# What's next for adversarial ML?

## (and why ad-blockers should care)

**Florian Tramèr**

EPFL

July 9th 2018

Joint work with Gili Rusak, Giancarlo Pellegrino and Dan Boneh

# The Deep Learning Revolution

First they came for images…

# The Deep Learning Revolution

And then everything else…



nature
International journal of science

Article | Published: 18 October 2017

Mast
know

The Download
What's up in emerging technology

November 16, 2017

Andrew Ng @AndrewYNg
Follow

Pretty much anything that a normal person can do in <1 sec, we can now automate with AI.

A New Algorithm Can Spot Pneumonia Better Than a Radiologist

Add diagnosing dangerous lung diseases to the growing list of things artificial intelligence can do better than humans.

# The ML Revolution

Including things that likely won't work…

# What does this mean for privacy & security?



Crypto, Trusted hardware

**Privacy & integrity**

Outsourced learning

Differential privacy

**Data inference Model theft** [**T**ZJRR16]

Test outputs

Training data

**Data poisoning**

Robust statistics

B ??? n

**Adversarial Examples**

Test data

dog cat bird

Outsourced inference

**Privacy & integrity**

Crypto, Trusted hardware
⇒ Check out Slalom! [**T**B18]

Adapted from (Goodfellow 2018)

# What does this mean for privacy & security?

Crypto, Trusted hardware

**Privacy & integrity**

Outsourced learning

Differential privacy

**Data inference**
**Model theft** [TZJRR16]

Test outputs

Training data

**Data poisoning**

Robust statistics

???

**Adversarial Examples**

Test data

dog cat bird

Outsourced inference

**Privacy & integrity**

Crypto, Trusted hardware
⇒ Check out Slalom! [TB18]

Adapted from (Goodfellow 2018)

# ML models make surprising mistakes

$$+ .007 \times$$

$$=$$

**Pretty sure this
is a panda**

**I'm certain this
is a gibbon**

(Szegedy et al. 2013, Goodfellow et al. 2015)

# Attacks on cyber-physical systems



(Sharif et al. 2016)

(Kurakin et al. 2016)

(Athalye et al. 2018)

Hi, how can I help?

(Carlini et al. 2016,
Cisse et al. 2017)

(Eykholt et al. 2017)

(Eykholt **et al.** 2018)

# Where are the defenses?

- ## Adversarial training
  Szegedy et al. 2013, Goodfellow et al. 2015, Kurakin et al. 2016, **T** et al. 2017, Madry et al. 2017, Kannan et al. 2018

  Prevent "all/most attacks" **for a given norm ball**

- ## Convex relaxations with provable guarantees
  Raghunathan et al. 2018, Kolter & Wong 2018, Sinha et al. 2018

- ## A lot of broken defenses…

**Adversarial Examples Are Not Easily Detected: Bypassing Ten Detection Methods**

Nicholas Carlini      David Wagner

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

Anish Athalye [*1]   Nicholas Carlini [*2]   David Wagner [2]

# Do we have a realistic threat model? (no…)

Current approach:
1. Fix a "toy" attack model (e.g., some $l_\infty$ ball)
2. Directly optimize over the robustness measure
   - $\Rightarrow$ Defenses do not generalize to other attack models
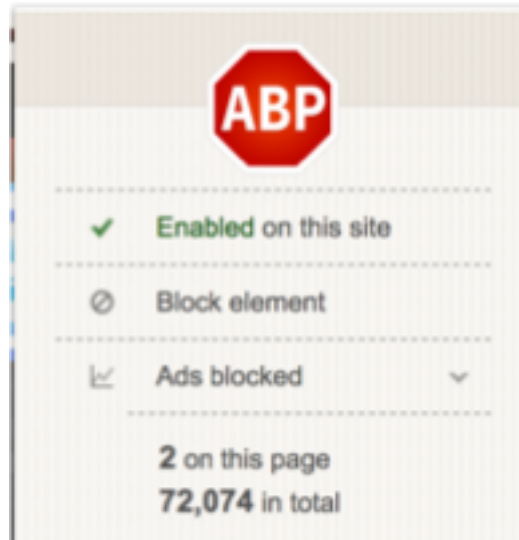   - $\Rightarrow$ Defenses are meaningless for applied security

What do we want?
- Model is "always correct" (sure, why not?)
- Model has blind spots that are "hard to find"
  - "Non-information-theoretic" notions of robustness?
  - CAPTCHA threat model is interesting to think about

# ADVERSARIAL EXAMPLES ARE HERE TO STAY!

For many things that humans can do "robustly", ML will fail miserably!

