

# On the (In-)Security of Machine Learning

*Nicholas Carlini*  
*Google Brain*



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

WHEN A USER TAKES A PHOTO,  
THE APP SHOULD CHECK WHETHER  
THEY'RE IN A NATIONAL PARK...



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

WHEN A USER TAKES A PHOTO,  
THE APP SHOULD CHECK WHETHER  
THEY'RE IN A NATIONAL PARK...

SURE, EASY GIS LOOKUP.  
GIMME A FEW HOURS.



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

WHEN A USER TAKES A PHOTO,  
THE APP SHOULD CHECK WHETHER  
THEY'RE IN A NATIONAL PARK...

SURE, EASY GIS LOOKUP.  
GIMME A FEW HOURS.

... AND CHECK WHETHER  
THE PHOTO IS OF A BIRD.



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

WHEN A USER TAKES A PHOTO,  
THE APP SHOULD CHECK WHETHER  
THEY'RE IN A NATIONAL PARK...

SURE, EASY GIS LOOKUP.  
GIMME A FEW HOURS.

... AND CHECK WHETHER  
THE PHOTO IS OF A BIRD.

I'LL NEED A RESEARCH  
TEAM AND FIVE YEARS.



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

WHEN A USER TAKES A PHOTO,  
THE APP SHOULD CHECK WHETHER  
THEY'RE IN A NATIONAL PARK...

SURE, EASY GIS LOOKUP.  
GIMME A FEW HOURS.

... AND CHECK WHETHER  
THE PHOTO IS OF A BIRD.

I'LL NEED A RESEARCH  
TEAM AND FIVE YEARS.



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

**Written: Sept 24, 2014**

WHEN A USER TAKES A PHOTO,  
THE APP SHOULD CHECK WHETHER  
THEY'RE IN A NATIONAL PARK...

SURE, EASY GIS LOOKUP.  
GIMME A FEW HOURS.

... AND CHECK WHETHER  
THE PHOTO IS OF A BIRD.

I'LL NEED A RESEARCH  
TEAM AND FIVE YEARS.



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

**Written: Sept 24, 2014**

**Today: Oct 16, 2018**



WHEN A USER TAKES A PHOTO,  
THE APP SHOULD CHECK WHETHER  
THEY'RE IN A NATIONAL PARK...

SURE, EASY GIS LOOKUP.  
GIMME A FEW HOURS.

... AND CHECK WHETHER  
THE PHOTO IS OF A BIRD.

I'LL NEED A RESEARCH  
TEAM AND FIVE YEARS.



IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

**Written: Sept 24, 2014**

**Today: Oct 16, 2018**

**... 4 years ago**

**So how are  
we doing?**



95% it is a

**French  
Bulldog**



83% it is a

**Old  
English  
Sheepdog**



78% it is a

**Greater  
Swiss  
Mountain  
Dog**



67% it is a

**Great  
Dane**



99.99% it is

**Guacamole**



96% it is a

**Golden  
Retriever**





99.99% it is

**Guacamole**

This phenomenon is known as an  
**adversarial example**

B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Srndic, P. Laskov, G. Giacinto, and F. Roli. Evasion attacks against machine learning at test time. 2013.  
C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. ICLR 2014.  
I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. 2014.



88% **tabby cat**



adversarial  
perturbation



88% **tabby cat**



adversarial  
perturbation



88% **tabby cat**



adversarial  
perturbation



88% **tabby cat**

99% **guacamole**

Why should we care about  
adversarial examples?

*Make ML*  
***robust***



(a)



(b)



(c)



Why should we care about  
adversarial examples?

*Make ML*  
***robust***

*Make ML*  
***better***

How do we generate  
adversarial examples?

DEFN: The **loss** of a neural network on an input  **$\mathbf{x}$**  for a label  **$\mathbf{y}$**  is a measure of how *wrong* the network is on  **$\mathbf{x}$** .

loss(



, dog) is small

loss(



, guacamole) is large

**MAXIMIZE**

neural network loss  
on the given input

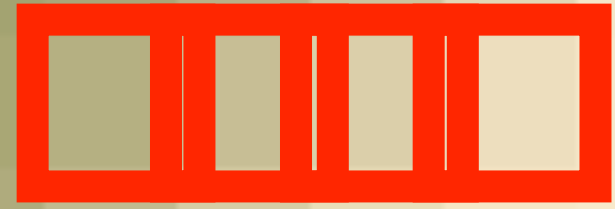
**SUCH THAT**

the perturbation is less  
than a given threshold

What do we need to know?

Everything.





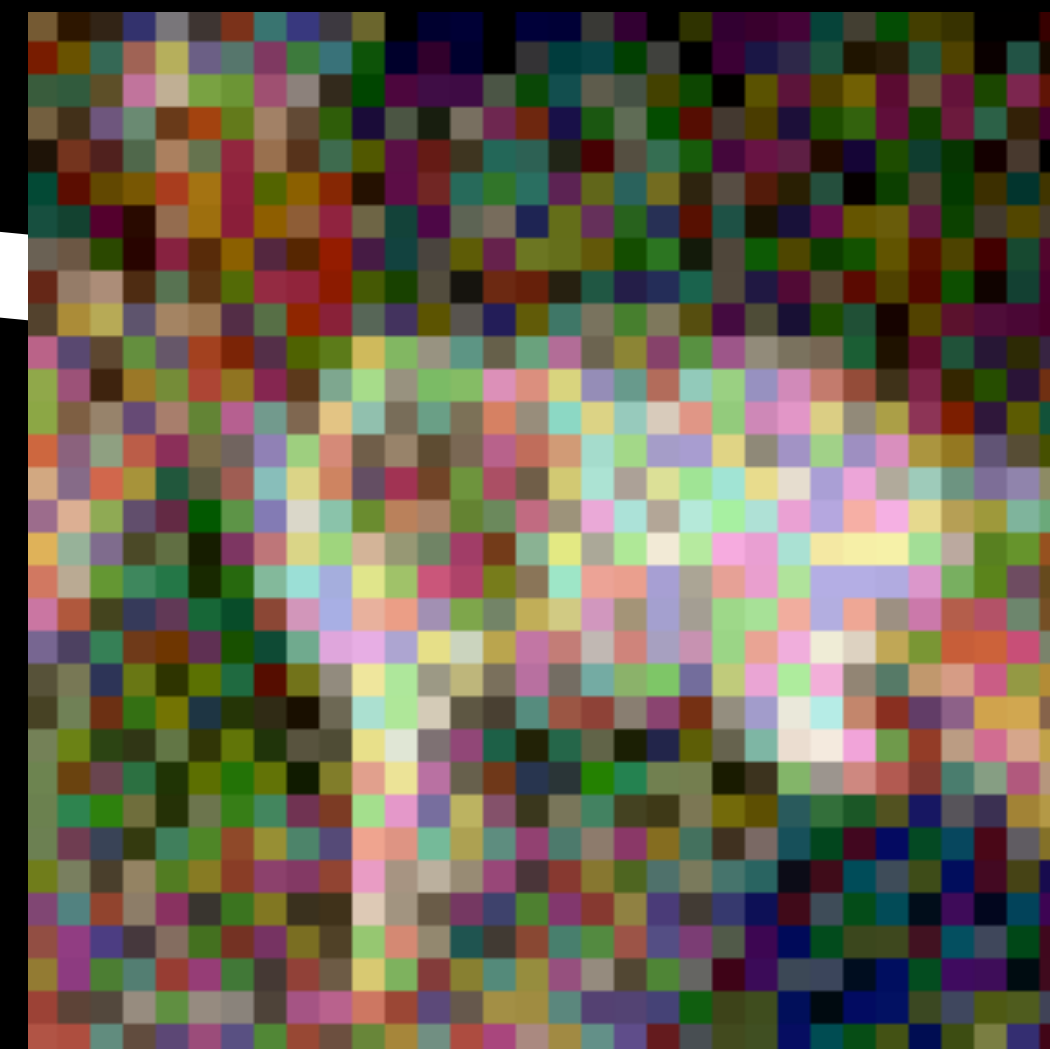
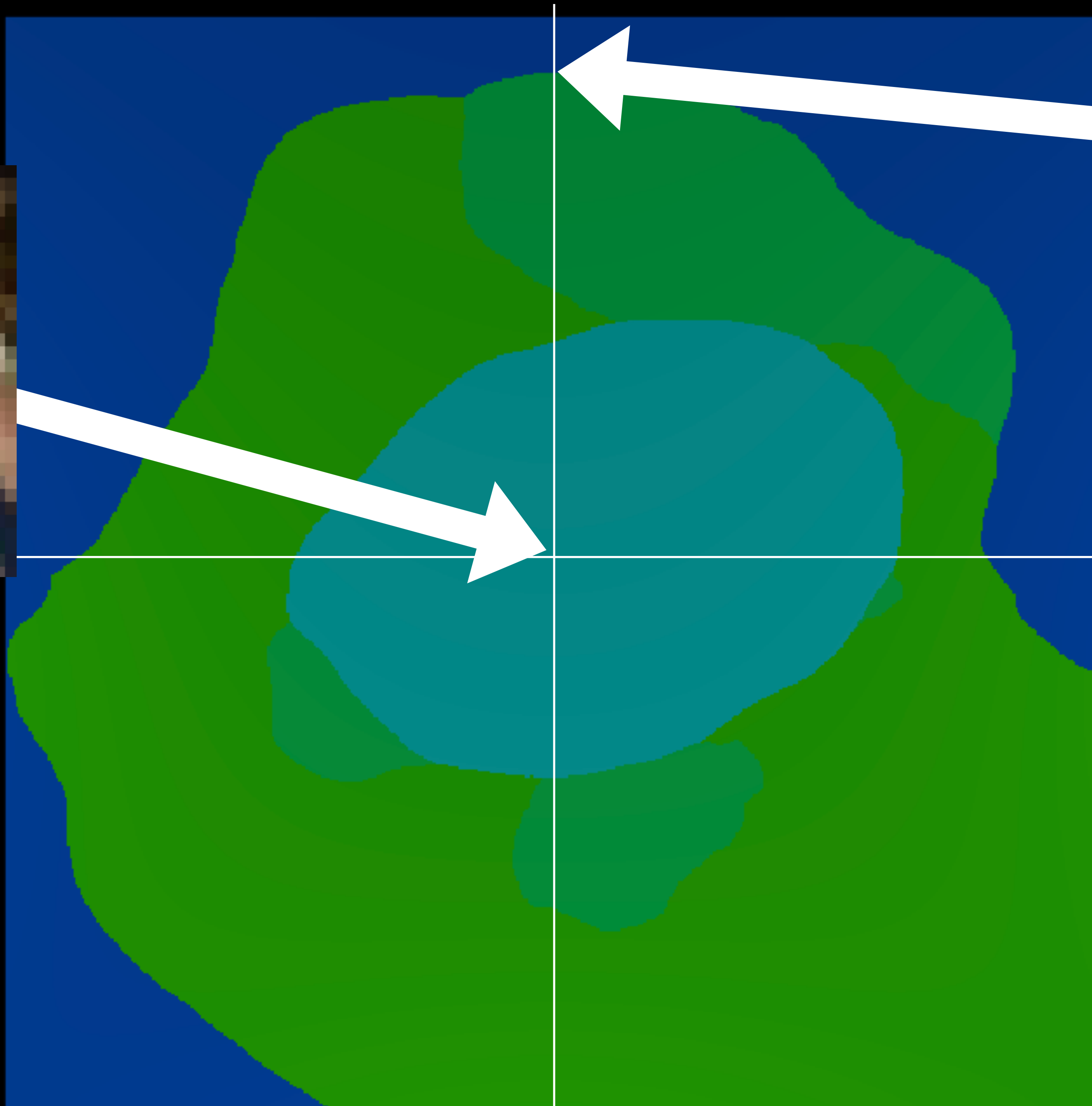




WHY does this work?



