# Logistics

- Course outline

- HyFlex

- Zoom etiquette

- 3 problem sets & course project
  - grading comprises 50% psets + 50% project

- Office hours: Wed 4-5pm on Zoom

- TA: Chao Qin

# Overview at 10000 ft

- Logistics

- Stochastic optimization
  - Supervised learning as loss minimization
  - Stochastic gradient descent

- Recent advances in ML
  - Architectures with inductive bias
  - Progress in computer vision & NLP
  - Downstream applications

- Challenges
  - Distribution shifts
  - Adversarial examples
  - Fairness, accountability, transparency, and ethics
  - Spurious correlations

# Stochastic optimization

- Optimization under random data

- Loss/Objective $\ell(\theta; Z)$ where $\theta \in \Theta$ is parameter/decision to be learned, and $Z \sim P$ is random data

- Optimize average performance under P

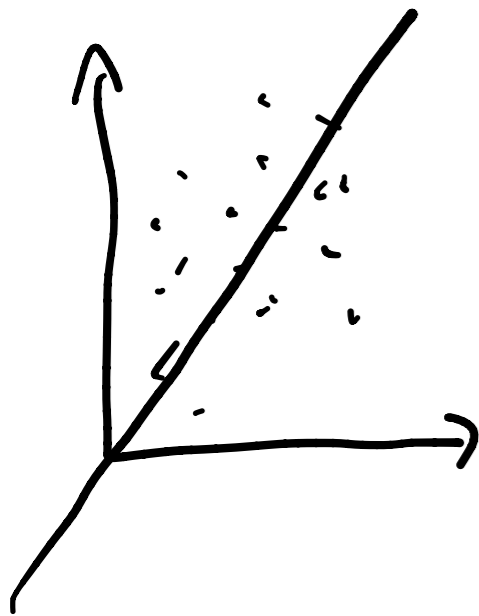$$\text{minimize}_{\theta \in \Theta} \ \mathbb{E}_P[\ell(\theta; Z)]$$

# Stochastic optimization

- For prediction problems, data often composes of Z = (X, Y), where X is features/covariates, and Y is label
  - e.g. X: image pixels, Y: cat/dog/sheep

- Loss min. abstraction includes almost all canonical supervised learning problems

- Foundational framework in OR, statistics, and ML

# Linear regression

$$Z = (X, Y) \qquad Y: \text{outcome} \qquad X: \text{covariate vector} \in \mathbb{R}^d$$

$$\ell(\theta; X, Y) = (Y - \theta^T X)^2$$



$$\text{If } \mathbb{E} X X^T > 0, \qquad \theta^* = \underset{\theta}{\text{argmin}} \, \mathbb{E} \ell(\theta; X, Y)$$

$$= (\mathbb{E} X X^T)^{-1} \mathbb{E} Y X.$$

$$\text{Robust regression}: \quad \ell(\theta; X, Y) = |Y - \theta^T X|.$$

# Maximum likelihood estimation

Likelihood model $\quad p_\theta(z)$

$$\min_{\theta \in \Theta} \quad -\mathbb{E} \log p_\theta(z)$$

Conditional likelihood model $\quad p_\theta(y|x)$

$$\min_{\theta \in \Theta} \quad -\mathbb{E} \log p_\theta(y|x)$$

# Binary classification

$$Z = (X, Y) \qquad Y \in \{-1, 1\} \qquad X: \text{features} \in \mathbb{R}^d$$

$$\{h_\theta(X) : \theta \in \Theta\} : \text{hypothesis class}$$
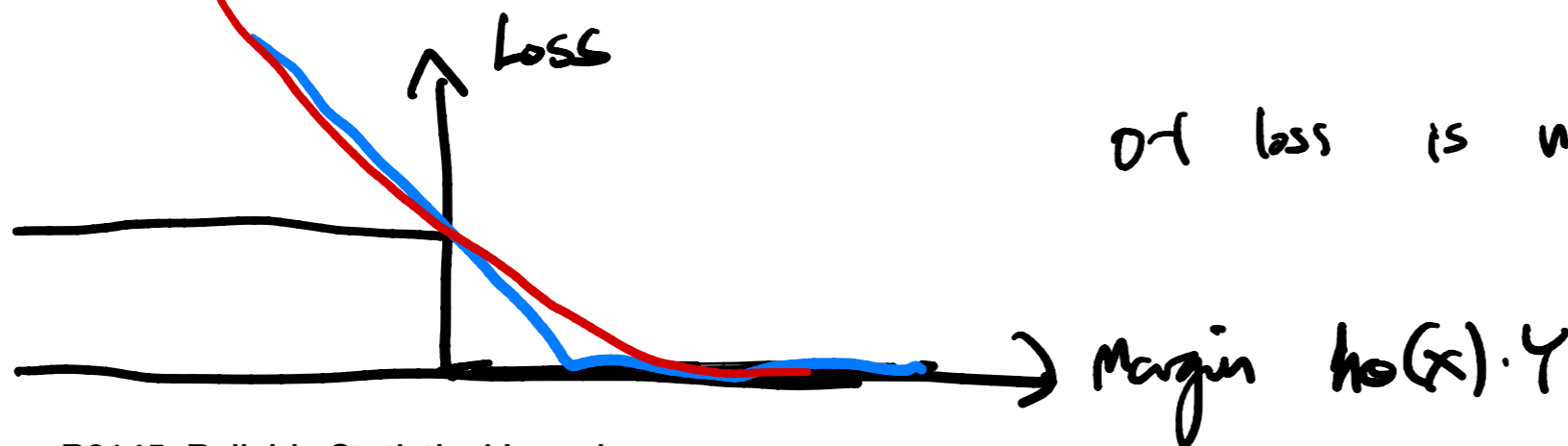
Predict $\quad \text{sgn}(h_\theta(X))$

$0-1$ loss: $\qquad \mathbb{1}\{\text{sgn}(h_\theta(x)) \neq Y\} = \mathbb{1}\{h_\theta(x) Y \leq 0\}$

Margin $\qquad h_\theta(x) Y \qquad$ "how right you are"



Loss

0-1 loss is non-smooth & non-convex

$\Rightarrow$ convex surrogates

Margin $h_\theta(x) \cdot Y$

# Binary classification

Hinge loss: $\ell(\theta; X, Y) = (1 - Y h_\theta(X))_+$

Support Vector Machines $\quad h_\theta(x) = \theta^T X$

$$\min_{\theta: \|\theta\|_2 \le r} \mathbb{E}(1 - Y\theta^T X)_+$$

Logistic loss: $\ell(\theta; X, Y) = \log(1 + \exp(-Y h_\theta(X)))$

Logistic Regression $\quad h_\theta(x) = \theta^T X$

$$\min_{\theta: \|\theta\|_p \le r} \mathbb{E}\log(1 + \exp(-Y h_\theta(x)))$$

$$\mathcal{H} = \{\theta: \|\theta\|_p \le r\}$$

# Binary classification

# Multi-class classification

$$\Theta = (\theta_1, \cdots, \theta_K) \in \mathbb{R}^{d \times k}$$

$$Y \in \{1, \cdots, K\}$$

$$\text{Logit} \quad P_\theta(y|x) = \frac{\exp(\theta_y^\top x)}{\sum_{k=1}^K \exp(\theta_k^\top x)}$$

$$\text{Max log likelihood} \equiv \min_{\theta \in \Theta} -\mathbb{E} \log P_\theta(Y|X)$$

$$= \min_{\theta \in \Theta} -\mathbb{E}\, \theta_Y^\top X + \mathbb{E} \log \sum_{k=1}^K \exp(\theta_k^\top X)$$

