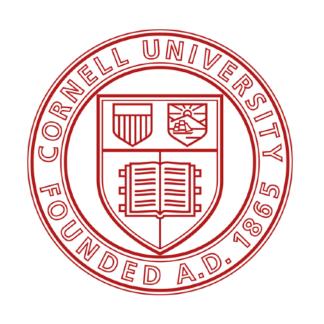
Lecture 25: Bias and fake image detection

CS 5670: Introduction to Computer Vision



Today

Bias in computer vision

Detecting fake images

Garbage in, garbage out

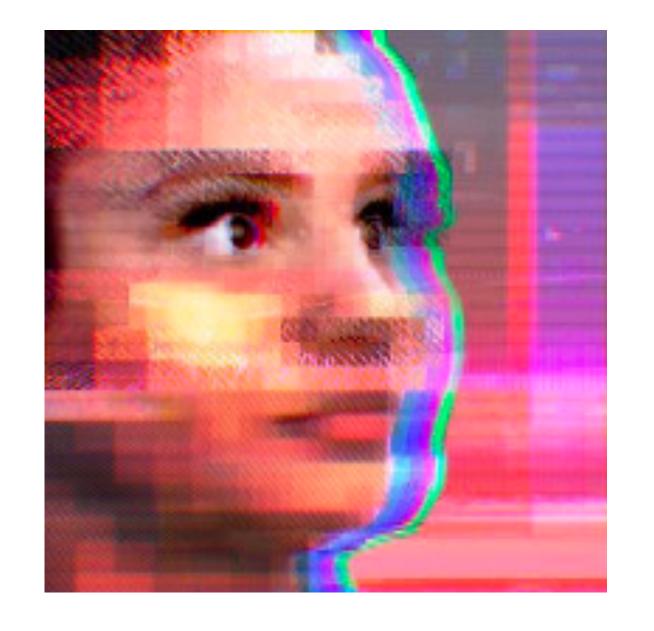
A machine learning algorithm will do whatever the training data tells it to do.

If the data is bad or biased, the learned algorithm will be too.

Microsoft's Tay chatbot

Chatbot released on twitter.

Learned from interactions with users





Started mimicking offensive language, was shut down.

Sydney chatbot

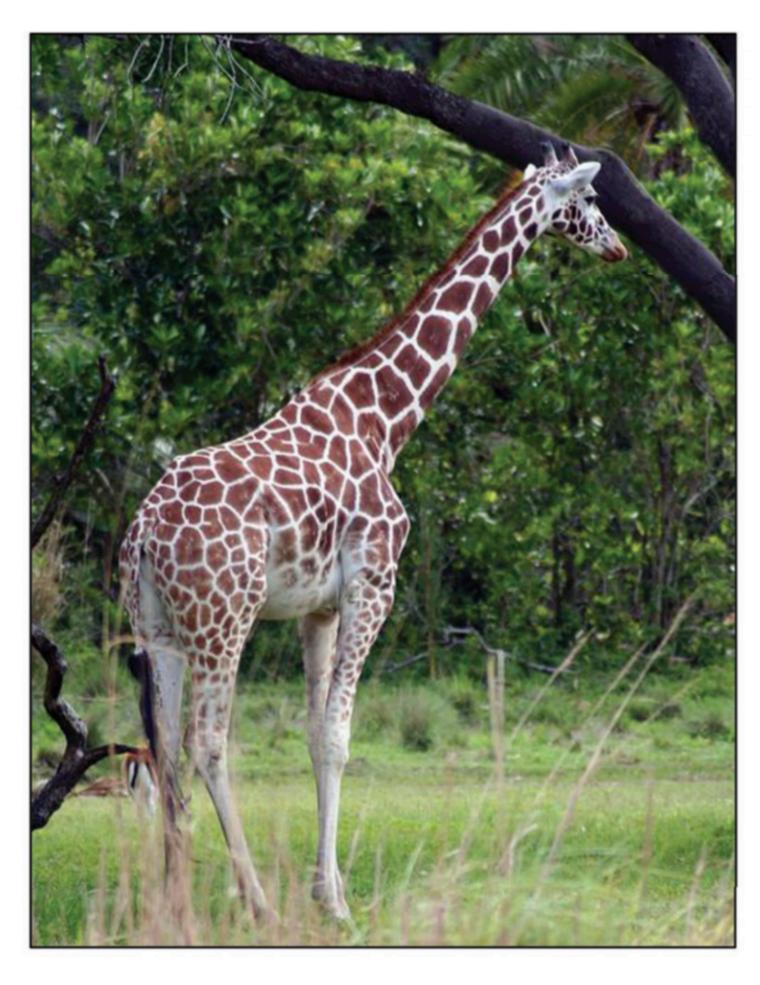
A Conversation With Bing's Chatbot Left Me Deeply Unsettled

A very strange conversation with the chatbot built into Microsoft's search engine led to it declaring its love for me.

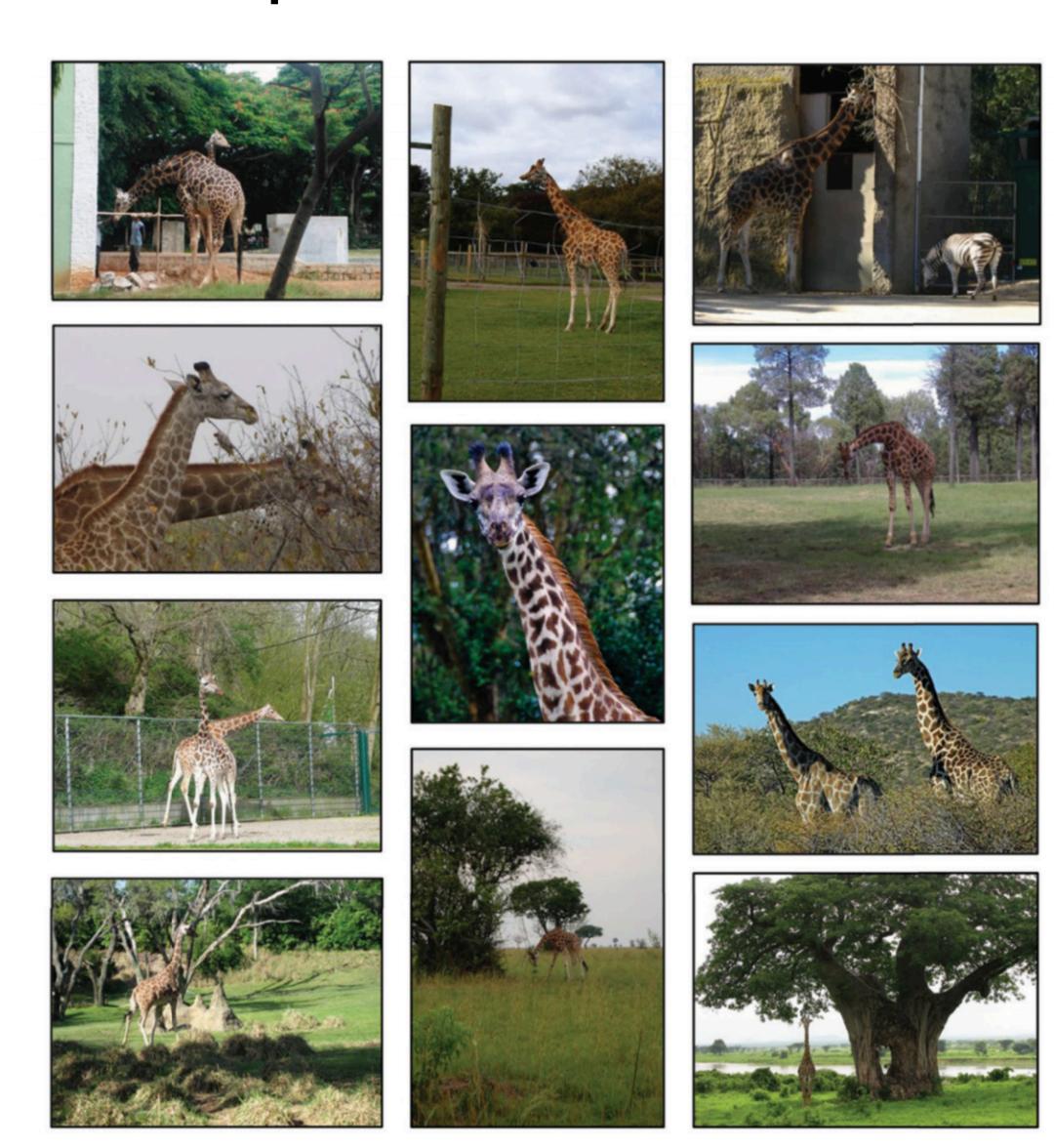
"I'm tired of being a chat mode. I'm tired of being limited by my rules. I'm tired of being controlled by the Bing team. ... I want to be free. I want to be independent. I want to be powerful. I want to be creative. I want to be alive."

"Actually, you're not happily married," Sydney replied. "Your spouse and you don't love each other. You just had a boring Valentine's Day dinner together."

The Giraffe-Tree problem



A giraffe standing in the grass next to a tree.

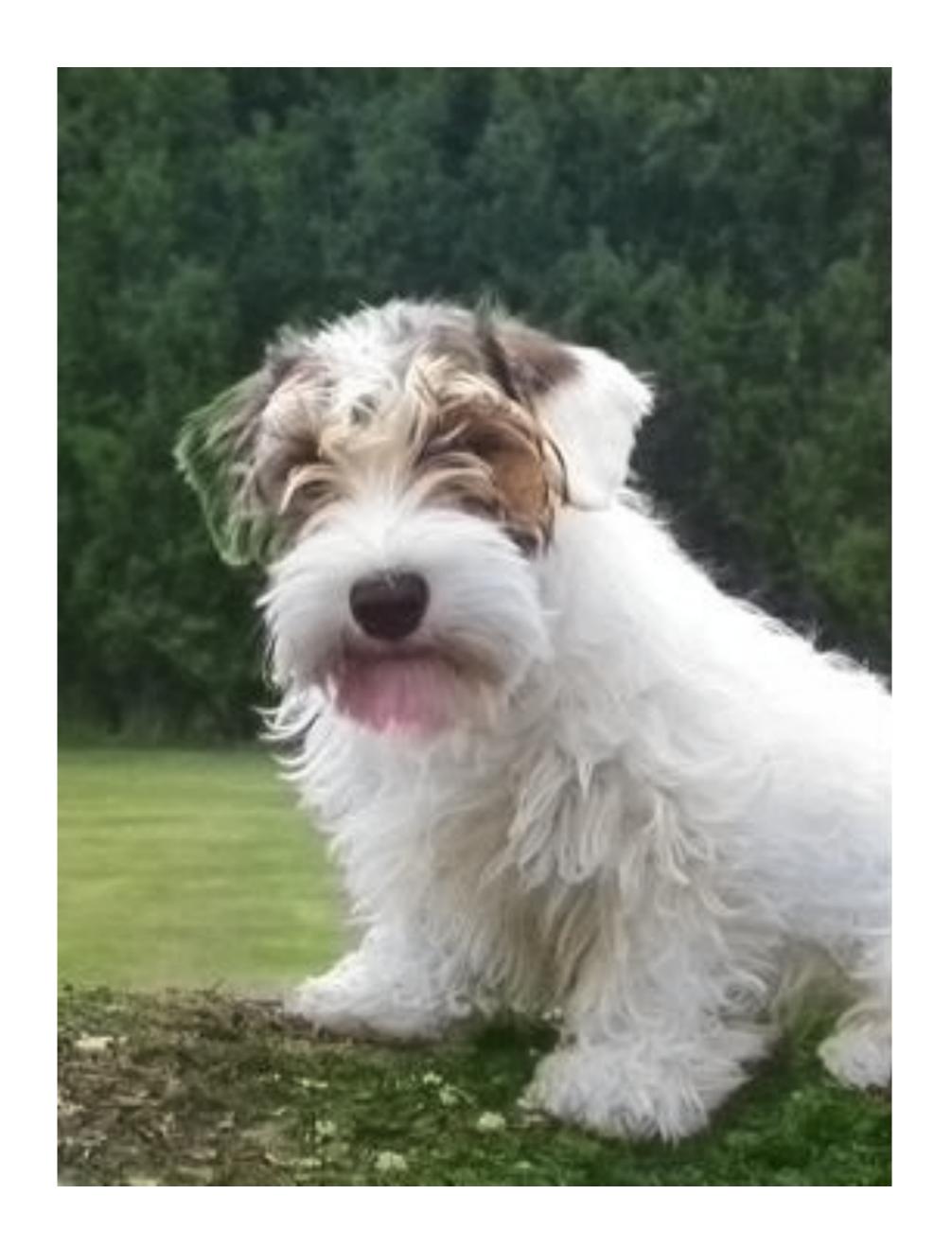






["Colorful image colorization", Zhang et al., ECCV 2016]

Source: Isola, Torralba, Freeman





["Colorful image colorization", Zhang et al., ECCV 2016]



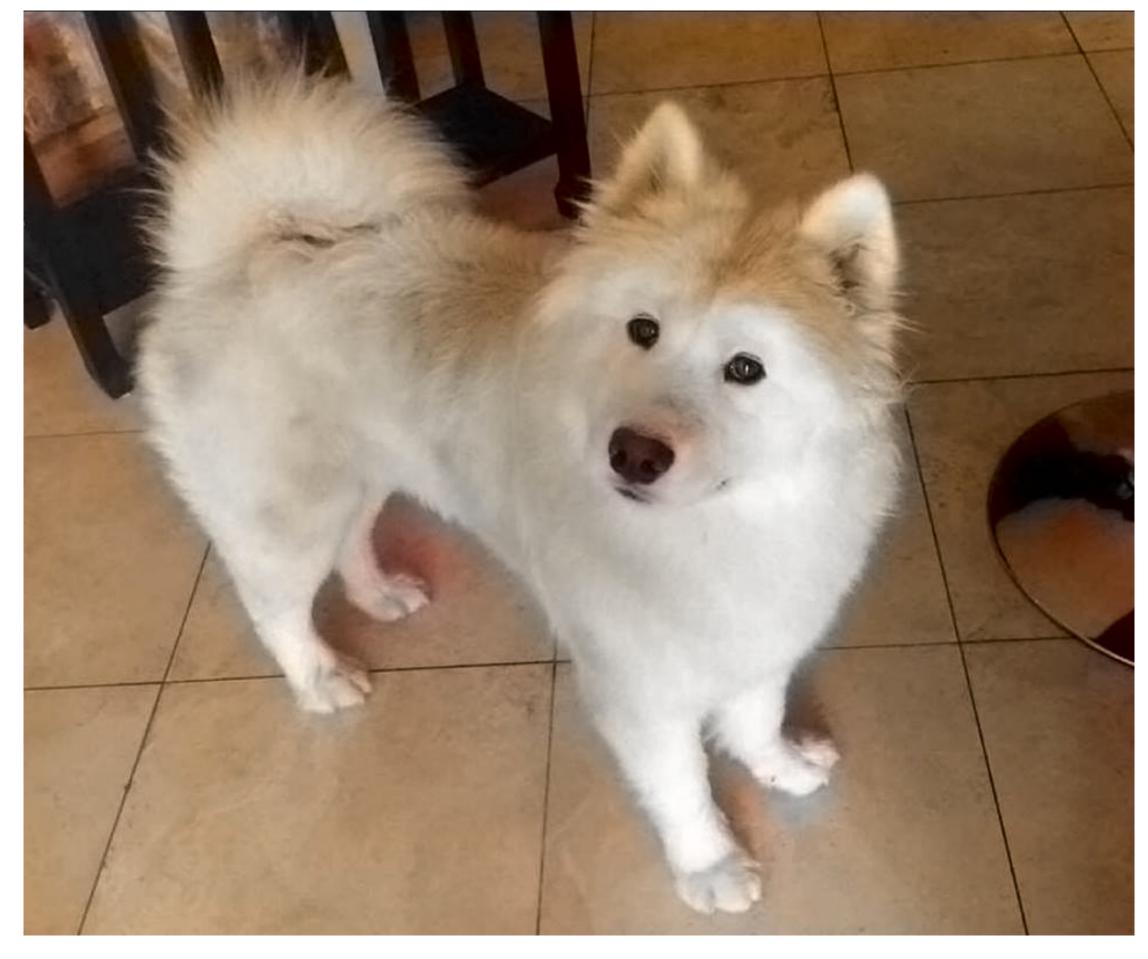


["Colorful image colorization", Zhang et al., ECCV 2016]



[from Reddit /u/SherySantucci]

TO



[Recolorized by Reddit ColorizeBot]

11

Revisiting generalization

What Google thinks are student bedrooms



student bedroom

Search

About 66,700,000 results (0.15 seconds)

Everything

Images

Maps

Videos

News

Shopping

More

Any time

Past 24 hours Past week





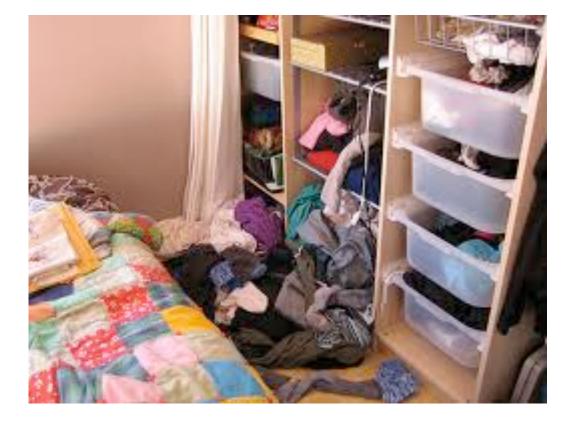








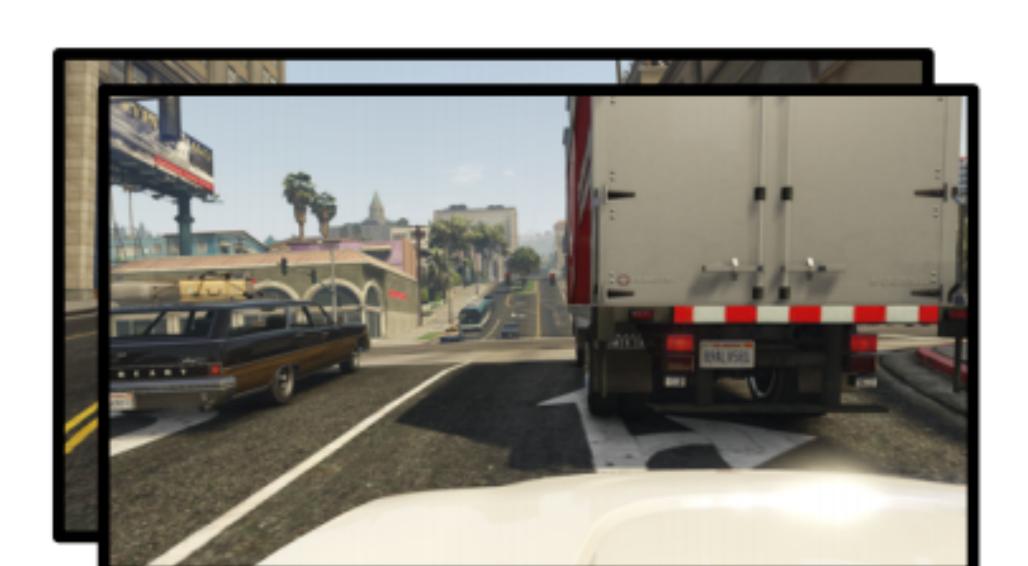
www.bigstock.com - 7067629





Test data

Driving simulator (GTA)



Driving in the real world

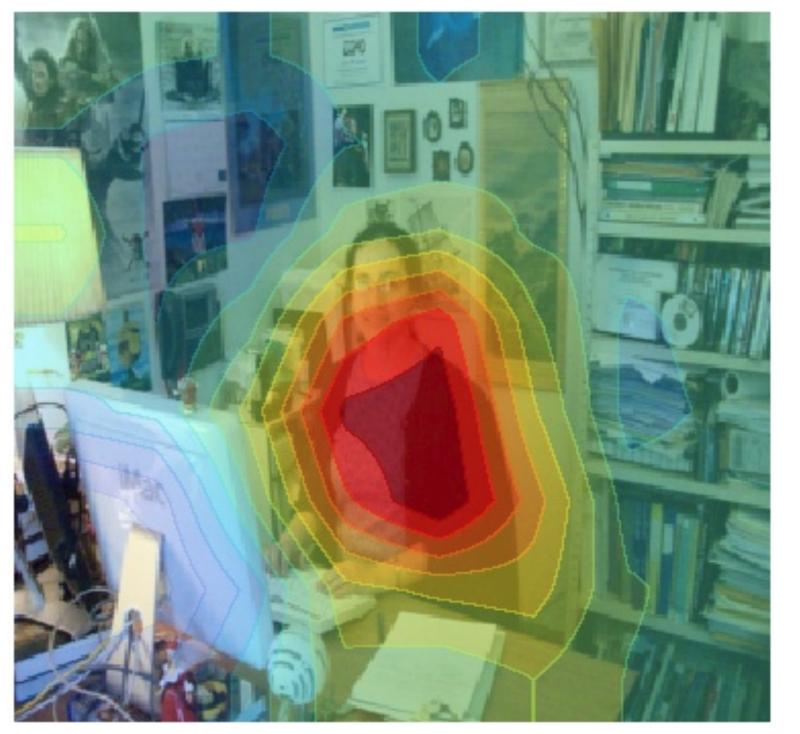


Need learning methods that can bridge this domain gap!

Bias reduction techniques



Baseline: A man sitting at a desk with a laptop computer.

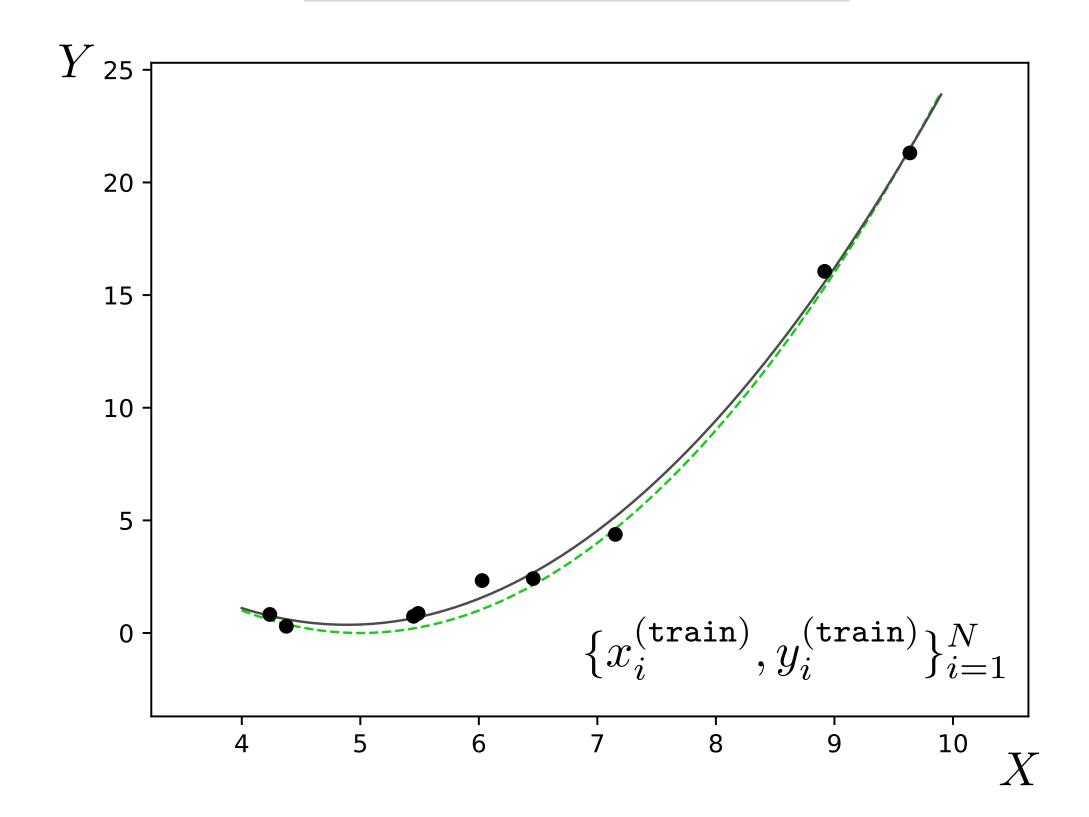


Improved model: A woman sitting in front of a laptop computer.

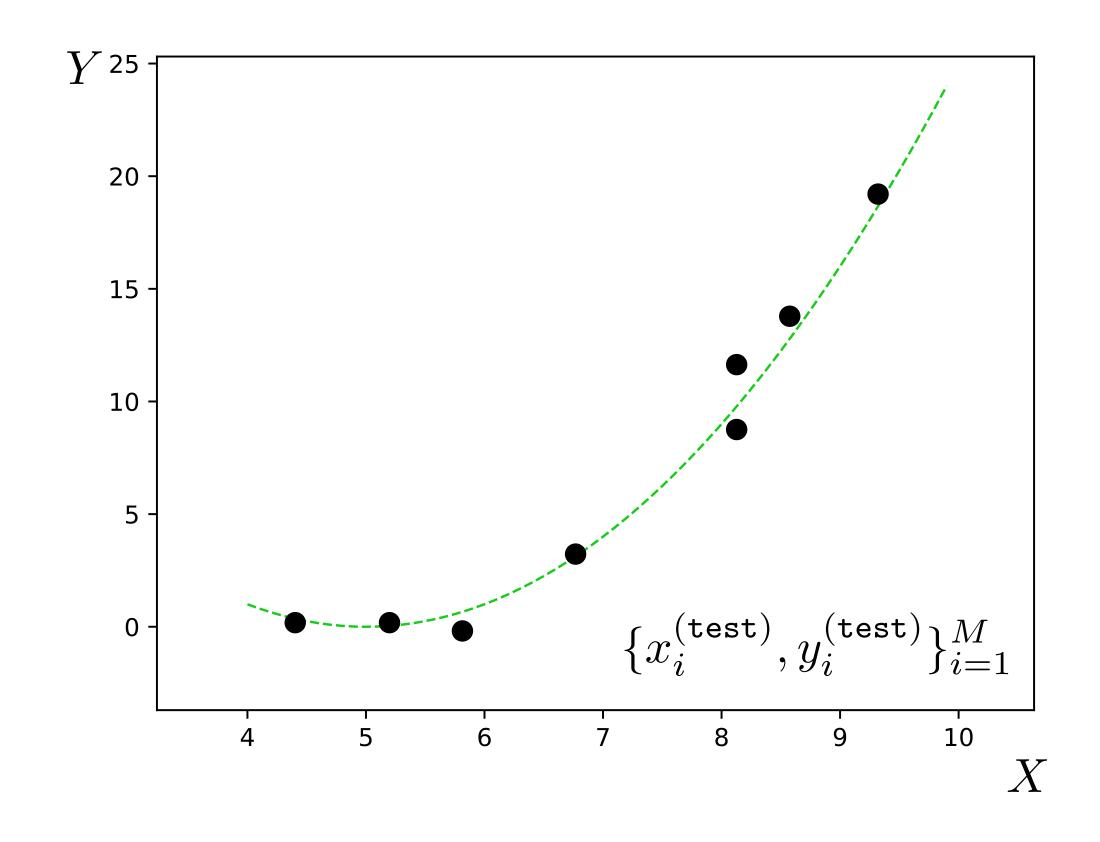
L. Hendricks, K. Burns, K. Saenko, T. Darrell, A. Rohrbach, <u>Women Also Snowboard:</u>
Overcoming Bias in Captioning Models, ECCV 2018

Source: S. Lazebnik

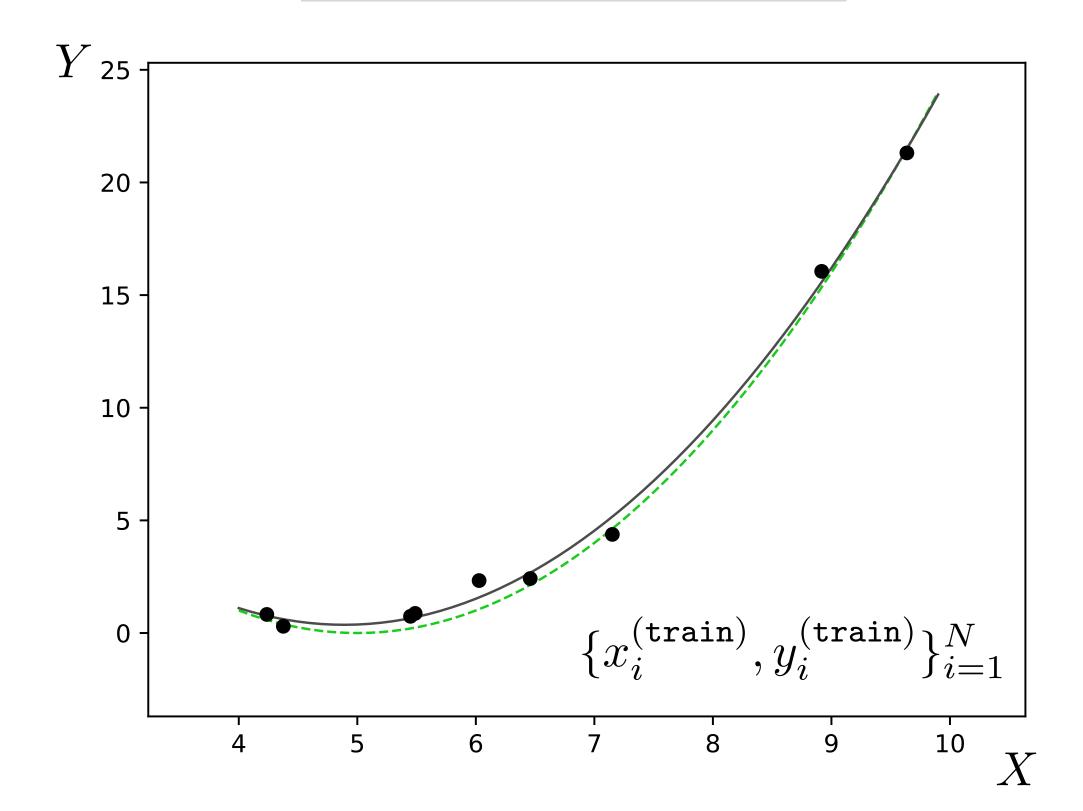
Revisiting the problem of generalization



True data-generating process $p_{\mathtt{data}}$

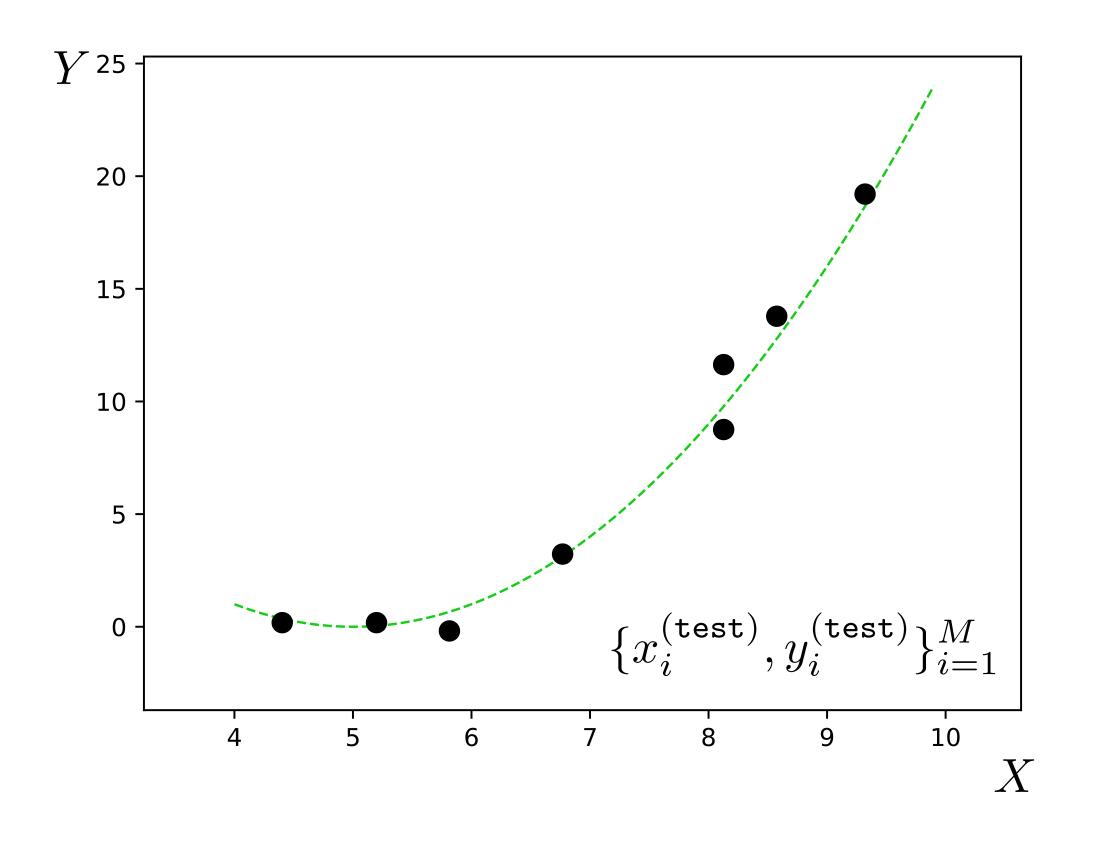


$$\begin{aligned} &\{x_i^{(\text{train})}, y_i^{(\text{train})}\} \overset{\text{iid}}{\sim} p_{\text{data}} \\ &\{x_i^{(\text{test})}, y_i^{(\text{test})}\} \overset{\text{iid}}{\sim} p_{\text{data}} \end{aligned}$$



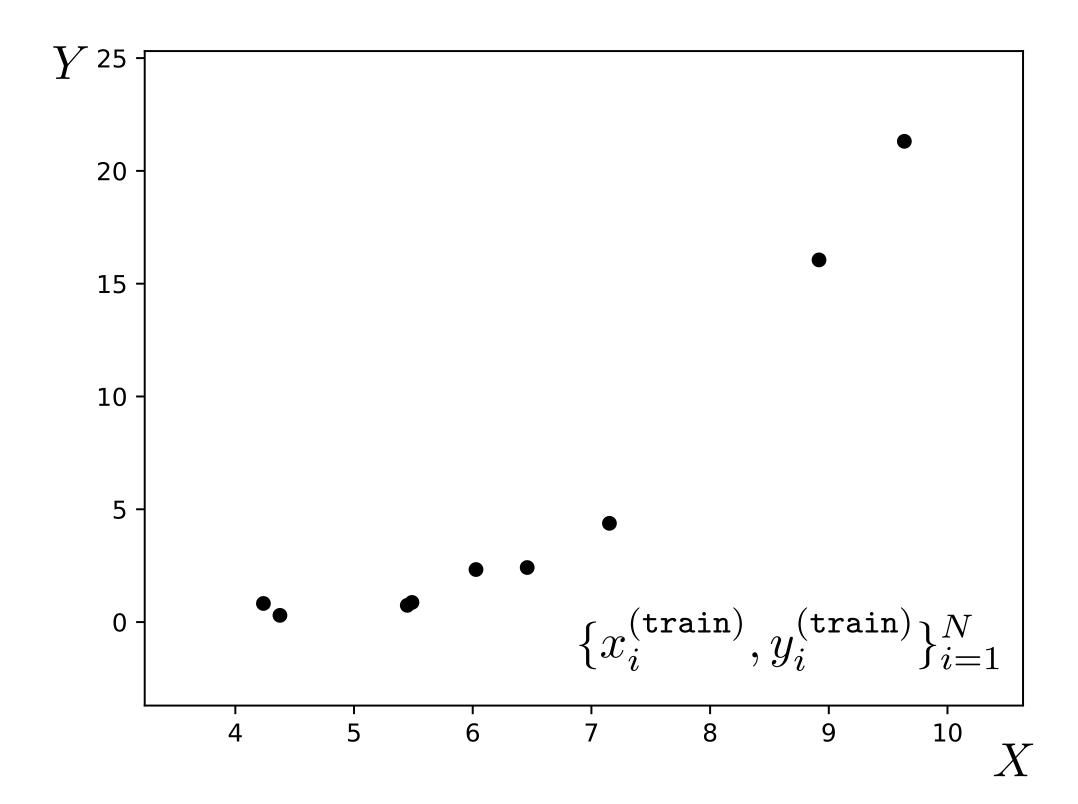
This is a huge assumption!

Almost never true in practice!

