial recognition models customized to their needs.

Gallery Database. In addition to a reliable feature extractor, an FR model requires a collection of gallery images covering a diverse set of identities, denoted as the gallery database $\mathcal{D}_g = (\mathcal{X}_g, \mathcal{Y}_g)$, where \mathcal{X}_g is a set of gallery images and \mathcal{Y}_g is their corresponding ground-truth identities. Given a target identity that a facial recognition model aims to recognize, an implicit requirement is that the gallery database \mathcal{D}_g contains at least one annotated facial image of the target identity, a prerequisite for successful facial identification. Including multiple facial images for the identities within the gallery database can enhance the robustness of facial recognition, particularly when the identity’s face images vary widely, which has been adopted in state-of-the-art facial recognition systems.

Query Matching. At inference time, a facial recognition model will output an identity prediction for any query image $\mathbf{x}_q \in \mathcal{X}$ based on the feature extractor f , the gallery database \mathcal{D}_g , and a query matching step. To be more specific, let $\mathcal{Z}_g = \{z_g = f(\mathbf{x}_g) : \mathbf{x}_g \in \mathcal{X}_g\}$ be the set of extracted latent features corresponding to the gallery images. Let $z_q = f(\mathbf{x}_q)$ be the extracted latent facial features of \mathbf{x}_q . The query matching step first computes a cosine similarity score between z_q and each of the gallery latent features $z_g \in \mathcal{Z}_g$ and then outputs the identity \hat{y}_g corresponding to the gallery image with the highest similarity score. Let \mathcal{Y} be the output space of identities, then a facial recognition model can be regarded as a function $\text{FR} : \mathcal{X} \rightarrow \mathcal{Y}$ such that: for any query image \mathbf{x}_q ,

$$\text{FR}(\mathbf{x}_q) = \hat{y}_g, \text{ where } (\hat{\mathbf{x}}_g, \hat{y}_g) = \underset{(\mathbf{x}_g, y_g) \in \mathcal{D}_g}{\operatorname{argmax}} \cos(f(\mathbf{x}_q), f(\mathbf{x}_g)). \quad (1)$$

Note that the FR model, based on the cosine distance from the query image, returns the most similar gallery identity, which is adopted without explicit mention in our paper. Using cosine similarity for the query matching step aligns with existing literature [12, 33, 54]. Other alternative prediction rules can also be used, depending on the application scenario. For instance, returning a set of top- k most similar identities is also considered in [47] as the output of the FR model, while Euclidean distance was used to compute the similarity scores in earlier works such as [46]. When the gallery database contains a large number of high-dimensional images, dimension-reduction techniques and tree-based structures are often employed in k -nearest neighbor search to improve the computational efficiency of query matching. In addition, facial verification [7, 42, 51]

is closely related to facial recognition, where similar feature extraction and query matching steps are employed. However, the key difference lies in that a single reference image of the target identity is used to compute the similarity score with a predefined threshold to determine the facial verification result, whereas the gallery database contains multiple images with diverse identities, and no thresholding step is involved in facial recognition.

2.2 Anti-Facial Recognition

Anti-facial recognition is an approach to avoid online tracking of user facial images based on adversarial examples [11, 19, 56]. An AFR method can be understood as a function $\text{AFR} : \mathcal{X} \rightarrow \mathcal{X}$ that maps any input facial image to a perturbed version of it. Existing AFR techniques can be divided into two categories, *gallery-target* and *query-target*, depending on whether AFR injects small adversarial perturbations into the gallery or query images.

Gallery-Target. Gallery-target AFR aims to fool FR models by injecting adversarial perturbations into the gallery images of the target identity. Given a target identity y_{targ} that AFR aims to protect, the collection of gallery images after perturbation $\tilde{\mathcal{X}}_g$ is defined as:

$$\begin{aligned} \tilde{\mathcal{X}}_g &= \mathcal{X}_g^{\text{nt}} \cup \tilde{\mathcal{X}}_g^{\text{t}}, \text{ where } \mathcal{X}_g^{\text{nt}} = \{\mathbf{x}_g \in \mathcal{X}_g \mid y_g \neq y_{\text{targ}}\}, \\ \text{and } \tilde{\mathcal{X}}_g^{\text{t}} &= \{\text{AFR}(\mathbf{x}_g) \mid \mathbf{x}_g \in \mathcal{X}_g, y_g = y_{\text{targ}}\}, \end{aligned} \quad (2)$$

where \mathcal{X}_g denotes the set of clean gallery images before perturbation. Note that only gallery images belonging to y_{targ} are adversarially perturbed by some gallery-target AFR techniques. Due to the injected adversarial features in $\tilde{\mathcal{X}}_g^{\text{t}}$, an FR model is expected to output an incorrect identity when a clean facial image of the target identity is queried.

Query-Target. The other category is known as query-target AFR, which crafts perturbations to the query images of the target. Let y_{targ} be the target and y_q be the ground-truth identity of \mathbf{x}_q , then the collection of (perturbed) query images $\tilde{\mathcal{X}}_q$ is defined as:

$$\begin{aligned} \tilde{\mathcal{X}}_q &= \mathcal{X}_q^{\text{nt}} \cup \tilde{\mathcal{X}}_q^{\text{t}}, \text{ where } \mathcal{X}_q^{\text{nt}} = \{\mathbf{x}_q \in \mathcal{X}_q \mid y_q \neq y_{\text{targ}}\} \\ \text{and } \tilde{\mathcal{X}}_q^{\text{t}} &= \{\text{AFR}(\mathbf{x}_q) \mid \mathbf{x}_q \in \mathcal{X}_q, y_q = y_{\text{targ}}\}, \end{aligned} \quad (3)$$

where $\tilde{\mathcal{X}}_q^{\text{t}}$ is the set of perturbed query images with respect to y_{targ} and some query-target AFR technique, and \mathcal{X}_q denotes the clean counterpart of $\tilde{\mathcal{X}}_q$. Due to the mismatched features between the unperturbed gallery and perturbed query images, the facial privacy of the target identity can be protected.

Related Work. In the prior literature on gallery-target AFR, Shan et al. proposed adding small perturbations that pretend to be non-target identities in the latent facial feature space [48], whereas Cherepanova et al. targeted to mislead the commercial facial recognition APIs by creating perturbations based on an ensemble of FR models [9]. Later, Radiya-Dixit et al. argued that there exists an asymmetry between users who posted online facial images and the trainers of FR models, revealing the limitations of gallery-target AFR techniques against future-released FR models [41]. Generally speaking, gallery-target AFR imposes stronger assumptions on both the construction procedure of the FR model’s gallery database and when the target identity posts their facial images online with or

without AFR protections. Initial attempts of query-target AFR focused on adversarial perturbations bounded in certain distance metrics such as ℓ_p -norm [13, 14, 32, 37, 58, 62], while other works considered patch-based approaches, which typically place an adversarial patch in localized regions [25, 49, 57]. Despite improved privacy protection, these methods suffer undesirable visual artifacts due to the restrictive constraint of the injected perturbations.

Recent works have shifted focus toward generative-based methods to produce more natural adversarial examples. For instance, Hu et al. proposed AMT-GAN that employs generative adversarial networks for makeup transfer [19], Shamshad et al. used text-guided generative models to craft adversarial makeup [47], while Sun et al. developed DiffAM by leveraging the strong generative capability of diffusion models [52]. These methods attain state-of-the-art privacy protection while largely preserving the naturalness of original facial images. Due to their superior performance and milder assumptions required, we primarily focus on generative-based AFR in this work.

3 Security Game Between Tracker and Trackee

Recall that our work aims to study AFR technology for facial privacy protection against unauthorized FR models. However, existing literature on AFR does not provide a unified definition of attackers and defenders due to their distinct research aims. To avoid potential confusion, we formally define the problem of anti-facial recognition as a two-party security game between *Tracker* and *Trackee*, which will be used throughout this paper.

Definition 1. In this security game, Tracker employs an FR model $\text{FR}(\cdot)$ to track the online images Trackee posts. Before posting, Trackee deploys a query-target AFR protection scheme $\text{AFR}(\cdot)$ on their images to dodge the tracking. Let y_{ee} denote the Trackee’s identity and \mathcal{X}_q be a collection of clean query images, then the goal of Tracker can be captured by the following optimization problem:

$$\begin{aligned} \max \quad & \sum_{\mathbf{x}_q \in \mathcal{X}_q} \mathbb{1}\{y_q = y_{\text{ee}}, \text{FR}(\text{AFR}(\mathbf{x}_q)) = y_{\text{ee}}\}, \\ \text{s.t.} \quad & \sum_{\mathbf{x}_q \in \mathcal{X}_q} \mathbb{1}\{y_q \neq y_{\text{ee}}, \text{FR}(\mathbf{x}_q) = y_{\text{ee}}\} \leq \gamma, \end{aligned} \quad (4)$$

where y_q stands for the ground-truth identity of query image \mathbf{x}_q , and $\gamma \in \mathbb{N}_+$ captures the tolerance threshold of the number of false positives. Correspondingly, the goal of Trackee can be defined as:

$$\min \quad \sum_{\mathbf{x}_q \in \mathcal{X}_q} \mathbb{1}\{y_q = y_{\text{ee}}, \text{FR}(\text{AFR}(\mathbf{x}_q)) = y_{\text{ee}}\}. \quad (5)$$

Equations 4 and Equation 5 encode the objectives of Tracker and Trackee, respectively. Tracker aims to accurately identify Trackee’s images (i.e., maximizing true positives) while avoiding incorrect recognitions of other identities’ query images as Trackee (i.e., keeping a low false positive rate). In contrast, Trackee aims to lower the identification rate as much as possible by slightly perturbing the posted images using AFR. Figure 1 illustrates the two-party security game between Tracker and Trackee. Note that we consider the query dataset’s setting to have multiple different images for each identity, which aligns with the characteristics of the dataset gathered by Tracker in practical scenarios. To distinguish from the clean query set \mathcal{X}_q , we use $\tilde{\mathcal{X}}_q$ (similar in Equation 3) to denote the set of actual query images collected by Tracker. In particular, AFR

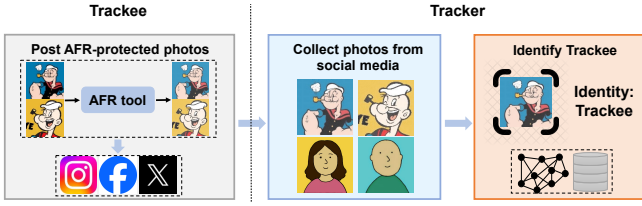


Figure 1: Illustration of the security game between Tracker and Trackee. Trackee uploads AFR-protected images to social media, while Tracker aims to track Trackee using FR models. To avoid potential ethical issues, we use the fictional cartoon character *Popeye* to represent the victim identity of Trackee.

perturbations are assumed to be involved in all the query images of Trackee, while other identities’ face images remain unperturbed since Trackees often only have access to their own images. This assumption largely simplifies our later analyses and is reasonable since, from Trackee’s perspective, perturbing all the images before posting is likely to minimize the risk of tracking by Tracker’s FR models and maximize protection success rates.

In this work, we focus on the dodging attack scenario when evaluating the success of AFR, whereas impersonating attacks are often primarily considered in the existing literature on generative-based AFR [19, 47, 52]. The key difference is that dodging aims for an untargeted attack goal through AFR, whereas impersonation is designed to mimic a specific identity different from Trackee, which is a targeted attack. We believe that dodging is a more suitable objective for protecting Trackee’s facial privacy than impersonating. As noted in [63], a successful impersonating attack does not necessarily imply high success rates in terms of dodging.

Threat Modeling of Tracker. We introduce the threat model of the Tracker in more detail. Recall that Tracker’s objective is to track Trackee’s online posted facial images, where the collection of Trackee’s query images is assumed to be all AFR-protected. To realize the tracking goal, Tracker is assumed to employ an FR model to recognize Trackee’s images from the collected query dataset automatically. We assume that Tracker initially has access to a single facial image of the Trackee, which will be included in the gallery database of the adopted facial recognition model. We believe these assumptions provide a realistic characterization of Tracker’s capabilities and align with the expectation that the constraints imposed on Tracker’s behavior should be as mild as possible.

In addition to the above assumptions, we do not impose any constraints on Tracker, aiming to capture the possible variabilities of FR schemes that Trackers may adopt in practice. In particular, we consider two scenarios for Tracker’s initial knowledge, starting with a clean or an AFR-protected image of Trackee, and evaluate the efficacy of AFR with respect to various publicly available FR models and even a black-box facial verification API. Existing literature on generative-based AFR only considers the initial knowledge of a clean Trackee gallery image for evaluation. However, since Trackee’s online posted AFR-protected image may also be spotted and downloaded by the Tracker, we consider both scenarios of Tracker initially holding a clean or protected Trackee image for comprehensiveness. In addition, we design a new dynamic FR strategy of

Tracker tracking Trackee in Section 4, which is proven to be much stronger than existing static FR schemes and can better capture the actual behavior of determined Trackees and their incentives.

