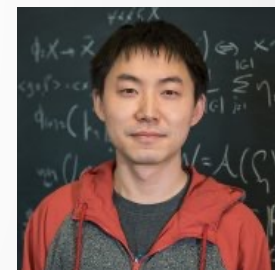


# Chasing the Long Tail: What Neural Networks Memorize and Why

Vitaly Feldman  
Apple\*



Chiyuan Zhang  
Google Research



*\*Part of the work done while at Google Research  
and while visiting the Simons Institute, UC Berkeley*

# Label memorization

For state-of-the-art deep learning algorithms on image datasets

- Training set error: 0-5%
- (Typical) test set error: 10-30%

Inception on ImageNet 1000 **with random labels**:

Training error: **9%**

Test error: **99.9%**

[Zhang,Bengio,Hardt,Recht,Vinyals '17]

**Not explained by existing theories**

# Defining memorization

**Def:** For dataset  $S = ((x_1, y_1), \dots, (x_n, y_n))$ ,  $i \in [n]$  and learning algorithm  $A$  let

$$\text{mem}(A, S, i) = \Pr_{\hat{h}=A(S)} [\hat{h}(x_i) = y_i] - \Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x_i) = y_i]$$

For data distribution  $P$  over  $X \times Y$

$$\mathbf{E}_{S \sim P^n} \left[ \frac{1}{n} \sum_{i \in [n]} \text{mem}(A, S, i) \right] \approx \mathbf{E}_{S \sim P^n, \hat{h}=A(S)} \left[ \underbrace{\text{err}_P(\hat{h}) - \text{err}_S(\hat{h})}_{\text{generalization gap}} \right]$$

generalization error      empirical error

Why is label memorization ubiquitous? Is it necessary?

# Privacy concerns

Black-box membership inference attack with high accuracy

[Shokri,Stronati,Song,Shmatikov 17; LongBWBWTGC 18; SalemZFHB 18,...]

Reconstruction of sensitive training data from language models

[Carlini et al. 20]



Is memorization necessary for achieving optimal error?

# Related work

## Generalization of interpolating methods

- Interpolation may be “harmless”
- Does not address the “**why?**” question
- Memorization also happens without interpolation

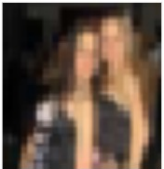
[Cover,Hart '67; Belkin,Hsu,Mitra '18;  
Liang,Rakhlin '18; Belkin,Rakhlin,Tsybakov '19;  
Belkin,Hsu,Xu '19; Bartlett,Long,Lugosi,Tsigler '19;  
Hastie,Montanari,Rosset,Tibshirani '19;  
Muthukumar,Vodrahalli,Sahai '19; Mei,Montanari 19;  
Montanari,Zhong 20; ....]