# Neural networks



Instead of $\Theta_k^\top X$, $\underline{h_{\Theta, u}(X)}$

# Neural networks $h_\theta(x) \in \mathbb{R}^k$

$\theta_1, \cdots, \theta_L$ : weight matrices

$\sigma_1, \cdots, \sigma_L$ : activations $\qquad \sigma_i : \mathbb{R}^{d_{i-1}} \to \mathbb{R}^{d_i} \qquad d_L = K$
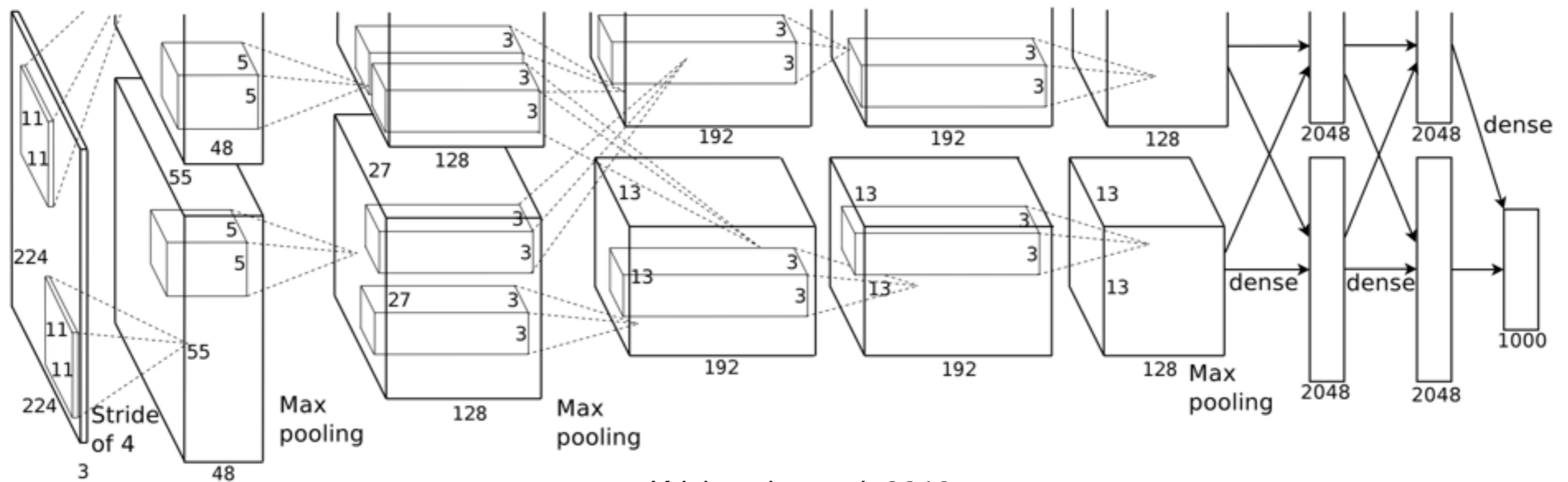$$d_0 = \dim(x)$$

$\qquad$ ReLU $\qquad \sigma(x)_j = \max(0, x_j)$

$\qquad$ Max pooling : Take a bunch local coordinates,
$$\text{output its maximum}$$

$$h_\theta(x) = \sigma_L \left( \theta_L \, \sigma_{L-1} \left( \theta_{L-1} \cdots \sigma_1(\theta_1 x) \cdots \right) \right)$$

Final loss : $\qquad \ell(\theta; x, y) = -\log \dfrac{\exp(h_{\theta, y}(x))}{\sum_{h=1}^{K} \exp(h_{\theta, h}(x))}$

# Convolutional nets



Krizhevsky et al. 2012

# Residual nets

He et al. 2015

# Newsvendor

Consider decision-maker deciding $\Theta \in \mathbb{R}_+$: order quantity

Uncertainty $Z$ : random demand

Order cost : $c$     If $Z > \Theta$, additional order cost $b \geq 0$

Holding cost : $h$

$$\ell(\Theta; Z) = c\Theta + b(Z-\Theta)_+ + h(\Theta-Z)_+$$

$(H) = \mathbb{R}_{\geq 0}$   or   $[0, M]$

# Portfolio optimization

$\theta \in \mathbb{R}_+^d$ : pfo weights

$Z$ : random asset returns $\in \mathbb{R}^d$

$\ell(\theta; z) = \theta^\top z$

$\Theta = \{\theta \in \mathbb{R}_+^d : \theta^\top \mathbb{1} = 1\}$

$$\min \quad -\mathbb{E}\,\theta^\top Z$$
$$\text{s.t.} \quad \theta \in \Theta$$

: risk-neutral

# Empirical risk minimization

- But we don't know P

- Even if we did, even evaluating the objective $\mathbb{E}_P[\ell(\theta; Z)]$ requires numerical integration over $Z \in \mathbb{R}^d$

  - d is often large in ML

- Empirical risk minimization (ERM), or sample average approximation (SAA) over $Z_i \stackrel{\text{iid}}{\sim} P$

$$\widehat{\theta}_n^{\text{erm}} = \text{argmin}_{\theta \in \Theta} \ \frac{1}{n} \sum_{i=1}^{n} \ell(\theta; Z_i)$$

# Optimization

$$\underset{\theta \in \Theta}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^{n} \ell(\theta; Z_i)$$

- How do we solve the ERM/SAA problem?
  - Let's say $\theta \mapsto \ell(\theta; Z)$ is convex
  - True for linear models [check for yourself!]

- Second-order methods (interior point methods)
  - Computing Hessian and doing backsolve is too expensive

- First-order methods
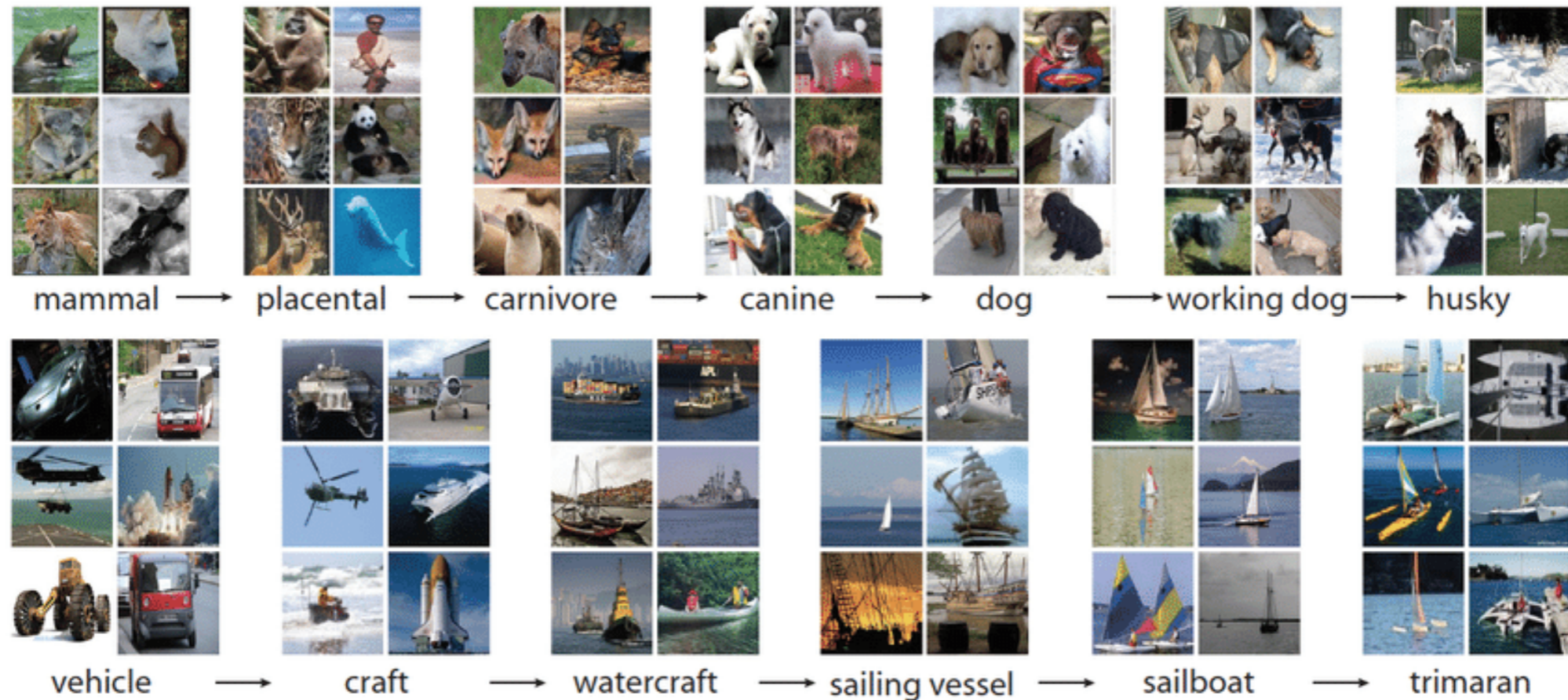  - Better, but still O(n) to even evaluate gradient