Expectation for AFR from Trackee. Trackee aims to dodge the tracking of unauthorized FR models by crafting AFR perturbations to their images before posting. Generally, two main objectives are expected to be attained by a desirable AFR protection scheme: *protection efficacy* and *visual quality*, which will be explained below.

The primary objective in designing AFR is protection efficacy, which characterizes how successfully AFR can avoid identification by FR models. Note that Trackee does not have the underlying ground-truth knowledge of the actual FR model employed by Tracker. Thus, it is important to consider the protection efficacy of AFR against various FR models and to promote transferability in the AFR design. Besides, AFR needs to preserve the visual quality of Trackee’s facial images to ensure a satisfactory user experience; otherwise, Trackee may not post the AFR-protected images online if they look unnatural or visually dissimilar to the initial images. We primarily focus on query-target AFR schemes that do not significantly change the visual appearance of Trackee’s images, while privacy-preserving tools that remove visually distinctive features of the identity, such as face obfuscation and anonymization [5, 35, 50], are considered out of scope. Since AFR-protected Trackee images are posted before Tracker starts tracking, developing strong AFR schemes is more challenging than crafting adversarial examples to fool ML systems. Trackee has very limited knowledge about the FR model or strategy employed by Tracker. In contrast, traditional adversarial attacks usually assume the victim model can be exploited at least in a black-box manner. As we will illustrate in Section 4, a simple dynamic FR strategy can render almost all of the existing competitive AFR protection schemes ineffective, confirming the challenges we anticipate for the development of AFR.

4 Dynamic Strategy of Facial Recognition

In this section, we first explain why the (static) FR strategies adopted in prior AFR literature are insufficient to characterize the capacity of determined Trackers (Section 4.1). This observation motivates us to develop a dynamic FR strategy, *DynTracker*, which iteratively enriches the gallery database of Tracker’s FR model with newly identified Trackee images (Section 4.2). In addition, we show preliminary results that *DynTracker* can significantly lower the protection efficacy of existing AFR methods (Section 4.3).

4.1 Motivation

As discussed in Section 3, a unique characteristic of the security game between Tracker and Trackee is that AFR protections are deployed on Trackee’s images, which are then posted online before Tracker starts the tracking. That said, the constraints imposed on modeling the Tracker’s behavior should be as less restrictive as possible. Otherwise, it is likely that we will underestimate the capability of Trackers who are determined to track the Trackee. Unfortunately, existing gallery-target and query-target AFR methods implicitly assume in their evaluation that the gallery database of the FR models used by Tracker remains fixed. We term such an evaluation setting as the *static FR strategy*. For example, gallery-target AFR injects small perturbations into gallery images of Trackee identity

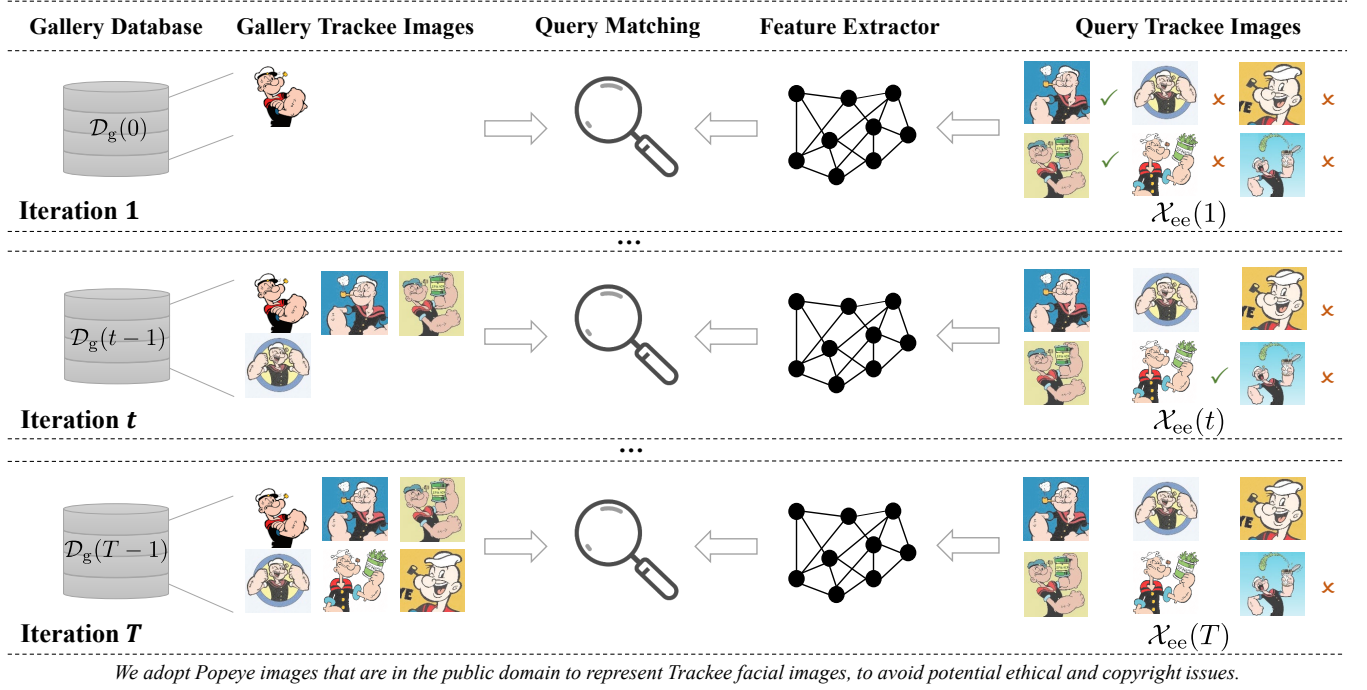


Figure 2: Illustration of the proposed dynamic facial recognition strategy, DynTracker, for Tracker tracking Trackee.

and assumes $\tilde{\mathcal{D}}_g$ is used as the gallery database, while query-target AFR uses a clean facial image dataset with known identities \mathcal{D}_g as the gallery database. Despite different choices, the gallery database will remain unchanged throughout the entire tracking process.

We argue that determined Trackers have both the *capability* and *strong incentives* to update the gallery database dynamically, particularly when the adopted FR model recognizes some of Trackee’s query images. First, given that many well-trained facial recognition models are available online, Trackers can easily download them and make necessary changes to the corresponding gallery database to realize their tracking goal. In fact, Trackee’s initial facial image needs to be inserted into the gallery database to ensure the tracking is specific to the target, which already justifies Tracker’s capability of updating the gallery database. Second, the strong incentives are due to the fact that abundant information about the Trackee might exist in the query images recognized as Trackee. Including these identified Trackee images in the gallery database will likely increase the accuracy of the targeted identity’s recognition, thereby matching the goal of determined Trackers. To support our argument with more concrete evidence, we are going to introduce DynTracker, a dynamic FR strategy featuring multi-round tracking and iterative updates of the gallery database, and provide empirical results showing that DynTracker is much stronger than the static FR strategy for realizing the tracking goal in the following sections.

4.2 DynTracker: Dynamic FR Strategy

So far, we’ve explained the limitations of existing static FR strategies used in AFR literature and motivated the need to consider dynamic FR strategies for better capturing the capabilities and incentives of

determined Trackers. In this section, we introduce the detailed design of our DynTracker, which is illustrated in Figure 2. Algorithm 1 in Appendix A depicts the pseudocode of our DynTracker. In particular, DynTracker consists of the following stages of tracking:

Preparation. Tracker initially holds a single Trackee image, either clean or protected, and includes it in the gallery database of the FR model. Initialize $\mathcal{D}_g(0) = \mathcal{D}_g$ as the gallery database at the start of tracking, and let $\tilde{\mathcal{X}}_q(0) = \tilde{\mathcal{X}}_q$ be the set of query images collected by Tracker, consisting of both AFR-protected Trackee images and clean images from other identities. Here, we use 0 to indicate the 0-th iteration before any facial recognition happens.

Iterative Update. As discussed in Section 4.1, a determined tracker is incentivized to update the gallery database of the FR model with the recognized query images and conduct dynamic facial recognition. Taking the t -th iteration as an example, Tracker first obtains $\mathcal{X}_{ee}(t)$, the set of query Trackee images predicted as y_{ee} by the FR model, then includes $\mathcal{X}_{ee}(t)$ in the gallery database and perform the t -th facial recognition. To be more specific, the update rule with respect to the t -th iteration of DynTracker is defined as: for any $t = 1, 2, 3, \dots$,

$$\begin{aligned} \mathcal{X}_{ee}(t) &= \{\mathbf{x}_q \in \tilde{\mathcal{X}}_q(t-1) \mid \text{FR}(\mathbf{x}_q; \mathcal{D}_g(t-1)) = y_{ee}\}, \\ \mathcal{D}_g(t) &= \mathcal{D}_g(t-1) \cup \{\mathbf{x}, y_{ee} \mid \mathbf{x} \in \mathcal{X}_{ee}(t)\}, \\ \tilde{\mathcal{X}}_q(t) &= \tilde{\mathcal{X}}_q(t-1) \setminus \mathcal{X}_{ee}(t), \end{aligned} \quad (6)$$

where identified images from the previous iteration are dynamically removed from the query dataset to avoid repetition, and we explicitly write out the dependence of $\text{FR}(\cdot)$ on the corresponding gallery database for clarity. Since $\mathcal{D}_g(t)$ is updated dynamically, the set

of Trackee query images $\mathcal{X}_{ee}(t)$ returned by the FR model will be different for different iterations. Technically speaking, there might exist other identities' images incorrectly recognized as Trackee in the returned set $\mathcal{X}_{ee}(t)$. However, the false positive rate turns out to be very low if the FR model is well-trained, which has been verified in our experiments (see Section 7.3 for more discussions on the negligible impact of false positive rates). This iterative update stage gradually inserts more Trackee images into the gallery database to improve the target identity's recognition.

Completion. As long as new query images can be recognized as Trackee's identity y_{ee} , Tracker will continue the iterative update stage, aiming to uncover as many Trackee images as possible. Finally, Tracker will complete the whole tracking process when no more Trackee query images are returned from the previous iteration since no benefits can be gained from an extra FR step. For ease of presentation, we use T to denote the total number of iterations.

Although the gallery dataset is dynamically updated during tracking, the query matching scheme remains the same, which still outputs the identity of the gallery image most similar to the query image. As more Trackee images are recognized and included in the gallery database, the likelihood of the remaining Trackee query images identified by the updated FR model is expected to increase. It is worth noting that the first iteration of the dynamic update phase corresponds to the static FR strategy used in existing query-target AFR methods; thus, our DynTracker is more powerful by design in identifying Trackee's query facial images.

4.3 DynTracker Defeats Existing AFR

Moreover, we conduct a preliminary study to evaluate the efficacy of the proposed dynamic strategy *DynTracker* with comparisons to the static FR strategy. In particular, we randomly select an identity from FaceScrub [36] as the Trackee. Tracker employs ResNet-101 [17] as the feature extractor, which is trained on MS-Celeb-1M [16] via MagFace loss [33], and adopts FaceScrub as the gallery database with a single clean or protected image of Trackee. Table 1 summarizes the tracking success rates against various AFR methods, including gallery-target schemes [9, 48], query-target approaches based on adversarial noises [13, 14, 32, 58] and generative models [19, 47, 60]. We strictly follow the iterative update rule specified in Section 4.2 to carry out DynTracker against query-target AFR schemes. In comparison, we follow a similar dynamic strategy to evaluate gallery-target AFR but iteratively update the protected gallery database $(\tilde{\mathcal{X}}_g, \mathcal{Y}_g)$ with identified clean query images recognized as Trackee. This aligns with the evaluation setup of existing works on gallery-target AFR [9, 48] (see Section 2.2 for details).

Table 1 shows that DynTracker significantly outperforms the static FR strategy in identifying Trackee images from the query dataset, regardless of the initial Tracker's knowledge being a single clean or protected Trackee image. For instance, starting with a single clean image, DynTracker can increase the tracking success rate from 33.83% to 98.04% against Clip2Protect, almost rendering the deployed AFR protection useless. Even against the static FR strategy, if Tracker initially holds a protected Trackee image, state-of-the-art generative-based AFR methods exhibit significantly decreased protection efficacy; not to say against our proposed dynamic FR strategy, all the existing AFR methods appear to be almost useless.

