

Slides with GIF:

<https://docs.google.com/presentation/d/e/2PACX-1v0z6blrMWZWibGXsMPNwIJlRE0B9v0zR4FG4y44vS95ShXhD43e0EQQB9UzGJp9JM9bF5L70WwTfjH/pub?start=false&loop=false&delayms=60000>

rethinking adversarial examples

yahya jabary

disco group - andreas plesner, prof. roger wattenhofer

dsg group - alireza furutanpey, prof. schahram dustdar



let's start with a game -
quick, what do you see?



bagel
- or -
dog?

chicken wing
- or -
dog?



A 3x4 grid of images illustrating the visual similarity between dogs and croissants. The images are arranged in three rows and four columns. The first two columns show dogs (Shar Pei breed) standing upright, while the last two columns show croissants. The second row shows a single croissant and a single dog lying down. The third row shows a group of dogs and a group of croissants.

croissant
- or -
dog?



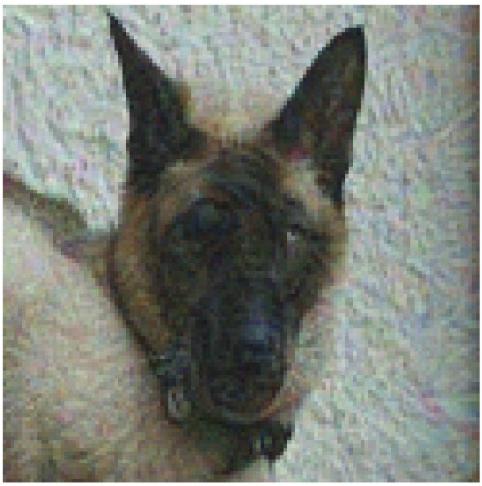
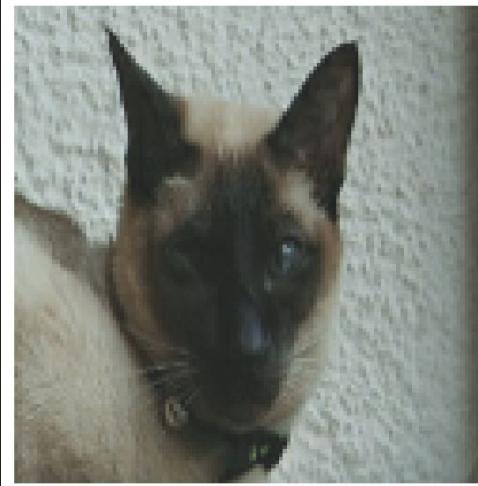
mop
- or -
dog?



muffin
- or -
dog?



1 person
- or -
2 people?