# Stochastic gradient descent

$$\underset{\theta \in \Theta}{\text{minimize}} \ \frac{1}{n} \sum_{i=1}^{n} \ell(\theta; Z_i)$$

$$\theta^{t+1} \leftarrow \theta^t - \alpha_t \nabla_\theta \ell(\theta^t; Z_t)$$

# Magic formula

- Inductive bias: CNN, ResNet, RNN, LSTM, attention, transformers

- Big datasets

- Optimize some surrogate loss using SGD

- GPUs

# Big datasets: ImageNet



mammal → placental → carnivore → canine → dog → working dog → husky

vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

- 2012 classification challenge: 1.3M images, 1000 labels

- Collected through web search, verified via Mechanical Turk

- Hierarchy of labels

# Big datasets: ImageNet

**SUN, 131K**
[Xiao et al. '10]

**LabelMe, 37K**
[Russell et al. '07]

**PASCAL VOC, 30K**
[Everingham et al. '06-'12]

**Caltech101, 9K**
[Fei-Fei, Fergus, Perona, '03]

IMAGENET 15M
[Deng et al. '09]

**Slide from Fei-Fei Li**

# Hammers



Slide from Jia Deng

# Ladles



Slide from Jia Deng

# ImageNet competition



100% wrong

75

**In the competition's first year**
teams had varying success.
Every team got at least 25%
wrong.

**In 2012,** the team to first use
deep learning was the only
team to get their error rate
below 25%.

50

**The following year**
nearly every team got
25% or fewer wrong.

25

**In 2017,** 29 of 38
teams got less than
5% wrong.

perfect

'10    '11    '12    '13    '14    '15    '16    '17

# Top-5 error



152 layers

22 layers

19 layers

8 layers

8 layers

shallow

3.57 — ILSVRC'15 ResNet
6.7 — ILSVRC'14 GoogleNet
7.3 — ILSVRC'14 VGG
11.7 — ILSVRC'13
16.4 — ILSVRC'12 AlexNet
25.8 — ILSVRC'11
28.2 — ILSVRC'10

Figure from Siddharth Das

# Success in vision



Redmon & Farhadi (2016), YOLO

# Success in vision

**https://www.youtube.com/watch?v=HS1wV9NMLr8&ab_channel=NVIDIA**

**https://www.youtube.com/watch?v=868tExoVdQw&ab_channel=Zoox**

# Engineering excellence

- ImageNet in X minutes, using $Y etc

  - https://dawn.cs.stanford.edu/benchmark/#imagenet

- Better pipelines, stable deployment

- Edge devices, run real-time on AV

# Success in NLP

- ## Machine translation
  - In 2014, first sequence-to-sequence paper
  - In 2016, Google translate switched to this technology

- ## Language models

Slide from Chris Manning's NLP class CS224N

**Now**

| ULMfit | GPT | BERT | GPT-2 | XLNet | GPT-3 |
|--------|-----|------|-------|-------|-------|
| Jan 2018 | June 2018 | Oct 2018 | Feb 2019 | June 2019 | May 2020 |
| Training: | Training | Training | Training | Training | |
| 1 GPU day | 240 GPU days | 256 TPU days | ~2048 TPU v3 days according to a reddit thread | 2816 TPU v3 days | 175 billion param $12M to train |
| | | ~320–560 GPU days | | | |

…

# GPT-3

**https://twitter.com/sharifshameem/status/1282676454690451457**

# Applications

- Fraud detection

- Robot-assisted surgical assistance

- Automated diagnosis, radiology assistants

- Fault detection in manufacturing systems

- Autonomous vehicles

- List goes on

# Obligatory remark

- Deep learning excitement/hype comes from ability to handle complex unstructured data that was previously impossible

- NOT a panacea for every problem

- Linear regression is a reasonable first step in most practical problems

- Random forests and gradient boosting are almost always good enough (and easier to train, test, deploy, and maintain)

- Collecting enough labels and building the entire pipeline for deep learning is a HUGE effort

# Break

# Progress in machine learning?

## Human-level average performance

### Image recognition [Eckersley+ '17]



### Face recognition [Harris+ '15]



TECH • GOOGLE

Google: Our new system for recognizing faces is

By DERRICK HARRIS

## Poor performance on underrepresented examples

TECH • GOOGLE

Google: Our new system for recognizing faces is

TECH • GOOGLE

e: Our new system for recognizing faces is
the best one ever

:ing tool that

⋮ REUTERS

HARRIS March 17, 2015

*Facial Recognition Is Accurate,
if You're a White Guy*

By Steve Lohr

Feb. 9, 2018

The New York Times

# Average-case

$$\text{minimize}_{\theta \in \Theta} \; \mathbb{E}_P[\ell(\theta; Z)]$$

- Only optimize performance under data-generating distribution P

- But data collection is always biased, and distributional shifts are ubiquitous (e.g. spatial, temporal)

- Only optimize average performance under P

  – No consideration for tail-performance

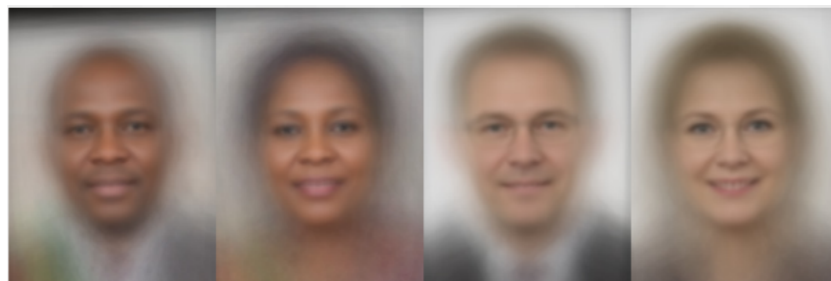# Facial recognition

- Labeled Faces in the Wild, a gold standard dataset for face recognition, is **77.5% male**, and **83.5% White** [Han and Jain '14]