Table 1: Comparisons of tracking success rates (%) between static and dynamic FR strategies against various AFR methods. Here, *Clean* and *Protected* indicate Tracker starting with a single clean and protected Trackee image, respectively.

AFR Method	Clean		Protected	
	Static	Dynamic	Static	Dynamic
No Protection	96.15	99.62	N/A	N/A
Fawkes [48]	N/A	N/A	38.75	98.87
Lowkey [9]	N/A	N/A	9.37	98.69
PGD [32]	28.47	98.76	48.32	97.66
MI-FGSM [13]	31.93	98.94	47.36	97.53
TI-DIM [14]	25.57	97.80	45.87	96.68
TIP-IM [58]	13.82	97.10	40.56	96.54
Adv-Makeup [60]	43.49	98.02	62.34	97.12
AMT-GAN [19]	44.25	99.87	83.50	98.85
Clip2Protect [47]	33.83	98.04	80.39	98.04
DiffAM [52]	26.74	98.10	82.36	98.42

Our preliminary results confirm that the static FR strategy is insufficient to capture the capability of determined Trackers, where the stronger DynTracker should be considered when evaluating the protection effectiveness of AFR methods (see Table 3 and Table 5 for more comprehensive results with similar trends).

5 Promoting Diversity to Combat DynTracker

As demonstrated in Section 4.3, existing AFR methods are reasonably effective against the static FR strategy but fail against our DynTracker. In this section, we hypothesize that a fundamental reason behind such a failure is the lack of diversity in the design of AFR protections (Section 5.1), which further motivates us to develop *DivTrackee*, a novel framework that explicitly promotes diversity in generating AFR protections for Trackee images (Section 5.2). Algorithm 2 in Appendix A presents the pseudocode of DivTrackee.

5.1 Lack of Diversity in Existing AFR

Before developing countermeasures, it is important to pinpoint the cause of failure for existing AFR against DynTracker. As explained in Section 4.2, the main difference between the static FR strategy and DynTracker lies in the set of gallery images corresponding to Trackee's identity. As DynTracker proceeds with more iterative updates, more query images recognized as Trackee will be included in the gallery database. Compared with $\mathcal{D}_g(0)$, the gallery database updated after the t -th iteration of DynTracker $\mathcal{D}_g(t)$ is expected to contain multiple AFR-protected images of Trackee, thus covering more variations of Trackee's facial features and may provide crucial information about the deployed AFR protection scheme, both of which increase the success rates of Tracker tracking Trackee.

We hypothesize that adversarial perturbations produced by existing AFR share a similar pattern across different facial images, causing the failure against dynamic FR strategies. Revisiting Table 1, we can observe that protection success rates across different

query-target AFR methods against the static FR strategy are consistently lower when the Tracker starts with a protected Trackee image, compared with the other clean settings. Such a uniformly decreased performance supports our hypothesis that adversarial perturbations generated by a specific AFR scheme are likely to share similar characteristics, which DynTracker can further exploit. For generative-based AFR [19, 47, 52, 60], an auxiliary facial image that belongs to a different person, denoted as \mathbf{x}_{aux} for simplicity, is usually employed to guide the targeted image generation process. However, the same auxiliary image is shared across different query images when generating AFR perturbations, potentially leading to low protection effectiveness against our DynTracker. For traditional adversarial noise-based AFR [13, 14, 32, 58], although no auxiliary image \mathbf{x}_{aux} is utilized, they rely on untargeted adversarial examples, which may inherently be prone to misclassification into a particular identity class close to Trackee in the latent space. We believe the failure of existing AFR methods in privacy protection against DynTracker can be mainly attributed to the lack of explicit diversity promotion in protected image generations. As shown in Table 1, existing AFR methods fail catastrophically against our DynTracker, which provides strong empirical support to our hypothesis.

5.2 DivTrackee: Diversity-Promoting AFR

Built upon previous analyses, we explore whether promoting diversity in AFR can result in more effective protection. In the following discussions, we will focus on generative-based AFR since it can preserve a better visual quality of facial images than adversarial noise-based AFR. Besides, unlike gallery-target AFR that assumes Trackers will construct their FR model’s gallery database in a certain period of time [41], generative-based AFR does not impose strong assumptions on Tracker’s behavior to be effective.

Initial Attempt. Our first attempt is to promote diversity by assigning different Trackee images with distinct auxiliary images in existing generative-based AFR. Specifically, let \mathcal{X}_{aux} be the auxiliary image dataset with distinct identities. For any clean Trackee image \mathbf{x}_{ee} , one can randomly select an auxiliary image $\mathbf{x}_{\text{aux}} \in \mathcal{X}_{\text{aux}}$, then apply some generative-based AFR method to craft the protection on \mathbf{x}_{ee} . Since randomly selected auxiliary images with different identities are used to guide the image generation process, the diversity of Trackee’s protected images is expected to increase.

Under the same setting as our preliminary experiments in Section 4.3, we evaluate whether diversifying the auxiliary images can help improve AFR performance. We select 2000 images of distinct identities from LFW [20] to construct the auxiliary dataset \mathcal{X}_{aux} . Table 2 shows the evaluation results for a few makeup-based AFR methods with diversified auxiliary images. Compared with fixing the auxiliary image, diversifying auxiliary images increases the protection success of generative-based AFR against DynTracker. In particular, tracking success rates of DynTracker against Clip2Protect are lowered by 15.39% and 12.23%, respectively, when Tracker initially holds a single clean and protected Trackee image. For DiffAM, diversifying auxiliary images improves performance by nearly 20% against DynTracker. Nevertheless, the limited overall improvement in protection effectiveness is unsatisfactory from Trackee’s perspective. Therefore, we take a further step to explore the concept of diversity and propose *DivTrackee*, a novel generative-based AFR

Table 2: Comparisons of DynTracker’s tracking success rates (%) between *fixed* and *diverse* auxiliary images for generative-based AFR. Here, *Clean* and *Protected* mean Tracker starts with a single clean and protected Trackee image, respectively.

AFR Method	Clean		Protected	
	Fixed	Diverse	Fixed	Diverse
Adv-Makeup [60]	98.02	96.23	97.12	95.82
AMT-GAN [19]	99.87	97.69	98.85	92.54
Clip2Protect [47]	98.04	82.65	98.04	85.81
DiffAM [52]	98.10	79.68	98.42	82.90

method with diversity-promoting modules and losses to enhance the protection efficacy of Trackee images against DynTracker.

Text-Guided AFR Generation. We adopt StyleGAN [23], which has been employed for various text-guided image generation and editing tasks [43, 45, 47, 53, 61], as the generative backbone of DivTrackee. Let $G_\theta : \mathcal{W} \rightarrow \mathcal{X}$ be a generative model that maps a latent vector $\mathbf{w} \in \mathcal{W}$ to a generated facial image $G_\theta(\mathbf{w})$. Given a Trackee’s clean image \mathbf{x}_{ee} and a target text prompt p_{targ} (e.g., “natural makeup”), the goal is to search for a latent code $\hat{\mathbf{w}} \in \mathcal{W}$ such that the generated image $G_\theta(\hat{\mathbf{w}})$ follows the specifications of the text prompt p_{targ} and satisfies the expectations for AFR from Trackee discussed in Section 3. Using text-guided generative models and makeup prompts allows ease and flexibility in creating natural AFR perturbations on Trackee’s facial images.

To obtain a desirable $\hat{\mathbf{w}}$, we start with finding a good initialization \mathbf{w}_{init} such that $G_\theta(\mathbf{w}_{\text{init}})$ is close to \mathbf{x}_{ee} to ensure visual similarity. Leveraging the inversion-based technique for StyleGAN image manipulation [53], we use a pre-trained e4e encoder $F_\phi : \mathcal{X} \rightarrow \mathcal{W}$ to infer the inverse of the clean Trackee image as the initialization $\mathbf{w}_{\text{init}} = F_\phi(\mathbf{x}_{\text{ee}})$. Inspired by [44, 47], we fine-tune the generator G_θ to improve the quality of the reconstructed image $G_\theta(F_\phi(\mathbf{x}_{\text{ee}}))$ without sacrificing the editing capabilities of the generative model (see Appendix A.1 for implementation details of such a generator fine-tuning step). After acquiring the fine-tuned generator, denoted as G_{θ^*} , we perform gradient descent on a carefully designed loss function \mathcal{L}_{tot} to optimize \mathbf{w}_{init} . To be more specific, the gradient update rule is defined as: for any $s = 0, 1, 2, \dots, S - 1$,

$$\mathbf{w}_{s+1} = \mathbf{w}_s - \lambda \nabla_{\mathbf{w}} \mathcal{L}_{\text{tot}}(\mathbf{w}_s). \quad (7)$$

Here \mathbf{w}_0 is initialized as \mathbf{w}_{init} , S is the total number of gradient updates, and λ denotes the step size. In particular, \mathcal{L}_{tot} represents the total loss function that encodes all Trackee’s expectations for text-guided AFR generation, which will be detailed next.

Diversity-Promoting Adversarial Loss. Now that we’ve determined the generative framework, the remaining task is to design the total loss \mathcal{L}_{tot} such that the final output from the gradient-based optimization $G_{\theta^*}(\mathbf{w}_S)$ achieves high protection success rates against DynTracker and maintains the visual quality with reference to the Trackee’s original image \mathbf{x}_{ee} and the given text prompt p . To achieve high protection efficacy against DynTracker, we design the following adversarial loss that explicitly promotes diversity:

$$\mathcal{L}_{\text{eff}}(\mathbf{w}) = \mathcal{L}_{\text{adv}}(\mathbf{w}) + \alpha_1 \cdot \mathcal{L}_{\text{guide}}(\mathbf{w}) + \alpha_2 \cdot \mathcal{L}_{\text{div}}(\mathbf{w}), \quad (8)$$

where $\alpha_1 \geq 0$ and $\alpha_2 \geq 0$ are hyperparameters that control the weights of each loss term. In particular, \mathcal{L}_{eff} consists of three terms: an adversarial loss \mathcal{L}_{adv} , a diverse guidance loss $\mathcal{L}_{\text{guide}}$ and a diversity-promoting loss \mathcal{L}_{div} , which are detailed below:

Adversarial Loss. The first loss term \mathcal{L}_{adv} in Equation 8 contributes to the primary goal of dodging: avoiding $G_{\theta^*}(\mathbf{w})$ being recognized as Trackee’s identity. Specifically, let f be a facial feature extractor, then the adversarial loss \mathcal{L}_{adv} is defined as:

$$\mathcal{L}_{\text{adv}}(\mathbf{w}) = -D_f(G_{\theta^*}(\mathbf{w}), \mathbf{x}_{\text{ee}}), \quad (9)$$

where $D_f : \mathcal{X} \times \mathcal{X} \rightarrow [0, 2]$ captures the cosine dissimilarity between the features of two input images extracted by the feature extractor f , i.e., $D_f(\mathbf{x}_1, \mathbf{x}_2) = 1 - \cos(f(\mathbf{x}_1), f(\mathbf{x}_2))$. Note that the feature extractor used in \mathcal{L}_{adv} does not need to be the same as that of DynTracker. As we optimize \mathbf{w} to decrease \mathcal{L}_{adv} through gradient descent, $G_{\theta^*}(\mathbf{w})$ is expected to be less similar to the original Trackee image \mathbf{x}_{ee} in the latent facial feature space.

Diverse Guidance Loss. The second term $\mathcal{L}_{\text{guide}}$ in Equation 8 denotes the diverse guidance loss, which encourages the generated image $G_{\theta^*}(\mathbf{w})$ to get closer to a randomly selected auxiliary image \mathbf{x}_{aux} up to some margin-based threshold $\delta \geq 0$. Specifically, the diverse guidance loss $\mathcal{L}_{\text{guide}}$ is defined as:

$$\mathcal{L}_{\text{guide}}(\mathbf{w}) = \max(0, D_f(G_{\theta^*}(\mathbf{w}), \mathbf{x}_{\text{aux}}) - \delta). \quad (10)$$

Similar to our initial attempt, different Trackee image \mathbf{x}_{ee} will be assigned with random auxiliary image \mathbf{x}_{aux} selected from a dataset \mathcal{X}_{aux} to guide the generation. We additionally employ a margin-based loss with a threshold δ to allow less restrictive directions to search for the desirable latent code.