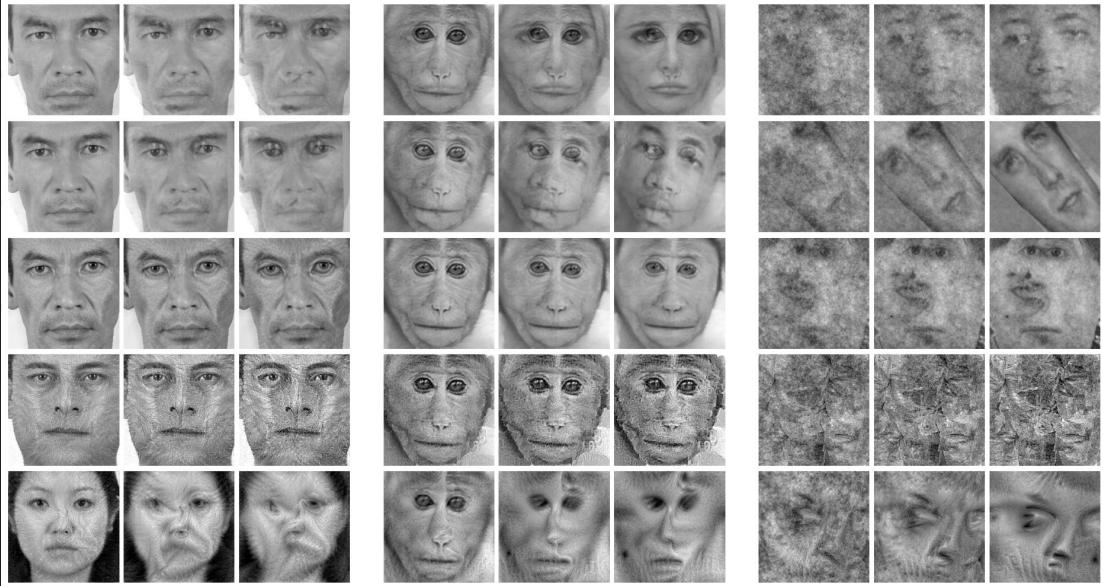


cat

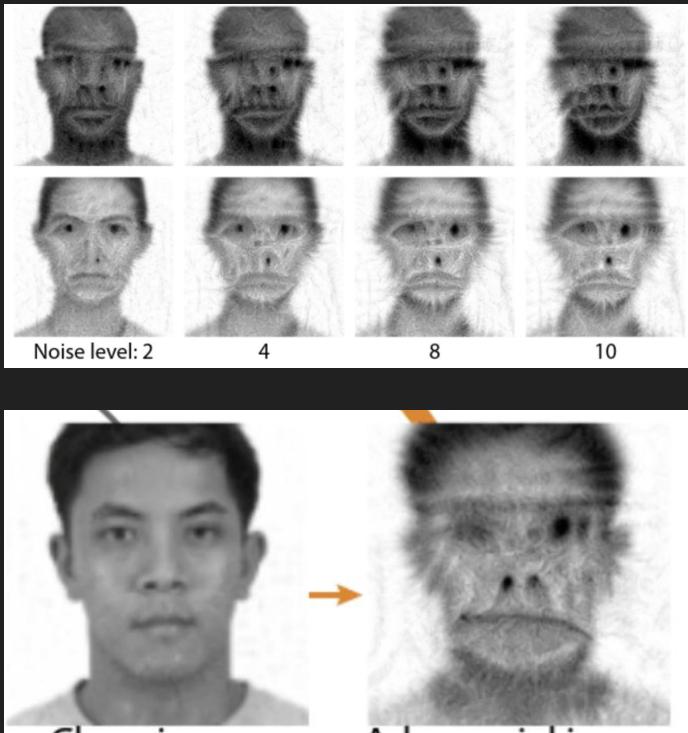
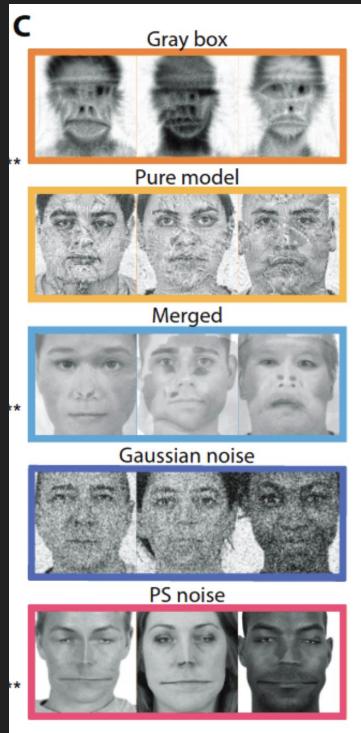
- or -

dog?

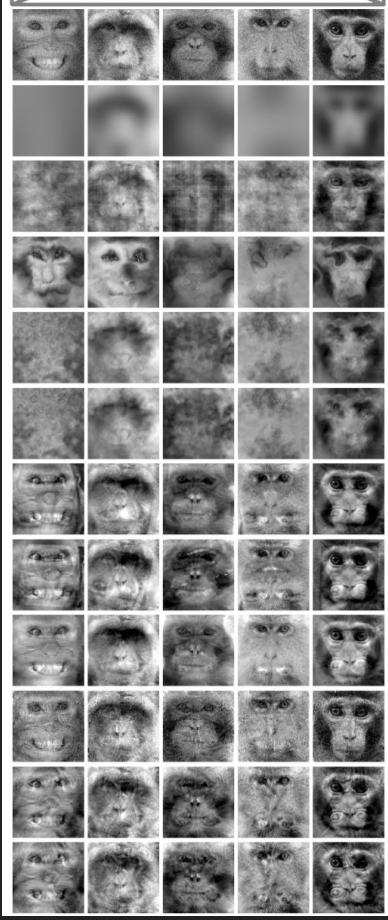




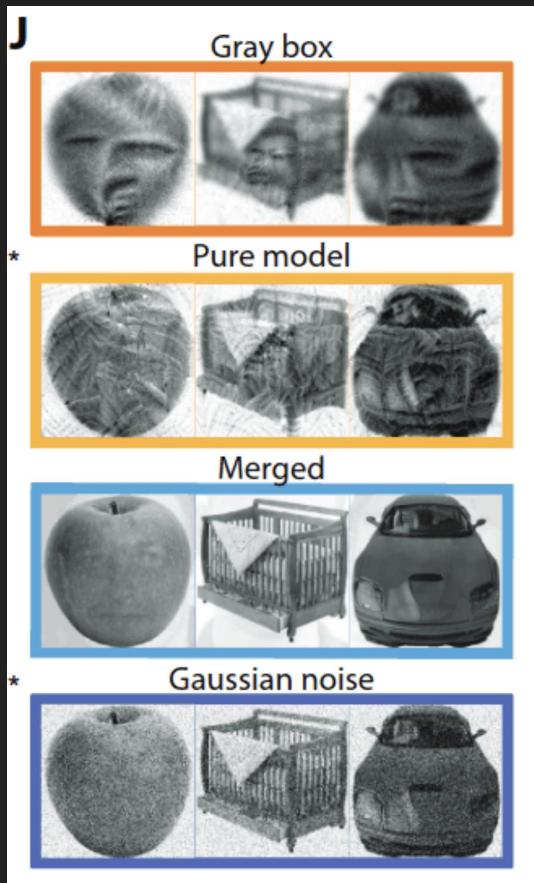
human
- or -
monkey?



human
- or -
monkey?



human
- or -
monkey?



face
- or -
no face?

Classifier Input



Classifier Input



banana

- or -

toaster?

Classifier Input



Classifier Input



too easy right?

not for vision models!

Classifier Input



Classifier Output



we call these "weird" examples
"adversarial" examples.

but you can call them however you want

statistics: out-of-distribution data / adversarial distribution / latent space collapse / information bottleneck failures / covariate shift outliers / data drift (covariate vs. concept) / distribution mismatch / complexity artifacts / overfitting residues / feature attribution misalignment / simpson's paradox / concept drift / black swan events / long tail events / shortcut learning / optimization instabilities

cybersecurity: jailbreaks / gradient-sign evasion attacks / zero-day model vulnerabilities

psychology: cognitive and reasoning biases / perceptual dissonance / deep learning pareidolia / anthropomorphism errors / semantic illusions / confabulation effects / hallucinations / illusory correlations / pattern misperceptions

what are they?