# A case study on ad-blocking



Ad blocking is a "cat & mouse" game
1. Ad blockers build crowd-sourced filter lists
2. Ad providers switch origins
3. Rinse & repeat
(4?) Content provider (e.g., Cloudflare) hosts the ads

# A case study on ad-blocking

## New method: perceptual ad-blocking (Storey et al. 2017)

- Industry/legal trend: ads have to be clearly indicated to humans

AdChoices ▶  E **The Economist** Sponsored · 🌐

**If humans can detect ads, so can ML!**

"[…] **we deliberately ignore all signals invisible to humans**, including URLs and markup. Instead we consider visual and behavioral information. […] **We expect perceptual ad blocking to be less prone to an "arms race."**

(Storey et al. 2017)

ADVERTISING
SPONSORED
**DETECTED**

**Meet Sentinel**
the artificial intelligence ad detector.

With your help, Sentinel could be the future of ad blocking.

Sentinel uses machine learning to detect Facebook ads visually. The more Facebook screenshots you submit, the faster Sentinel will learn.

Team up with Sentinel for the future of ad blocking!

**FEED SENTINEL**

# Detecting ad logos is not trivial

## No strict guidelines, or only loosely followed:



**Fuzzy hashing + OCR** (Storey et al. 2017)

$\Rightarrow$ Fuzzy hashing is very brittle (e.g., shift all pixels by 1)

$\Rightarrow$ OCR has adversarial examples (Song & Shmatikov, 2018)