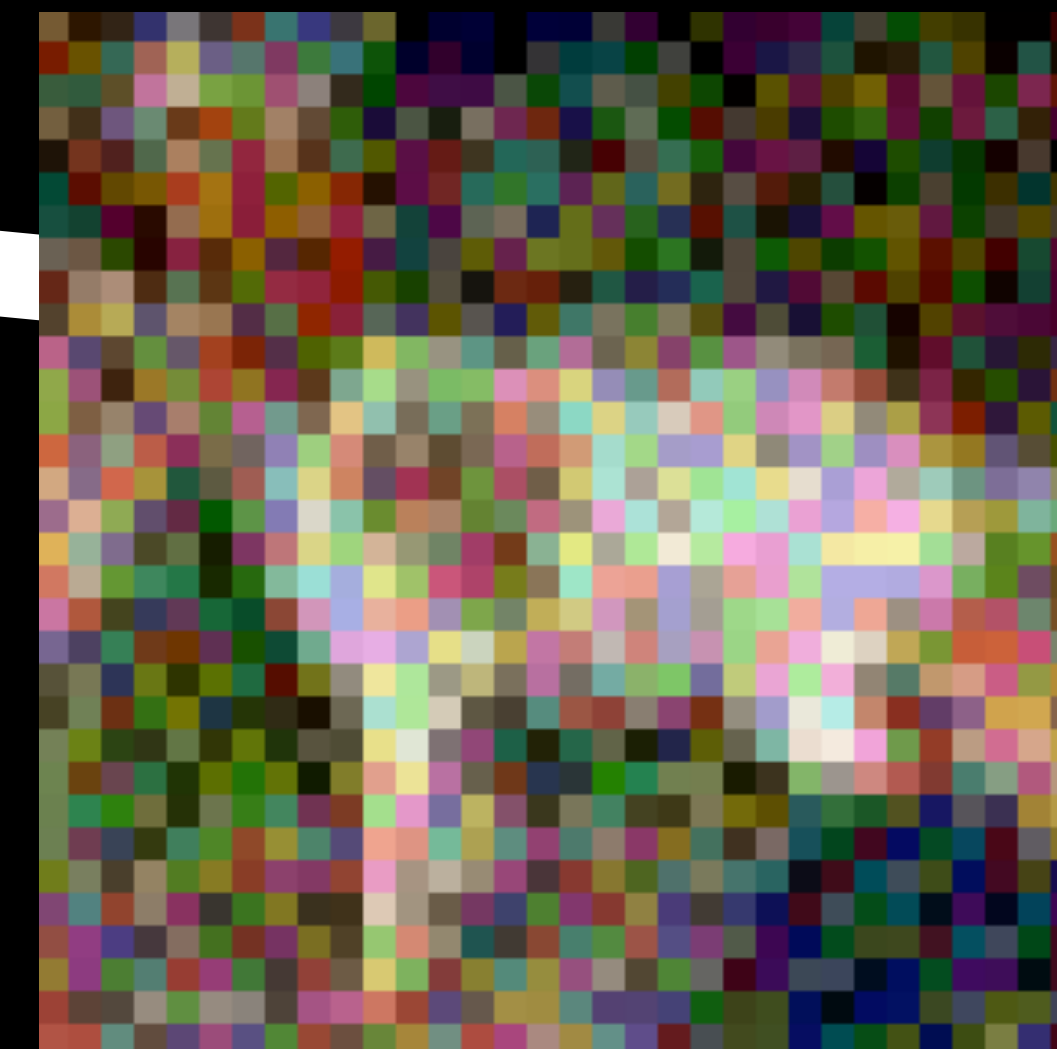
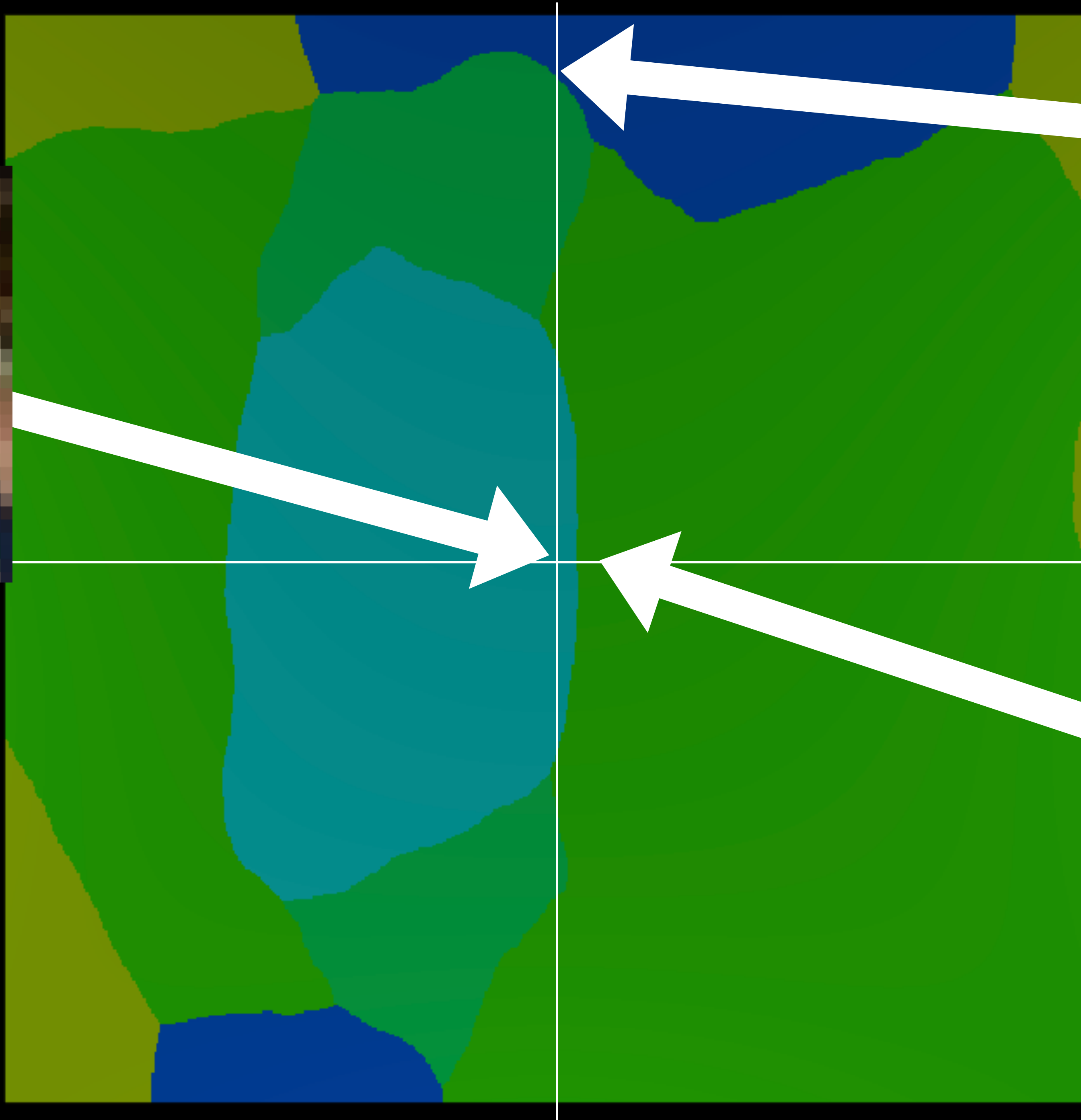
**Dog**



**Truck**



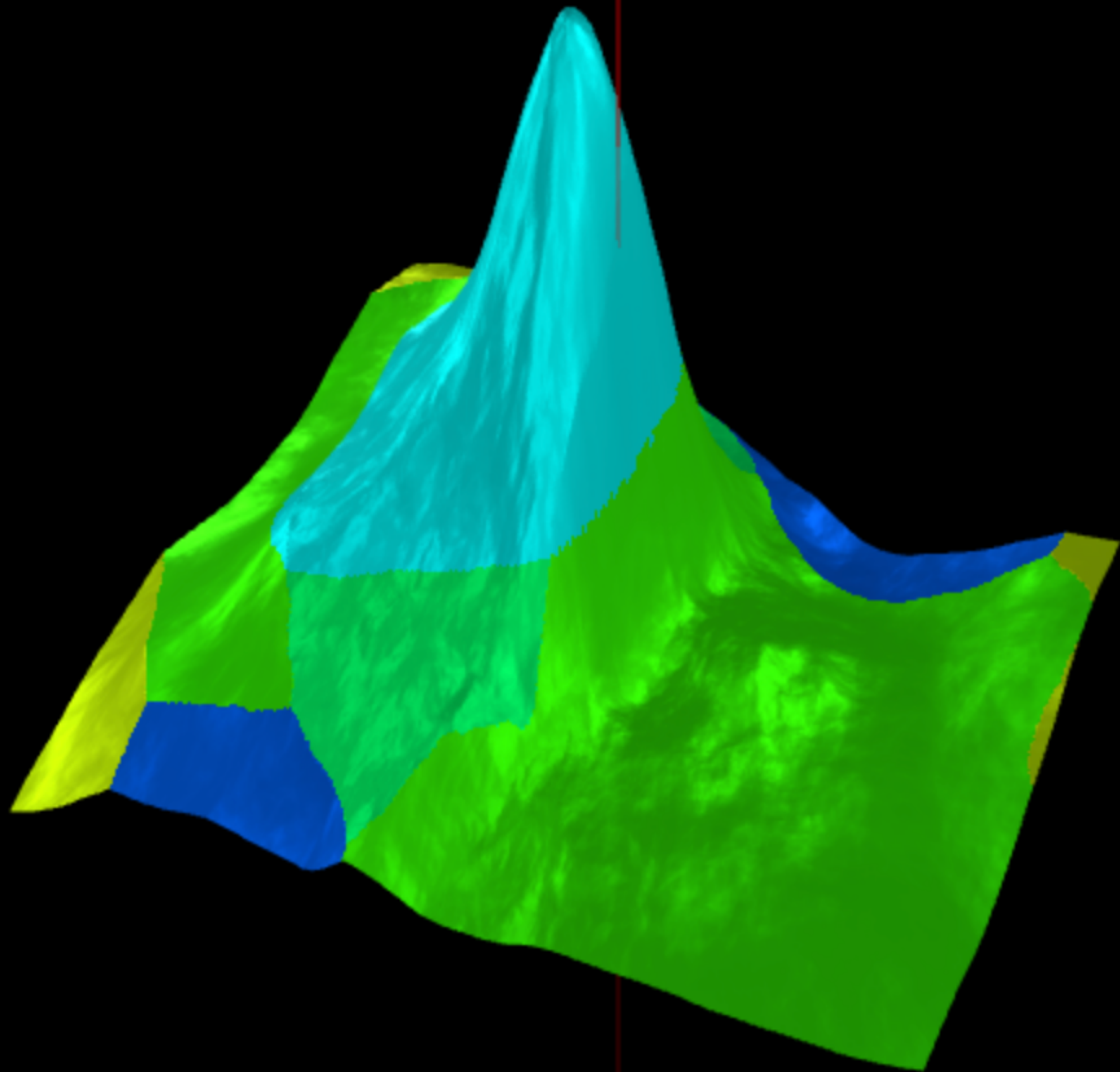
**Dog**

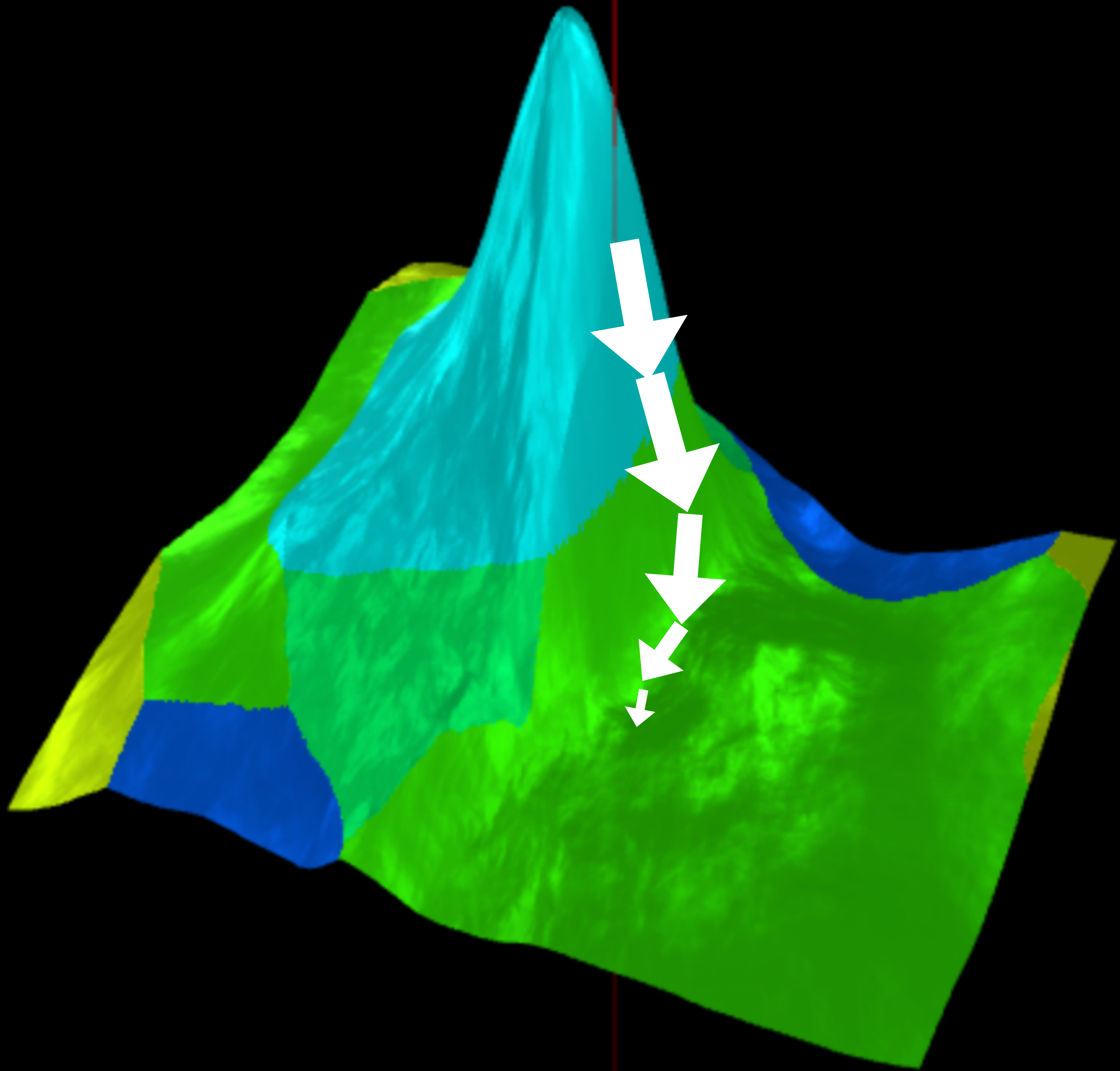


**Truck**



**Airplane**





Okay, lesson learned.

Okay, lesson learned.

Don't classify dogs with  
neural networks.





99.99% it is a

**School  
Bus**

Okay, lesson learned.

Okay, lesson learned.

*images*

Don't classify ~~dogs~~ with  
neural networks.