Diversity-Promoting Loss. The last loss term \mathcal{L}_{div} in Equation 8 explicitly promotes diversity within the generated protected images of Trackee. Specifically, a queue Q with maximum length $m > 0$ is maintained using the *first-in, first-out* (FIFO) principle. For any newly arrived Trackee’s query image \mathbf{x}_{ee} , \mathcal{L}_{div} is defined as the averaged cosine dissimilarity to previously stored images in Q :

$$\mathcal{L}_{\text{div}}(\mathbf{w}) = -\frac{1}{m} \sum_{\mathbf{x} \in Q} D_f(G_{\theta^*}(\mathbf{w}), \mathbf{x}). \quad (11)$$

Note that Q is sequentially expanded as we generate Trackee’s protected images. When the first Trackee image arrives, Q is empty, so the initial $\mathcal{L}_{\text{div}} = 0$. After the first protected image is generated and stored in Q , the generation of the new image needs to account for previous ones. When $|Q|$ reaches the maximum allowable length m , the queue Q will be updated in a FIFO manner. In our experiments, we set $m = 10$ to achieve the diversity-promoting goal without incurring high computational overhead, where further increasing the value of m does not improve the performance much. Minimizing \mathcal{L}_{div} encourages the newly generated image $G_{\theta^*}(\mathbf{w})$ to be dissimilar to the previously-stored protected Trackee images, thus promoting diversity in AFR and enhancing the protection against DynTracker. As will be illustrated in Section 7.2, the proposed diversity-promoting loss terms, $\mathcal{L}_{\text{guide}}$ and \mathcal{L}_{div} , are essential for DivTrackee to achieve high protection success rates.

Loss Design for Visual Quality. In addition to designing diversity-promoting adversarial loss necessary to achieve effective protection against DynTracker, we also need to ensure the visual quality of

generated images. To achieve this goal, we adopt the following loss inspired by prior work [15, 28, 47]:

$$\mathcal{L}_{\text{vis}}(\mathbf{w}) = \alpha_3 \cdot \mathcal{L}_{\text{align}}(\mathbf{w}) + \alpha_4 \cdot \mathcal{L}_{\text{latent}}(\mathbf{w}), \quad (12)$$

where $\alpha_3 \geq 0$ and $\alpha_4 \geq 0$ are hyperparameters that controls the trade-off. In particular, the first term $\mathcal{L}_{\text{align}}$ in Equation 12 captures the alignment of the embeddings regarding the text-image pairs between the initial and generated images based on a pre-trained CLIP model [40]. Let E_I and E_T be the pre-trained CLIP image and text encoders, respectively, then $\mathcal{L}_{\text{align}}$ is defined as:

$$\mathcal{L}_{\text{align}}(\mathbf{w}) = 1 - \cos(\Delta I, \Delta T), \quad (13)$$

where $\Delta I = E_I(G_{\theta^*}(\mathbf{w})) - E_I(G_{\theta^*}(\mathbf{w}_{\text{init}}))$ is the distance between the image embeddings, $\Delta T = E_T(p_{\text{targ}}) - E_T(p_{\text{src}})$ denotes the distance between the text embeddings, and p_{src} is the semantic text prompt of the query image \mathbf{x}_{ee} . Following [47], we simply set p_{src} as “face” since we are generally working with facial images. Such an alignment loss promotes the generated images to follow the specifications of the target text prompt more closely. Finally, the last term $\mathcal{L}_{\text{latent}}$ in Equation 12 penalizes the Euclidean distance between the current latent code \mathbf{w} and the initial latent code \mathbf{w}_{init} :

$$\mathcal{L}_{\text{latent}}(\mathbf{w}) = \|\mathbf{w} - \mathbf{w}_{\text{init}}\|_2. \quad (14)$$

By keeping \mathbf{w} close to \mathbf{w}_{init} , $\mathcal{L}_{\text{latent}}$ ensures the visual similarity between initial and generated images. Putting pieces together, we design the total loss function as $\mathcal{L}_{\text{tot}}(\mathbf{w}) = \mathcal{L}_{\text{eff}}(\mathbf{w}) + \mathcal{L}_{\text{vis}}(\mathbf{w})$ and perform gradient descent on \mathcal{L}_{tot} to obtain the optimized latent code \mathbf{w}_S based on Equation 7. Finally, $G_{\theta^*}(\mathbf{w}_S)$ will be returned as the protected Trackee image corresponding to the original \mathbf{x}_{ee} .

6 Experiments

To systematically study the strength of our DynTracker and validate the superiority of our DivTracker regarding state-of-the-art AFR protection schemes, we conduct comprehensive experiments to compare the protection efficacy and visual aspects of different AFR methods across various image benchmarks and facial recognition settings (Section 6.1). Besides, we also test the generalizability of our methods to facial verification models (Section 6.2). Below, we start with explaining the main setup of our experiments, while more detailed descriptions are provided in Appendix B.

Dataset. We consider 4 widely-used facial image benchmarks, FaceScrub [36], PubFig [27], UMDFaces [4] and CelabA-HQ [22], among which CelabA-HQ is the dataset with high-resolution facial images. All these datasets consist of images with annotated identities, each associated with multiple facial images.

Configuration. For each benchmark dataset, we first randomly select 5 identities as the target identities of Trackees. Then, for each Trackee identity, we generate AFR protections for all the Trackee images, including one image in the gallery database owned by Tracker, while treating the remaining images as the query images of Trackees. We follow Lowkey [9] to set up the query and gallery images for the remaining identities other than Trackee. We compare our DivTrackee with state-of-the-art generative-based methods, such as Adv-Makeup [60], AMT-GAN [19], Clip2Protect [47] and DiffAM [52]. For DivTrackee, we use the first 2000 images from

Table 3: Comparisons of tracking success rates (%) of static and dynamic FR strategies against different generative-based AFR methods across various image benchmarks and feature extractors. For each entry, the first number stands for TSR against the static FR strategy, and the second represents TSR against our DynTracker. Here, *Clean* and *Protected* mean Tracker starting with a single clean and protected Trackee image in the gallery database, respectively. For DynTracker, the lowest tracking success rate is bolded to highlight the best-performing AFR method in terms of Trackee’s facial privacy protection.

Dataset	AFR Method	MobileFace		WebFace		VGGFace		MagFace	
		Clean	Protected	Clean	Protected	Clean	Protected	Clean	Protected
FaceScrub	No Protection	85.89/87.76	N/A	98.17/98.29	N/A	98.26/98.29	N/A	99.29/99.36	N/A
	Adv-Makeup	15.93/85.52	45.93/82.62	31.54/91.02	58.66/88.84	27.96/92.58	59.42/89.16	43.49/98.02	62.34/97.12
	AMT-GAN	13.38/85.42	24.26/82.64	29.87/92.35	56.64/91.17	30.29/92.58	57.18/91.26	44.25/99.87	83.50/98.85
	Clip2Protect	2.46/83.33	31.37/83.33	4.68/91.40	42.64/89.65	4.32/90.87	43.85/88.60	33.83/98.04	80.39/98.04
	DiffAM	14.54/84.96	34.65/85.20	27.64/92.72	52.50/91.46	24.86/91.25	53.90/90.82	39.65/98.26	83.74/98.06
	DivTrackee	2.67/26.58	1.86/22.16	7.82/30.62	2.37/29.49	5.94/31.69	1.98/32.68	6.82/31.94	2.53/28.57
PubFig	No Protection	86.29/88.07	N/A	96.85/97.08	N/A	98.01/98.15	N/A	99.15/99.24	N/A
	Adv-Makeup	24.37/83.38	55.87/79.19	34.46/86.12	58.36/82.08	32.18/88.54	56.73/83.42	35.68/94.57	65.24/92.85
	AMT-GAN	27.64/83.83	56.44/79.67	40.78/87.01	49.64/84.82	34.15/88.53	51.07/85.44	39.34/94.74	60.85/93.16
	Clip2Protect	3.45/81.93	32.37/79.53	4.83/84.69	39.69/82.71	5.35/86.28	40.60/83.46	42.63/93.70	68.18/91.89
	DiffAM	22.10/83.90	52.46/80.40	27.65/85.38	46.09/84.46	28.24/87.32	48.70/83.98	38.10/94.65	63.25/93.18
	DivTrackee	3.23/28.57	2.88/29.58	4.82/32.17	2.63/29.72	5.25/35.76	3.27/32.54	5.88/37.06	4.71/34.53
UMDFaces	No Protection	86.89/87.42	N/A	96.01/96.27	N/A	98.19/98.33	N/A	99.30/99.30	N/A
	Adv-Makeup	18.93/84.28	53.58/80.62	29.23/87.38	57.32/84.21	25.03/88.77	56.79/85.82	33.68/96.08	59.60/95.36
	AMT-GAN	22.37/82.83	47.68/80.60	31.23/88.17	57.72/83.83	30.78/89.48	56.28/84.25	33.48/97.08	61.34/94.88
	Clip2Protect	3.01/81.57	33.01/80.62	4.54/88.61	43.60/87.45	5.38/90.18	48.52/89.87	31.46/95.48	63.56/95.48
	DiffAM	18.62/83.44	45.18/81.36	27.52/87.95	51.46/87.02	23.70/91.04	52.26/91.32	32.19/96.68	60.46/96.30
	DivTrackee	3.69/29.17	3.22/27.54	5.63/31.64	3.89/28.79	4.86/32.61	3.46/29.57	5.18/33.64	4.10/31.15
CelebA-HQ	No Protection	82.30/85.83	N/A	95.62/95.80	N/A	96.90/97.08	N/A	99.29/99.37	N/A
	Adv-Makeup	42.38/90.26	45.54/91.42	46.68/91.39	52.62/91.25	47.70/92.40	52.16/90.94	49.85/90.62	54.62/92.46
	AMT-GAN	36.74/89.42	38.36/90.50	41.69/90.42	46.78/92.40	40.42/91.86	43.60/90.74	42.65/93.42	44.55/91.84
	Clip2Protect	16.34/86.62	31.75/88.36	19.42/92.38	33.68/91.52	21.40/91.98	29.64/90.75	22.56/94.64	36.72/93.18
	DiffAM	28.05/89.60	36.42/88.16	37.48/93.25	48.20/92.66	36.44/92.50	43.84/91.94	43.36/94.90	46.72/93.56
	DivTrackee	18.62/45.82	15.34/42.98	20.68/49.26	16.54/47.38	18.75/51.32	15.68/44.83	21.24/52.37	17.95/50.60

LFW [20] to construct the auxiliary dataset \mathcal{X}_{aux} in the diverse guidance loss and consider the text prompt p_{targ} as “natural makeup”. We choose the hyperparameters involved in DivTrackee as follows: $m = 10$, $\alpha_1 = 0.6$, $\alpha_2 = 1.2$, $\alpha_3 = 0.5$, $\alpha_4 = 0.02$, $\delta = 0.2$, $\lambda = 0.01$, and $S = 60$ (see Section 7.1 for our sensitivity analysis experiments).

Feature Extractor. To set the FR model of Tracker, we consider 4 pre-trained facial feature extractors, including MobileFace [8], WebFace [59], VGGFace [6], and MagFace [33]. To generate AFR protections for Trackee images, we use IR-50 [18], IR-152 [12], and FaceNet [46] as the substitute feature extractors to compute the cosine dissimilarity function D_f for the adversarial losses in Equation 8. Aligned with existing literature [19, 47, 52], we consider the scenarios where feature extractors used for AFR generation differ from those employed by Tracker. This black-box transferability setting simulates the scenario of having imprecise knowledge of Tracker’s FR model when generating AFR protections.

Evaluation Metric. We adopt *tracking success rate* (TSR) to evaluate the strength of a specific Tracker’s FR strategy. Specifically, TSR is defined as the ratio of the number of query images correctly identified as Trackee to the total number of Trackee images in the

query dataset. The higher the TSR is, the stronger the employed FR strategy is. For Trackee, we use *protection success rate* (PSR), defined as $1 - \text{TSR}$, to evaluate the efficacy of an AFR protection scheme. In addition, we visualize the generated AFR-protected images to provide qualitative comparison results, which is considered the primary criterion for evaluating the visual quality of generated images (see Appendix C.3 for additional quantitative comparisons).

6.1 Main Results on Facial Recognition

Table 3 summarizes the tracking success rates of Tracker’s FR strategies against various generative-based AFR methods. We report the average tracking success rate for each AFR scheme over the 5 selected Trackee identities against static FR strategy and our DynTracker, respectively. We consider two settings of Tracker’s initial knowledge of Trackee, either including a clean or an AFR-protected image in the gallery database gathered by Tracker. As a reference, we evaluate the performance of tracking strategies when no AFR protection is employed on the set of Trackee query images, depicted as *No Protection* in Table 3. For completeness, we also evaluate the performance of existing adversarial noise-based AFR protection schemes and provide the results in Table 5 in Appendix C.1.