**Unsupervised feature detector (SIFT)**

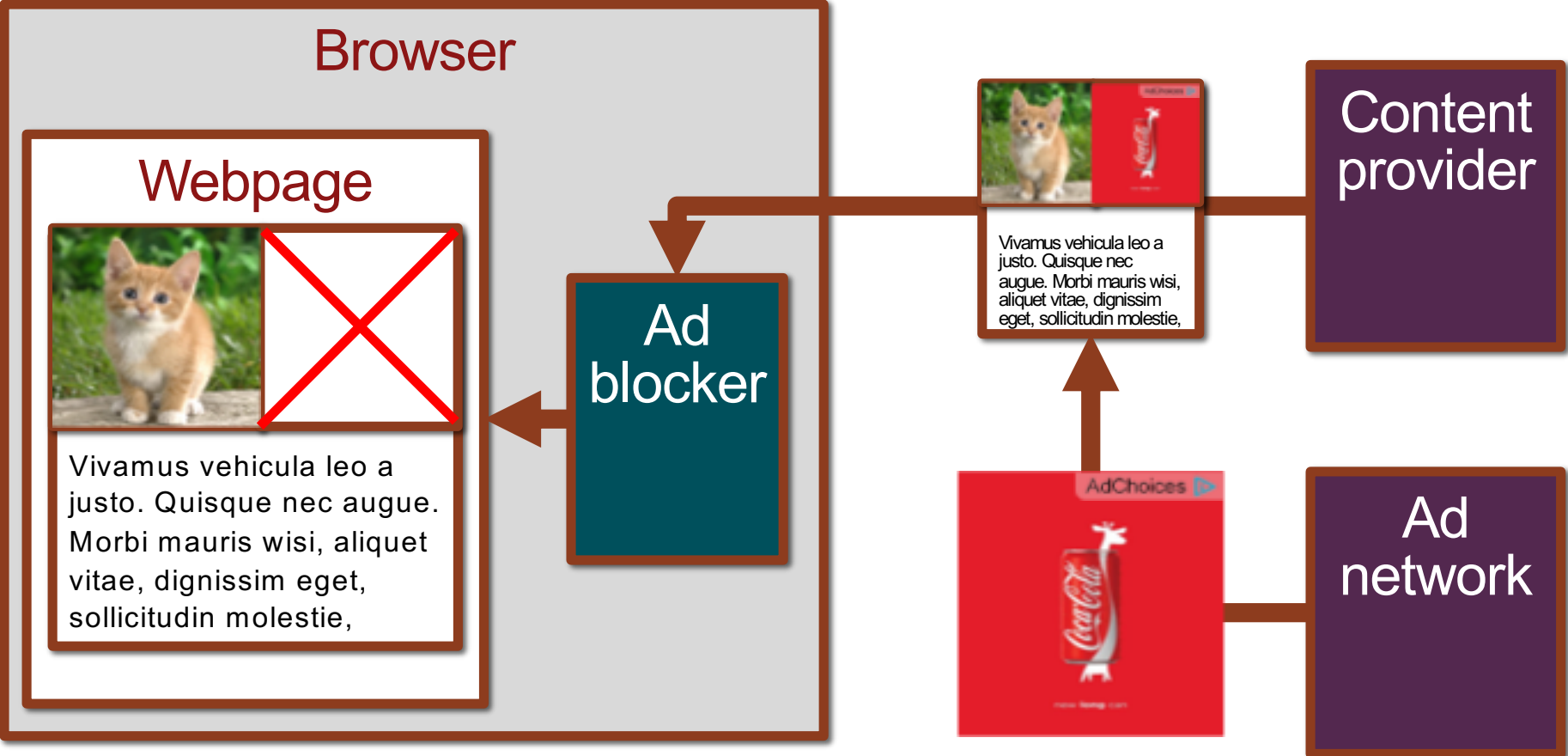$\Rightarrow$ More robust method for matching object features ("keypoints")