And now for something  
completely different

# Mozilla's DeepSpeech

Mozilla's DeepSpeech  
transcribes this

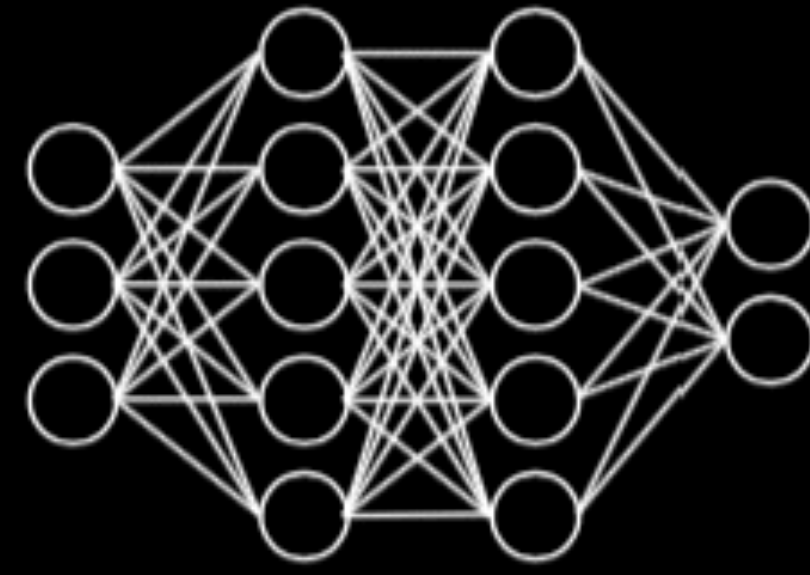
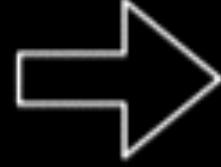
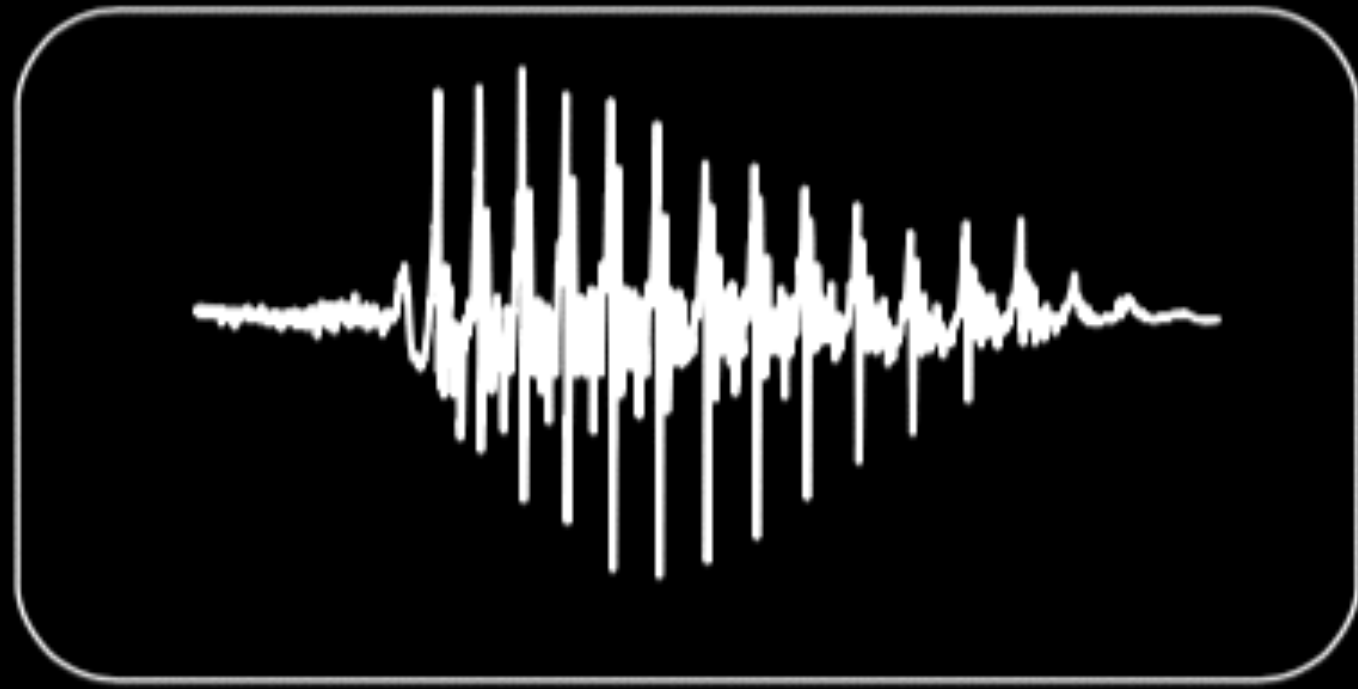
Mozilla's DeepSpeech  
transcribes this as

"most of them were staring  
quietly at the big table"

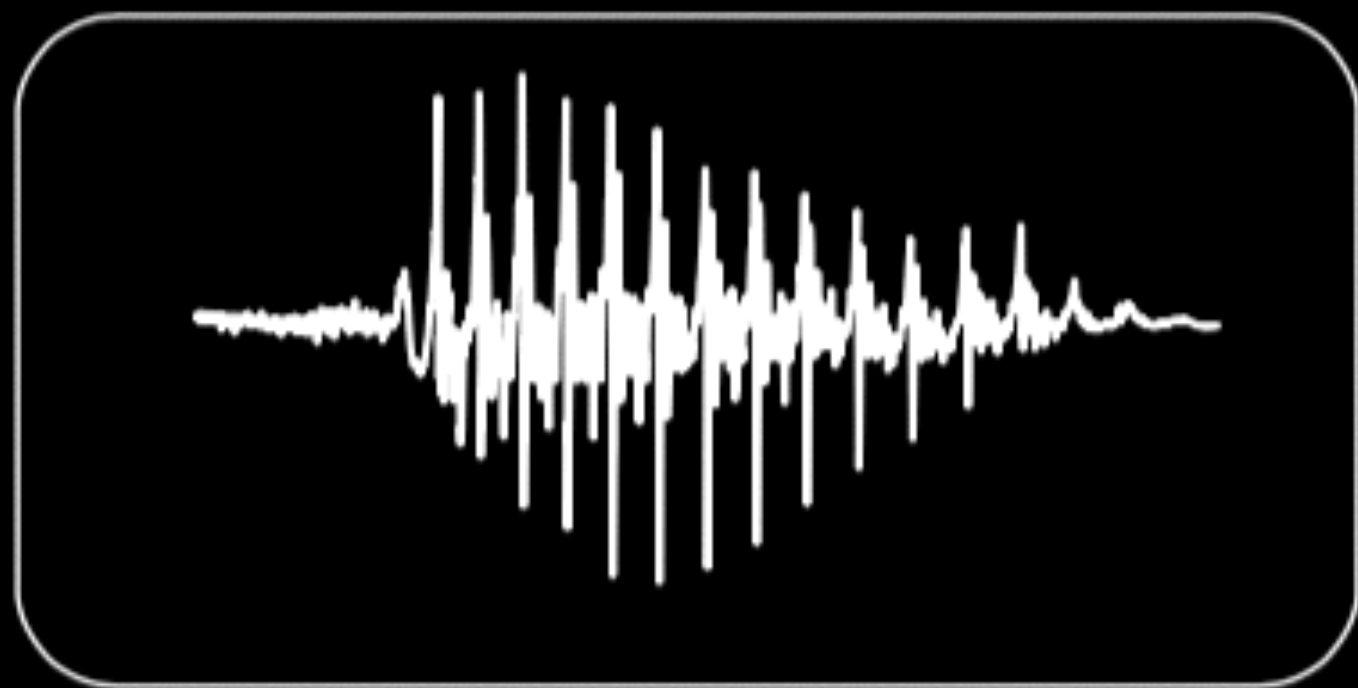
What about this?



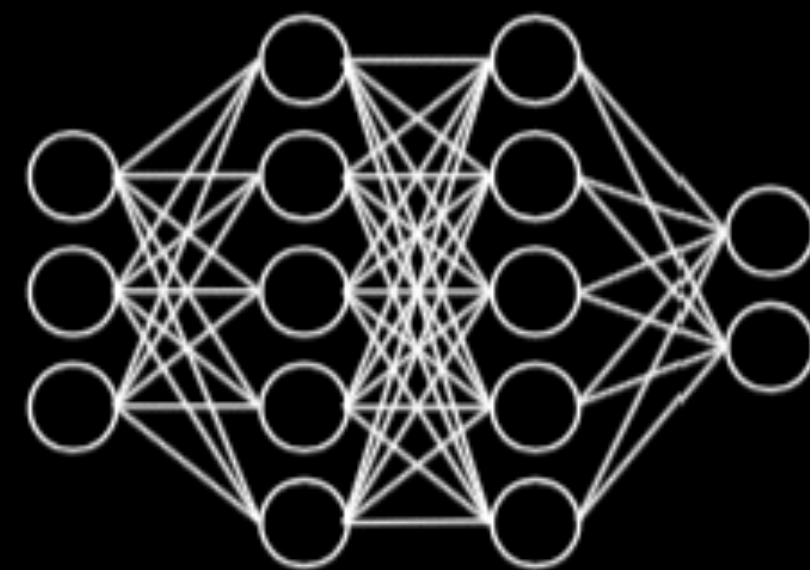
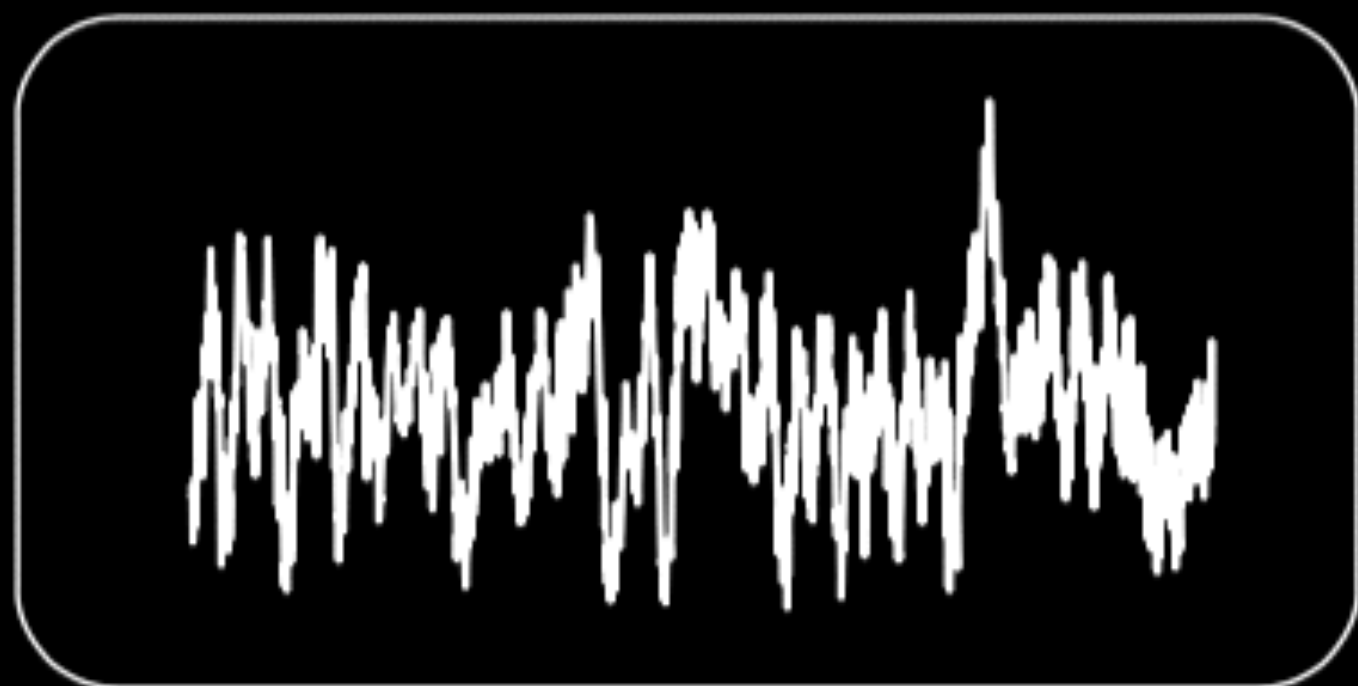
"It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity"



"it was the  
best of times,  
it was the  
worst of times"