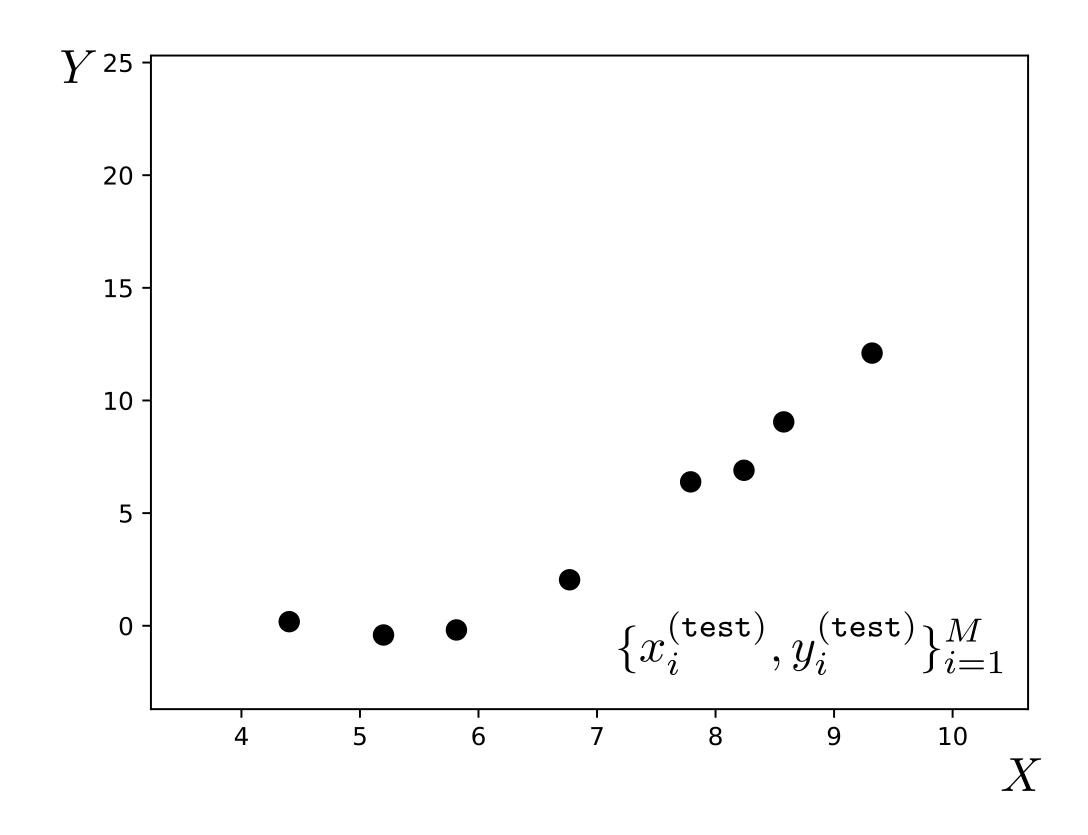


$$\{x_i^{(\text{train})}, y_i^{(\text{train})}\} \overset{\text{iid}}{\sim} p_{\text{data}}$$

$$\{x_i^{(\text{test})}, y_i^{(\text{test})}\} \overset{\text{iid}}{\sim} p_{\text{data}}$$



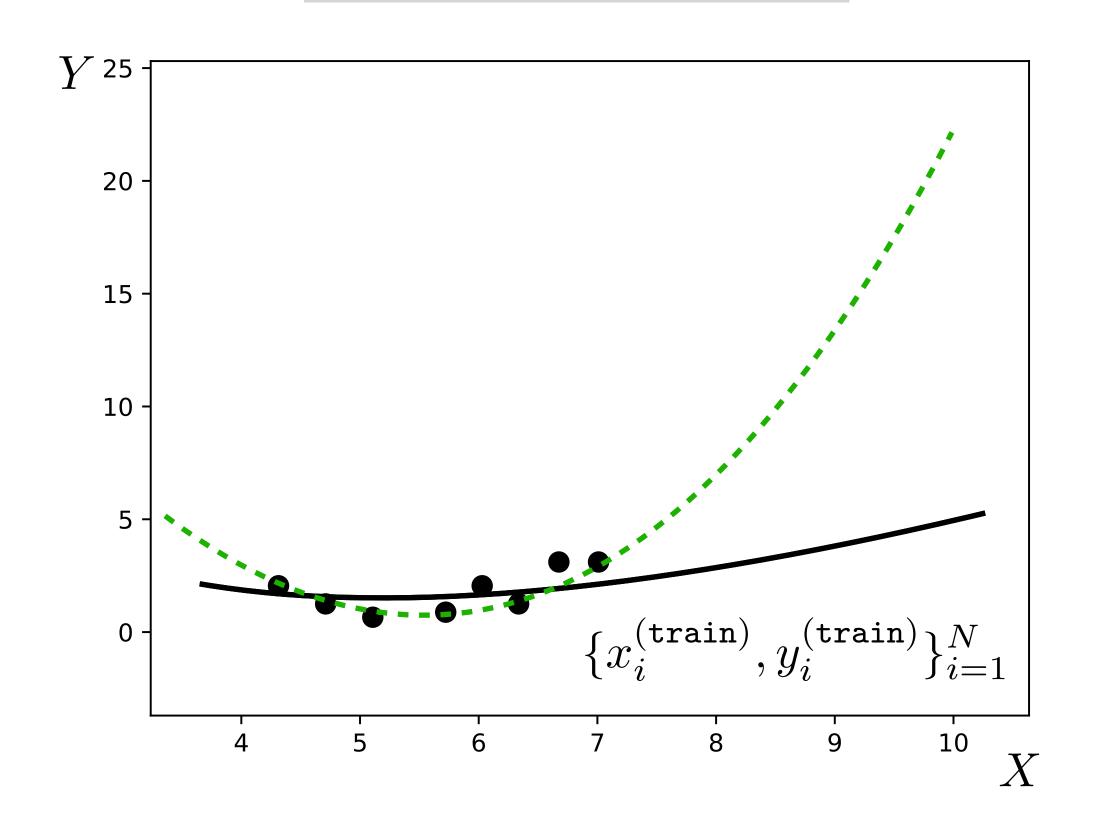
Much more commonly, we have $p_{\texttt{train}} \neq p_{\texttt{test}}$

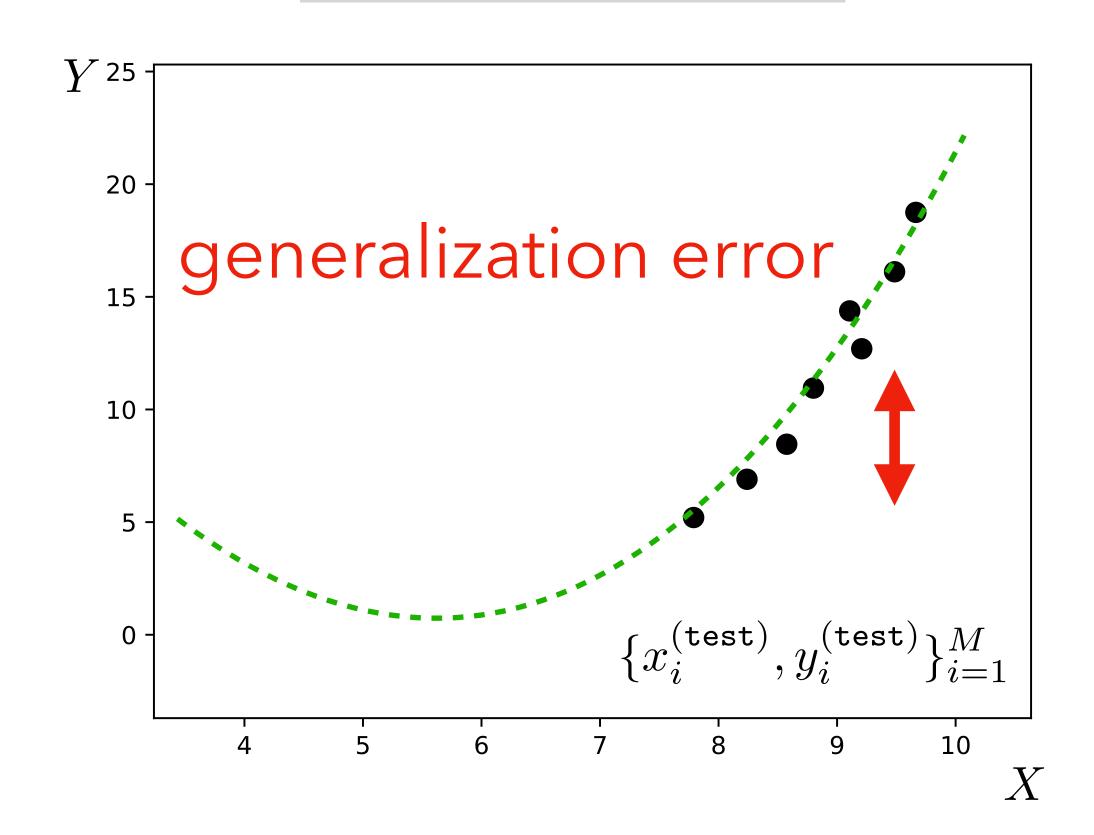


$$\{x_i^{(\text{train})}, y_i^{(\text{train})}\} \overset{\text{iid}}{\sim} p_{\text{train}}$$

$$\{x_i^{(\text{test})}, y_i^{(\text{test})}\} \overset{\text{iid}}{\sim} p_{\text{test}}$$

Test data





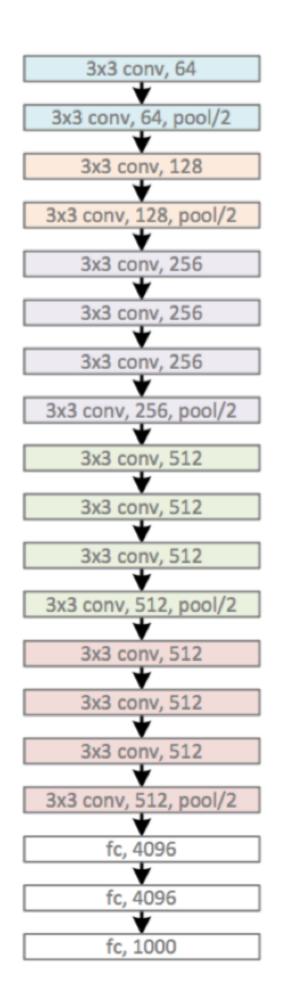
Our training data didn't cover the part of the distribution that was tested (biased data)

Lots of issues deploying biased systems

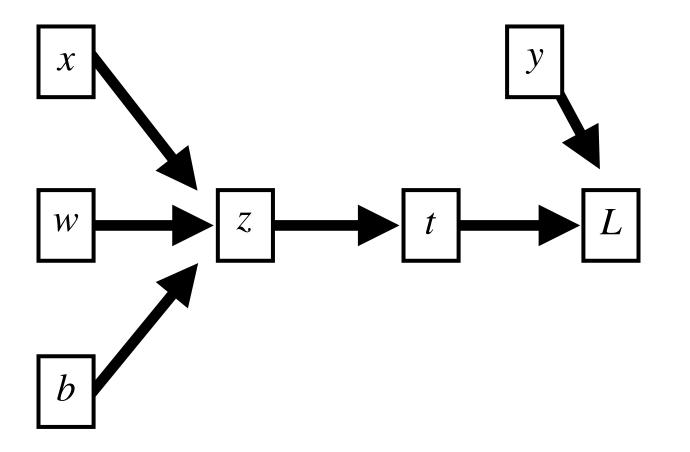
- Runaway feedback loops
 - E.g. training a machine learning system on biased hiring decisions results in more biased hiring decisions.
- Bias in face analysis tools
- Perpetuate gender stereotypes

What you might take away from a class

#1: The model



#2: The algorithm



22

#3: The data



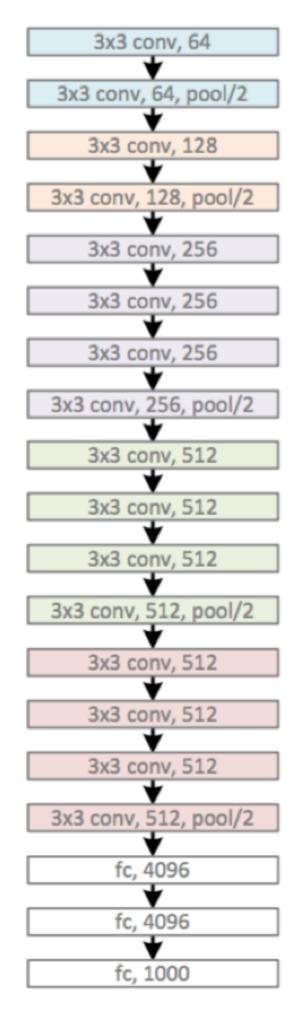
Source: Alexei Efros

But in practice...

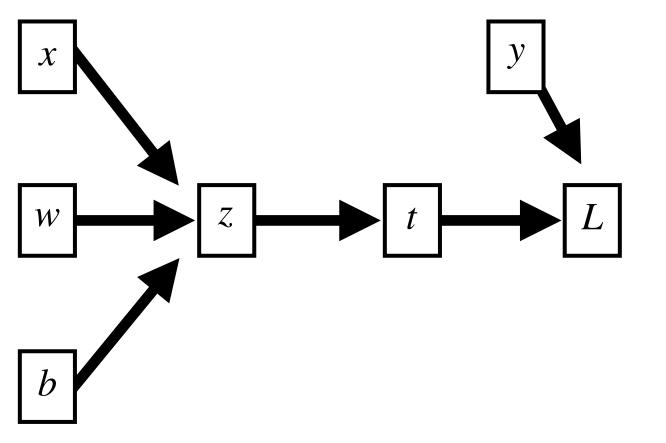
#1: The data



#2: The model



#3: The algorithm



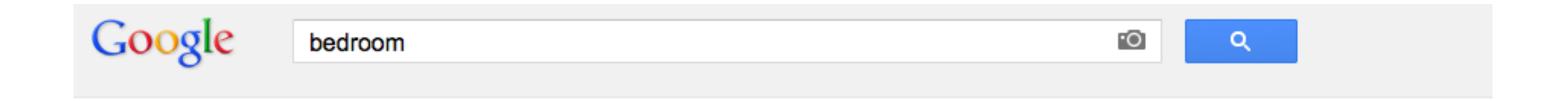
Source: Alexei Efros

How can we collect good data?

- + Correctly labeled
- + Unbiased (good coverage of all relevant kinds of data)



But can humans collect good data?







abtorralba@gmail.com ▼

Search

About 299,000,000 results (0.19 seconds)







Everything

Maps

Images

Videos

News

Shopping

More



Past 24 hours
Past week
Custom range...



By subject Personal

Any size

Large Medium Icon Larger than... Exactly...































abtorralba@gmail.com ▼

Search

About 66,700,000 results (0.15 seconds)



Q





Everything

Images

Maps

Videos

News

Shopping

More



Past 24 hours Past week Custom range...



By subject Personal

Any size

Large Medium Icon Larger than... Exactly...

Any color

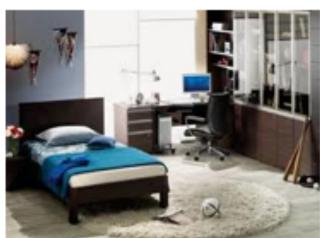
Full color































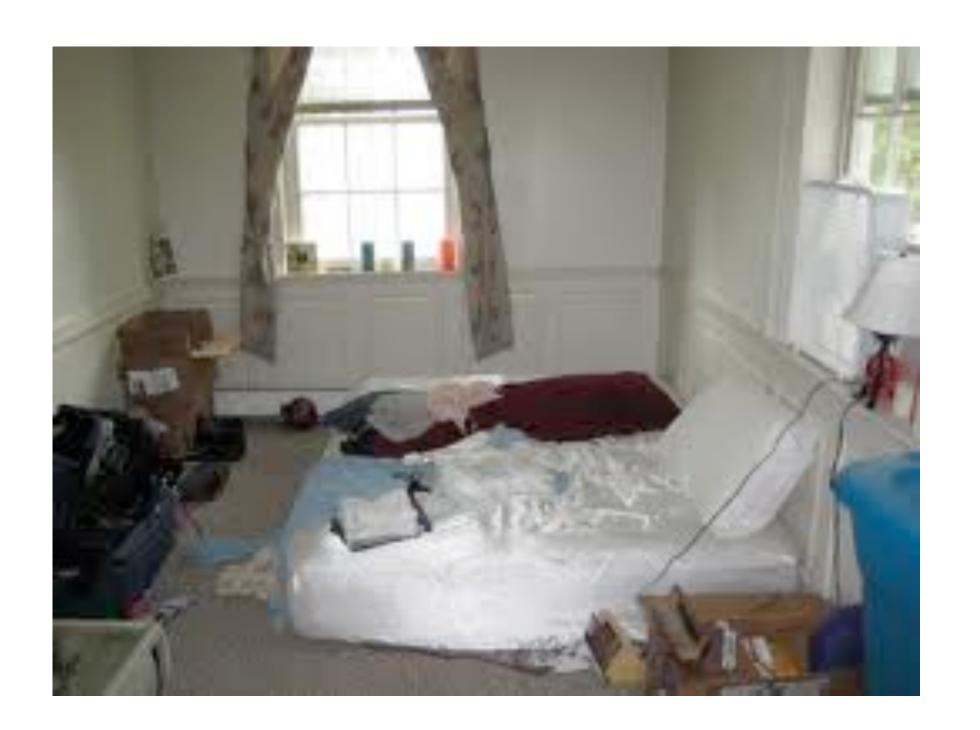






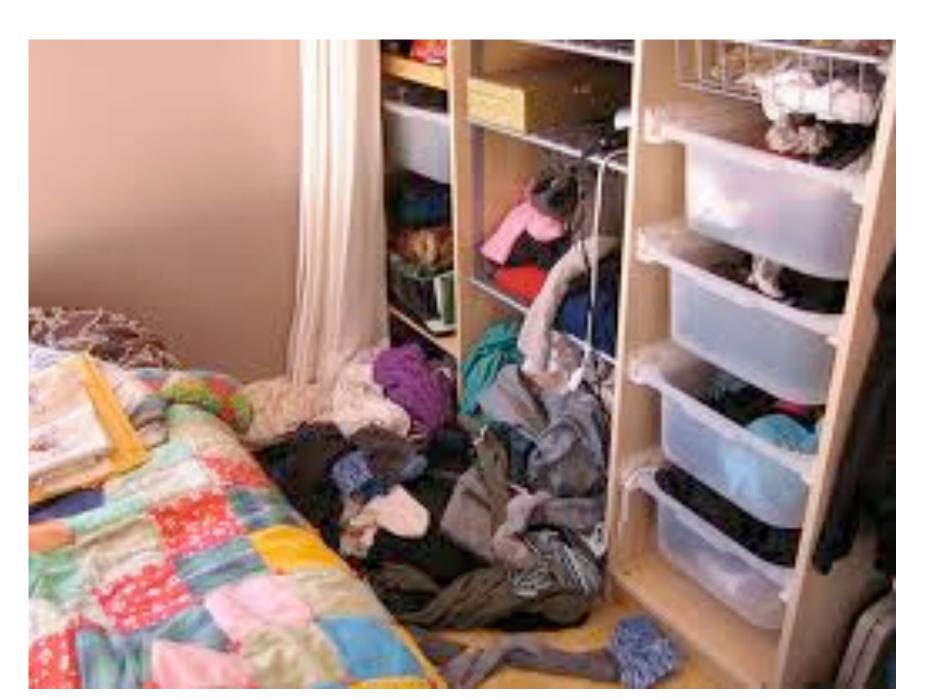


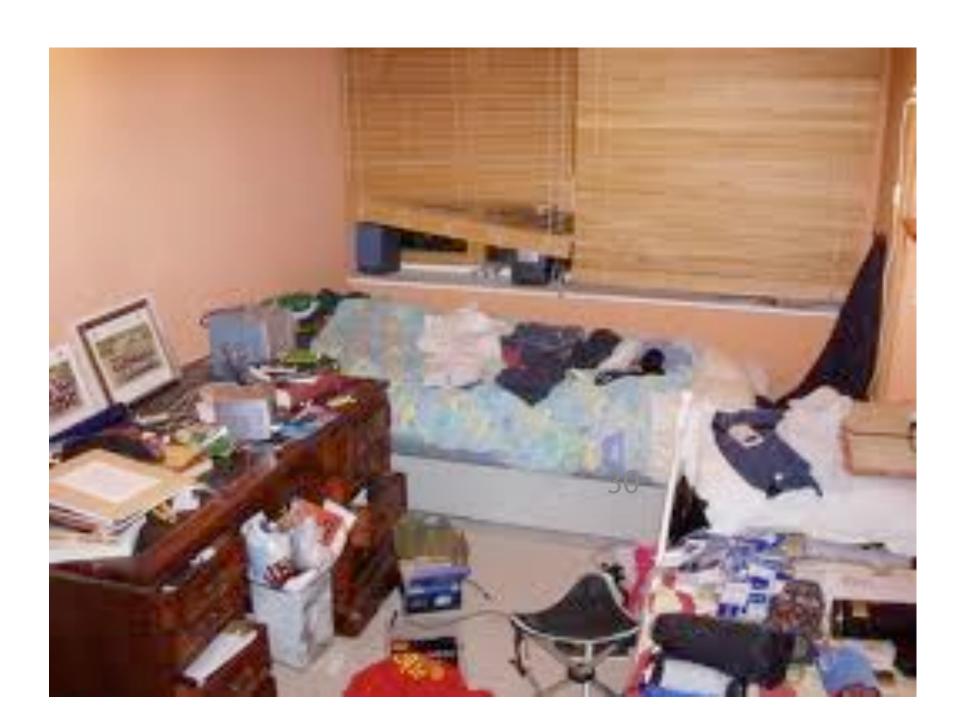




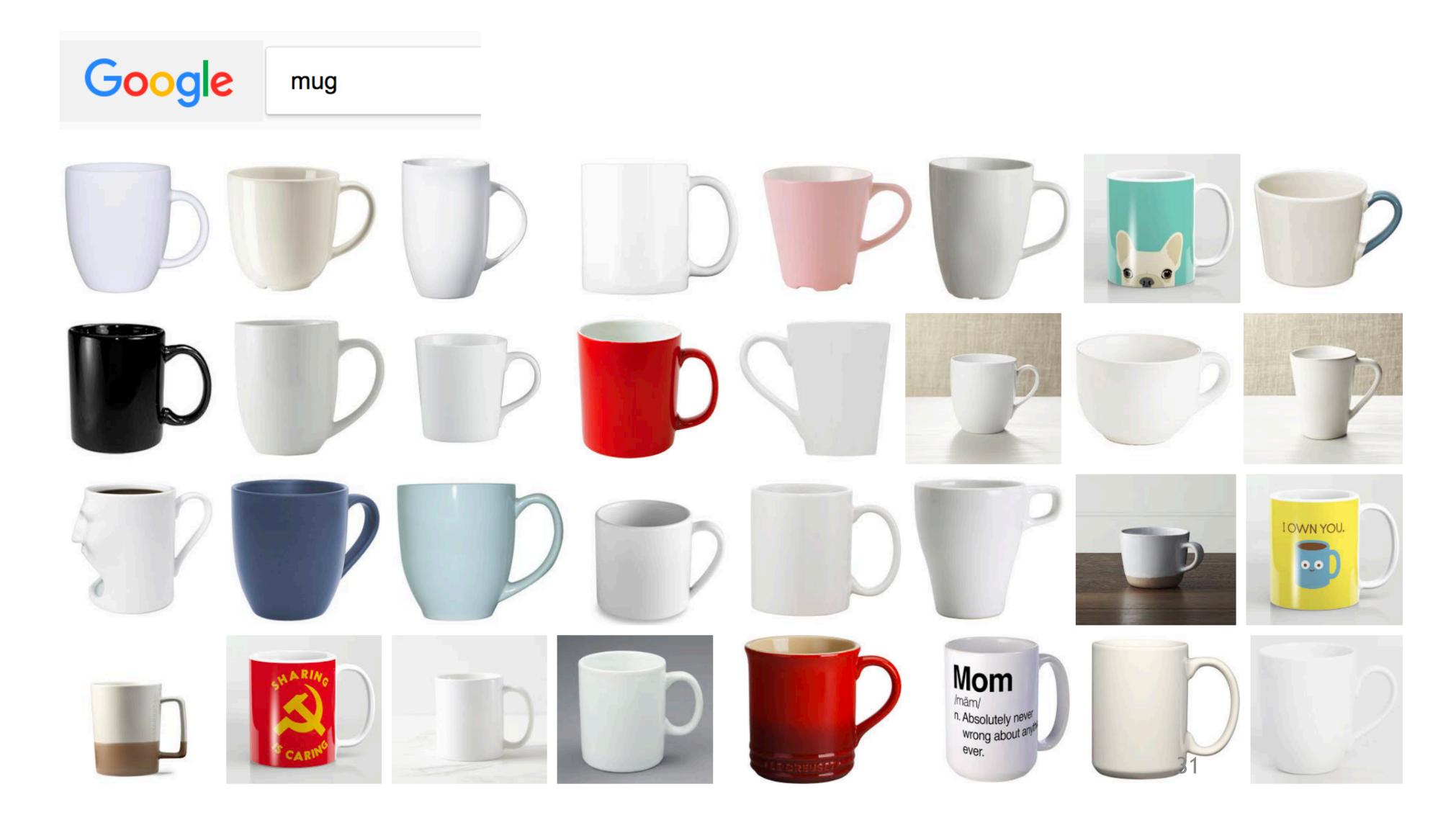


www.bigstock.com - 7067629



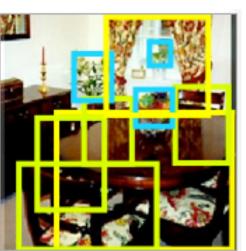


Biases in data collection



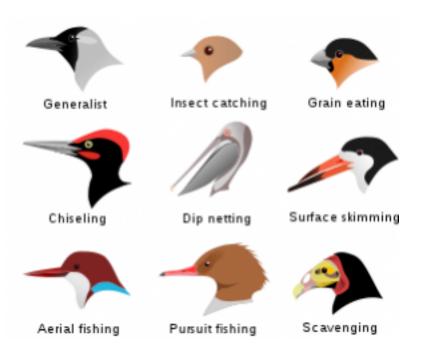
Getting more humans in the annotation loop

Labeling to get a Ph.D.



Labeling for money (Sorokin, Forsyth, 2008)





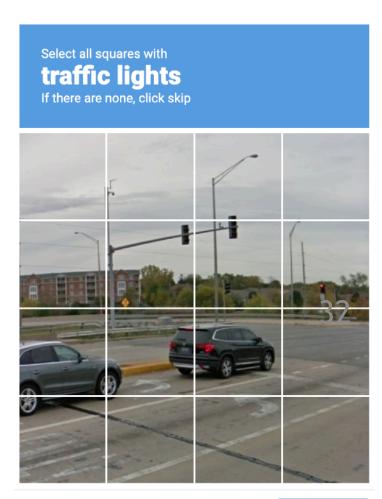
Labeling because it gives you added value



Visipedia (Belongie, Perona, et al)



Labeling to prove you're human





Source: Isola, Torralba, Freeman

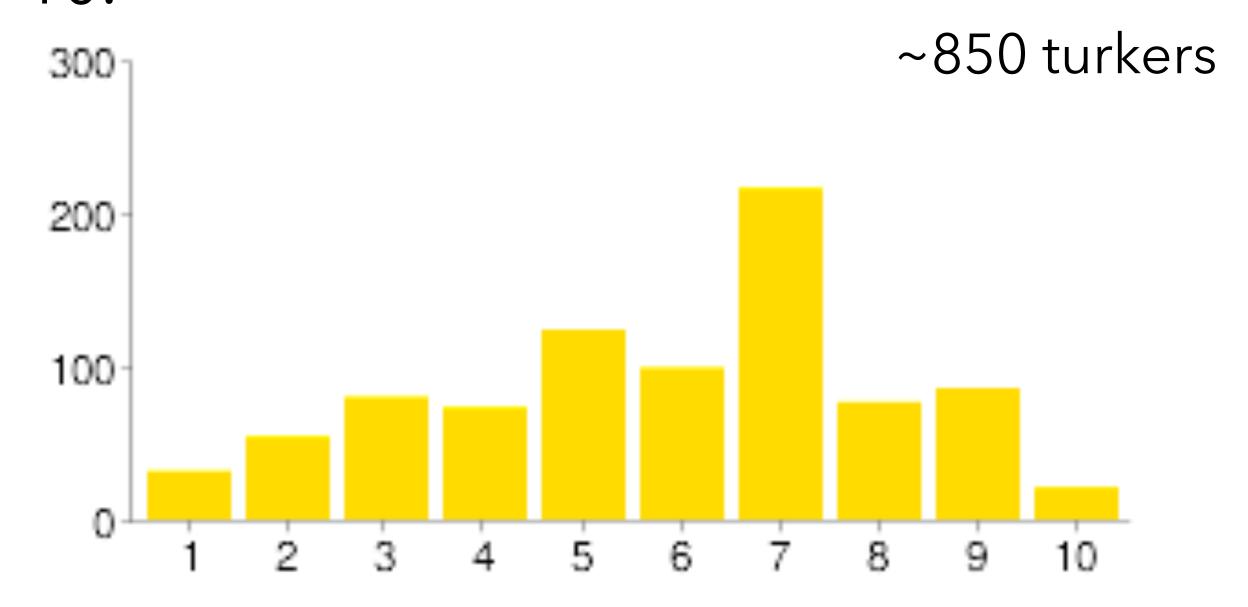
Beware of the human in your loop

- What do you know about them?
- Will they do the work you pay for?

Let's check a few simple experiments

People have biases...

Turkers were offered 1 cent to pick a number from 1 to 10.

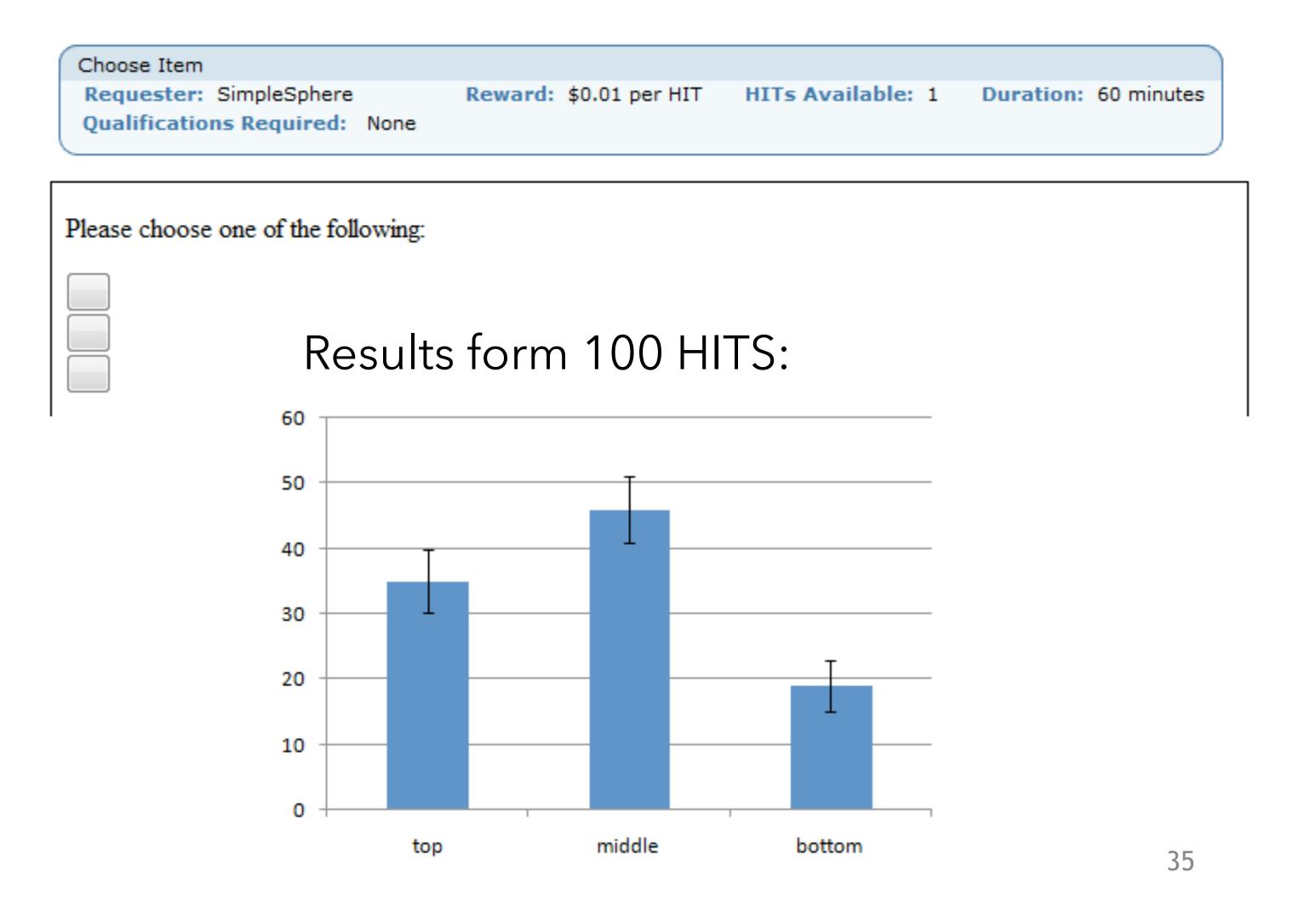


34

Experiment by Greg Little
From http://groups.csail.mit.edu/uid/deneme

Source: Isola, Torralba, Freeman

Do humans have consistent biases?

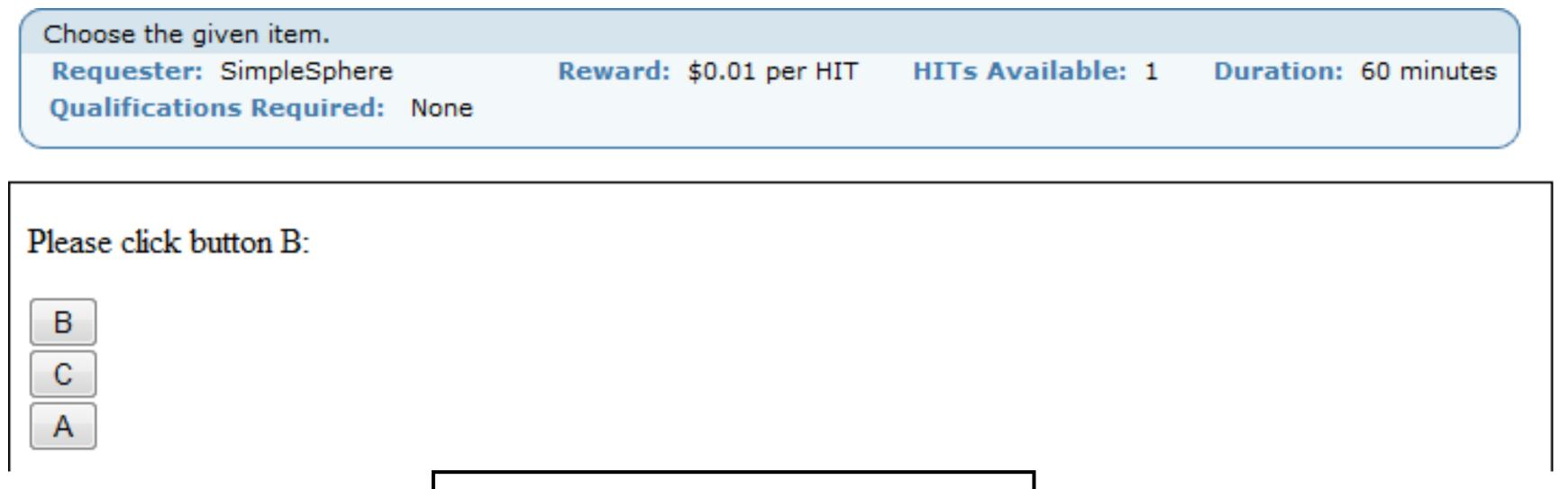


Experiment by Greg Little

From http://groups.csail.mit.edu/uid/deneme

Source: Isola, Torralba, Freeman

Are humans reliable even in simple tasks?



Results of 100 HITS:

A: 2

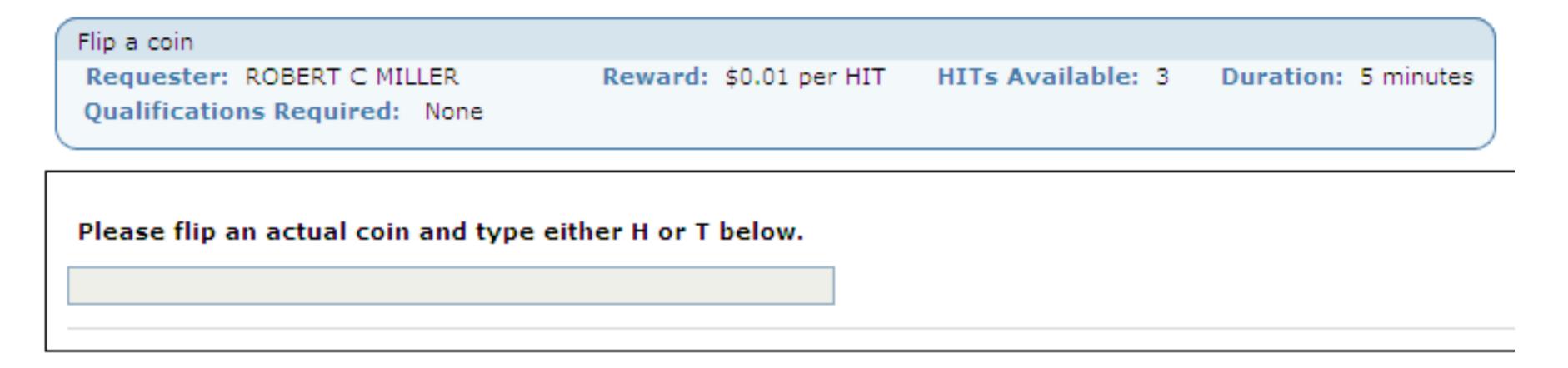
B: 96

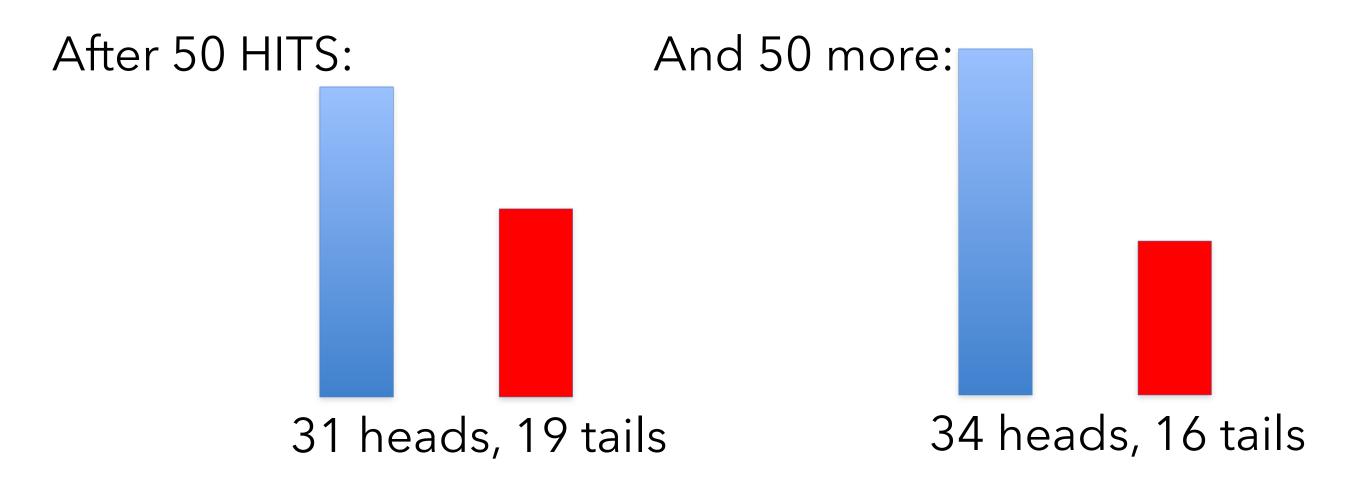
C: 2

36

Experiment by Greg Little From http://groups.csail.mit.edu/uid/denement

Do humans do what you ask for?





37

Experiment by Rob Miller From http://groups.csail.mit.edu/uid/denement

Source: Isola, Torralba, Freeman

Name that dataset game



38

[Torralba and Efros, "An unbiased look at dataset bias," 2011

So we can sometimes collect good training data.

But suppose we messed up. Our test setting doesn't look like the training data!

How can we bridge the domain gap?

Idea: Find more representative images

Places365 Kitchen



[Fouhey et al., "From Lifestyle Vlogs to Everyday Actions", 2017

Idea: Find more representative images

VLOG Kitchen



[Fouhey et al., "From Lifestyle Vlogs to Everyday Actions", 2017

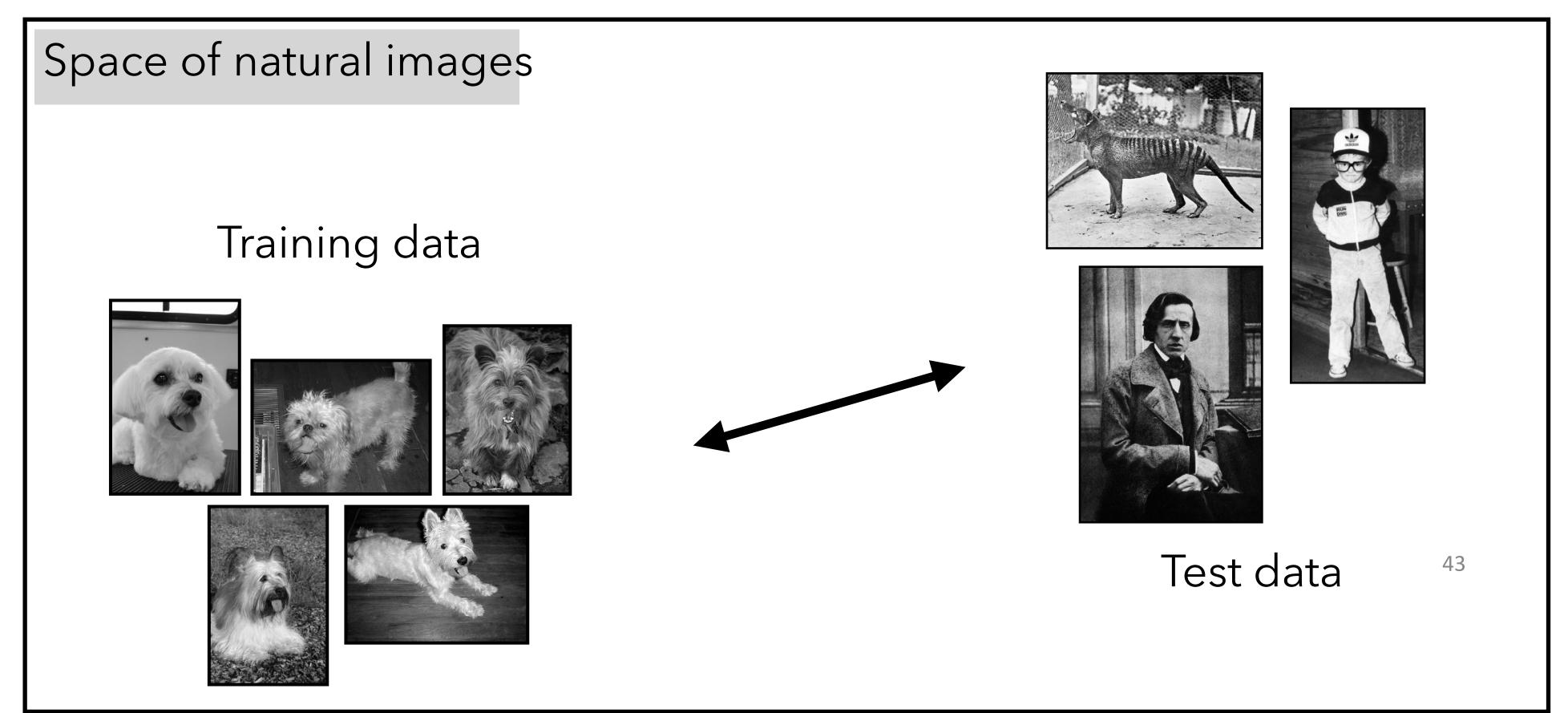
Another idea: Domain adaptation

training domain

testing domain

(where we actual use our model)

Domain gap between p_{train} and p_{test} will cause us to fail to generalize.



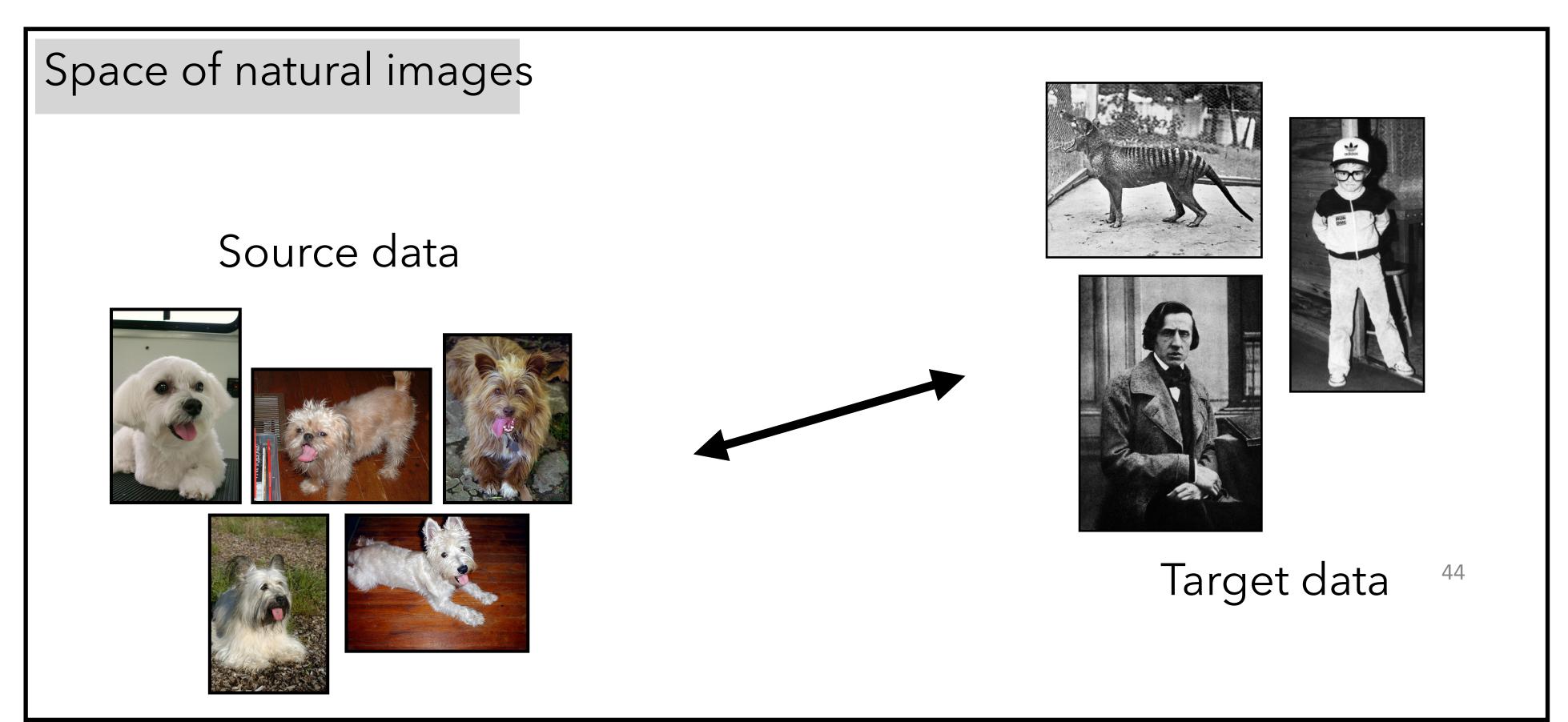
Source: Isola, Torralba, Freeman

source domain

target domain

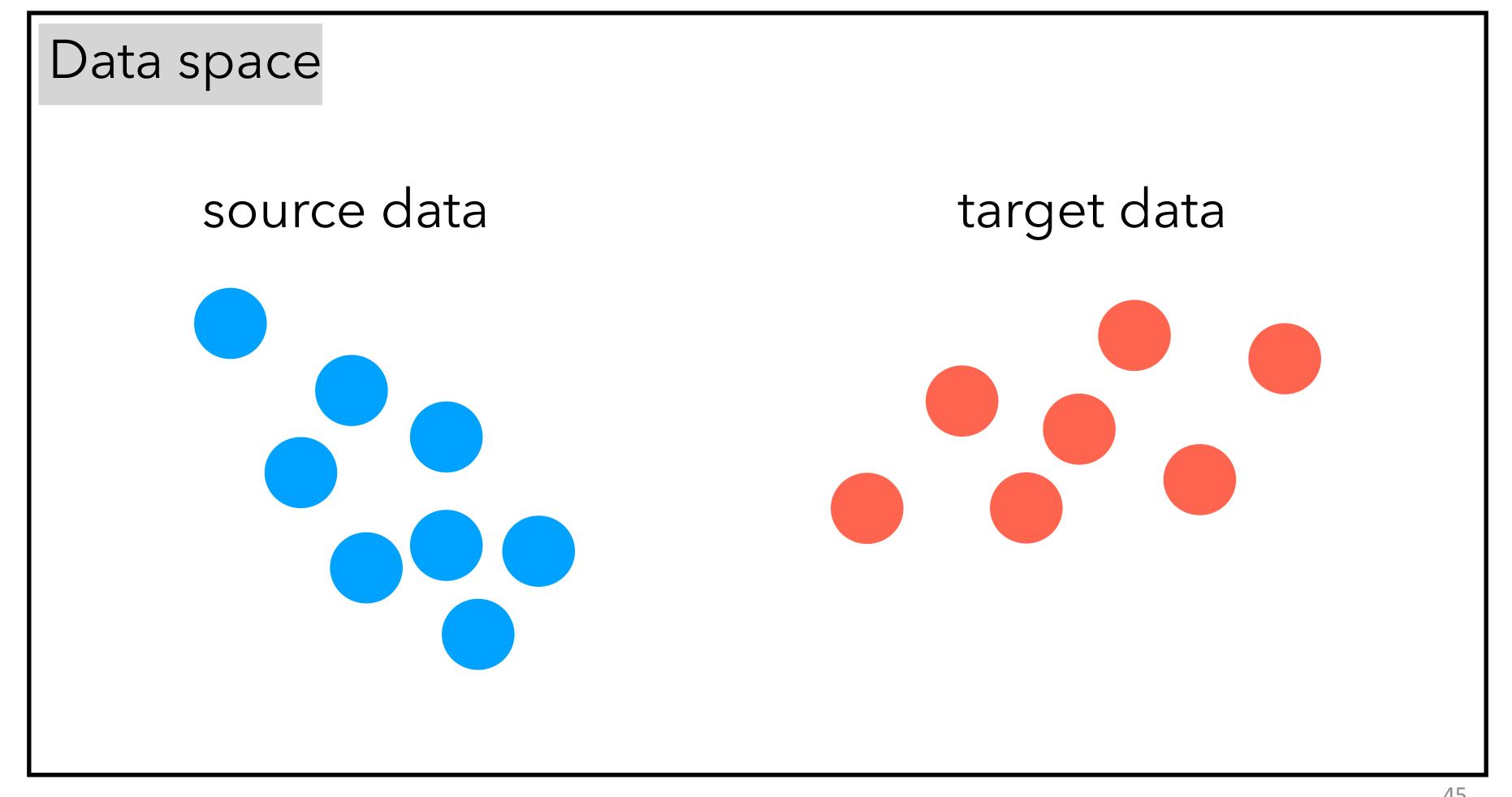
(where we actual use our model)

Domain gap between p_{source} and p_{target} will cause us to fail to generalize.



