# D4: <u>D</u>etection of Adversarial <u>D</u>iffusion <u>D</u>eepfakes Using <u>D</u>isjoint Ensembles

**Ashish Hooda**[1*]    **Neal Mangaokar**[2*]    **Ryan Feng**[2]
**Kassem Fawaz**[1]    **Somesh Jha**[1]    **Atul Prakash**[2]
[1]University of Wisconsin-Madison    [2]University of Michigan
{ahooda,kfawaz,jha}@wisc.edu
{nealmgkr,rtfeng,aprakash}@umich.edu

## Abstract

Detecting diffusion-generated deepfake images remains an open problem. Current detection methods fail against an adversary who adds imperceptible adversarial perturbations to the deepfake to evade detection. In this work, we propose **D**isjoint **D**iffusion **D**eepfake **D**etection (D4), a deepfake detector designed to improve black-box adversarial robustness beyond de facto solutions such as adversarial training. D4 uses an ensemble of models over disjoint subsets of the frequency spectrum to significantly improve adversarial robustness. Our key insight is to leverage a redundancy in the frequency domain and apply a saliency partitioning technique to disjointly distribute frequency components across multiple models. We formally prove that these disjoint ensembles lead to a reduction in the dimensionality of the input subspace where adversarial deepfakes lie, thereby making adversarial deepfakes harder to find for black-box attacks. We then empirically validate the D4 method against several black-box attacks and find that D4 significantly outperforms existing state-of-the-art defenses applied to diffusion-generated deepfake detection. We also demonstrate that D4 provides robustness against adversarial deepfakes from unseen data distributions as well as unseen generative techniques.

## 1 Introduction

Significant advances in deep learning are responsible for the advent of "deepfakes", which can be misused by bad actors for malicious purposes. Deepfakes broadly refer to digital media that has been synthetically generated or modified by deep neural networks (DNNs). Modern DNNs such as diffusion models [1, 2, 3, 4, 5, 6] and generative adversarial networks (GANs) [7, 8, 9, 10, 11, 12, 13] are now capable of synthesizing hyper-realistic deepfakes, which can then be used to craft fake social media profiles [14], generate pornography [15], spread political propaganda, and manipulate elections.

The deepfake detection problem asks the defender to classify a given image as deepfake or real. While recent work has made remarkable efforts towards solving the deepfake detection problem, many of these detectors are rendered ineffective by *adversarial examples* (Fig. 1). Specifically, these state-of-the-art detectors often leverage DNNs, and Carlini et al. [17] (amongst others) have shown that such DNNs are vulnerable — the attacker can simply use adversarial perturbation techniques to evade detection [18, 19, 20, 21]. Recent work has shown that these "adversarial deepfakes" can even be crafted in a black-box setting, where the attacker only has query access to the detector [22, 23, 24, 25]. Defending against adversarial examples, in general, has been shown to be a difficult task [26], and is a critical problem in the deepfake detection setting.

---

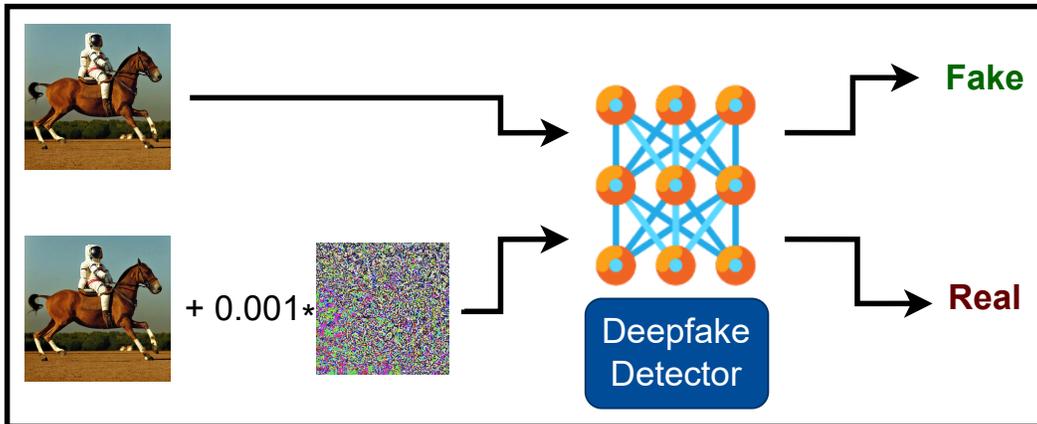*Indicates equal contribution.

Preprint.

Figure 1: Under the non-adversarial setting, the deepfake detector correctly classifies the image produced by Stable Diffusion [16] on the text prompt "a photograph of an astronaut riding a horse." as fake. However, one can flip the detector's prediction by adding imperceptible adversarial perturbations to the deepfake.

Our key intuition to mitigate this problem is to utilize *redundant information in the frequency feature space* of deepfakes to generate *disjoint ensembles* for adversarial deepfake detection. Specifically, we show in Sec. 3.1 that we can achieve good detection performance with only a subset of the features, particularly in the frequency domain. This enables us to build an ensemble of performant classifiers, each using a disjoint set of frequencies. In contrast to traditional ensembles (where each model shares the same set of frequencies), a key advantage of this design is that non-robust frequencies are partitioned across all the models in the ensemble (see Sec. 3.2). Thus, an attacker is no longer able to perturb a single non-robust frequency to evade all models — rather, they must find perturbations to evade multiple sets of disjoint frequencies. D4 thus aims to thwart the attacker's generation of adversarial deepfakes.

To summarize, our key contributions are as follows:

1. We propose D4, a diffusion-generated deepfake detection framework designed to be adversarially robust. D4 builds an ensemble of models that use disjoint partitions of the input features. This is achieved by leveraging redundancy in the feature space. D4 achieves robustness while still maintaining natural deepfake detection average precision scores as high as 93%. (see Sec. 4 for details).

2. Extending the theoretical results by Tramer et al. [27] on dimensionality of adversarial subspaces, we prove new bounds on the maximum number of adversarial directions that can be found under an ensemble with disjoint inputs. Our bounds are tight for both the $\ell_2$ and $\ell_\infty$ perturbation norms (Lemmas 3.1 and 3.2 in Sec. 3.3) and indicate that D4 reduces the dimension of the adversarial subspace, i.e., there are simply fewer adversarial examples to be found.

3. We evaluate D4 against query-based black-box attacks, as well as adaptive frequency and post-processing attacks. Across a variety of diffusion-generated deepfake images, we find that D4 significantly outperforms state-of-the-art defenses such as the recently proposed EnsembleDet [28, 29, 30], suggesting that D4 indeed provides a reduction in dimension of the adversarial subspace. For example, as indicated by our evaluation in Sec. 4, D4 reduces the attack success rate to 28%, whereas baselines incur attack success rates of more than 90%. These improvements also extend to unseen image domains and deepfake generation models.

## 2   Background and Related Work

**Notation.** We consider a distribution $\mathcal{D}$ over $\mathcal{X} \times \mathcal{Y}$, where $\mathcal{X} \subseteq \mathbb{R}^d$ is the input space and $\mathcal{Y} \subseteq \mathbb{Z}^c$ is the finite class-label space. We denote vectors in boldface (e.g., $\mathbf{x}$). We denote a trained classifier as a

function $\mathcal{F} : \mathcal{X} \to \mathcal{Y}$ (the classifier is usually parameterized by its weights $w$, omitted for brevity). We denote the loss function as $\mathcal{L}(\mathbf{x}, y)$. An ensemble classifier is a function $M_{(\mathcal{F}_1, \mathcal{F}_2, ..., \mathcal{F}_n)} : \mathcal{X} \to \mathcal{Y}$ that combines the logit outputs $l_1, l_2, ..., l_n$ of multiple classifiers $\mathcal{F}_1, \mathcal{F}_2, ..., \mathcal{F}_n$ with a voting aggregation function $\mathcal{A} : \mathbb{R}^{n \times c} \to \mathcal{Y}$.

For a classifier $\mathcal{F}$ and input-label pair $(\mathbf{x}, y)$, an adversarial example is a perturbed input $\mathbf{x}'$ such that (1) $\mathbf{x}'$ is misclassified, i.e., $\mathcal{F}(\mathbf{x}') \neq y$ and (2) $||\mathbf{x} - \mathbf{x}'||$ is within a small $\epsilon$ ball, where $||.||$ is a given norm. The value of $\epsilon$ is chosen to be small so that the perturbation is imperceptible to humans.