Strength of DynTracker. We start by looking at the tracking success rates of DynTracker when Trackee’s protected images are crafted using existing AFR techniques. Table 3 shows that our dynamic FR strategy, DynTracker, is very effective in breaking the protections crafted by existing AFR generative-based AFR methods that do not promote diversity explicitly, achieving more than 80% tracking success rates in almost all the tested scenarios. For example, DiffAM can only achieve protection success rates of around 10% on the high-resolution CelebA-HQ dataset against our dynamic FR strategy. Similar trends can be observed in Table 5 for traditional adversarial noise-based AFR. Compared with the scenario where no protections are deployed, existing AFR methods can only achieve a marginal improvement in protection success against DynTracker, confirming their limitations against determined trackers.

In sharp contrast, the tracking success rates of static FR strategy are found to be significantly lower than DynTracker, especially when Tracker holds an initial clean Trackee image. This is the typical evaluation setting adopted in the prior AFR literature. Even if Tracker initially includes a single Trackee’s protected image in the FR model’s gallery database, the protection efficacy of existing AFR methods clearly drops in most cases. For instance, Clip2Protect is very effective on the FaceScrub dataset against the static FR strategy with MobileFace and a clean Trackee image, while TSR increases from 2.46% to 31.37% when the initial Trackee image is changed from clean to protected. These empirical results align with our insights and preliminary results discussed in Section 4, suggesting a crucial overlooked aspect in AFR evaluations. That said, if we use the static FR strategy to evaluate and compare AFR methods, we will obtain a largely overestimated protection efficacy due to the strong assumptions imposed on Tracker’s behavior of not updating the FR model’s gallery database with newly identified Trackee’s protected images. This overestimation can lead to wrong concluding arguments and a false sense of security for AFR protections against real-world determined Trackers. Thus, we strongly recommend the community consider dynamic FR strategies like DynTracker when evaluating the efficacy of AFR protection schemes.

Effectiveness of DivTrackee. Next, we examine the protection efficacy of our DivTrackee. Table 3 shows that DivTrackee consistently improves protection success rates against DynTracker over existing generative-based AFR methods. In particular, DivTrackee maintains tracking success rates of around 30% against DynTracker across FaceScrub, PugFig, and UMDFaces. For the high-resolution CelebA-HQ image dataset, DivTrackee achieves around 50% protection success rates against DynTracker, which significantly improves protection efficacy over existing AFR methods with PSR usually less than 10%. As we have explained in Section 5, the crucial difference lies in the diverse guidance and the diversity-promoting losses involved in DivTrackee, reducing the similarity between AFR-protected Trackee images during generation, which is the reason why DivTrackee can be much more resilient to dynamic FR strategies. In addition, looking at the TSRs of the static FR strategy, we note that the protection efficacy with respect to DivTrackee is even better when the initial Trackee gallery image is protected, which is in sharp contrast to the decreased PSRs for existing AFR methods. Our evaluation results highlight the usefulness of explicit diversity promotion in AFR protections against Trackers who update the



Figure 3: Visualizations of protected Trackee images produced by different generative-based AFR methods. The target text prompt p_{targ} in DivTrackee is set as “natural makeup”.

gallery database. We believe that our exploration of diversifying the AFR-protected images in DivTrackee can serve as an important initial step toward more effective AFR protection schemes against determined Trackers in real-world application scenarios.

Visual Quality. Besides focusing on the primary goal of enhancing the protection efficacy of AFR, we also empirically study whether the protected images produced by DivTrackee can preserve the visual quality of the original facial images, which is a secondary but still important pursuit from Trackee’s perspective. In particular, we visualize the AFR-protected Trackee images generated by different query-target AFR techniques on high-resolution images from the CelebA-HQ dataset.

Figure 3 compares DivTrackee with state-of-the-art generative-based AFR methods, including AMT-GAN, Clip2Protect, and DiffAM, to compare their capabilities of preserving the visual quality of original facial images. We can observe that the protected facial images generated by DivTrackee look natural, without visual distortion of the original identity. In contrast, due to its severe impact at the pixel level, protected images generated by AMT-GAN can be easily distinguished from their clean version. Despite preserving facial characteristics, Clip2Protect tends toward more sallow skin colors. DiffAM has a prominent makeup effect on the eyes and mouth, and the overall skin tone appears brighter, while protected images generated by DivTrackee are more vivid. From the human perceptual view, DivTrackee can largely satisfy the expectations from Trackees for the visual quality pursuit of their images while maintaining strong protection performance.

6.2 Generalization to Facial Verification

Due to the similarity to facial recognition, Trackers have the freedom and capability to use a facial verification model to realize a similar goal of identifying Trackee images. Therefore, we further investigate whether the proposed DynTracker and DivTrackee can generalize to facial verification scenarios. Unlike facial recognition

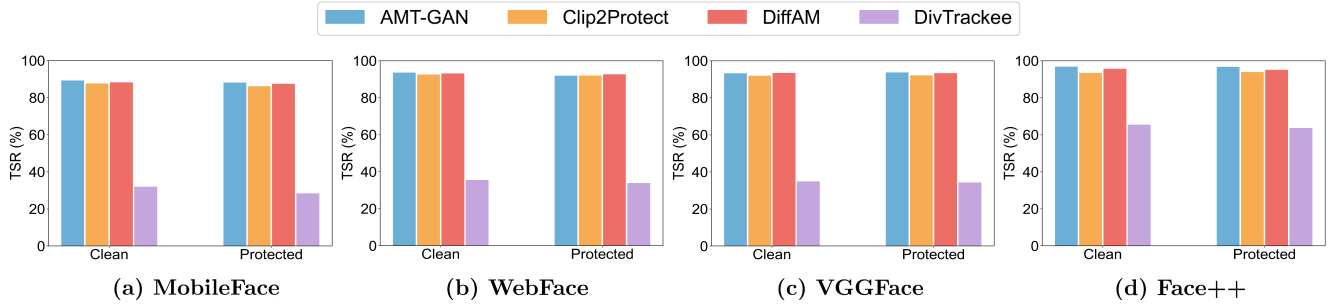


Figure 4: Comparisons of tracking success rates (TSRs) of DynTracker adapted to facial verification settings against various generative-based AFR methods. Here, Figures (a)-(c) consider the scenarios where Tracker adopts some publicly available facial feature extractor, while a facial verification API, Face++, is tested in Figure (d).

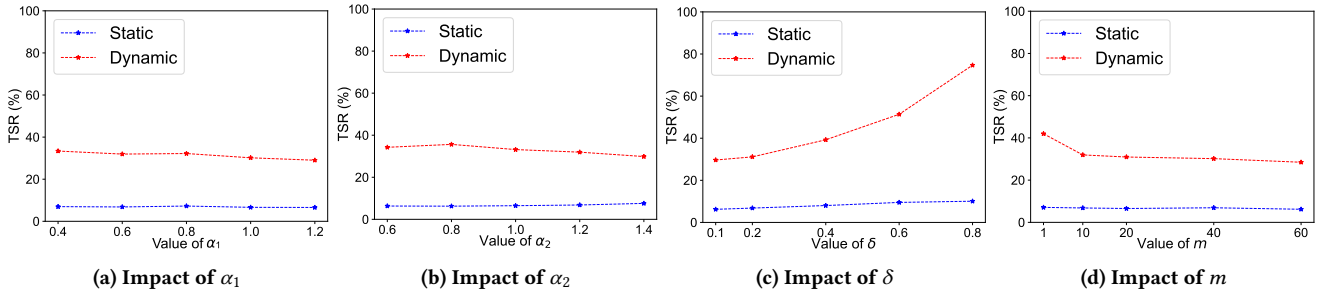


Figure 5: Visualizations of sensitivity analysis results of key hyperparameters regarding the protection efficacy of DivTrackee.

models that require the construction of a gallery database, a reference image of a known identity (e.g., Trackee) and a predefined threshold $\gamma \in [-1, 1]$ are sufficient to verify whether a given query image belongs to the reference identity. Facial images from other individuals are not involved in the facial verification process. Similar to the settings of DynTracker in Section 4, we assume Tracker initially holds a Trackee image, either clean or protected, as the reference image for facial verification. During the dynamic updating stage, Tracker iteratively verifies the facial images in the query dataset based on the similarity scores to the reference image. As long as a successful match exists, the newly identified images will be regarded as reference images in the following iterations. When there are multiple reference images, a new query image is compared with each of them, which will be considered successful if any reference image yields a similarity score exceeding the threshold γ . The tracking process will be terminated when no new query matches can be returned in the current iteration.

Result. Figure 4 summarizes the comparisons of the tracking success rates of DynTracker adapted to facial verification across various generative-based AFR schemes. In this experiment, we consider the FaceScrub image dataset and different facial verification models, including 3 publicly available feature extractors and Face++, a commercial facial verification API.¹ We set the threshold γ to 0.542, which is fixed throughout the tracking process. Similar to our facial recognition experiments, we report the tracking success rate averaged over the same 5 selected identities to compare the efficacy of various AFR methods. Figure 4 shows that DynTracker can be successfully applied to different facial verification frameworks, which

all render the existing generative-based AFR methods ineffective. In sharp contrast, DivTrackee achieves significantly higher protection success rates against the facial verification adapted version of DynTracker. For instance, DivTrackee lowers DynTracker’s tracking success rates by more than 50% compared to existing generative-based AFR methods. These results again validate the strength of our dynamic tracking strategy and suggest that DivTrackee is superior in achieving high protection efficacy, particularly when the Tracker updates the reference Trackee images.

Analysis. Compared with publicly available feature extractors, we observe a performance drop in the protection efficacy achieved by our DivTrackee when the Tracker employs the Face++ API. We suspect that Face++ may have adopted data augmentation or adversarial training techniques to enhance the robustness of the employed facial feature extractor, making it easier to identify the protected Trackee images generated by DivTrackee. This suggests that real-world APIs are powerful tools to implement DynTracker from Tracker’s perspective, where developing better schemes to enhance the protection efficacy of DivTrackee against dynamic tracking strategy using API tools would be an interesting future direction. In addition, we note that the TSRs of DynTracker under facial verification setup are slightly higher than those in Table 3. Such a difference is likely due to the difference in the query matching step: a query match is regarded as successful in our facial verification experiments as long as there exists a reference image with similarity exceeding the threshold γ , whereas only the most similar gallery identity is returned for facial recognition. As will be illustrated in Section 7.3, like what we observe under facial recognition, DynTracker under facial verification settings also induces only

¹Face++ is available online at <https://www.faceplusplus.com/>.

negligible false positive rates. Given the improved identification accuracy of Trackee images and negligible impact on false positive rates, dynamic tracking strategies with facial verification models pose significant threats to user facial privacy, where our diversity-promoting AFR method, DivTrackee, offers a viable solution to mitigate the privacy risks induced by DynTracker.

7 Further Analysis and Discussion

In this section, we study the impact of hyperparameter choices and designed loss terms on the performance of DivTrackee (Sections 7.1 and 7.2), as well as how dynamically updating gallery database affects false positive rates (Section 7.3). All the experiments are conducted on the FaceScrub dataset without loss of generality.

7.1 Hyperparameter Sensitivity

We conduct sensitivity analyses of the key hyperparameters used in DivTrackee. Specifically, we consider the scenario where Tracker initially holds a clean Trackee image and employs MagFace as the facial feature extractor. We record the tracking success rates of DynTracker by varying the values of $\alpha_1, \alpha_2, \delta$, and m , where we observe similar trends for other experimental settings. Figure 5 visualizes our sensitivity analysis results, showing that DivTrackee is highly sensitive to the threshold parameter δ involved in $\mathcal{L}_{\text{guide}}$, i.e., the TSR of DynTracker against DivTrackee sharply increases as we increase the value of δ . In contrast, other hyperparameters have more minor impacts on DivTrackee’s efficacy, and increasing their values decreases DynTracker’s TSR slightly. Also, we note that all the tested hyperparameters have almost no impact on the tracking success rates of the static FR strategy. This phenomenon aligns with our expectations, as these hyperparameters are associated with the diversity-promoting loss terms in \mathcal{L}_{tot} and thus are expected to be more influential to the protection success rates of DivTrackee when DynTracker is employed. We chose the values of these hyperparameters for DivTrackee based on our sensitivity analyses, considering their impact on image generation quality.

