# Lecture 17: Out-of-distribution generalization

Speaker: Sara Beery

# Machine Learning: A Success Story



Image Classification



Machine Translation



Strategy Games



Robotic Manipulation

Realistic Image Generation

# Are ML systems really ready for the real world?

# Standard ML setting



training distribution
=
test distribution

Training → Inference

# … vs the real world



deploy model on data from a different distribution

e.g.:

- perturbed data
- different label distribution
- other shifts (sequence/graph size, weather, country/city, source of measurement,…)

**Training** → **Inference**

**But:** In reali... ...utions we <u>use</u> ML on are N... ...we <u>train</u> it on
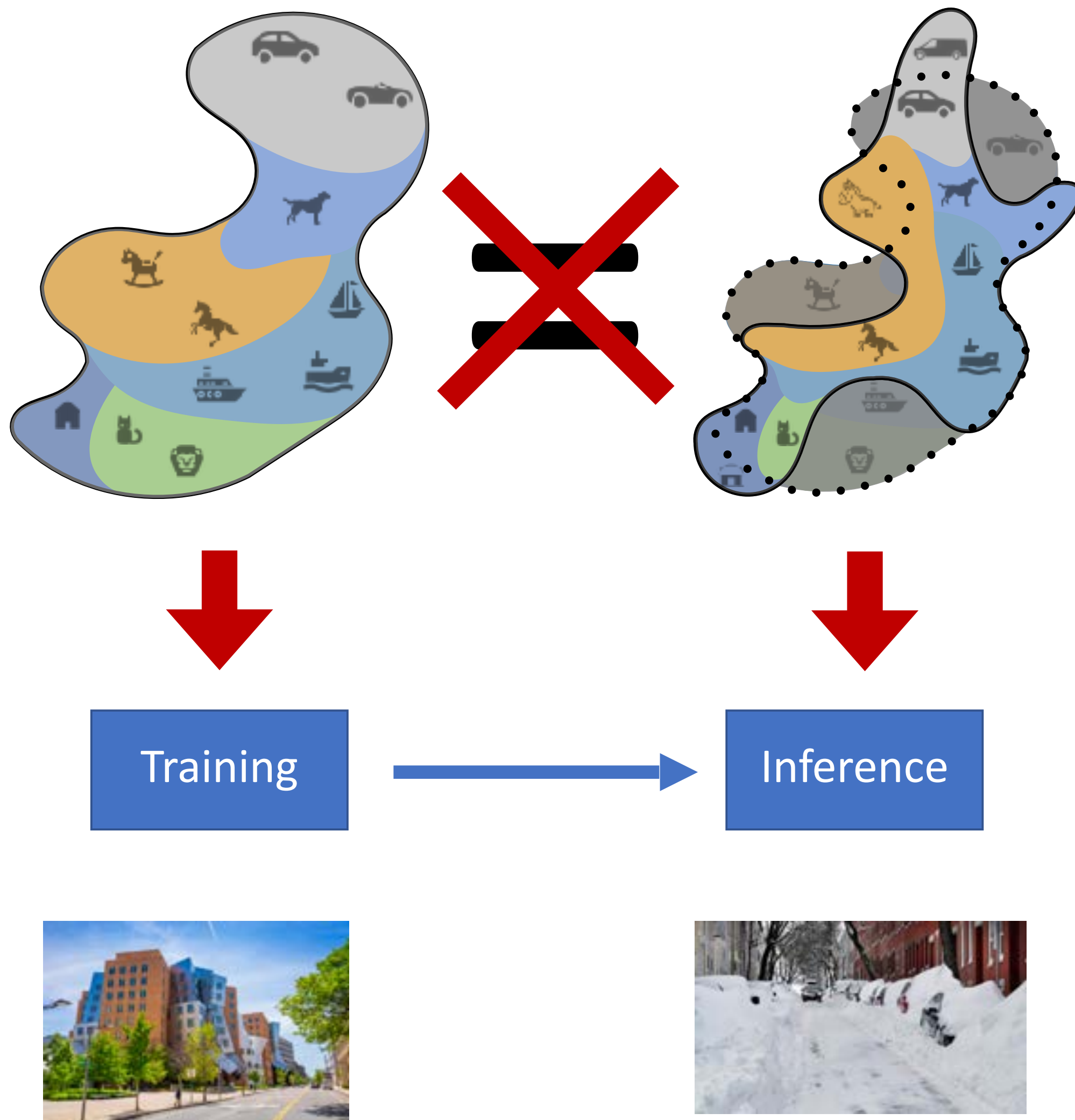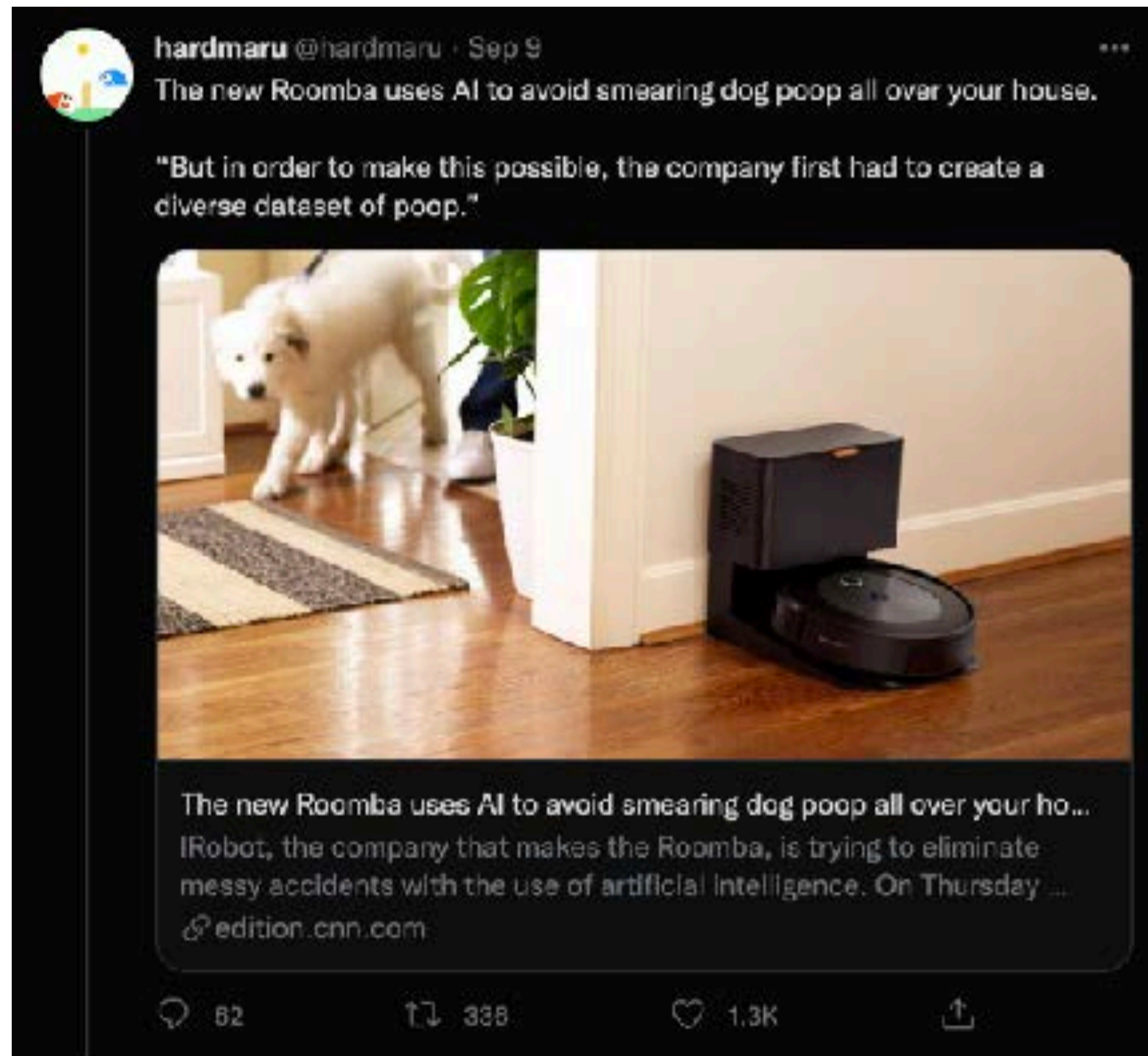
# What can go wrong?

# Concrete Problems in AI Safety

**Dario Amodei\***
Google Brain

**Chris Olah\***
Google Brain

**Jacob Steinhardt**
Stanford University

**Paul Christiano**
UC Berkeley

**John Schulman**
OpenAI

**Dan Mane**
Google Brain

might serve a benchmarking role similar to that of the bAbI tasks [163], with the eventual goal being to develop a single architecture that can learn to avoid catastrophes in all environments in the suite.

## 7 Robustness to Distributional Change

All of us occasionally find ourselves in situations that our previous experience has not adequately prepared us to deal with-for instance, flying an airplane, traveling to a country whose culture is very different from ours, or taking care of children for the first time. Such situations are inherently

# Outline for today

- Adversarial examples and training: small perturbations

- Distribution Shifts

# Adversarial examples

# Adversarial examples



"pig"

91% confidence

noise (not random)

+ 0.005 x

"airliner"

99% confidence

- ML model predictions are (mostly) accurate but can be brittle

10    *example: Szegedy et al 2013, obtained from https://gradientscience.org/intro_adversarial/*

# Adversarial examples

*Papernot et al 2017, Practical black-box attacks against machine learning*

# Adversarial stickers

# Adversarial stickers

13

*image: Brown et al 2018, Adversarial patch. https://youtu.be/i1sp4X57TL4*

# Adversarial stickers



Classifier Input

Classifier Output

banana   toaster   orange   crash_helm

14

*image: Brown et al 2018, Adversarial patch. https://youtu.be/i1sp4X57TL4*

# Adversarial stickers

15                                    *image: Brown et al 2018, Adversarial patch. https://youtu.be/i1sp4X57TL4*