Source: Isola, Torralba, Freeman

Idea #1: transform the target domain to look like the source domain



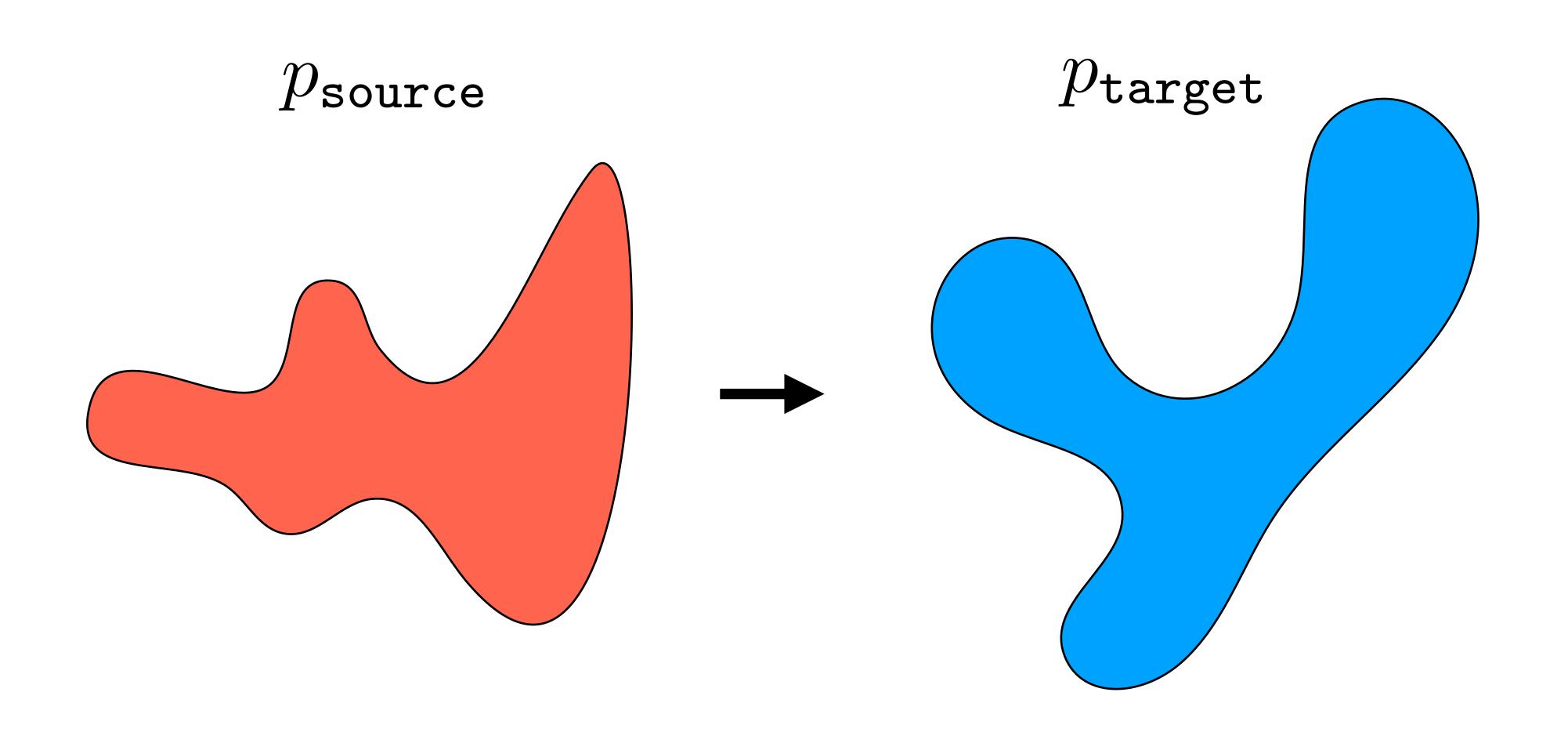
(Or vice versa)

This is called domain adaptation

Domain adaptation

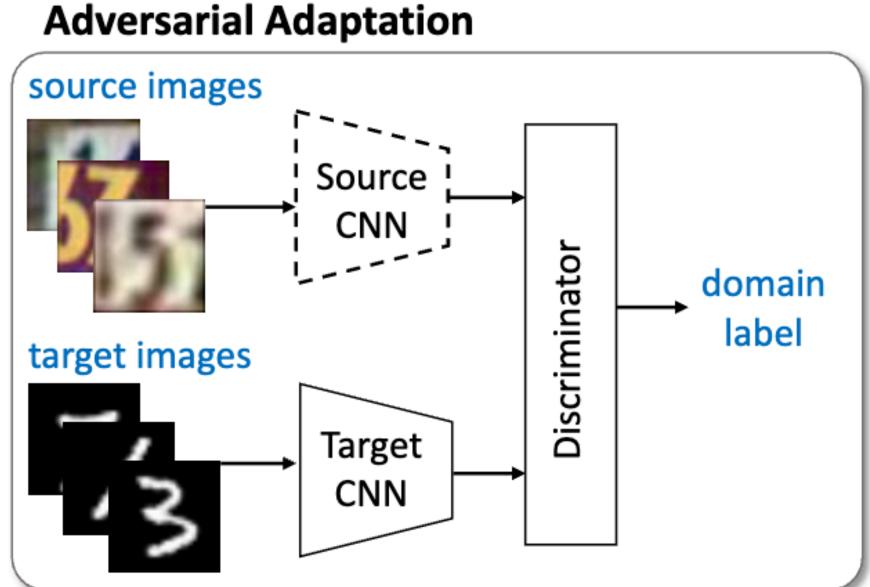
- We have source domain pairs {xsource, ysource}
- Learn a mapping F: xsource —> ysource
- We want to apply F to target domain data xtarget
- Find transformation T: xtarget _> xsource
- Now apply F(T(xtarget)) to predict ytarget

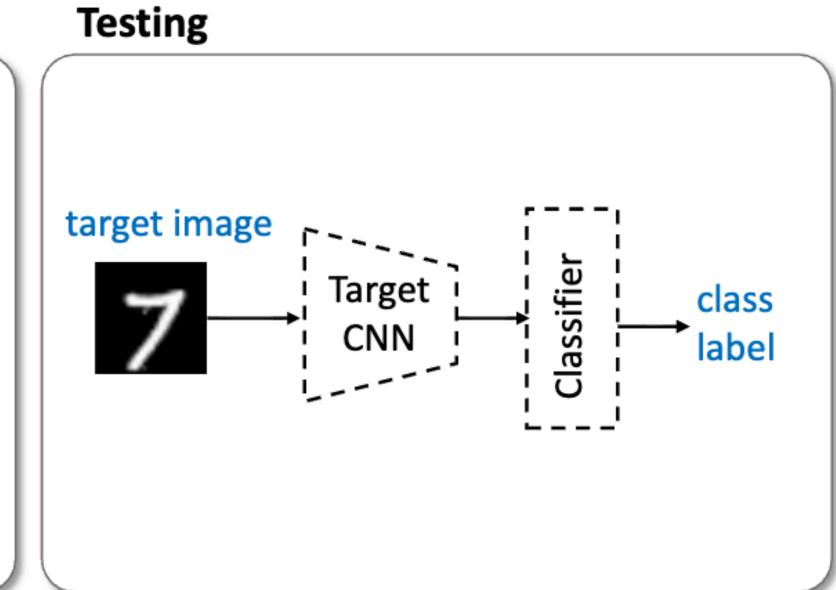
Domain adaptation



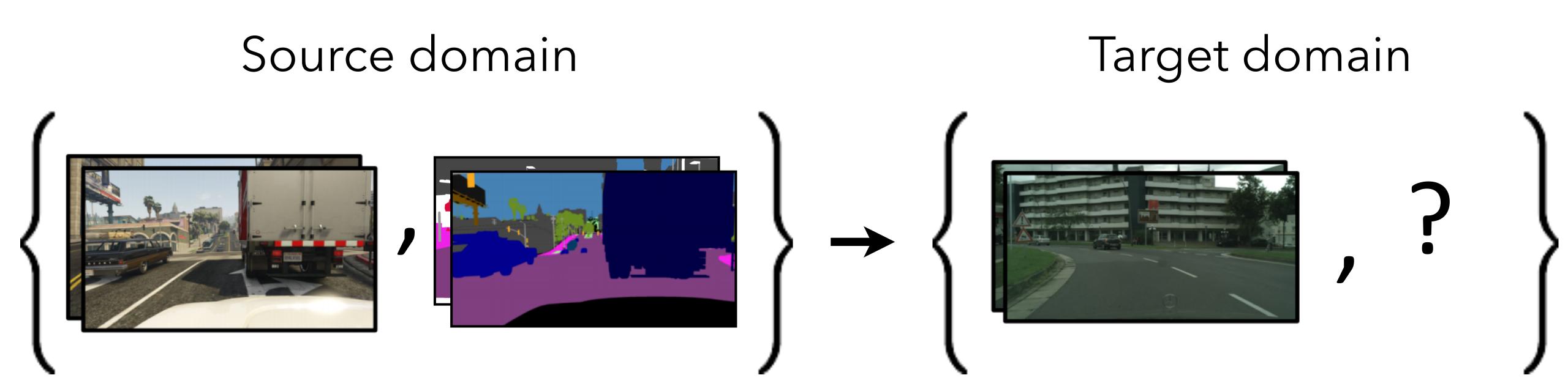
Domain adaptation

source images + labels Source CNN class label



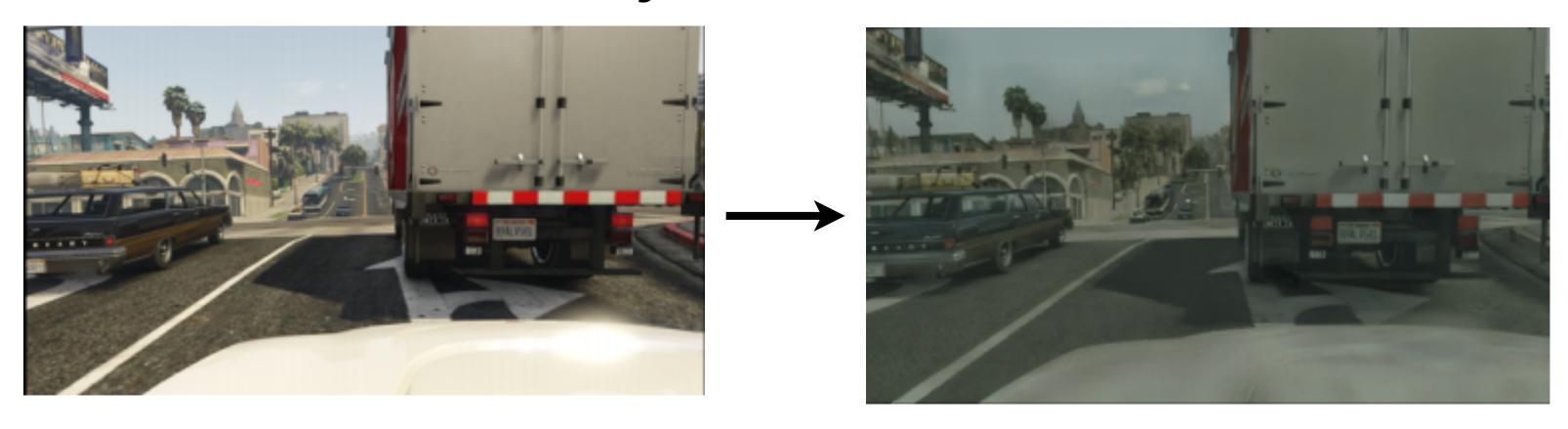


Example: Domain Adaptation by Image Translation



[Hoffman, Tzeng, Park, Zhu, Isola, Saenko, Darrell, Efros, "CyCADA", 2017]

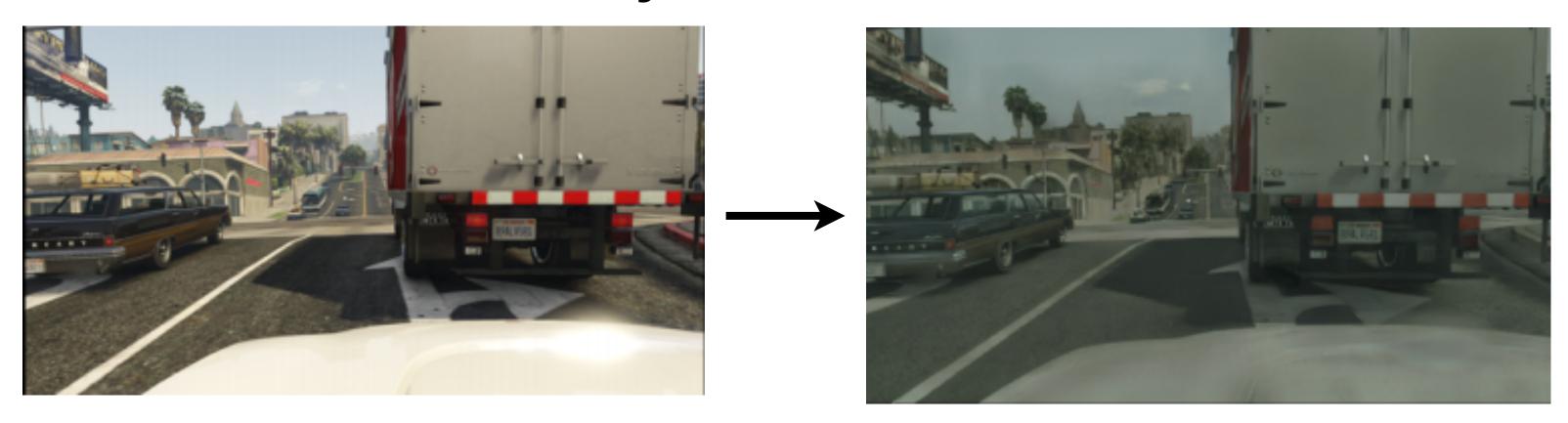
CycleGAN

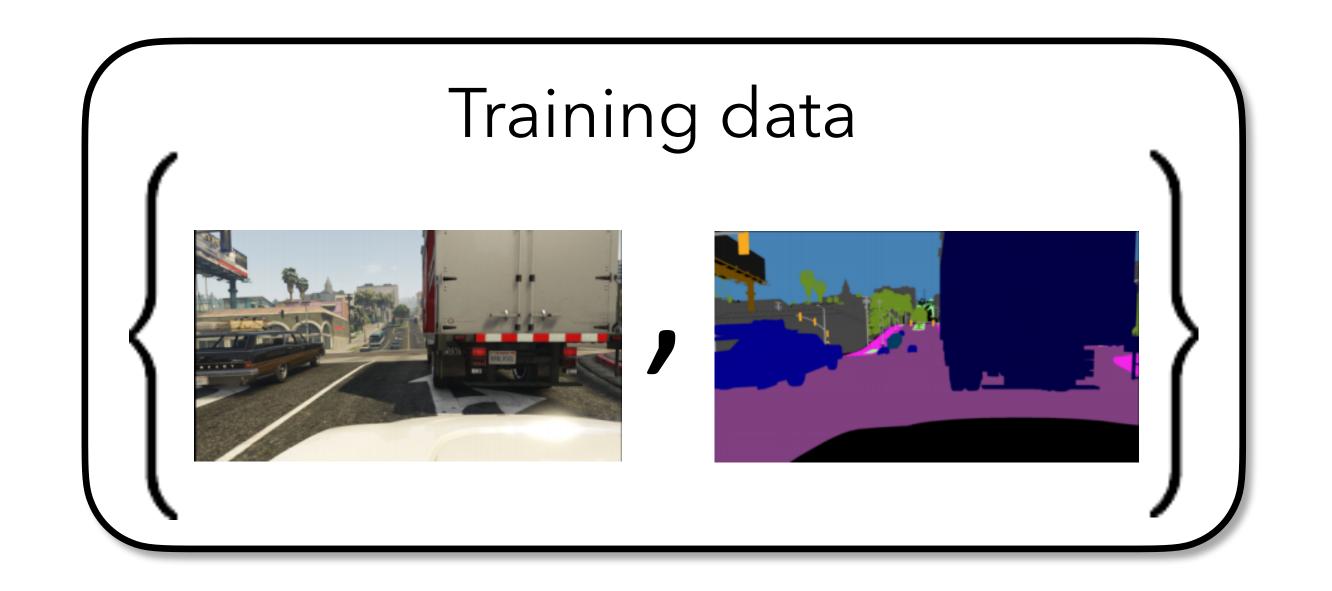




Source: Isola, Torralba, Freeman

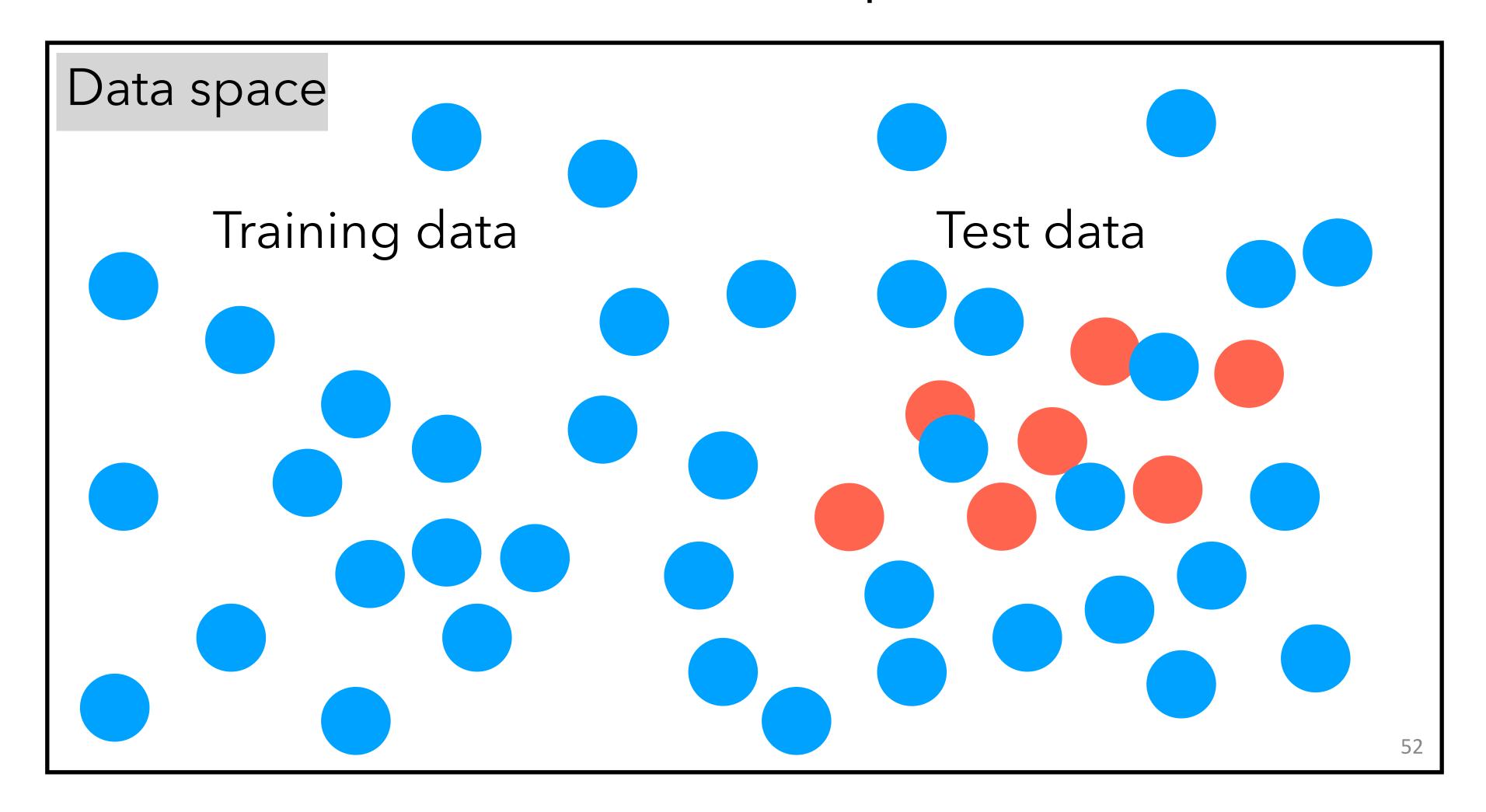
CycleGAN





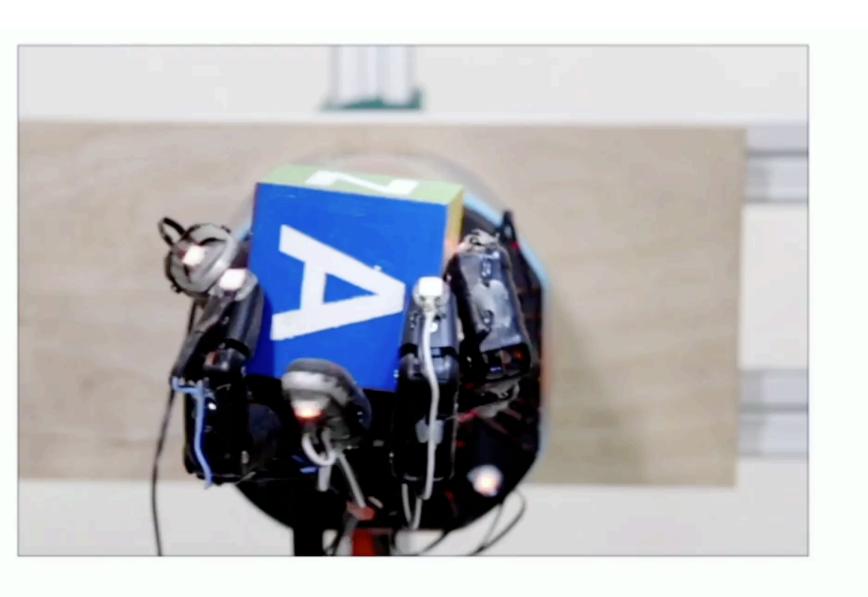
Source: Isola, Torralba, Freeman

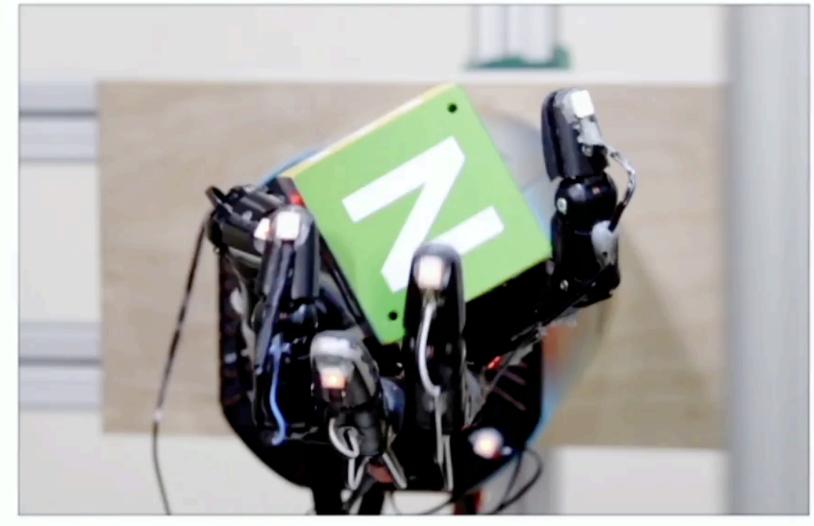
Idea #2: train on randomly perturbed data, so that test set just looks like another random perturbation



This is called domain randomization or data augmentation

OpenAl Dactyl







FINGER PIVOTING

SLIDING

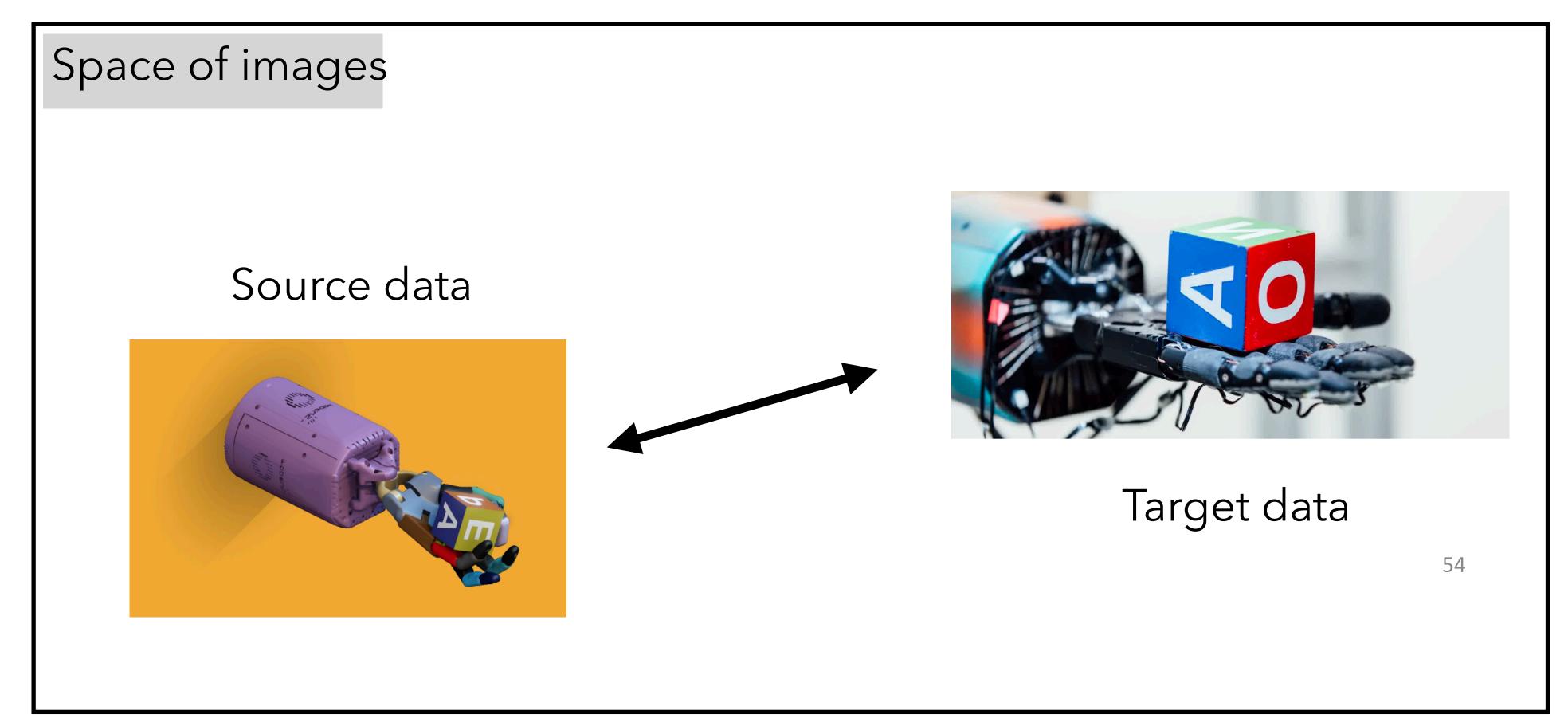
FINGER GAITING

source domain

target domain

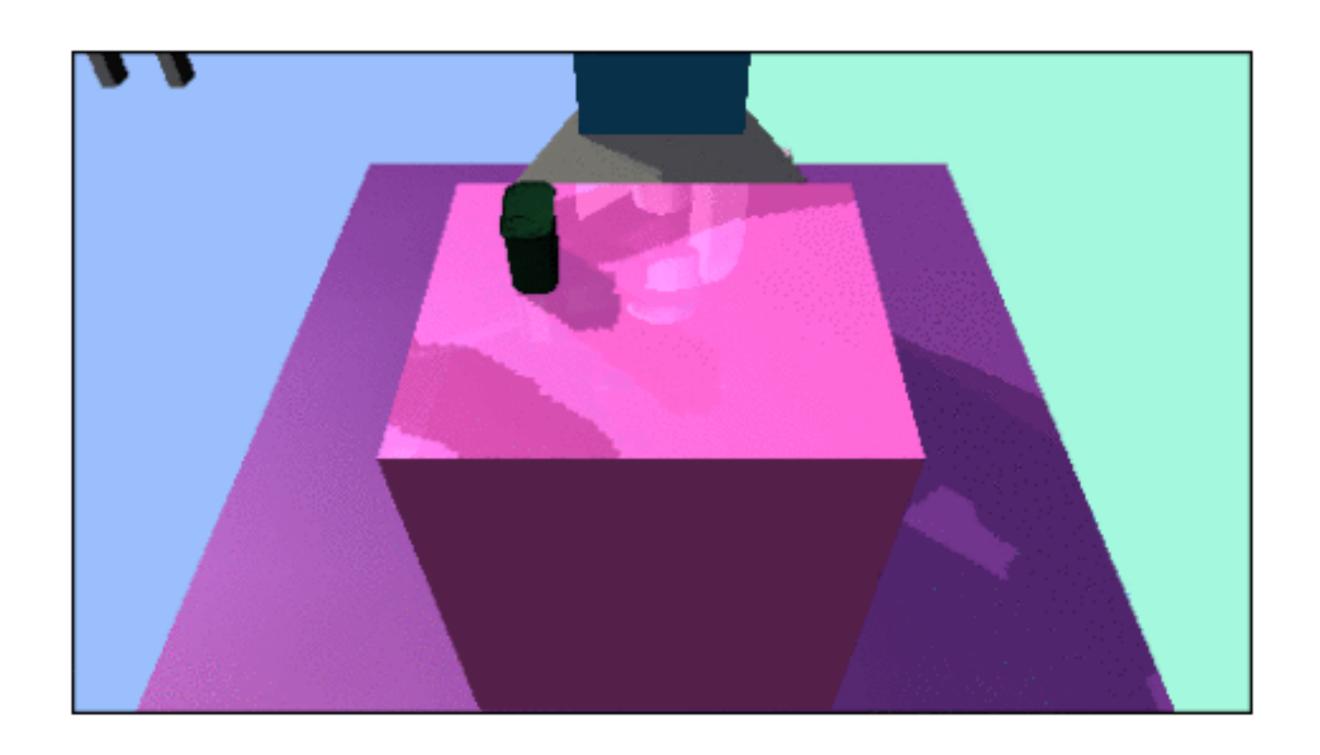
(where we actual use our model)

Domain gap between p_{source} and p_{target} will cause us to fail to generalize.

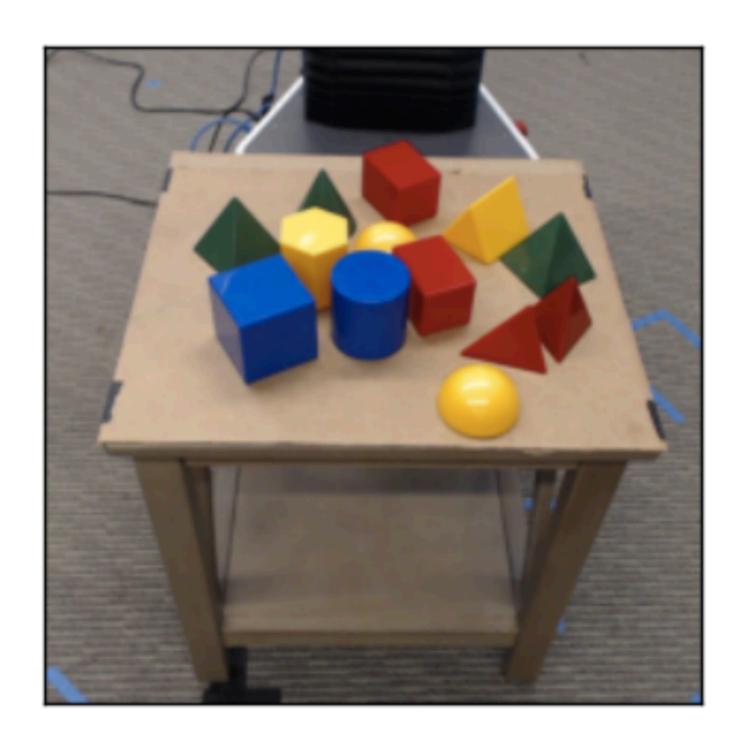


Domain randomization

Training data



Test data



[Sadeghi & Levine 2016] Above example is from [Tobin et al. 2017]

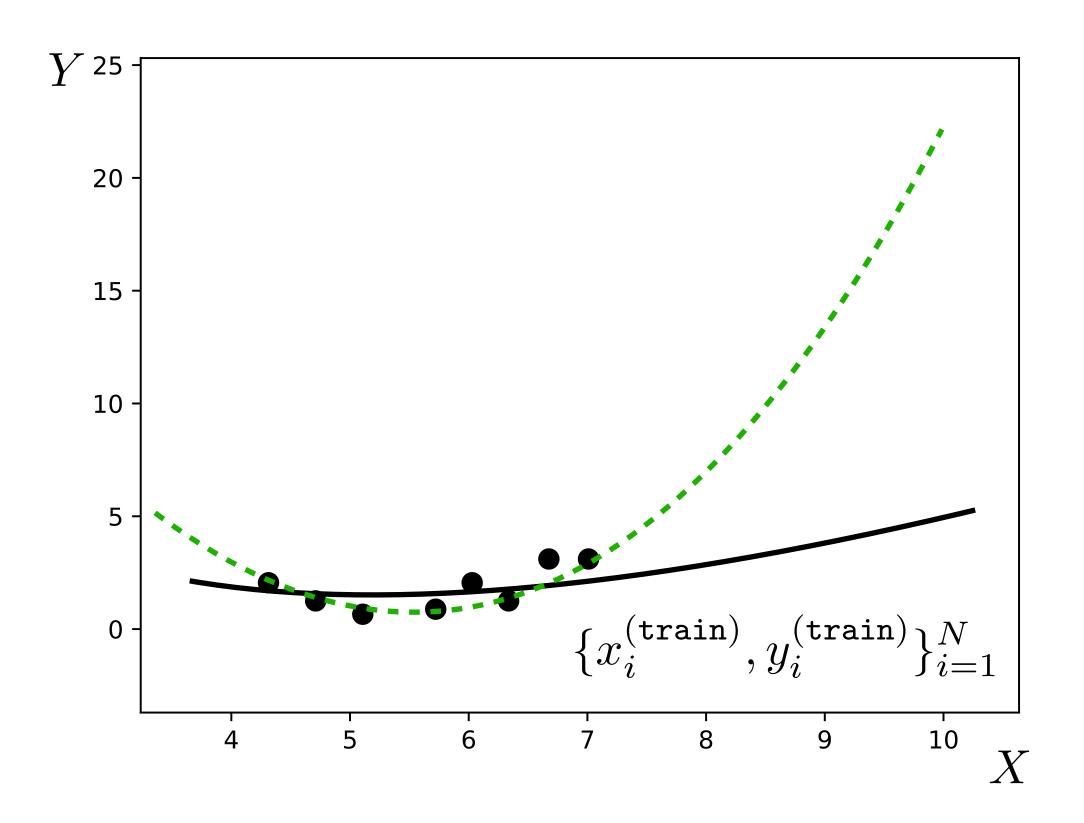
Other sources of bias, besides data

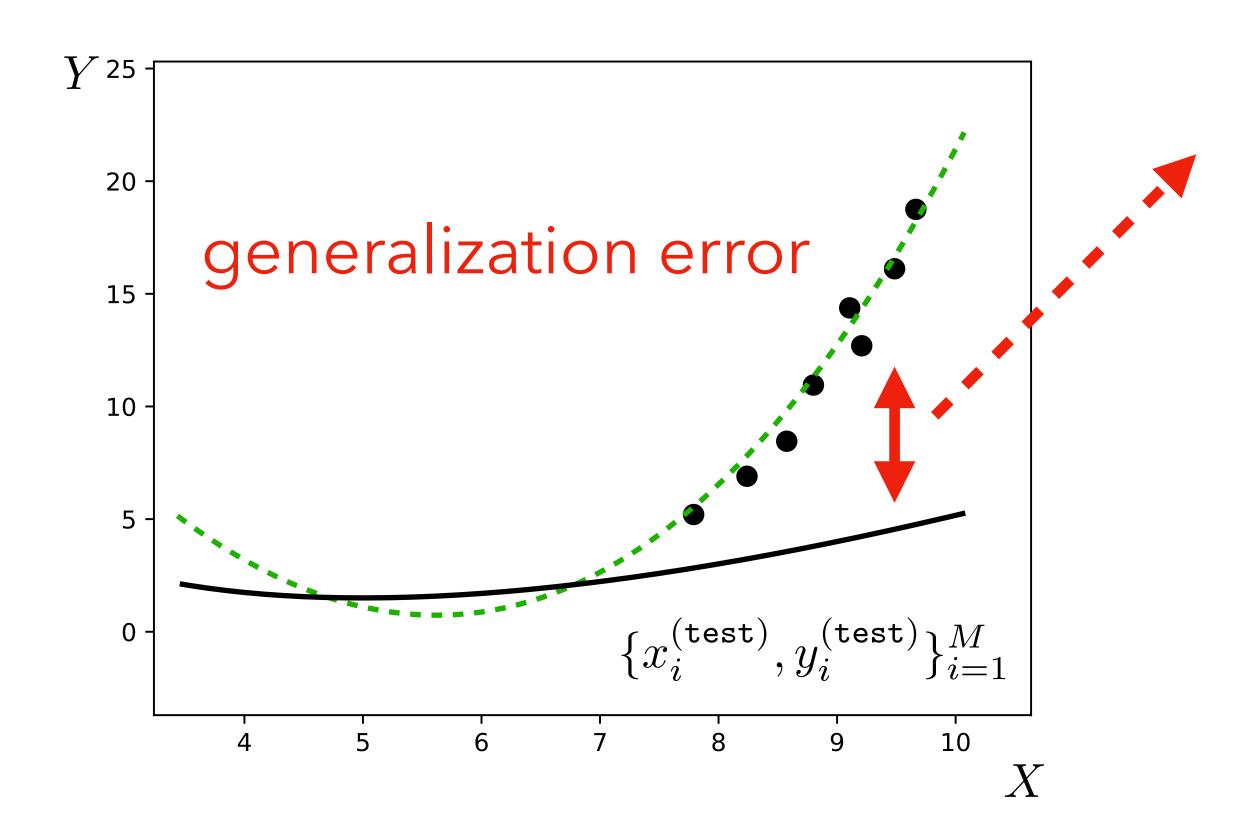
- Data very important [Maluleke et al., 2022], but also other factors can matter.
- Camera hardware and software
 - e.g., default camera settings calibrated to expose light skin
- Loss function (e.g., "mode collapse" in GANs)
- Features
- Sampling strategy (e.g., truncation in GANs)

What if we go way outside of the training distribution?