7.2 Design Choice of DivTrackee

To validate the key design choices of DivTrackee, we conduct ablation studies to examine the effectiveness of the proposed diversity-promoting adversarial loss terms in Equation 8. We are particularly interested in how the diverse guidance loss $\mathcal{L}_{\text{guide}}$ and the diversity-promoting loss \mathcal{L}_{div} affect the protection success rates of DivTrackee against DynTracker. Figures 6a and 6b illustrate the results under facial recognition and verification, respectively, where Tracker employs MagFace as the feature extractor of the FR model and involves an initial clean or protected Trackee image in the gallery database. To be more concrete, we compare the tracking success rates of DynTracker when $\mathcal{L}_{\text{guide}}$, \mathcal{L}_{div} or both are separately excluded from the total loss function used in DivTrackee for optimizing the latent code in Equation 7. We set the same remaining loss terms and choose the same hyperparameters for DivTrackee.

Result Analysis. The blue bars in Figures 6a and 6b correspond to the performance of our DivTrackee with the total loss employed, where DynTracker exhibits the lowest tracking success rates. When $\mathcal{L}_{\text{guide}}$ is solely excluded from \mathcal{L}_{tot} , the TSR of DynTracker increases by approximately 15%, which suggests the effectiveness of

diversifying the auxiliary image guidance with a properly selected thresholding parameter δ . In addition, when we exclude \mathcal{L}_{div} from the total loss function, the tracking success rate increases more drastically with nearly 30% on average. In comparison, the TSR of DynTracker can be as high as 95% if both $\mathcal{L}_{\text{guide}}$ and \mathcal{L}_{div} loss terms are ignored. Our results highlight the effectiveness of the proposed losses in DivTrackee against determined Trackers, confirming the importance of explicitly promoting diversity in the design of AFR.

7.3 Impact on False Positive Rate

Although we expect well-trained FR models to be capable of accurately recognizing the identity from clean facial images, it is theoretically possible for clean query images of other identities to be misclassified as Trackee. If the employed FR model induces many false positive identifications, the Tracker will need to manually verify the recognized query images, which is particularly undesirable for DynTracker since it involves multiple rounds of facial recognition. Therefore, we conduct experiments to investigate the impact of false positives during the tracking process of DynTracker. Specifically, we first adopt a 7/3 split on FaceScrub to construct training and testing sets and then record the testing identities into which clean images of other individuals are misclassified. Two such identities are selected as Trackee, denoted as Trackee A and Trackee B. Figures 6c and 6d visualize how the number of true and false positives changes during the dynamic update stage of DynTracker, where we consider both facial recognition and verification settings. Similar to our experiments in Section 7.1, Tracker is assumed to adopt MagFace and initially hold a clean Trackee image. We assume Trackee employs our DivTrackee as the AFR protection scheme.

Result Analysis. According to our experimental design, both selected Trackees incur two false positives at the initial iteration of DynTracker. Note that these false positives will be included and paired with the Trackee identity in the FR model’s gallery database after the first iteration of DynTracker. However, as the gallery database is updated with newly identified Trackee query images, surprisingly, 0 (resp. 1) additional false positives are introduced for Trackee A (resp. Trackee B) during the subsequent tracking process. In comparison, the number of true positives for both Trackees steadily increases until it plateaus. These results not only confirm the effectiveness of dynamically updating the gallery database in boosting the tracking success rate but also validate the negligible impact of DynTracker on false positive rates. One possible explanation for the negligible increase in FPR is that since the most similar gallery identity is always returned and the gallery database contains many images, it is unlikely that the query matching step picks the false positive samples included in the gallery database.

7.4 Discussion and Future Work

Computation. For DynTracker, even when the gallery database contains over 40,000 images, the employed FR model can still quickly predict the identity of a query facial image in an average of 28ms. Each update of the nearest neighbor classifier generated from the gallery database takes an average of 8ms. The dynamic update stage of DynTracker typically takes 4-8 iterations to complete. Since our DynTracker is not intended for real-time prediction, we

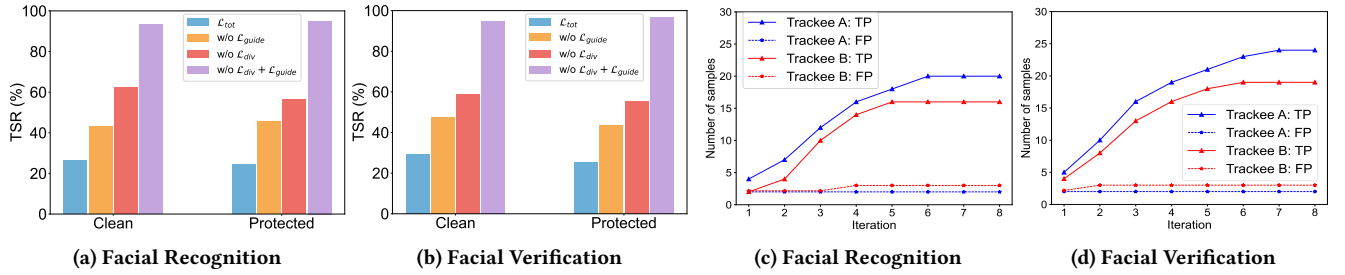


Figure 6: (a)-(b) visualizes the effectiveness of the terms designed in our diversity-promoting adversarial loss, and (c)-(d) plots the curves of the number of true and false positives versus the number of dynamic updates in DynTracker.

Table 4: TSR comparisons of DynTracker with various data preprocessing schemes against DivTrackee on FaceScrub.

Preprocessing Scheme	Clean	Protected
No Preprocessing	31.94%	28.57%
JPEG Compression	32.87%	29.16%
Gaussian Smoothing	34.52%	36.10%

regard the additionally-incurred computational costs as negligible from Tracker’s perspective. In addition, the most time-consuming step for DivTrackee is generator fine-tuning, which is designed to improve the generator’s generalizability to individual Trackee facial images but often takes 1-2 minutes. Exploring how to reduce the computation overhead or even avoid the costly generative fine-tuning step in AFR could be an interesting future direction.

Effectiveness. Despite achieving high PSRs on FaceScrub, PubFig, and UMDFaces, DivTrackee’s performance degrades on the CelebA-HQ dataset against dynamic FR strategies. This may be because high-resolution images contain more facial feature details, and achieving adversarial effects requires adding more perturbations, thereby impacting the visual quality of the generated AFR-protected images. Future work can explore ways to improve DivTrackee’s protection efficacy on high-resolution datasets. In addition, this paper does not consider potential time shifts in Trackee’s facial images, such as the differences in appearance between childhood and adulthood. It remains an open question whether DynTracker can effectively track recently posted Trackee images when the initial gallery image is collected from much earlier periods.

Image Preprocessing. Tracker might adopt state-of-the-art facial recognition models or more robust models in the future, or apply preprocessing techniques (such as image compression or makeup removal tools) to compromise the adversarial perturbations added to the Trackee facial images. Our experiments in Section 6.2 suggest that when DynTracker is implemented based on a more robust FR model (e.g., the Face++ facial verification API), the protection efficacy of DivTrackee can deteriorate. As a preliminary exploration, we conduct experiments to examine the sensitivity of DivTrackee’s performance to two simple data preprocessing schemes, including JPEG compression and Gaussian Smoothing. These preprocessing steps were also studied in existing literature [9]. The results are demonstrated in Table 4, where we consider the FaceScrub dataset

and implement our DynTracker based on the MagFace feature extractor. We can observe that DynTracker’s tracking success rate slightly improves when Trackee images are preprocessed before feeding them into the FR model. Nevertheless, Tracker needs to ensure that the preprocessing methods do not remove excessive facial features of the target identity; otherwise, integrating preprocessing techniques might adversely affect DynTracker’s tracking success.

Adaptive Variations of FR Tracking Strategies. Our proposed DynTracker can be regarded as a specific adaptive FR strategy that is likely to be adopted by determined trackers. However, there are also other variants of tracking strategies that are worth discussing. In this work, Tracker is assumed to include a single Trackee facial image, either clean or protected, in the gallery database to launch DynTracker. In reality, however, Tracker may have already collected multiple clean or protected images by following Trackee’s social media accounts. Including more Trackee images might increase DynTracker’s tracking success rates, aligned with the incentives of determined trackers and real-world application scenarios. In addition, it is worth studying the scenarios in which the FR model’s feature extractor is further fine-tuned using the additional recognized Trackee query images. Note that the feature extractor is assumed to be fixed in DynTracker to simplify the tracking strategy. Similar to the idea of DynTracker versus the static FR strategy in terms of the gallery database, dynamically updating the feature extractor using additional Trackee images may also improve tracking success, compared with fixing the feature extractor. However, the challenge lies in the fact that Tracker does not know precisely about Trackee’s AFR strategy. Instead, Tracker needs to rely on a small amount of AFR-protected images to implement adaptive fine-tuning or un-learning, which is more difficult than building adversarial defenses against specific (known types of) perturbations. Nevertheless, we consider exploring stronger adaptive modifications to DynTracker that may break our protections as interesting future work. On the flip side, developing effective schemes based on our DivTrackee’s insight to be robust to the aforementioned (adaptive) variations of FR tracking strategies would be an important future direction for advancing anti-facial recognition technology.

Game between Tracker and Trackee. Again, we want to emphasize that in the security game between Tracker and Trackee (Definition 1), Trackee is perpetually at a disadvantage since Trackee cannot modify the facial images once they are published. In contrast, Tracker can iteratively refine the tracking strategy, aiming to identify more Trackee posted images [41]. Such asymmetry

does not imply that the Trackee is at a “fundamental” disadvantage but rather represents a challenging task that the community can work on to find better solutions. At a higher level, such asymmetry between Tracker and Trackee is similar to the problem of building robust models against adversarial perturbations, where the defender (like Trackee) is at a disadvantage compared with the attacker (like Tracker). The fundamental difference lies in that we do not have a clear picture of how to model Tracker’s behavior, as opposed to the explicit constraint that we usually impose on adversarial examples. Our work makes an initial step by introducing the precise working pipeline of DynTracker, which is strictly more potent than its static counterpart, relies only on necessary assumptions, and is easy to implement. This enables us to characterize the actual behaviors of determined Trackers better and establish an essential foundation for developing more reliable AFR technologies.

8 Conclusion

We designed a simple but highly effective dynamic FR strategy, DynTracker, where the model’s gallery database is progressively updated with newly recognized target identity images. Our work reveals the limitations of existing AFR, particularly highlighting the importance of using dynamic FR strategies for more rigorous AFR evaluations. Besides, we develop DivTrackee, an innovative text-guided generative AFR method that leverages diversity-promoting modules and adversarial losses. It explicitly encourages more diverse generations of AFR-protected images as a promising solution to defend against DynTracker. We hope that our comprehensive investigations of dynamic FR tracking strategies and their countermeasures can inspire more reliable designs of anti-facial recognition techniques that can potentially be deployed for practical use.

Availability

Our DynTracker and DivTrackee implementations, as well as all our experiments, are available as open-source code at [this url](#).

Acknowledgments

This work is supported by the National Key R&D Program of China under Grant 2022YFB3103500, the National Natural Science Foundation of China under Grant 62402087 and 62020106013, the Chengdu Science and Technology Program under Grant 2023-XT00-00002-GX, the Sichuan Science and Technology Program under Grant 2024ZHCG0188 and 2025ZNSFSC1490, the Fundamental Research Funds for Chinese Central Universities under Grant ZYGX2024J019, the China Postdoctoral Science Foundation under Grant BX20230060, BX20240053 and 2024M760356.