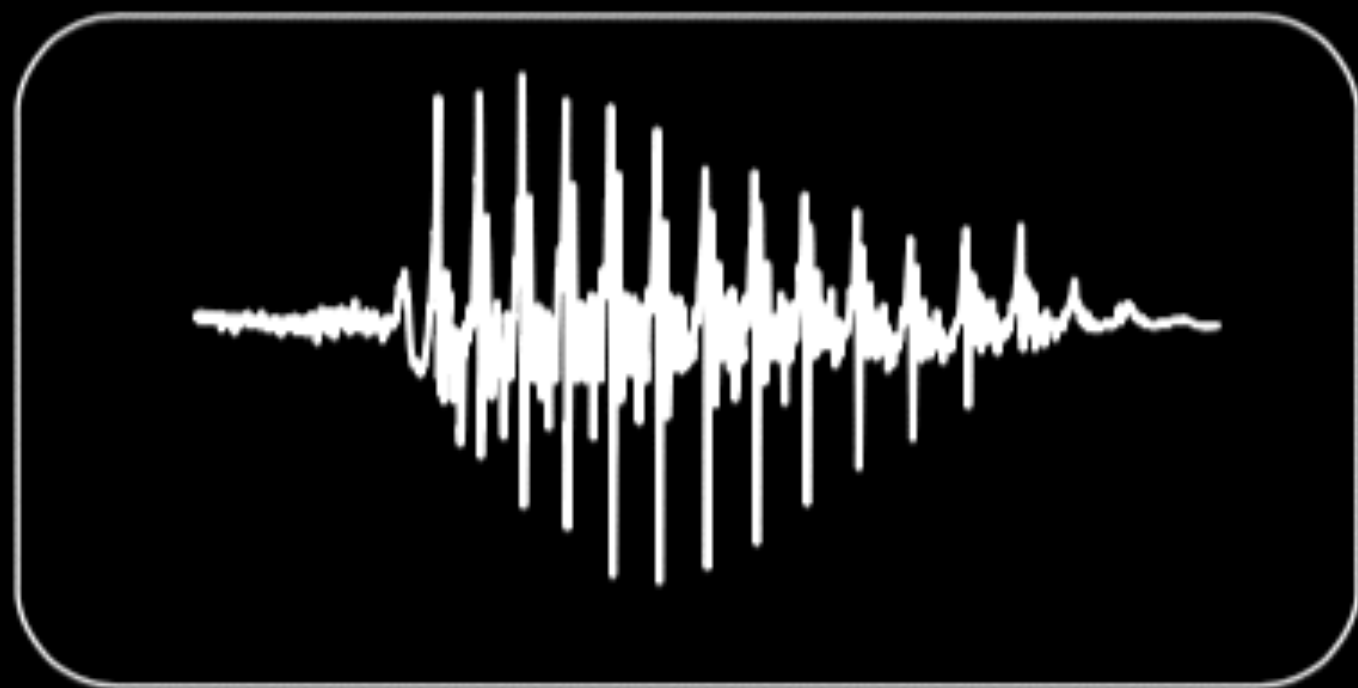


+

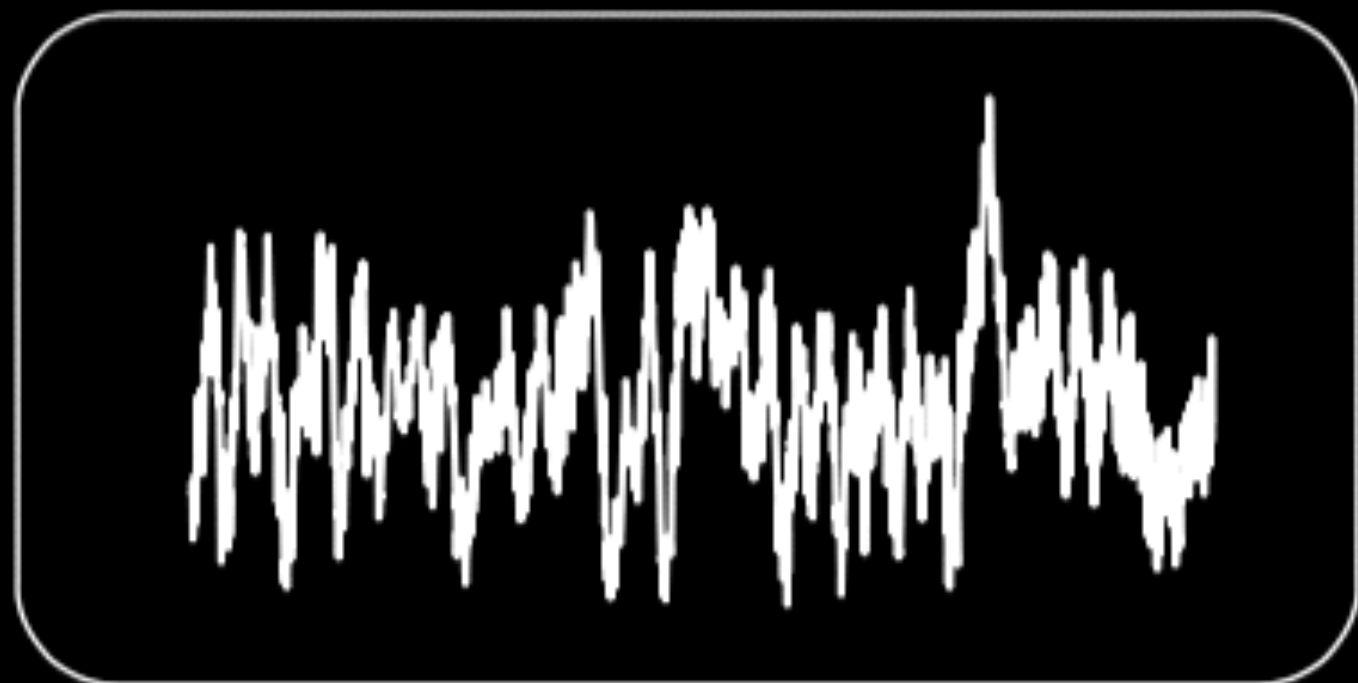


"it was the  
best of times,  
it was the  
worst of times"

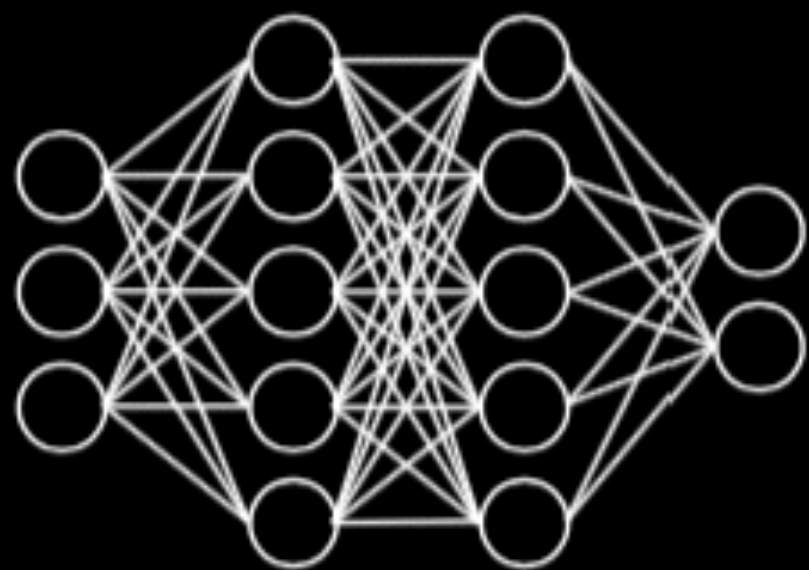
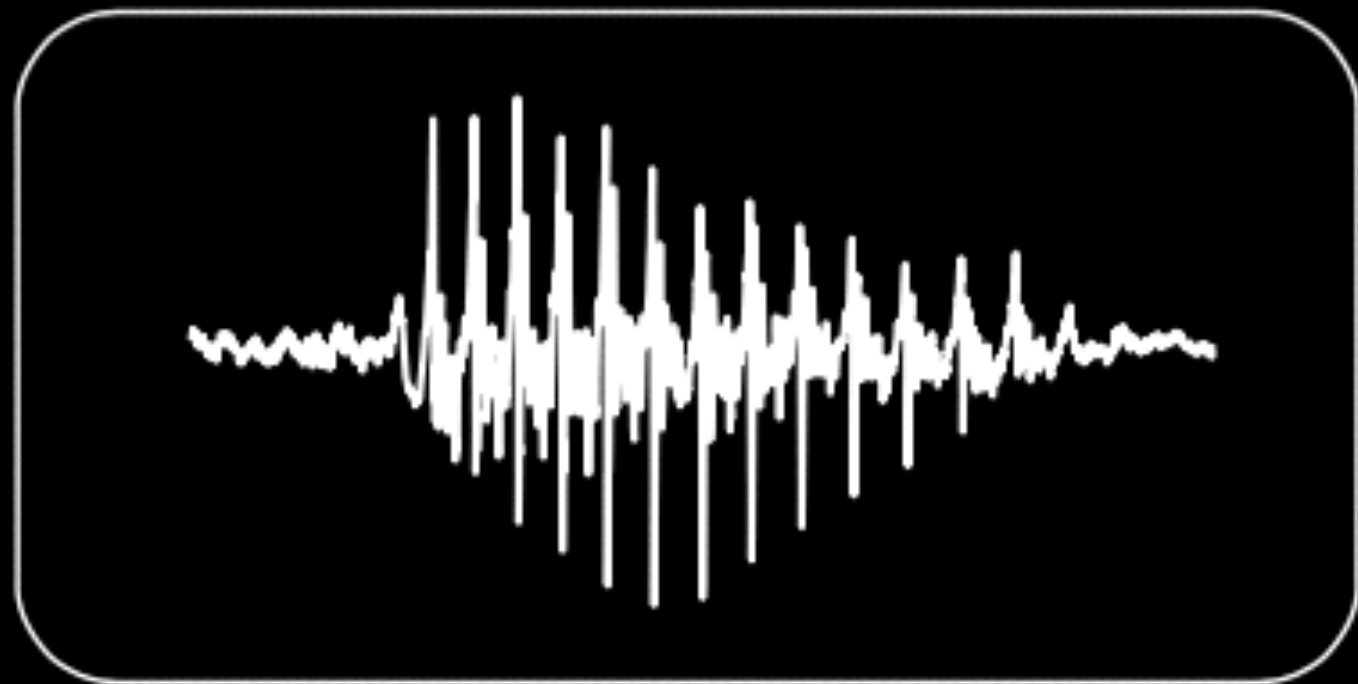
× 0.001



+

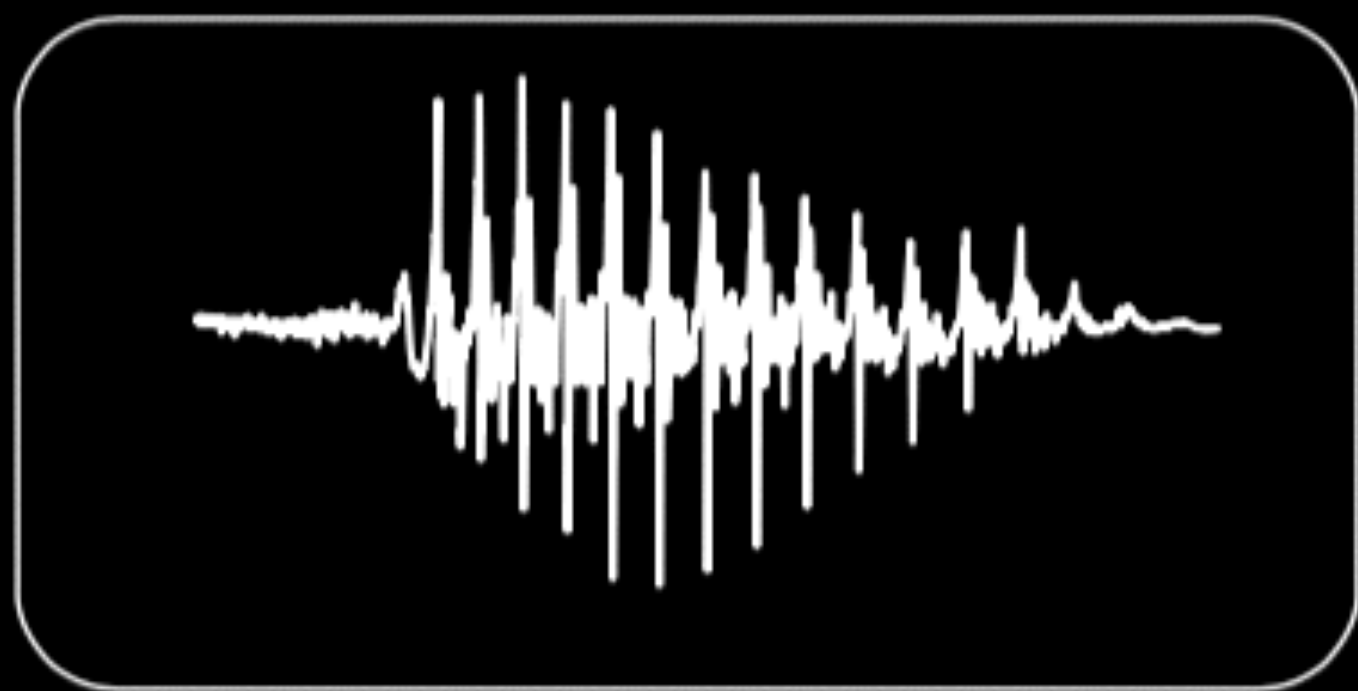


=

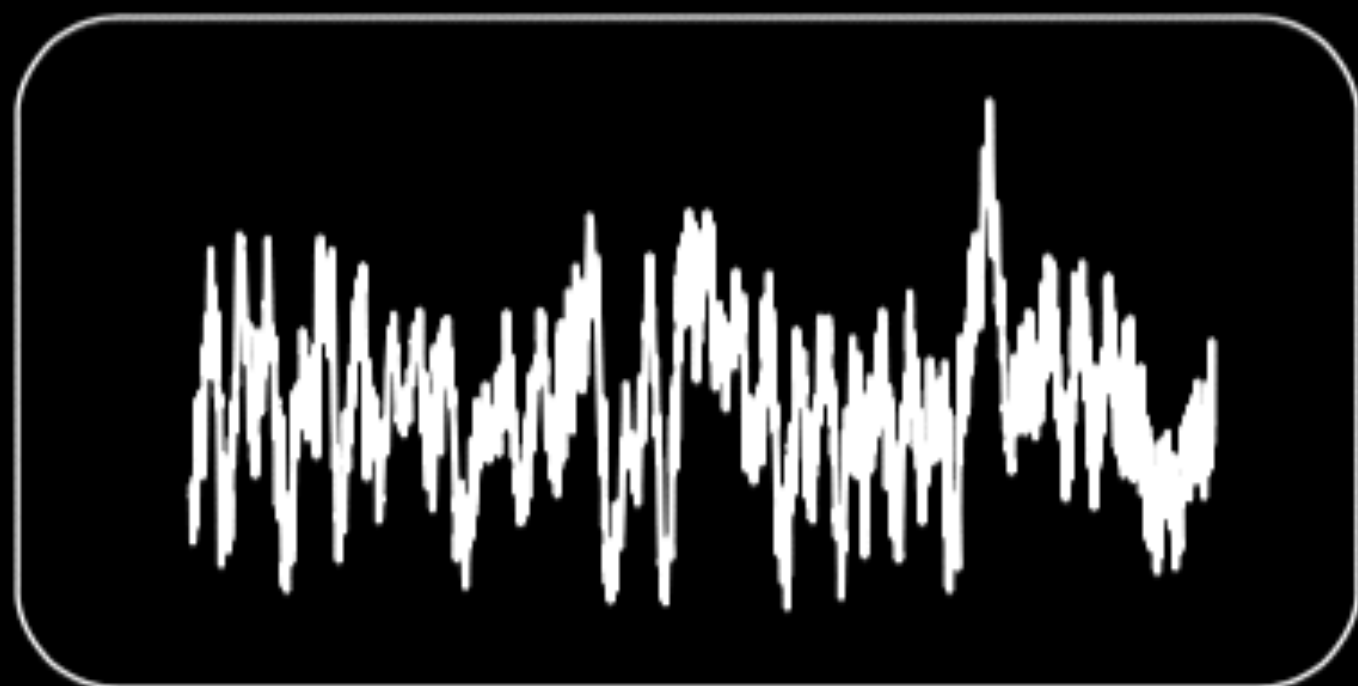


"it was the  
best of times,  
it was the  
worst of times"