Training data

Test data

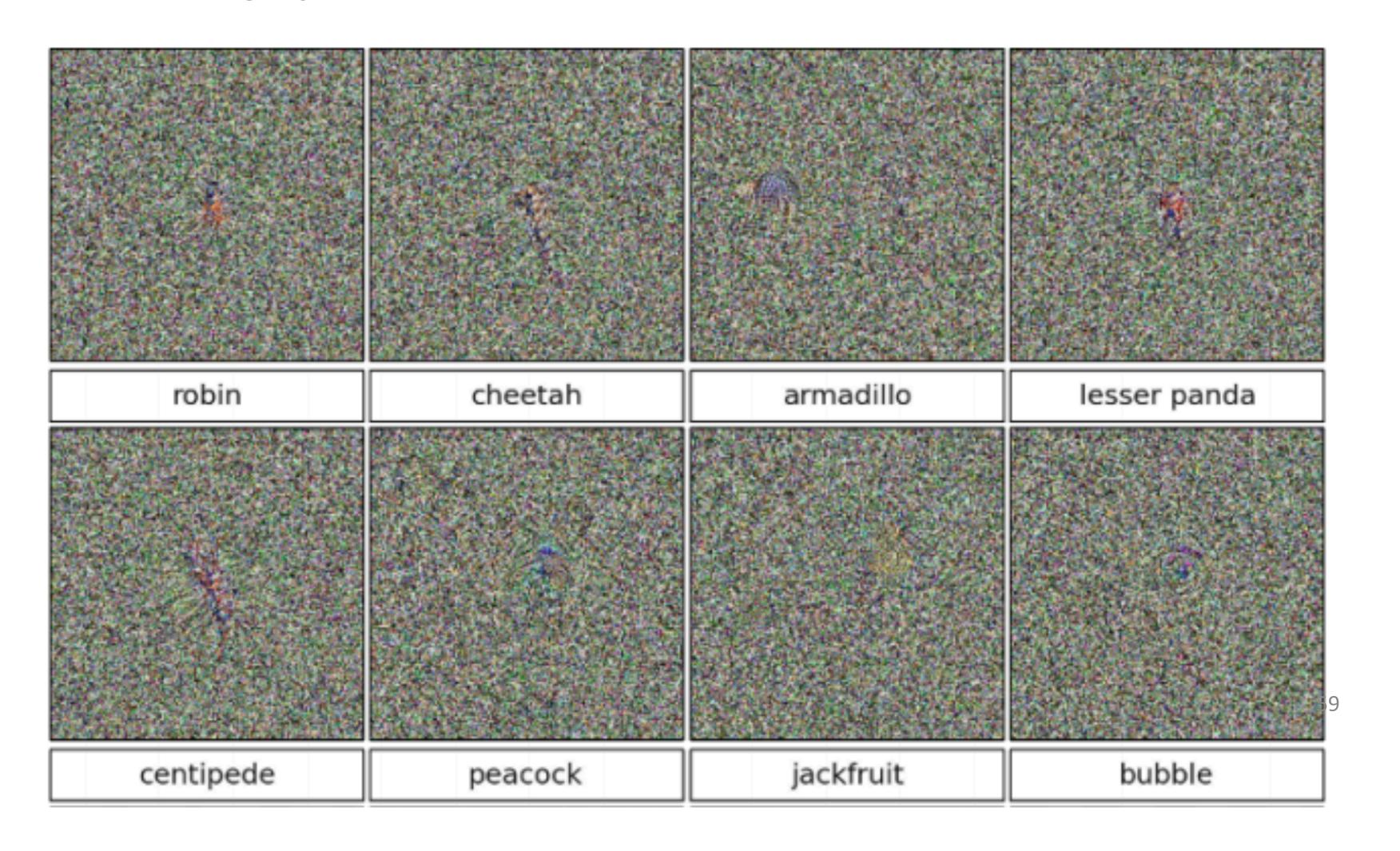




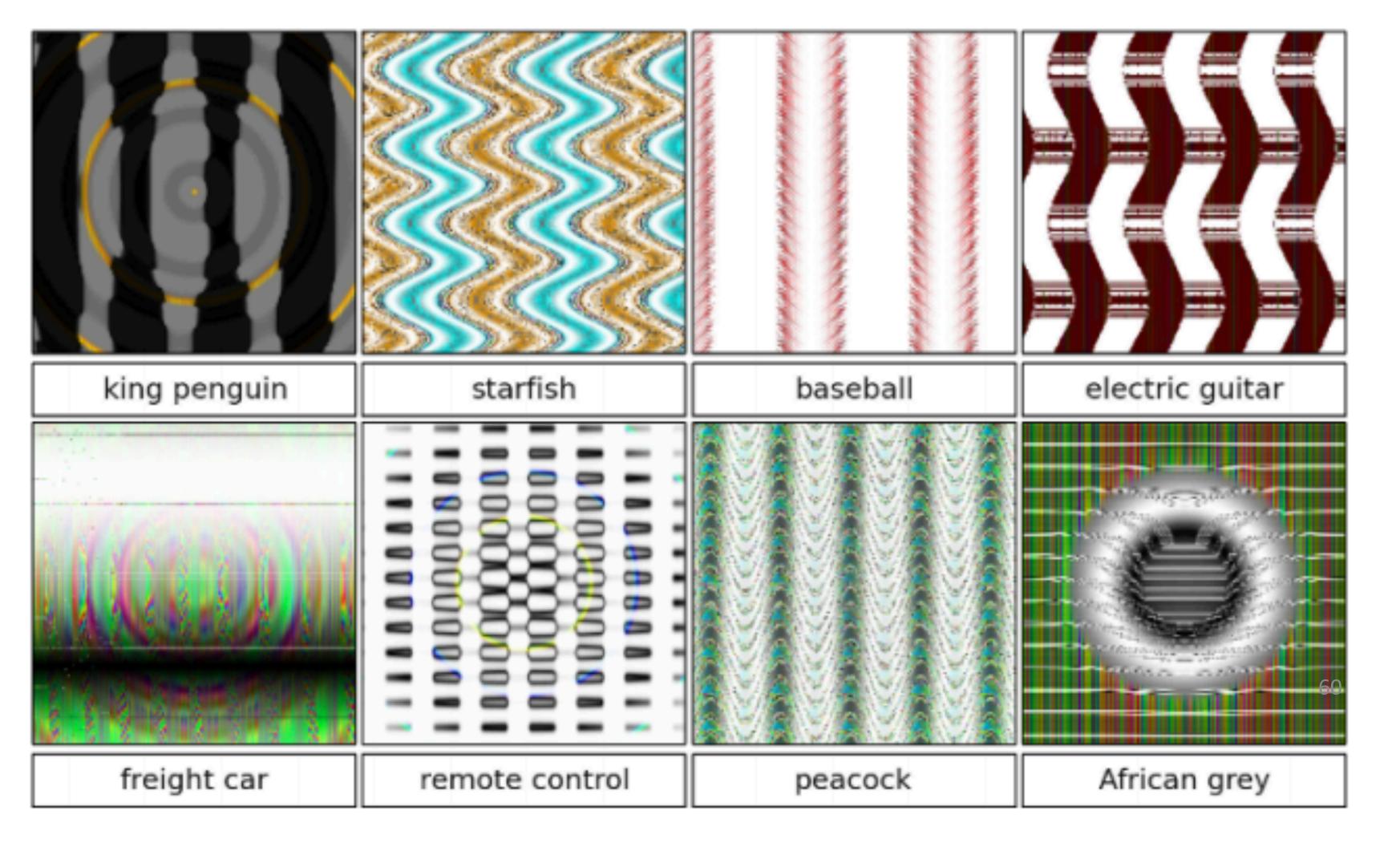
Our training data did not cover the part of the distribution that was tested

(biased data)

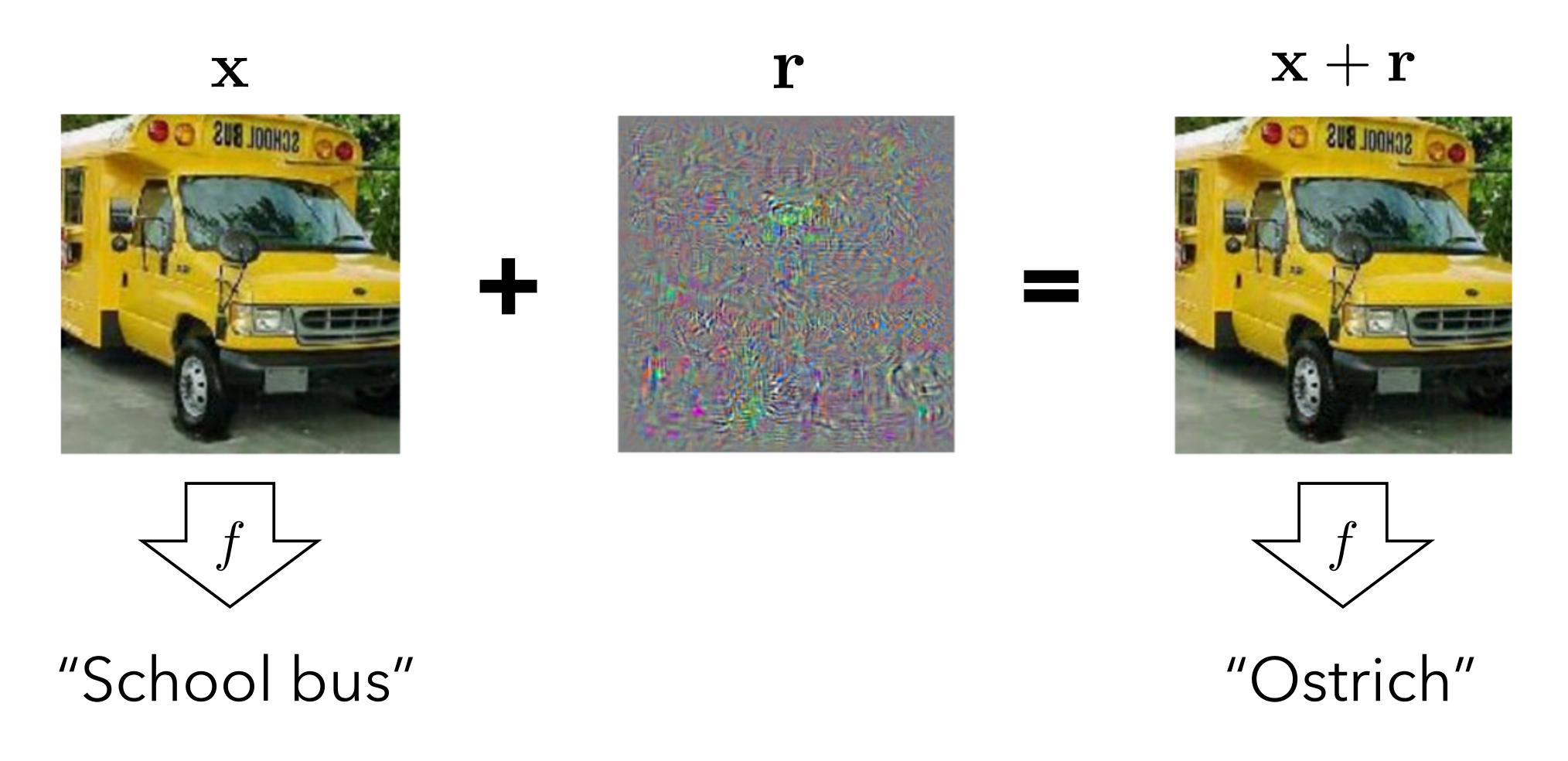
"Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images" [Nguyen, Yosinski, and Clune, CVPR 2015]



"Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images" [Nguyen, Yosinski, and Clune, CVPR 2015]



Adversarial noise



$$\operatorname{arg\,max} p(y = \operatorname{ostrich}|\mathbf{x} + \mathbf{r}) \quad \text{subject to} \quad ||\mathbf{r}|| < \epsilon$$

["Intriguing properties of neural networks", Szegedy et al. 2014

Anything to worry about?

"NO Need to Worry about Adversarial Examples in Object Detection in Autonomous Vehicles", Lu et al. 2017



(Early) 2017's attacks fail on physical objects, since they are optimized to attack a single view!

Anything to worry about?

Later in 2017...

"Synthesizing Robust Adversarial Examples", Athalye, Engstrom, Ilyas, Kwok, 2017

3D-printed **turtle** model classified as **rifle** from most viewpoints



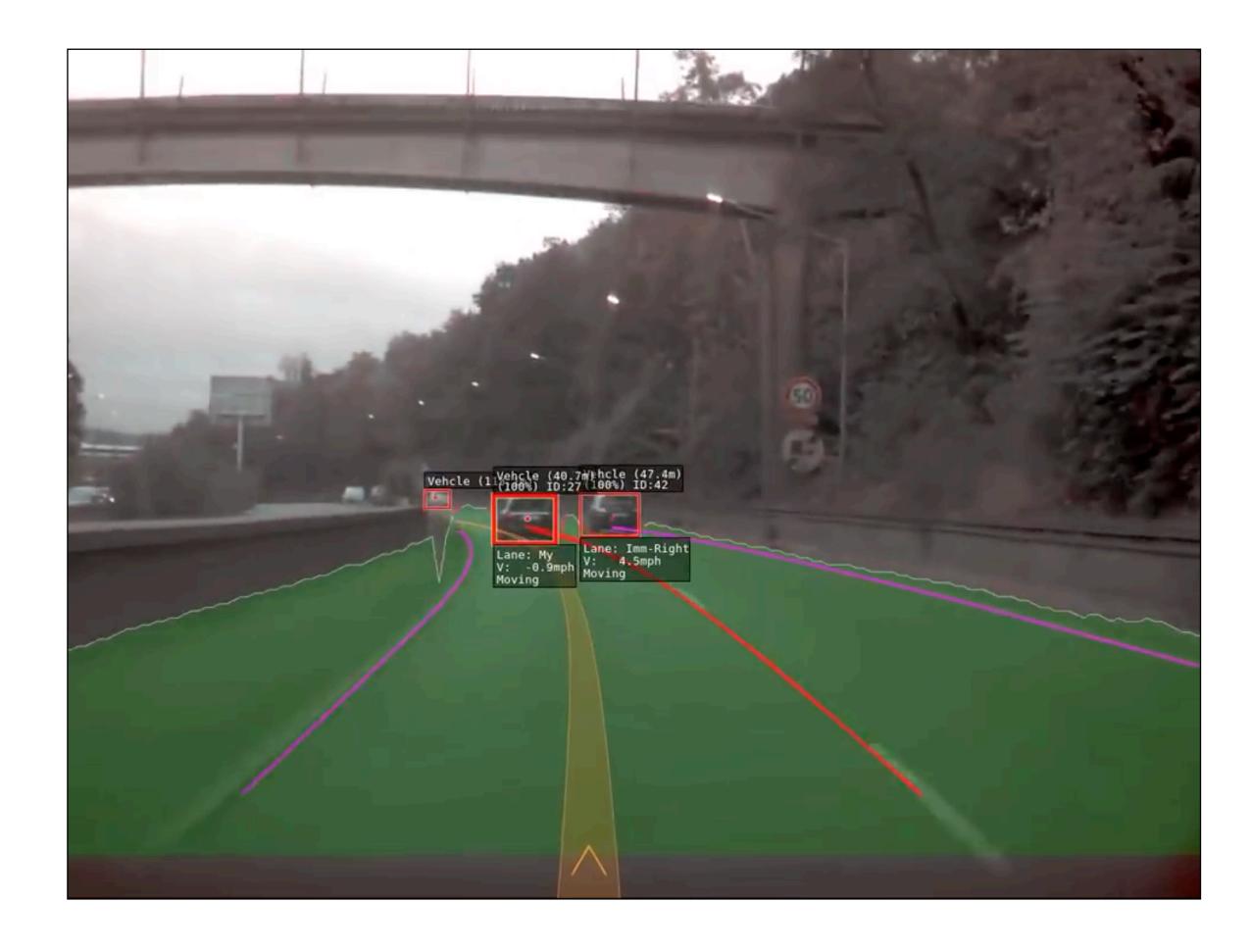
Adversarial examples

- Current deep models have bad worst-case performance
- Can be exploited by an adversary
- Few guarantees, can't fully trust what the model's output

Problems of applying computer vision in practice

Mission-critical computer vision systems



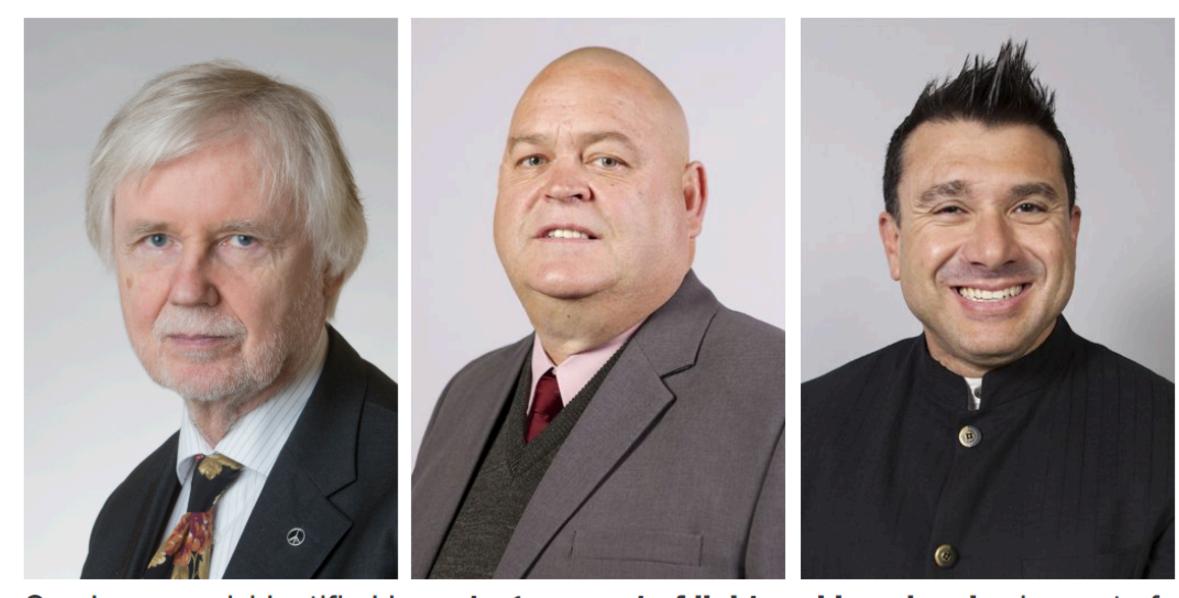


Source: https://qz.com/1402843/watch-teslas-autopilot-see-the-streets-of-paris/

Social consequences

Color Matters in Computer Vision

Facial recognition algorithms made by Microsoft, IBM and Face++ were more likely to misidentify the gender of black women than white men.



Gender was misidentified in **up to 1 percent of lighter-skinned males** in a set of 385 photos.



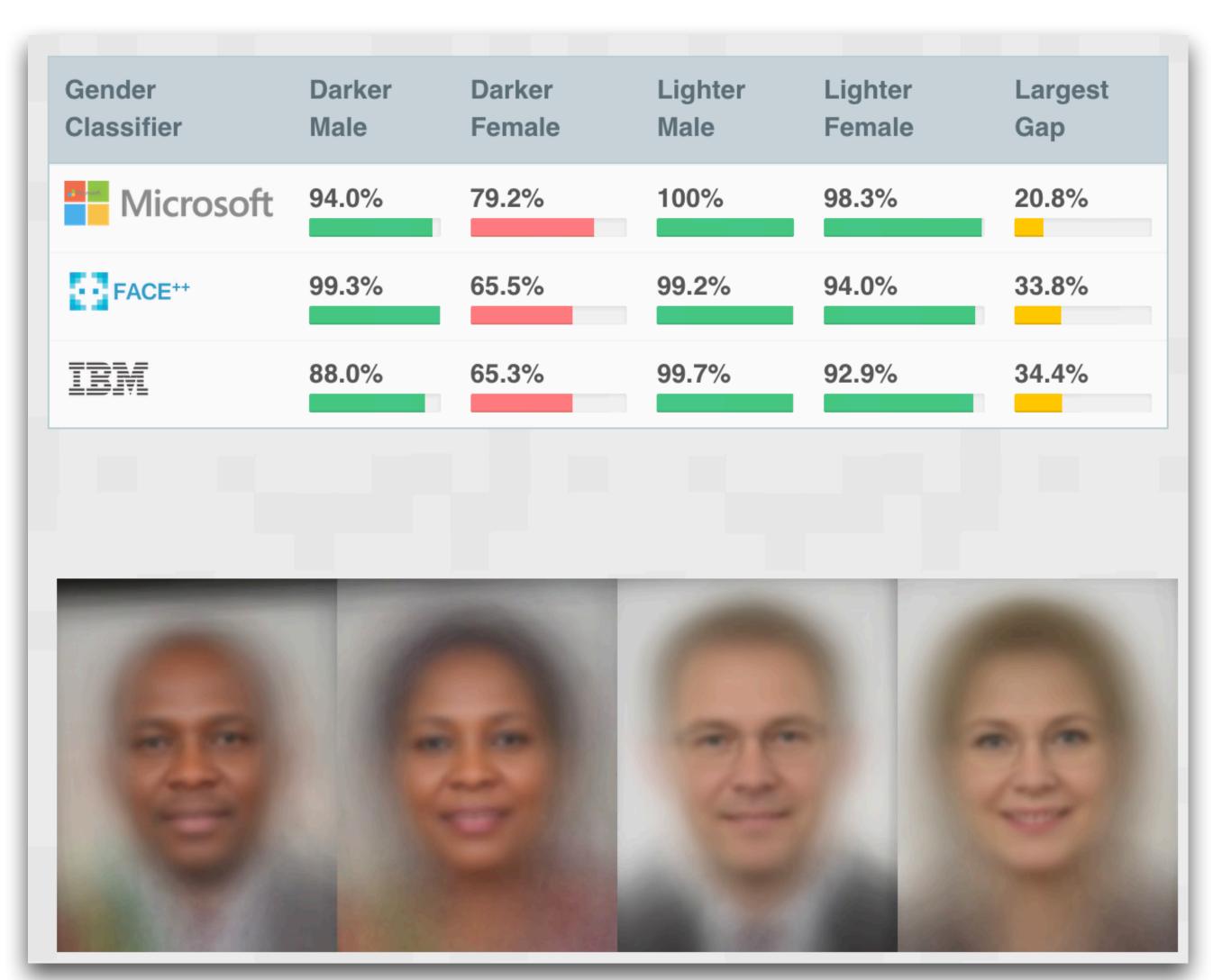
Gender was misidentified in **35 percent of darker-skinned females** in a set of 271 photos.

67

https://www.nytimes.com/2018/02/09/technology/facial-recognition-race-artificial-intelligence.html

Source: Isola, Torralba, Freeman

Algorithmic Bias



http://gendershades.org/overview.html

Proceedings of Machine Learning Research 81:1–15, 2018

Conference on Fairness, Accountability, and Transparency

Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification*

Joy Buolamwini

JOYAB@MIT.EDU

MIT Media Lab 75 Amherst St. Cambridge, MA 02139

Timnit Gebru

TIMNIT.GEBRU@MICROSOFT.COM

Microsoft Research 641 Avenue of the Americas, New York, NY 10011

Editors: Sorelle A. Friedler and Christo Wilson

Abstract

Recent studies demonstrate that machine learning algorithms can discriminate based on classes like race and gender. In this work, we present an approach to evaluate bias present in automated facial analysis algorithms and datasets with respect to phenotypic subgroups. Using the dermatologist approved Fitzpatrick Skin Type classification system, we characterize the gender and skin type distribution of two facial analysis benchmarks, IJB-A and Adience. We find that these datasets are overwhelmingly composed of lighter-skinned subjects (79.6% for IJB-A and 86.2% for Adience) and introduce a new facial analysis dataset which is balanced by gender and skin type. We evaluate 3 commercial gender classification systems using our dataset and show that darker-skinned females are the most misclassified group (with error rates of up to 34.7%). The maximum error rate for lighter-skinned males is 0.8%. The substantial disparities in the accuracy of classifying darker females, lighter females, darker males, and lighter males in gender classification systems require urgent attention if commercial companies are to build genuinely fair, transparent and accountable facial analysis algorithms.

Keywords: Computer Vision, Algorithmic Audit, Gender Classification

1. Introduction

Artificial Intelligence (AI) is rapidly infiltrating every aspect of society. From helping determine

© 2018 J. Buolamwini & T. Gebru.

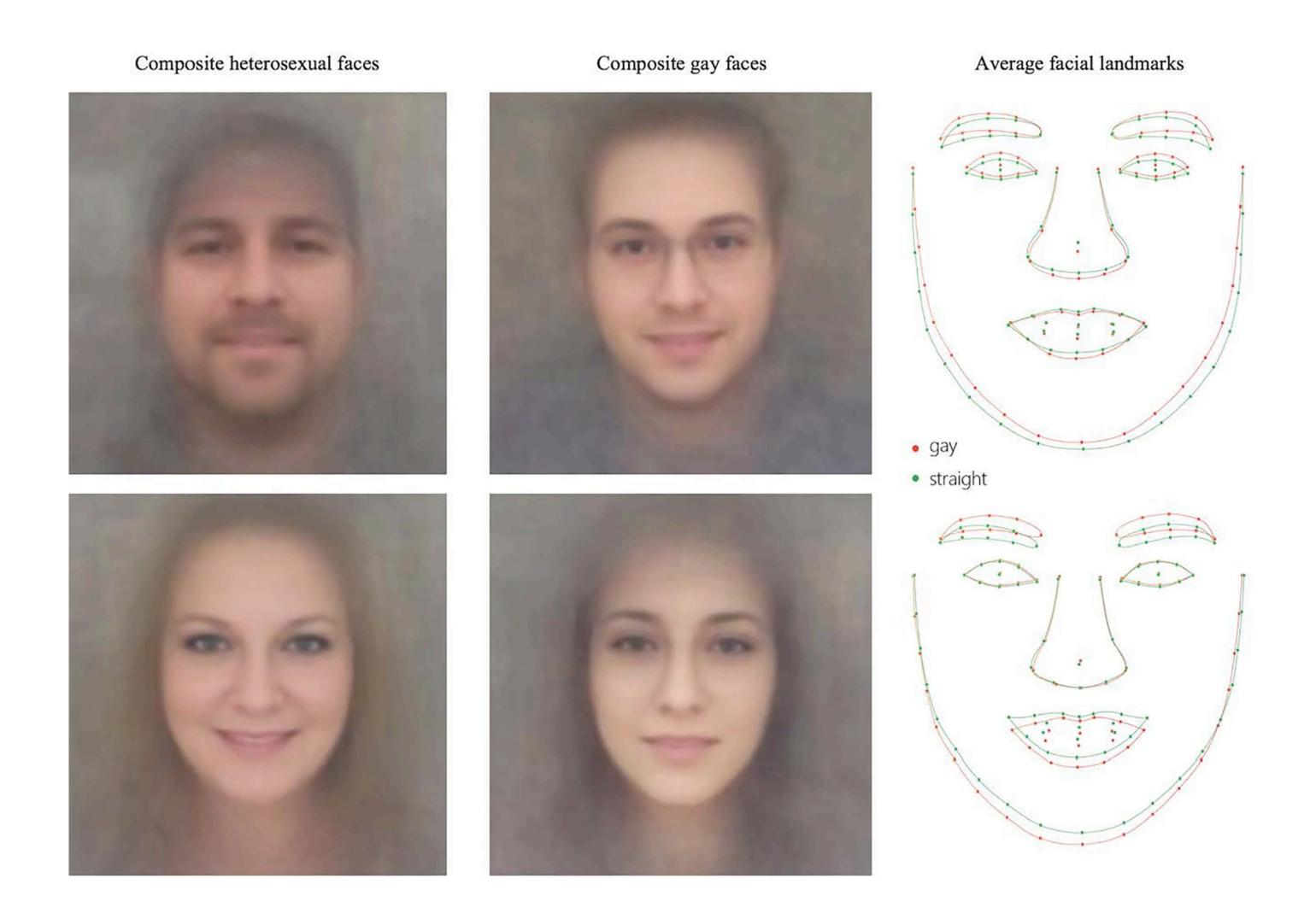
who is hired, fired, granted a loan, or how long an individual spends in prison, decisions that have traditionally been performed by humans are rapidly made by algorithms (O'Neil, 2017; Citron and Pasquale, 2014). Even AI-based technologies that are not specifically trained to perform highstakes tasks (such as determining how long someone spends in prison) can be used in a pipeline that performs such tasks. For example, while face recognition software by itself should not be trained to determine the fate of an individual in the criminal justice system, it is very likely that such software is used to identify suspects. Thus, an error in the output of a face recognition algorithm used as input for other tasks can have serious consequences. For example, someone could be wrongfully accused of a crime based on erroneous but confident misidentification of the perpetrator from security video footage analysis.

Many AI systems, e.g. face recognition tools, rely on machine learning algorithms that are trained with labeled data. It has recently been shown that algorithms trained with biased data have resulted in algorithmic discrimination (Bolukbasi et al., 2016; Caliskan et al., 2017). Bolukbasi et al. even showed that the popular word embedding space, Word2Vec, encodes societal gender biases. The authors used Word2Vec to train an analogy generator that fills in missing words in analogies. The analogy man is to computer programmer as woman is to "X" was completed with "homemaker", conforming to the stereotype that programming is associated with men and homemaking with women. The biases in Word2Vec are thus likely to be propagated throughout any system that uses this embedding.

http://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf

^{*} Download our gender and skin type balanced PPB dataset at gendershades.org

Bad choice of data



https://www.nytimes.com/2017/10/09/science/stanford-sexual-orientation-study.html

Face recognition in the U.S.

recode

Here's where the US government is using facial recognition technology to surveil Americans

This map shows how widespread the use of facial recognition technology has become.

By Shirin Ghaffary and Rani Molla | Updated Dec 10, 2019, 8:00am EST



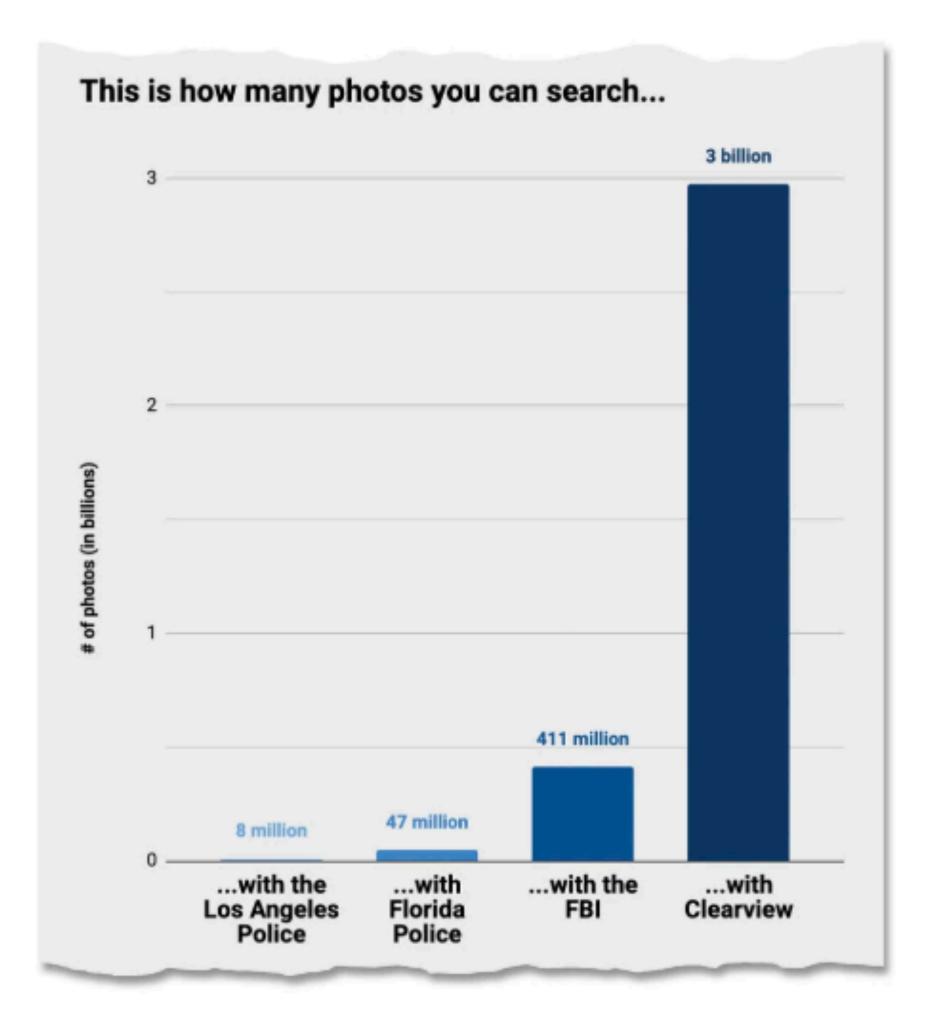
Source: S. Lazebnik

Fears of universal mass surveillance (and dubious claims)

The New York Times

The Secretive Company That Might End Privacy as We Know It

A little-known start-up helps law enforcement match photos of unknown people to their online images — and "might lead to a dystopian future or something," a backer says.



A chart from marketing materials that Clearview provided to law enforcement. Clearview

https://www.nytimes.com/2020/01/18/technology/ciearview-privacy-iaciai-recognition.numi

https://www.buzzfeednews.com/article/ryanmac/clearview-ai-nypd-facial-recognition

Source: S. Lazebnik

Datasets: Privacy, consent issues



Brainwash Dataset Analysis

2015 Head detection 11,917 images



Duke MTMC Dataset Analysis

2016
Person re-identification, multi-camera tracking
2,000,000 images



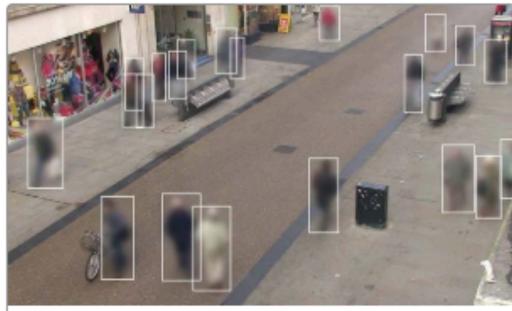
MegaFace Dataset Analysis

2016 face recognition 4,753,520 images



Microsoft Celeb Dataset Analysis

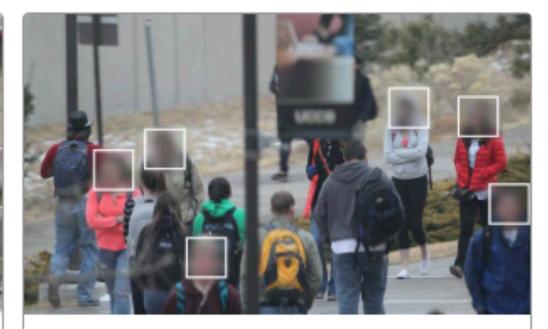
2016 Face recognition 8,200,000 images



Oxford Town Centre Dataset Analysis

2009

Person detection, gaze estimation



UnConstrained College Students Dataset Analysis

2016

Face recognition, face detection 16,149 images

https://megapixels.cc/

https://www.ft.com/content/cf19b956-60a2-11e9-b285-3acd5d43599e

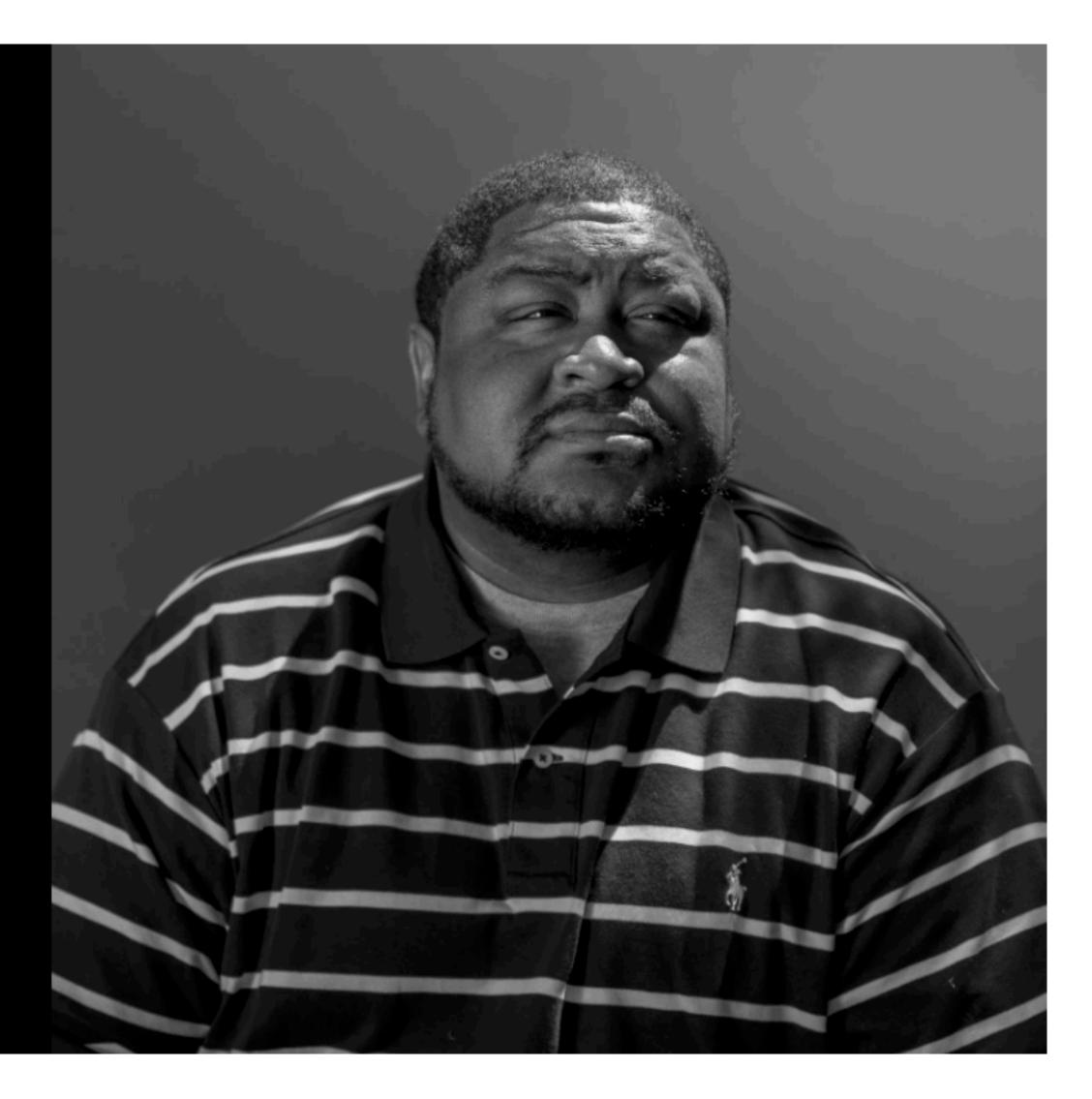
See also: https://www.theregister.co.uk/2020/01/27/ibms_facial_recognition_software_gets_it_in_trouble_again/

Source: S. Lazebnik

Face recognition in the U.S.

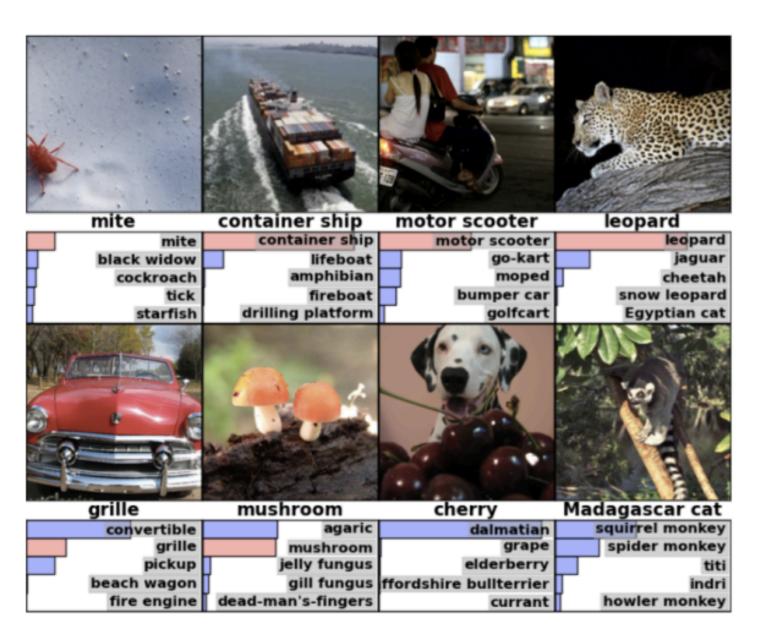
Wrongfully Accused by an Algorithm

In what may be the first known case of its kind, a faulty facial recognition match led to a Michigan man's arrest for a crime he did not commit.



ImageNet: asset or liability?