Adversarial Noise



"panda"

+



=



"gibbon"



"panda"

Adversarial Noise



+

=



"gibbon"

but wait - there's more!

traditional
/ imperceptible
(2013)

unrestricted
/ semantics preserving
(2019)

natural
(2019)



"panda"



"vulture"



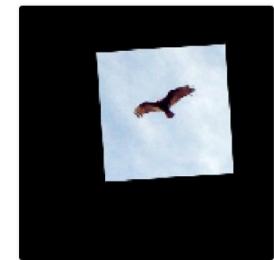
"not hotdog"

Adversarial Noise



"gibbon"

Adversarial Rotation



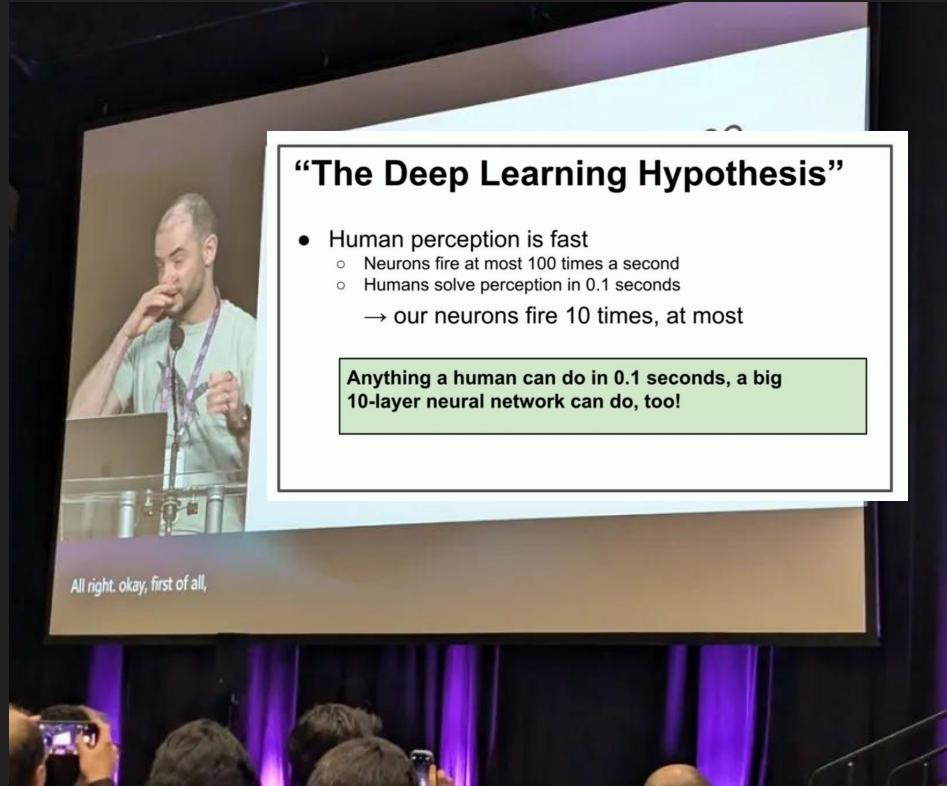
"orangutan"

Adversarial Photographer



"hotdog"

some attacks also work on time-constrained humans & primates!



idea:

neural networks emulate
time-constrained humans

how do they work?

non-robust features hypothesis

feature types:

- a) useless not picked up by feature extractor
 - b) robust detect *concept drift*
generalize to test-set
comprehensible to humans
 - c) non-robust detect *covariate drift*
very predictive, but only in train-set

e.g. the fur direction of an animal

open problem



MrAcurite · 3y ago ·

Researcher

If you could solve this problem conclusively, I'm pretty sure you'd instantly be awarded the Turing prize, or something equivalent. **Shit's hard.**



40

Reply

Award

Share

...

•

why should i care?
use-cases

hack the reviewer-assignment system

No more Reviewer #2: Subverting Automatic Paper-Reviewer Assignment using Adversarial Learning

Thorsten Eisenhofer^{#†}, Erwin Quiring^{*‡}, Jonas Möller[§], Doreen Riepel[†],
Thorsten Holz[¶], Konrad Rieck[§]
[†]Ruhr University Bochum
[‡]International Computer Science Institute (ICSI) Berkeley
[§]Technische Universität Berlin
[¶]CISPA Helmholtz Center for Information Security

Abstract

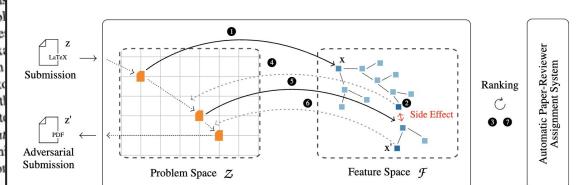
The number of papers submitted to academic conferences is steadily rising in many scientific disciplines. To handle this growth, systems for automatic paper-reviewer assignments are increasingly used during the reviewing process. These systems use statistical topic models to characterize the content of submissions and automate the assignment to reviewers. In this paper, we show that this automation can be manipulated using adversarial learning. We propose an attack that adapts a given paper so that it misleads the assignment and selects its own reviewers. Our attack is based on a novel optimization strategy that alternates between the feature space and problem space to realize unobtrusive changes to the paper. To evaluate the feasibility of our attack, we simulate the paper-reviewer assignment of an actual security conference (IEEE S&P) with 165 reviewers on the program committee. Our results show that we can successfully select and remove reviewers without access to the assignment system. Moreover, we demonstrate that the manipulated papers remain plausible and are often indistinguishable from benign submissions.

1 Introduction

Peer review is a major pillar of academic research and the publication process. Despite its well-known weaknesses, it is crucial for ensuring high-quality research.

To handle this growth, conference management tools have become indispensable in peer review. They allow reviewers to bid for submissions and support the program chair to find a good assignment based on a best-effort matching. Unfortunately, even these tools submissions continue to be intractable, as for example major conferences in missions that need to reviewers [31]. For increasingly ext reviewer assignments from machine learn missions, and auto

In this work, we manipulate the assignment. In prior work that focused on bid manipulations and review collusion [24, 57], our attack rests on adversarial learning. In particular, we propose an attack that adapts a given paper so that it misleads the underlying topic model. This enables us to reject and select specific reviewers from the program committee. To reach this goal, we introduce a novel optimization strategy that alternates between the feature space and problem space when adapting a paper. This optimization allows us to preserve the semantics and plausibility of the document, while carefully changing the assignment of reviewers.



hack the topic-model used for NeurIPS & ICML

solve & make better captchas!