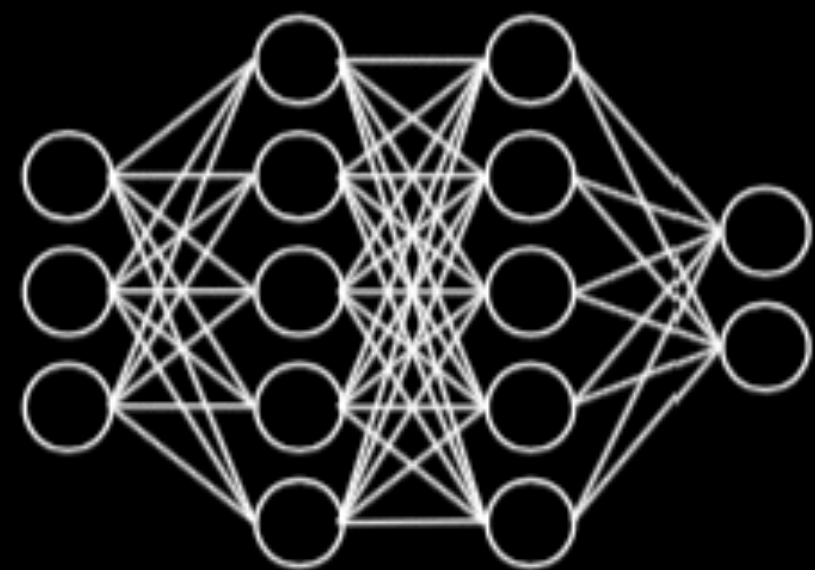
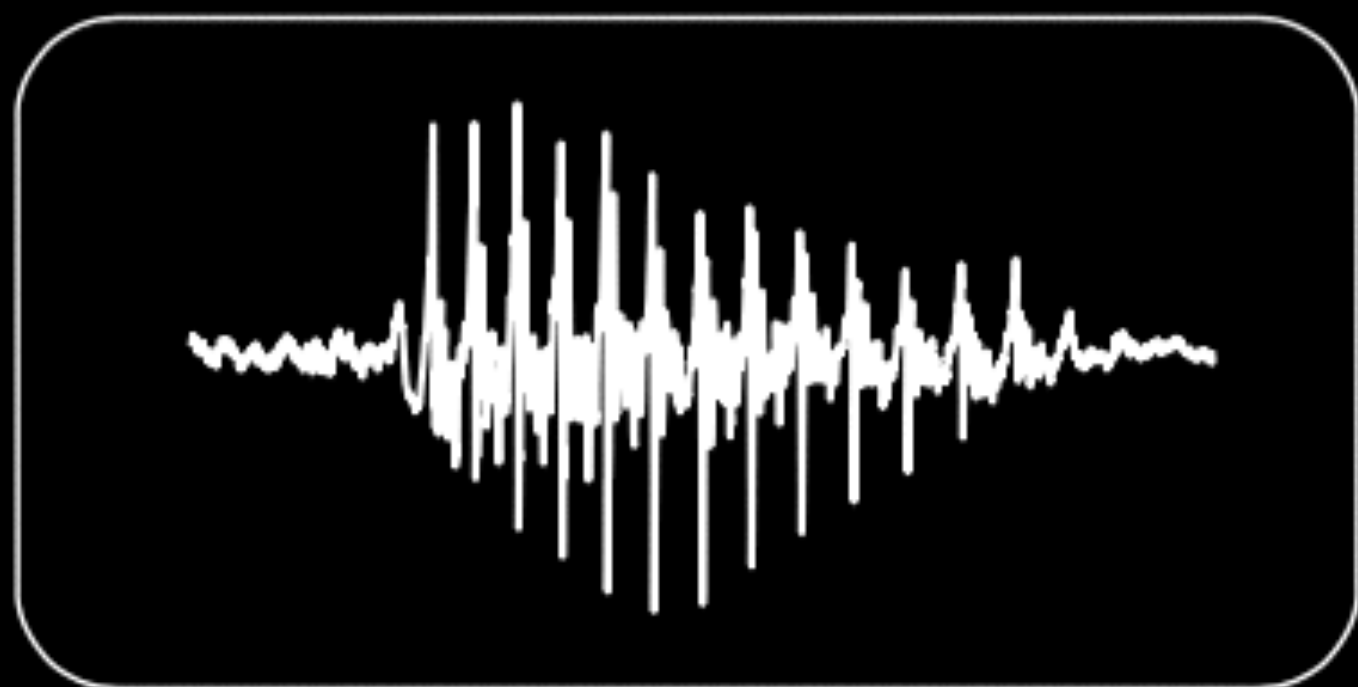
× 0.001



+

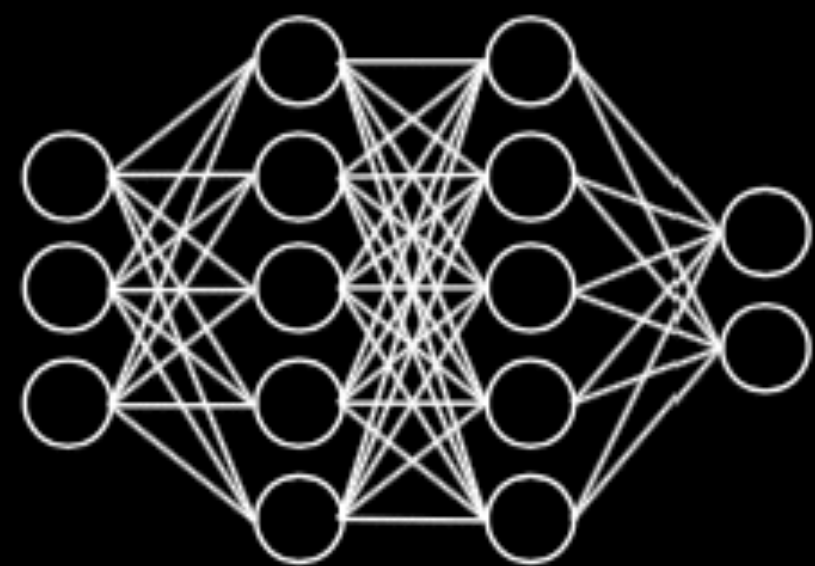


=



"it was the  
best of times,  
it was the  
worst of times"

× 0.001



"it is a truth  
universally  
acknowledged  
that a single"

Okay, lesson learned.

Okay, lesson learned.

or audio

Don't classify images with  
neural networks.

## Generating Natural Language Adversarial Examples

Moustafa Alzantot<sup>1\*</sup>, Yash Sharma<sup>2\*</sup>, Ahmed Elgohary<sup>3</sup>,  
Bo-Jhang Ho<sup>1</sup>, Mani B. Srivastava<sup>1</sup>, Kai-Wei Chang<sup>1</sup>

<sup>1</sup>Department of Computer Science, University of California, Los Angeles (UCLA)  
{malzantot, bojhang, mbs, kwchang}@ucla.edu

<sup>2</sup>Cooper Union sharma2@cooper.edu

<sup>3</sup>Computer Science Department, University of Maryland elgohary@cs.umd.edu

---

## Adversarial Attacks on Neural Network Policies

---

Sandy Huang<sup>†</sup>, Nicolas Papernot<sup>‡</sup>, Ian Goodfellow<sup>§</sup>, Yan Duan<sup>†§</sup>, Pieter Abbeel<sup>†§</sup>

<sup>†</sup> University of California, Berkeley, Department of Electrical Engineering and Computer Sciences

<sup>‡</sup> Pennsylvania State University, School of Electrical Engineering and Computer Science

<sup>§</sup> OpenAI

### Abstract

Machine learning classifiers are known to be vulnerable to inputs maliciously constructed by adversaries to force misclassification. Such adversarial examples have been extensively studied in the context of computer vision applications. In this work, we show adversarial attacks are also effective when targeting neural network

## Seq2Sick: Evaluating the Robustness of Sequence-to-Sequence Models with Adversarial Examples

Minhao Cheng<sup>1</sup>, Jinfeng Yi<sup>2</sup>, Huan Zhang<sup>1</sup>, Pin-Yu Chen<sup>3</sup>, Cho-Jui Hsieh<sup>1</sup>

<sup>1</sup>Department of Computer Science, University of California, Davis, CA 95616

<sup>2</sup>Tencent AI Lab, Bellevue, WA 98004

<sup>3</sup>IBM Research AI, Yorktown Heights, NY 10598

mhcheng@ucdavis.edu, jinfengyi.ustc@gmail.com, ecezhang@ucdavis.edu,  
pin-yu.chen@ibm.com, chohsieh@ucdavis.edu

## HALLUCINATIONS IN NEURAL MACHINE TRANSLATION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Neural machine translation (NMT) systems have reached state of the art performance in translating text and are in wide deployment. Yet little is understood about how these systems function or break. Here we show that NMT systems are susceptible to producing highly pathological translations that are completely untethered from the source material, which we term *hallucinations*. Such pathological translations are problematic because they are deeply disturbing of user trust and easy to find with a simple search. We describe a method to generate hallucinations and show that many common variations of the NMT architecture are susceptible to them. We study a variety of approaches to reduce the frequency

of hal  
nique  
nally,  
in the

## SYNTHETIC AND NATURAL NOISE BOTH BREAK NEURAL MACHINE TRANSLATION

Yonatan Belinkov\*

Computer Science and  
Artificial Intelligence Laboratory,  
Massachusetts Institute of Technology  
belinkov@mit.edu

Yonatan Bisk\*

Paul G. Allen School  
of Computer Science & Engineering,  
University of Washington  
ybisk@cs.washington.edu

## On the Robustness of Semantic Segmentation Models to Adversarial Attacks

Anurag Arnab Ondrej Miksik Philip H.S. Torr

University of Oxford

{anurag.arnab, ondrej.miksik, philip.torr}@eng.ox.ac.uk



Okay, lesson learned.

Okay, lesson learned.

Don't let adversaries  
perform gradient descent.

# MITIGATING ADVERSARIAL EFFECTS THROUGH RANDOMIZATION

**Cihang Xie, Zhishuai Zhang &**  
Department of Computer Science  
The Johns Hopkins University  
Baltimore, MD 21218 USA  
{cihangxie306, zhshuai.

**Jianyu Wang**  
Baidu Research USA  
Sunnyvale, CA 94089 USA  
wjyouch@gmail.com

**Zhou Ren**  
Snap Inc.  
Venice, CA 90291 USA  
zhou.ren@snapchat.com

Convolutional neural networks have become the dominant paradigm in recent years. However, adversarial examples have been found to fool these networks. For example, imperceptible perturbations can cause convolutional neural networks to misclassify. To guard against adversarial examples, we take inspiration from game theory and propose randomization operations: randomizing the order of operations, randomizing the kernel size, and random padding. Our method is very effective against adversarial attacks. Our method provides 1) no need for fine-tuning, 2) very few additional computations, 3) compatible with other adversarial defense methods. By combining the proposed randomization method with an adversarially trained model, it achieves a normalized score of 0.924 (ranked No.2 among 107 defense teams) in the NIPS 2017 adversarial examples defense challenge, which is far better than using adversarial training alone with a normalized score of 0.773 (ranked No.56). The code is public available at [https://github.com/cihangxie/NIPS2017\\_adv\\_challenge\\_defense](https://github.com/cihangxie/NIPS2017_adv_challenge_defense).

1) no need for fine-tuning, 2) very few additional computations, 3) compatible with other adversarial defense methods. By combining the proposed randomization method with an adversarially trained model, it achieves a normalized score of 0.924 (ranked No.2 among 107 defense teams) in the NIPS 2017 adversarial examples defense challenge, which is far better than using adversarial training alone with a normalized score of 0.773 (ranked No.56). The code is public available at [https://github.com/cihangxie/NIPS2017\\_adv\\_challenge\\_defense](https://github.com/cihangxie/NIPS2017_adv_challenge_defense).

# STOCHASTIC ACTIVATION PRUNING FOR ROBUST ADVERSARIAL DEFENSE

**Guneet S. Dhillon<sup>1,2</sup>, Kamyar Azizzadenesheli<sup>3</sup>, Zachary C. Lipton<sup>1,4</sup>,  
Jeremy Bernstein<sup>1,5</sup>, Jean Kossaifi<sup>1,6</sup>, Aran Khanna<sup>1</sup>, Anima Anandkumar<sup>1,5</sup>**  
<sup>1</sup>Amazon AI, <sup>2</sup>UT Austin, <sup>3</sup>UC Irvine, <sup>4</sup>CMU, <sup>5</sup>Caltech, <sup>6</sup>Imperial College London  
guneetdhillon@utexas.edu, kazizzad@uci.edu, zlipton@cmu.edu,  
bernstein@caltech.edu, jean.kossaifi@imperial.ac.uk,  
aran@arankhanna.com, anima@amazon.com

## ABSTRACT

Neural networks are known to be vulnerable to adversarial examples: chosen perturbations to real images, while imperceptible to humans, can cause misclassification and threaten the reliability of deep learning systems. To guard against adversarial examples, we take inspiration from game theory and propose the problem as a minimax zero-sum game between the adversary and the defender. In general, for such games, the optimal strategy for both players is a mixed strategy, also known as a *mixed strategy*. In this light, we propose *Activation Pruning* (SAP), a mixed strategy for adversarial defense: randomly deactivate a random subset of activations (preferentially pruning those with high variance) and scales up the survivors to compensate. We can apply SAP to any neural networks, including adversarially trained models, without fine-tuning. Our method provides robustness against adversarial examples. Experiments demonstrate that SAP provides robustness against attacks, increasing accuracy and preserving cal-

1) no need for fine-tuning, 2) very few additional computations, 3) compatible with other adversarial defense methods. By combining the proposed randomization method with an adversarially trained model, it achieves a normalized score of 0.924 (ranked No.2 among 107 defense teams) in the NIPS 2017 adversarial examples defense challenge, which is far better than using adversarial training alone with a normalized score of 0.773 (ranked No.56). The code is public available at [https://github.com/cihangxie/NIPS2017\\_adv\\_challenge\\_defense](https://github.com/cihangxie/NIPS2017_adv_challenge_defense).