**Deepfake Image Generation.** In this work, we focus on detection of deepfake images that are generated entirely from scratch using a deep generative model. Two prominent techniques that have achieved state-of-the-art for such generation are GANs [7] and diffusion models [1]. GANs comprise two DNNs: a generator and a discriminator. The generator synthesizes images, while the discriminator attempts to distinguish between real and deepfake samples. Through this adversarial training procedure, GANs learn to generate increasingly realistic and high-quality outputs, and have achieved remarkable success [31]. However, GANs have recently been superceded by diffusion models. These models iteratively add noise to data samples and then remove it, thereby learning to generate images from randomly sampled noise. Diffusion models have now also achieved new state-of-the-art FID scores for image generation benchmarks [5]. Given the quality of images generated, detection of diffusion model deepfakes poses a pressing concern.

**Deepfake Image Detection.** The research community has made rapid progress towards detecting deepfake images. Some efforts propose classification DNNs that operate directly on pixel features [32, 29, 33, 34]. Others have instead trained DNNs using features extracted from the deepfakes, i.e., co-occurrence matrices [35], color-space anomalies [36], convolutional traces [37], texture representations [38], pixel-patches [39], or more recently using neural features extracted from foundation models such as CLIP [40].

One particular line of work that has shown great promise for deepfake detection is leveraging *frequency* features. Specifically, Frank et al. [30] proposed the idea of detecting deepfakes with the Discrete Cosine Transform (DCT) as a pre-processing transform before the DNN. Similar work has also achieved remarkable performance using frequency features — [41, 42], and more recent efforts have emphasized their utility in detecting deepfakes from both GANs and diffusion models [43]. D4 also leverages frequency features, but through a unique disjoint ensembling approach.

**Adversarial Deepfakes.** Unfortunately, regardless of the chosen feature space, the aforementioned detectors have been rendered ineffective in adversarial settings. Specifically, Carlini et al. [17] showed that DNN detectors are vulnerable to adversarial examples — an adversary can add imperceptible adversarial perturbations to a deepfake that evade such detectors, rendering them ineffective. Others have corroborated this observation [24, 25, 44, 23, 22, 18, 19, 20, 21].

An adversary can construct these "adversarial deepfakes" in either a white-box setting (with complete knowledge of the detector's weights and parameters), or black-box setting (with only query access to the detector). In this work, we focus on the black-box setting. The DNN models of deepfake detectors in the real world are often hidden from the users and likely to be black-box services, such as those already offered by Intel [45], Deepware [46], Reality Defender [47] and Sensity AI [48]. These services are typically available through web-based platforms or through API access. Moreover, defending against white-box attacks is a challenging open problem on all vision tasks. We leave defense in the white-box setting to future work.

**Adversarial Deepfake Detection.** Existing defenses for adversarial deepfakes can be broadly classified into (a) training time defenses (which adjust the training process), and (b) inference-time defenses.

(a) Training time. The original training time defense is adversarial training, in which the model is trained on adversarial examples generated during training [49]. However, Carlini et al. [17] suggest that adversarial training alone is unlikely to achieve significant improvement in robustness in the difficult deepfake detection setting (confirmed by our experiments in Sec. 4.2).

Instead of adversarial training, recent work has also proposed using an ensemble of models — in principle, the adversary is then forced to attack multiple models, instead of just one. However, He et al. [50] have shown that arbitrarily ensembling models does not necessarily lead to more robustness.
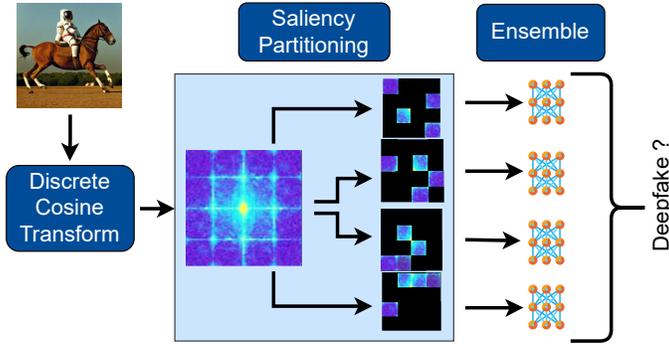
Figure 2: The processing pipeline of D4. It partitions the DCT spectrum of an image into disjoint partitions using a saliency-based approach. Each frequency partition is fed to a separate model that is adversarially trained. A voting mechanism over the ensemble decides the output.
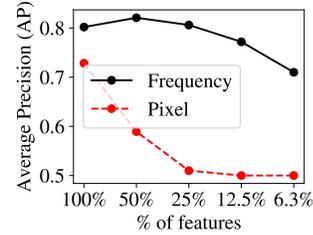


Figure 3: Average Precision (AP) values of a single CNN classifier trained on a fixed subset (randomly selected) of the input features. Redundancy in frequency spectra permits good deepfake detection even when using $\sim 6\%$ of the components (see black line).

Prior work suggests that each model in an ensemble tends to learn the same non-robust features, i.e., an adversary is able to perturb the same features to evade all models [51].

Most recently, Dutta et al. [28] and Khan et al. [52] proposed EnsembleDet and ARDD respectively as ensemble defenses against adversarial deepfakes. Both defenses ensemble multiple different DNN architectures, and posit that the different architectures will learn different features, thereby providing improved robustness. Devasthale et al. [53] take this approach one step further and also adversarially train each model in the ensemble.

Note that while D4 is also an ensemble-based detector, it is fundamentally different than existing approaches since it does not require different architectures to avoid learning the same features — we instead leverage artifact redundancy to design frequency-partitioned ensembles (see Sec. 3 for a more detailed explanation). Furthermore, D4's partitioning approach to achieve disjoint ensembles is novel, and it offers theoretical and empirical advantages where such feature redundancy is available, e.g., generated deepfakes, allowing it to outperform prior work (see Sec. 4).

(b) Inference time. Another class of defenses focuses on removing the effect of the adversarial perturbation without any changes made to the underlying detection technique — unfortunately, these approaches are computationally intensive, e.g., upto 30 minutes per image [23, 54], in comparison to no additional overhead in D4. Nevertheless, both approaches are complementary in that pre-processing could be combined with training-time defenses such as D4.

Finally, we note that there are other inference-time defenses that are not specific to deepfakes. For example, stateful defenses [55] were proposed to defend against black-box adversarial examples, but have been recently shown to be vulnerable [56]. Another group of defenses post-process the detector's response by manipulating the detector's confidence scores [57, 58], but these do not work against hard-label black-box attacks.

## 3   Our Approach

We now present D4 (Fig. 2), a deepfake detection framework that leverages an ensemble of disjoint frequency models to achieve robust detection of diffusion-generated deepfake images. Sec. 3.1 presents our observation of redundant information in the frequency space of deepfakes. Sec. 3.2 details how redundancy allows frequencies to be partitioned between multiple models for robust ensembling and explains our exact frequency partitioning schemes. Finally, Sec. 3.3 provides theoretical insight into why D4 improves adversarial robustness.

### 3.1   Redundant Information in Deepfake Image Spectra

As discussed in Sec. 2, ensembles are a promising approach to adversarial robustness, so long as perturbing the same set of features does not simultaneously evade the individual models. We propose

designing an ensemble that avoids this shortcoming by *disjointly partitioning the input feature space* amongst individual models.

As mentioned earlier, the frequency spectrum of images facilitates deepfake detection because generation techniques leave discriminating artifacts that are more prominent in the frequency space. We additionally observe that these artifacts are spread throughout the frequency feature space, and they are more uniformly spread as compared to pixel space. We confirm this in Fig. 3, which plots the performance of a simple convolutional classifier [2] on an increasingly smaller random subset of the input features.

Our first key insight is thus that disjoint partitioning is feasible for deepfake detection. Specifically, we observe a "redundancy" in frequency-space artifacts — signals relevant for deepfake detection[3] are distributed throughout the frequency spectrum. This observation is best exemplified by the black line in Fig. 3, which plots the deepfake detection performance using increasingly small subsets of the spectrum using the same classifier. Using as few as $\sim 6\%$ of the frequency components yields good deepfake detection performance. We emphasize that this does not hold for subsets of pixels (red line), as signals for detection are not well-distributed in the RGB space. Overall, these findings suggest that the frequency-space contains plenty of redundancy, which can potentially be leveraged to design a robust ensemble.