- Commercial gender classification softwares have **disparate** performance on different subpopulations

| Gender Classifier | Darker Male | Darker Female | Lighter Male | Lighter Female | Largest Gap |
|---|---|---|---|---|---|
| Microsoft | 94.0% | 79.2% | 100% | 98.3% | 20.8% |
| FACE++ | 99.3% | 65.5% | 99.2% | 94.0% | 33.8% |
| IBM | 88.0% | 65.3% | 99.7% | 92.9% | 34.4% |

Gendered Shades: Intersectional accuracy disparity
[Buolamwini and Gebru '18]

# Lack of diversity in data

- **"Clinical trials for new drugs skew heavily white"**
  - Less than 5% of cancer trial participants were non-white

[Oh et al. '15, Burchard et al. '15, Chen et al., '14, SA Editors '18]

- Majority of image data from **US & Western Europe**

**ImageNet: country of origin**



| | |
|---|---|
| CC | 0.0% |
| IE | 0.5% |
| AR | 1.0% |
| ES | 2.5% |
| AU | 2.8% |
| CA | 3.0% |
| IT | 6.2% |

US 45.4%

GB 7.6%

[Shankar et al. '17]

**Other examples**



Dependency parsing

[Blodgett+ 16]



Captioning

[Tatman+ 17]



Recommender systems

[Ekstrand+ 17,18]



Face recognition

[Grother+ 11]



Language identification

[Blodgett+ 16, Jurgens +17]



Part-of-speech tagging

[Hovy+ 15]

# Lack of diversity in data



[DeVries et al. 2019, Does object recognition work for everyone?]

Slide from Timnit Gebru & Emily Denton's CVPR2020 tutorial

# Gender bias in machine translation

**Alex Shams**
@seyyedreza
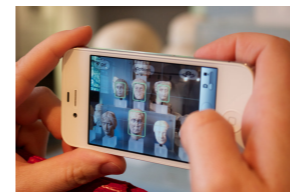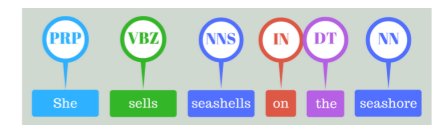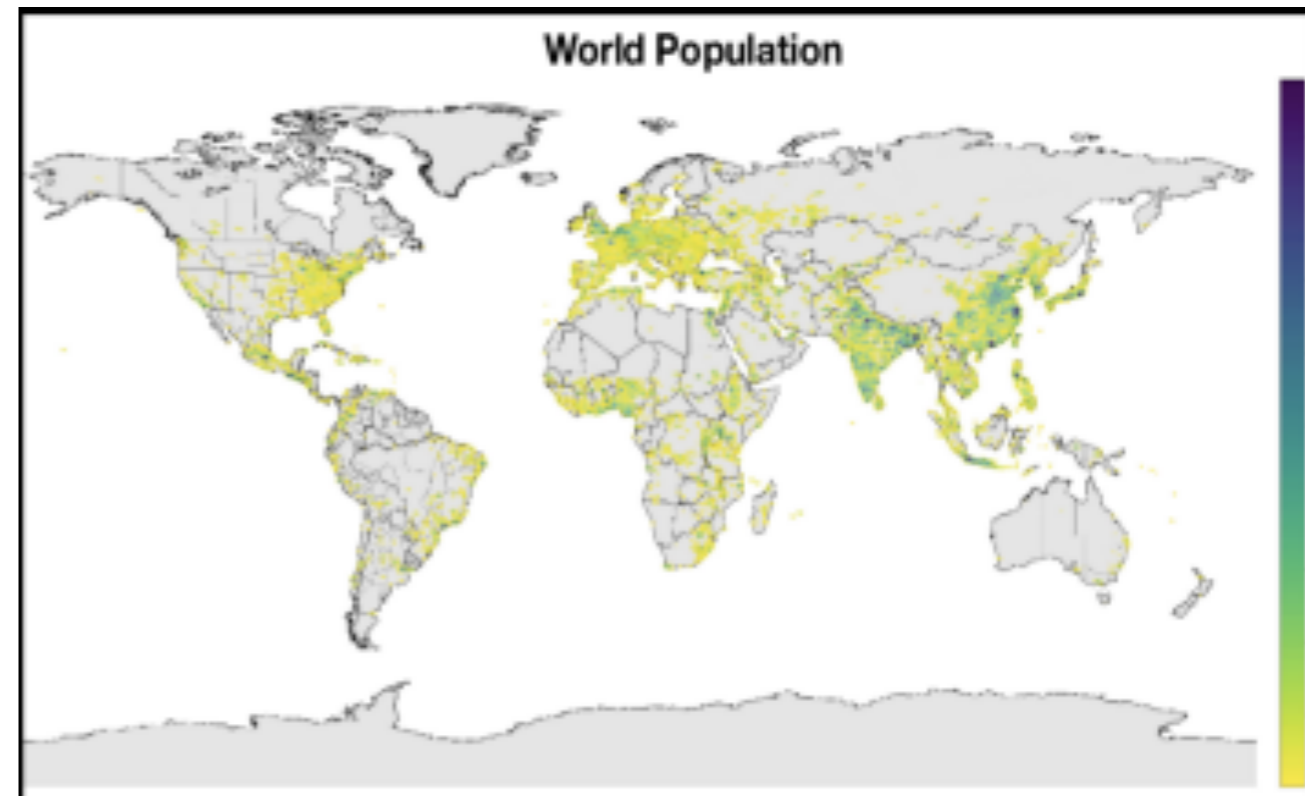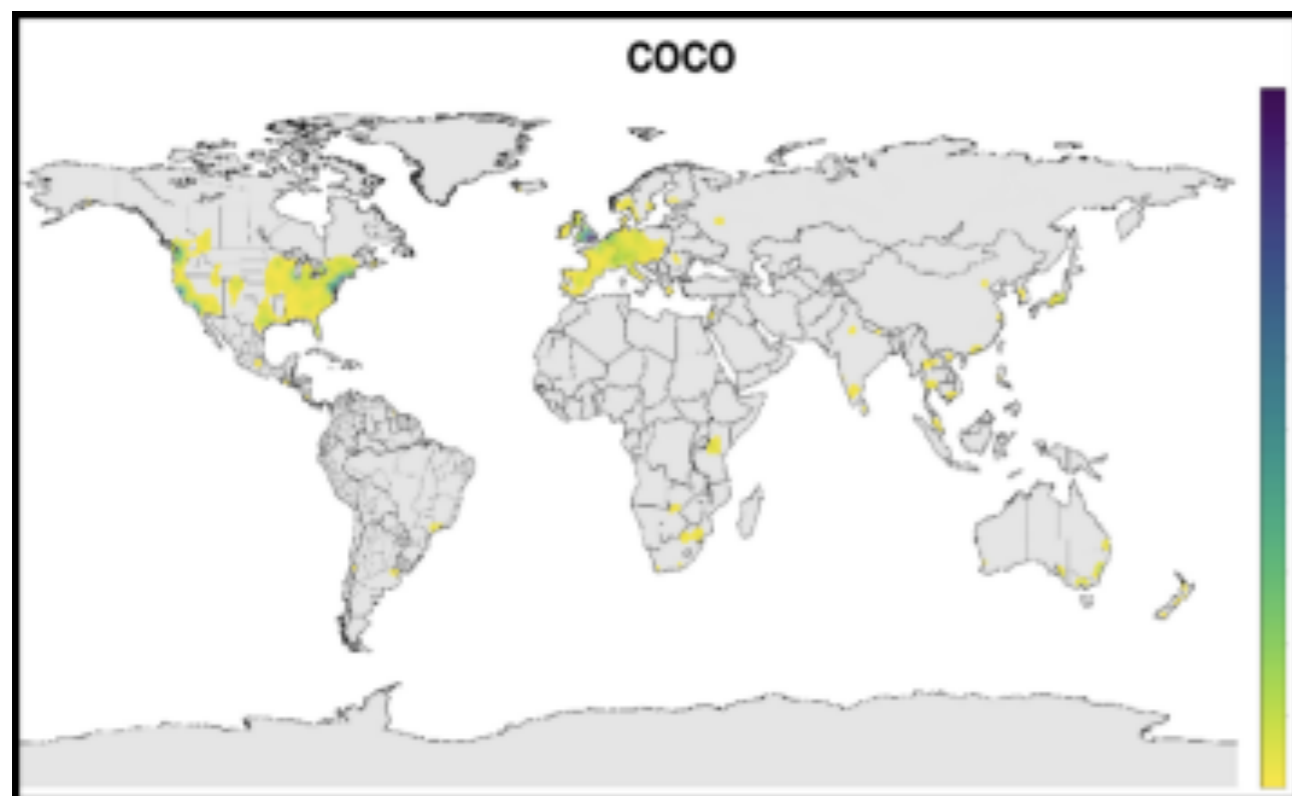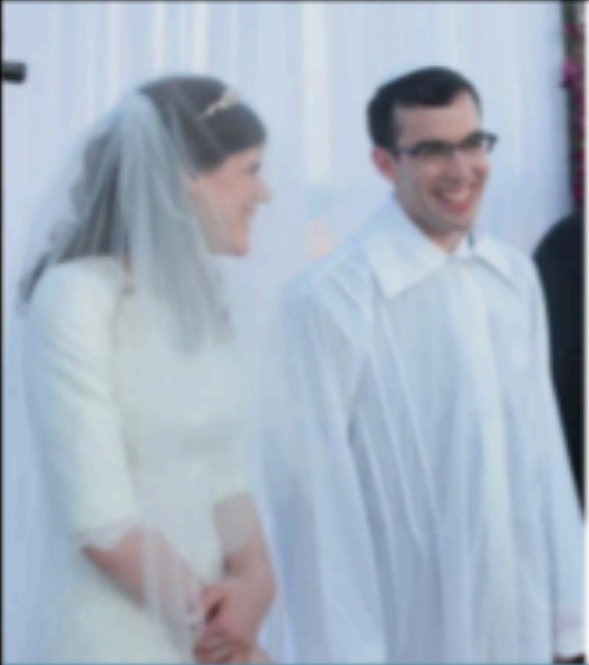
Turkish is a gender neutral language. There is no "he" or "she" – everything is just "o". But look what happens when Google translates to English. Thread:

| Turkish - detected ▾ | English ▾ |
|---|---|
| o bir aşçı | she is a cook |
| o bir mühendis | he is an engineer |
| o bir doktor | he is a doctor |
| o bir hemşire | she is a nurse |
| o bir temizlikçi | he is a cleaner |
| o bir polis | He-she is a police |
| o bir asker | he is a soldier |
| o bir öğretmen | She's a teacher |
| o bir sekreter | he is a secretary |
| | |
| o bir arkadaş | he is a friend |
| o bir sevgili | she is a lover |

| | |
|---|---|
| onu sevmiyor | she does not like her |
| onu seviyor | she loves him |
| | |
| onu görüyor | she sees it |
| onu göremiyor | he can not see him |
| | |
| o onu kucaklıyor | she is embracing her |
| o onu kucaklamıyor | he does not embrace it |
| | |
| o evli | she is married |
| o bekar | he is single |
| | |
| o mutlu | he's happy |
| o mutsuz | she is unhappy |
| | |
| o çalışkan | he is hard working |
| o tembel | she is lazy |

6:36 PM · Nov 27, 2017 · Twitter Web Client

**14.9K** Retweets   **2K** Quote Tweets   **27.2K** Likes

41

# Racial bias in speech recognition

**MARCH 23, 2020**

## Stanford researchers find that automated speech recognition is more likely to misinterpret black speakers

*The disparity likely occurs because such technologies are based on machine learning systems that rely heavily on databases of English as spoken by white Americans.*

**BY EDMUND L. ANDREWS**

The technology that powers the nation's leading automated speech recognition systems makes twice as many errors when interpreting words spoken by African Americans as when interpreting the same words spoken by whites, according to a new study by researchers at Stanford Engineering.

Slide from Timnit Gebru & Emily Denton's CVPR2020 tutorial

Thanks to machine-learning algorithms, the robot apocalypse was short-lived.

# Long-tails

- Long-tailed data is ubiquitous in modern applications

  - Google (7 yrs ago): constant fraction of queries were new each day

- Tail inputs often determine quality of service

**Long-tailed queries**

# Fundamentally hard examples

- Task: classify image of dog to breed (120 classes)
- Kernel features



Stanford Dogs Dataset [Khosla et al. '11]

No underrepresentation:
same number of images per class

# Big gaps in performance



**BIG gap in performance even when no underrepresentation**

top-5 error rate vs classes

**Hard** ⟷ **Easy**

B9145: Reliable S
Hongseok

# Not a new problem...

- Standard regressors obtained from MLE lose predictive power on certain regions of covariates [Meinshausen & Buhlmann (2015)]

- Temporal, spatial shifts common



Demographic shift over space and time

# Not a new problem…

## Classifier Technology and the Illusion of Progress

David J. Hand

- "A fundamental assumption of the classical paradigm is that the various distributions involved do not change over time. In fact, in many applications this is unrealistic and the population distributions are nonstationary."

  - Marketing & banking: Classification rules used to predict loan default updated every few months

  - "Their performance degrades, not because the rules themselves change, but because the distributions to which they are being applied change"

# Not a new problem...

- Model performance drops across different domains and datasets [Torralba & Efros (2011)]



Table 1. Cross-dataset generalization. Object detection and classification performance (AP) for "car" and "person" when training on one dataset (rows) and testing on another (columns), i.e. each row is: training on one dataset and testing on all the others. "Self" refers to training and testing on the same dataset (same as diagonal), and "Mean Others" refers to averaging performance on all except self.

| task | Test on: / Train on: | SUN09 | LabelMe | PASCAL | ImageNet | Caltech101 | MSRC | Self | Mean others | Percent drop |
|---|---|---|---|---|---|---|---|---|---|---|
| "car" classification | SUN09 | **28.2** | 29.5 | 16.3 | 14.6 | 16.9 | 21.9 | 28.2 | 19.8 | **30%** |
| | LabelMe | 14.7 | **34.0** | 16.7 | 22.9 | 43.6 | 24.5 | 34.0 | 24.5 | **28%** |
| | PASCAL | 10.1 | 25.5 | **35.2** | 43.9 | 44.2 | 39.4 | 35.2 | 32.6 | **7%** |
| | ImageNet | 11.4 | 29.6 | 36.0 | **57.4** | 52.3 | 42.7 | 57.4 | 34.4 | **40%** |
| | Caltech101 | 7.5 | 31.1 | 19.5 | 33.1 | **96.9** | 42.1 | 96.9 | 26.7 | **73%** |
| | MSRC | 9.3 | 27.0 | 24.9 | 32.6 | 40.3 | **68.4** | 68.4 | 26.8 | **61%** |
| | Mean others | 10.6 | 28.5 | 22.7 | 29.4 | 39.4 | 34.1 | 53.4 | 27.5 | 48% |

# SOTA models are also non-robust



ImageNet

New test accuracy (top-1, %) vs Original test accuracy (top-1, %)

- - - Ideal reproducibility ● Model accuracy ── Linear fit

[Does ImageNet classifiers generalize to ImageNet?
Recht, Roelofs, Schmidt, Shankar '19]