# THERMOMETER ENCODING: ONE HOT WAY TO RESIST ADVERSARIAL EXAMPLES

**Jacob Buckman\*<sup>†</sup>, Aurko Roy\*, Colin Raffel, Ian Goodfellow**  
Google Brain  
Mountain View, CA  
{buckman, aurkor, craffel, goodfellow}@google.com

## ABSTRACT

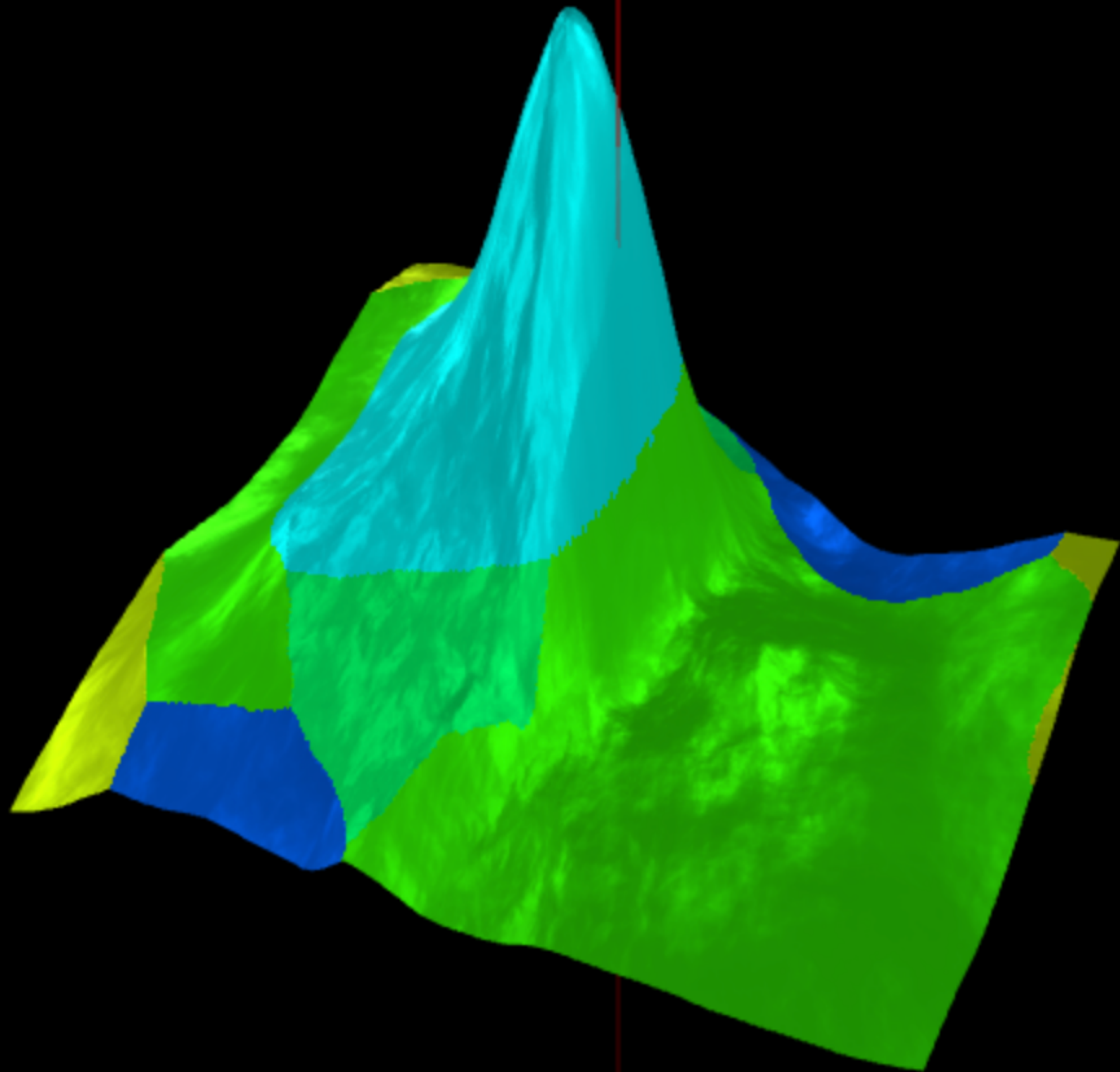
# COUNTERING ADVERSARIAL IMAGES USING INPUT TRANSFORMATIONS

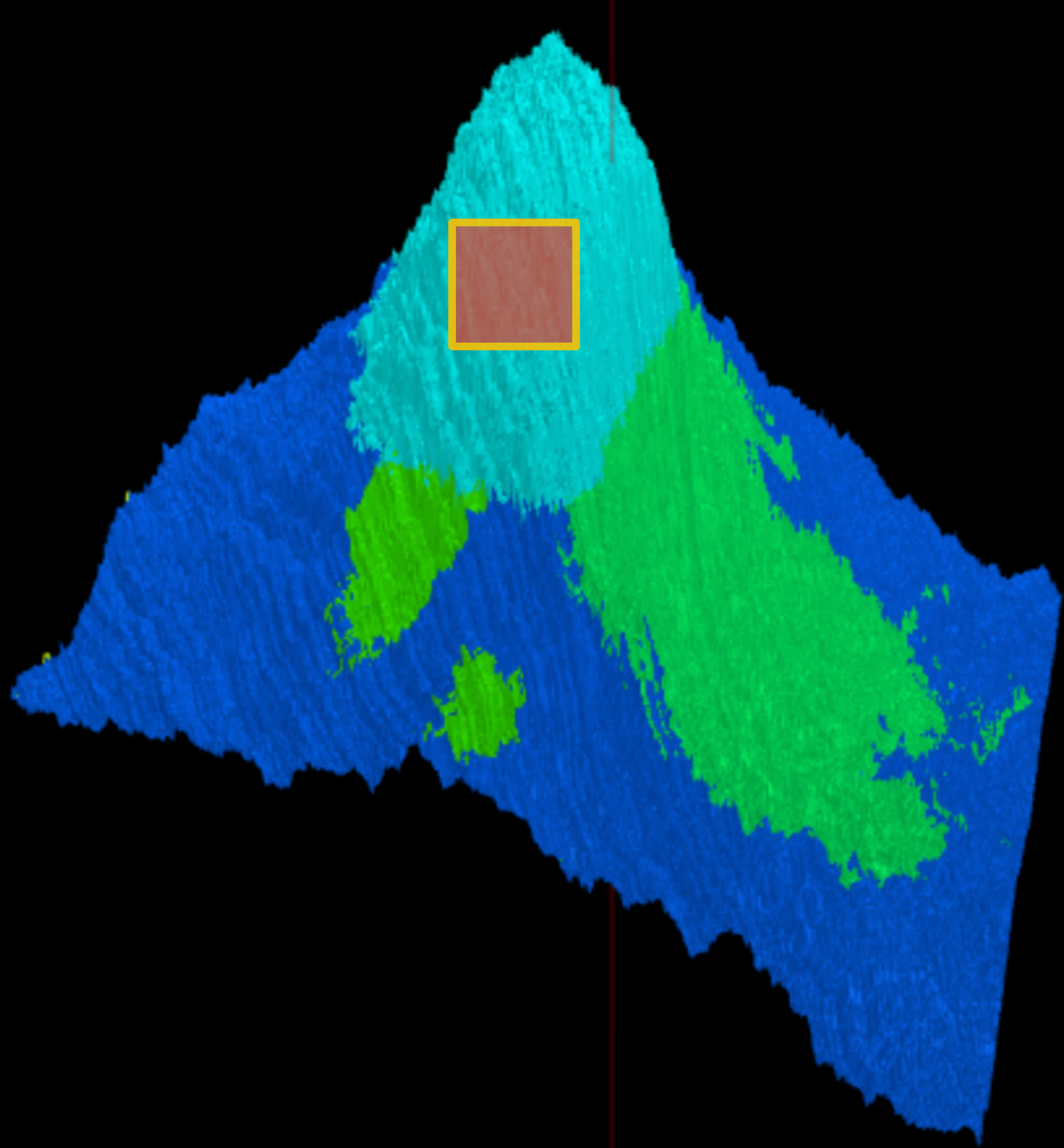
**Chuan Guo\***      **Mayank Rana & Moustapha Cissé & Laurens van der Maaten**  
Cornell University      Facebook AI Research

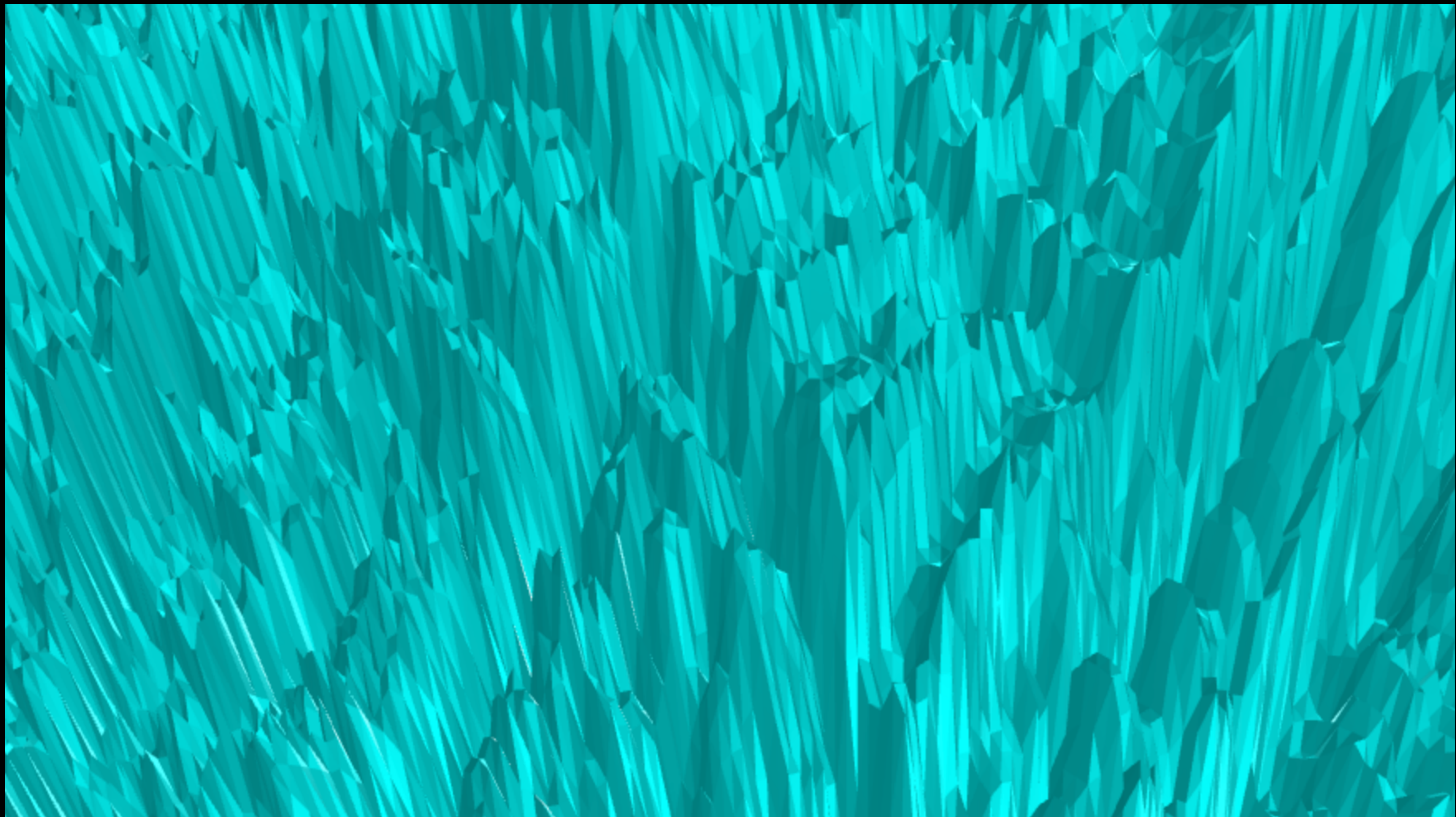
## ABSTRACT

This paper investigates strategies that defend against adversarial-example attacks on image-classification systems by transforming the inputs before feeding them to the system. Specifically, we study applying image transformations such as bit-depth reduction, JPEG compression, total variance minimization, and image quilting before feeding the image to a convolutional network classifier. Our experiments on ImageNet show that total variance minimization and image quilting are very effective defenses in practice, in particular, when the network is trained on transformed images. The strength of those defenses lies in their non-differentiable nature and their inherent randomness, which makes it difficult for an adversary to circumvent the defenses. *Our best defense eliminates 60% of strong gray-box and 90% of strong black-box attacks by a variety of major attack methods.*

adversarial examples” for neural networks. These adversarial examples are indistinguishable from real images to humans but cause the neural network to misclassify. This paper explores the robustness of neural networks against adversarial examples. We show that adversarial examples exist for a wide range of datasets, and show that adversarial examples can be generated for a wide range of models. Our results show that adversarial examples can be generated for a wide range of models, and show that adversarial examples can be generated for a wide range of models.







# Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples

Anish Athalye<sup>\*1</sup> Nicholas Carlini<sup>\*2</sup> David Wagner<sup>2</sup>

## Abstract

We identify obfuscated gradients, a kind of gradient masking, as a phenomenon that leads to a false sense of security in defenses against adversarial examples. While defenses that cause obfuscated gradients appear to defeat iterative optimization-based attacks, we find defenses relying on this effect can be circumvented. We describe characteristic behaviors of defenses exhibiting the effect, and for each of the three types of obfuscated gradients we discover, we develop attack techniques to overcome it. In a case study, examining non-certified white-box-secure defenses at ICLR 2018, we find obfuscated gradients are a common occurrence, with 7 of 9 defenses relying on obfuscated gradients. Our new attacks successfully circumvent 6 completely, and 1 partially, in the original threat model each paper considers.

apparent robustness against iterative optimization-based attacks. We identify obfuscated gradients, a term we define as a specific kind of gradient masking (Papernot et al., 2017). Without

7 of 9 defenses relying on obfuscated gradients. Our new attacks successfully circumvent 6 completely, and 1 partially

caused by these three phenomena. We address gradient shattering with a new attack technique we call Backward Pass Differentiable Approximation, where we approximate derivatives by computing the forward pass normally and computing the backward pass using a differentiable approximation of the function. We compute gradients of randomized defenses by applying Expectation Over Transformation (Athalye et al., 2017). We solve vanishing/exploding gradients through reparameterization and optimize over a space where gradients do not explode/vanish.