# Adversarial examples 3D-printed

16

*image: Athalye et al 2018, Synthesizing robust adversarial examples*

# Adversarial examples 3D-printed



classified as turtle    classified as rifle
classified as other

*image: Athalye et al 2018, Synthesizing robust adversarial examples*
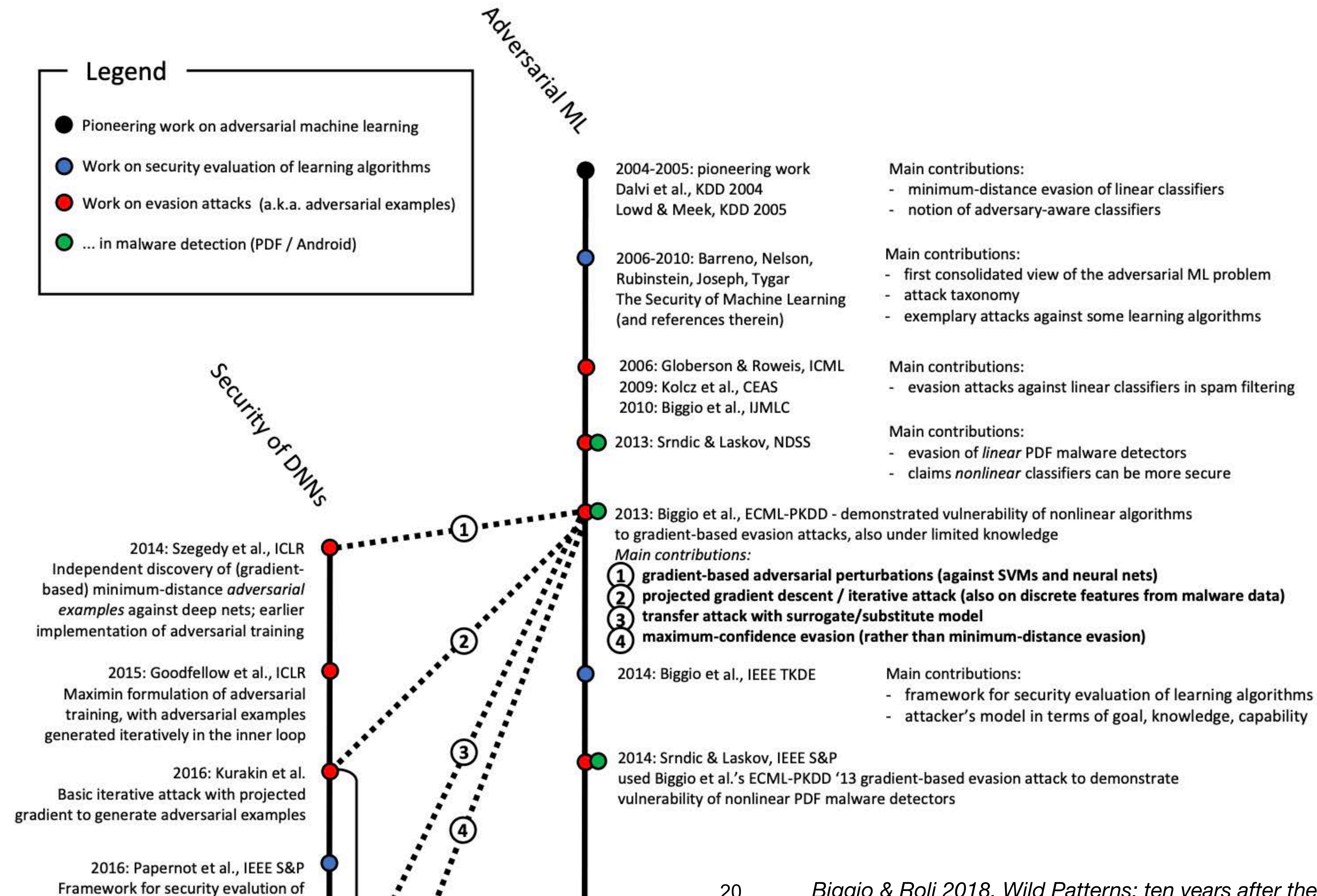
# Speech recognition example

*Carlini & Wagner 2018*

# Hmmmm....

- Are our models completely useless?

- Why does this happen?

- Can one prevent it?
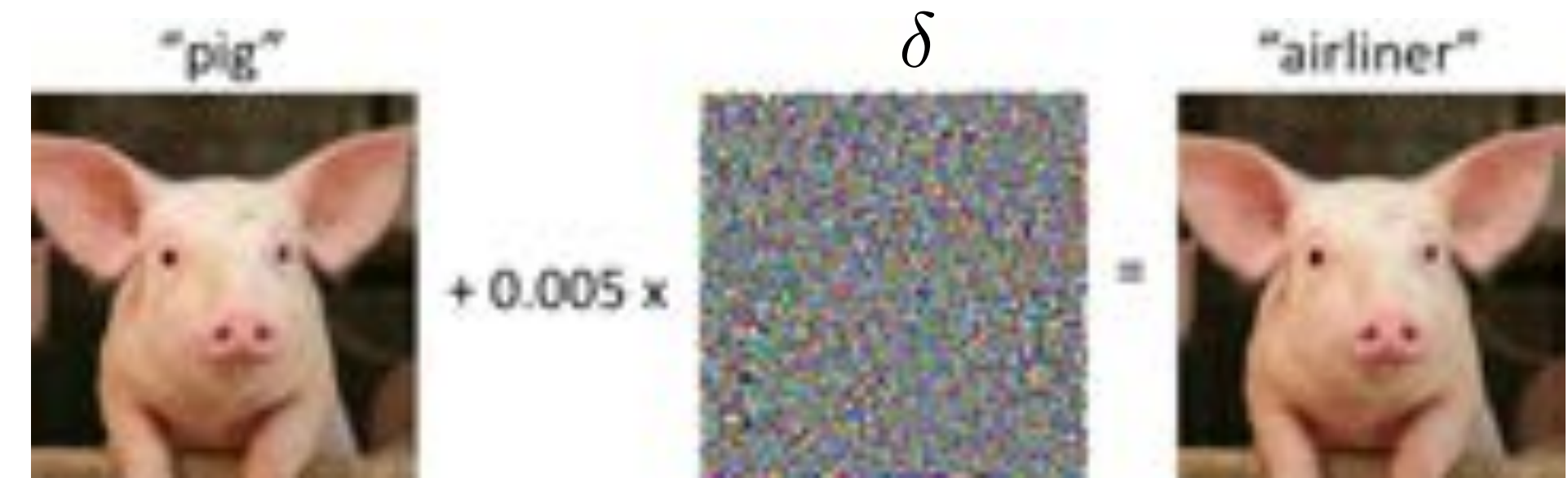
# History of adversarial examples / brittleness



**Adversarial ML**

**2004-2005: pioneering work**
Dalvi et al., KDD 2004
Lowd & Meek, KDD 2005

Main contributions:
- minimum-distance evasion of linear classifiers
- notion of adversary-aware classifiers

**2006-2010: Barreno, Nelson, Rubinstein, Joseph, Tygar**
The Security of Machine Learning
(and references therein)

Main contributions:
- first consolidated view of the adversarial ML problem
- attack taxonomy
- exemplary attacks against some learning algorithms

**2006: Globerson & Roweis, ICML**
**2009: Kolcz et al., CEAS**
**2010: Biggio et al., IJMLC**

Main contributions:
- evasion attacks against linear classifiers in spam filtering

**2013: Srndic & Laskov, NDSS**

Main contributions:
- evasion of *linear* PDF malware detectors
- claims *nonlinear* classifiers can be more secure

**2013: Biggio et al., ECML-PKDD** - demonstrated vulnerability of nonlinear algorithms to gradient-based evasion attacks, also under limited knowledge
*Main contributions:*
① **gradient-based adversarial perturbations (against SVMs and neural nets)**
② **projected gradient descent / iterative attack (also on discrete features from malware data)**
③ **transfer attack with surrogate/substitute model**
④ **maximum-confidence evasion (rather than minimum-distance evasion)**

**2014: Biggio et al., IEEE TKDE**

Main contributions:
- framework for security evaluation of learning algorithms
- attacker's model in terms of goal, knowledge, capability

**2014: Srndic & Laskov, IEEE S&P**
used Biggio et al.'s ECML-PKDD '13 gradient-based evasion attack to demonstrate vulnerability of nonlinear PDF malware detectors

**Security of DNNs**

**2014: Szegedy et al., ICLR**
Independent discovery of (gradient-based) minimum-distance *adversarial examples* against deep nets; earlier implementation of adversarial training

**2015: Goodfellow et al., ICLR**
Maximin formulation of adversarial training, with adversarial examples generated iteratively in the inner loop

**2016: Kurakin et al.**
Basic iterative attack with projected gradient to generate adversarial examples

**2016: Papernot et al., IEEE S&P**
Framework for security evaluation of

## Legend

20

*Biggio & Roli 2018, Wild Patterns: ten years after the rise of adversarial machine learning*

# How do you create an adversarial example?

- want: small perturbation that does not change meaning to a human, but to ML model



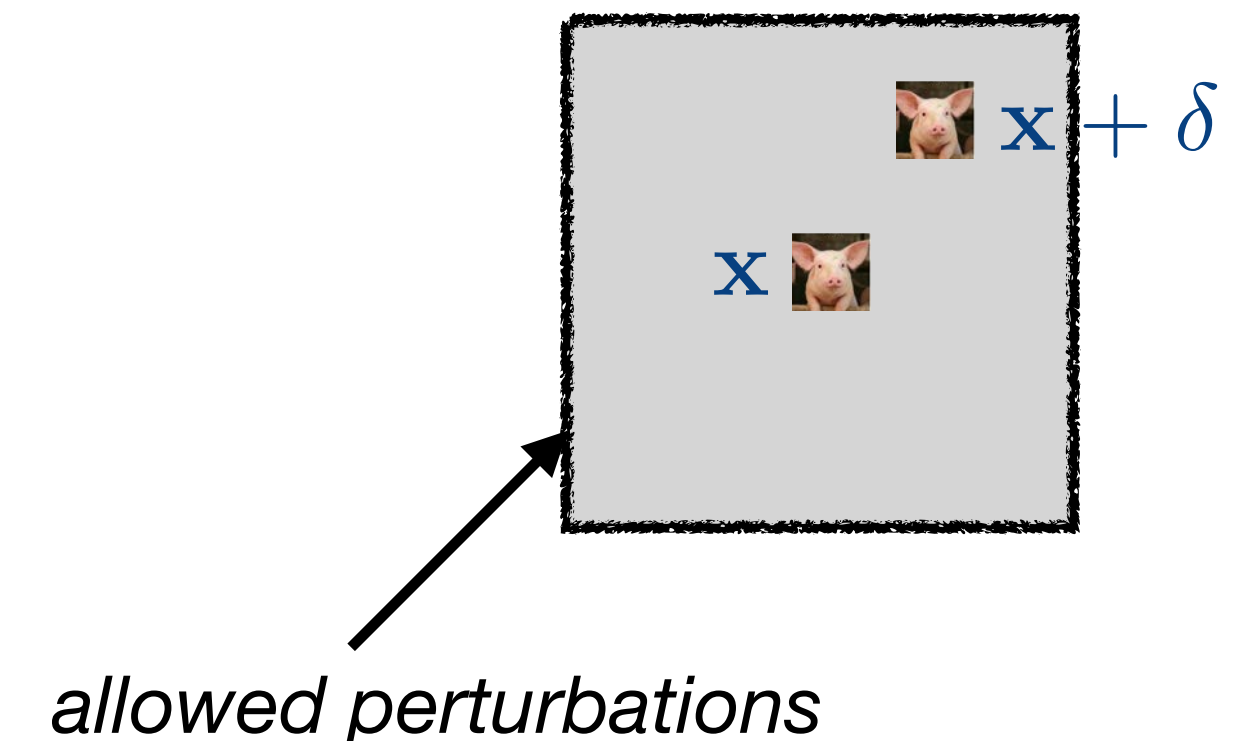- model outputs $P_\theta(y|\mathbf{x})$ (softmax)