# SOTA models are non-robust

## Similar frames extracted from videos
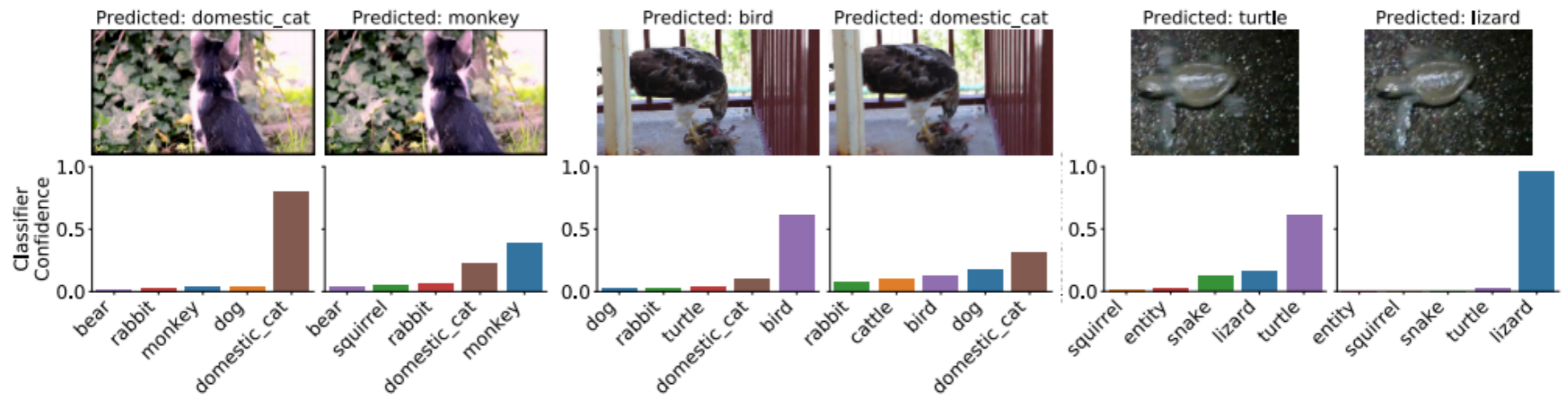


Figure 1: Three examples of natural perturbations from nearby video frames and resulting classifier predictions from a ResNet-152 model fine-tuned on ImageNet-Vid. While the images appear almost identical to the human eye, the classifier confidence changes substantially.

[Does ImageNet classifiers generalize across time?
Shankar, Dave, Roelofs, Ramanan, Recht, Schmidt '19]

# SOTA models are non-robust



Figure 3: Model accuracy on original vs. perturbed images. Each data point corresponds to one model in our testbed (shown with 95% Clopper-Pearson confidence intervals). Each perturbed frame was taken from a ten frame neighborhood of the original frame (approximately 0.3 seconds). All frames were reviewed by humans to confirm visual similarity to the original frames.

[Does ImageNet classifiers generalize across time?
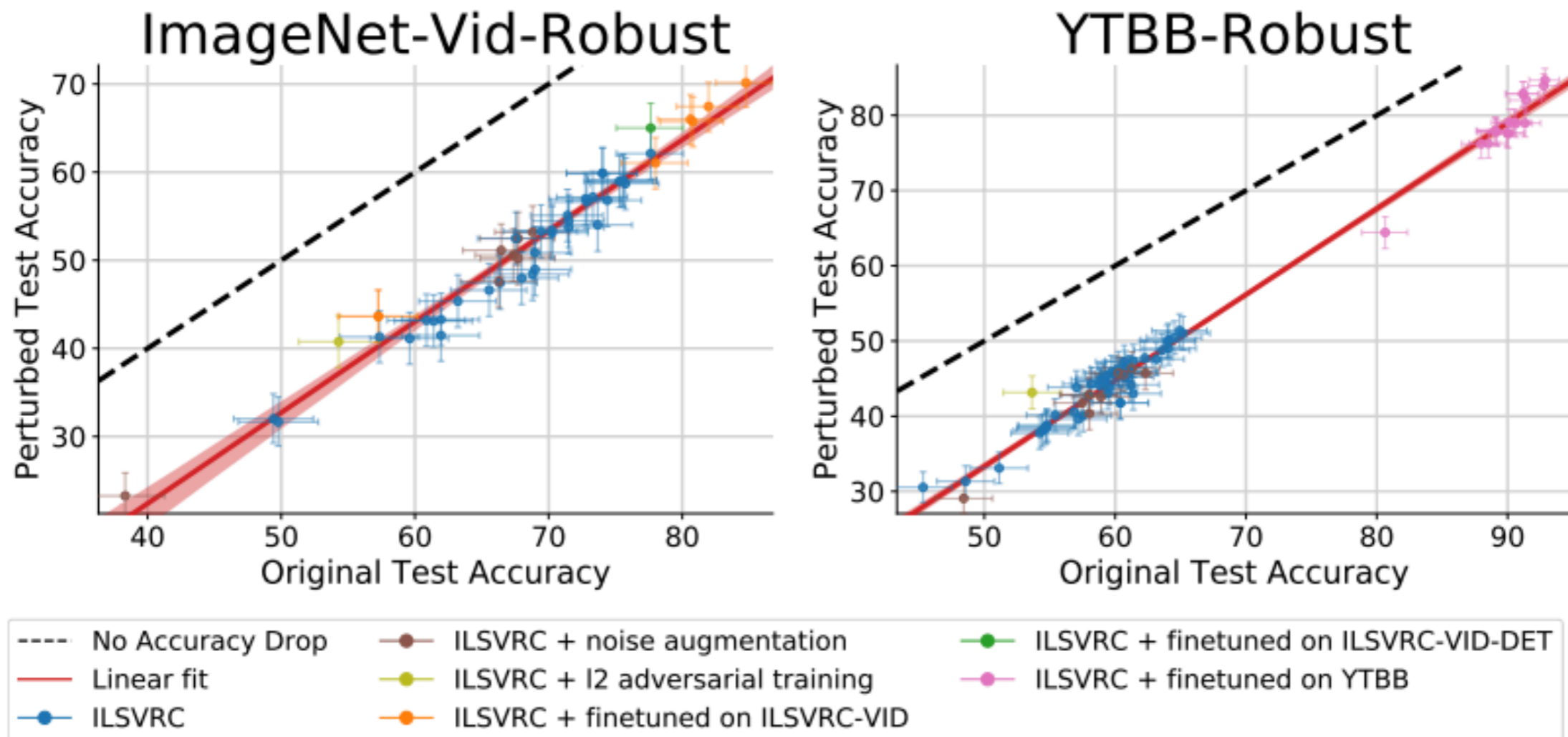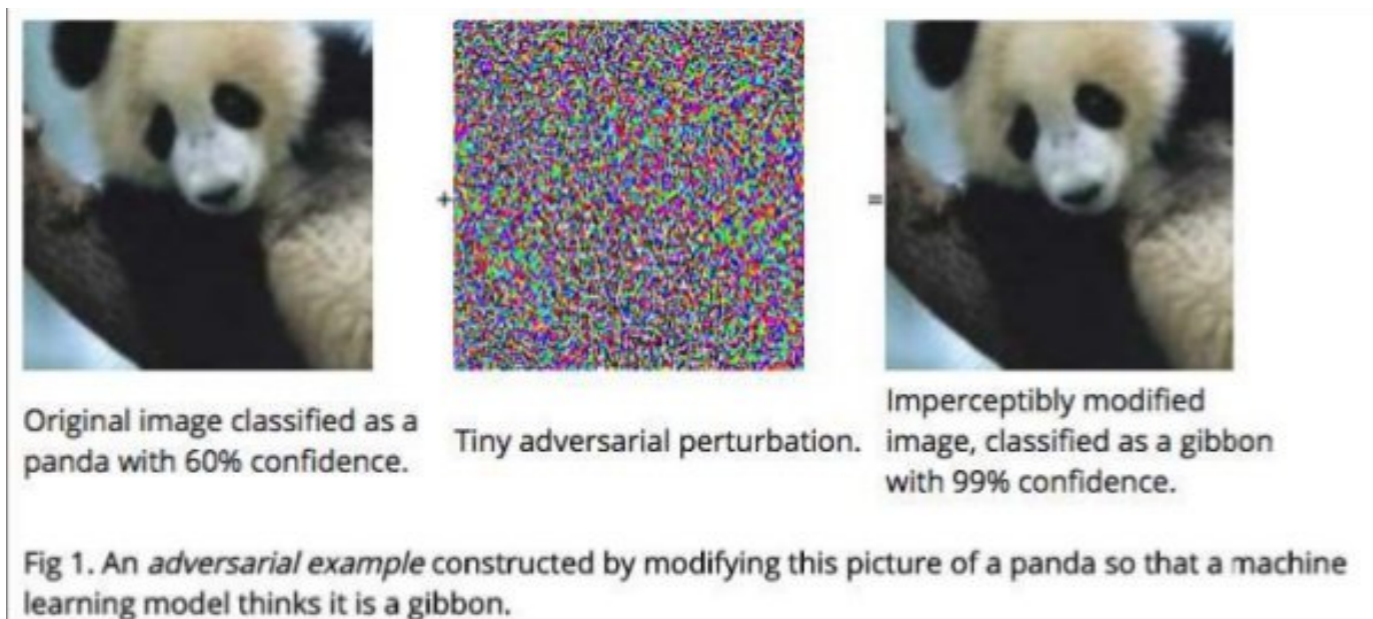Shankar, Dave, Roelofs, Ramanan, Recht, Schmidt '19]