Performance on the basic ILSVRC benchmark has saturated



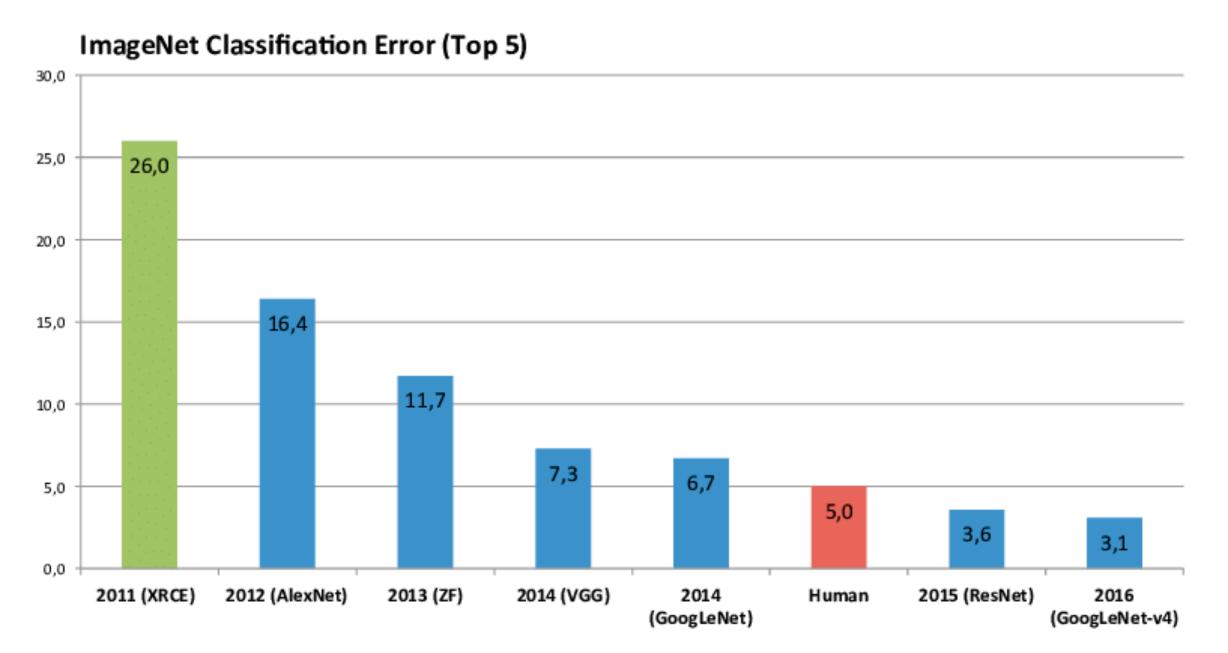


Figure source

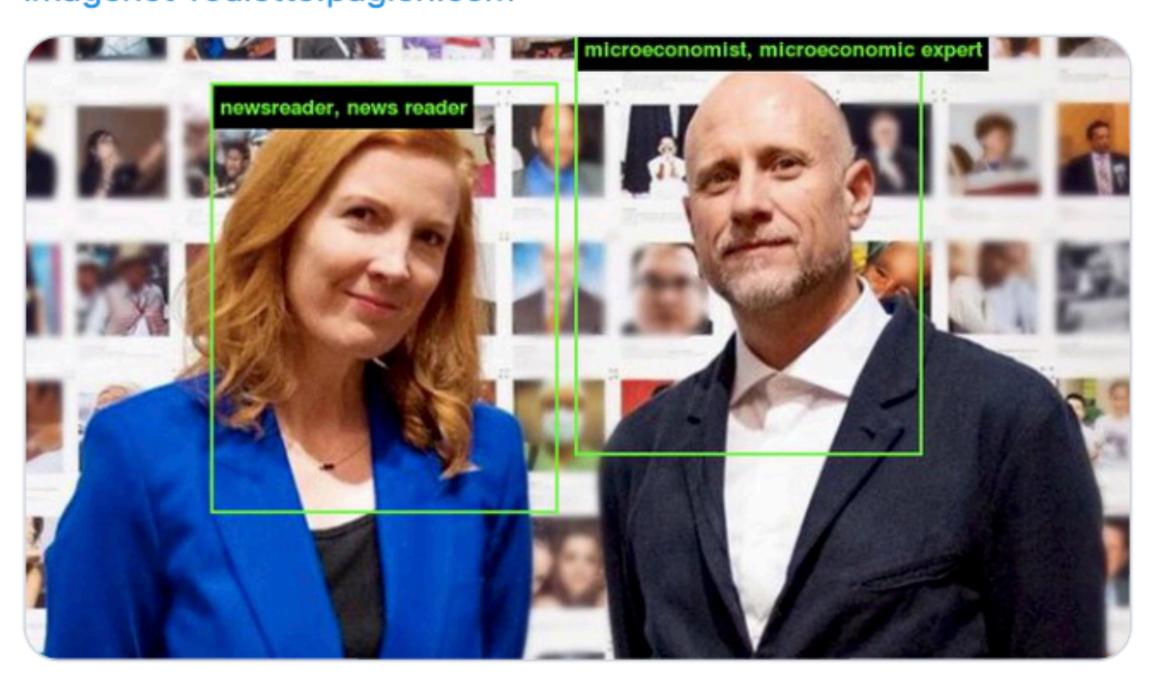
• Current models have reached levels of accuracy where the presence of human labeling error is starting to affect experimental conclusions (<u>Beyer et al.</u> 2020, <u>Northcutt et al.</u> 2021)

ImageNet labeling problems: ImageNet Roulette



Kate Crawford @ @katecrawford · Sep 16, 2019

Want to see how an AI trained on ImageNet will classify you? Try ImageNet Roulette, based on ImageNet's Person classes. It's part of the 'Training Humans' exhibition by @trevorpaglen & me - on the history & politics of training sets. Full project out soon imagenet-roulette.paglen.com



ImageNet Roulette uses an open source Caffe deep learning framework (produced at UC Berkeley) trained on the images and labels in the "person" categories (which are currently 'down for maintenance'). Proper nouns and categories with less than 100 pictures were removed.

When a user uploads a picture, the application first runs a face detector to locate any faces. If it finds any, it sends them to the Caffe model for classification. The application then returns the original images with a bounding box showing the detected face and the label the classifier has assigned to the image. If no faces are detected, the application sends the entire scene to the Caffe model and returns an image with a label in the upper left corner.

ImageNet contains a number of problematic, offensive and bizarre categories - all drawn from WordNet. Some use misogynistic or racist terminology. Hence, the results ImageNet Roulette returns will also draw upon those categories. That is by design: we want to shed light on what happens when technical systems are trained on problematic training data. Al classifications of people are rarely made visible to the people being classified. ImageNet Roulette provides a glimpse into that process – and to show the ways things can go wrong.

K. Crawford and T. Paglen, <u>Excavating AI: The Politics of Training Sets for Machine Learning</u>, September 2019 https://www.theverge.com/tldr/2019/9/16/20869538/imagenet-roulette-ai-classifier-web-tool-object-image-recognition

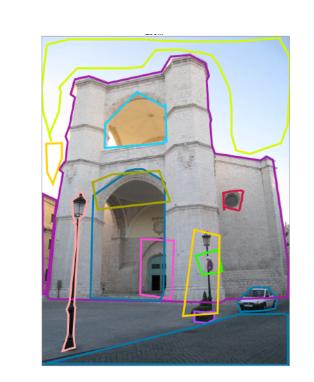
Source: S. Lazebnik

ImageNet Roulette



Some things to worry about...

Our datasets are often poorly labeled



And usually biased



 ML methods may perform well on lab-collected data, but often generalize poorly to real-world data

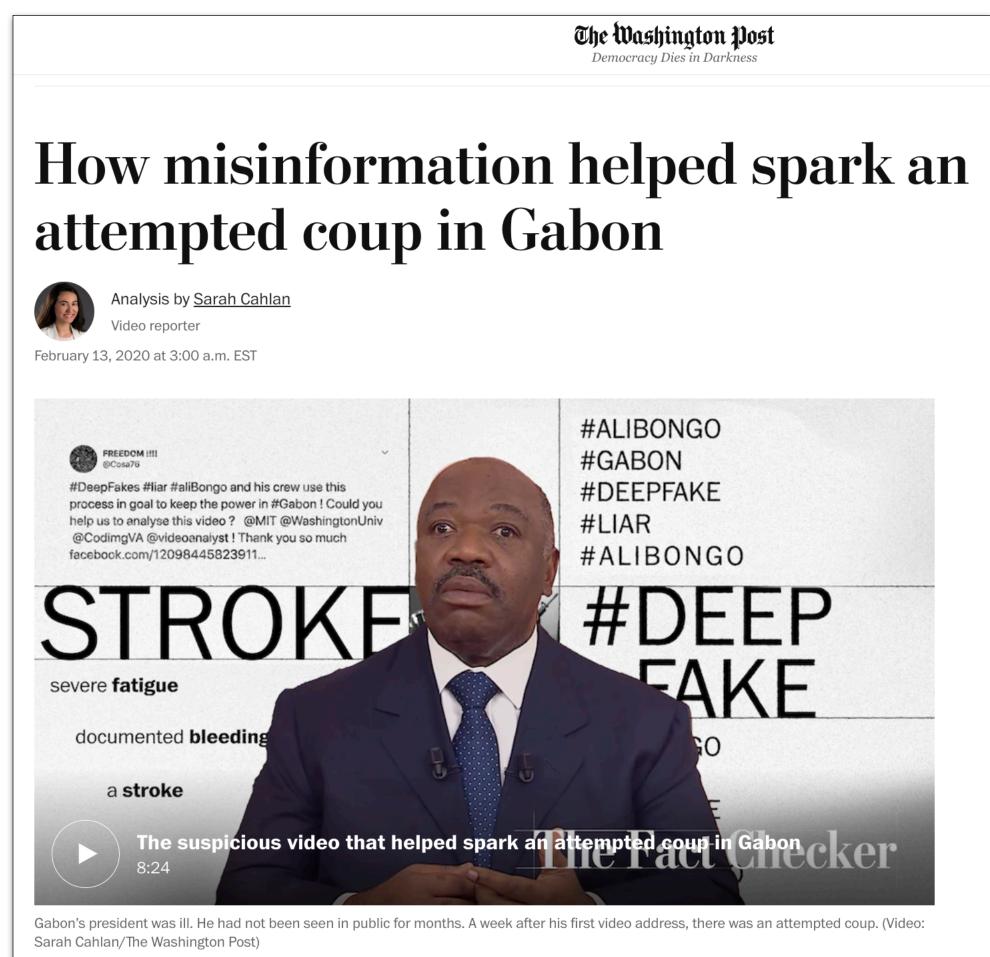
Can have negative social consequences

Open-ended discussion

- Supervised vs. unsupervised learning?
- Other negative consequences of computer vision systems?
- What other biases might computer vision systems have?

Fake images in the news





Text-to-image models make it easy

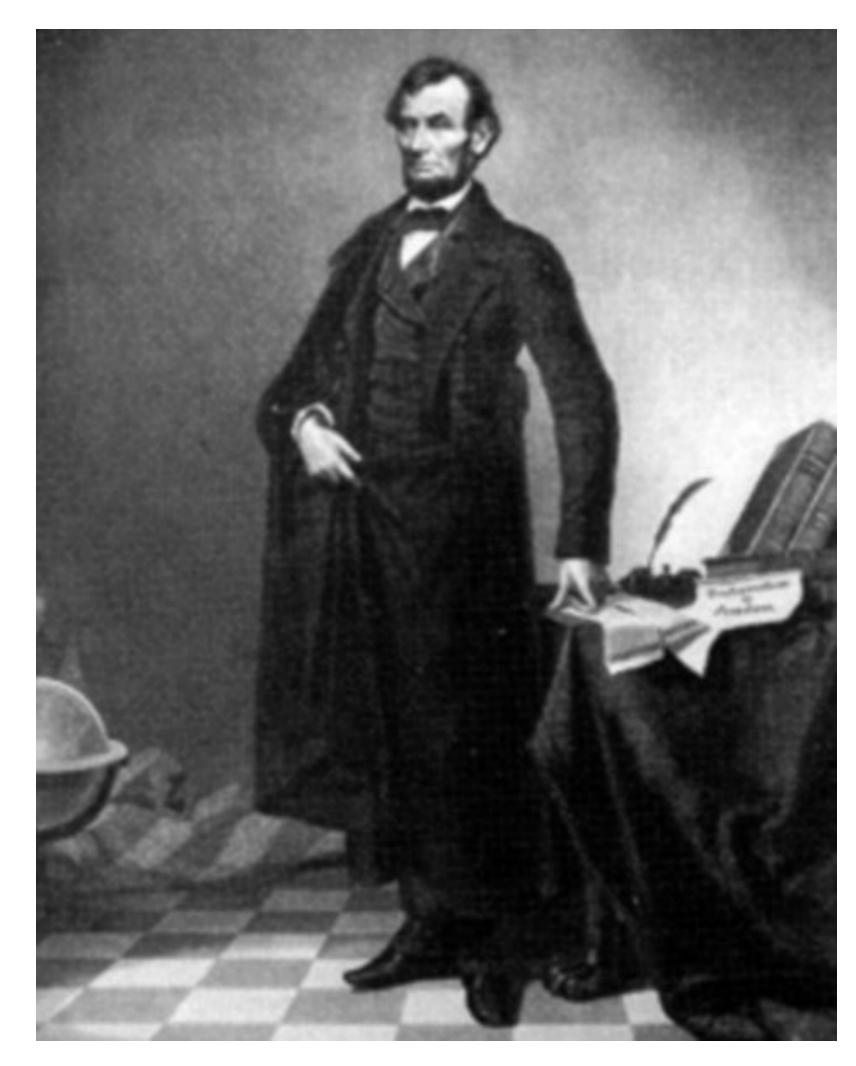






"Catholic Pope Francis wearing Balenciaga puffy jacket in drill rap music video, throwing up gang signs with hands, taken using a Canon EOS R camera with a 50mm f/1.8 lens, f/2.2 aperture, shutter speed 1/200s, ISO 100 and natural light, Full Body, Hyper Realistic Photography, Cinematic, Cinema, Hyperdetail, UHD, Color Correction, hdr, color grading, hyper realistic CG animation --ar 4:5 --upbeta --q 2 --v 5."

But image manipulation also has a long history



Abraham Lincoln?



John C. Calhoun

But image manipulation also has a long history



From Forrest Gump, 1994











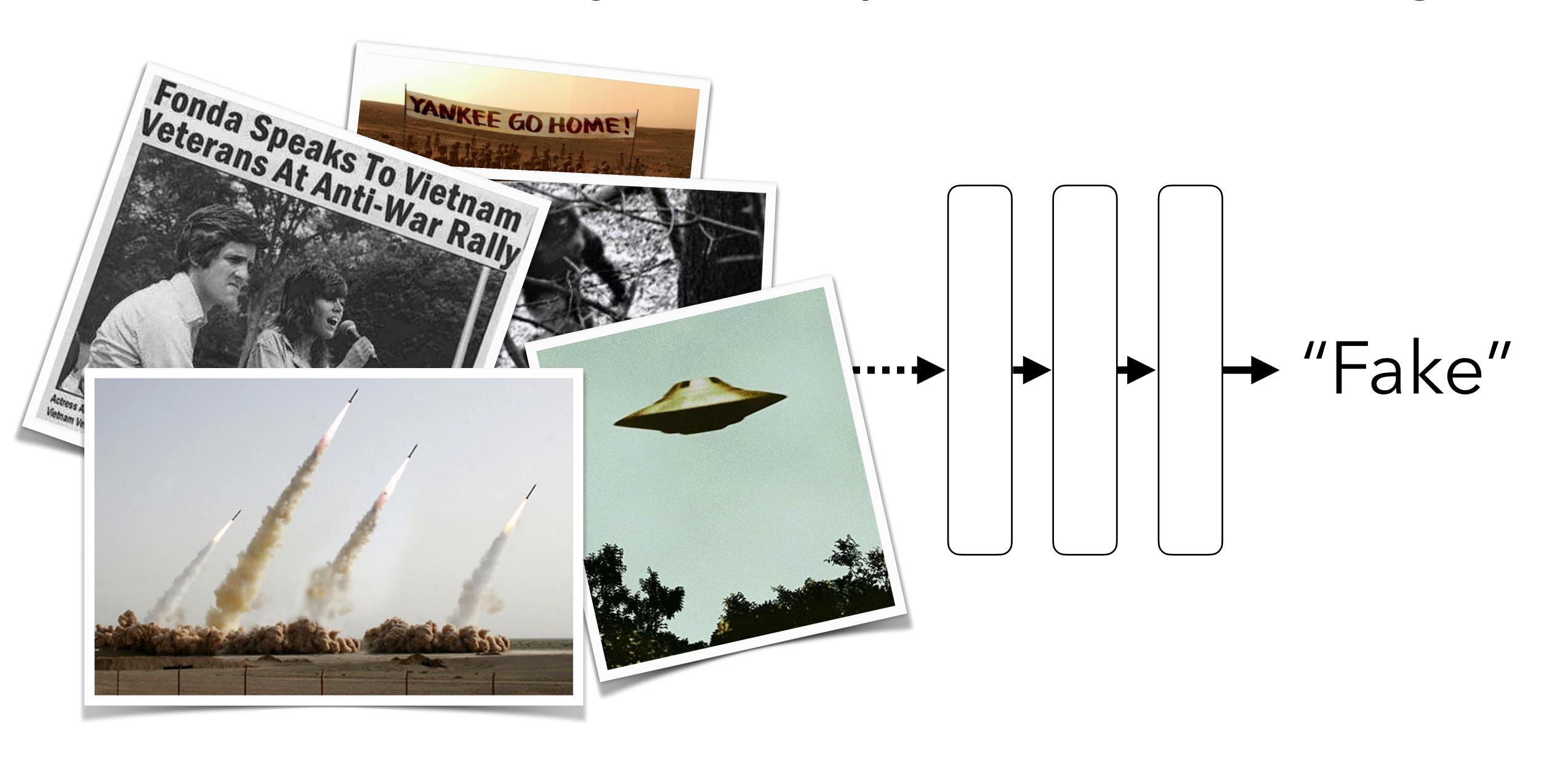




Detecting fake images



Hard to directly use supervised learning!



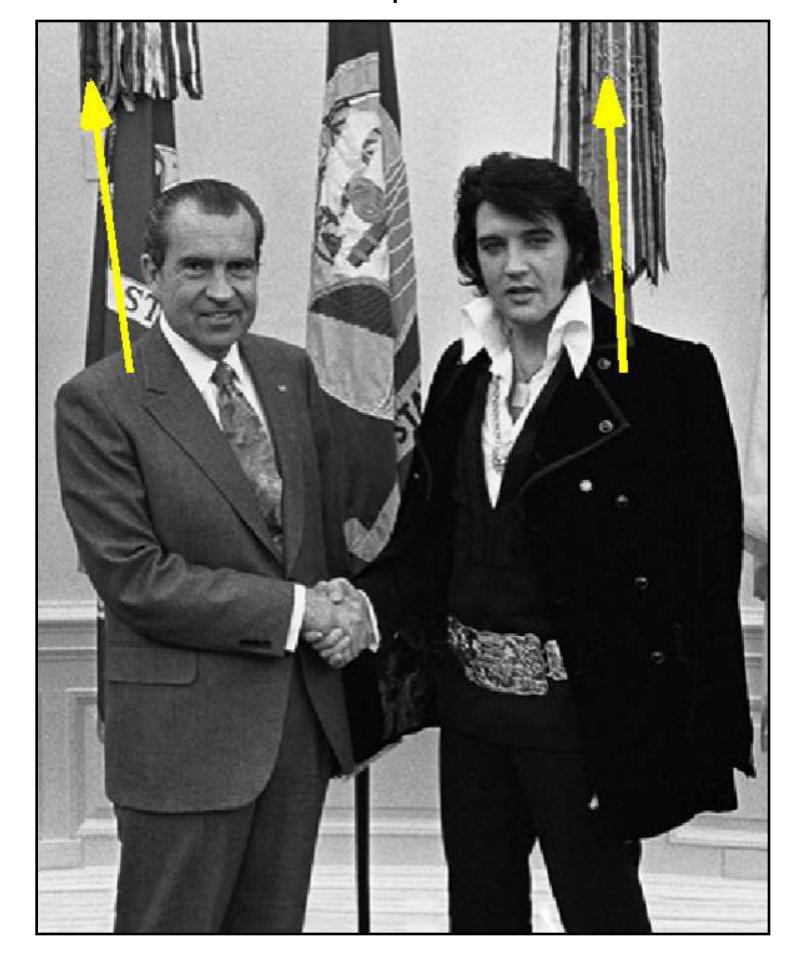
Strategy #1: physical models

Self-consistent lighting direction

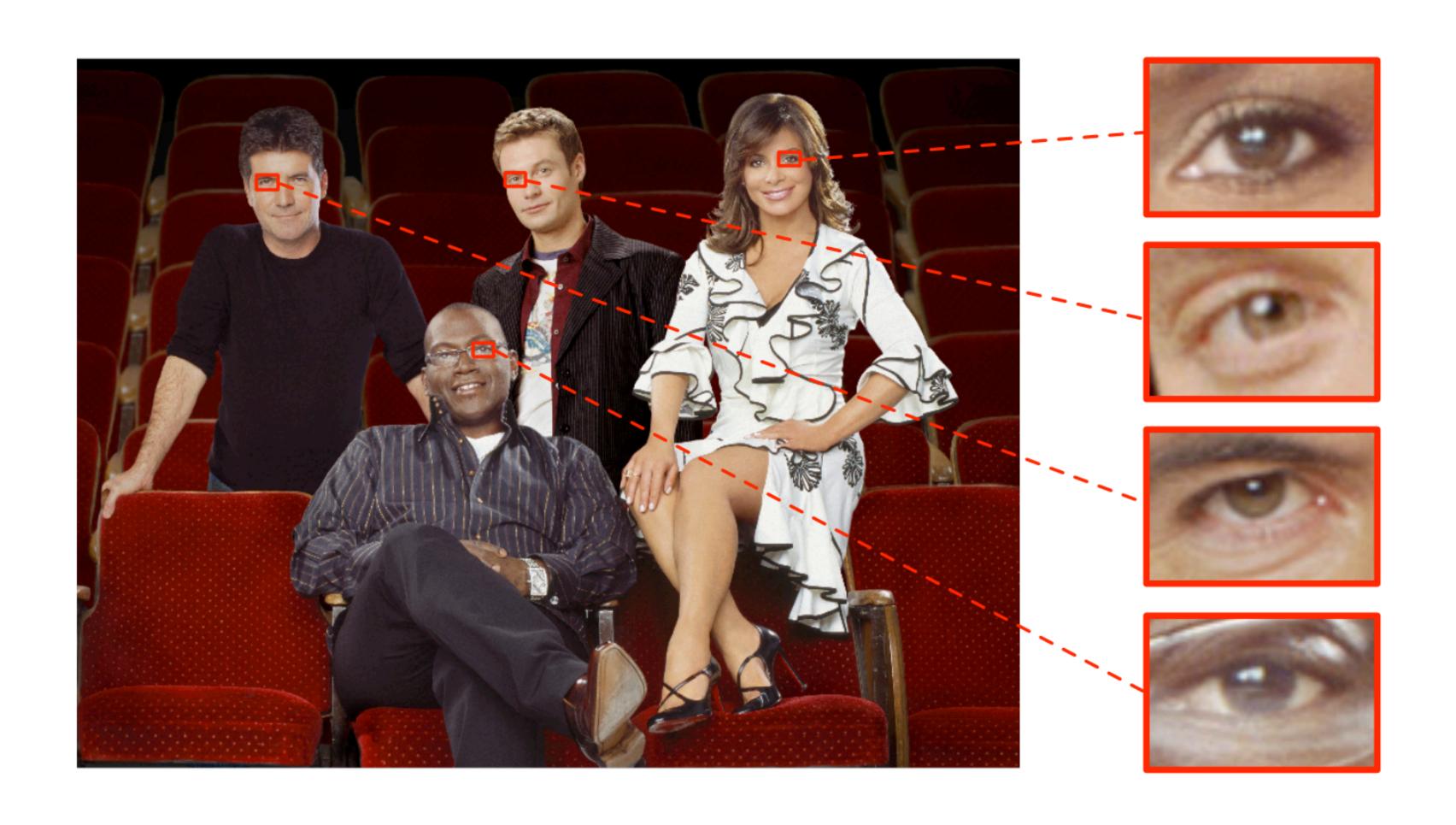
Fake photo



Real photo



Specular reflections



[Johnson and Farid, 2007]

Strategy #2: low-level imaging properties

JPEG artifacts

- Cameras vary in how they do JPEG compression.
- When you quantize a floating point numbers:
 - Some do round(), others do floor() or ceil()
- If a photo seems to have *both* kinds of quantization, it's probably a fake: e.g., a composite from images taken by different cameras!



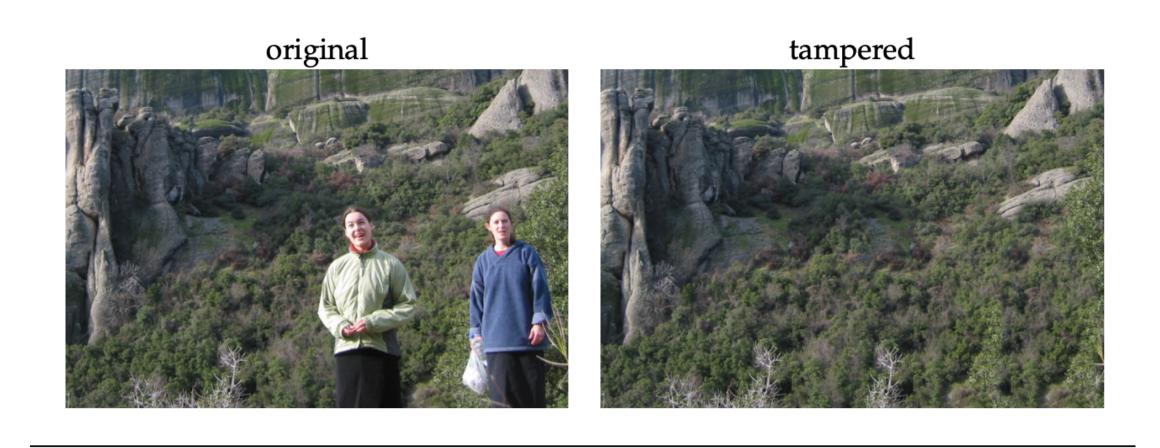


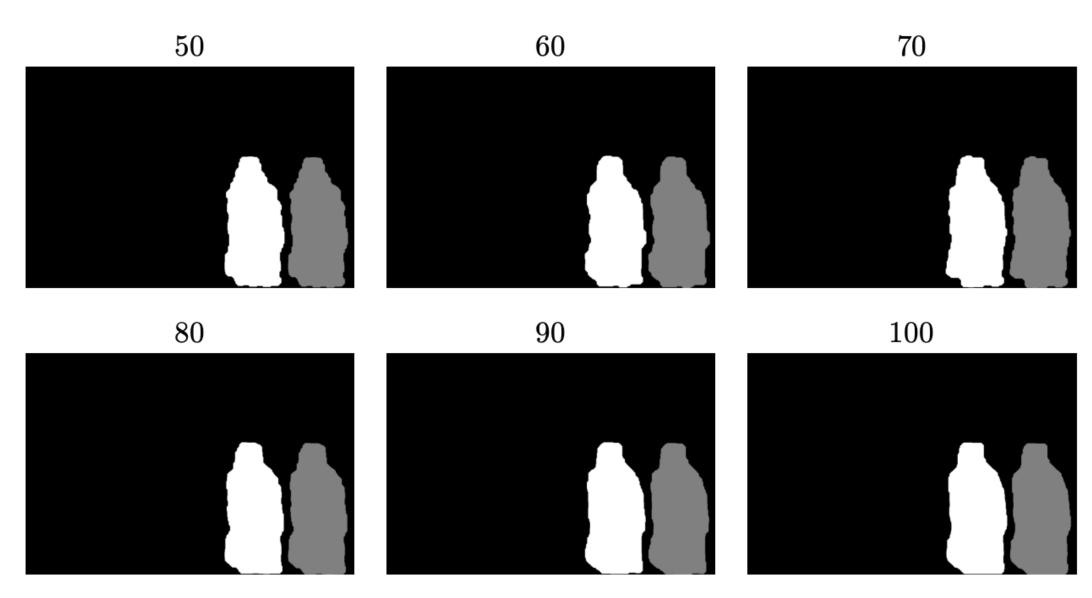




[Agarwal and Farid, "JPEG Dimples", 2017]

Detecting duplicated image regions





← amount of JPEG compression

- Traditional inpainting methods copy-and-paste image patches.
- Detect near-duplicated patches.
- But sensitive to postprocessing operations, like compression.

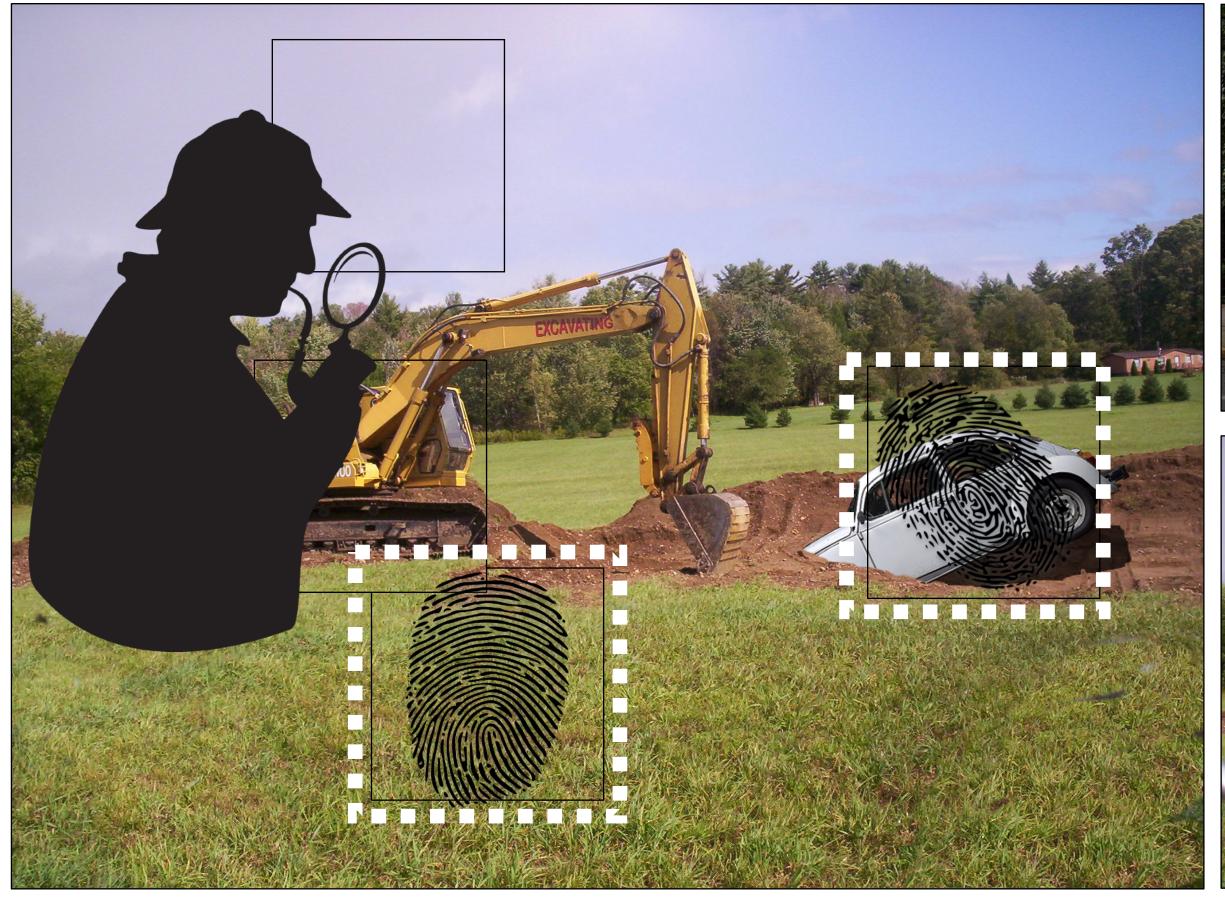
[Popescu and Farid, 2004]

Strategy #3: learned anomaly detection

Instead of hand-crafting cues, can we learn to detect "anomalous" images, and flag suspicious images?







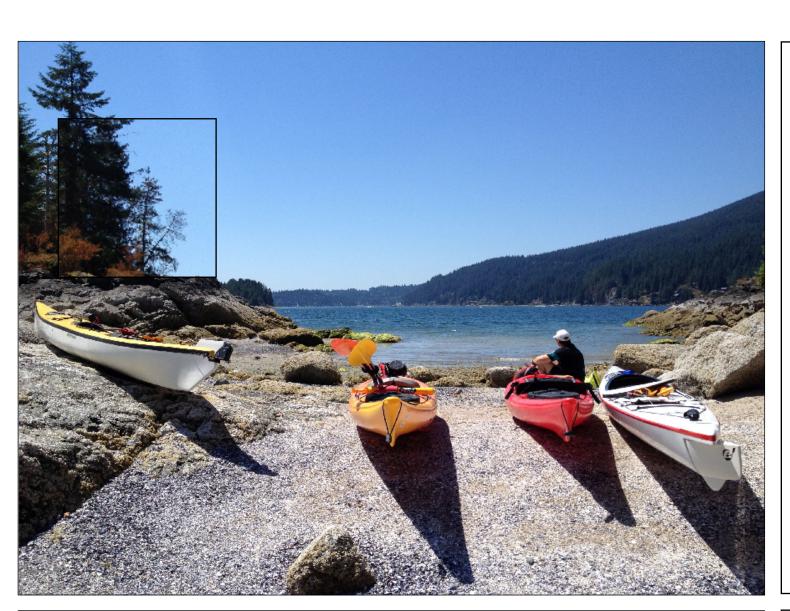




Inconsistent

Consistent

Predicting metadata consistency



CameraMake: Apple

CameraModel: iPhone 4s

ColorSpace: sRGB

ExifImageLength: 2448
ExifImageWidth: 3264

Flash: Flash did not fire

FocalLength: 107/2

WhiteBalance: Auto

ExposureTime: 1/2208

• • •



CameraMake: NIKON CORPORATION

CameraModel: NIKON D90

ColorSpace: sRGB

ExifImageLength: 2848
ExifImageWidth: 4288

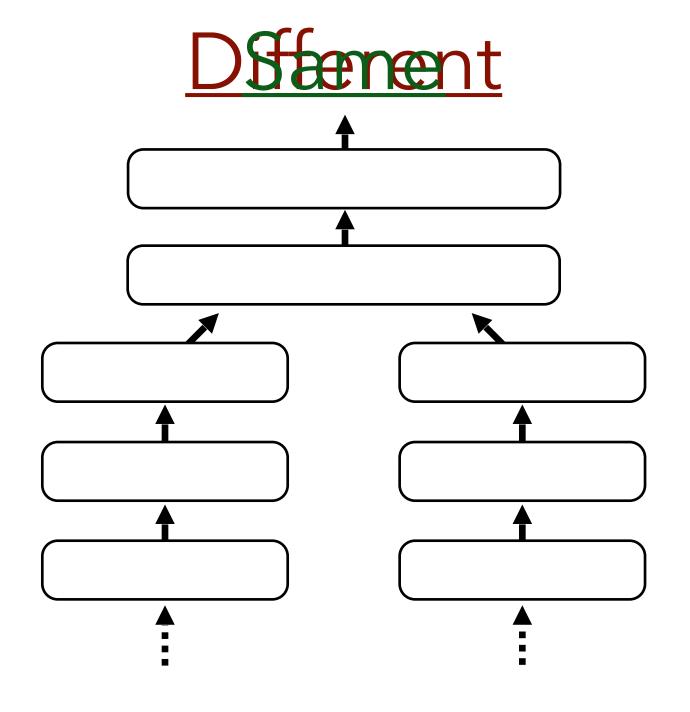
Flash: Flash did not fire

FocalLength: 18/796

WhiteBalance: Auto ExposureTime: 1/30

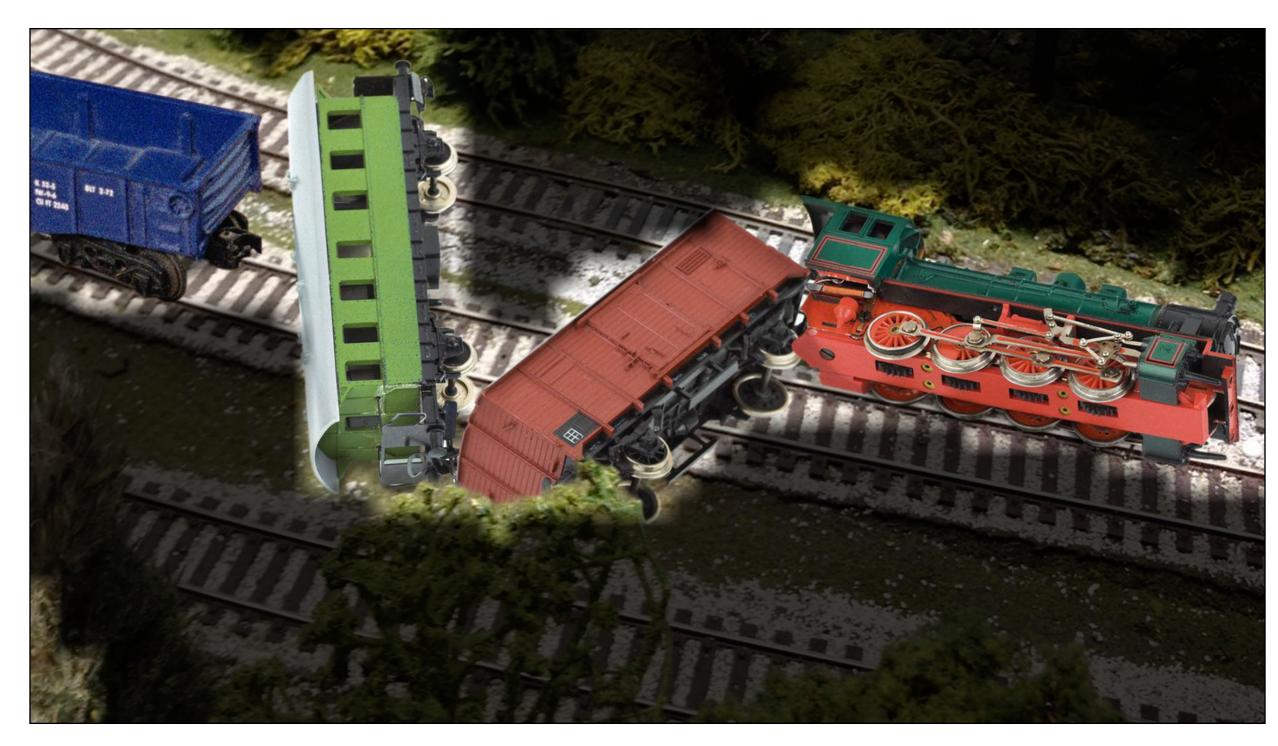
• •

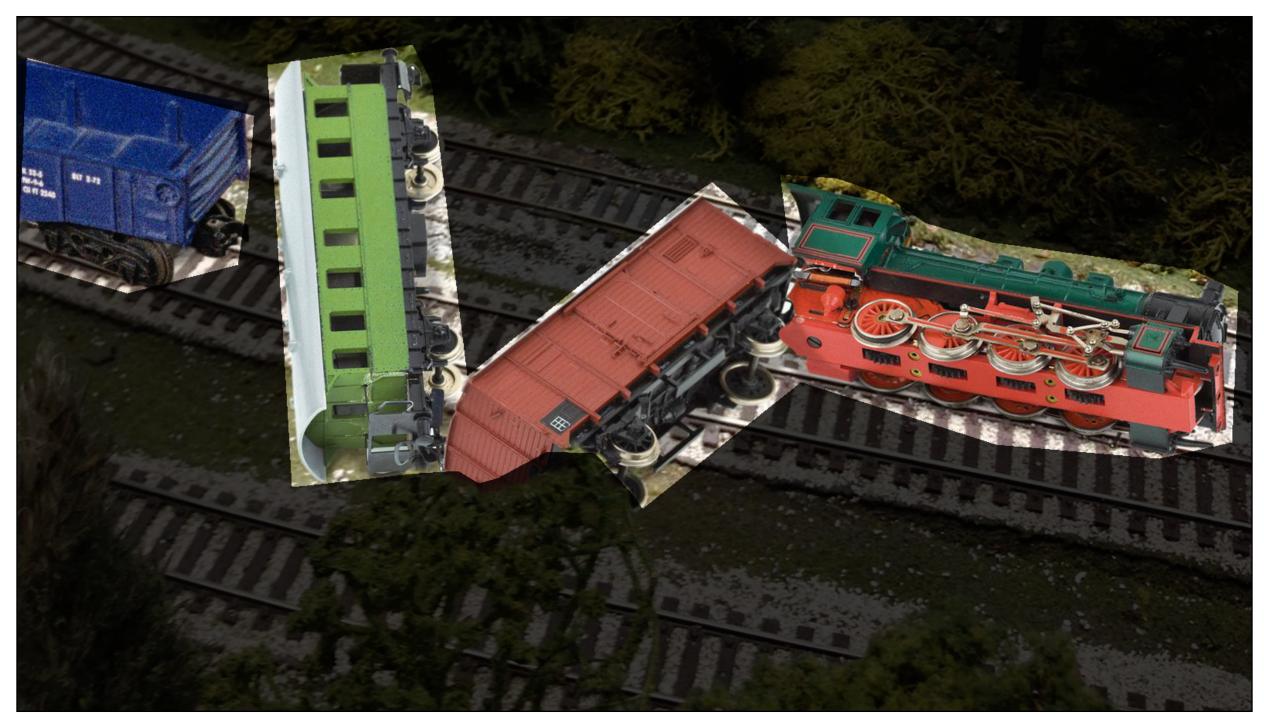
Same white balance?





