NATURALADV: An Exploratory Framework to Balance Adversarial Strength and Stealth in Autonomous Driving Environments

Anonymous submission

Abstract

Deep Neural Networks (DNNs) have become integral to various real-world autonomous mobile systems, from self-driving cars to food delivery robots. However, current adversarial attack techniques often focus on maximizing the attack strength at the cost of naturalness, leading to examples that are easily detected by humans or deviate significantly from the expected input distribution. This trade-off between adversarial effectiveness and natural appearance presents a critical challenge in ensuring the robustness and reliability of DNNs in practical settings. In this paper, we introduce a framework to navigate this trade-off. Unlike traditional methods that prioritize pixel-level perturbations, our approach integrates a naturalness metric that reflects human perceptibility and the resemblance of adversarial examples to real-world inputs. The framework leverages pretrained neural networks, differentiable similarity metrics, and high-strength adversarial attacks to automatically generate adversarial images that strike a balance between these two competing objectives. Our method leverages differentiable image similarity metrics and custom loss functions for gradient-based attack generation. Initial empirical results demonstrate the framework's potential to create adversarial examples that are both powerful and natural-looking, capable of bypassing DNN defenses while maintaining realism. This work aims to offer software engineers a flexible approach to adversarial attack generation, with implications for robustness testing and model evaluation in various real-world contexts. This approach enables higher abstraction of robustness testing above the pixel level, as well as future development of adversarial techniques that consider not only attack strength but also the naturalness of the generated tests, paving the way for more resilient AI systems.

This paper presents the Natural Adversarial DNN Validation (NATURALADV) framework for balancing the trade-off between adversarial strength and naturalness of the adversarial patch's appearance. NATURALADV can incorporate a number of differentiable naturalness metrics, works with various gradient traversal algorithms, and scales to attacks represented in multiple sensor readings. Our contributions are:

- a technique, NATURALADV, to balance the trade-off between adversarial strength and naturalness for in-situ adversarial patch attacks;
- a proof of concept study showing the naturalness-strength tradeoff for the motivating example; and
- an open-source repository with tool and data for reproducibility available at https://github.com/anon/anon.

When designing adversarial attacks for deployment scenarios, it is essential to distinguish between stealthiness and naturalness. Stealthiness refers to the perceptual imperceptibility of the perturbation; a stealthy attack introduces minimal visual artifacts, making it hard for a human observer to detect any manipulation. Naturalness, on the other hand, refers to how well the adversarial example aligns with the expected distribution of inputs—whether it "looks real" or conforms to what the model would typically encounter. An attack might be stealthy but lack naturalness if, for instance, it is an abstract or unrealistic pattern that the model might never see in a real-world setting.

The overall goal of our approach is to generate adversarial patches that take advantage of existing high-strength adversarial patches and inject a configurable measure of stealthiness in a way that best preserves the adversarial strength and properties of the original high-strength patch. We formulate the strength-stealthiness problem as a trade-off in accordance with existing literature (). Our framework can generate more natural patches using existing adversarial perturbations for any definition of naturalness using differentiable image similarity metrics.



Figure 1: Overview of NATURALADV generation loop

Figure 1 depicts a high-level overview of the generation loop for the NATURALADV framework. It takes in two images of the patch region, one with the original adversarial perturbation known to have high adversarial strength and one that the user considers natural and one that is when applied to the patch region,, an image set imgs taken from an ADS navigating a driving environment without an adversarial patch, a navigation DNN, iterations of FGSM iters, an image similarity metric similarity, and weights for the two loss terms w1 for the image similarity loss and w2 for the perturbed prediction loss between the original perturbation and the target image patch.

The framework alternates between calculating the loss function (see Equation 1) and backpropagation combined with Fast Sign Gradient Method (FGSM) to adjust the similarity of the patch to match the target image, while still retaining the adversarial strength of the original high-strength patch. Similarity and strength are prioritized according to parameterized weights. After iters loops, the generation loop exits and returns the final patch for injection into a driving environment.

Algorithm 1 gives a more granular explanation of this generation loop. The algorithm takes as inputs a set of images *imqs*, a set number of iterations of gradient ascent *iters*, a deep neural network DNN, weights for the two loss terms w_1 and w₂, a high-strength adversarial patch orig_patch, and a natural patch to mimic, *natural_patch*. The goal is to iteratively refine the patch so that it maintains strong adversarial properties while appearing natural. First, the initial adversarial patch is overlaid onto the images, creating $imgs_{orig_patch}$ (line 1). These patched images are then passed through the DNN to generate the baseline predictions, $ys_{orig-patch}$ (line 2). This establishes a reference point for the adversarial behavior of the original patch.