## 3.2    Leveraging Redundancy to Build a Robust Ensemble

Using the observations from Sec. 3.1, one can craft a robust ensemble by partitioning the frequency components amongst multiple detector models, without hurting natural detection performance. Fig. 2 visualizes this partitioning as part of our ensembling pipeline. Specifically, for each individual model we mask (i.e., zero out) the frequencies not used. For example, consider an ensemble of two such "disjoint" models $F_A$ and $F_B$, with the full-spectrum frequency feature vector $f = [f_1; f_2]$. Then, the input feature vector to $F_A$ is $[f_1; 0]$ and to $F_B$ is $[0; f_2]$. Since the input feature space (frequencies) is not shared amongst the individual models, the adversary cannot simply attack the ensemble by targeting common frequencies.

Furthermore, we note that choice of partitioning scheme, i.e., *how* the frequencies are partitioned plays an important role in robustness of the ensemble. Specifically, the chosen scheme should aim to design all models as "equals" — if some models are less robust than others, then the adversary can target them to overturn the ensemble's decision.

**Saliency Partitioning.** While signals for deepfake detection are distributed throughout the spectrum, there still exists an *adversarial saliency* ordering of these frequencies that determines their robustness for the deepfake detection task. For an ensemble of size $n$, our saliency partitioning technique is aimed at ensuring each model receives a fair proportion of salient frequencies. To this end, we follow [59] and [60] to compute saliency values for all frequencies. This is achieved by adversarially perturbing deepfake $x$ to $x + \delta^x$, where $\delta^x$ is the perturbation computed with the Carlini-Wagner $\ell_2$ attack for 1000 steps. Then, we compute saliency $s_i$ for the $i^{th}$ frequency as

$$s_i = \mathbb{E}_{x \in \mathcal{X}} \nabla f(x + \delta^x)_i \cdot \delta_i^x \tag{1}$$

where subscript $i$ denotes the $i^{th}$ component. Intuitively, higher gradients and larger perturbation magnitudes imply larger saliencies. Frequencies are then ordered by their saliencies, and distributed in a round-robin fashion amongst the models.

## 3.3    Adversarial Subspace of Disjoint Ensembles

Given the partitioning approach in Sec. 3.2, we now show that an ensemble of such disjoint frequency models increases robustness against adversarial examples by reducing the dimension of the adversarial subspace. For a single model $\mathcal{F}$ and input $\mathbf{x}$, Tramer et al. [27] approximate the $k-$dimensional adversarial subspace as the span of orthogonal perturbations $\mathbf{r_1}, \cdots, \mathbf{r_k} \in \mathbb{R}^d$ such that $\mathbf{r_i^\intercal g} \geq \gamma \; \forall \; 1 \leq i \leq k$ where $\mathbf{g} = \nabla_{\mathbf{x}} \mathcal{L}_{\mathcal{F}}(\mathbf{x}, y)$, $\mathcal{L}_{\mathcal{F}}$ is the loss function used to train $\mathcal{F}$, and $\gamma$ is the

---

[2]We use the same architecture as Frank et al. [30]

[3]Frank et al. [30] show that these artifacts manifest in the form of a grid-like spectral pattern, and attribute their presence to the upsampling process in generative models. Existing work proposes detectors that leverage the *entire frequency spectrum* for deepfake detection.

increase in loss sufficient to cause a mis-classification. For perturbations satisfying the $\ell_2$-norm, i.e. $||\mathbf{r_i}||_2 \leq \epsilon \; \forall \; 1 \leq i \leq k$, the adversarial dimension $k$ is bounded by $\frac{\epsilon^2||\mathbf{g}||_2^2}{\gamma^2}$ (tight). In what follows, we extend this result and provide bounds for dimensionality of the shared adversarial subspace between $n$ disjoint models. We provide tight bounds for both $\ell_2$ and $\ell_\infty$ norms in Lemma 3.1 and Lemma 3.2 respectively (with detailed proofs in Appendix 7.1). We consider these bounds for two voting mechanisms: (1) majority, where the ensemble outputs deepfake if at least $\lceil n/2 \rceil$ classifiers predict deepfake, and (2) at-least-one, where the classifier outputs deepfake if at least one classifier predicts deepfake, otherwise it outputs real.

**Lemma 3.1.** *Given $n$ disjoint models, $\mathcal{F}_1, ..., \mathcal{F}_n$, having gradients $\mathbf{g_1}, \cdots, \mathbf{g_n} \in \mathbb{R}^d$ for input-label pair $(\mathbf{x}, y)$ (where $\mathbf{g_j} = \nabla_\mathbf{x}\mathcal{L}_{\mathcal{F}_j}(\mathbf{x}, y)$), the maximum number $k$ of orthogonal vectors $\mathbf{r_1}, \mathbf{r_2}, \cdots, \mathbf{r_k} \in \mathbb{R}^d$ satisfying $||\mathbf{r_i}||_2 \leq \epsilon$ and, $\mathbf{r_i}^\intercal\mathbf{g_j} \geq \gamma_j$ for all $1 \leq j \leq n$ (at-least-one voting) or for at least $\lceil \frac{n}{2} \rceil$ models (majority voting), for all $1 \leq i \leq k$ is given by:*

$$k = \min\left(d, \left\lfloor \frac{\epsilon^2}{\left(\sum\limits_{j=1}^{n} \gamma_j\right)^2} \sum_{j=1}^{n} ||\mathbf{g_j}||_2^2 \right\rfloor\right) \tag{2}$$

*(at-least-one voting)*

$$k \leq \min\left(d, \left\lfloor \max_{|K|=\lceil \frac{n}{2}\rceil} \frac{\epsilon^2}{\left(\sum\limits_{j \in K} \gamma_j\right)^2} \sum_{j \in K} ||\mathbf{g_j}||_2^2 \right\rfloor\right) \tag{3}$$

*(majority voting)*

$$k \geq \min\left(d, \left\lfloor \min_{|K|=\lceil \frac{n}{2}\rceil} \frac{\epsilon^2}{\left(\sum\limits_{j \in K} \gamma_j\right)^2} \sum_{j \in K} ||\mathbf{g_j}||_2^2 \right\rfloor\right) \tag{4}$$

*(majority voting)*

**Lemma 3.2.** *Given $n$ disjoint models, $\mathcal{F}_1, ..., \mathcal{F}_n$, having gradients $\mathbf{g_1}, \cdots, \mathbf{g_n} \in \mathbb{R}^d$ for input-label pair $(\mathbf{x}, y)$ (where $\mathbf{g_j} = \nabla_\mathbf{x}\mathcal{L}_{\mathcal{F}_j}(\mathbf{x}, y)$), the maximum number $k$ of orthogonal vectors $\mathbf{r_1}, \mathbf{r_2}, \cdots, \mathbf{r_k} \in \mathbb{R}^d$ satisfying $||\mathbf{r_i}||_\infty \leq \epsilon$ and $\mathbb{E}[\mathbf{g_j}^\intercal\mathbf{r_i}] \geq \gamma_j$ for all $1 \leq j \leq n$ (at-least-one voting) or for at least $\lceil \frac{n}{2} \rceil$ models (majority voting), for all $1 \leq i \leq k$*

$$k = \min\left(d, \left\lfloor \min\left(\frac{\epsilon^2||\mathbf{g_1}||_1^2}{n^2\gamma_1^2}, ..., \frac{\epsilon^2||\mathbf{g_n}||_1^2}{n^2\gamma_n^2}\right)\right\rfloor\right) \tag{5}$$

*(at-least-one voting)*

$$k = \min\left(d, \left\lfloor median\left(\frac{\epsilon^2||\mathbf{g_1}||_1^2}{n^2\gamma_1^2}, ..., \frac{\epsilon^2||\mathbf{g_n}||_1^2}{n^2\gamma_n^2}\right)\right\rfloor\right) \tag{6}$$

*(majority voting)*

**Implications.** If all $n$ disjoint models in the ensemble are "near-identical" (as expected per our saliency partitioning scheme), i.e., $||\mathbf{g_1}||_2^2 \approx \cdots \approx ||\mathbf{g_n}||_2^2$ and $\gamma_1 \approx \cdots \approx \gamma_n$, then Lemma 3.1 for at-least-one voting reduces to $k \approx \min\left(d, \left\lfloor \frac{\epsilon^2||\mathbf{g_1}||_2^2}{n\gamma_1^2} \right\rfloor\right)$. This implies that an ensemble of $n$ disjoint models offers potential reduction in dimensionality of the adversarial subspace by a factor of $n$ compared to any individual constituent disjoint model. Similar interpretation holds for Lemma 3.2, where reductions are now by a factor of $n^2$. Next, in Sec. 4, we empirically demonstrate that this reduction in dimension of adversarial subspace leads to improved performance against black-box adversarial examples.