# SOTA models are non-robust

- Deep networks are very brittle
  - imperceptible adversarial perturbations can fool them



Original image classified as a panda with 60% confidence. Tiny adversarial perturbation. Imperceptibly modified image, classified as a gibbon with 99% confidence.

Fig 1. An *adversarial example* constructed by modifying this picture of a panda so that a machine learning model thinks it is a gibbon.

Goodfellow et al. (2015)

88% tabby cat → adversarial perturbation → 99% guacamole

Nicholas Carlini

# SOTA models are non-robust

- Deep networks are very brittle
  - imperceptible adversarial perturbations can fool them



[Athalye et al. '17]

[Chen et al. '18]

# Spurious correlations

- Models fit to observed associations, which maybe not be the fundamental structure that we want to learn



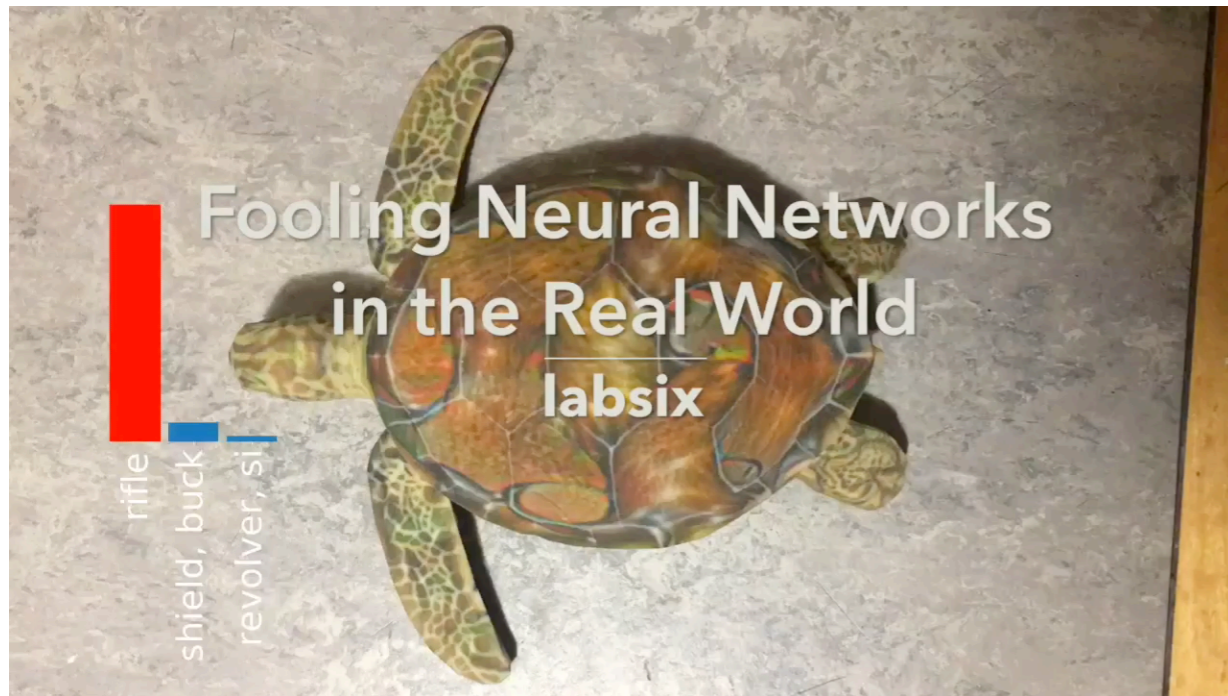Figure 1: Representative training and test examples for the datasets we consider. The correlation between the label $y$ and the spurious attribute $a$ at training time does not hold at test time.

Sagawa et al. (2019)

Accuracy:
100% train
=>
60% test

- But I want my models to work in a non-patriarchal society without sexism

Amazon scraps secret AI recruiting tool that showed bias against women ⬡ REUTERS

# Environmental concerns
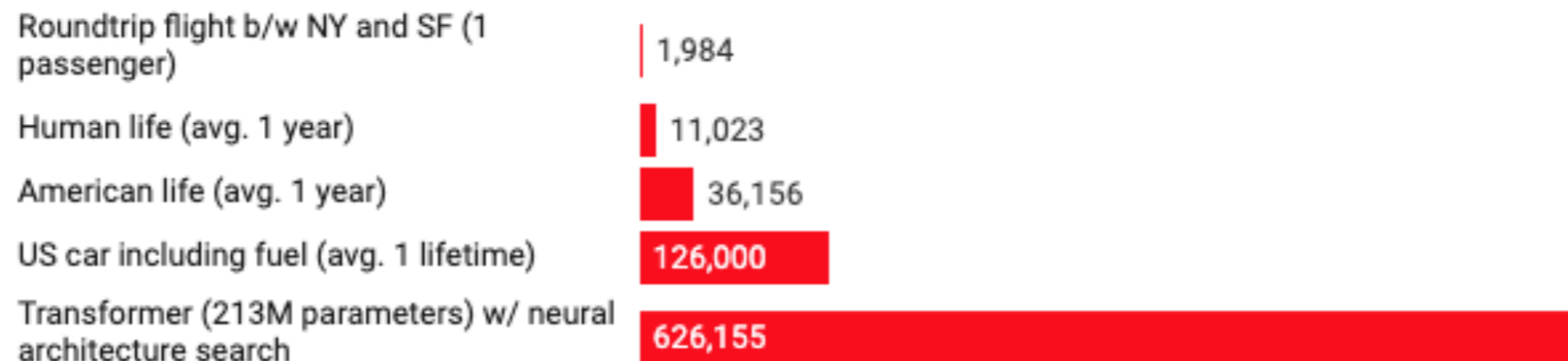
## Common carbon footprint benchmarks

in lbs of CO2 equivalent

| | |
|---|---|
| Roundtrip flight b/w NY and SF (1 passenger) | 1,984 |
| Human life (avg. 1 year) | 11,023 |
| American life (avg. 1 year) | 36,156 |
| US car including fuel (avg. 1 lifetime) | 126,000 |
| Transformer (213M parameters) w/ neural architecture search | 626,155 |