References

- [1] [n. d.]. <https://gdpr-info.eu/>.
- [2] [n. d.]. <https://oag.ca.gov/privacy/ccpa>.
- [3] Insaf Adjabi, Abdeldjalil Ouahabi, Amir Benzaoui, and Abdelmalik Taleb-Ahmed. 2020. Past, present, and future of face recognition: A review. *Electronics* 9, 8 (2020), 1188.
- [4] Ankan Bansal, Anirudh Nanduri, Carlos D Castillo, Rajeev Ranjan, and Rama Chellappa. 2017. Umdfaces: An annotated face dataset for training deep networks. In *2017 IEEE international joint conference on biometrics (IJCB)*. IEEE, 464–473.
- [5] Jingyi Cao, Bo Liu, Yunqian Wen, Rong Xie, and Li Song. 2021. Personalized and invertible face de-identification by disentangled identity information manipulation. In *Proceedings of the IEEE/CVF international conference on computer vision*. 3334–3342.
- [6] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. 2018. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*. IEEE, 67–74.
- [7] Xudong Cao, David Wipf, Fang Wen, Genquan Duan, and Jian Sun. 2013. A practical transfer learning algorithm for face verification. In *Proceedings of the IEEE international conference on computer vision*. 3208–3215.
- [8] Sheng Chen, Yang Liu, Xiang Gao, and Zhen Han. 2018. Mobilefacenet: Efficient cnns for accurate real-time face verification on mobile devices. In *Biometric Recognition: 13th Chinese Conference, CCBR 2018, Urumqi, China, August 11-12, 2018, Proceedings 13*. Springer, 428–438.
- [9] Valeriia Cherepanova, Micah Goldblum, Harrison Foley, Shiyuan Duan, John P Dickerson, Gavin Taylor, and Tom Goldstein. 2021. LowKey: Leveraging Adversarial Attacks to Protect Social Media Users from Facial Recognition. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=hJmtwocEqzc>
- [10] Ka-Ho Chow, Sihao Hu, Tiansheng Huang, Fatih Ilhan, Wenqi Wei, and Ling Liu. 2024. Diversity-driven privacy protection masks against unauthorized face recognition. *Proceedings on Privacy Enhancing Technologies* (2024).
- [11] Ka-Ho Chow, Sihao Hu, Tiansheng Huang, and Ling Liu. 2025. Personalized Privacy Protection Mask Against Unauthorized Facial Recognition. In *European Conference on Computer Vision*. Springer, 434–450.
- [12] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. 2019. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4690–4699.
- [13] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. 2018. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 9185–9193.
- [14] Yinpeng Dong, Tianyu Pang, Hang Su, and Jun Zhu. 2019. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4312–4321.
- [15] Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–13.
- [16] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. 2016. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*. Springer, 87–102.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [18] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7132–7141.
- [19] Shengshan Hu, Xiaogeng Liu, Yechao Zhang, Minghui Li, Leo Yu Zhang, Hai Jin, and Libing Wu. 2022. Protecting facial privacy: Generating adversarial identity masks via style-robust makeup transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15014–15023.
- [20] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. 2008. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Workshop on faces in Real-Life Images: detection, alignment, and recognition*.
- [21] Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, and Yisen Wang. 2021. Unlearnable examples: Making personal data unexploitable. *arXiv preprint arXiv:2101.04898* (2021).
- [22] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2017. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196* (2017).
- [23] Tero Karras, Samuli Laine, and Timo Aila. 2019. A Style-Based Generator Architecture for Generative Adversarial Networks. In *Conference on Computer Vision and Pattern Recognition CVPR*. IEEE, 4401–4410.
- [24] Paramjit Kaur, Kewal Krishan, Suresh K Sharma, and Tanuj Kanchan. 2020. Facial-recognition algorithms: A literature review. *Medicine, Science and the Law* 60, 2 (2020), 131–139.
- [25] Stepan Komkov and Aleksandr Petiushko. 2021. Advhat: Real-world adversarial attack on arcface face id system. In *2020 25th international conference on pattern recognition (ICPR)*. IEEE, 819–826.
- [26] Yassin Kortli, Maher Jridi, Ayman Al Falou, and Mohamed Atri. 2020. Face recognition systems: A survey. *Sensors* 20, 2 (2020), 342.
- [27] Neeraj Kumar, Alexander C Berg, Peter N Belhumeur, and Shree K Nayar. 2009. Attribute and simile classifiers for face verification. In *2009 IEEE 12th international conference on computer vision*. IEEE, 365–372.
- [28] Gihyun Kwon and Jong Chul Ye. 2022. Clipstyler: Image style transfer with a single text condition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18062–18071.
- [29] Minh-Ha Le and Niklas Carlsson. 2024. StyleAdv: A Usable Privacy Framework Against Facial Recognition with Adversarial Image Editing. *Proceedings on Privacy Enhancing Technologies* (2024).

- [30] Lixiang Li, Xiaohui Mu, Siying Li, and Haipeng Peng. 2020. A review of face recognition technology. *IEEE access* 8 (2020), 139110–139120.
- [31] Minghui Li, Jiangxiong Wang, Hao Zhang, Ziqi Zhou, Shengshan Hu, and Xiaobing Pei. 2024. Transferable adversarial facial images for privacy protection. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 10649–10658.
- [32] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*.
- [33] Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. 2021. Magface: A universal representation for face recognition and quality assessment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 14225–14234.
- [34] Dongbin Na, Sangwoo Ji, and Jong Kim. 2022. Unrestricted Black-Box Adversarial Attack Using GAN with Limited Queries. In *European Conference on Computer Vision*. Springer, 467–482.
- [35] Elaine M Newton, Latanya Sweeney, and Bradley Malin. 2005. Preserving privacy by de-identifying face images. *IEEE transactions on Knowledge and Data Engineering* 17, 2 (2005), 232–243.
- [36] Hong-Wei Ng and Stefan Winkler. 2014. A data-driven approach to cleaning large face datasets. In *2014 IEEE international conference on image processing (ICIP)*. IEEE, 343–347.
- [37] Seong Joon Oh, Mario Fritz, and Bernt Schiele. 2017. Adversarial Image Perturbation for Privacy Protection – A Game Theory Perspective. In *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [38] Omkar Parkhi, Andrea Vedaldi, and Andrew Zisserman. 2015. Deep face recognition. In *BMVC 2015-Proceedings of the British Machine Vision Conference 2015*. British Machine Vision Association.
- [39] Divyavarsinh N Parmar and Brijesh B Mehta. 2014. Face recognition methods & applications. *arXiv preprint arXiv:1403.0485* (2014).
- [40] 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*. PMLR, 8748–8763.
- [41] Evani Radiya-Dixit, Sanghyun Hong, Nicholas Carlini, and Florian Tramer. 2022. Data Poisoning Won't Save You From Facial Recognition. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=B5XahNLmna>
- [42] Narayanan Ramanathan and Rama Chellappa. 2006. Face verification across age progression. *IEEE transactions on image processing* 15, 11 (2006), 3349–3361.
- [43] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. 2021. Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2287–2296.
- [44] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. 2022. Pivotal tuning for latent-based editing of real images. *ACM Transactions on graphics (TOG)* 42, 1 (2022), 1–13.
- [45] Axel Sauer, Katja Schwarz, and Andreas Geiger. 2022. Stylegan-xl: Scaling stylegan to large diverse datasets. In *ACM SIGGRAPH 2022 conference proceedings*. 1–10.
- [46] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 815–823.
- [47] Fahad Shamshad, Muzammal Naseer, and Karthik Nandakumar. 2023. CLIP2Protect: Protecting Facial Privacy using Text-Guided Makeup via Adversarial Latent Search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20595–20605.
- [48] Shawn Shan, Emily Wenger, Jiayun Zhang, Huiying Li, Haitao Zheng, and Ben Y Zhao. 2020. Fawkes: Protecting privacy against unauthorized deep learning models. In *29th USENIX security symposium (USENIX Security 20)*. 1589–1604.
- [49] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K. Reiter. 2019. A general framework for adversarial examples with objectives. *ACM Transactions on Privacy and Security* 22, 3 (2019).
- [50] Qianru Sun, Liqian Ma, Seong Joon Oh, Luc Van Gool, Bernt Schiele, and Mario Fritz. 2018. Natural and effective obfuscation by head inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5050–5059.
- [51] Yi Sun, Xiaogang Wang, and Xiaoou Tang. 2013. Hybrid deep learning for face verification. In *Proceedings of the IEEE international conference on computer vision*. 1489–1496.
- [52] Yuhao Sun, Lingyun Yu, Hongtao Xie, Jiaming Li, and Yongdong Zhang. 2024. DiffAM: Diffusion-based Adversarial Makeup Transfer for Facial Privacy Protection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 24584–24594.
- [53] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. 2021. Designing an encoder for stylegan image manipulation. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–14.
- [54] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. 2018. Cosface: Large margin cosine loss for deep face recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5265–5274.
- [55] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* (2004), 600–612.
- [56] Emily Wenger, Shawn Shan, Haitao Zheng, and Ben Y Zhao. 2023. Sok: Anti-facial recognition technology. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 864–881.
- [57] Zihao Xiao, Xianfeng Gao, Chilin Fu, Yinpeng Dong, Wei Gao, Xiaolu Zhang, Jun Zhou, and Jun Zhu. 2021. Improving transferability of adversarial patches on face recognition with generative models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 11845–11854.
- [58] Xiao Yang, Yinpeng Dong, Tianyu Pang, Hang Su, Jun Zhu, Yuefeng Chen, and Hui Xue. 2021. Towards face encryption by generating adversarial identity masks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3897–3907.
- [59] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. 2014. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923* (2014).
- [60] Bangjie Yin, Wenxuan Wang, Taiping Yao, Junfeng Guo, Zelun Kong, Shouhong Ding, Jilin Li, and Cong Liu. 2021. Adv-makeup: A new imperceptible and transferable attack on face recognition. *arXiv preprint arXiv:2105.03162* (2021).
- [61] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujia Yang. 2022. Styleheat: One-shot high-resolution editable talking face generation via pre-trained stylegan. In *European conference on computer vision*. Springer, 85–101.
- [62] Yaoyao Zhong and Weihong Deng. 2022. Opom: Customized invisible cloak towards face privacy protection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 3 (2022), 3590–3603.
- [63] Fengfan Zhou, Qianyu Zhou, Bangjie Yin, Hui Zheng, Xuequan Lu, Lizhuang Ma, and Hefei Ling. 2024. Rethinking Impersonation and Dodging Attacks on Face Recognition Systems. In *ACM Multimedia 2024*.

A Algorithm Pseudocode

Algorithm 1 DynTracker: Dynamic Facial Recognition Strategy

```

1: function DYNTRACKER(facial recognition model FR, initial gallery database  $\mathcal{D}_g$ , query images  $\tilde{\mathcal{X}}_q$ )
2:   Initialize  $t \leftarrow 0$ ,  $\mathcal{X}_{ee} \leftarrow \{\}$ ,  $\mathcal{D}_g(0) \leftarrow \mathcal{D}_g$  and  $\tilde{\mathcal{X}}_q(0) \leftarrow \tilde{\mathcal{X}}_q$ 
3:   while True do
4:      $t \leftarrow t + 1$ 
5:      $\mathcal{X}_{ee}(t) \leftarrow \{x_q \in \tilde{\mathcal{X}}_q(t-1) \mid \text{FR}(x_q; \mathcal{D}_g(t-1)) = y_{ee}\}$  ▷ Recognize Trackee query images
6:      $\mathcal{D}_g(t) \leftarrow \mathcal{D}_g(t-1) \cup \{(x, y_{ee}) \mid x \in \mathcal{X}_{ee}(t)\}$  ▷ Update FR model's gallery database
7:      $\tilde{\mathcal{X}}_q(t) \leftarrow \tilde{\mathcal{X}}_q(t-1) \setminus \mathcal{X}_{ee}(t)$ 
8:     if  $|\mathcal{X}_{ee}(t)| > 0$  then
9:        $\mathcal{X}_{ee} \leftarrow \mathcal{X}_{ee} \cup \mathcal{X}_{ee}(t)$  ▷ Merge recognized Trackee query images
10:    else
11:      Break ▷ Terminate if no more images can be recognized
12:    end if
13:  end while
14:  return  $\mathcal{X}_{ee}$ 
15: end function

```

Algorithm 2 DivTrackee: Diversity-Promoting Anti-Facial Recognition

```

1: function DIVTRACKEE(generator  $G_\theta$ , clean Trackee image  $x_{ee}$ , auxiliary dataset  $\mathcal{X}_{aux}$ , target text prompt  $p_{targ}$ , dissimilarity function  $D_f$ , e4e encoder  $F_\phi$ , CLIP image encoder  $E_I$ , CLIP text encoder  $E_T$ , current queue  $Q$ , hyperparameters  $\alpha_1, \alpha_2, \alpha_3, \alpha_4, \delta, m, \lambda, S$ )
2:    $G_{\theta^*} \leftarrow \text{FineTune}(G_\theta, F_\phi, x_{ee})$  ▷ Fine-tune the StyleGAN generator according to Appendix A.1
3:   Initialize  $w_{init} \leftarrow F_\phi(x_{ee})$  and  $w_0 \leftarrow w_{init}$  ▷ Initialize the latent code based on GAN inversion
4:   for  $s = 0, 1, 2, \dots, S-1$  do
5:     // Compute diversity-promoting adversarial loss terms in Equation 8
6:      $x_{aux} \leftarrow$  Randomly select an image from  $\mathcal{X}_{aux}$ 
7:      $\mathcal{L}_{adv}(w_s) = -D_f(G_{\theta^*}(w_s), x_{ee})$ 
8:      $\mathcal{L}_{guide}(w_s) = \max(0, D_f(G_{\theta^*}(w_s), x_{aux}) - \delta)$ 
9:     if  $|Q| > 0$  then
10:       $\mathcal{L}_{div}(w_s) = -\frac{1}{m} \sum_{x \in Q} D_f(G_{\theta^*}(w_s), x)$  ▷ Note that if  $|Q| = 0$ ,  $\mathcal{L}_{div} = 0$ 
11:    end if
12:    // Compute visual quality loss terms in Equation 12
13:     $\mathcal{L}_{align}(w_s) = 1 - \cos(E_I(G_{\theta^*}(w_s)) - E_I(G_{\theta^*}(w_{init})), E_T(p_{targ}) - E_T(p_{src}))$  ▷ Source text prompt  $p_{src}$  is set as "face"
14:     $\mathcal{L}_{latent}(w_s) = \|w_s - w_{init}\|_2$ 
15:    // Optimize latent code using gradient descent based on the total loss
16:     $\mathcal{L}_{tot}(w_s) = \mathcal{L}_{adv}(w_s) + \alpha_1 \cdot \mathcal{L}_{guide}(w_s) + \alpha_2 \cdot \mathcal{L}_{div}(w_s) + \alpha_3 \cdot \mathcal{L}_{align}(w_s) + \alpha_4 \cdot \mathcal{L}_{latent}(w_s)$ 
17:     $w_{s+1} = w_s - \lambda \nabla_w \mathcal{L}_{tot}(w_s)$ 
18:  end for
19:  if  $|Q| < m$  then
20:    Append  $G_{\theta^*}(w_S)$  to the end of  $Q$ 
21:  else
22:    Update  $Q$  with  $G_{\theta^*}(w_S)$  in a FIFO manner ▷ Ensure the length of  $Q$  not exceeding  $m$ 
23:  end if
24:  return  $G_{\theta^*}(w_S)$ 
25: end function