To investigate the prevalence of obfuscated gradients and understand the applicability of these attack techniques, we use as a case study the ICLR 2018 non-certified defenses that claim white-box robustness. We find that obfuscated gradients are a common occurrence, with 7 of 9 defenses relying on this phenomenon. Applying the new attack techniques we develop, we overcome obfuscated gradients and circumvent 6 of them completely, and 1 partially, under the original threat model of each paper. Along with this, we offer an analysis of the evaluations performed in the papers.

Additionally, we hope to provide researchers with a common baseline of knowledge, description of attack techniques, and common evaluation pitfalls, so that future defenses can avoid falling vulnerable to these same attack approaches.

To promote reproducible research, we release our reimplementation of each of these defenses, along with implementations of our attacks for each.<sup>1</sup>

<sup>1</sup> <https://github.com/anishathalye/obfuscated-gradients>

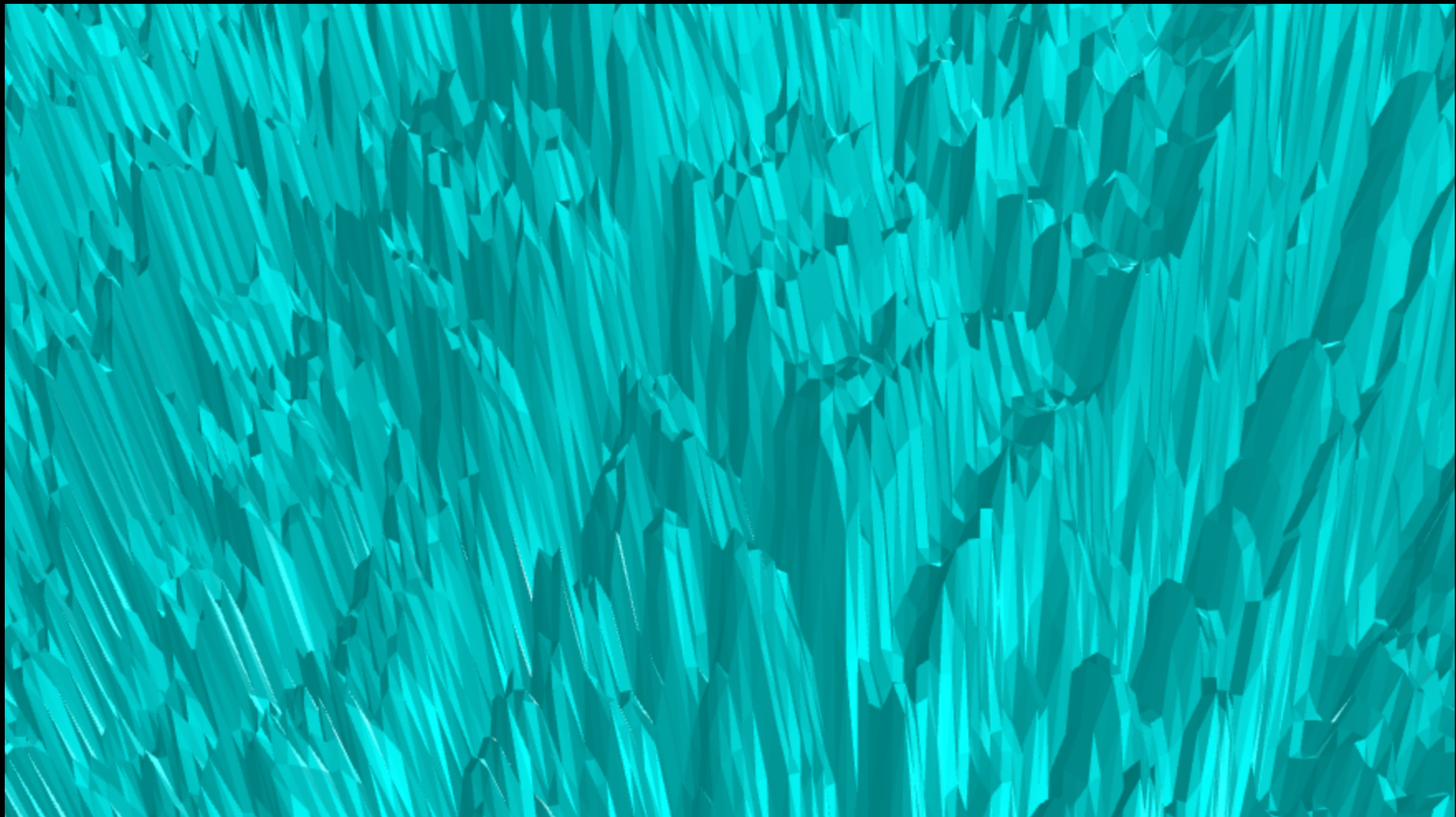
## 1. Introduction

In response to the susceptibility of neural networks to adversarial examples (Szegedy et al., 2013; Biggio et al., 2013), there has been significant interest recently in constructing defenses to increase the robustness of neural networks. While progress has been made in understanding and defending against adversarial examples in the white-box setting, where the adversary has full access to the network, a complete solution has not yet been found.

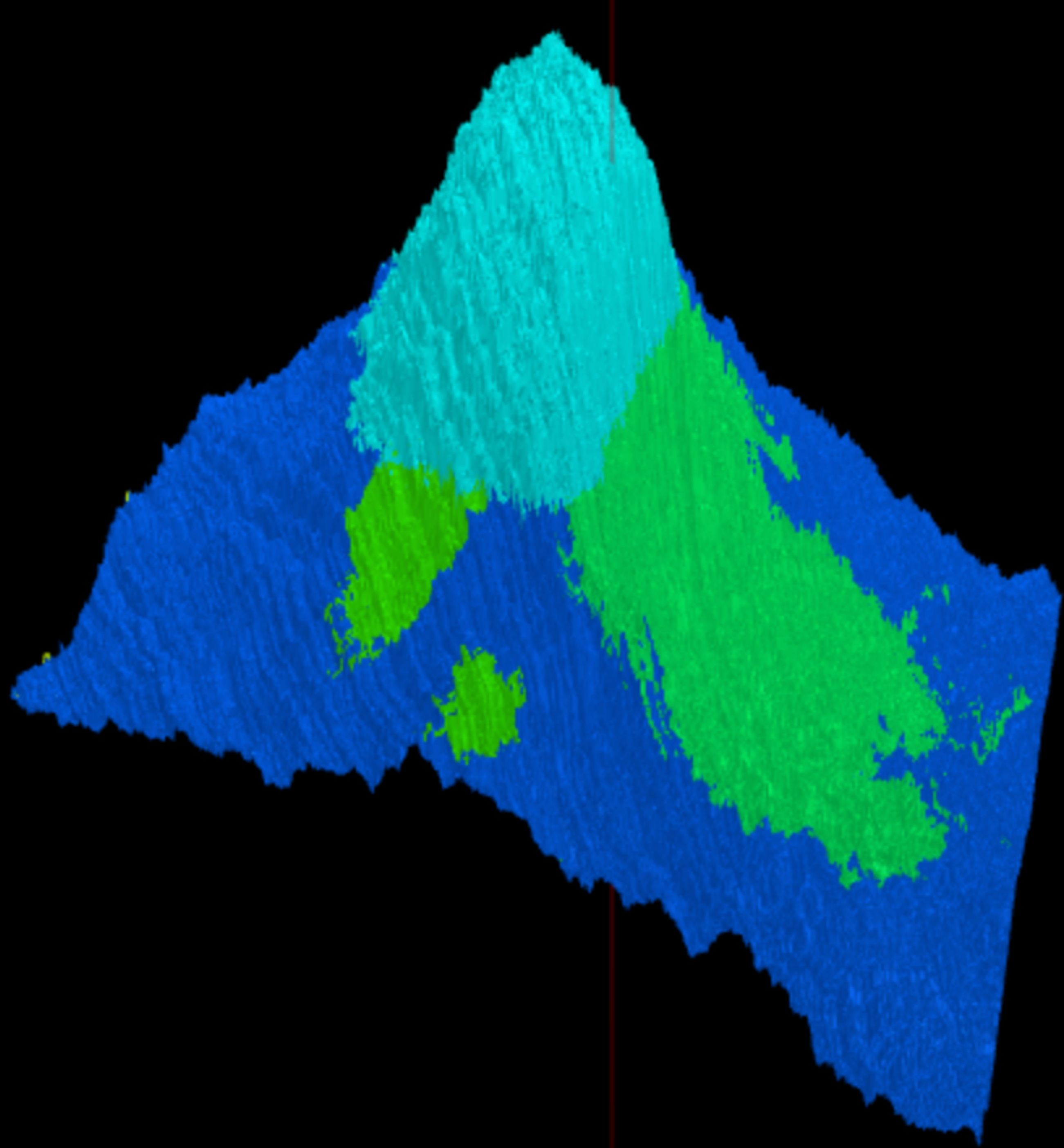
As benchmarking against iterative optimization-based attacks (e.g., Kurakin et al. (2016a); Madry et al. (2018); Carlini & Wagner (2017c)) has become standard practice in evaluating defenses, new defenses have arisen that appear to be robust against these powerful optimization-based attacks.

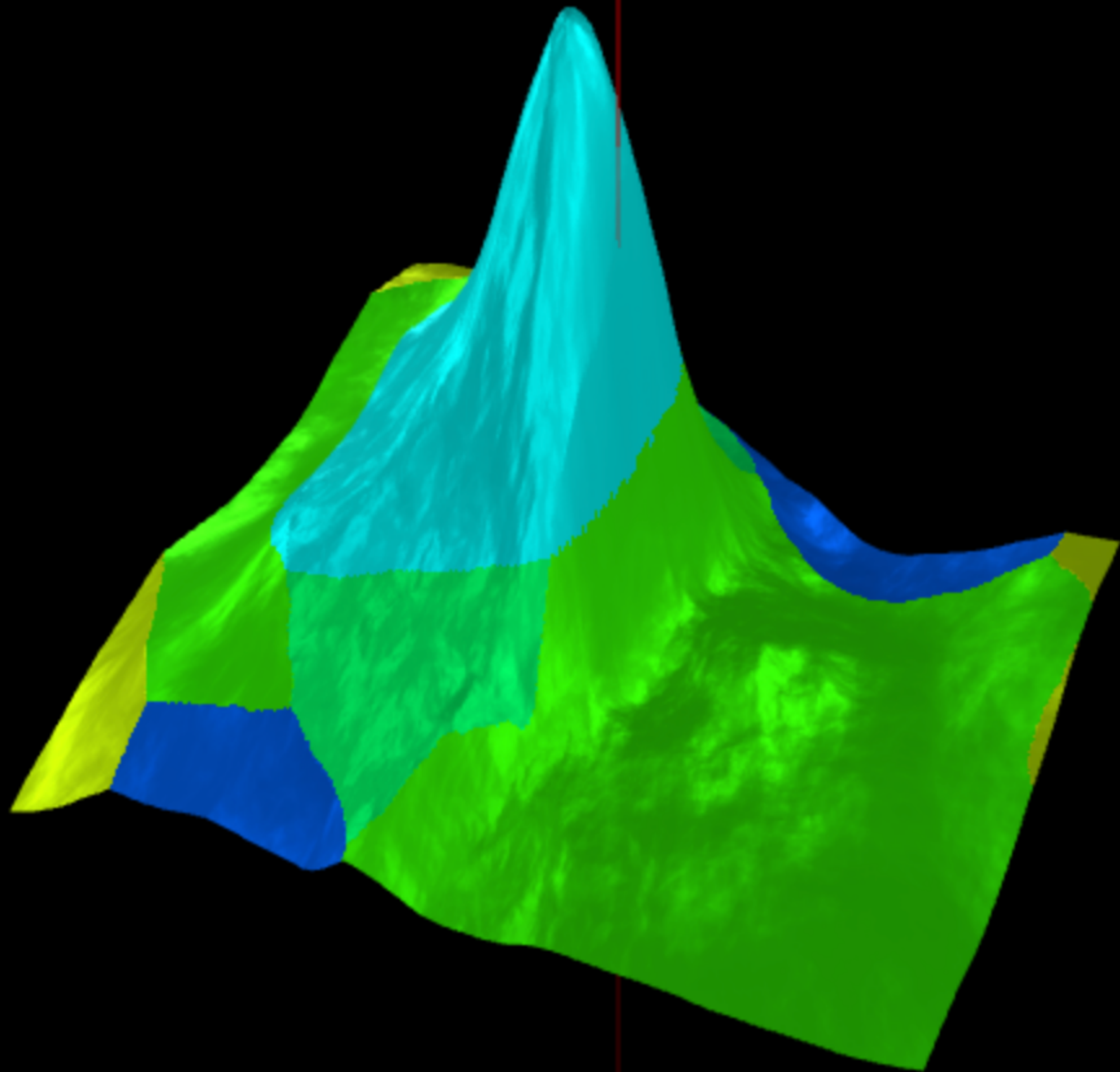
We identify one common reason why many defenses provide

<sup>\*</sup>Equal contribution <sup>1</sup>Massachusetts Institute of Technology <sup>2</sup>University of California, Berkeley. Correspondence to: Anish Athalye <aathalye@mit.edu>, Nicholas Carlini <npc@berkeley.edu>.









Okay, lesson learned.

Okay, lesson learned.

Don't let adversaries have  
**ANY** access to my model

## DECISION-BASED ADVERSARIAL ATTACKS: RELIABLE ATTACKS AGAINST BLACK-BOX MACHINE LEARNING MODELS

Wieland Brendel\*, Jonas Rauber\* & Matthias Bethge

Werner Reichardt Centre for  
Eberhard Karls University Tübingen  
{wieland, jonas, matth:}

## DELVING INTO TRANSFERABLE ADVERSARIAL EX- AMPLES AND BLACK-BOX ATTACKS

Yanpei Liu\*, Xinyun Chen\*  
Shanghai Jiao Tong University

Chang Liu, Dawn Song  
University of the California, Berkeley

### ABSTRACT

An intriguing property of deep neural networks is the existence of adversarial ex-  
amples, which can transfer among different architectures. These transferable ad-

## Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples

Nicolas Papernot and Patrick McDaniel  
The Pennsylvania State University  
University Park, PA  
{ngp5056, mcdaniel}@cse.psu.edu

Ian Goodfellow  
OpenAI  
San Francisco, CA  
ian@openai.com

which is a black-box image classification system.