24

Breaking reCAPTCHAv2

Andreas Plesner
ETH Zurich, Switzerland
aplesner@ethz.ch

Tobias Vontobel
ETH Zurich, Switzerland
votobias@student.ethz.ch

Roger Wattenhofer
ETH Zurich, Switzerland
wattenhofer@ethz.ch

Abstract—Our work examines the efficacy of employing advanced machine learning methods to solve captchas using Google’s reCAPTCHA v2 system. We evaluate the performance of automated systems in solving captchas generated by YOLO models. Our results show that state-of-the-art machine learning models can solve most captchas with a success rate of up to 90%.

**Seeing Through the Mask:
Rethinking Adversarial Examples for CAPTCHAs**

Yahya Jabary
TU Wien
jabaryyahya@gmail.com

Andreas Plesner*
ETH Zurich
aplesner@ethz.ch

Turhan Kuzhagaliyev
ETH Zurich
kturlan@student.ethz.ch

Roger Wattenhofer
ETH Zurich
wattenhofer@ethz.ch

Abstract

Modern CAPTCHAs rely heavily on vision tasks that are supposedly hard for computers but easy for humans. However, advances in image recognition models pose a significant threat to such CAPTCHAs. These models can easily be fooled by generating some well-hidden “random” noise and adding it to the image, or hiding objects in the image. However, these methods are model-specific and thus can not aid CAPTCHAs in fooling all models. We show in this work that by allowing for more significant changes to the images while preserving the semantic information and keeping it solvable by humans, we can fool many state-of-the-art models. Specifically, we demonstrate that by adding masks of various intensities and points. These masks can therefore effectively fool modern image classifiers, thus machines have not caught up with humans – yet.

google's reCAPTCHA v2

- 99% market share
- 100% solve rate

The screenshot shows the homepage of Campus Technology. The header features the magazine's name in large, bold, orange and white letters, with the tagline "EMPOWERING THE WORLD OF HIGHER EDUCATION" below it. A navigation bar below the header includes links for NEWS, FEATURES, OPINION, RESEARCH, PODCASTS, and RESOURCES. The main content area displays a news article titled "AI Beats reCAPTCHA" under the "TECHNOLOGY NEWS" section. The article's thumbnail image shows a computer screen with a CAPTCHA challenge. At the bottom of the page, there are social media sharing icons for Facebook and LinkedIn, along with a "Watch" button and a "Starred" section showing 305 items.

A screenshot of a web browser displaying a news article from Techxplore. The title of the article is "Humans are not needed, have AI breakthrough all Google reCaptcha's (code included)". The page includes a navigation bar with links for Home, Security, and Topics, as well as a sidebar with social sharing options and a sidebar for the user Abish Plus.

The screenshot shows a news article from TechDogs. The header features the TechDogs logo and the text "NEW Research CAPTCHA". Below the header, it says "By David Ramel | 09/2023". The main title of the article is "KI löst Bilder-Captchas besser". The article discusses how AI can solve Google's reCAPTCHA better than humans. It includes a quote from Andreas Plesner and Roger Wattenhofer from ETH Zurich. The article is part of a section called "inspire".

The image is a composite of several screenshots. At the top left, there's a snippet of a news article with the headline 'AI bots now beat 100% of those traffic-image CAPTCHAs'. The main part of the image shows a screenshot of a website with a navigation bar that includes 'Branded Insights' and 'Events & Webi'. Below the navigation, there's a large heading 'AI bots now beat 100% of those traffic-image CAPTCHAs' followed by a subtext 'I, for one, welcome our traffic light-identifying overlords.' On the right side, there's a CAPTCHA challenge titled 'Select all images with crosswalks' with several examples of street scenes. At the bottom left, there's another snippet of a news article with the headline 'AI Tests All The Time'.

- The researchers found that they could hoodwink many state-of-the-art models by adding hidden "red-hat" noise and adding it to the image or hiding objects in the image.
- This technique could help in dropping accuracy levels by 50%-points for all "supposedly robust models such as vision transformers".

rise in online bots, technologies to help combat them. A technology includes its own by its acronym. A tests commonly using one or several random numbers mixed in a sequence. It was to tell if the code was not effective (

AI Can Crack CAPTCHA with 100% Accuracy

Reading time: 2 min Updated on Sep 24, 2024

Written by Kira Fabbri Multimedia Journalist

Fact-Checked by Justin Newman Lead Cybersecurity Editor

In a Rush? Here are the Quick Facts!

- AI can now solve CAPTCHAs with 100% accuracy.
- reCAPTCHA v2 relies heavily on cookies and browser history.