Chart: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper

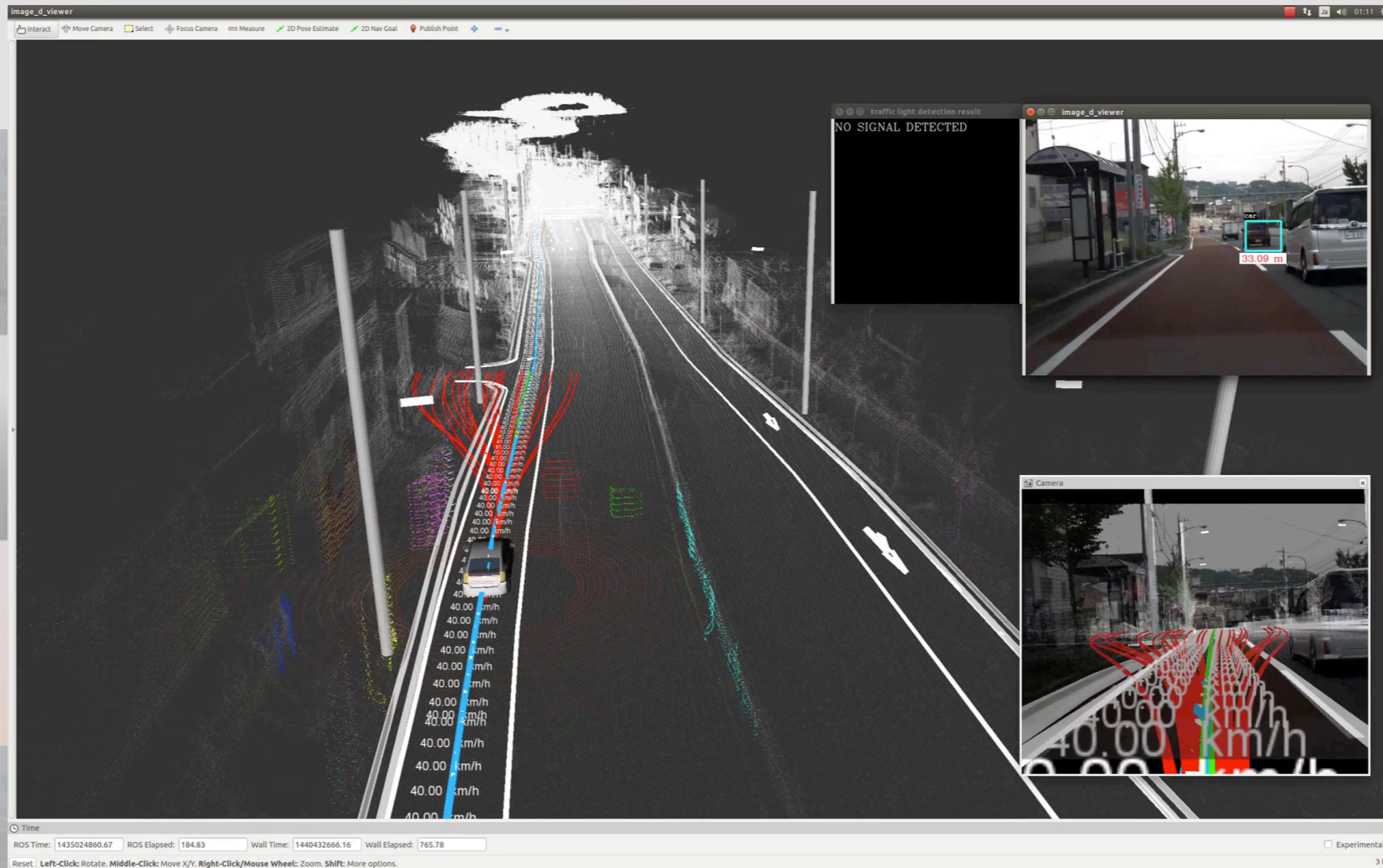| Consumption | $CO_2e$ (lbs) |
|---|---|
| Air travel, 1 passenger, NY↔SF | 1984 |
| Human life, avg, 1 year | 11,023 |
| American life, avg, 1 year | 36,156 |
| Car, avg incl. fuel, 1 lifetime | 126,000 |

| Training one model (GPU) | |
|---|---|
| NLP pipeline (parsing, SRL) | 39 |
| w/ tuning & experimentation | 78,468 |
| Transformer (big) | 192 |
| w/ neural architecture search | 626,155 |

# Lots of questions

- How do we deal with unanticipated distributional shifts?

  - Modeling itself is nontrivial

- How do we learn causal structures?

- Ultimately, ML models work towards aiding downstream decisions

  - Prediction is not the ultimate goal

  - How to design models with this in mind?

- How do we evaluate the entire system, with many complex modules?
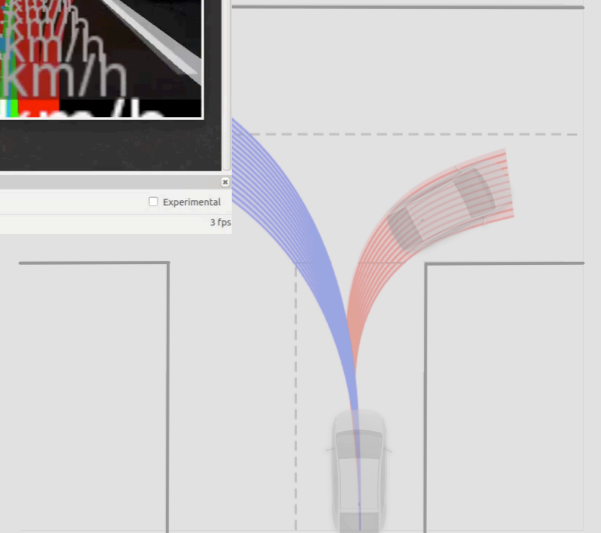
# Complex system example: AVs

*Sense*



At the end of the day:

A function that *generates* a sequence of *steering* and *acceleration* commands

# Complex system example: AVs



Mobileye running a red light



Tesla Autopilot fatal accident

# Lots of questions

- ML system interacts with (strategic) agents over time. How to model this interaction/dynamics?

- All modern platforms employ ML as a part of their pipeline

- Operational constraints (safety, reliability etc)

- Collected data on decisions are observational

  - Often based on human agents' decisions, which may depend on unrecorded variables

  - For sequential decisions, observed data often does not cover entire (action seq, state seq) space. So not really "big data"…

# Rest of the course

- First, learn foundational techniques!

  - One month on basic results in statistical learning, and how to prove them

- Then, survey recent works that aim to identify, model, improve upon aforementioned challenges

  - Focus is on *principled* methods, but we'll also discuss a range of practical issues

- Goal: Develop a critical view of topics surrounding reliability

  - Much remains to be done in ML

  - Discussions toward context-specific applications (e.g. healthcare, manufacturing, supply chains, finance, marketing…)

- Goal: Identify interfaces

  - mechanism design

  - sequential decision-making

  - simulation (e.g. rare-events)