## 4 Experimental Evaluation

Our experiments broadly aim to answer the following questions:

| Detector | No Attack | Attack (ASR) | | | | | |
|---|---|---|---|---|---|---|---|
| | (AP) | SurFree | HSJA | QEBA | Triangle | Boundary | SignOPT |
| **EnsembleDet** | 100% | 100% | 100% | 100% | 77% | 100% | 100% |
| **ARDD** | 100% | 100% | 100% | 100% | 99% | 90% | 100% |
| **ARDD-AT** | 100% | 91% | 23% | 97% | 18% | 19% | 43% |
| D4 (**SIZE=1**) | 98% | 93% | 11% | 53% | 51% | **6%** | 10% |
| D4 (**SIZE=4**) | 93% | **28%** | **3%** | **2%** | **0%** | 8% | **8%** |

Table 1: An attacker achieves lower attack success rates (ASRs) when attacking D4 (SIZE=4) as compared to the baselines. Attacks are launched on LSUN bedroom deepfake images from an LDM diffusion model, with a query budget of 50k queries and $\ell_2$ perturbation budget of $\epsilon = 10$.
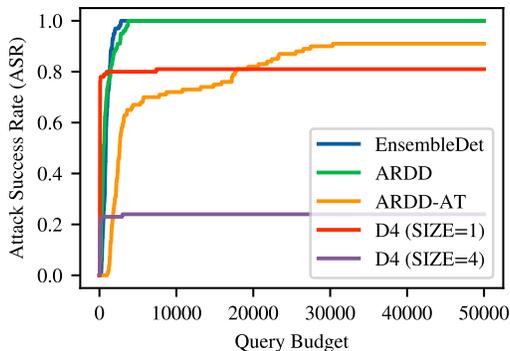


Figure 4: SurFree ASR vs. query budget against all detectors, on LSUN bedroom deepfakes from the LDM diffusion model.

**Q1.** How robust are D4 ensembles against black-box attackers, and how does D4 raise the attacker's cost of creating adversarial deepfakes?

**Q2.** Does D4's robustness hold for deepfakes from different diffusion models? How well does D4 generalize to deepfakes from unseen image domains, diffusion models, or even different generative architectures, i.e., GANs?

**Q3.** How does D4 fare against traditional post-processing, and more adaptive frequency-based attacks?

In the following sections, we address these questions by detailing our setup and experiments.

## 4.1 Experimental Setup

We adopt the following settings to evaluate D4 for detecting deepfakes and adversarial deepfakes.

**Datasets and Pre-processing.** We perform our experiments using real images from the LSUN Bedroom dataset [61], and the CelebaHQ dataset [8]. For each of these datasets, we obtain deepfake images from the LDM [4], DDIM [3], and PNDM [6] diffusion models. Deepfakes for bedroom are sourced from those already generated by prior work [43], and we generate deepfakes for CelebaHQ ourselves using models from the HuggingFace public model repository [62]. For our generalization evaluation with GANs, we source ProGAN [8], StyleGAN [9], and Diff-StyleGAN2 [11] images again from [43]. For any given diffusion model and dataset, the training set comprises 39k training images for each of the deepfake and real classes, 1k validation, and 10k testing (as per [43]). All images are resized to 256x256, and then center-cropped to 224x224.

**Baselines.** We implement five baselines to evaluate D4. The first two baselines are the pixel-space ensemble defenses *EnsembleDet* [28] and *ARDD* [52] proposed by prior work. EnsembleDet comprises three models with the EfficientNet, XCeption, and Resnet50 architectures. ARDD is similar, but instead comprises the VGG16, InceptionV3, and Xception architectures. Our third baseline *ARDD-AT* is also from prior work [53], which changes the architectures of ARDD to VGG19, Vision Transformer, Wide-ResNet, and also adversarially trains each individual model using the PGD-150 attack on deepfake images only. Our fourth baseline is simply an adversarially trained version of a full-spectrum frequency space detector — this is equivalent to a D4 ensemble of

| Model | Detector | No Attack (AP) | Attack (ASR) |
|-------|----------|----------------|--------------|
| LDM | D4 (SIZE=1) | 98% | 95% |
|     | D4 (SIZE=4) | 93% | **28%** |
| DDIM | D4 (SIZE=1) | 98% | 97% |
|      | D4 (SIZE=4) | 78% | **15%** |
| PNDM | D4 (SIZE=1) | 97% | 100% |
|      | D4 (SIZE=4) | 97% | **0%** |

Table 2: D4 (SIZE=4) only trained on LSUN Bedroom is able to generalize its robustness to CelebaHQ, an entirely different domain unseen during training time.

| Model | Detector | Attack (ASR) | | | | | |
|-------|----------|---------|------|------|----------|----------|---------|
|       |          | SurFree | HSJA | QEBA | Triangle | Boundary | SignOPT |
| LDM | D4 (SIZE=1) | 93% | 11% | 53% | 51% | **6%** | 10% |
|     | D4 (SIZE=4) | **28%** | **3%** | **2%** | **0%** | 8% | **8%** |
| DDIM | D4 (SIZE=1) | 95% | 18% | 99% | **4%** | 26% | 24% |
|      | D4 (SIZE=4) | **9%** | **13%** | **5%** | 6% | **11%** | **15%** |
| PNDM | D4 (SIZE=1) | 100% | 87% | 77% | 16% | 91% | 89% |
|      | D4 (SIZE=4) | **0%** | **0%** | **0%** | **0%** | **0%** | **0%** |

Table 3: D4 (SIZE=4) continues to reduce the attacker's ASR more than the baselines, even when trained and tested on adversarial deepfakes from diffusion models other than LDM.

size 1, i.e., *D4 (SIZE=1)* without multiple models or disjoint partitioning. We perform adversarial training for D4 (SIZE=1) with PGD-50 attacks [49] using the TRADES loss [63]. Finally, we use the pretrained CNNDet detector from [29] as an additional baseline for our generalization experiments in Sec. 4.3.3, since this detector was designed for detection of unseen generative models [4]. All baseline models (except CNNDet) are trained using the Adam optimizer [64] with an initial learning rate of 0.0001, batch size of 32, and a maximum of 20 epochs.

**Attacks.** We evaluate D4 and baselines against six popular black-box attacks spanning both the gradient-based and random search-based categories. For gradient-based attacks, we select HSJA [65] and QEBA [66]. For random search-based attacks, we select SurFree [67], Triangle [68], Boundary [69], and SignOPT [70]. We employ the default hyperparameters for each attack, and focus our attack evaluation on the $\ell_2$ norm with a standard perturbation $\epsilon = 10$. We impose a query budget of 50,000 queries on each attack following [57] (roughly equivalent to between $50 - $75 as per modern MLaaS platforms such as Clarifai).

D4 **Architecture and Training.** We implement a D4 ensemble *D4 (SIZE=4)* comprising four models. Each model follows a ResNet50 architecture, and the 2D-Discrete Cosine Transform (DCT) is used to convert images to the frequency space before distributing frequencies amongst the models. Each individual model is also adversarially trained using the same procedure as the D4 (SIZE=1) baseline.

**Metrics.** We follow prior work [29, 43, 40] and use average precision (AP) to measure the natural, i.e., non-adversarial, unperturbed deepfake detection performance of D4. We then consider an adversary that attempts to perturb deepfakes to the "real" class, and employ *attack success rate*, i.e., ASR (fraction of successfully perturbed deepfakes) as our performance metric for robustness. Lower ASR implies that the detector is more effective against adversarial deepfakes.