**Deep object detector (YOLO)**

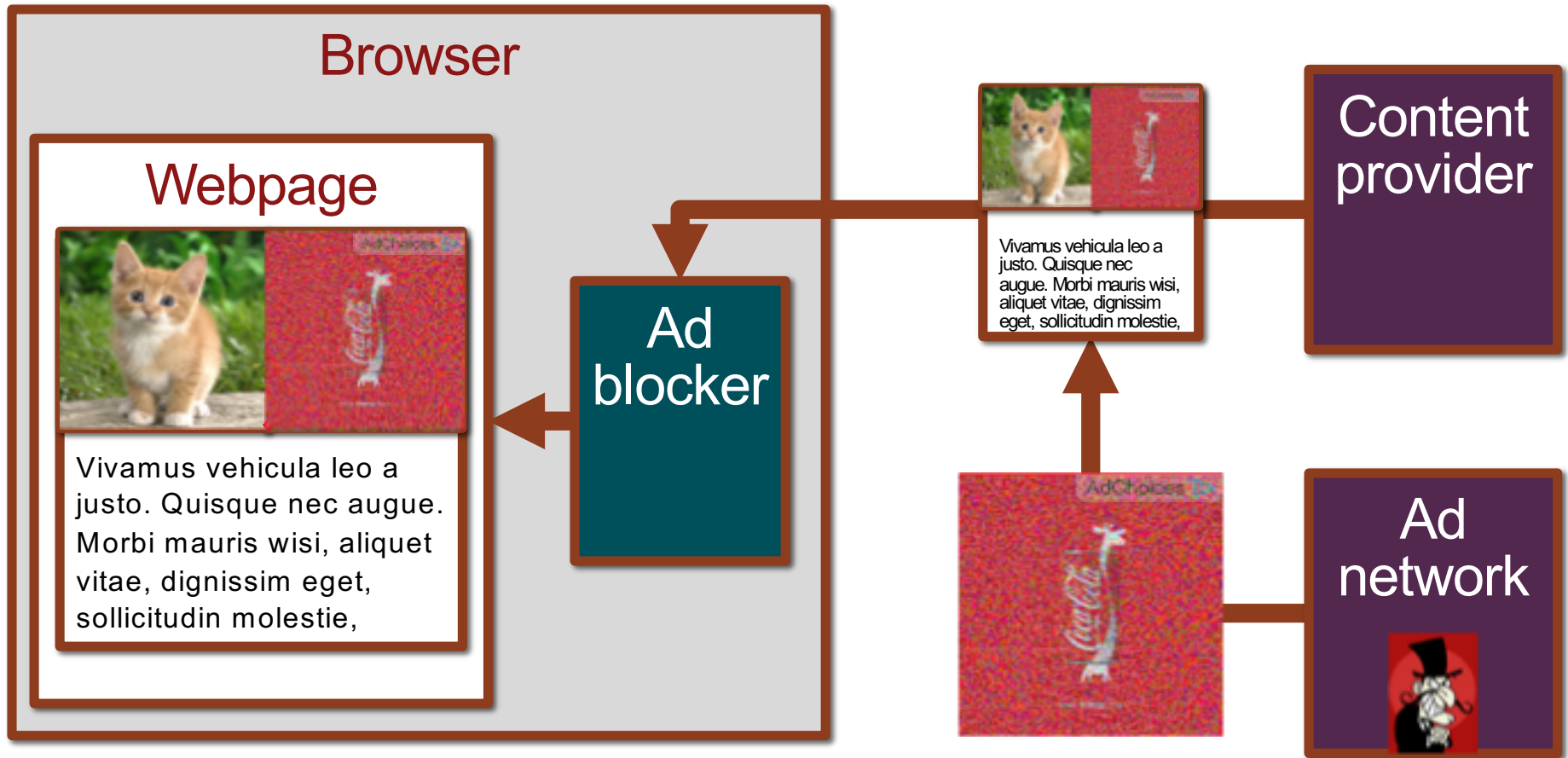