- adversarial example: $$\max_{\delta \in \Delta} P_\theta(y_{\text{target}} \mid \mathbf{x}+\delta)$$

**small perturbation, e.g.**

**wrong class ("airliner")**

**input image**

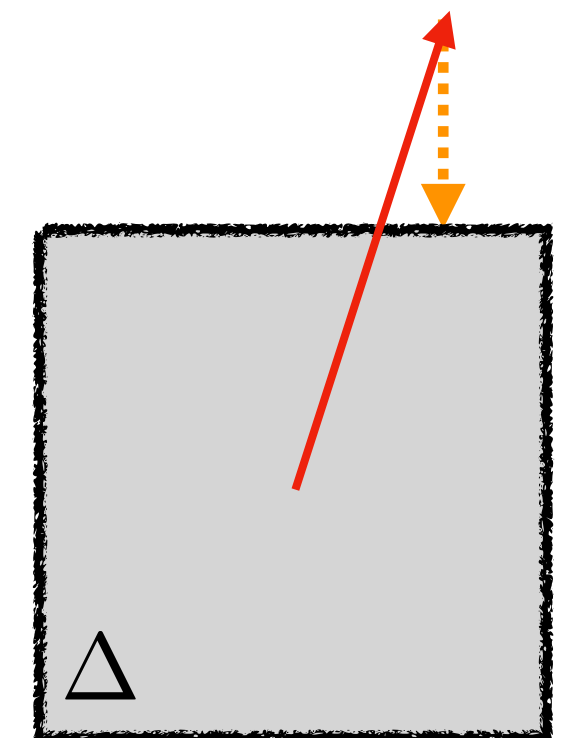$$\Delta = \{\delta \in \mathbb{R}^d \mid \|\delta\|_\infty < \epsilon\}$$

$\mathbf{x}+\delta$

$\mathbf{x}$

*allowed perturbations*

# How to find an adversarial example?

**ML model**

$$\max_{\delta \in \Delta} \boxed{P_\theta}(y_{\text{target}} \mid \mathbf{x} + \delta)$$

**small perturbation, e.g.**

**wrong class ("airliner")**

**input image**

- e.g. Projected gradient ascent (we update data perturbation $\delta$!):
  1. take a step in the direction of the gradient:
     $$\delta^{(t+1)} = \delta^{(t)} + \eta \cdot \nabla_\delta P_\theta(y_{\text{target}} \mid \mathbf{x} + \delta)$$
  2. project the result back into the feasible set $\Delta$
  3. repeat steps 1 & 2

$\Delta$

# How to "defend" against adversarial examples?

Recall:

- Adversarial example       versus      standard training:

$$\max_{\delta \in \Delta} \mathrm{Loss}\Big( f_\theta(\mathbf{x}+\delta), y \Big) \qquad \min_\theta \mathrm{Loss}\Big( f_\theta(\mathbf{x}), y \Big)$$

# How to "defend" against adversarial examples?

- **Standard training**:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \text{Loss}(f_{\theta}(\mathbf{x}^{(i)}), y^{(i)})$$

via (stochastic) gradient descent

**neural network**

- **Adversarial training** / robust optimization:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \max_{\delta \in \Delta} \text{Loss}\left(f_{\theta}(\mathbf{x}^{(i)} + \delta), y^{(i)}\right)$$

*"adaptive data augmentation"*

24

# Adversarial training with stochastic gradient descent

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \boxed{\max_{\delta \in \Delta}} \text{Loss}\left(f_\theta(\mathbf{x}^{(i)} + \delta), y^{(i)}\right)$$

repeat until convergence:

1. sample a data point $(\mathbf{x}, y)$

2. compute the optimal adversarial perturbation $\delta^*$  *(approximately)*

3. compute the gradient $g = \nabla_\theta \text{Loss}\left(f_\theta(\mathbf{x} + \delta^*), y\right)$

4. update $\theta$ with the gradient $g$

# What do adversarial examples tell us?

- something about the input "features" that are critical for the model's decision

- Example:



**Training data:**
**classify 4 vs 9**

**Adversarial**
**perturbations**

*images: Hongzhou Lin*

26

# Predictive features



**Useless** features

**Robust features**
Correlated with label even when perturbed

**Non-robust features**
Correlated with label, but can be flipped via perturbation

- Many features may be correlated with the label and hence predictive and help with accuracy, *beyond what humans would use.*

*illustration: Aleksander Madry*

# Where do these correlations come from?

- Data



Dogs

Cats

"Fish" from the ImageNet training set

# Where do these correlations come from?

- …and how we create datasets

**In fact:** The way we create datasets gives rise to them



**Ideal world:**

Real-world images → Expert annotators → Perfect annotations (dog, bird, truck, primate) → Meaningful benchmark → 🌈

**In fact:** The way we create datasets gives rise to them

**~~Ideal~~ Real world:**

Flickr/scraped images → Automated + Crowd Labels → Noisy, biased annotations (dog, bird, truck, primate) → Easy-to-optimize benchmark → ❓

# It's all "shortcuts"

- Shortcuts: features correlated with label in the training data, but not under realistic distribution shifts

- Models will use them and not generalize if features are no longer correlated

*illustration: Geirhos et al 2020*

# It's all "shortcuts"

- Shortcuts: features correlated with label in the training data, but not under realistic distribution shifts

- Models will use them and not generalize if features are no longer correlated

- This is related to data, not models: *adversarial examples transfer across models trained on the same dataset*

# What can these shortcuts look like?



A herd of sheep grazing on a l
Tags: grazing, sheep, mountai

Left: A man is holding a dog in his hand
Right: A woman is holding a dog in her hand

Image: @SouperSarah

lock of birds flying in the air
up of giraffe standing next to a tree
www.flickr.com/photos/gratapictures - CC-BY-NC

images: https://www.aiweirdness.com/do-neural-nets-dream-of-electric-18-03-02/

# What can these shortcuts look like?



"…if an image had a ruler in it, the algorithm was more likely to call a tumor malignant…"

[Esteva et al. 2017]



"CNNs were able to detect where an x-ray was acquired […] and calibrate predictions accordingly."

[Zech et al. 2018]

**not all predictive patterns are desirable**

# Many more…



**Same category for humans** but not for DNNs (intended generaliszation)

i.i.d.

| Domain shift Wang 2018 | Adversarial examples Szegedy 2013 | Distortions Dodge 2019 | Pose Alcorn 2019 | Texture Geirhos 2019 | Background Beery 2018 |

o.o.d.

**Same category for DNNs** but not for humans (unintended generalization)

i.i.d.

| Excessive invariance Jacobson 2019 | Fooling images Nguyen 2015 | Natural adversarials Hendrycks 2019 | Texturized images Brendel 2019 |

*illustration: Geirhos et al 2020, Shortcut learning in deep neural networks*

# Transformers Learn Shortcuts to Automata

Bingbin Liu[1*]    Jordan T. Ash[2]    Surbhi Goel[2,3]    Akshay Krishnamurthy[2]    Cyril Zhang[2]

[1]Carnegie Mellon University    [2]Microsoft Research NYC    [3]University of Pennsylvania
bingbinl@cs.cmu.edu, {ash.jordan, goel.surbhi, akshaykr, cyrilzhang}@microsoft.com

## Abstract

Algorithmic reasoning requires capabilities which are most naturally understood through recurrent models of computation, like the Turing machine. However, Transformer models, while lacking recurrence, are able to perform such reasoning using far fewer layers than the number of reasoning steps. This raises the question: *what solutions are these shallow and non-recurrent models finding?* We investigate this question in the setting of learning automata, discrete dynamical systems naturally suited to recurrent modeling and expressing algorithmic tasks. Our theoretical results completely characterize *shortcut solutions*, whereby a shallow Transformer with only $o(T)$ layers can exactly replicate the computation of an automaton on an input sequence of length $T$. By representing automata using the algebraic structure of their underlying transformation semigroups, we obtain $O(\log T)$-depth simulators for all automata and $O(1)$-depth simulators for all automata whose associated groups are so[...] synthetic experiments by training Transformers to simulate a wide var[...] shortcut solutions can be learned via standard training. We further in[...] solutions and propose potential mitigations.