---

[4]The pre-trained models for the more recent generalization detector from [40] are currently unavailable.

| Detector | LDM | DDIM | PNDM | StyleGAN | ProGAN | Diff-StyleGAN2 |
|----------|-----|------|------|----------|--------|----------------|
| **CNNDet** | 62% (-) | 68% (100%) | 64% (100%) | **95** (100%) | **100%** (100%) | **100%** (100%) |
| D4 **(SIZE=4)** | **93%** **(28%)** | **70%** **(9%)** | **79%** **(19%)** | 67% (**14%**) | 59% (28%) | 67% (11%) |
| **Both** | **93%** **(28%)** | **70%** **(9%)** | **79%** **(19%)** | 80% (54%) | **100%** (**63%**) | 94% (**67%**) |

Table 4: Generalization of CNNDet (trained on ProGAN) and D4 (SIZE=4) (trained on LDM) to diffusion and GAN deepfakes that were unseen during training. Results are presented in the following format: non-adversarial AP (ASR).

## 4.2 Robustness Against Adversarial Examples

We now present performance of D4 and baselines under the six attacks described in Sec. 4.1. For each detector, we present ASR over 100 images under a 50k query budget.

Results for each baseline detector are presented in rows 1-4 of Tab. 1. Notably, these baselines achieve excellent AP scores of $\sim 100\%$ on non-adversarial deepfakes. However, we find that for each detector at least one attack achieves ASR $> 90\%$, rendering it entirely ineffective. Interestingly, the random-search based SurFree attack is particularly effective against all the baselines, e.g., 91% and 93% against the the ARDD-AT and D4 (SIZE=1) baselines respectively. In contrast D4 (SIZE=4) (presented in row 5) is not vulnerable to any ASRs over 28% (again, achieved by SurFree). In some cases, it can even reduce ASR to $< 3\%$. Furthermore, D4 is able to withstand these attacks without much drop in performance on non-adversarial deepfakes (93% AP). We re-emphasize that this is only achievable due to the redundancy observation from Sec. 3.1.

To better visualize how D4 raises the cost of an attack, we also plot SurFree ASR against attack query budget for all detectors in Fig. 4. An attacker can typically achieve over 80% ASR against all baselines within 20k queries. However, D4 (SIZE=4) continues to prevent ASR over 28% even at over double, i.e. 50k queries.

## 4.3 Generalization of Robustness

We now expand beyond the standard setting discussed in Sec. 4.2 and instead consider the generalization capabilities of D4 in different contexts. First, we repeat the experiments from Sec. 4.2 for other diffusion models. Second, we consider a more difficult setting and extend our evaluation of D4 to adversarial deepfakes from models and domains *unseen at training time*. While this generalization is known to be possible for non-adversarial deepfake detection [29, 40], to the best of our knowledge generalization for adversarial deepfake detection has not yet been explored.

### 4.3.1 Different Diffusion Models

Tab. 3 presents the results of repeating the experiments in Sec. 4.2, but now for adversarial deepfakes from the DDIM and PNDM diffusion models. As expected, we observe that the trends continue to hold — in fact, D4 (SIZE=4) on PNDM is able to completely prevent all six attacks, i.e., ASR=0 across all images. For reference, the D4 (SIZE=1) baseline (one of the stronger ones from Tab. 1) is again vulnerable to at least one attack with $> 90\%$ ASR.

### 4.3.2 Unseen Image Domains

Tab. 2 presents the results of evaluating the D4 (SIZE=4) models from Tab. 3 (which were trained only on bedroom images) on adversarial deepfakes from an entirely different data domain, i.e., human faces. We focus on the strongest SurFree attack. Again, D4 reduces ASR significantly more than the baselines, with minimal cost to non-adversarial deepfake detection. One exception is DDIM, for which non-adversarial AP drops to 78%. However, the decrease in an attacker's ASR from 97% to 15% likely offsets this drop in adversarial settings.

On potential explanation for D4's success over the baselines here is as follows: any detector using the full feature set may overfit to a small set of non-robust features unique to the specific training image domain. On the other hand, D4's disjoint partitioning of non-robust features forces each model to learn more from the robust artifacts caused by the diffusion model. Overall, D4 suggests that generalizing adversarial robustness to a variety of domains is possible.

### 4.3.3 Unseen Generative Models

We now evaluate D4's generalization to adversarial deepfakes from models unseen during training time, including different diffusion models, as well as other architectures, i.e., GANs. To this end, we select the D4 (SIZE=4) model from Tab. 1 trained on LSUN bedroom only. Since prior work has only focused on generalization in the non-adversarial context, our baseline is the popular CNNDet detector [29] renown for its detection capabilities across a wide variety of GANs. CNNDet is a pixel-space detector trained only on ProGAN deepfake images, using heavy JPEG and blurring data augmentation.

| Detector | Noise | Blur | JPEG | Freq-Peaks |
|----------|-------|------|------|------------|
| D4 (SIZE=4) | 78% | 93% | 83% | 92% |

Table 5: AP scores on LSUN Bedroom LDM deepfakes that are subject to post-processing, or an adaptive frequency-peaks attack.

Row 1 of Tab. 4 presents CNNDet's AP scores for non-adversarial deepfakes and ASR for adversarial deepfakes, across the variety of diffusion and GAN models listed in Sec. 4.1. As expected, CN-NDet performs well for detection of non-adversarial GAN deepfake images (all AP scores $> 80\%$). However, it performs poorly for diffusion deepfakes which can be explained by prior work's observation that the high frequency artifacts differ between GANs and diffusion models [43]. Under the adversarial setting, it is rendered completely ineffective with 100% ASR for both GAN and diffusion deepfakes. This also suggests that simultaneously acheiving both generalization and robustness against adversarial examples is a challenging problem. Row 2 presents D4 (SIZE=4) scores, which exhibit the opposite trend — it is able to generalize better (both non-adversarial and adversarial) for diffusion deepfakes, but worse for GANs (since it is trained on diffusion images).

The above observations suggest that *combining the two detectors* may yield improvements. To this end, row 3 presents the results of ensembling D4 (SIZE=4) and CNNDet using an at-least-one voting scheme. The resulting aggregate detector "merges" the benefits to an extent, presenting AP scores $>= 80\%$ and $>= 70\%$ for non-adversarial GAN and diffusion model deepfakes respectively. Furthermore, ASRs are significantly reduced from the 100% of CNNDet. Overall, this indicates that D4 is complementary to existing detectors, and can be combined to improve adversarial (or even non-adversarial) generalization.

## 4.4 Post-Processing and Adaptive Attacks

We now evaluate D4 (SIZE=4)'s robustness to standard post-processing image transforms that are common, e.g., when distributed through social media. We focus on standard transforms leveraged by prior work, including additive Gaussian noise ($\sigma = 2$), blurring ($\sigma = 2$), and JPEG compression (80%). We also evaluate against the adaptive frequency-peaks attack proposed by recent work [71], designed to evade frequency-based deepfake detectors by removing frequency artifacts. At a high level, the attack manipulates frequency coefficients to remove "peaks" in the spectrum. Specifically, it computes a fingerprint as the difference between the log-scaled mean spectra of deepfake and real images. To attack a deepfake image, the attack then intensifies and subtracts this fingerprint from the image.

Tab. 5 presents results of evaluating D4 (SIZE=4) under the above settings. We observe that D4 is generally robust to post-processing, with the largest drop happening for Gaussian noise (78% AP). This is to be expected, as prior work has shown that the frequency space is generally robust to all standard transforms except noise [30]. Furthermore, the frequency-peaks attack does not appear to hurt performance. This is likely because a disjoint partitioning of features ensures that many frequency components are used for detection, and not just the peaks.

## 5  Discussion and Future Work

**Societal impact and limitations.** Modern deepfakes raise several societal and security threats; D4 is a step towards mitigating that. Nonetheless, adversarial deepfakes also have benign use cases, e.g., anonymization of an end-user on an online network; D4 could prevent such anonymization. Additionally, D4 is focused on diffusion-generated deepfake images — since it relies upon redundancy in the frequency space, it may not be effective against against future types of deepfakes that avoid these artifacts. We believe that the benefits of D4 outweigh such potential concerns.