```

A.1 Generator Fine-Tuning

To achieve better editing performance for out-of-distribution (OOD) face images while maintaining image quality, we adopt the techniques proposed by [44] to fine-tune the generator for each Trackee query image. More specifically, the generator fine-tuning step can be cast into

Table 5: Comparisons on TSRs (%) of static and dynamic FR strategies against various noise-based AFR methods. For each entry, the first number stands for TSR against the static FR strategy, and the second represents TSR against our DynTracker.

Dataset	Method	MobileFace		WebFace		VGGFace		MagFace	
		Clean	Protected	Clean	Protected	Clean	Protected	Clean	Protected
FaceScrub	PGD	7.42/84.54	30.22/81.69	10.69/91.24	38.84/90.35	10.48/92.15	39.62/90.24	28.47/98.76	48.32/97.66
	MI-FGSM	6.15/83.68	29.15/79.25	9.35/90.28	40.84/87.69	10.06/91.41	39.62/90.24	31.93/98.94	47.36/97.53
	TI-DIM	4.62/83.26	27.56/79.54	6.76/90.22	41.80/88.06	11.34/90.23	40.53/88.15	25.57/97.80	45.87/96.68
	TIP-IM	1.98/83.96	20.28/80.86	4.26/91.38	28.65/89.10	5.13/91.81	30.82/88.30	13.82/97.10	40.56/96.54
PubFig	PGD	10.79/80.14	29.97/77.12	8.34/82.65	35.71/79.56	11.13/85.08	39.25/80.93	28.68/94.46	52.38/93.34
	MI-FGSM	8.63/81.27	31.45/75.26	8.42/85.82	32.63/78.94	9.56/87.13	33.47/81.06	29.78/94.37	51.88/93.52
	TI-DIM	6.87/81.69	32.12/76.54	7.25/82.53	33.64/79.32	10.18/85.06	34.51/82.15	21.73/94.43	40.69/93.57
	TIP-IM	3.07/81.89	31.76/76.02	5.18/80.28	35.62/78.08	6.13/83.04	34.74/80.76	13.48/94.11	41.51/92.34
UMDFaces	PGD	9.68/83.27	32.58/79.97	12.83/88.87	38.65/85.66	10.28/89.12	39.82/86.32	26.74/96.73	49.74/95.42
	MI-FGSM	7.87/82.76	32.14/78.52	9.48/87.37	37.81/84.45	13.77/88.63	39.67/85.88	21.43/96.13	45.17/95.26
	TI-DIM	6.68/82.87	30.90/77.06	7.78/85.07	29.27/82.29	9.87/87.23	36.20/83.15	18.78/95.34	49.28/94.18
	TIP-IM	2.78/79.37	24.87/74.25	5.77/86.23	25.92/80.04	3.78/87.54	27.94/81.62	15.87/93.34	24.83/92.30
CelebA-HQ	PGD	24.37/89.60	23.28/87.66	26.10/90.62	22.68/87.10	29.60/89.65	27.32/87.54	31.02/91.37	29.62/87.23
	MI-FGSM	24.88/88.72	22.16/86.48	26.72/90.62	26.90/91.82	27.44/91.37	25.92/90.69	28.73/92.05	24.80/91.60
	TI-DIM	21.64/90.58	23.62/91.35	22.36/92.60	25.42/90.82	23.70/90.85	27.28/91.06	26.10/92.52	25.60/91.38
	TIP-IM	19.72/88.60	27.96/86.32	22.78/89.65	28.64/90.45	24.66/91.38	26.80/91.52	26.74/90.56	25.60/91.37

an optimization problem with respect to the weight parameters of G_θ as follows:

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathcal{L}_{\text{LPIPS}}(\mathbf{x}_{\text{ee}}, G_\theta(\mathbf{w}_{\text{init}})) + \lambda \cdot \|\mathbf{x}_{\text{ee}} - G_\theta(\mathbf{w}_{\text{init}})\|_2,$$

where $\mathcal{L}_{\text{LPIPS}}$ denotes the LPIPS perceptual loss, \mathbf{w}_{init} is the inverse latent code of Trackee’s clean image \mathbf{x}_{ee} , and $\lambda > 0$ is a hyperparameter. In our experiments, we set λ to 0.5 and fine-tuned the generator for 450 epochs.

B Other Experimental Details

In addition to the main experimental setup described in Section 6, we lay out all other necessary implementation details for completeness.

Configuration. For all the other non-Trackee identities, 10% of their facial images are used as query images, while the remaining images are added to the gallery database of the FR model employed by Tracker. During the process of the AFR-protected image generation, we use the Adam optimizer for the gradient descent steps in Equation 7, where the hyperparameters β_1 and β_2 are set to 0.9 and 0.999, respectively. For the face verification experiments, in addition to the trackee’s other identities, each identity selects one image as the query image.

Dataset. The FaceScrub dataset contains 106, 863 photos of 530 celebrities with annotated names and genders, which are retrieved from the Internet and captured in real-world scenarios. The PubFig dataset contains 58, 797 images of 200 individuals with annotated ages, races, and occupations, which are largely various in pose, lighting, and scene. The UMDFaces dataset contains 367, 888 images of 8, 277 individuals with annotated genders and postures, where faces are located by human-curated bounding boxes. We use the CelebA-HQ dataset organized by identity [34], which contains a total of 307 individuals and 5, 478 images, with each identity having more than 15 images.

Feature Extractor. For the FR model employed by Tracker, we consider 4 pre-trained facial feature extractors in our evaluation, including MobileFace [8] trained on MS1MV2 using MobileFaceNet, WebFace [59] that is trained on CASIA-WebFace using Inception ResNet [59], VGGFace [6], which is pre-trained on VGG-Face2 using Inception ResNet [59], and MagFace [33] trained on MS1MV2 using MagFace loss and ResNet. Note that these feature extractors are set to be different from those used for generating AFR protections.

C Additional Experiments

In this section, we conduct additional experiments to study DynTracker’s strength in breaking existing AFR protection schemes and the effectiveness of our DivTrackee in enhancing protection efficacy against dynamic FR strategies while preserving image visual quality.

C.1 Evaluation of Adversarial Noise-Based AFR

In Section 6.1, we evaluate the efficacy of different generative-based AFR protection schemes against static and dynamic FR strategies. To more comprehensively examine the effectiveness of our DynTracker, we further evaluate the performance of existing adversarial noise-based AFR methods, including PGD [32], MI-FGSM [13], TI-DIM [14], and TIP-IM [58], under the same experimental settings. The evaluation

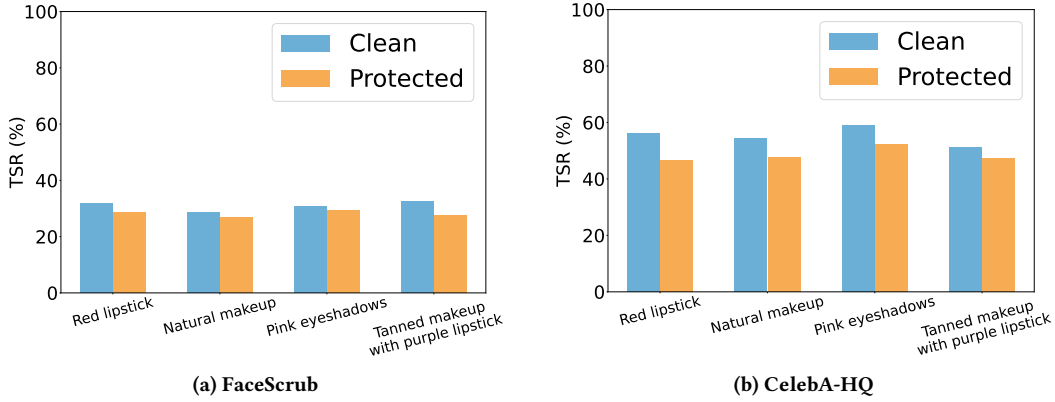


Figure 7: Comparisons of tracking success rates of DynTracker across various target text prompts, which are used to guide the generation process of AFR-protected Trackee images, on two face datasets: (a) FaceScrub, and (b) CelebA-HQ.

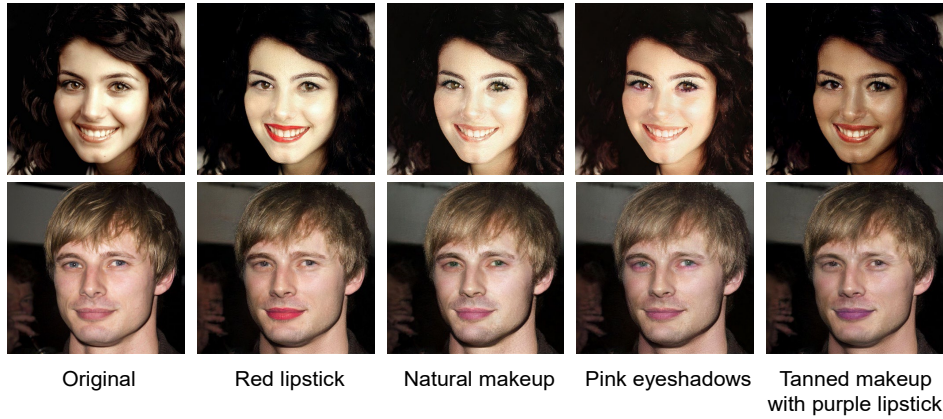


Figure 8: Visualizations of the generated AFR-protected Trackee facial images on CelebA-HQ across different text prompts.

Table 6: Quantitative comparisons of visual quality between different AFR methods on FaceScrub.

Method	AMT-GAN	Clip2Protect	DiffAM	DivTrackee
SSIM (\uparrow)	0.43	0.38	0.40	0.38
PSNR (\uparrow)	12.09	11.59	11.78	11.55

results are demonstrated in Table 5, which again confirms the limited protection success of existing AFR protection schemes against our DynTracker, regardless of Tracker initially holding a clean or protected Trackee image.

C.2 Influence of Target Text prompt

We investigate the impact of different target text prompts p_{targ} on the efficacy of DivTrackee. Figure 7 indicates that for both the “Clean” and “Protected” cases, the tracking success rates against our DynTracker exhibit slight variations across different text prompts, while the visualizations depicted in Figure 8 illustrate that text prompts primarily influence the final makeup effect of the protected images generated by DivTrackee, largely aligned with the expectations of Trackee for AFR in terms of visual quality.

C.3 Quantitative Results on Visual Quality

We adopt the metrics of *Structural Similarity Index Measure* (SSIM) [55] and *Peak Signal to Noise Ratio* (PSNR) to measure the visual quality of generated protected facial images quantitatively. Specifically, SSIM measures the similarity between two images by comparing the conditions of luminance, contrast, and structure; meanwhile, PSNR defines the similarity as the peak signal-to-noise ratio between two images. Besides, a higher value of SSIM/PSNR indicates better visual quality. In Table 6, while our method yields slightly lower SSIM and PSNR values compared to other approaches, the visual quality remains acceptable for users. We present the protected images generated using various prompts in Figure 8. We can observe that the generated images preserve the original facial features while incorporating the appropriate makeup effects.