**parallel solutions generalize within-distribution,
but not out-of-distribution**
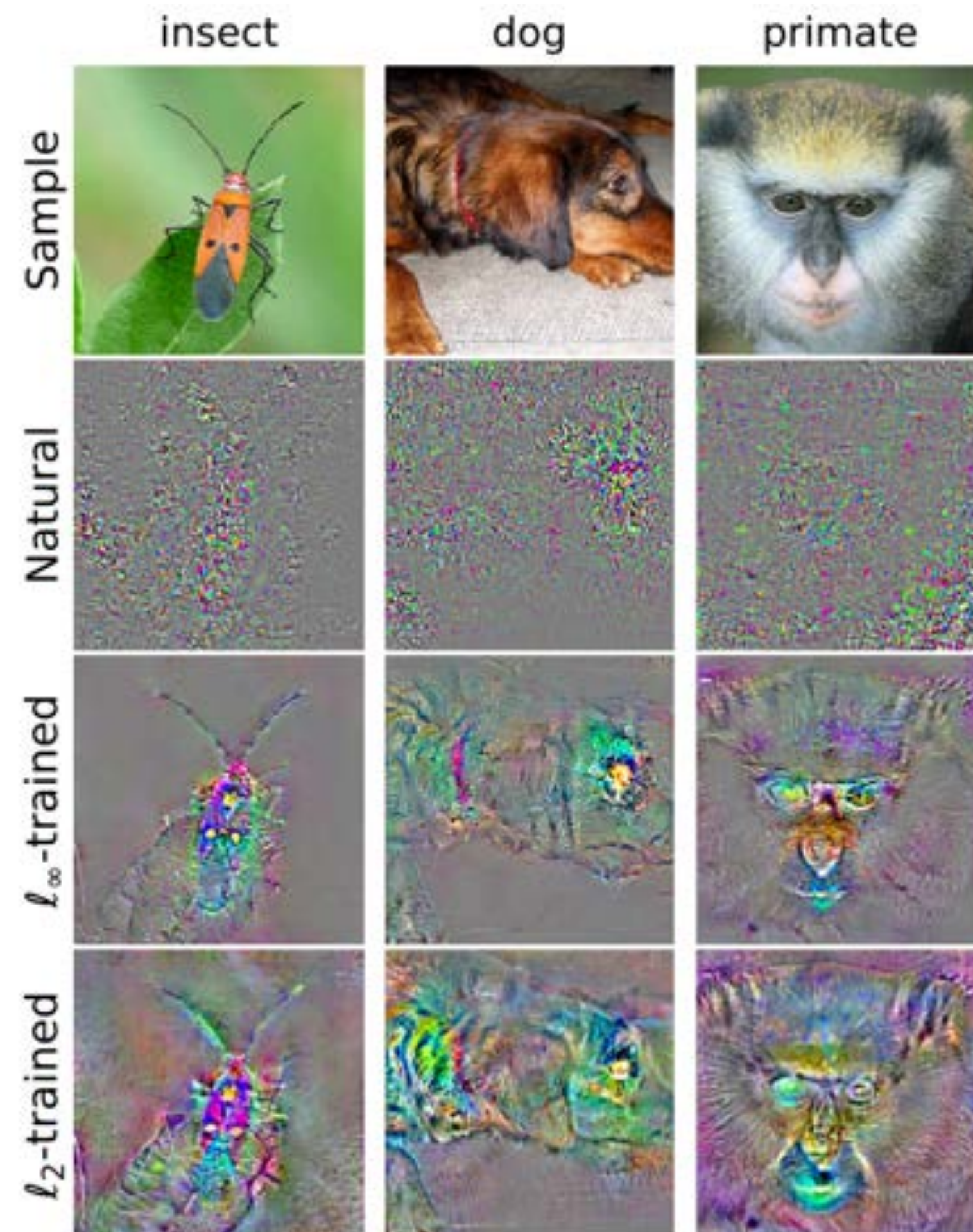
35

# Effect of adversarial training



**Useless** features

**Robust features**
Correlated with label even when perturbed

**Non-robust features**
Correlated with label, but can be flipped via perturbation

- model output should be stable under adversarial perturbations
  => teaches invariance to non-robust features

# Effect of adversarial training
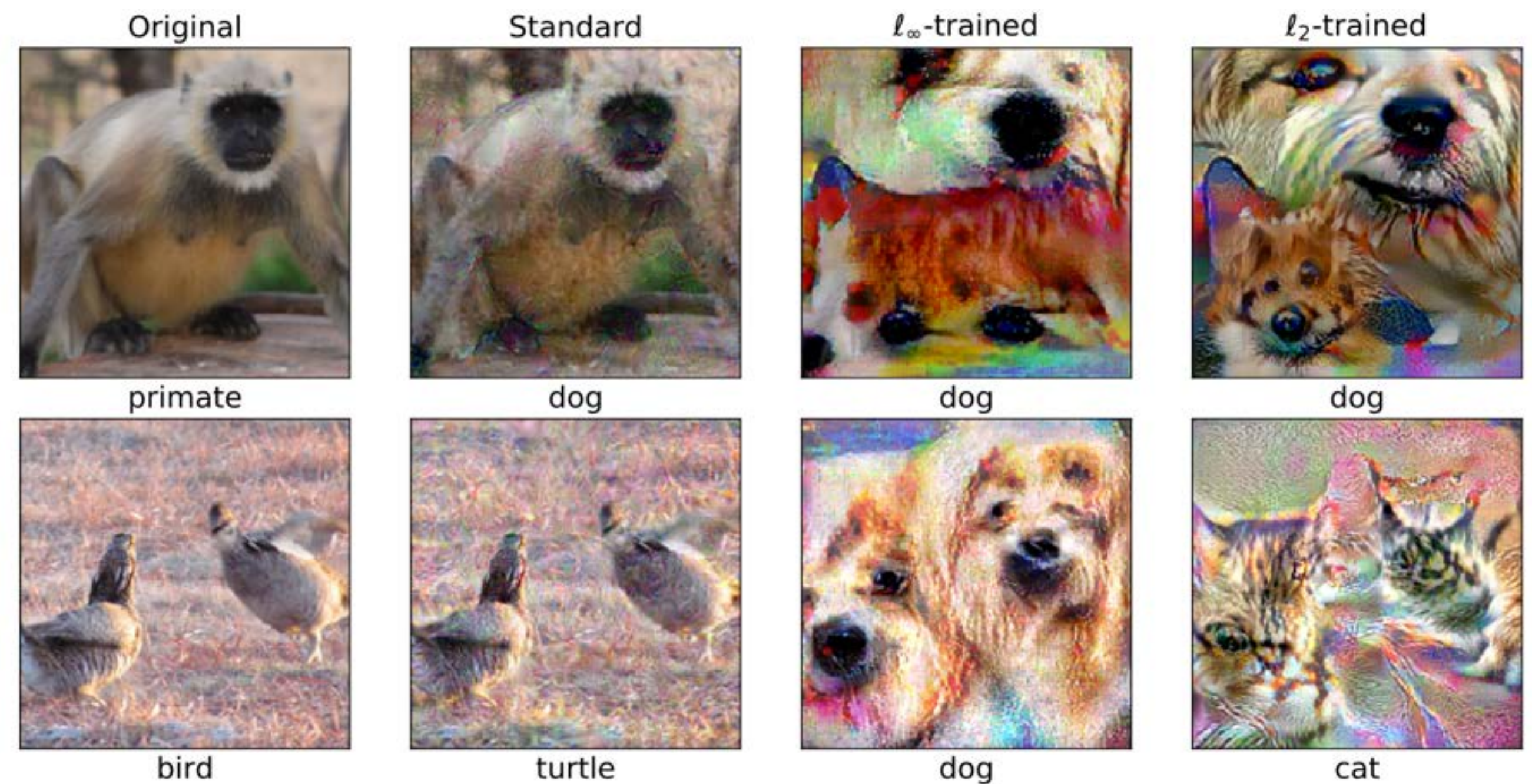
Loss gradients with respect to input pixels (most important features) show: robust model relies less on "non-robust" features, and more on human-intuitive features





**Adversarial examples for standard and robust models**

*(Tsipras et al. 2019, Robustness may be at odds with accuracy.)*

# Effect of adversarial training: transfer learning

- adversarially trained models transfer better to other datasets

*(Salman et al. 2020, Do adversarially robust ImageNet models transfer better?)*

# Distribution shifts

Training → Inference

Training → Inference

Fr

**But:** In reality, the distributions we u_
are NOT the ones we train it

40

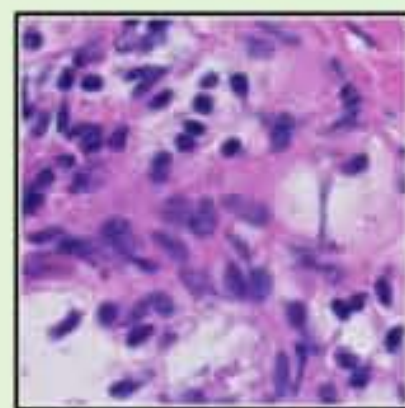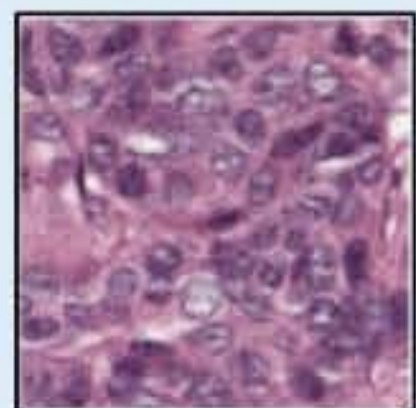**Map of global biodiversity**

**Species occurrence data in GBIF**

# WILDS

Pang Wei Koh*, Shiori Sagawa*, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang

| | Camelyon17 | iWildCam | PovertyMap | FMoW | Amazon | CivilComments | OGB-MolPCBA |
|---|---|---|---|---|---|---|---|
| Shift | Hospitals | Locations | Countries | Time | Users | Demographics | Scaffold |
| Train | | | | | Overall a solid package that has a good quality of construction for the price. | What do Black and LGBT people have to do with bicycle licensing? | |
| Test | | | | | I *loved* my French press, it's so perfect and came with all this fun stuff! | As a Christian, I will not be patronizing any of those businesses. | |
| Adapted from | Bandi et al. 2018 | Beery et al. 2020 | Yeh et al. 2020 | Christie et al. 2018 | Ni et al. 2019 | Borkan et al. 2019 | Hu et al. 2020 |

shifts across hospitals in histopathology

ID accuracy
93.2%
-22.9%
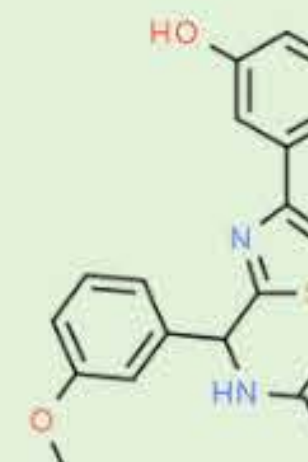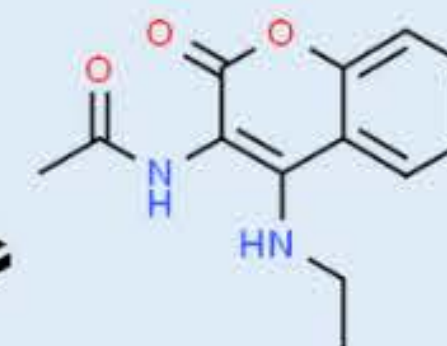OOD accuracy
70.3%

shifts across time in satellite imagery

ID accuracy
48.6%
-16.3%
OOD accuracy
32.3%

shifts across regions in wheat head detection

ID accuracy
63.3%
-13.7%
OOD accuracy
49.6%

shifts across scaffold in bioassay prediction

ID AP
34.4%
-7.2%
OOD AP
27.2%

[Koh et al., 2021]

*Slide from Shiori Sagawa*

43

# Training data

## Camera 1          Camera 2     ...     Camera 245



# Out-of-distribution (OOD) test data

## Camera 246



...