**Alternate Partitioning Strategies.** Recall from Section 3.2 that saliency of a feature may be viewed as a heuristic measure of its "robustness", as larger saliencies imply that the model is more sensitive to perturbations of that frequency. Figure 5 plots the distribution of absolute saliency values for D4 (SIZE=1) and D4 (SIZE=4) ensembles. We observe that saliencies for saliency-partitioning configurations D4 (SIZE=4) are of relatively lower values, and are *sharply concentrated around their mode*, implying higher feature robustness. This can be attributed to the round-robin, equal distribution
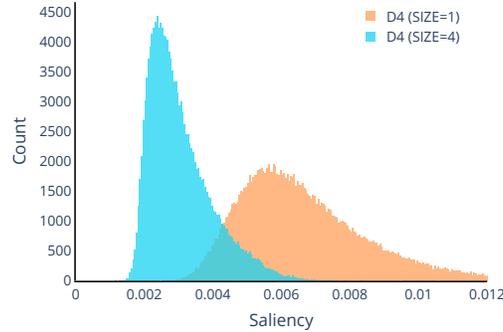
Figure 5: Distribution of frequency feature saliencies for D4 (SIZE=1) and D4 (SIZE=4) ensembles. Lower saliency implies higher feature robustness.

of robust frequencies amongst the constituent models. Improved approaches to saliency partitioning could increase this separation, improving model robustness even further. We leave this exploration to future work.

**Applicability to domains other than deepfake detection.** We presented D4 as a framework for adversarially robust deepfake detection. However, we hypothesize that this approach could apply to other classification tasks that exhibit redundancy in a feature space. While we are unaware of such a space for the popular CIFAR10 and ImageNet classification tasks, there are several classification tasks in, say, the audio domain that exhibit redundancy in features, e.g., keyword spotting and fake speech detection. Exploring this hypothesis is an interesting future research direction.

## 6 Conclusions

In this paper, we present D4, an ensemble approach to deepfake detection that exploits redundancy in frequency feature space by partitioning the frequencies across multiple models. We show theoretical advantages to such disjoint partitioning of input features, that reduces the dimensionality of the adversarial subspace. We empirically validate that D4 offers significant gains in robustness under black-box attacks, reducing attack success rates to as low as 0%.

## References

[1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 1, 3

[2] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 1

[3] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 1, 7

[4] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 1, 7

[5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 1, 3

[6] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022. 1, 7

[7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1, 3

[8] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *International Conference on Learning Representations*, 2018. 1, 7

[9] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1, 7

[10] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 1

[11] Zhendong Wang, Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Diffusion-gan: Training gans with diffusion. *arXiv preprint arXiv:2206.02262*, 2022. 1, 7

[12] Behnam Neyshabur, Srinadh Bhojanapalli, and Ayan Chakrabarti. Stabilizing gan training with multiple random projections. *arXiv preprint arXiv:1705.07831*, 2017. 1

[13] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 1

[14] Paris Martineau. Facebook removes accounts with ai-generated profile photos. `https://www.wired.com/story/facebook-removes-accounts-ai-generated-photos/`, 2019. 1

[15] Samantha Cole. Ai-assisted fake porn is here and we're all f****d. `https://www.vice.com/en_us/article/gydydm/gal-gadot-fake-ai-porn`, 2017. 1

[16] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2

[17] Nicholas Carlini and Hany Farid. Evading deepfake-image detectors with white-and black-box attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 658–659, 2020. 1, 3

[18] Shaikh Akib Shahriyar and Matthew Wright. Evaluating robustness of sequence-based deepfake detector models by adversarial perturbation. In *Proceedings of the 1st Workshop on Security Implications of Deepfakes and Cheapfakes*, pages 13–18, 2022. 1, 3

[19] Steven Lawrence Fernandes and Sumit Kumar Jha. Adversarial attack on deepfake detection using rl based texture patches. In *European Conference on Computer Vision*, pages 220–235. Springer, 2020. 1, 3

[20] Quanyu Liao, Yuezun Li, Xin Wang, Bin Kong, Bin Zhu, Siwei Lyu, Youbing Yin, Qi Song, and Xi Wu. Imperceptible adversarial examples for fake image detection. *arXiv preprint arXiv:2106.01615*, 2021. 1, 3

[21] Shehzeen Hussain, Todd Huster, Chris Mesterharm, Paarth Neekhara, Kevin An, Malhar Jere, Harshvardhan Sikka, and Farinaz Koushanfar. Reface: Real-time adversarial attacks on face recognition systems. *arXiv preprint arXiv:2206.04783*, 2022. 1, 3

[22] Ngan Hoang Vo, Khoa D Phan, Anh-Duy Tran, and Duc-Tien Dang-Nguyen. Adversarial attacks on deepfake detectors: A practical analysis. In *International Conference on Multimedia Modeling*, pages 318–330. Springer, 2022. 1, 3

[23] Apurva Gandhi and Shomik Jain. Adversarial perturbations fool deepfake detectors. In *2020 International joint conference on neural networks (IJCNN)*, pages 1–8. IEEE, 2020. 1, 3, 4

[24] Shehzeen Hussain, Paarth Neekhara, Malhar Jere, Farinaz Koushanfar, and Julian McAuley. Adversarial deepfakes: Evaluating vulnerability of deepfake detectors to adversarial examples. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3348–3357, 2021. 1, 3

[25] Shehzeen Hussain, Paarth Neekhara, Brian Dolhansky, Joanna Bitton, Cristian Canton Ferrer, Julian McAuley, and Farinaz Koushanfar. Exposing vulnerabilities of deepfake detection systems with robust attacks. *Digital Threats: Research and Practice (DTRAP)*, 3(3):1–23, 2022. 1, 3

[26] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *Proceedings of International Conference on Machine Learning*, 2018. 1

[27] Florian Tramèr, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. The space of transferable adversarial examples. *arXiv preprint arXiv:1704.03453*, 2017. 2, 5, 15, 16

[28] Himanshu Dutta, Aditya Pandey, and Saurabh Bilgaiyan. Ensembledet: ensembling against adversarial attack on deepfake detection. *Journal of Electronic Imaging*, 30(6):063030, 2021. 2, 4, 7

[29] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8695–8704, 2020. 2, 3, 8, 9

[30] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. Leveraging frequency analysis for deep fake image recognition. In *International Conference on Machine Learning*, pages 3247–3258. PMLR, 2020. 2, 3, 5, 10

[31] Zhaoqing Pan, Weijie Yu, Xiaokai Yi, Asifullah Khan, Feng Yuan, and Yuhui Zheng. Recent progress on generative adversarial networks (gans): A survey. *IEEE Access*, 7:36322–36333, 2019. 3

[32] Zeyang Sha, Zheng Li, Ning Yu, and Yang Zhang. De-fake: Detection and attribution of fake images generated by text-to-image diffusion models. *arXiv preprint arXiv:2210.06998*, 2022. 3

[33] Francesco Marra, Cristiano Saltori, Giulia Boato, and Luisa Verdoliva. Incremental learning for the detection and classification of gan-generated images. In *2019 IEEE international workshop on information forensics and security (WIFS)*, pages 1–6. IEEE, 2019. 3

[34] Francesco Marra, Diego Gragnaniello, Davide Cozzolino, and Luisa Verdoliva. Detection of gan-generated fake images over social networks. In *2018 IEEE conference on multimedia information processing and retrieval (MIPR)*, pages 384–389. IEEE, 2018. 3

[35] Lakshmanan Nataraj, Tajuddin Manhar Mohammed, BS Manjunath, Shivkumar Chandrasekaran, Arjuna Flenner, Jawadul H Bappy, and Amit K Roy-Chowdhury. Detecting gan generated fake images using co-occurrence matrices. *Electronic Imaging*, 2019(5):532–1, 2019. 3

[36] Scott McCloskey and Michael Albright. Detecting gan-generated imagery using color cues. *arXiv preprint arXiv:1812.08247*, 2018. 3

[37] Luca Guarnera, Oliver Giudice, and Sebastiano Battiato. Deepfake detection by analyzing convolutional traces. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 666–667, 2020. 3

[38] Zhengzhe Liu, Xiaojuan Qi, and Philip HS Torr. Global texture enhancement for fake face detection in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8060–8069, 2020. 3

[39] Lucy Chai, David Bau, Ser-Nam Lim, and Phillip Isola. What makes fake images detectable? understanding properties that generalize. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVI 16*, pages 103–120. Springer, 2020. 3