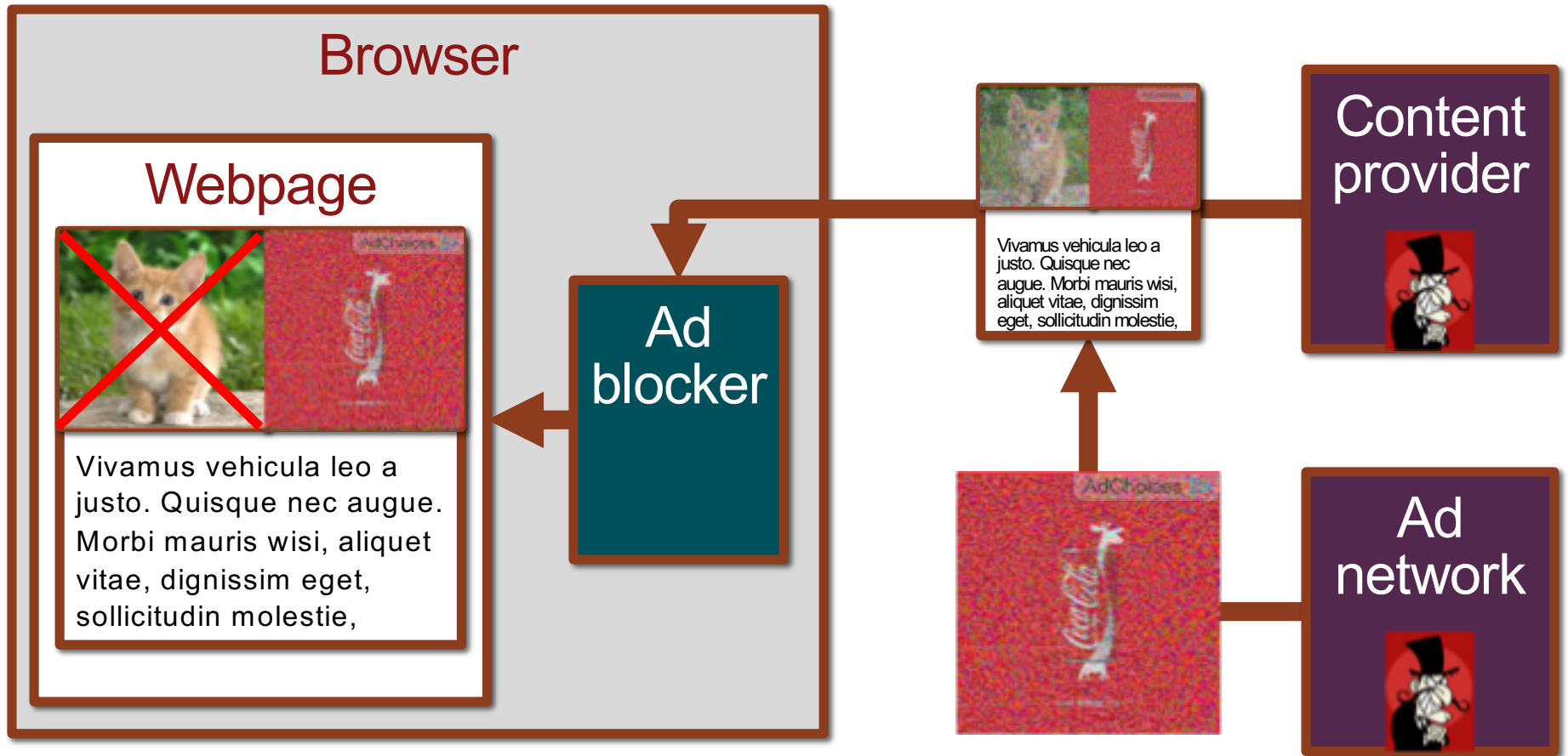$\Rightarrow$ Supervised learning