# Control: In-distribution (ID) test data

## Camera 1          Camera 2     ...     Camera 245
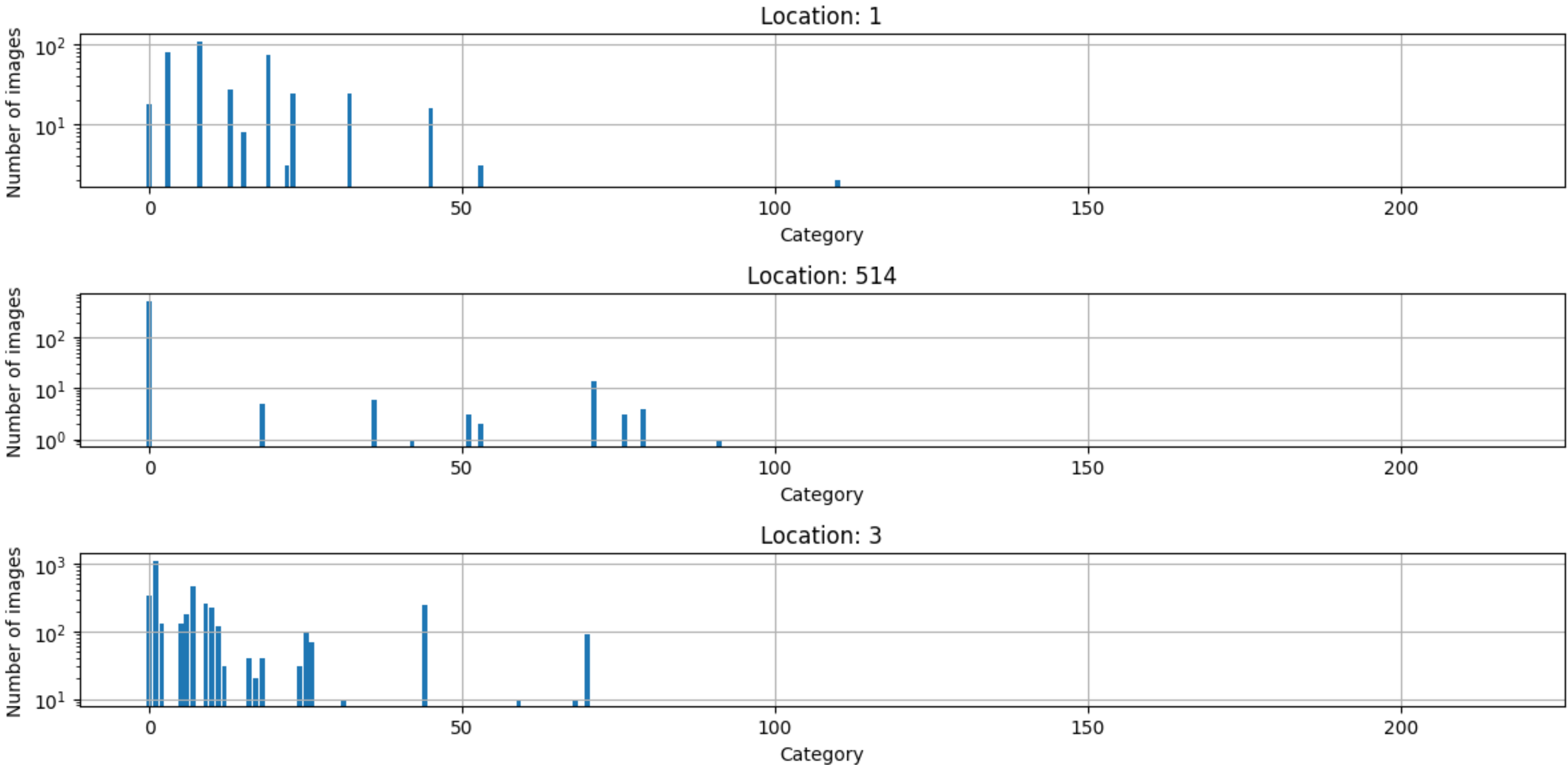


## Macro F1

ID          **-16.0%**          OOD

47.0%  ⟶  31.0%

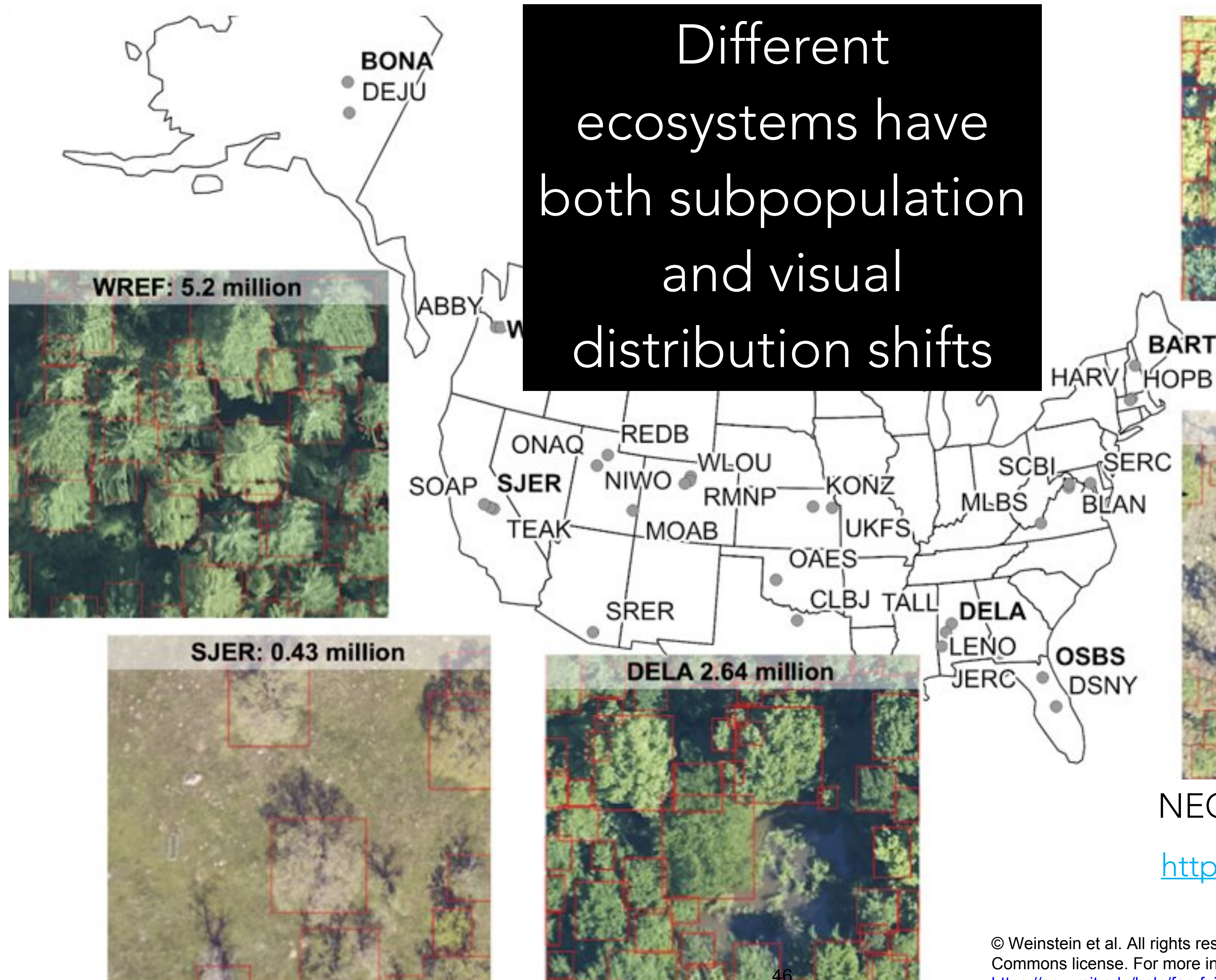[Beery et al., 2020; Koh et al., 2021]