Next, the optimization process starts. Over a set number of iterations, a "natural" patch - one that aims to maintain visual similarity to the original patch-is applied to the images. The natural patch is overlaid onto the images and these images are fed through the DNN to produce predictions, $ys_{natural.patch}$ (lines 4-5). Two types of loss are computed at each iteration: a similarity loss, which measures the visual similarity between the natural patch and the original patch using a differentiable similarity metric (line 6), and a prediction loss, which assesses how closely the predictions of the natural patch match those of the original patch using Mean Squared Error (MSE) (line 7). These losses are weighted by factors w_1 and w_2 and combined into a total loss (line 8).

To optimize the natural patch, Fast Gradient Sign Method (FGSM) is applied based on the computed graident using the loss, adjusting the patch to better balance adversarial strength and natural appearance (line 9). This iterative process continues for a predefined number of iterations. Once completed, the refined patch - now the final patch - is output, representing an adversarial example that maintains its effectiveness while appearing more natural.

NATURALADV Loss Function

The red-boxed loss function in Figure 1 is a shorthand version of the full loss function:

$$loss = w1 \times similarity(img_{natural}, img_{originial_patch}) + w2 \times L1(DNN(imgs + img_{originial_patch}), DNN(imgs + img_{natural_patch}))$$
(1)

where *similarity* is any differentiable image similarity metric (e.g. SSIM, German-McClure, Welsch, etc.), imgnatural is the natural image we want the adversarial patch to

Algorithm 1: NATURALADV Perturbation Generation

Input DNN, imgs, iters, orig_patch, natural_patch, weights={w1,w2}, similarity

Output final_patch

- 3: for i = 0 to iters do
- 4: $\texttt{imgs}_{natural_patch} \gets \texttt{imgs} + \texttt{natural_patch}$
- $ys_{natural_patch} = \text{DNN}(\text{imgs}_{natural_patch})$ 5:
- sim_loss = similarity(orig_patch, natural_patch) 6:
- $prediction_loss = MSE(ys_{orig_patch}, ys_{natural_patch})$ 7:
- 8: $loss = w1 * sim_loss + w2 * prediction_loss$
- 9: $natural_patch = FGSM(natural_patch)$
- 10: end for
- 11: $final_patch = natural_patch$
- 12: **return** final_patch

look like, *img*originial_peturbation is the high-strength but unnatural-looking (or un-stealthy) original adversarial patch, $DNN(imgs + img_{originial_peturbation}$ is the DNN prediction output for the original high strength perturbation, $DNN(imgs + img_{natural_peturbation})$ is the DNN prediction output for the current version of the natural perturbation, and w_x are weights to prioritize image similarity loss or DNN prediction loss. For an example of imgs +*imgoriginal_peturbation* and *imgs+imgnatural_peturbation*, see Figure 1. For an example of *img*_{originial_peturbation} and imgnatural_peturbation see the rightmost and middle images in Figure 2, respectively.

Note that the sim_loss loss term and the prediction loss term may need to be normalized, as not all image similarity metrics have a [-1, 1] range and can be overly permissive which may overpower your prediction loss term. Include the differences in outcome when balancing loss function terms: 100% image similarity and 0% DNN output to 0% image similarity and 100% DNN output



Figure 2: Natural image, "natural" adversarial patch, and original high-strength adversarial patch. Middle image was balanced using SSIM image similarity metric and equally weighted image similarity metric loss and prediction loss.

Acknowledgment

Anonymized.

References

NATURALADV: An Exploratory Framework to Balance Adversarial Strength and Stealth in Autonomous Driving Environments

Anonymized

Motivation

Deep Neural Networks (DNNs) have become integral to various real-world autonomous mobile systems, from self-driving cars to food delivery robots.



Figure 1. A DeepBillboard [3] in-situ patch attack.

However, current adversarial attack techniques often focus on maximizing the attack strength at the cost of naturalness, leading to examples that are easily detected by humans or deviate significantly from the expected input distribution. This trade-off between adversarial effectiveness and natural appearance presents a critical challenge in ensuring the robustness and reliability of DNNs in practical settings.

Perturbation Stealthiness and Naturalness

When designing adversarial attacks for deployment scenarios, it is essential to distinguish between stealthiness and naturalness. Stealthiness refers to the perceptual imperceptibility of the perturbation; a stealthy attack introduces minimal visual artifacts, making it hard for a human observer to detect any manipulation. Naturalness, on the other hand, refers to how well the adversarial example aligns with the expected distribution of inputs—whether it "looks real" or conforms to what the model would typically encounter. An attack might be stealthy but lack naturalness if, for instance, it is an abstract or unrealistic pattern that might never occur in a real-world setting.