This talk

# What's the threat model for perceptual ad-blockers?

# What's the threat model for perceptual ad-blockers?

# What's the threat model for perceptual ad-blockers?

# What's the threat model for perceptual ad-blockers?

Pretty much the worst possible!

1. **Adblocker is white-box** (browser extension)

   $\Rightarrow$ Alternative would be a privacy & bandwidth nightmare

2. **Adblocker operates on (large) digital images**

3. **Adblocker needs to resist adversarial examples and "DOS" attacks**

   $\Rightarrow$ Perturb ads to evade ad blocker

   $\Rightarrow$ Punish ad-block users by perturbing benign content

4. **Updating is more expensive than attacking**
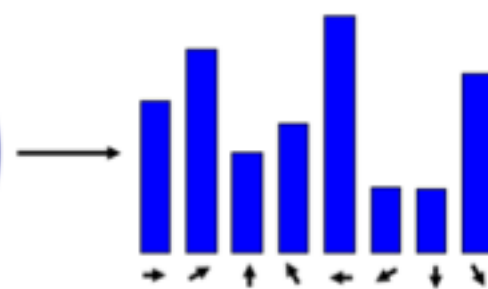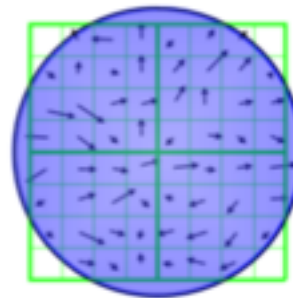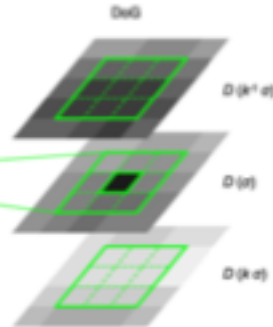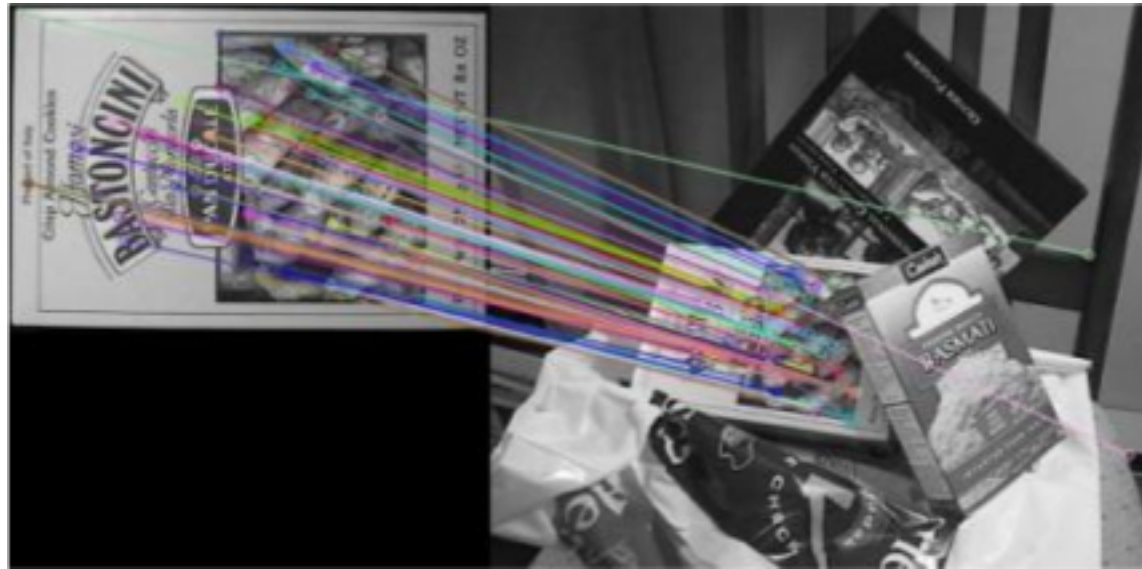
# An interesting contrast: CAPTCHAs



Deep ML models can solve text CAPTCHAs

$\Rightarrow$ Why don't CAPTCHAs use adversarial examples?
$\Rightarrow$ CAPTCHA $\simeq$ adversarial example for OCR systems

| | Model access | Vulnerable to DOS | Model Distribution |
|---|---|---|---|
| Ad blocker | White-box | Yes | Expensive |
| CAPTCHA | "Black-box" (not even query access) | No | Cheap (None) |

# BREAKING PERCEPTUAL AD-BLOCKERS WITH ADVERSARIAL EXAMPLES

# SIFT: How does it work? (I don't know exactly either)
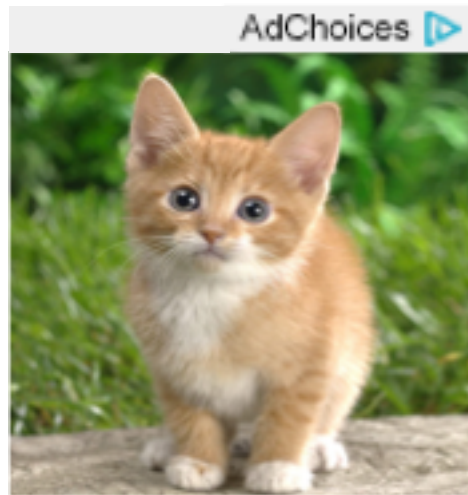
# Attack examples: SIFT detector



original ad



perturbed logo

- No keypoint matches between the two logos
- Attack uses standard black-box optimization
  $\Rightarrow$ Gradient descent with black-box gradient estimates
  $\Rightarrow$ There's surely more efficient attacks but SIFT is complicated…