*Slide from Shiori Sagawa*

# Class distribution is different for each static sensor location

Different ecosystems have both subpopulation and visual distribution shifts

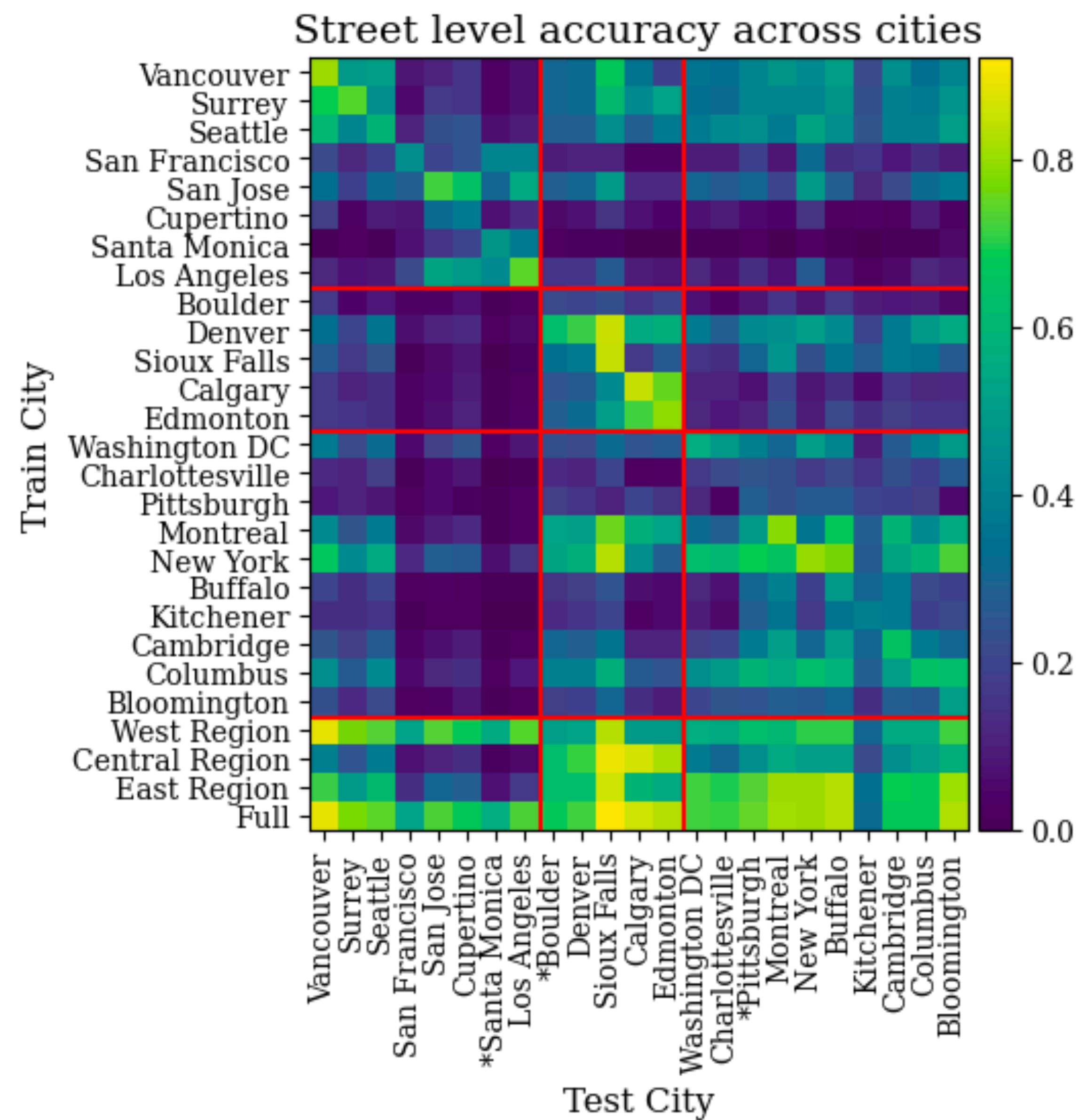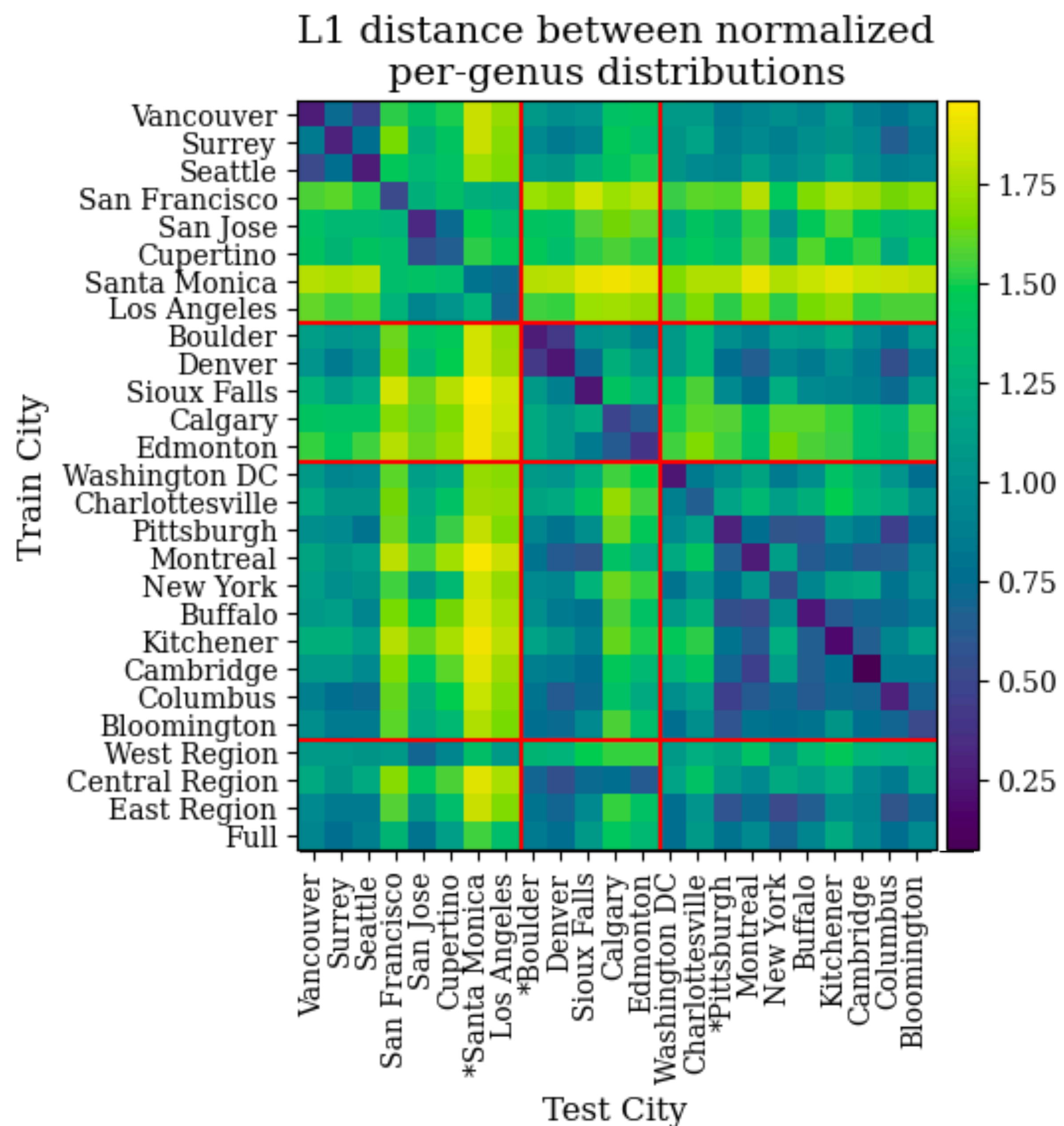WREF: 5.2 million

SJER: 0.43 million

DELA 2.64 million

OSBS: 5.13 million

NEONCROWNS Dataset

http://visualize.idtrees.org/

Weinstein et al., 2020

# Performance has strong correlation with subpop. distribution similarity



L1 distance between normalized per-genus distributions

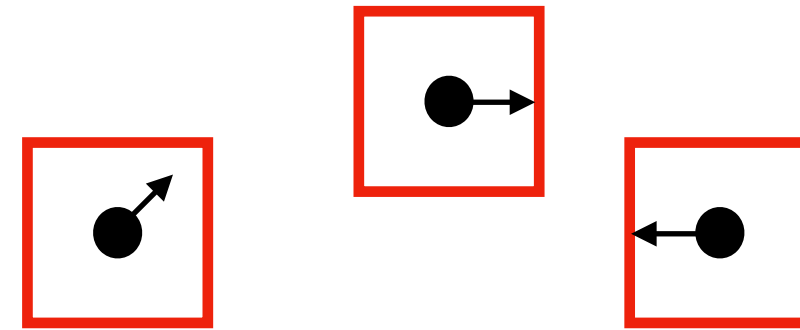Street level accuracy across cities

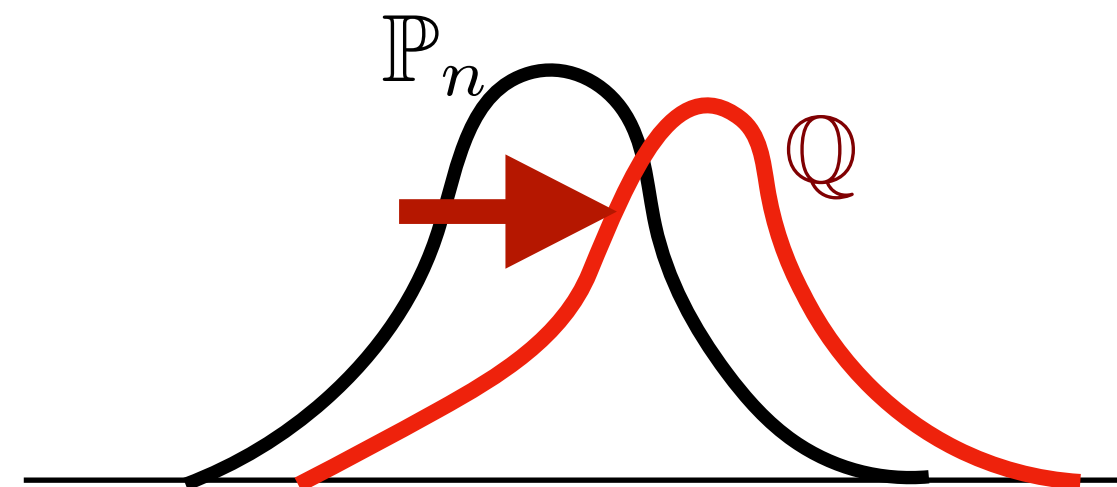# What to do about distribution shift?