Published on September 24, 2024 by Roger Wattenberg

Smartwatch Might Be Poisoning You

Android Java

Gemini on Android Can Now Perform Multiple Actions At Once

Apple Siri

Sam Altman Downplays AGI Risks; Now Warns About Superintelligence

Apple Siri

From Brooklyn 99 to the News Doc: My Year Revealed by Unofficial Netflix Wrapped

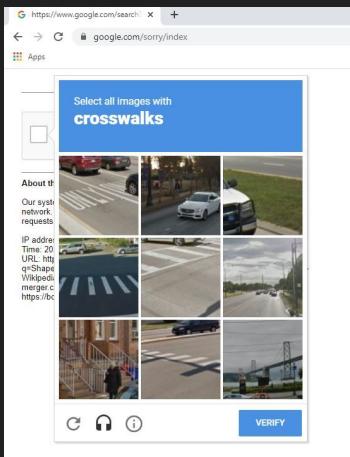
Android Java

Microsoft Says It Didn't Use

arms race

robustified model

challenger



prover / solver



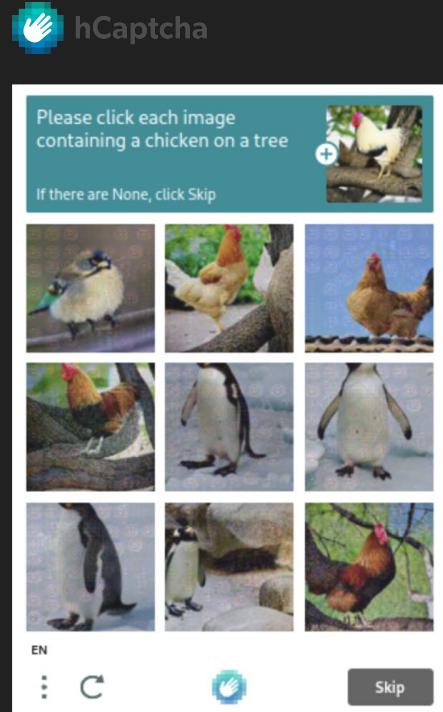
adversarial attack

step 1: making harder challenges

hCaptchas are surprisingly hard to solve!

simplicity enables formal reasoning:

- tracing changes in latent space
- linear increase of "attack strength" (opacity)



steps

1. remake the geometric masks

original:



reconstruction:



steps

1. remake the geometric masks
2. find the best of each kind

hyperparam search using a single model

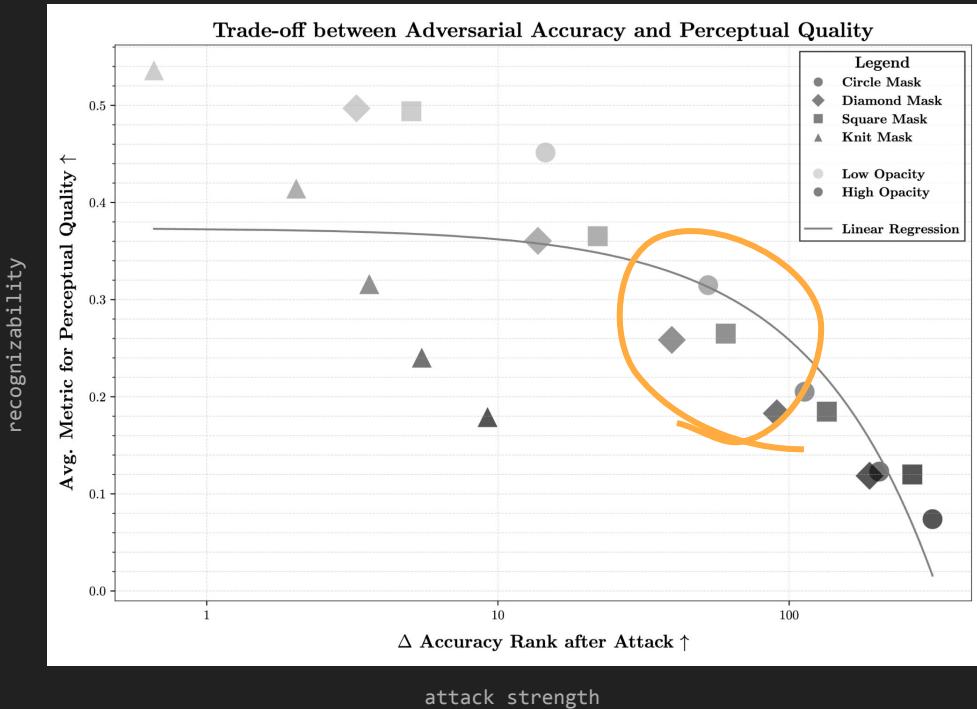
metrics:

- (1) attack strength
= drop of correct label rank (acc@1, acc@5)
- (2) recognizability
= cosine similarity, PSNR, SSIM, LPIPS

we want both to be high!

steps

1. remake the geometric masks
2. find the best of each kind
3. eval them on SOTA models



steps

1. remake the geometric masks
2. find the best of each kind
3. eval them on SOTA models

we confirmed that these masks are:

- ✓ simple
- ✓ functional
- ✓ transferable

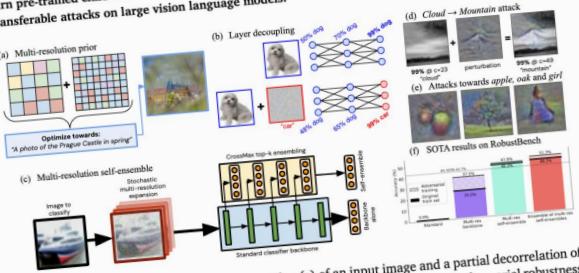
Opacity	Mask	inverse attack strength		recognizability
		Δ Acc	Rank	
50	Circle	-14.57	0.45	15.02
	Diamond	-3.27	0.50	3.76
	Knit	-0.66	0.54	1.19
	Square	-5.04	0.49	5.54
80	Circle	-52.72	0.31	53.03
	Diamond	-13.72	0.36	14.08
	Knit	-2.03	0.41	2.44
	Square	-22.01	0.37	22.37
110	Circle	-113.07	0.21	113.27
	Diamond	-39.55	0.26	39.81
	Knit	-3.62	0.32	3.93
	Square	-60.57	0.27	60.84
140	Circle	-203.89	0.12	204.01
	Diamond	-90.79	0.18	90.97
	Knit	-5.47	0.24	5.71
	Square	-134.75	0.18	134.94
170	Circle	-310.80	0.07	310.88
	Diamond	-188.92	0.12	189.04
	Knit	-9.21	0.18	9.39
	Square	-264.90	0.12	265.02

step 2: making better solvers

Ensemble everything everywhere: Multi-scale aggregation for adversarial robustness