[40] Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. Towards universal fake image detectors that generalize across generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24480–24489, 2023. 3, 8, 9

[41] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. Detecting and simulating artifacts in gan fake images. In *2019 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2019. 3

[42] Ricard Durall, Margret Keuper, Franz-Josef Pfreundt, and Janis Keuper. Unmasking deepfakes with simple features. *arXiv preprint arXiv:1911.00686*, 2019. 3

[43] Jonas Ricker, Simon Damm, Thorsten Holz, and Asja Fischer. Towards the detection of diffusion model deepfakes. *arXiv preprint arXiv:2210.14571*, 2022. 3, 7, 8, 10

[44] Paarth Neekhara, Brian Dolhansky, Joanna Bitton, and Cristian Canton Ferrer. Adversarial threats to deepfake detection: A practical perspective. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 923–932, 2021. 3

[45] Intel. Intel introduces real-time Deepfake Detector. https://www.intel.com/content/www/us/en/newsroom/news/intel-introduces-real-time-deepfake-detector.html, 2023. 3

[46] Deepware: Scan & detect deepfake videos. https://deepware.ai/, 2023. 3

[47] Reality Defender. https://realitydefender.com/, 2023. 3

[48] Sensity AI: Speed up online customer identity verification guaranteeing high anti-fraud standards. https://sensity.ai/, 2023. 3

[49] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. 3, 8

[50] Warren He, James Wei, Xinyun Chen, Nicholas Carlini, and Dawn Song. Adversarial example defense: Ensembles of weak defenses are not strong. In *11th {USENIX} workshop on offensive technologies ({WOOT} 17)*, 2017. 3

[51] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. *Advances in neural information processing systems*, 32, 2019. 4

[52] Sohail Ahmed Khan, Alessandro Artusi, and Hang Dai. Adversarially robust deepfake media detection using fused convolutional neural network predictions. *arXiv preprint arXiv:2102.05950*, 2021. 4, 7

[53] Aditya Devasthale and Shamik Sural. Adversarially robust deepfake video detection. In *2022 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 396–403. IEEE, 2022. 4, 7

[54] Jianguo Jiang, Boquan Li, Shuai Yu, Chao Liu, Shaohua An, Mingqi Liu, and Min Yu. A residual fingerprint-based defense against adversarial deepfakes. In *2021 IEEE 23rd Int Conf on High Performance Computing & Communications; 7th Int Conf on Data Science & Systems; 19th Int Conf on Smart City; 7th Int Conf on Dependability in Sensor, Cloud & Big Data Systems & Application (HPCC/DSS/SmartCity/DependSys)*, pages 797–804. IEEE, 2021. 4

[55] Huiying Li, Shawn Shan, Emily Wenger, Jiayun Zhang, Haitao Zheng, and Ben Y Zhao. Blacklight: Defending black-box adversarial attacks on deep neural networks. *arXiv preprint arXiv:2006.14042*, 2020. 4

[56] Ryan Feng, Ashish Hooda, Neal Mangaokar, Kassem Fawaz, Somesh Jha, and Atul Prakash. Investigating stateful defenses against black-box adversarial examples, 2023. 4

[57] Zeyu Qin, Yanbo Fan, Hongyuan Zha, and Baoyuan Wu. Random noise defense against query-based black-box attacks. *Advances in Neural Information Processing Systems*, 34:7650–7663, 2021. 4, 8

[58] Sizhe Chen, Zhehao Huang, Qinghua Tao, Yingwen Wu, Cihang Xie, and Xiaolin Huang. Adversarial attack on attackers: Post-process to mitigate black-box score-based query attacks. *arXiv preprint arXiv:2205.12134*, 2022. 4

[59] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017. 5

[60] Ryan Feng, Neal Mangaokar, Jiefeng Chen, Earlence Fernandes, Somesh Jha, and Atul Prakash. Graphite: Generating automatic physical examples for machine-learning attacks on computer vision systems. *arXiv preprint arXiv:2002.07088*, 2022. 5

[61] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 7

[62] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019. 7

[63] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I. Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, 2019. 8

[64] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 8

[65] Jianbo Chen, Michael I Jordan, and Martin J Wainwright. Hopskipjumpattack: A query-efficient decision-based attack. In *2020 ieee symposium on security and privacy (sp)*, pages 1277–1294. IEEE, 2020. 8

[66] Huichen Li, Xiaojun Xu, Xiaolu Zhang, Shuang Yang, and Bo Li. Qeba: Query-efficient boundary-based blackbox attack. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1221–1230, 2020. 8

[67] Thibault Maho, Teddy Furon, and Erwan Le Merrer. Surfree: a fast surrogate-free black-box attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10430–10439, 2021. 8

[68] Xiaosen Wang, Zeliang Zhang, Kangheng Tong, Dihong Gong, Kun He, Zhifeng Li, and Wei Liu. Triangle attack: A query-efficient decision-based adversarial attack. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part V*, pages 156–174. Springer, 2022. 8

[69] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *arXiv preprint arXiv:1712.04248*, 2017. 8

[70] Minhao Cheng, Simranjit Singh, Patrick Chen, Pin-Yu Chen, Sijia Liu, and Cho-Jui Hsieh. Sign-opt: A query-efficient hard-label adversarial attack. *arXiv preprint arXiv:1909.10773*, 2019. 8

[71] Vera Wesselkamp, Konrad Rieck, Daniel Arp, and Erwin Quiring. Misleading deep-fake detection with gan fingerprints. In *2022 IEEE Security and Privacy Workshops (SPW)*, pages 59–65. IEEE, 2022. 10

[72] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defenses. *arXiv preprint arXiv:1705.07204*, 2017. 16

[73] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International conference on machine learning*, pages 2206–2216. PMLR, 2020. 17

[74] Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In *International Conference on Machine Learning*, pages 2196–2205. PMLR, 2020. 17

[75] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square Attack: a query-efficient black-box adversarial attack via random search . In *Proceedings of European Conference on Computer Vision*, 2020. 17

14

# 7 Appendix

In Appendix 7.1, we provide proofs for our bounds on dimensionality of the adversarial subspace of disjoint ensembles, such as D4. In Appendix 7.2, we provide details of our baseline configurations and attack settings for reproducibility purposes, as well as additional experimental results on other datasets.

## 7.1 Proofs

### 7.1.1 Orthogonal Gradients

Prior to providing proofs for the adversarial dimensions, we demonstrate that gradients for disjoint classifiers are always orthogonal to each other. We will use this for our later results. Given input space $\mathcal{X}$ and class-label space $\mathcal{Y}$, we have $n$ disjoint classifiers $\mathcal{F}_1, ..., \mathcal{F}_n$. If $T$ is the DCT transformation matrix, we can define $T_i$ to be the transformation matrix for the classifier $\mathcal{F}_i$. Each disjoint transformation $T_i$ has a lot of zeros. Only the rows corresponding to the unmasked frequencies of classifier $\mathcal{F}_i$ have non-zero entries. Moreover, since no frequency is shared by any two classifiers, the $j^{th}$ row will have non-zero entries in exactly one of the $n$ disjoint transformation matrices, i.e. $T_i T_j^{\mathsf{T}} = O \ \forall i \neq j$.