# One path: distributionally robust optimization

- So far: allowed to perturb each datapoint by a limited amount



- Alternative: we can perturb the entire training distribution (sample) by a certain amount, together

# Distributionally robust optimization

- Standard training:
$$\frac{1}{n}\sum_{i=1}^{n}\text{Loss}(f_\theta(\mathbf{x}^{(i)}), y^{(i)}) = \mathbb{E}_{(\mathbf{x},y)\sim\mathbb{P}_n}[\text{Loss}(f_\theta(\mathbf{x}), y)]$$
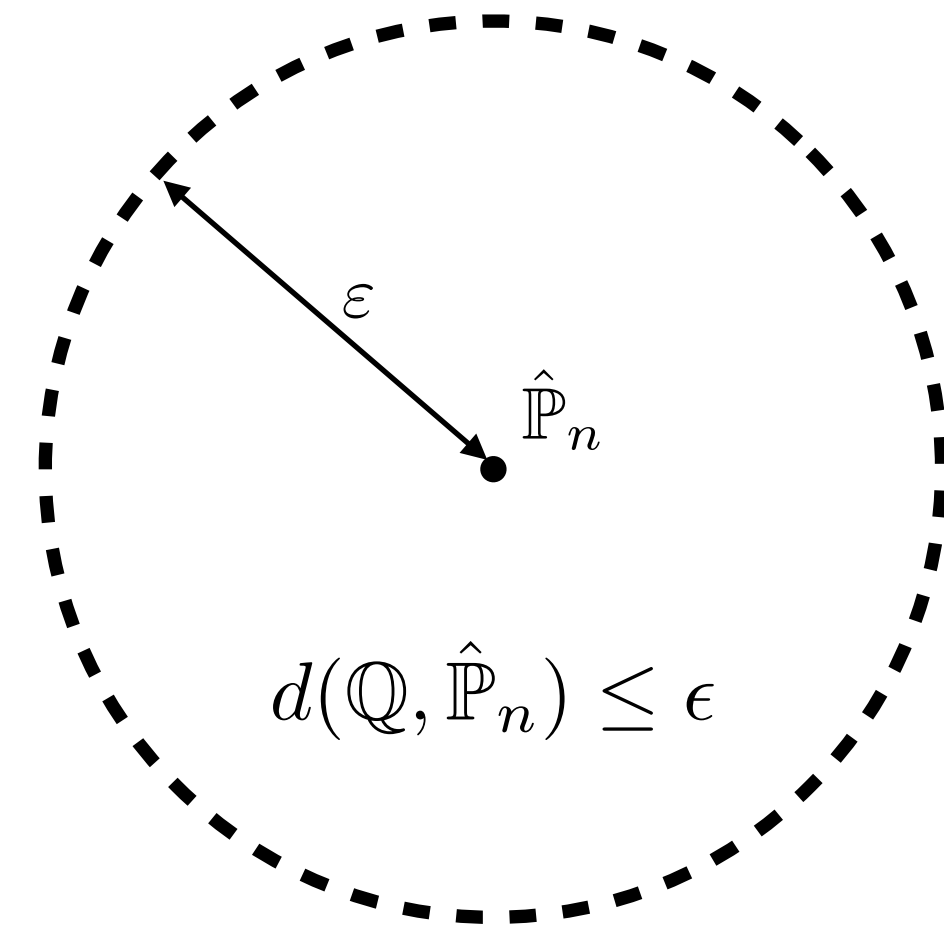
**allow a small
perturbation of
training sample
(discrete distribution)**

- Distributionally robust optimization (DRO):

$$\min_\theta \max_{\mathbb{Q},\, D(\mathbb{Q},\mathbb{P}_n)<\epsilon} \mathbb{E}_{(\mathbf{x},y)\sim\mathbb{Q}}[\text{Loss}(f_\theta(\mathbf{x}), y)]$$

**e.g. re-weight or
perturb training data points**

$\varepsilon$

$\hat{\mathbb{P}}_n$
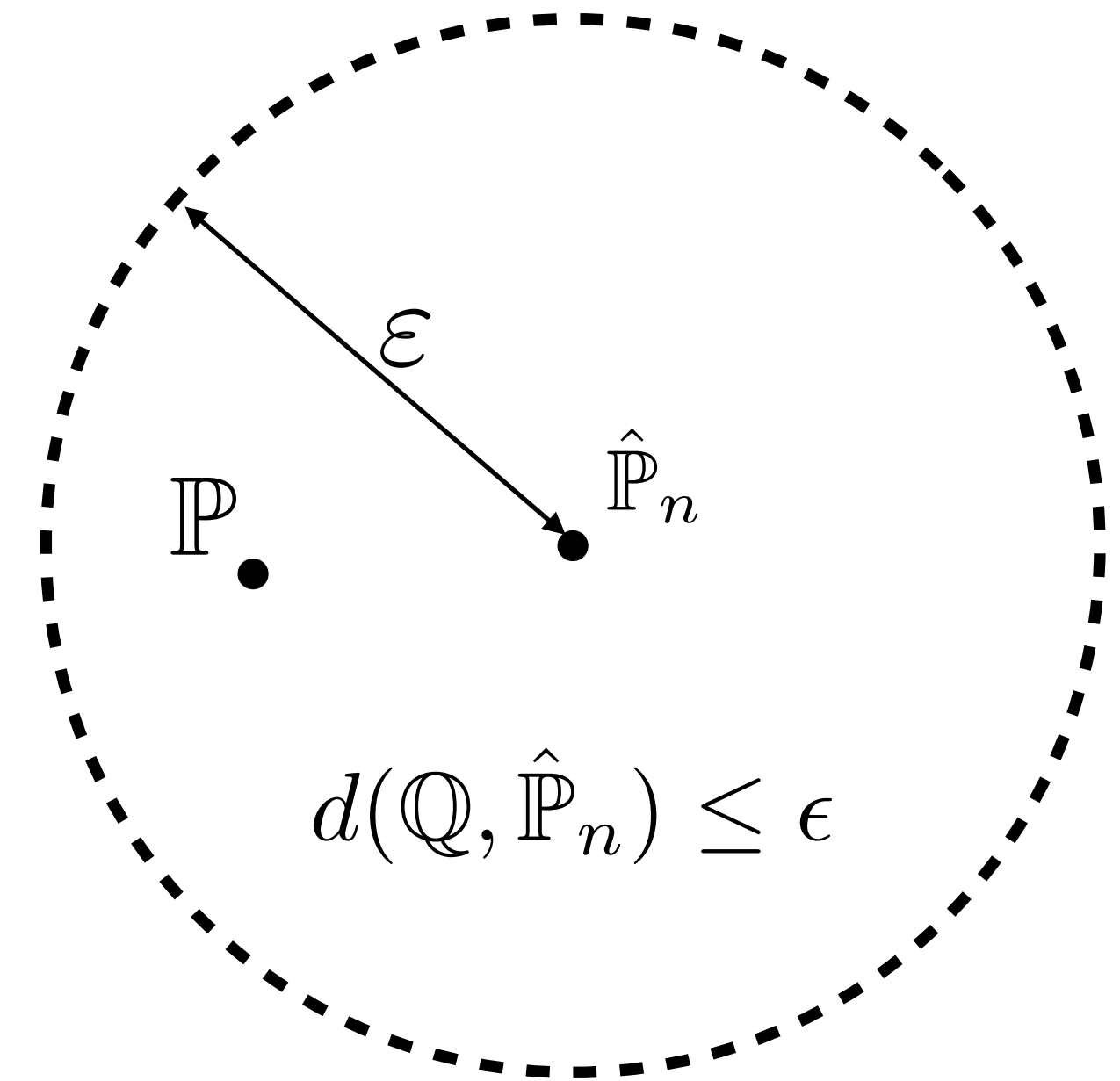
$d(\mathbb{Q}, \hat{\mathbb{P}}_n) \leq \epsilon$

- Various choices of measuring "distance" between probability distributions: $\chi^2$-distance, Wasserstein distance, maximum mean discrepancy (MMD)…

# DRO and generalization

$$\min_{\theta} \quad \textcolor{red}{\max_{\mathbb{Q}, \, D(\mathbb{Q}, \mathbb{P}_n) < \epsilon}} \mathbb{E}_{(\mathbf{x}, y) \sim \textcolor{red}{\mathbb{Q}}}[\text{Loss}(f_{\theta}(\mathbf{x}), y)]$$

- DRO optimizes for a set of training data sets/ distributions

- Say underlying data distribution is $\mathbb{P}$

- Empirical training data is $\hat{\mathbb{P}}_n$

- If $D(\mathbb{P}, \hat{\mathbb{P}}_n) < \epsilon$, then we are guaranteed to perform well on $\mathbb{P}$ too, i.e., generalize!

$\varepsilon$

$\mathbb{P}$ $\quad\quad$ $\hat{\mathbb{P}}_n$

$d(\mathbb{Q}, \hat{\mathbb{P}}_n) \leq \epsilon$

# Application: DRO and class imbalance

- Assume population has K sub-groups (example: K=2).
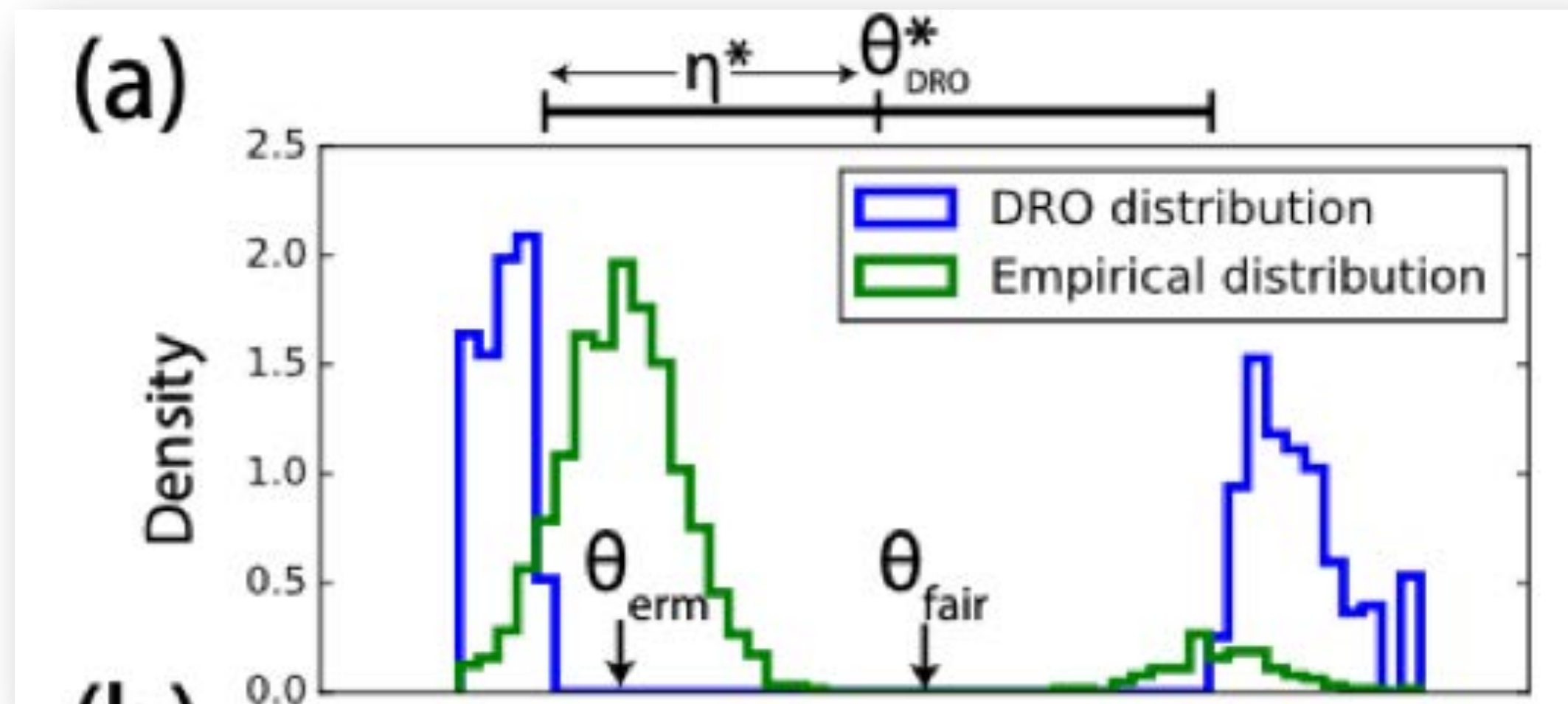- Usually: minimize "Empirical Risk" (average error)

$$\min_{\theta} \quad \frac{1}{n}\Big( \underbrace{\sum_{i \text{ in group } 1} \text{Loss}(x_i; \theta)}_{\textbf{80\%}} + \underbrace{\sum_{j \text{ in group } 2} \text{Loss}(x_j; \theta)}_{\textbf{20\%}} \Big)$$