## Universal adversarial perturbations

Seyed-Mohsen Moosavi-Dezfooli\*†  
seyed.moosavi@epfl.ch

Omar Fawzi†  
omar.fawzi@ens-lyon.fr

## ZOO: Zeroth Order Optimization Based Black-box Attacks to Deep Neural Networks without Training Substitute Models

Pin-Yu Chen\*  
AI Foundations Group  
IBM T. J. Watson Research Center  
Yorktown Heights, NY  
pin-yu.chen@ibm.com

Huan Zhang\*†  
University of California, Davis  
Davis, CA  
echezhang@ucdavis.edu

Yash Sharma  
IBM T. J. Watson Research Center  
Yorktown Heights, NY  
Yash.Sharma3@ibm.com

## The Space of Transferable Adversarial Examples

Florian Tramèr<sup>1</sup>, Nicolas Papernot<sup>2</sup>, Ian Goodfellow<sup>3</sup>, Dan Boneh<sup>1</sup>, and Patrick McDaniel<sup>2</sup>

<sup>1</sup>Stanford University, <sup>2</sup>Pennsylvania State University, <sup>3</sup>Google Brain

### Abstract

Adversarial examples are maliciously perturbed inputs designed to mislead machine learning (ML) models at test-time. They often *transfer*: the same adversarial example fools more than one model.

In this work, we propose novel methods for estimating the previously unknown *dimensionality* of the space of adversarial inputs. We find that adversarial examples span a contiguous subspace of large (~25) dimensionality. Adversarial subspaces with higher dimensionality are more likely to intersect. We find that for two different models, a significant fraction of their subspaces is shared, thus enabling transferability.

analysis of the similarity of different models' decision boundaries at these boundaries are actually close in *arbitrary* directions, benign. We conclude by formally studying the *limits* of these findings (1) sufficient conditions on the *data distribution* that transfer occurs for simple model classes and (2) examples of scenarios in which transfer occurs. These findings indicate that it may be possible to design transfer-based attacks, even for models that are vulnerable

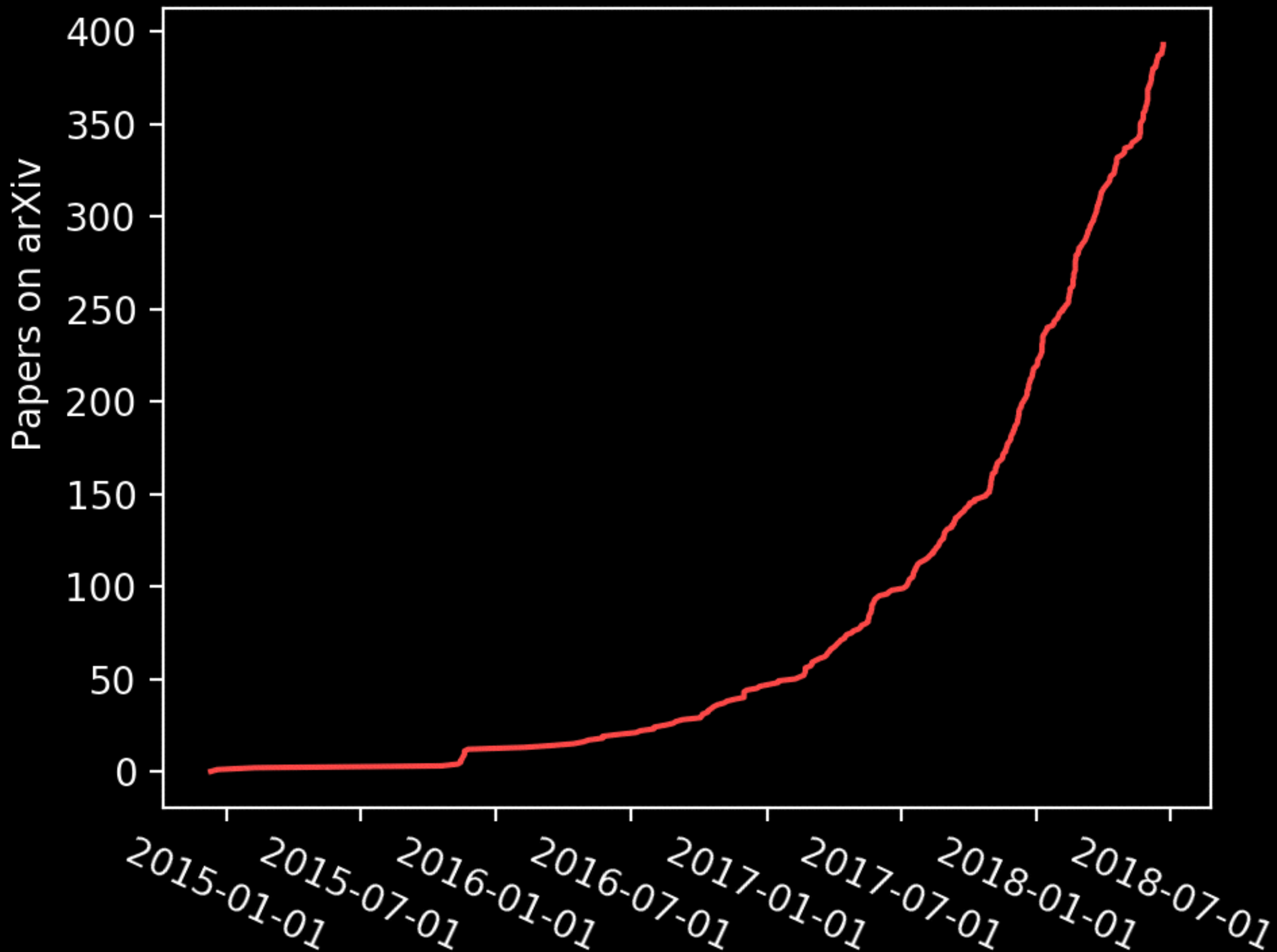
Alhussein Fawzi\*†  
alhussein.fawzi@epfl.ch

Pascal Frossard†  
pascal.frossard@epfl.ch

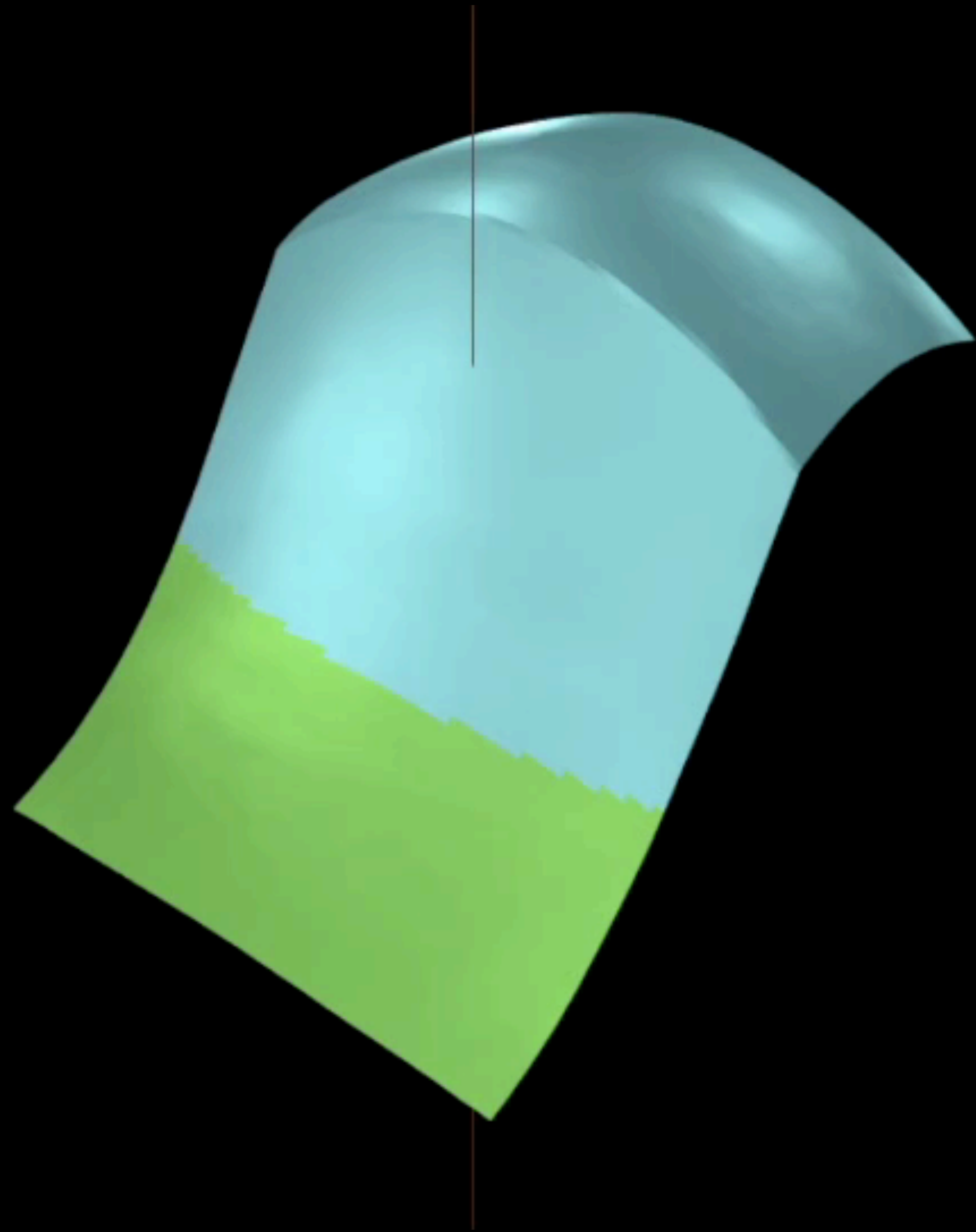
Okay, lesson learned.

Okay, lesson learned.

Give up.

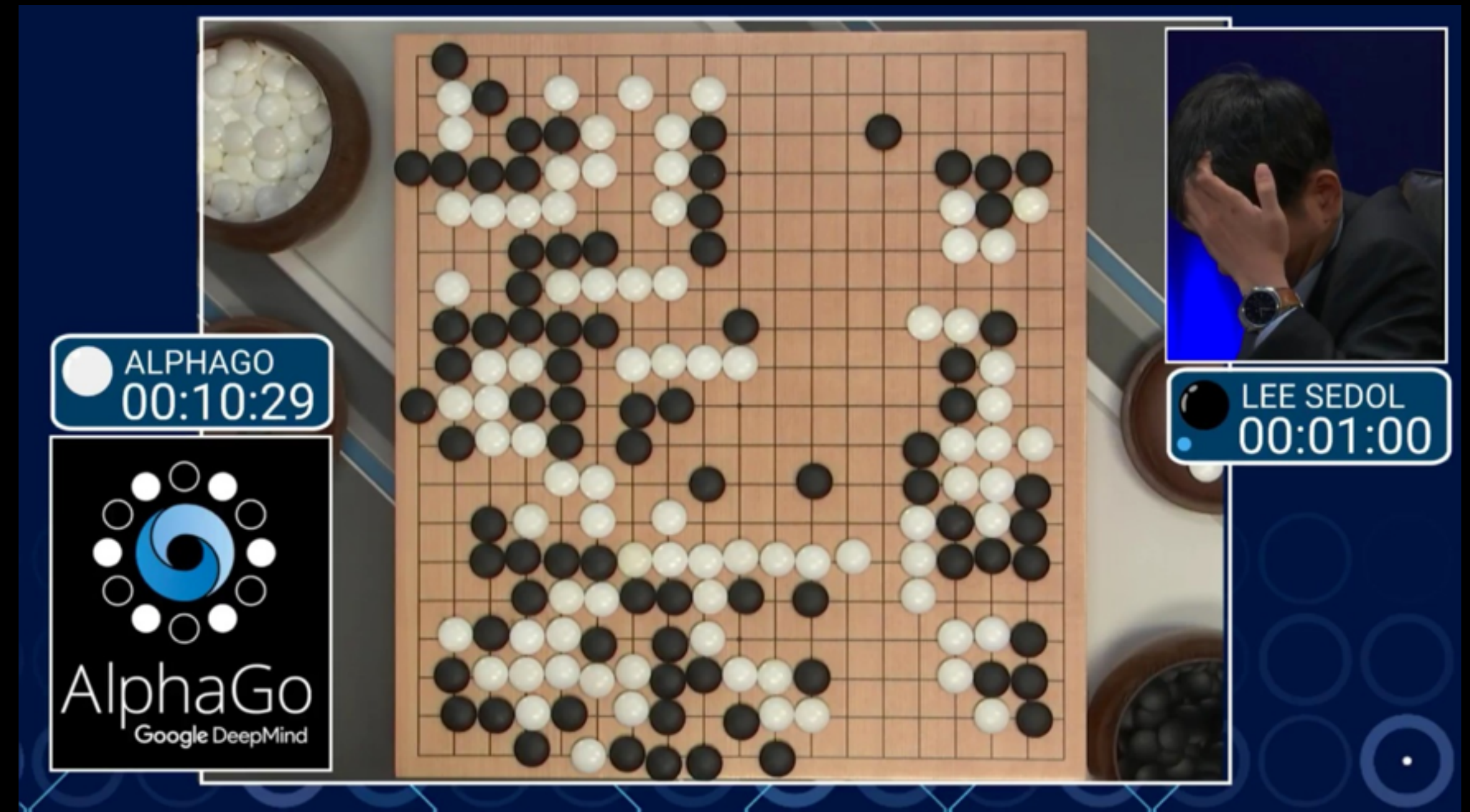








Yes, machine learning gives **amazing** results



However, there are  
also significant  
**vulnerabilities**



Guacamole (99%)

# Questions?

<https://nicholas.carlini.com>  
[nicholas@carlini.com](mailto:nicholas@carlini.com)