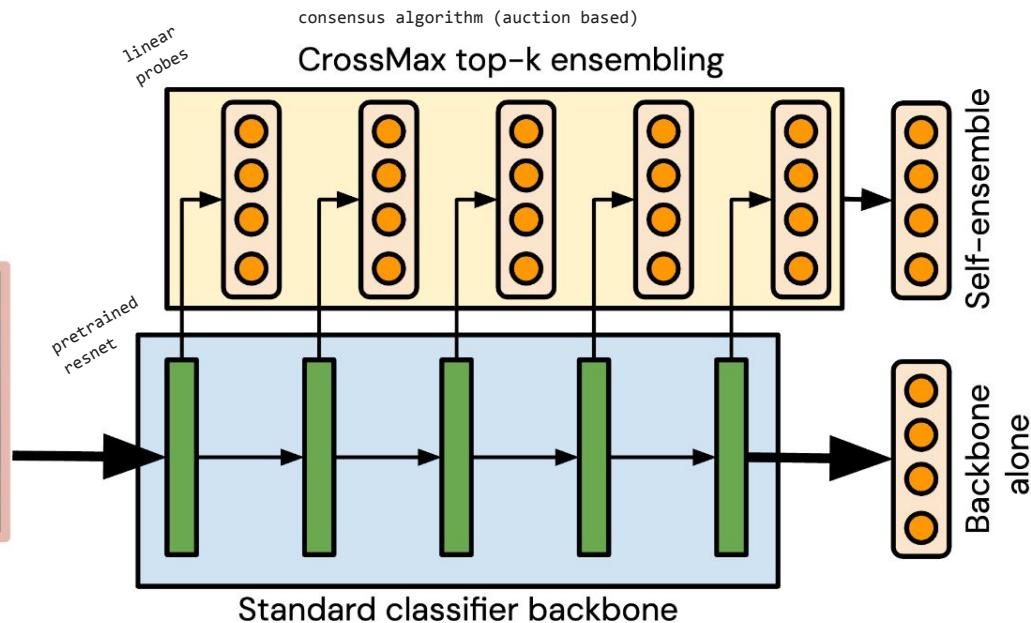
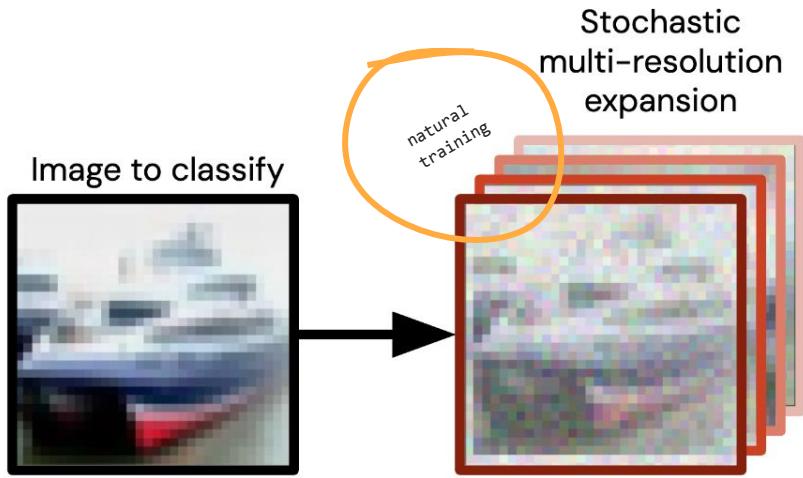
Stanislav Fort¹ and Balaji Lakshminarayanan¹
¹Google DeepMind

Adversarial examples pose a significant challenge to the robustness, reliability and alignment of deep neural networks. We propose a novel, easy-to-use approach to achieving high-quality representations that lead to adversarial robustness through the use of multi-resolution input representations and dynamic self-ensembling of intermediate layer predictions. We demonstrate that intermediate layer predictions exhibit inherent robustness to adversarial attacks crafted to fool the full classifier, and propose a robust aggregation mechanism based on Vickrey auction that we call CrossMax to dynamically ensemble them. By combining multi-resolution inputs and robust ensembling, we achieve significant adversarial robustness on CIFAR-10 and CIFAR-100 datasets without any adversarial training or extra data, reaching an adversarial accuracy of $\approx 72\%$ (CIFAR-10) and $\approx 48\%$ (CIFAR-100) on the RobustBench AutoAttack suite ($L_\infty = 8/255$) with a finetuned ImageNet-pretrained ResNet152. This represents a result comparable with the top three models on CIFAR-10 and a $+5\%$ gain compared to the best current dedicated approach on CIFAR-100. Adding simple adversarial training on top, we get $\approx 78\%$ on CIFAR-10 and $\approx 51\%$ on CIFAR-100, improving SOTA by 5% and 9% respectively and seeing greater gains on the harder dataset. We validate our approach through extensive experiments and provide insights into the interplay between adversarial robustness, and the hierarchical nature of deep representations. We show that simple gradient-based attacks against our model lead to human-interpretable images of the target classes as well as interpretable image changes. As a byproduct, using our multi-resolution prior, we turn pre-trained classifiers and CLIP models into controllable image generators and develop successful transferable attacks on large vision language models.



self-enssembled resnet:

- beats SOTA in robustness by 5%
- increased interpretability
- robustness through architecture, just not adversarial training



multires, noise, jitter, shuffle, light fgsm
(we found out: work best when combined)