# Attack examples: SIFT Denial Of Service



- Logos are similar in gray scale but not in color space



- Alternative: high confidence matches for visually close —yet semantically different—objects
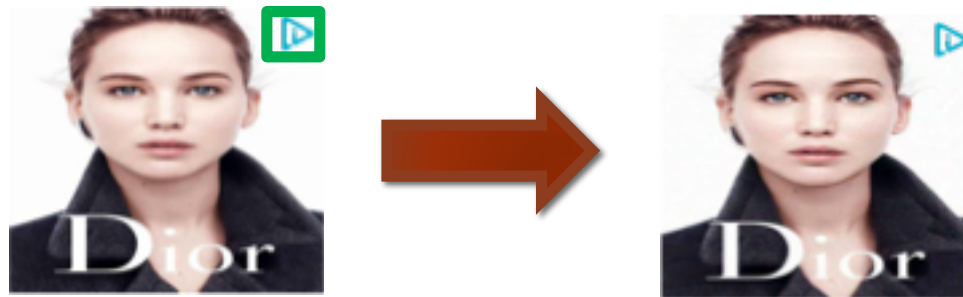
# Attack examples: YOLO object detector

Object detector trained to recognize AdChoice logo

$\Rightarrow$ Test accuracy is >90%

$\Rightarrow$ 0% accuracy with l∞ perturbations ≤ 8/256



Similar but simpler task than Sentinel (Adblock Plus)

$\Rightarrow$ Sentinel tries to detect ads in a whole webpage
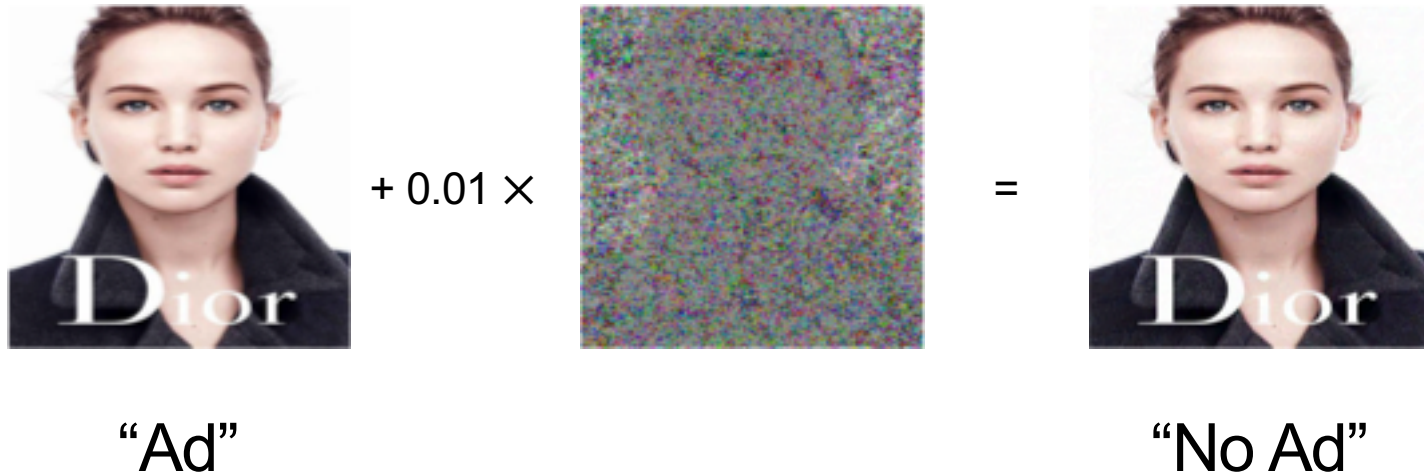
$\Rightarrow$ For now, it breaks even on non-adversarial inputs…

# Perceptual ad-blockers without ad-indicators

Hussain et al. 2017: Train a generic ad/no-ad classifier (for sentiment analysis)

$\Rightarrow$ Accuracy around 88% !

$\Rightarrow$ 0% accuracy with $l_\infty$ perturbations $\leq 4/256$



"Ad"          + 0.01 ×          =          "No Ad"

# Conclusion

Adversarial examples are here to stay

- No defense can address realistic attacks
- A truly robust defense likely implies a huge breakthrough in non-secure ML as well

Security-sensitive ML seems hopeless if adversary has white-box model access

- Ad-blocking ticks most of the "worst-case" boxes
- ML is unlikely to change the ad-blocker cat & mouse game

THANKS