Next, the $n$ disjoint classifiers $\mathcal{F}_1, ..., \mathcal{F}_n$, where $\mathcal{F}_i : T_i \mathbf{x} \to y$, are trained using loss functions $\mathcal{L}_{\mathcal{F}_1}, ..., \mathcal{L}_{\mathcal{F}_n}$ respectively. Now, the dot product between the gradients of classifiers $\mathcal{F}_i$ and $\mathcal{F}_j$ is given by

$$
\begin{aligned}
(\nabla_{\mathbf{x}} \mathcal{L}_{\mathcal{F}_i})^{\mathsf{T}} \left(\nabla_{\mathbf{x}} \mathcal{L}_{\mathcal{F}_j}\right) &= (T_i^{\mathsf{T}} \nabla_{T_i \mathbf{x}} \mathcal{L}_{\mathcal{F}_i})^{\mathsf{T}} \left(T_j^{\mathsf{T}} \nabla_{T_j \mathbf{x}} \mathcal{L}_{\mathcal{F}_j}\right) \\
&= (\nabla_{T_i \mathbf{x}} \mathcal{L}_{\mathcal{F}_i})^{\mathsf{T}} T_i T_j^{\mathsf{T}} \left(\nabla_{T_j \mathbf{x}} \mathcal{L}_{\mathcal{F}_j}\right) \\
&= 0
\end{aligned}
\tag{7}
$$

### 7.1.2 Proof of Lemma 3.1

From [27], we know that for a classifier $\mathcal{F} : \mathcal{X} \to \mathcal{Y}$ where $\mathcal{X} \in \mathbb{R}^d$ is the input space and $\mathcal{Y}^c \in \mathbb{Z}$ is the finite class label space, the dimension of the adversarial subspace around input-label pair $(\mathbf{x}, y)$ where $\mathbf{x} \in \mathcal{X}$ and $y \in \mathcal{Y}$, is approximated by the maximal number of orthogonal perturbations $\mathbf{r_1}, \mathbf{r_2}, ..., \mathbf{r_k}$ such that $||\mathbf{r_i}||_2 \leq \epsilon$ and $\mathbf{g}^{\mathsf{T}} \mathbf{r_i} \geq \gamma \ \forall \ 1 \leq i \leq k$. Here, $\mathbf{g} = \nabla_{\mathbf{x}} L(\mathcal{F}(\mathbf{x}), y)$ and $\gamma$ is the increase in loss function $L$ sufficient to cause a mis-classification. [27] provide a tight bound for $k$:

$$
k = \min \left( d, \left\lfloor \frac{\epsilon^2 ||\mathbf{g}||_2^2}{\gamma^2} \right\rfloor \right)
\tag{8}
$$

We now extend this result for $n$ disjoint classifiers. Let $\mathbf{g}' = \frac{\sum_{j=1}^{n} \mathbf{g_j}}{n}$.

Now, for *at-least-one voting*,

$$
\mathbf{g}'^{\mathsf{T}} \mathbf{r_i} = \frac{\sum_{j=1}^{n} \mathbf{g_j}^{\mathsf{T}} \mathbf{r_i}}{n} \geq \frac{\sum_{j=1}^{n} \gamma_j}{n} \quad \forall \ 1 \leq i \leq k
\tag{9}
$$

Applying the result from [27] (Equation 8) on the above inequality (Equation 9), we get:

$$k = \min\left(d, \left\lfloor \left| \frac{\epsilon^2 n^2 ||\mathbf{g}'||_2^2}{\left(\sum\limits_{j=1}^{n} \gamma_j\right)^2} \right| \right\rfloor \right)$$

$$= \min\left(d, \left\lfloor \left| \frac{\epsilon^2 \sum\limits_{j=1}^{n} ||\mathbf{g_j}||_2^2}{\left(\sum\limits_{j=1}^{n} \gamma_j\right)^2} \right| \right\rfloor \right).$$

(10)

(since $\mathbf{g_i}^\mathsf{T}\mathbf{g_j} = 0 \ \forall i \neq j$, using Equation 7)

Now, for *majority voting*, we again apply the results from [27] (Equation 8). However, the derivation now depends on the selection of $\left\lceil \frac{n}{2} \right\rceil$ models that the adversary chooses to target. To obtain the lower and upper bounds, we can select $\left\lceil \frac{n}{2} \right\rceil$ with the most and least adversarial dimensions respectively. Following a similar derivation as before, we get :

$$k \geq \min\left(d, \left\lfloor \left| \min_{|K|=\lceil \frac{n}{2} \rceil} \frac{\epsilon^2 \sum\limits_{j=1}^{n} ||\mathbf{g_j}||_2^2}{\left(\sum\limits_{j=1}^{n} \gamma_j\right)^2} \right| \right\rfloor \right)$$

(11)

$$k \leq \min\left(d, \left\lfloor \left| \max_{|K|=\lceil \frac{n}{2} \rceil} \frac{\epsilon^2 \sum\limits_{j=1}^{n} ||\mathbf{g_j}||_2^2}{\left(\sum\limits_{j=1}^{n} \gamma_j\right)^2} \right| \right\rfloor \right)$$

(12)

### 7.1.3 Proof of Lemma 3.2

Follow up work from [72] also provides a tight bound for the adversarial dimension in the $\ell_\infty$ case. They provide a tight bound for the number of $k$ orthogonal perturbations $\mathbf{r_1}, ..., \mathbf{r_k} \in \mathbb{R}^d$ such that $||\mathbf{r_i}||_\infty \leq \epsilon$, given by $sign(\mathbf{g})^\mathsf{T}\mathbf{r_i} = \frac{\epsilon d}{\sqrt{k}} \ \forall 1 \leq i \leq k$ where $sign(\mathbf{g})$ is the signed gradient.

We now extend this result for $n$ disjoint classifiers. For $\mathbf{g}' = \frac{\sum\limits_{j=1}^{n} \mathbf{g_j}}{n}$, since $\mathbf{g_j}'s$ are non-zero only on non-overlapping dimensions, we can see that $sign(\mathbf{g}')^\mathsf{T}r = \sum\limits_{j=1}^{n} sign(\mathbf{g_j})^\mathsf{T}r \ \forall \mathbf{r} \in \mathbb{R}^d$. Applying the above results here, we get

$$\sum_{j=1}^{n} sign(\mathbf{g_j})^\mathsf{T}\mathbf{r_i} = \frac{\epsilon d}{\sqrt{k}} \ \forall 1 \leq i \leq k$$

(13)

Now, similar to [72], we compute the perturbation magnitude along a random permutation of the signed gradient. For each $1 \leq j \leq n$ and $1 \leq i \leq k$, we get :

16

$$\mathbb{E}[\mathbf{g_j}^\top \mathbf{r_i}] = \mathbb{E}\left[\sum_{p=1}^{d} |g_j^{(p)}| \cdot sign(g_j^{(p)}) \cdot r_i^{(p)}\right]$$

$$= \sum_{p=1}^{d} |g_j^{(p)}| \mathbb{E}\left[sign(g_j^{(p)}) \cdot r_i^{(p)}\right] \qquad (14)$$

$$= \frac{\epsilon ||\mathbf{g_j}||_1}{n\sqrt{k}}$$

## 7.2 Additional Evaluation Details

### 7.2.1 Details for Attacks

We consider the following attacks for our evaluation. These attacks constitute the entire ensemble of attacks used in the AutoAttack Benchmark. We tune the attacks where necessary to get the strongest attack setting.

**APGD-CE/CW [73]** is a step-size free variant of the standard PGD attack. It adjusts the step size during the attack based on the convergence of the loss and the overall perturbation budget. We optimize the adaptive PGD-attack on the Cross Entropy (CE) and Carlini Wagner (CW) loss functions. We use the same set of parameters for the attack as mentioned in AutoAttack Benchmark other than the step size decay parameter $\alpha$ which we set to $0.1$ instead of $2$.

**FAB [74]** is a iterative first order attack that utilizes geometry of the decision boundary to minimize the perturbation required to cause mis-classification. We use the same set of parameters as AutoAttack.

**Square [75]** is an efficient black-box attack that uses random square shaped updates to approximate the decision boundary. We use the same set of parameters as AutoAttack.

### 7.2.2 Evaluation on CIFAR10

Table 6 presents results for attacking the D3-S(4) ensemble and AT baseline, for the CIFAR10 classification task. While D3-S(4) offers some improvements at the smallest perturbation budgets, it quickly drops off. We suspect that a more carefully chosen feature space with redundancies for animal/vehicle classification would improve these results.

Table 6: CIFAR10 adversarial accuracies for $\ell_\infty$ and $\ell_2$ perturbation attacks for same configurations as in Table **??**.

| Attacks | $\epsilon$ | D4 (SIZE=1) | D4 (SIZE=4) |
|---------|------------|-------------|-------------|
| | $\ell_\infty$ | | |
| APGD-CE(50) | 1/255 | 1.7 | **59.8** |
| | 4/255 | 0.1 | **1.2** |
| | 8/255 | **0.1** | 0 |
| | 16/255 | **0.1** | 0 |
| | $\ell_2$ | | |
| APGD-CE(50) | 0.5 | 15.1 | **63.5** |
| | 1 | 13.6 | **39.7** |
| | 5 | 0.9 | **8.2** |
| | 10 | 0.8 | **6.2** |