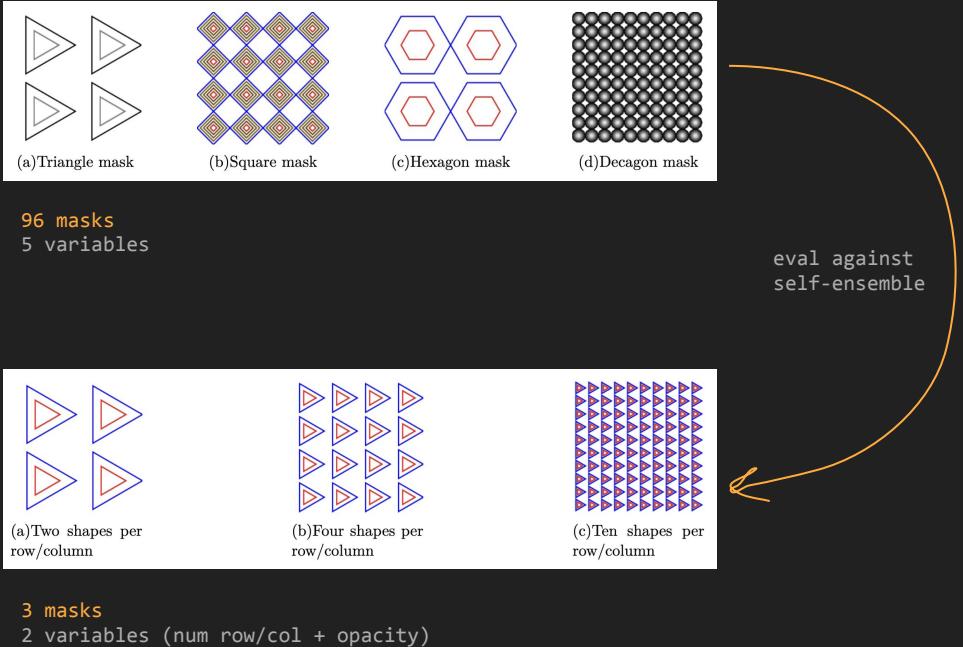
steps

1. choose dataset → CIFAR10

steps

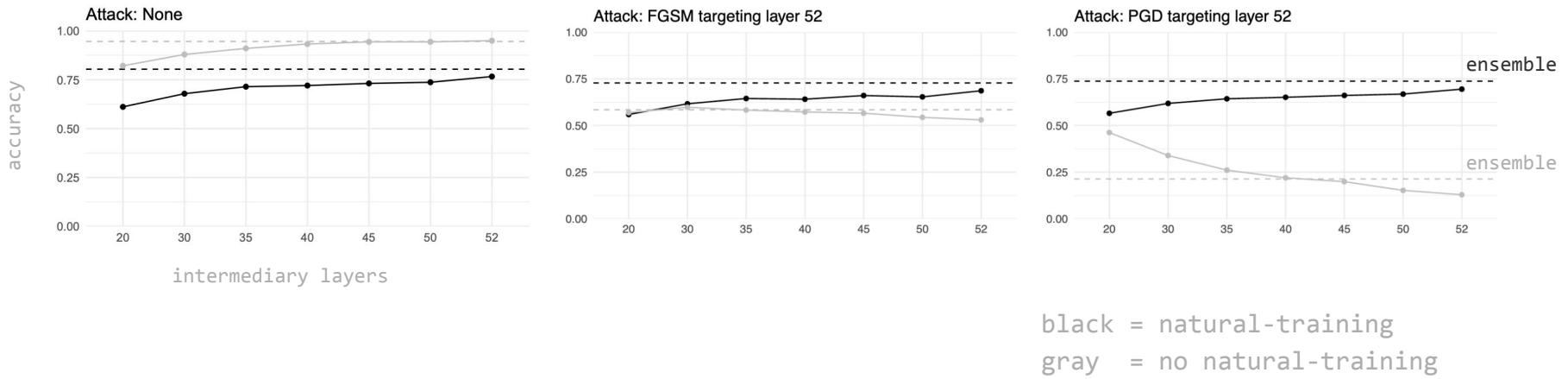
1. choose dataset
2. choose attacks

- 3x geometric masks
- fgsm attack
- pgd attack



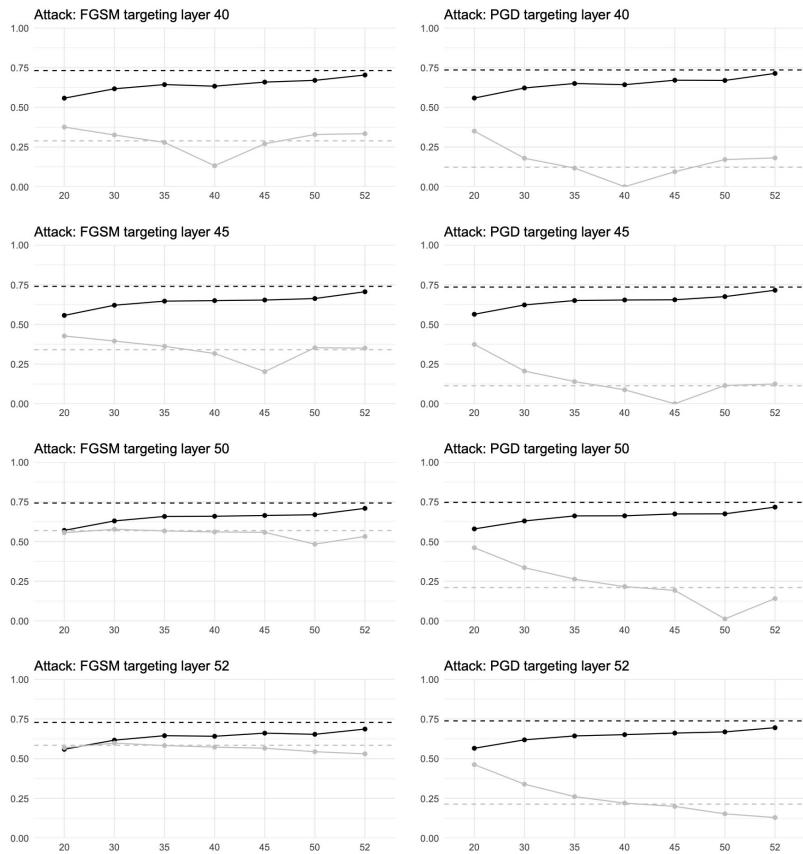
steps

1. choose dataset
2. choose attacks
3. eval



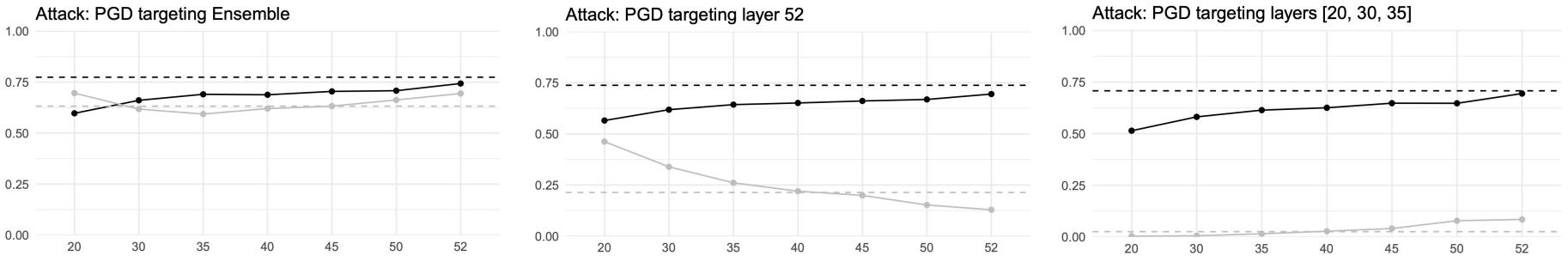
findings:

- natural-training > just the architecture
- natural-training worsens accuracy without a threat ("robustness-accuracy-tradeoff")
- PGD > FGSM



findings:

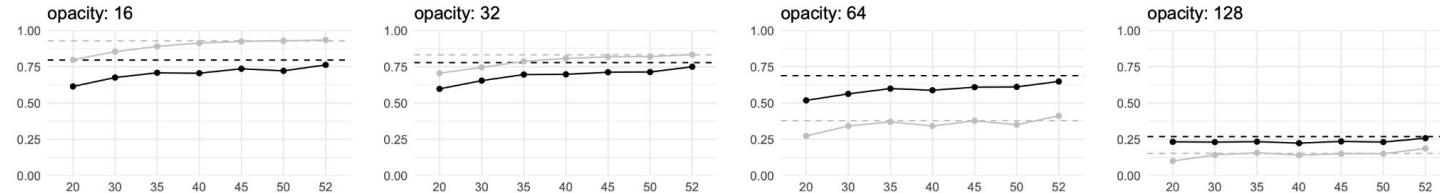
- FGSM loses strength towards final layers
- PGD does not (as much)



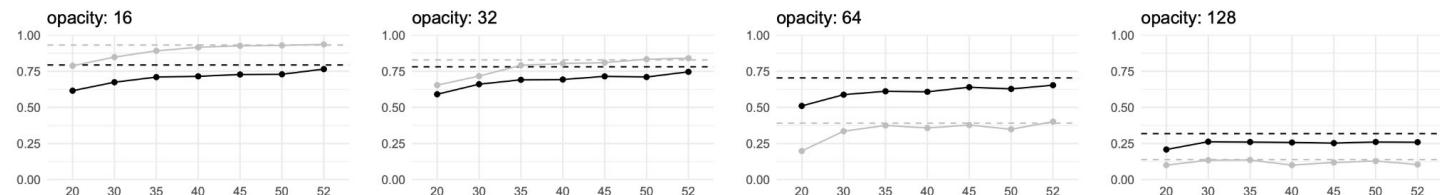
findings:

- targeting mid layers > targeting final layer > targeting ensemble

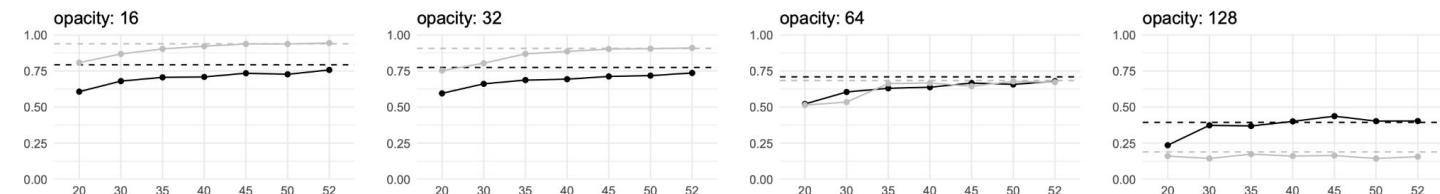
Attack: Geometric Mask
(3 sides, 2 per row/col, 2 concentric shapes, colors enabled)



(3 sides, 4 per row/col, 2 concentric shapes, colors enabled)



(3 sides, 10 per row/col, 2 concentric shapes, colors enabled)



findings:

- opacity \approx attack strength
- masks target all layers at once
- natural-training is less effective for masks (than FGSM / PGD)

what did we learn?

key takeaways

geometric masks:

- **good:** simple, allows formal reasoning
- **bad:** not "natural"

self-ensemble:

- **good:** robustness through architecture
- **bad:** but still needs natural training for SOTA

what could be next?

outlook

attacks: natural attacks

- generative models → synthetic data
- NeRF models → 3D data

defenses: robust architectures

- fermi-bose machine architecture
- depth estimation / NeRF model (data augmentation)
- map latent shifts to concept drifts (mechanistic interpretability)
- mutual information / MINE (entropy for intermediary layers)

thank you!

github.com/sueszli/thesis