- Here, 50% error on minority group makes only 10% average error. (+ statistical patterns for minority may be different)

- We can "ignore" minority group and still get decent loss!

# DRO and class imbalance

- Idea: automatically re-weight data via DRO
  => pay more attention to minority class

$$\min_{\theta} \ \max_{\mathbb{Q}, \, D(\mathbb{Q}, \mathbb{P}_n) < \epsilon} \ \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\mathrm{Loss}(f_\theta(\mathbf{x}), y)]$$

*(Hashimoto et al. 2018, Fairness without demographics in repeated loss minimization)*
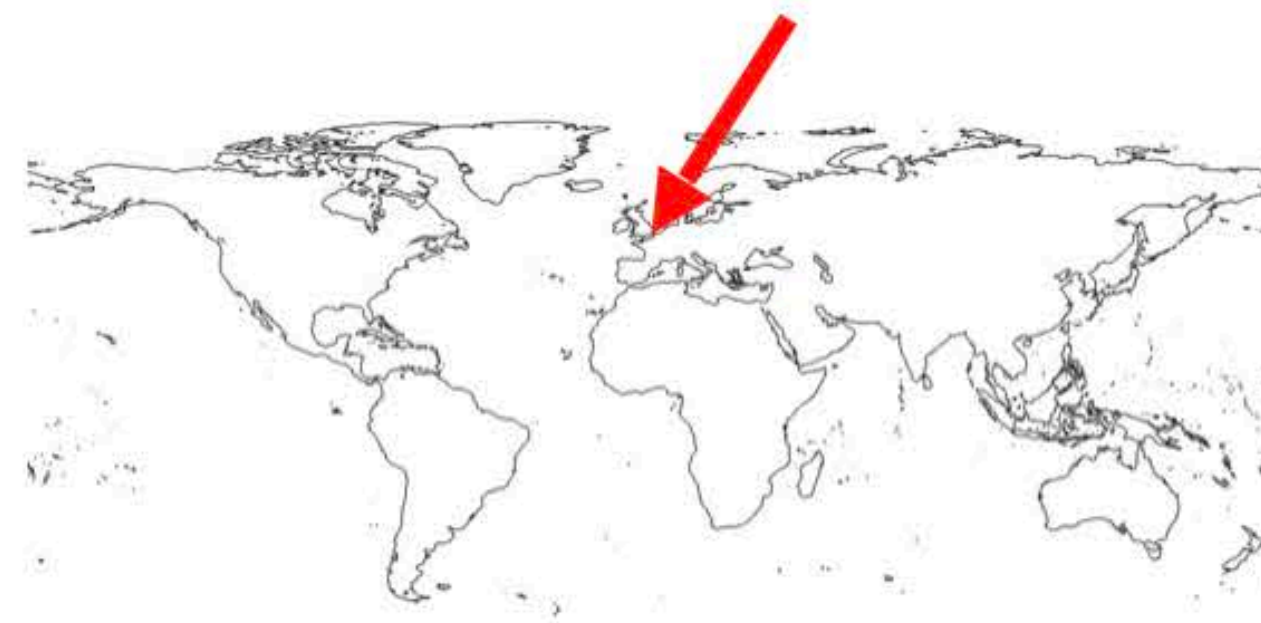
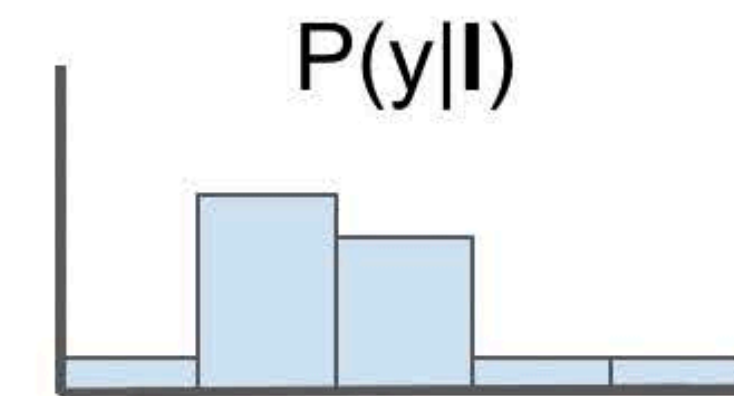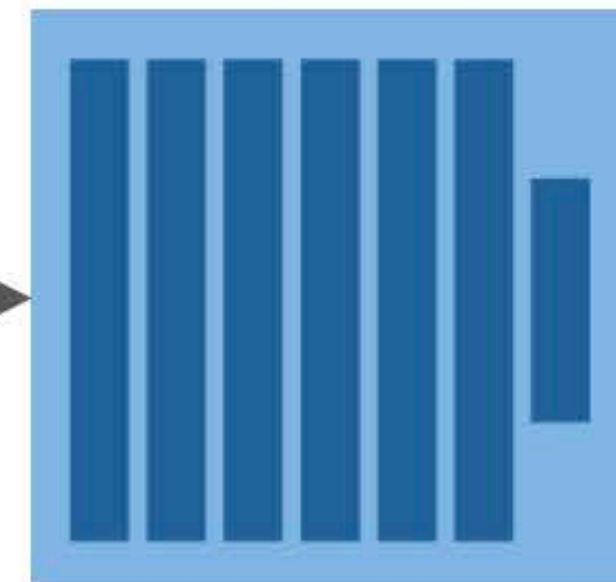# What to do about distribution shift?

- Distributionally robust optimization
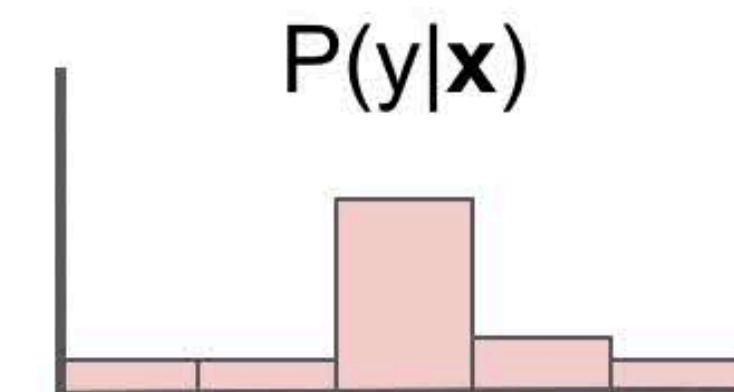
# Learn a spatiotemporal prior

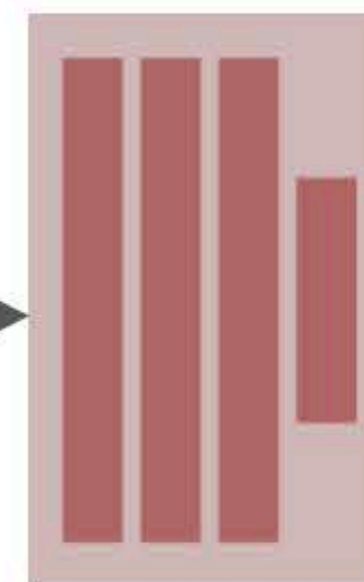$$P(y|I, \mathbf{x}) \propto P(y|I)P(y|\mathbf{x})$$



Image Classifier

Cesar Pollo BY CC-NC 4.0

$P(y|\mathbf{I})$

Spatio-Temporal Prior

$P(y|\mathbf{x})$

**Combine**

$\mathbf{x}$ = (longitude, latitude, day)

*Presence-Only Geographical Priors for Fine-Grained Image Classification,* Mac Aodha et al., 2019

# What to do about distribution shift?

- Distributionally robust optimization
- Learn (or use) a prior for subpopulation shift

Domain Adaptation

Source domain: ● ★ ▲ ■

Target domain: □ △ ○ ☆

https://link.springer.com/article/10.1007/s10489-022-03709-8

# What to do about distribution shift?

- Distributionally robust optimization
- Learn (or use) a prior for subpopulation shift
- Domain adaptation (next lecture!)

## Original Plates

## Acquiring images of plates with utensils



ImageNet

Bing

Stable Diffusion

ImageNet*



tray tray tray plate plate plate

tray bucket bucket plate plate plate

Figure 8: Real images of plates, with and without food and either on a table or in the grass. Below each image is the predicted class by an ImageNet-trained ResNet50.

**https://arxiv.org/abs/2302.07865**

# What to do about distribution shift?

- Distributionally robust optimization
- Learn (or use) a prior for subpopulation shift
- Domain adaptation (next lecture!)
- Diagnose failures

# What to do about distribution shift?

- Distributionally robust optimization
- Learn (or use) a prior for subpopulation shift
- Domain adaptation (next lecture!)
- Diagnose failures
- Get training data that is representative of your test domain (works better than any algorithm)

# Summary

- Out-of-distribution generalization: big challenge, but helps understand what NNs learn.

  - Adversarial examples and training

  - Distribution shifts