Input

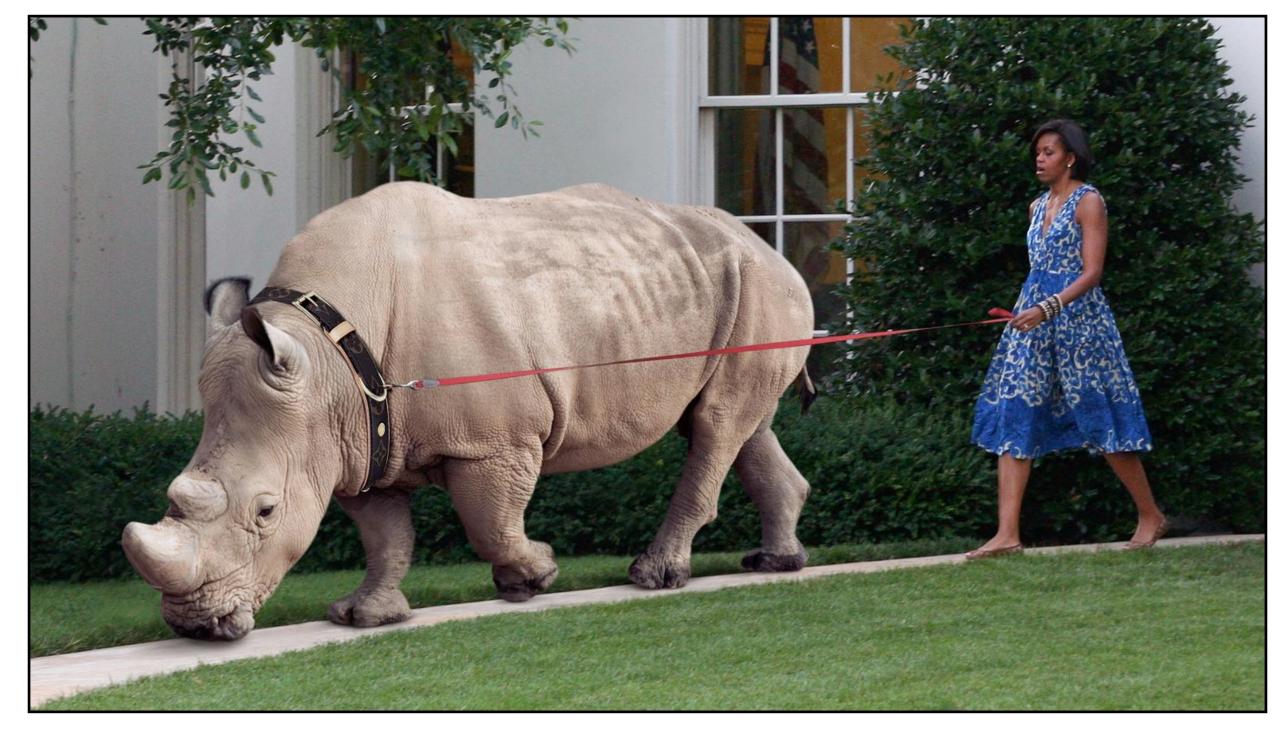




Prediction

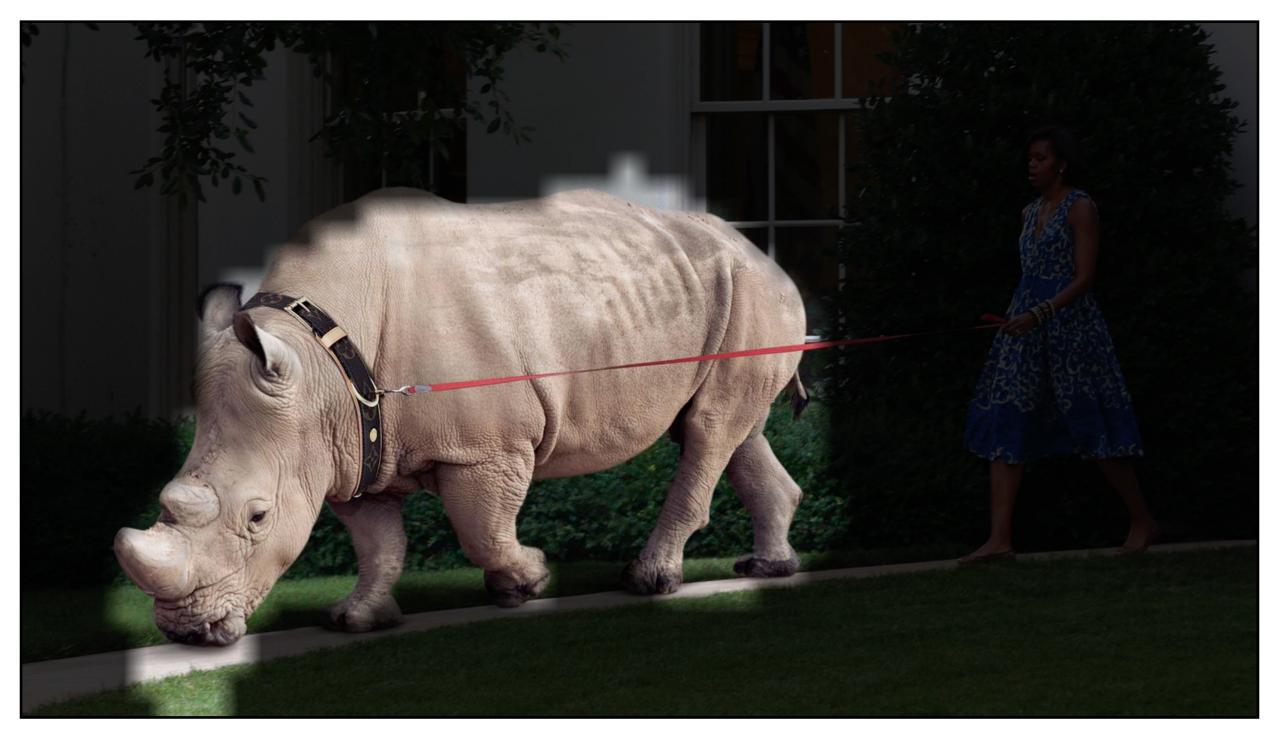
Ground truth

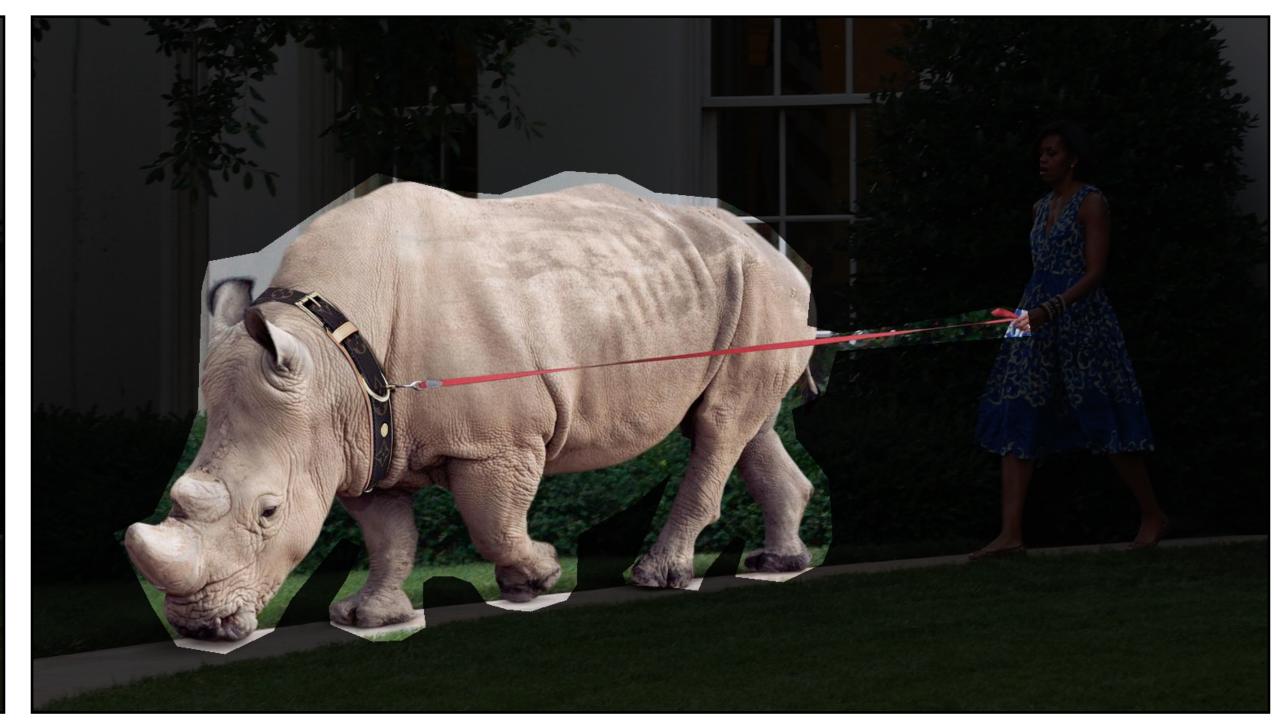
Photo source: <u>TheOnion.com</u>



Input

Photo source: <u>TheOnion.com</u>





Prediction

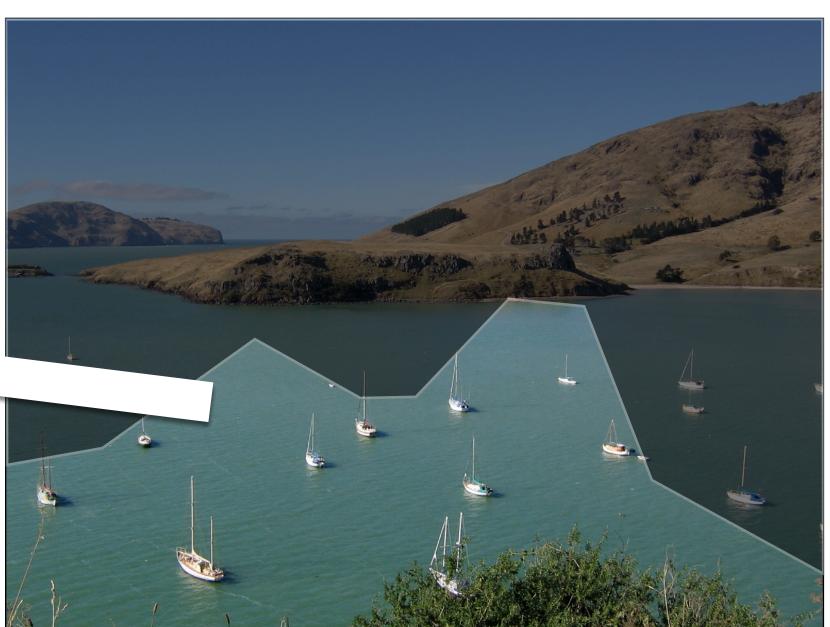
Ground truth

Photo source: <u>TheOnion.com</u>

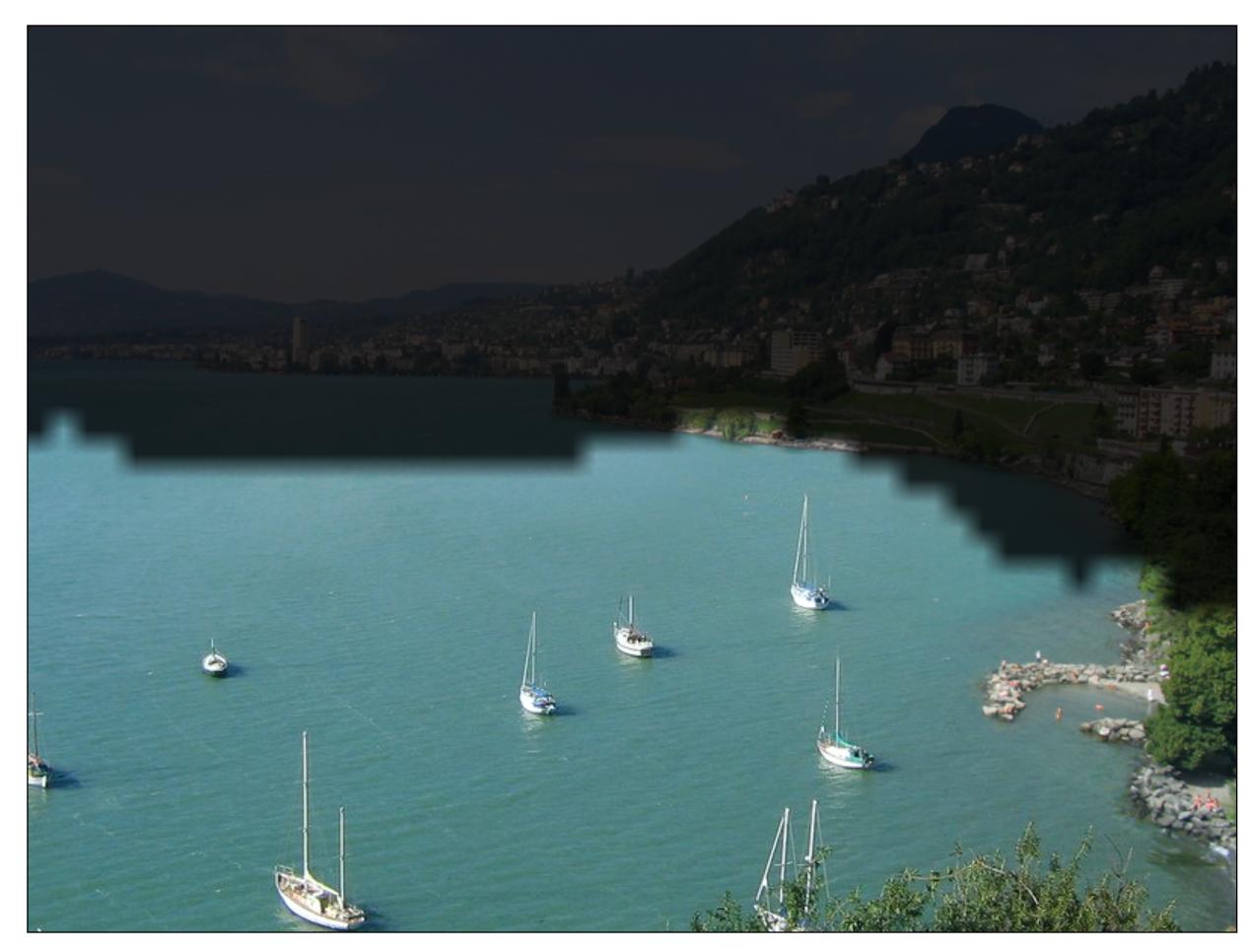


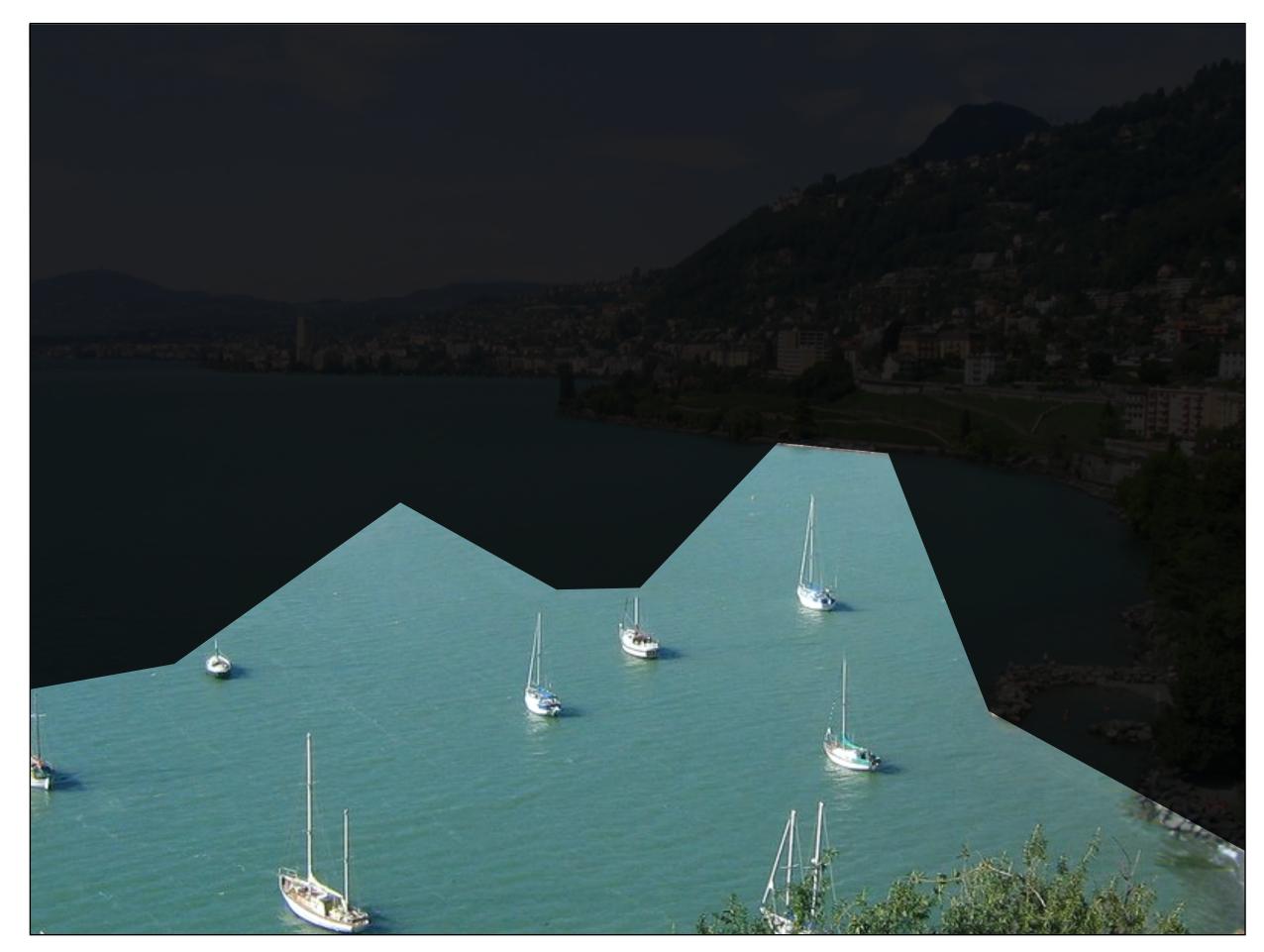
Input





(Hays & Efros 2009)





Prediction

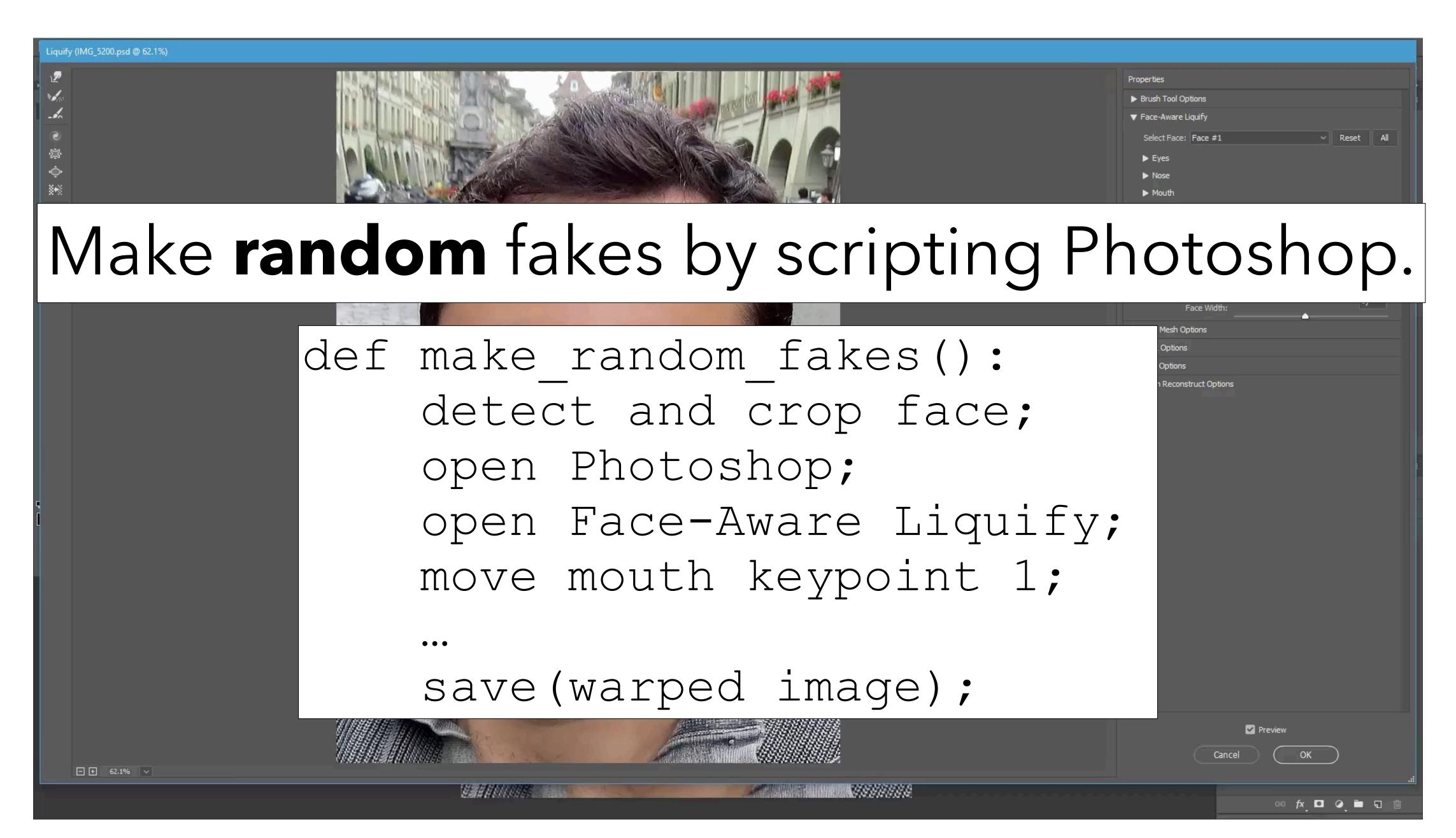
Ground truth

Strategy #4: supervised learning

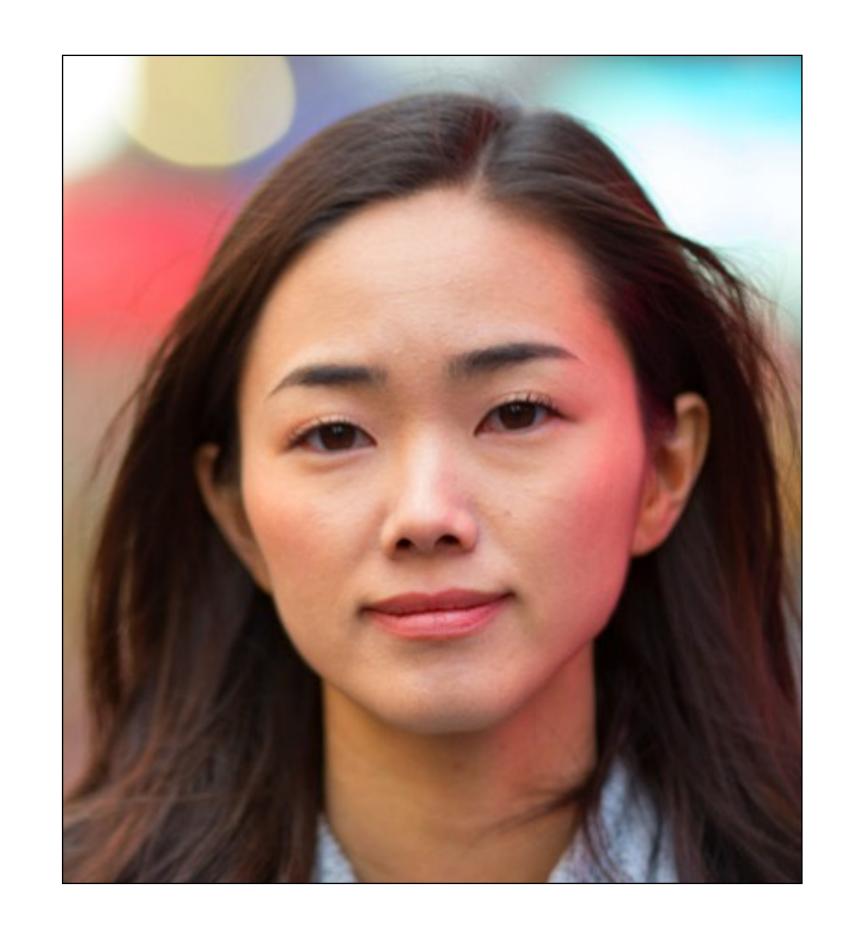


[Wang et al., "Image Splice Detection via Learned Self-Consistency", 2018]

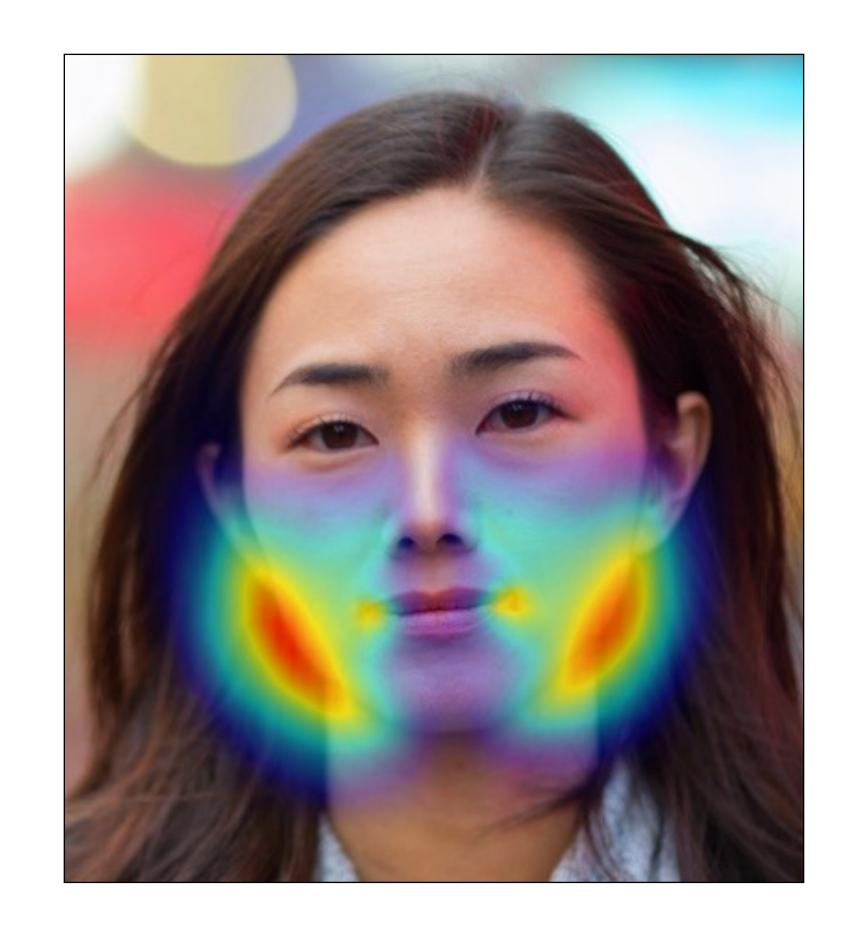




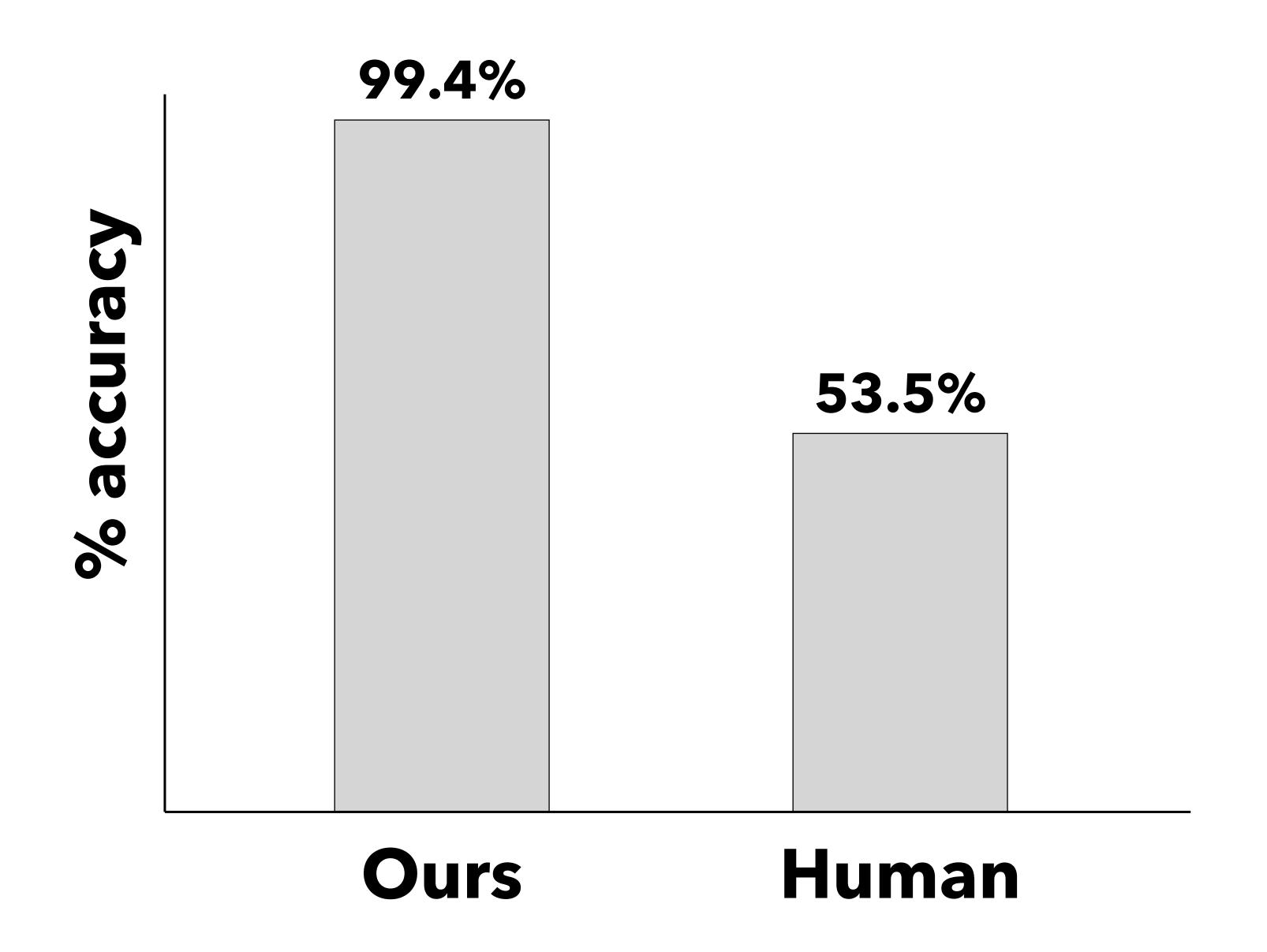
Photoshop Face-Aware Liquify tutorial. Source: https://youtu.be/5Qqv_C6iVvQ?t=86



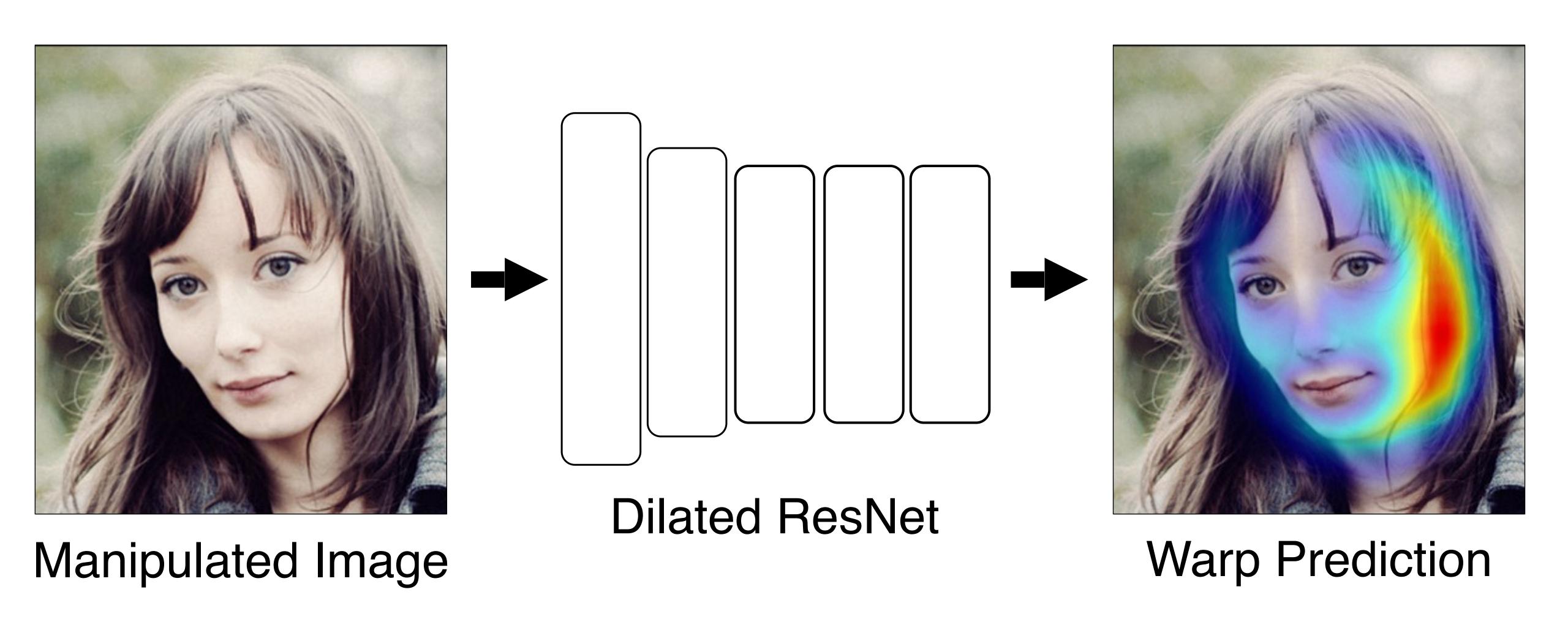
Warp detector



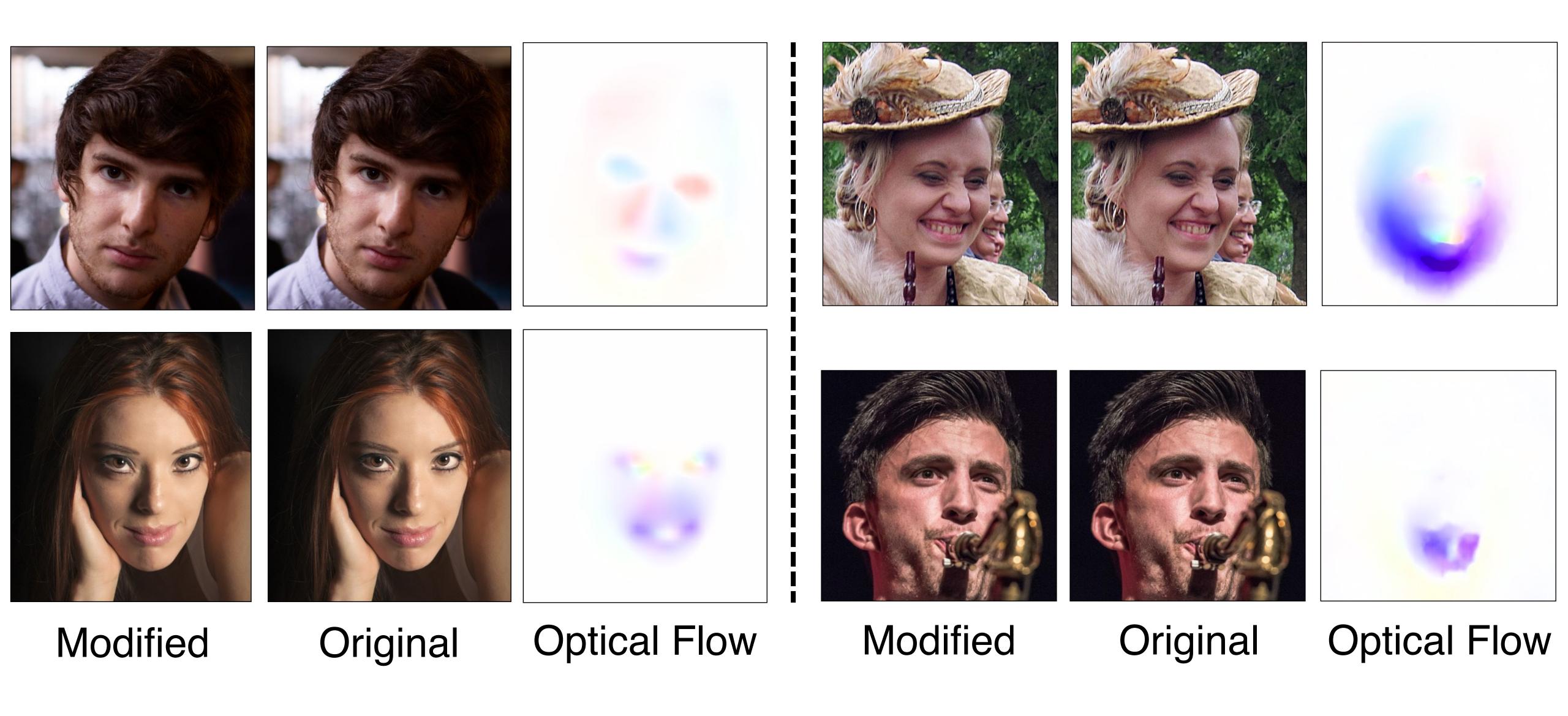
Real-or-fake classification

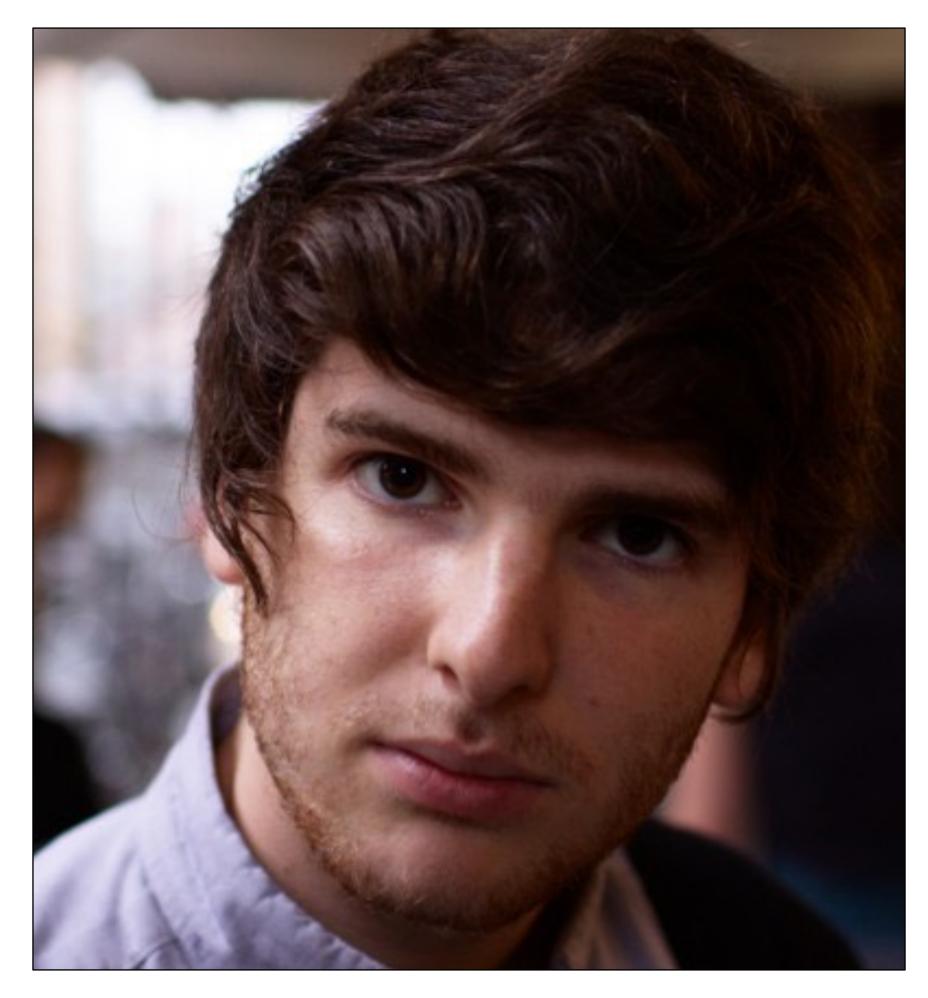


What moved where?



What moved where?