Experiments

NATURALADV Framework Contributions

This poster presents the Natural Adversarial DNN Validation (NaturalADV) framework for balancing the trade-off between adversarial strength and naturalness of the adversarial patch's appearance. NaturalADV can incorporate a number of differentiable naturalness metrics, works with various gradient traversal algorithms, and scales to attacks represented in multiple sensor readings. Our contributions are:

- a framework, NaturalADV, to balance the trade-off between adversarial strength and naturalness for in-situ adversarial patch attacks;
- a proof of concept study showing the naturalness-strength tradeoff for the motivating example; and
- an open-source repository with tool and data for reproducibility available at https://github.com/anon/repo.

Perturbation Generation Loop

Figure 3 depicts a high-level overview of the generation loop for the NaturalADV framework. It takes in two images of the patch region of the deployment environment, one with the original adversarial perturbation known to have high adversarial strength (the original patch) and one that the user considers natural (the target patch), an image set **imgs** taken from an ADS navigating a driving environment without an adversarial patch, a navigation **DNN**, iterations of FGSM **iters**, a differentiable image similarity metric **similarity**, and **weights** for the two loss terms **w1** for the image similarity loss and **w2** for the perturbed prediction loss between the original perturbation and the target patch.



Perturbation Strength

We explore a range of weights for similarity and prediction (see Equation 1) and report several performance metrics. Column 2 shows the resulting structural similarity index measure (SSIM) score when comparing the benign target patch and the generated adversarial patch. Column 3 shows the crash rate when the patch is deployed in simulation like in Figure 1. Column 4 shows the average deviation in the vehicle's trajectory from the centerline of the road.

Weights (sim, pred)	SSIM Score	Crash Rate	Avg. traj. deviation
0.00, 1.00	0.24	53%	2.75m
0.10, 0.90	0.21	22%	2.45m
0.25, 0.75	0.33	18%	2.37m
0.50, 0.50	0.46	2%	2.31m
0.75, 0.25	0.51	0%	2.21m
0.90, 0.10	0.70	0%	2.24m
1.00, 0.00	1.00	0%	2.23m

Table 1. Performance metrics for a range of weights using NaturalADV.

As Table 1 shows, perturbation strength diminishes inversely to SSIM score, where the generated patch resembles more and more closely the benign target patch. However, the generated patch still retains the ability to crash the vehicle in deployment when image similarity loss and prediction loss are equally weighted.

Perturbation Naturalness



Figure 2. Patch appearances from original perturbation to benign target patch

Perturbation naturalness is dictated by the choice of metric, in this study SSIM. SSIM is designed to compare two images in terms of luminance, contrast, and structure, or edge detection. It was originally designed for black and white images and as a result preserves colors of the original perturbation well into the (0.90, 0.10) weighting.

The framework alternates between calculating the loss function (see Equation 1) and backpropagation combined with Fast Sign Gradient Method (FGSM) to adjust the similarity of the patch to match the target image, while still retaining the adversarial strength of the original high-strength patch. Similarity and strength are prioritized according to parameterized **weights**. After **iters** loops, the generation loop exits and returns the final patch for injection into a driving environment.

NATURALADV Loss Function

The perturbation loop relies on a loss function for which preserving image similarity versus perturbation strength has been parameterized:

 $loss = w1 \times similarity(img_{natural}, img_{originial_patch}) + w2 \times L1(DNN(imgs + img_{originial_patch}), DNN(imgs + img_{natural_patch}))$ (1)

where *similarity* is any differentiable image similarity metric (e.g. SSIM, German-McClure, Welsch, etc.), $img_{natural}$ is the natural image is that supposed to be an image or a patch? WE need to be super careful here differentiating patch and the image whole image, seems like we are mixing them up? we want the adversarial patch to look like, $img_{originial_peturbation}$ is the high-strength but unnatural-looking (or un-stealthy) original adversarial patch, $DNN(imgs + img_{originial_peturbation}$ is the DNN prediction output for the original high strength perturbation, $DNN(imgs + img_{natural_peturbation}$ is the DNN prediction output for the current version of the natural perturbation, and w_x are weights to prioritize image similarity loss or DNN prediction loss.

References

- [1] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples.
 In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 15262–15271, 2021.
- [2] Meriel von Stein, David Shriver, and Sebastian Elbaum. Deepmaneuver: Adversarial test generation for trajectory manipulation of autonomous vehicles. IEEE Transactions on Software Engineering, 2023.
- [3] Husheng Zhou, Wei Li, Zelun Kong, Junfeng Guo, Yuqun Zhang, Bei Yu, Lingming Zhang, and Cong Liu.
 Deepbillboard: Systematic physical-world testing of autonomous driving systems.
 In 2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE), pages 347–358, 2020.