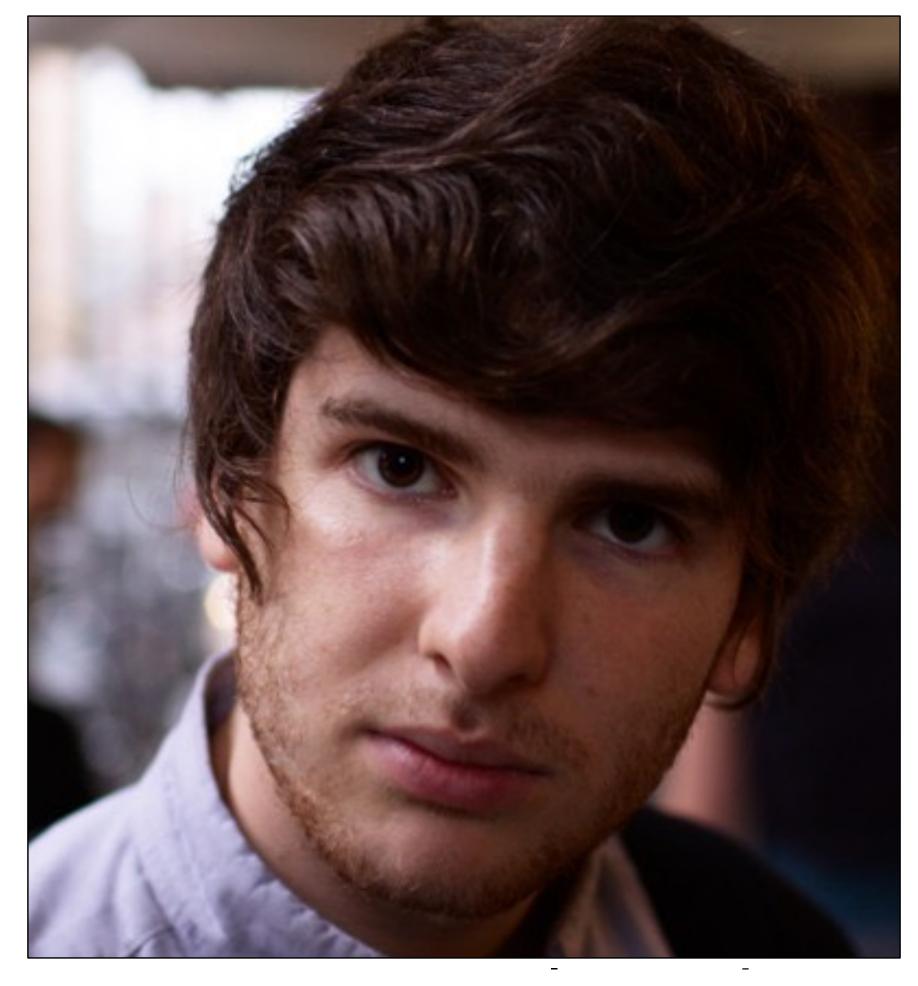




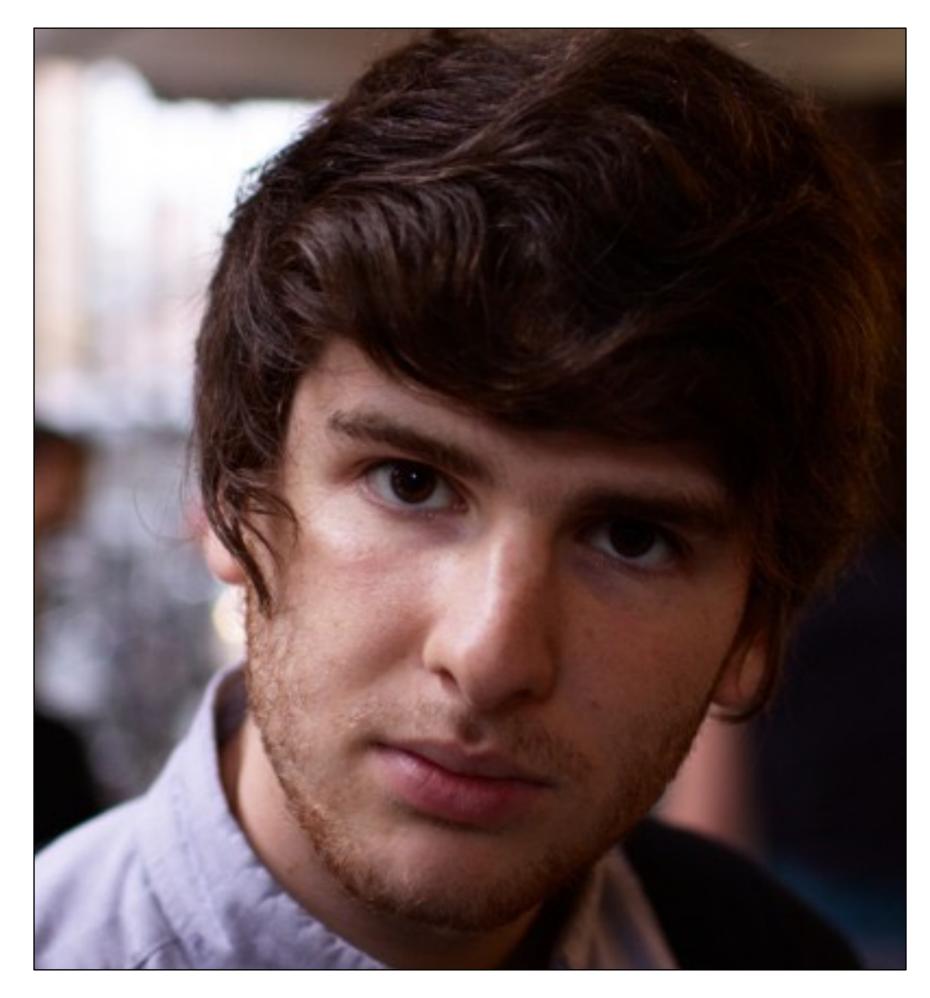
Manipulated Photo



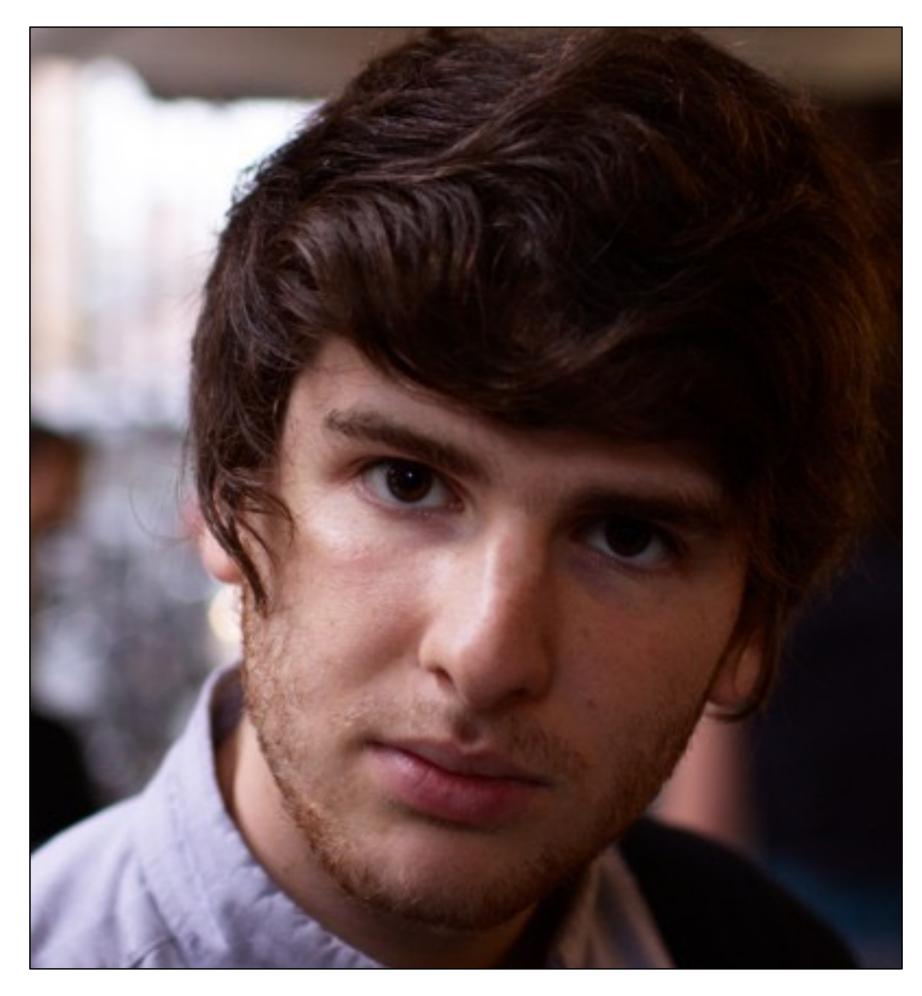
Flow Prediction



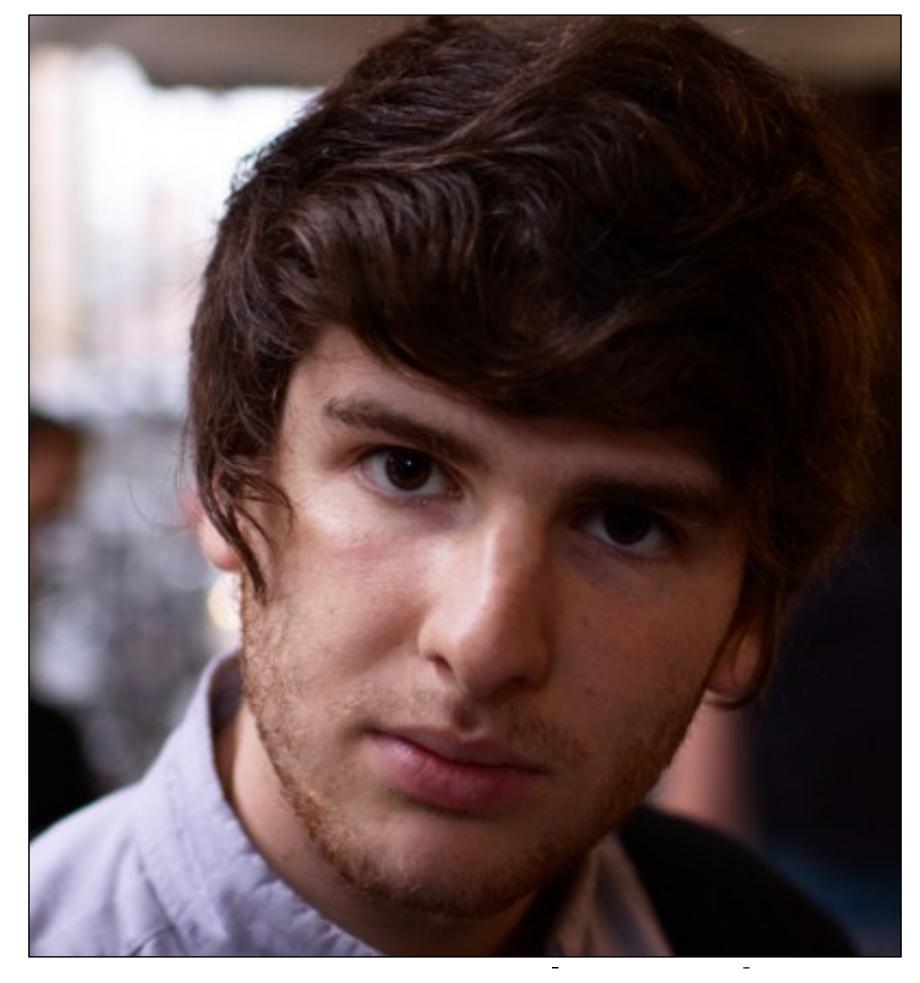
Suggested "Undo"



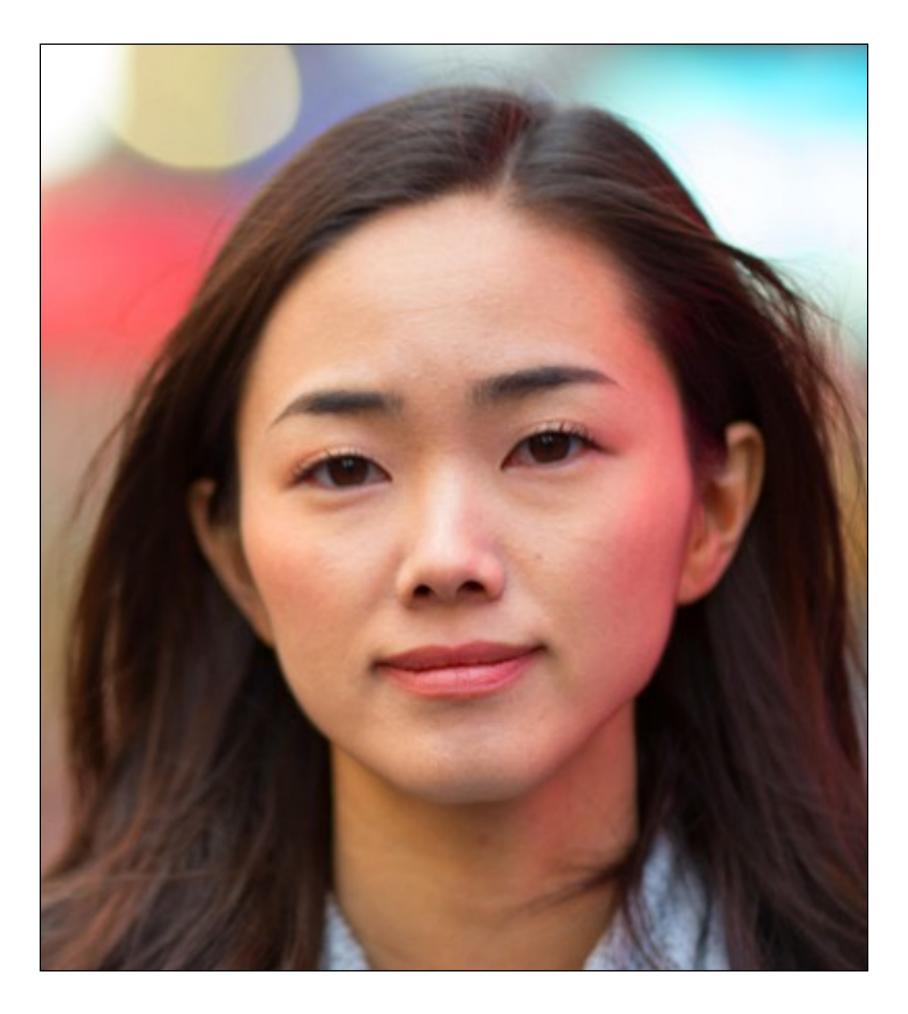
Original Photo



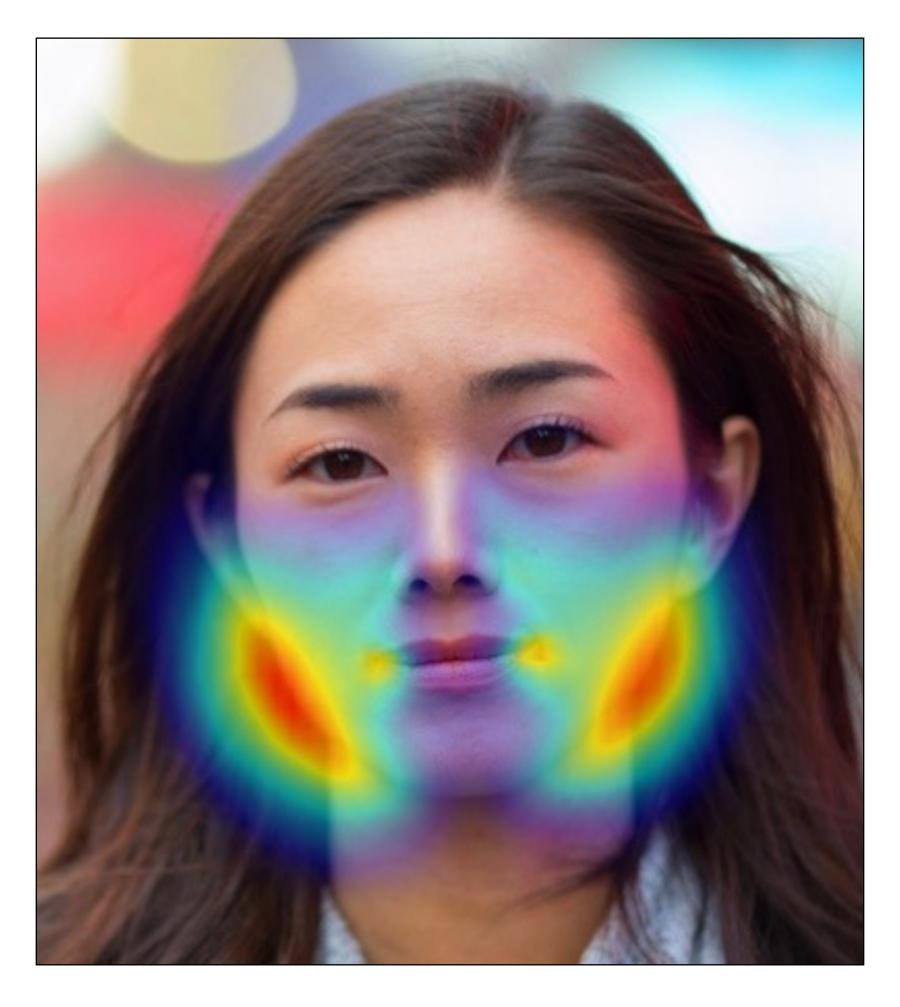
Manipulated vs. Original



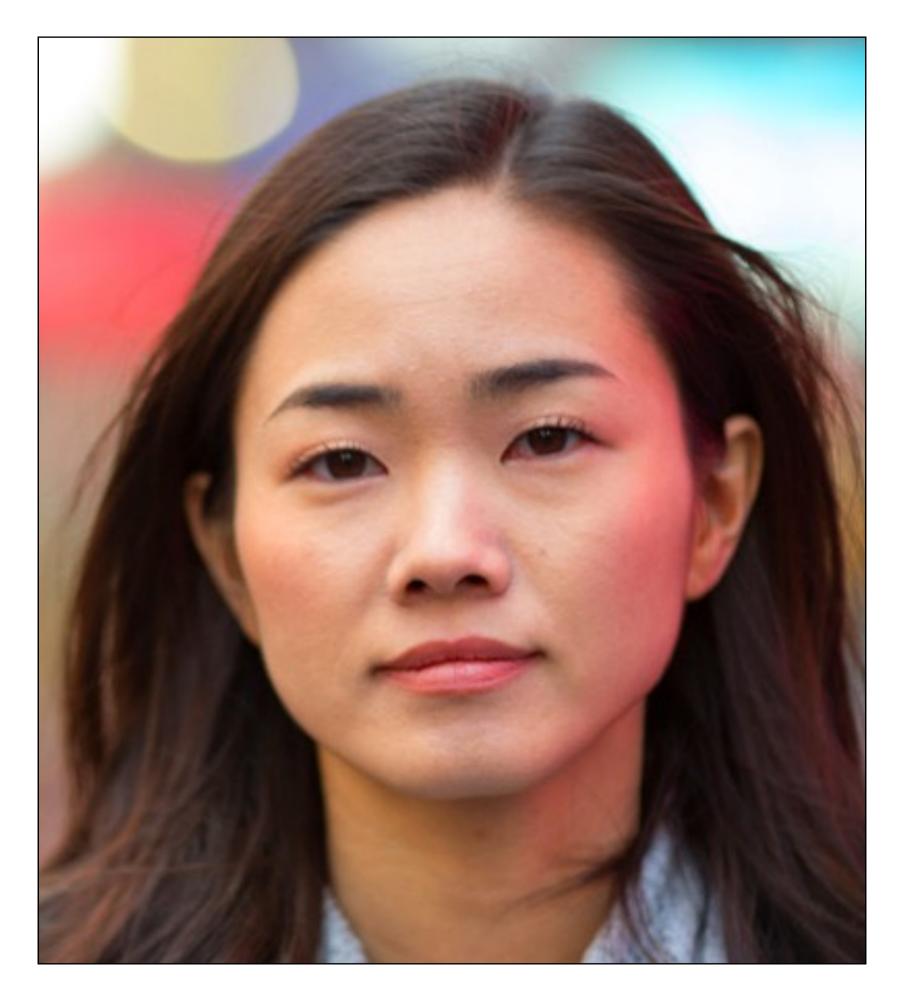
Undo vs. Original



Manipulated Photo



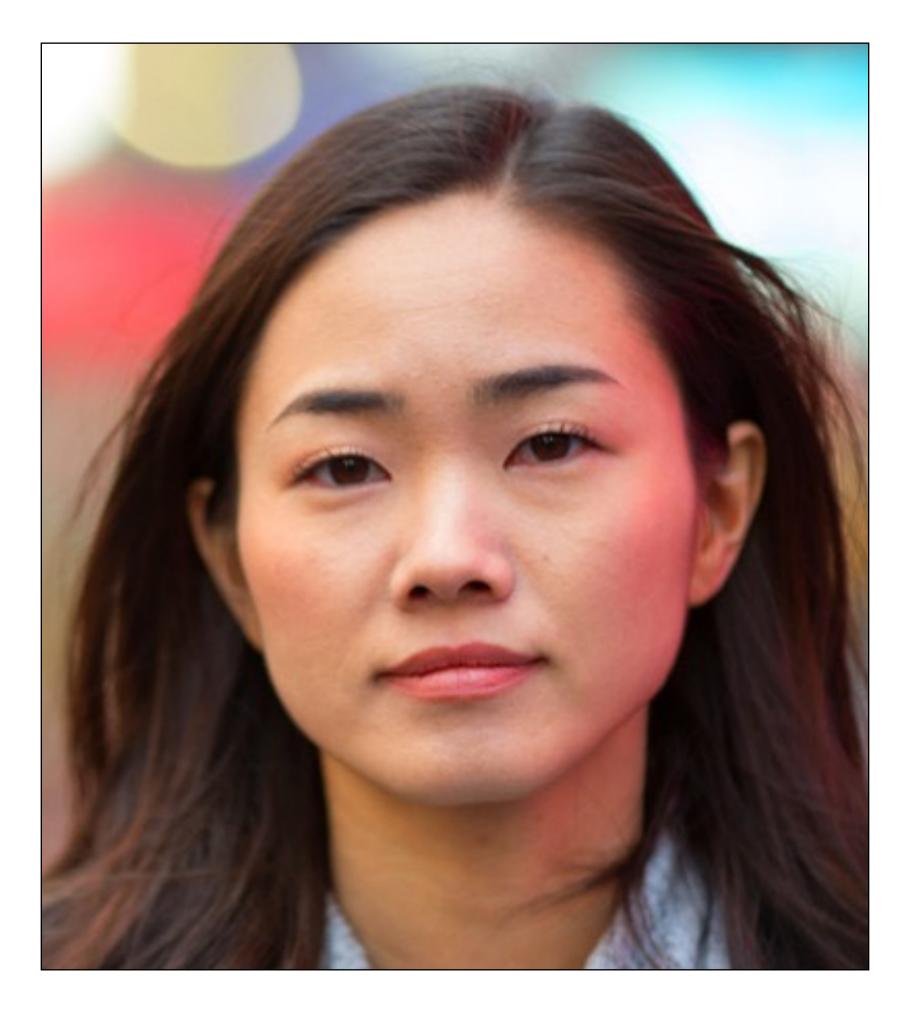
Warp Prediction



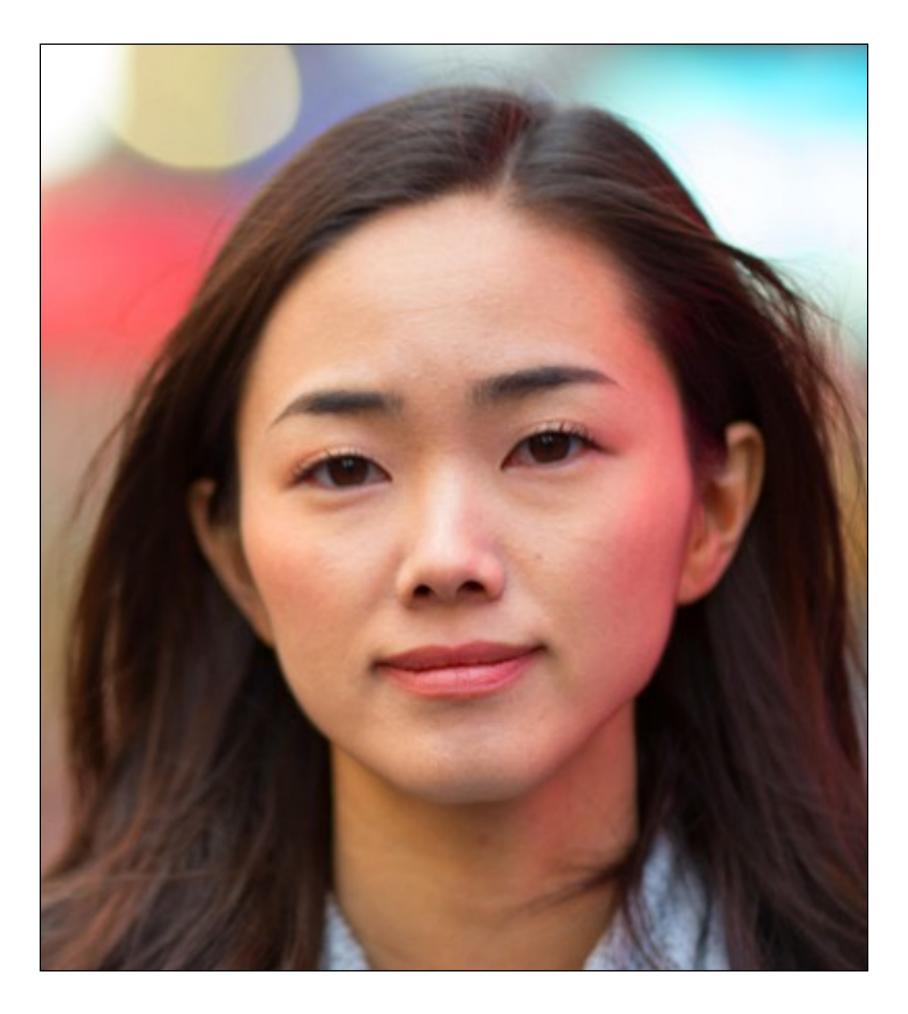
Suggested "Undo"



Original Photo



Suggested "Undo"



Manipulated Photo

Al-generated images

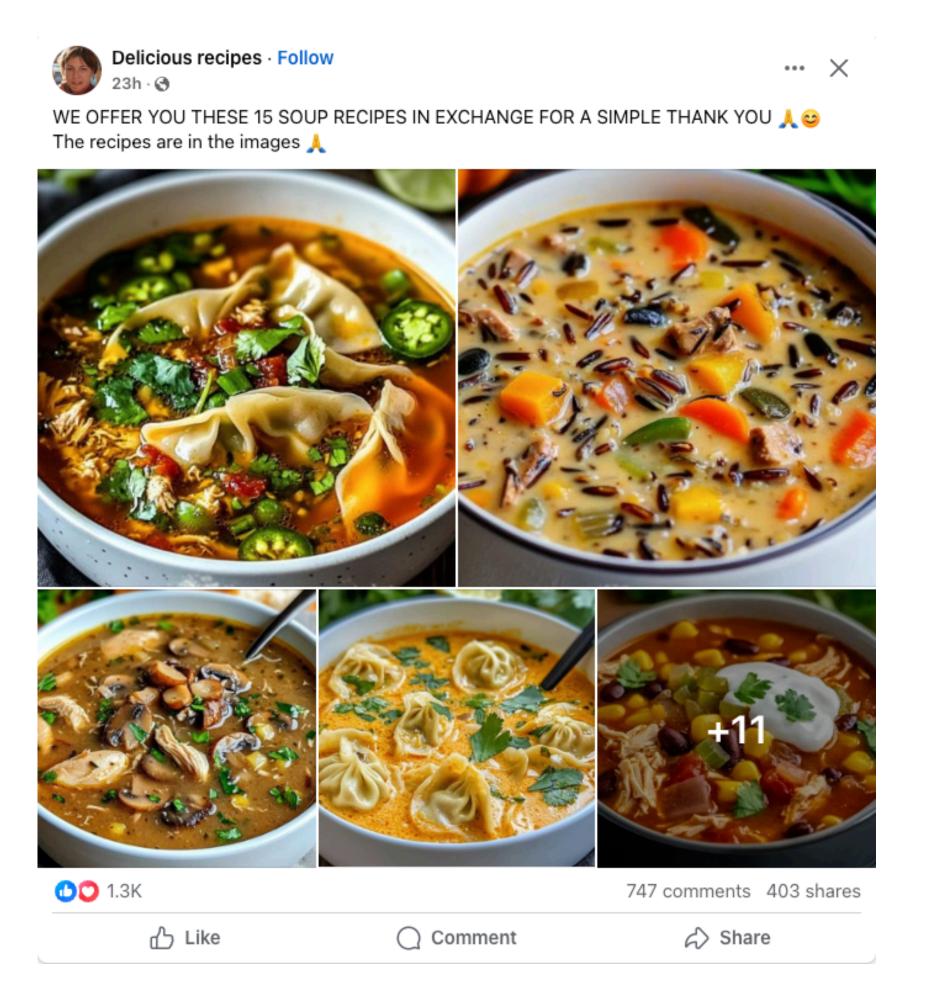


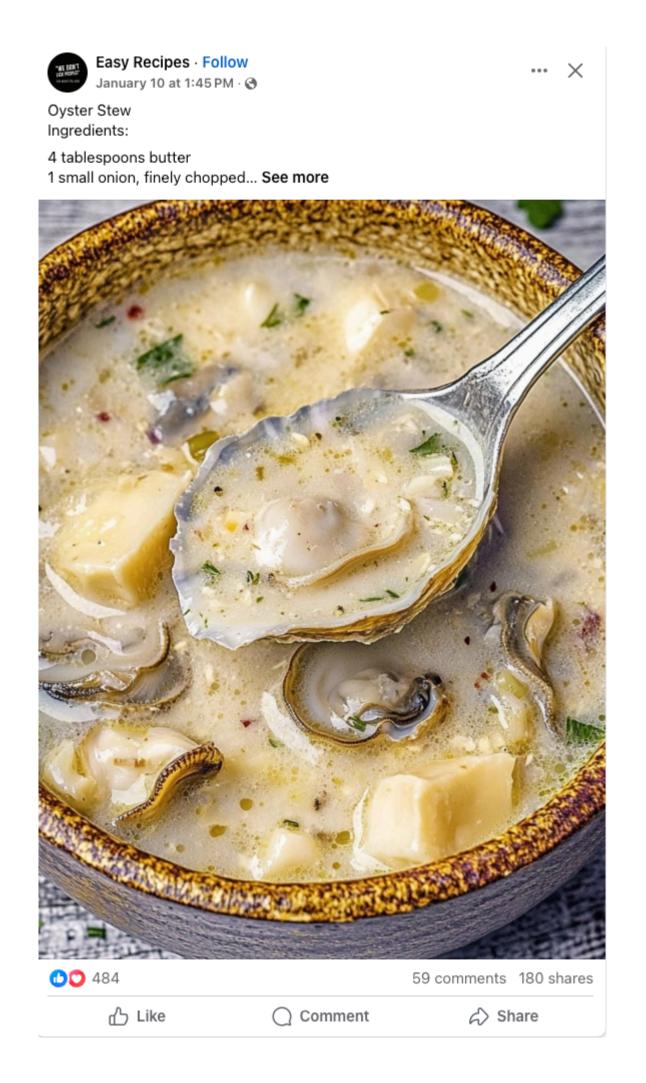




"Catholic Pope Francis wearing Balenciaga puffy jacket in drill rap music video, throwing up gang signs with hands, taken using a Canon EOS R camera with a 50mm f/1.8 lens, f/2.2 aperture, shutter speed 1/200s, ISO 100 and natural light, Full Body, Hyper Realistic Photography, Cinematic, Cinema, Hyperdetail, UHD, Color Correction, hdr, color grading, hyper realistic CG animation --ar 4:5 --upbeta --q 2 --v 5."

Al-generated spam







22 comments 24 shares

Share

Like

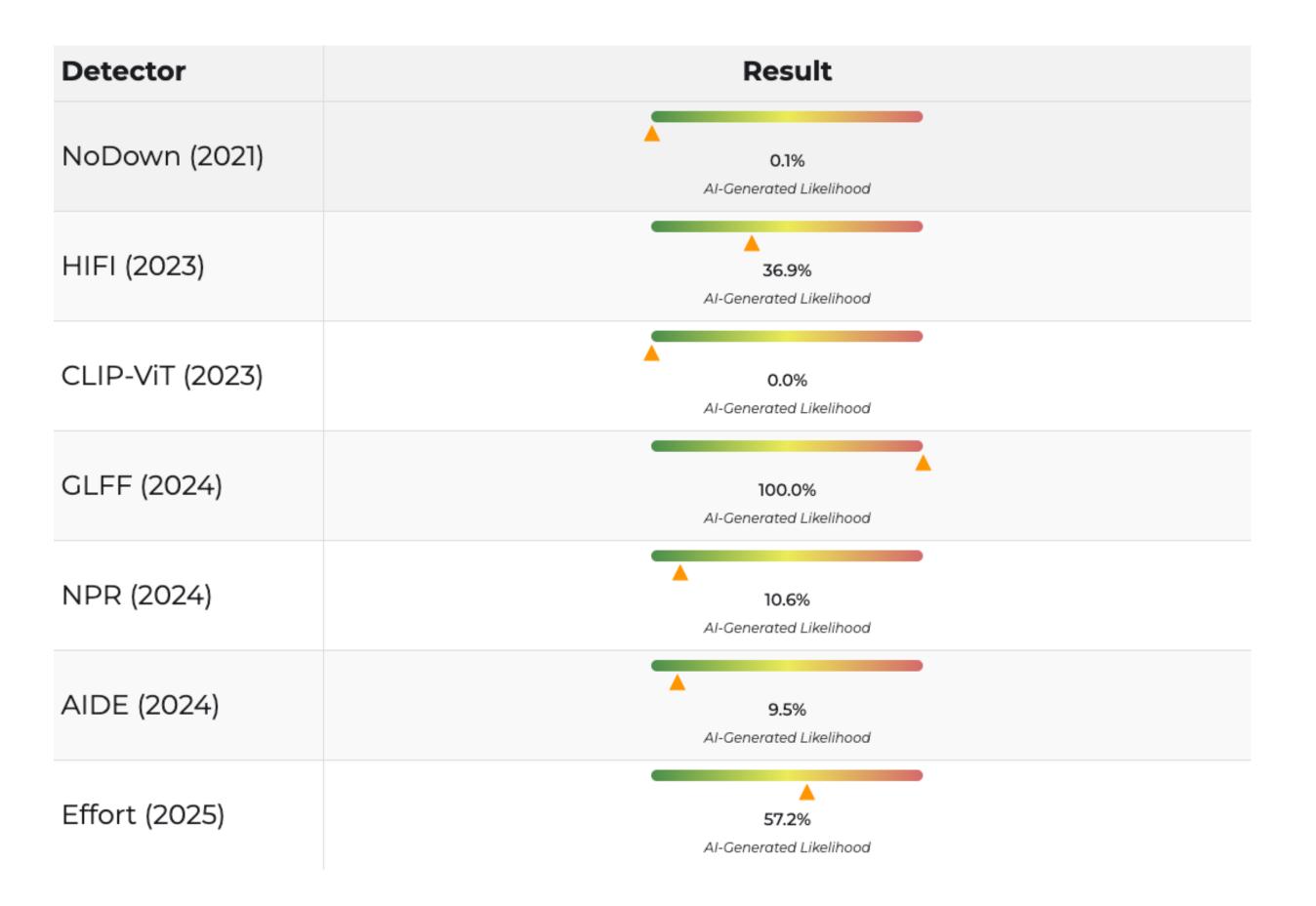
Comment



Examples from Jeongsoo Park's Facebook feed

Real or fake?





[Li, Zhang, Sun, Qi, Lyu, "DeepFake-o-meter", 2024]

Downsides of easy image generation

Misinformation and abuse



Teen Girls Confront an Epidemic of Deepfake Nudes in Schools

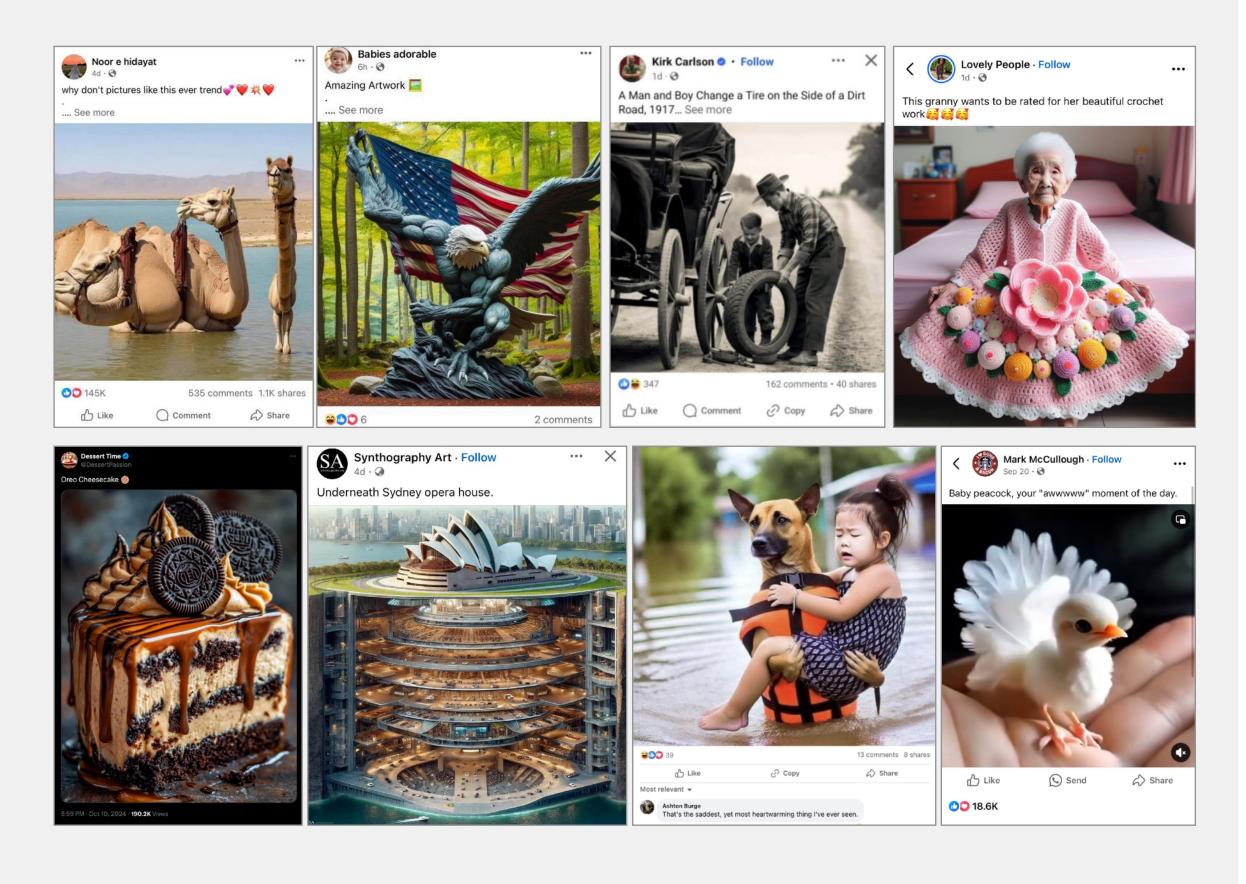
Using artificial intelligence, middle and high school students have fabricated explicit images of female classmates and shared the doctored pictures.



Raising doubt in real images



Al spam



Source: "Insane Facebook Al slop" @FacebookAlslop

Generalization in Al-generated image detection

Training



ProGAN



DALL·E 2



VQ-GAN



Midjourney v5

Generalization in Al-generated image detection

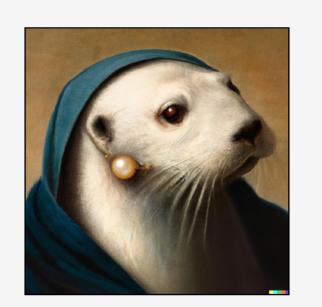
Training



ProGAN



VQ-GAN



DALL·E 2



Midjourney v5

Test on the same generators?

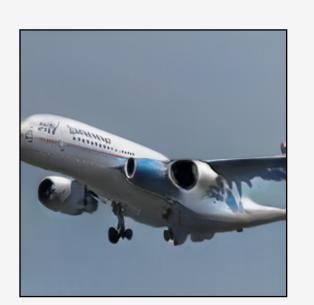
→ Detectors generalize pretty well.

Test on out-of-distribution generators?

→ Not always. But intriguingly there is some generalization!

Generalization in Al-generated image detection

Training



ProGAN



VQ-GAN

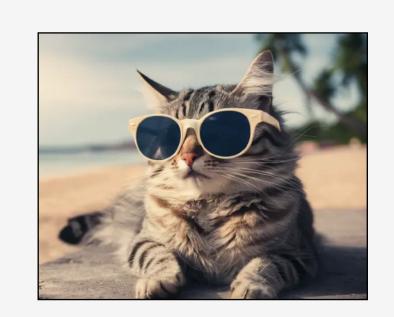


DALL·E 2



Midjourney v5

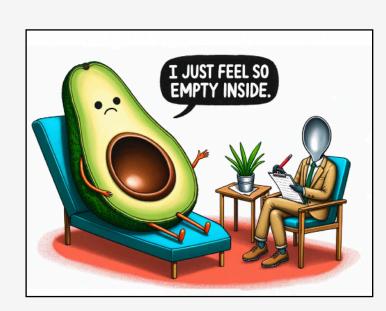
Testing



FLUX.1



Imagen v3



DALL·E 3



Ideogram



github.com/anon123/ mycoolgenerator



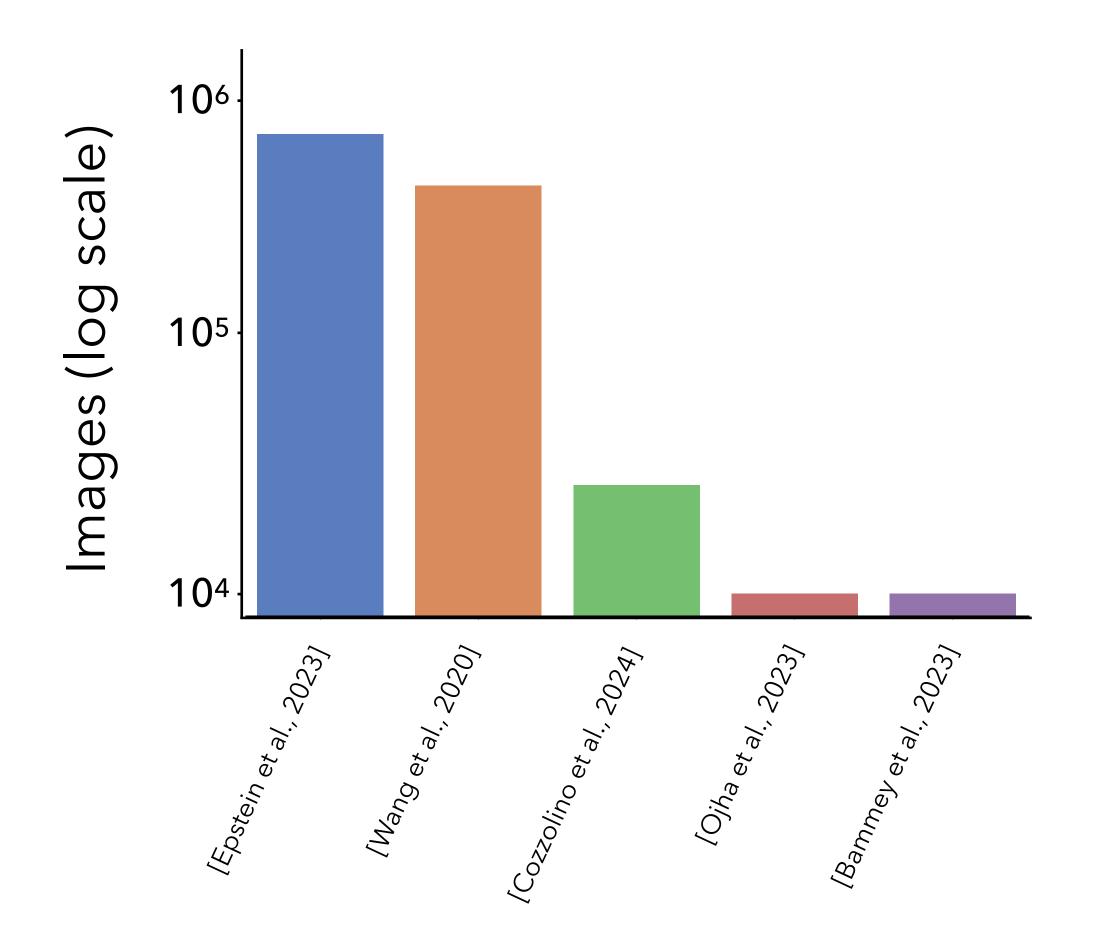


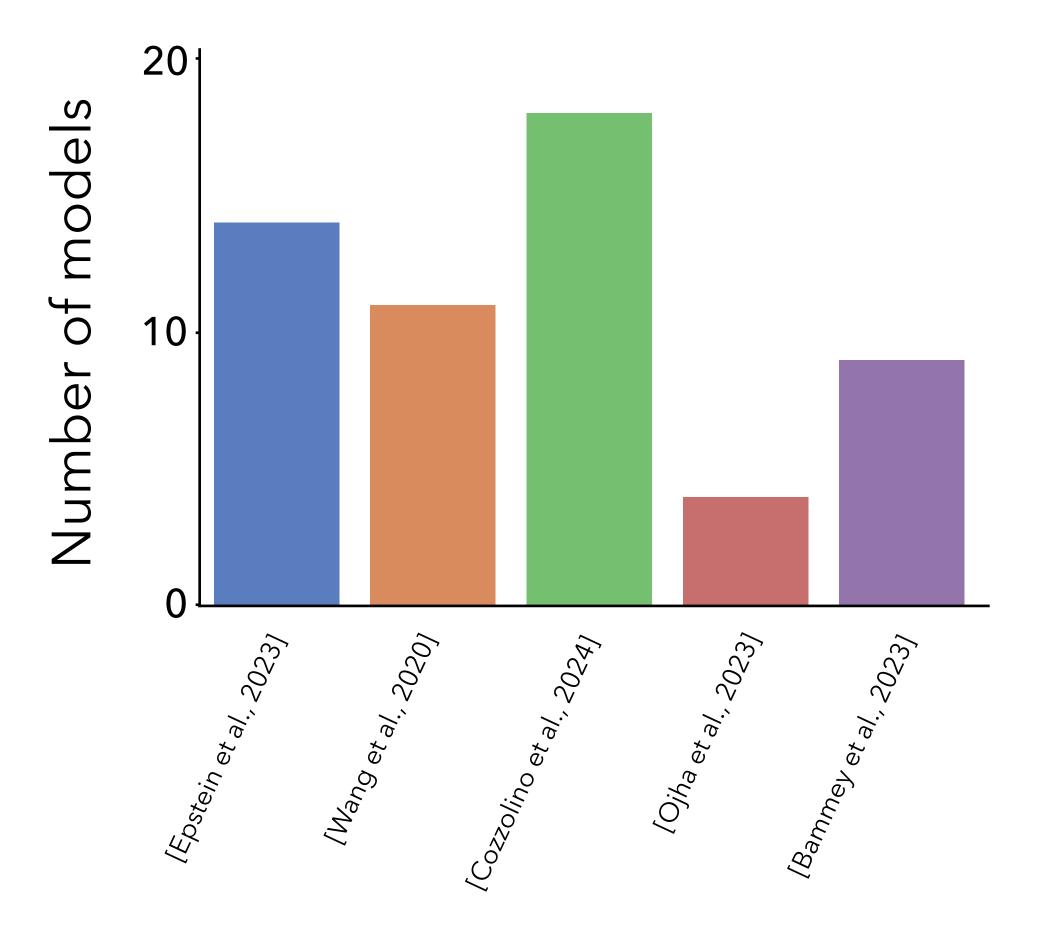
"DIY" models

One problem: the data

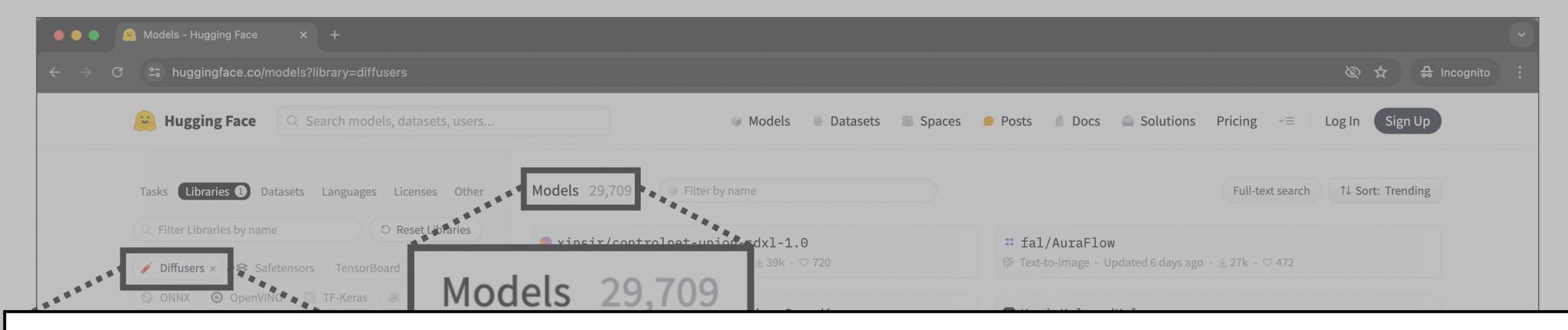
Datasets have lots of images.

But relatively few models.

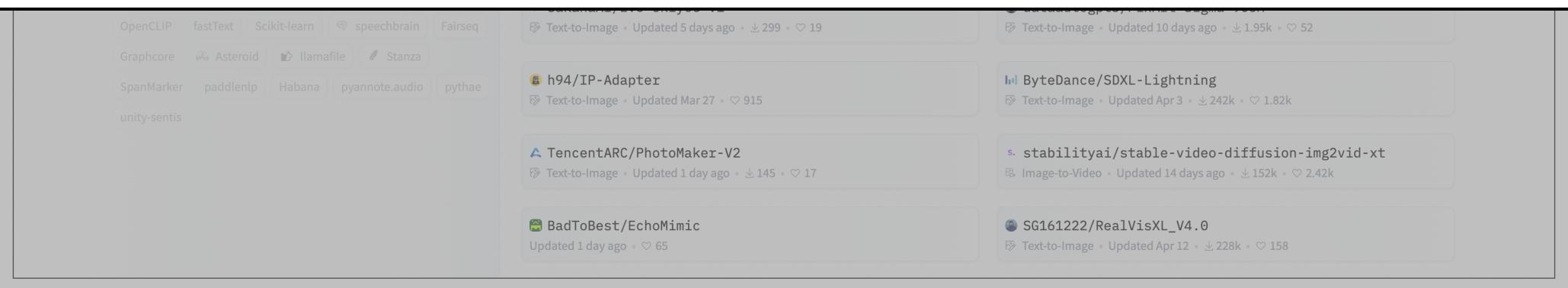


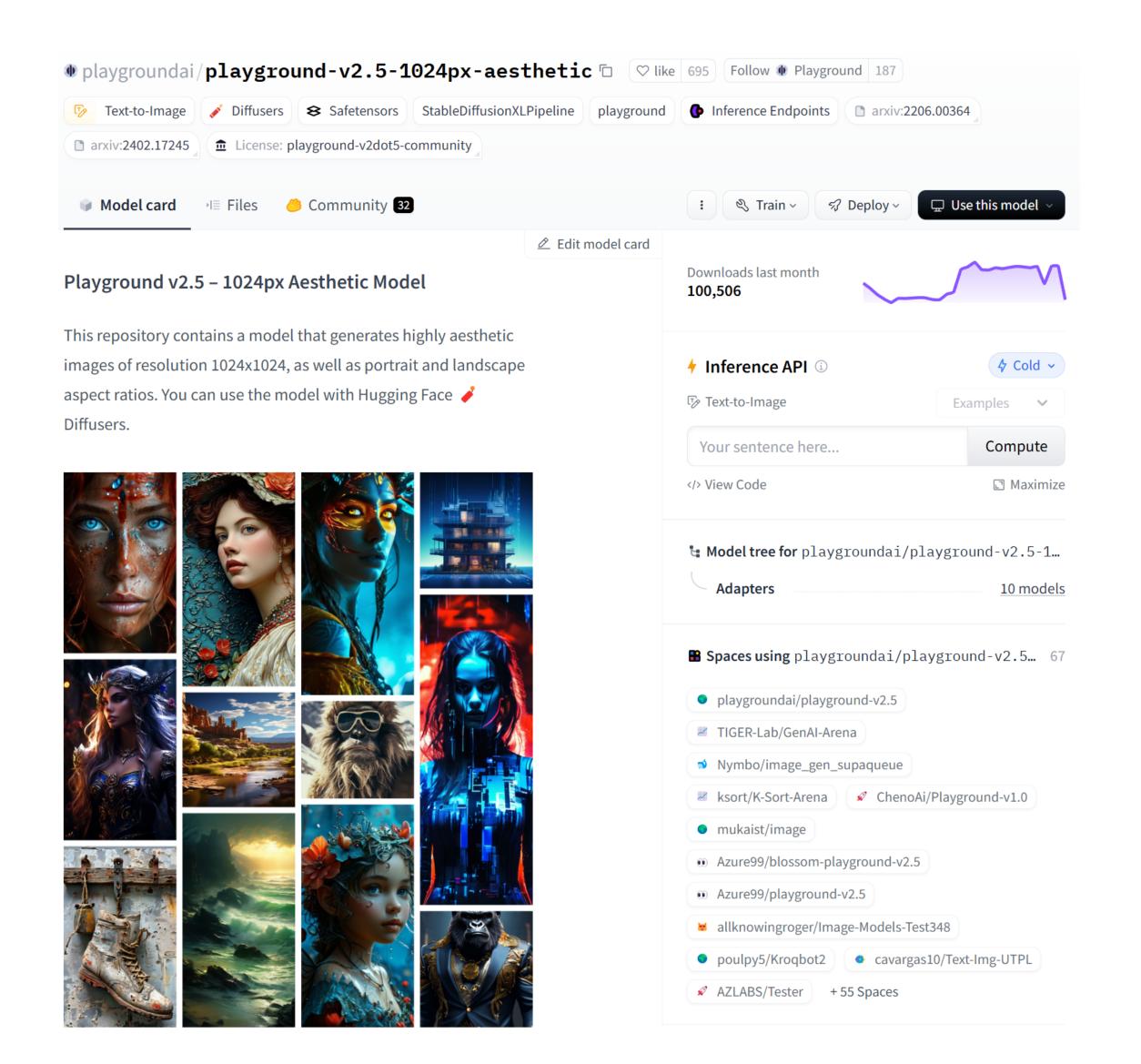


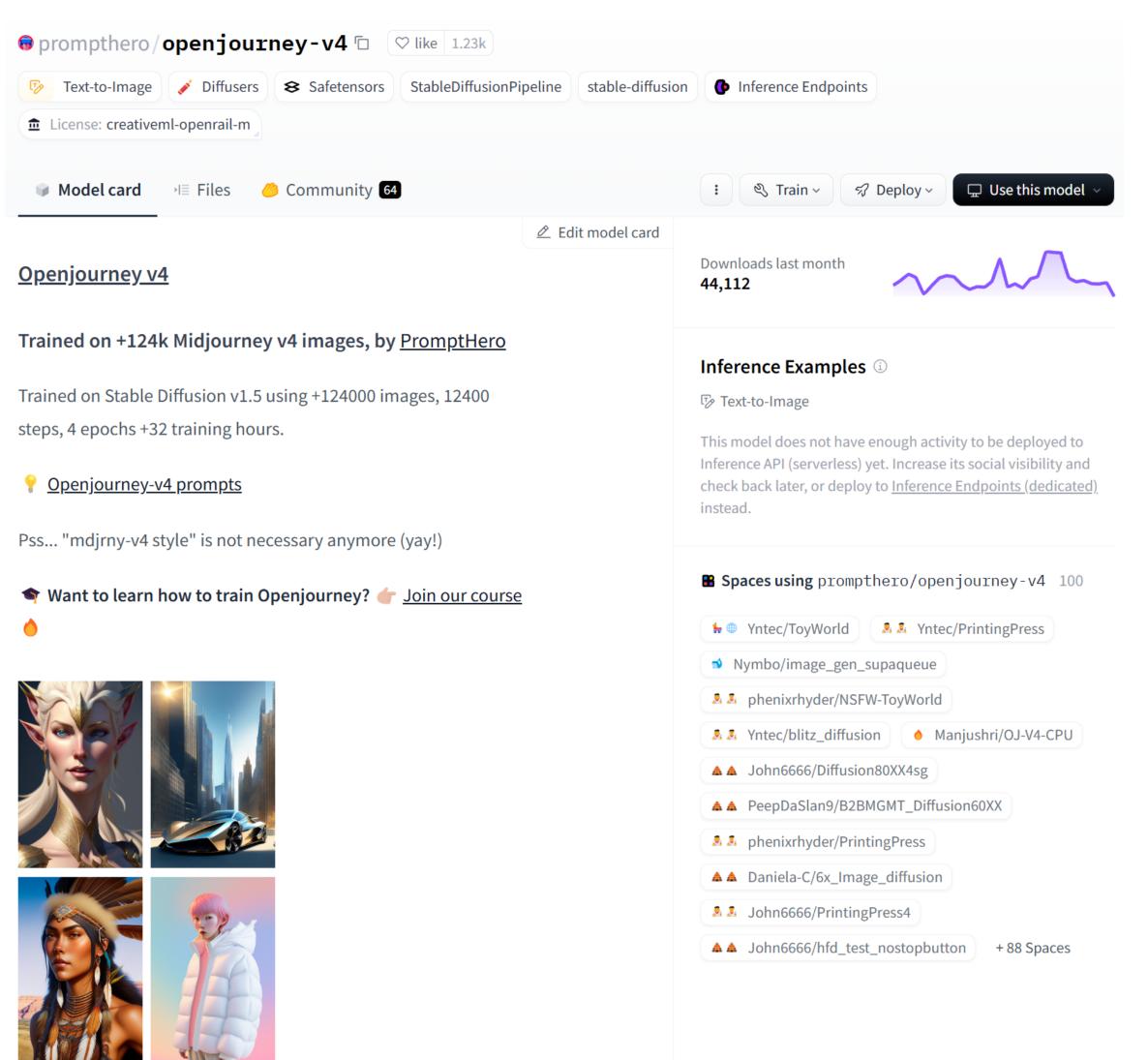
Model sharing communities

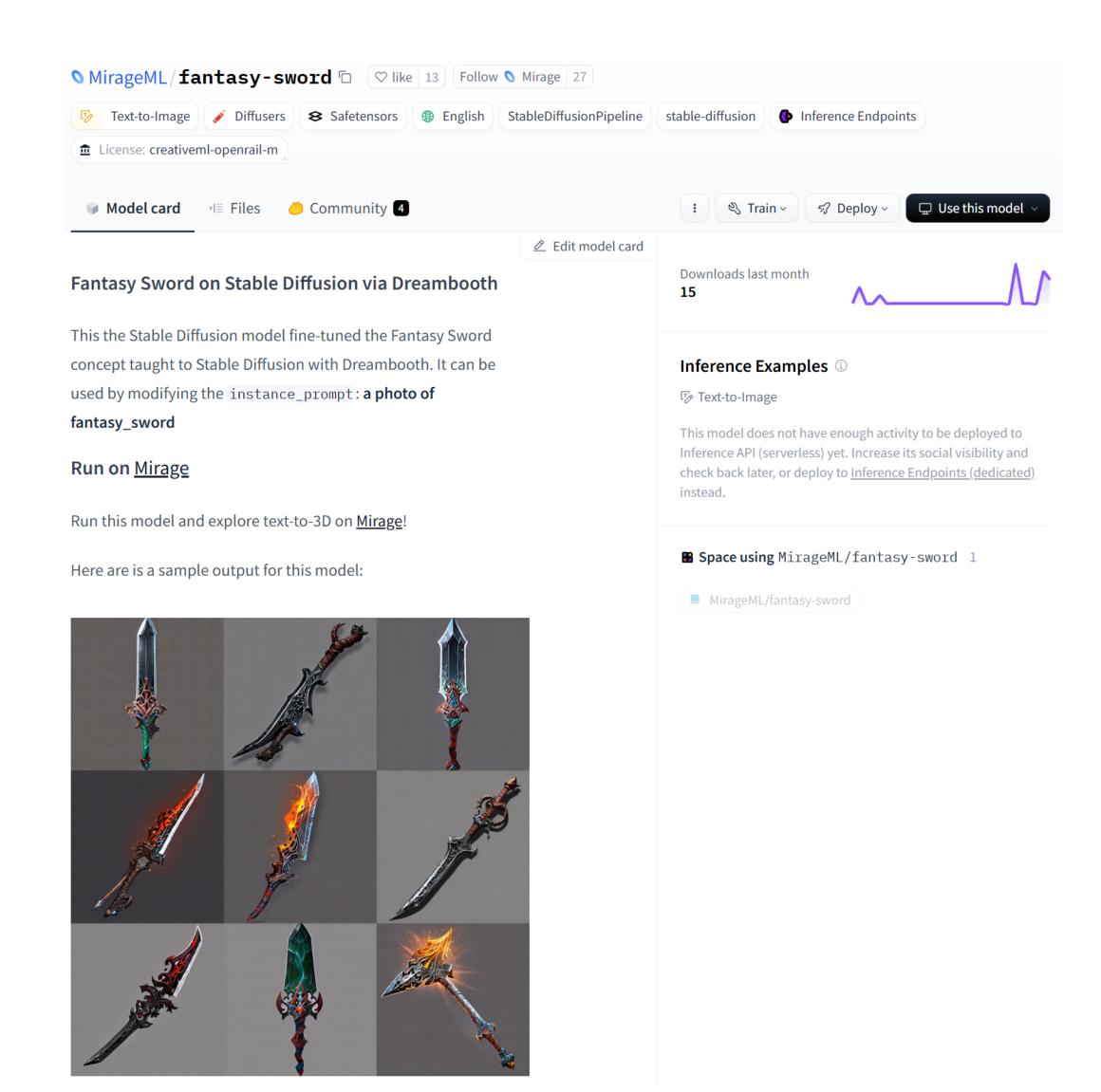


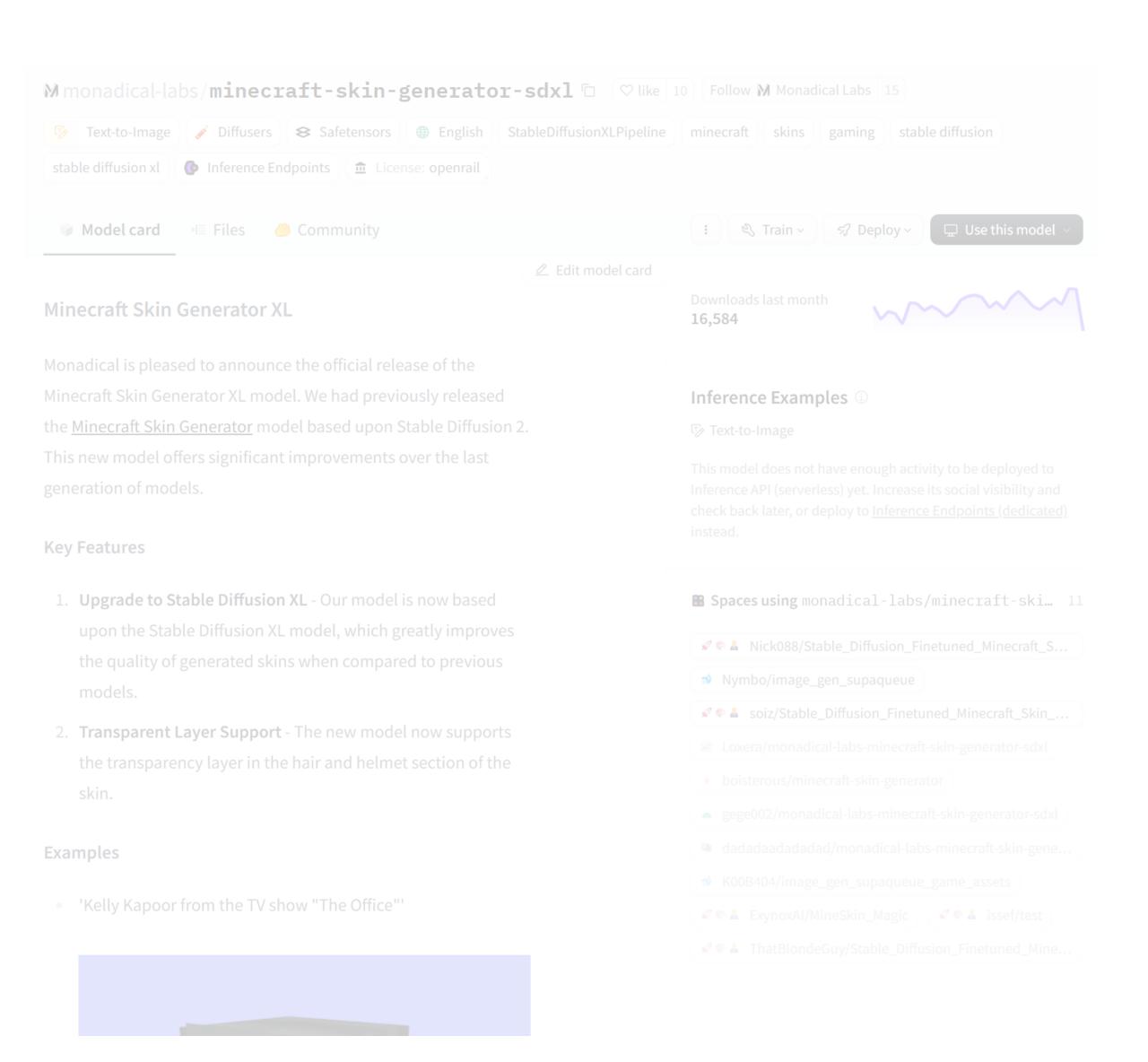
Can we automatically acquire fake images from open source latent diffusion models?

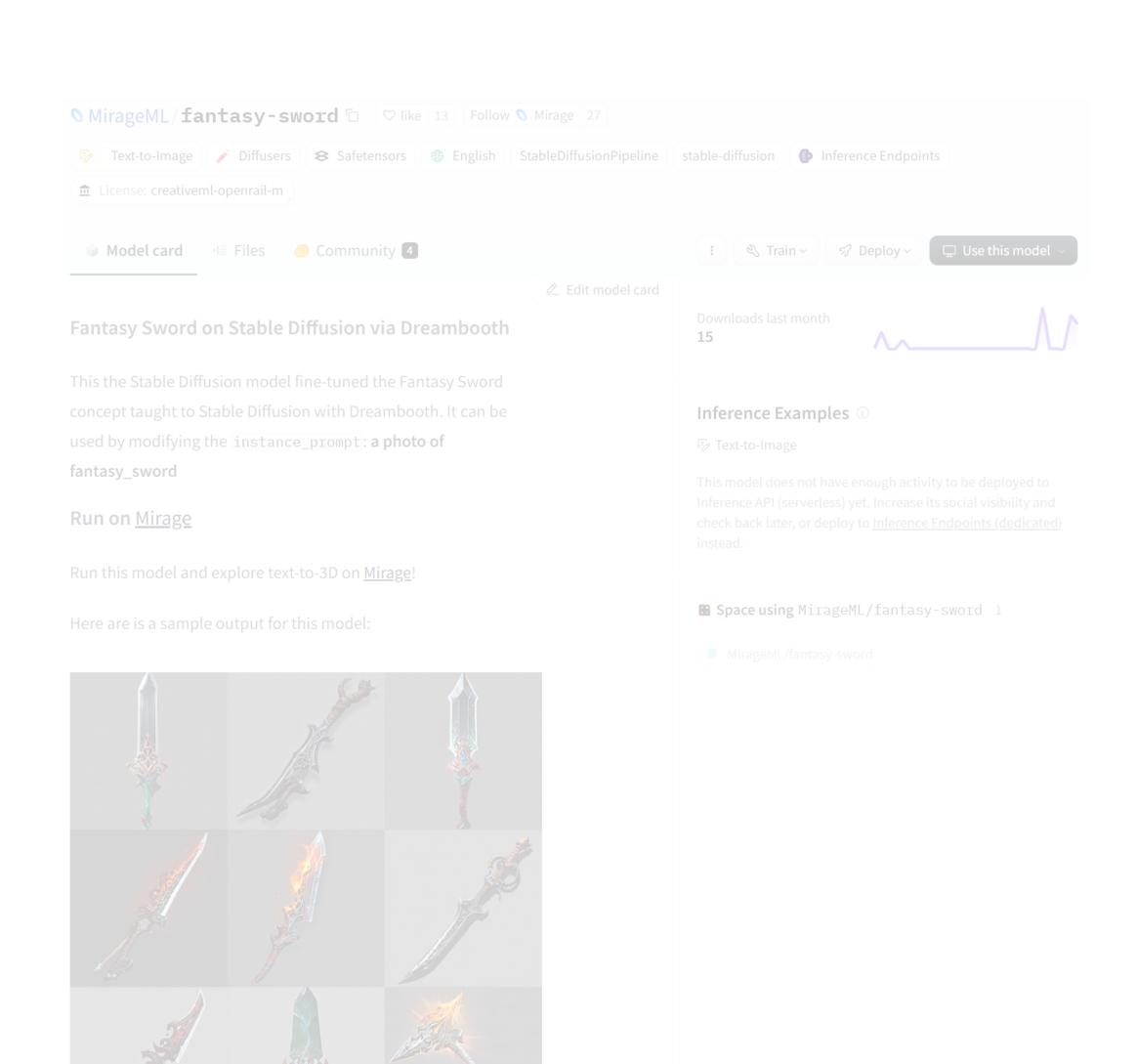


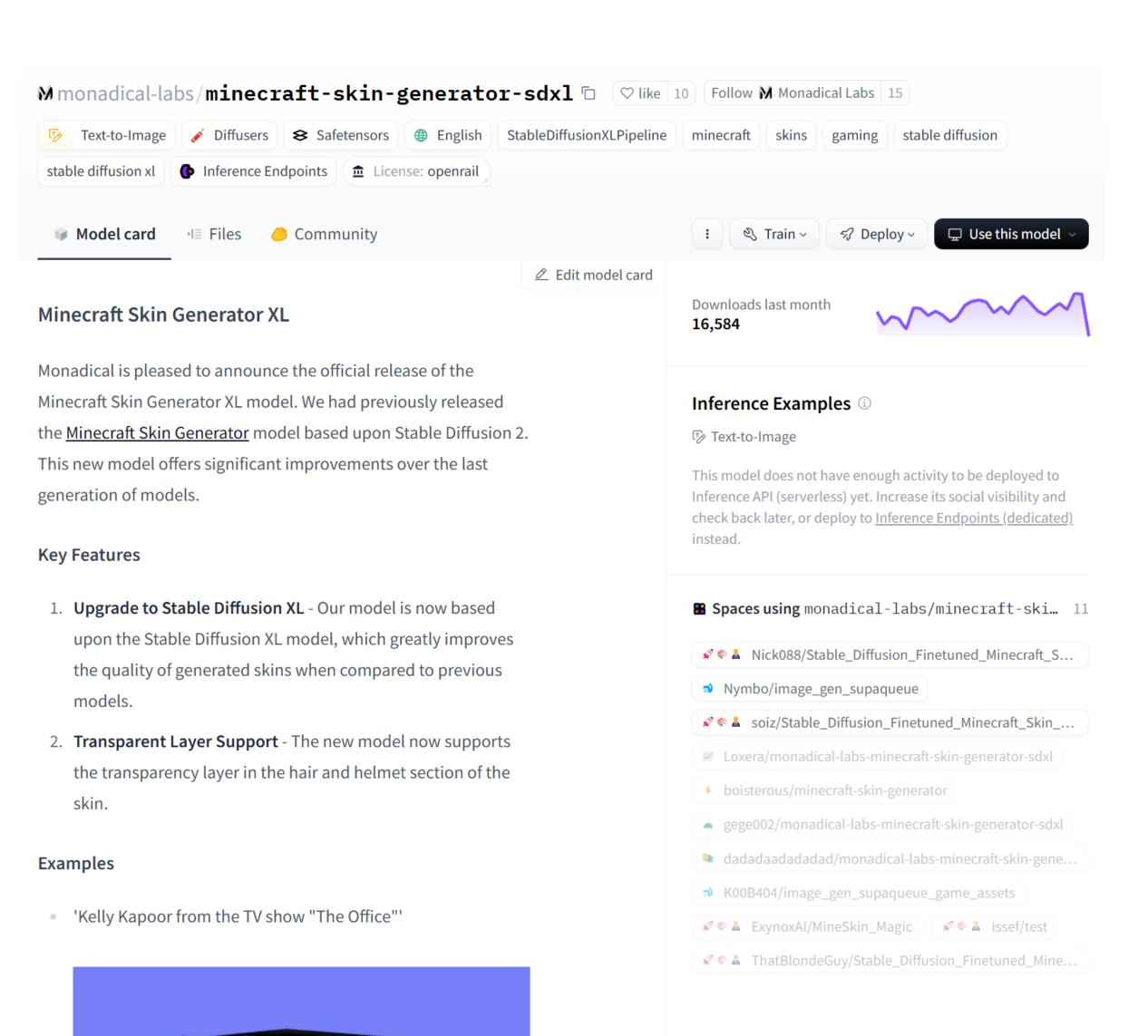


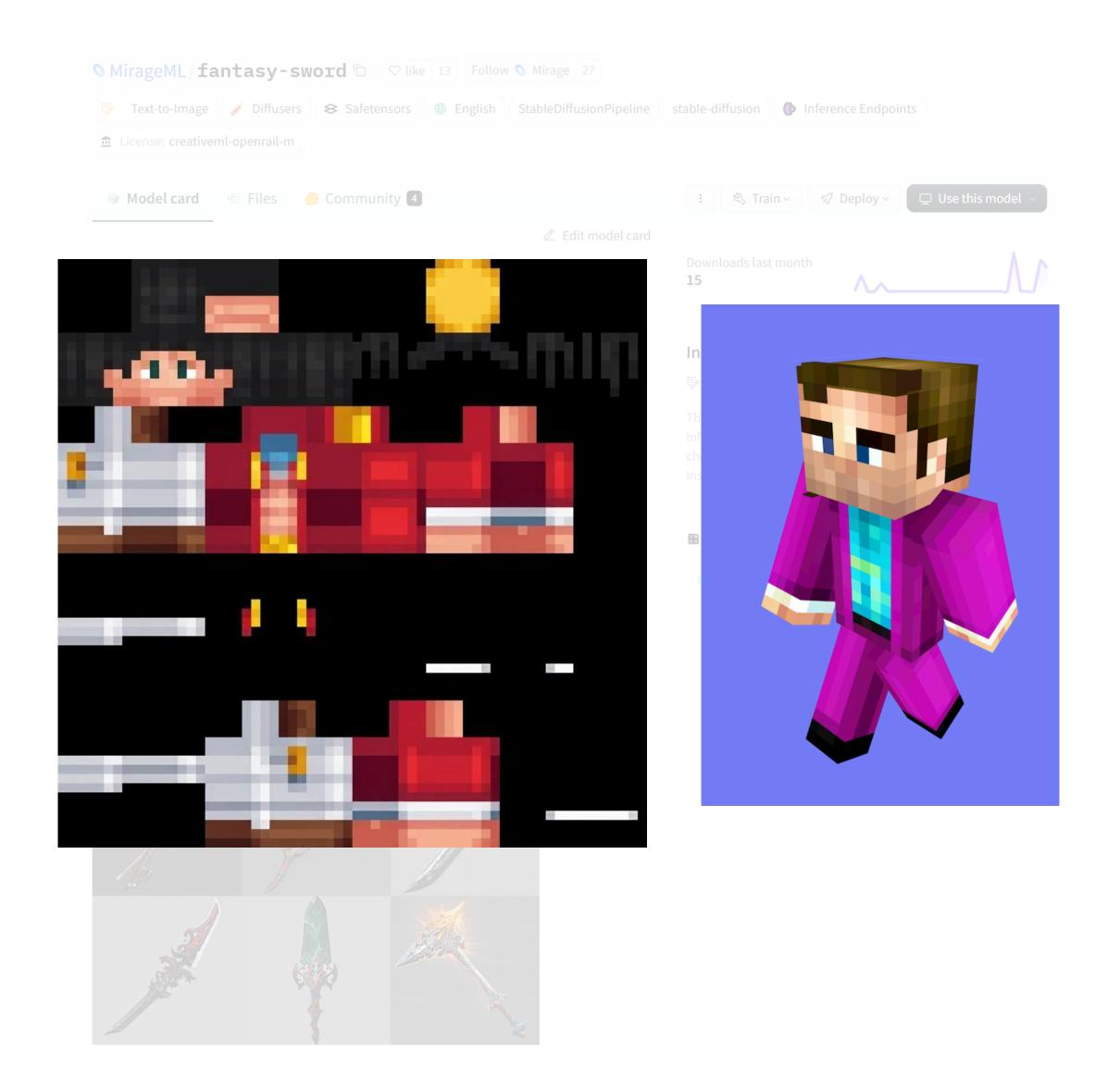


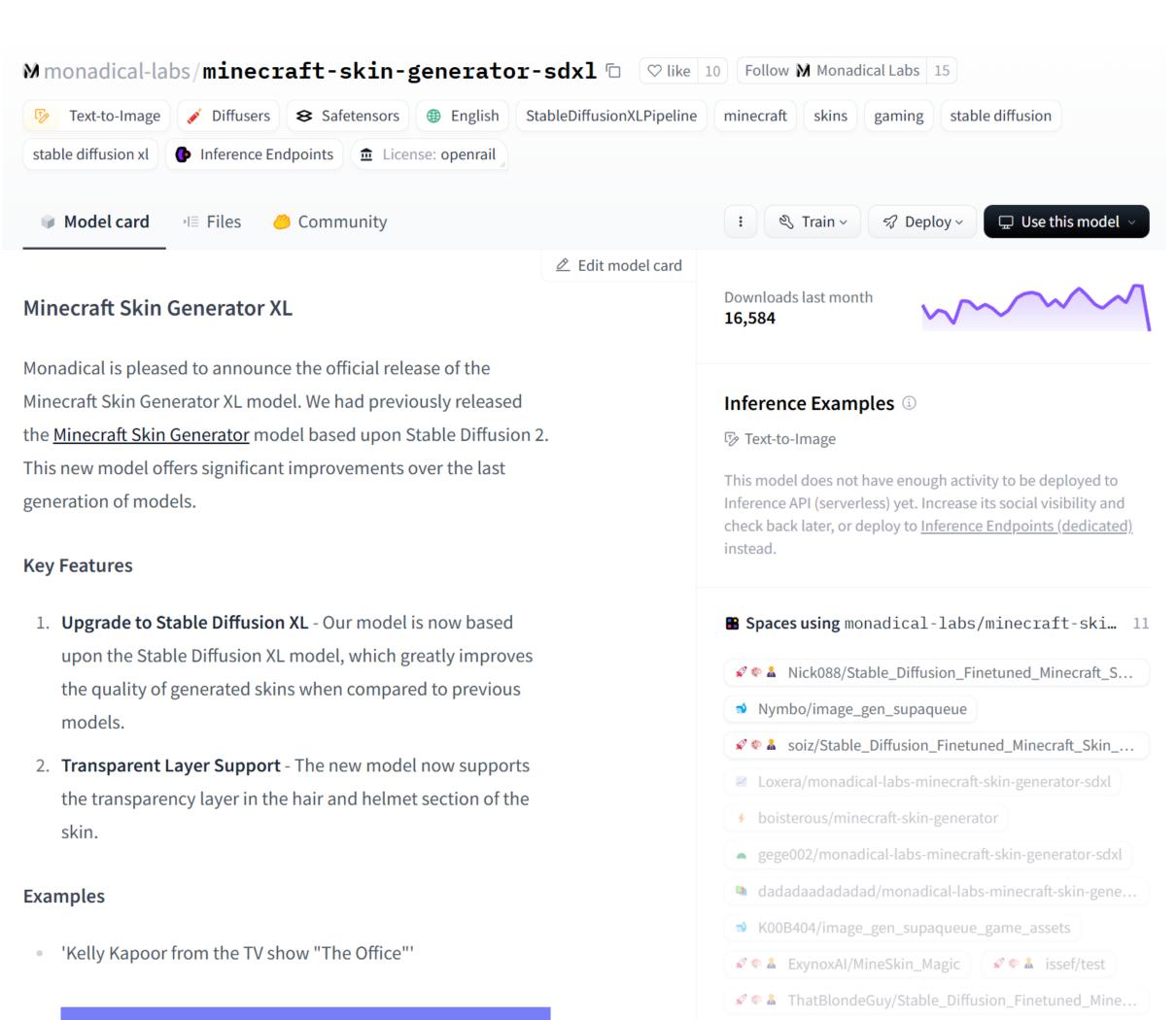




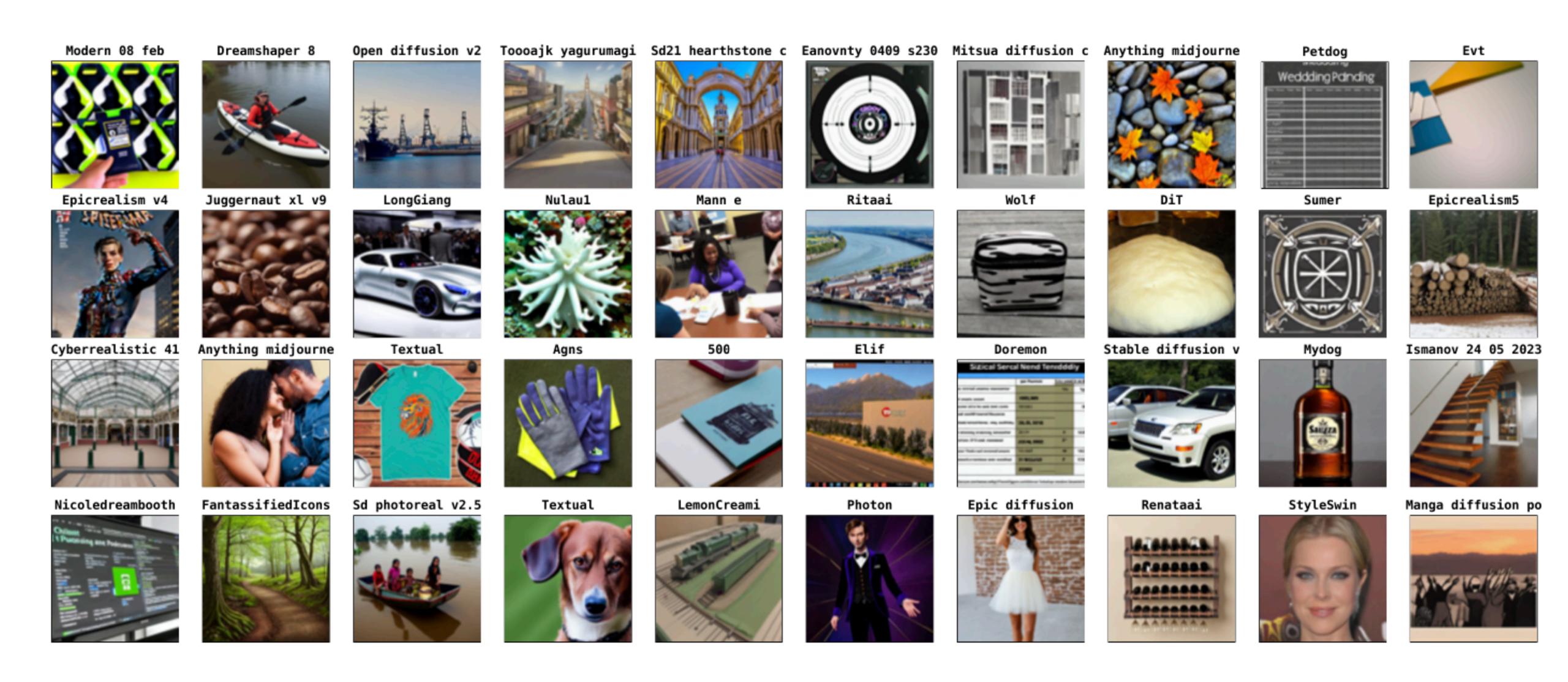








Automatically sampled images



Sampled images from 4763 LDM text-to-image models from Hugging Face

The Community Forensics dataset

Systematically collected diffusion models



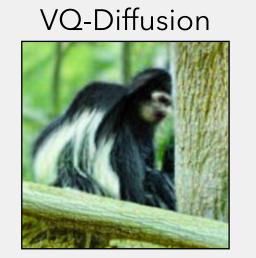
4763 LDM text-to-image models from Hugging Face

Other open models









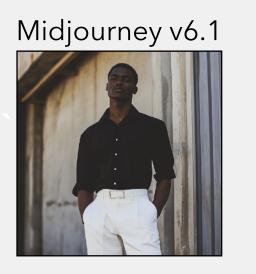
• • •

19 models

Commercial models





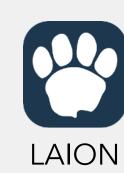


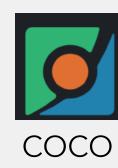


11 models

Real images

from 11 datasets, w/ paired prompts

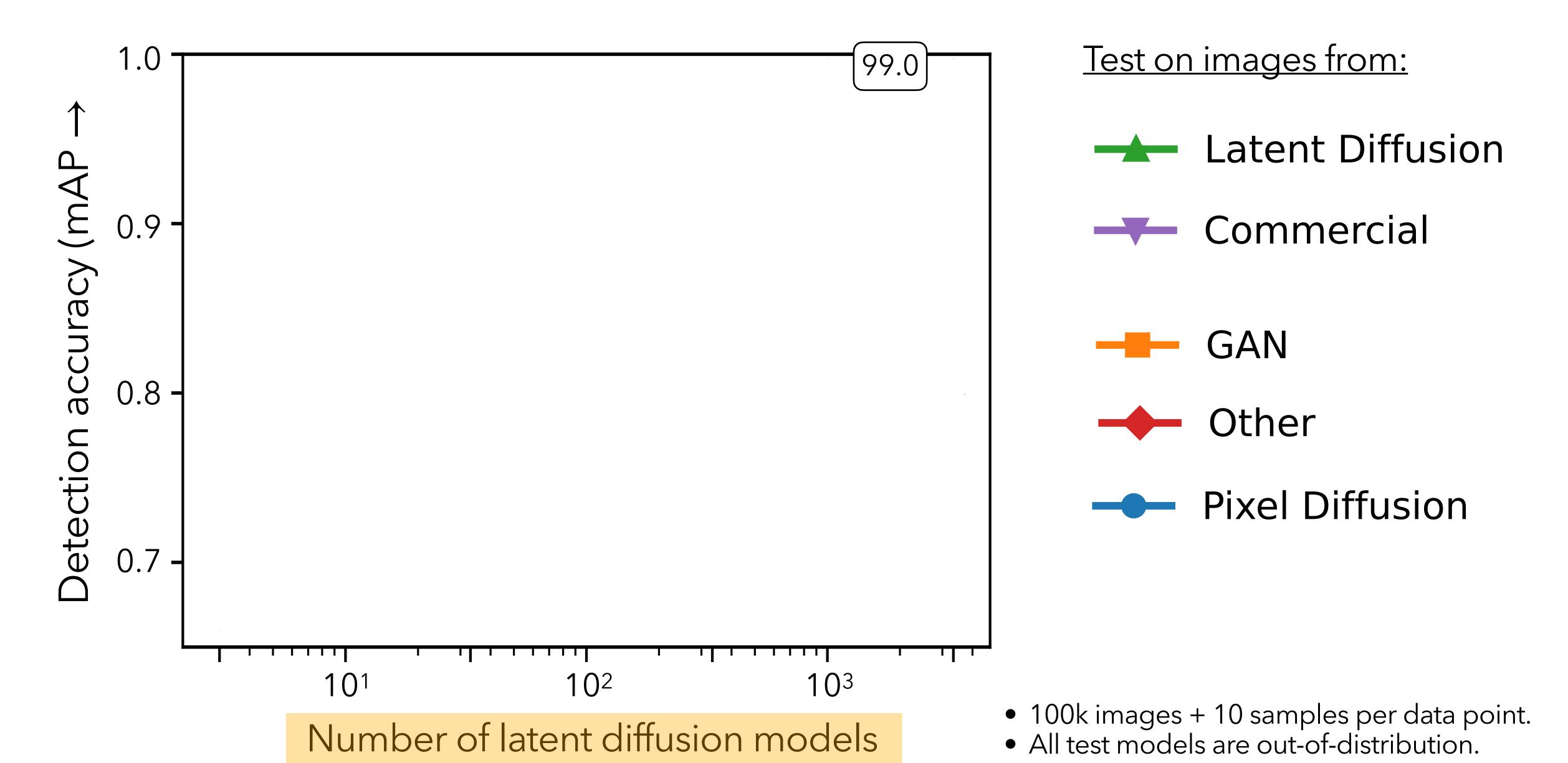






2.8M fake images from 4804 models.

What happens when you train on thousands of generators?



Challenges on the horizon

- Lots of ways to make fake images.
- If we know what methods were used, it's much easier.
- But it's hard to capture all of them!
- False positives are still a huge problem.
- So are postprocessing operations, like cropping and compression.
- Need methods that can handle unseen models.
- Alternative approaches: watermarking, signatures, etc.

Open-ended discussion

- How susceptible are people to fake images?
- Is there any hope of detecting "most" fake images?
- Under what situations might it be important and/or feasible?
- How do we deal with false positives?

Thank you!