# Hard atypical examples

CIFAR-10 **truck** class

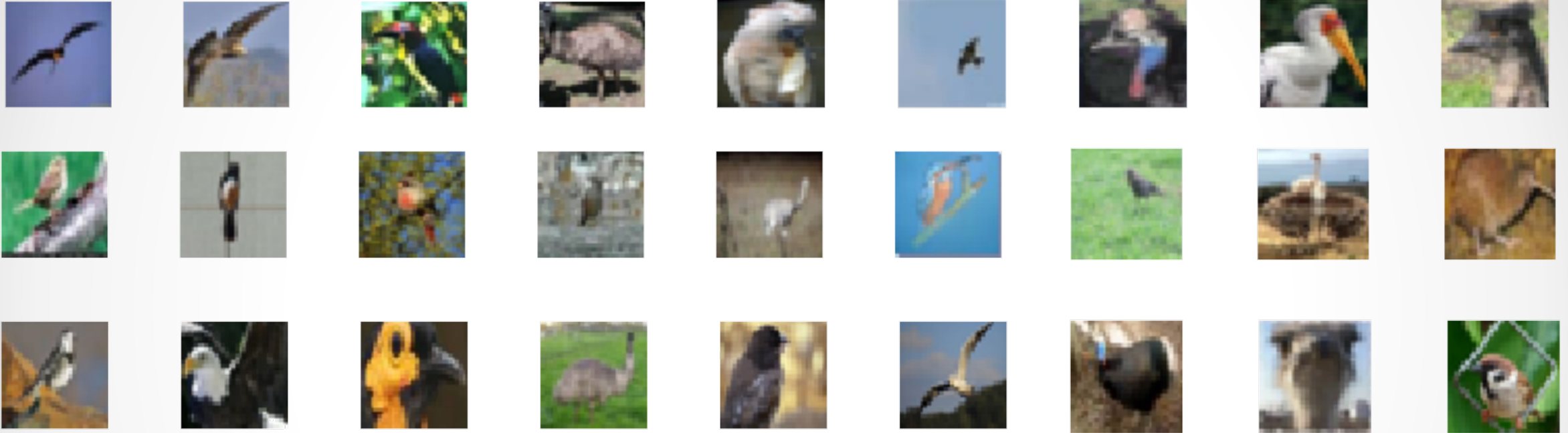
training set

test set



# Subpopulations

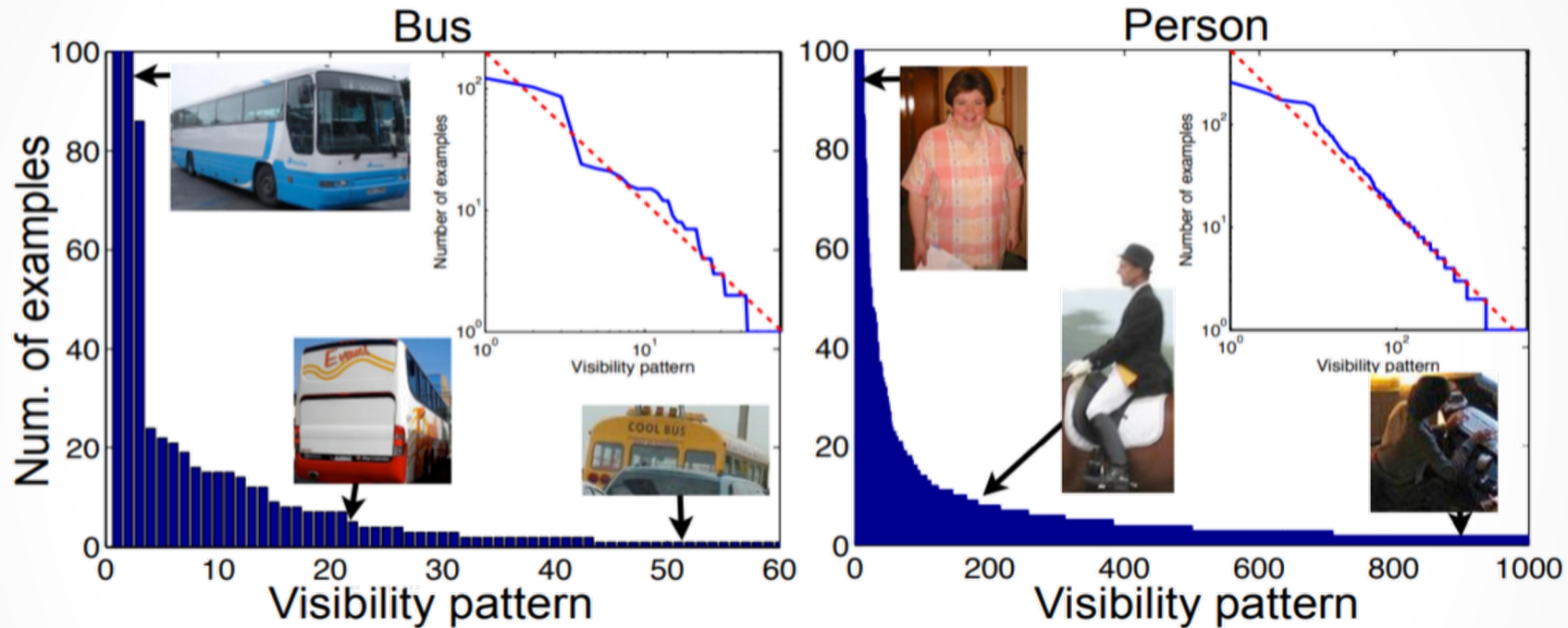
CIFAR-10 birds



Examples from low ( $\approx \frac{1}{n}$ ) frequency subpopulations may be hard to distinguish from mislabeled ones and outliers

# Long-tailed data

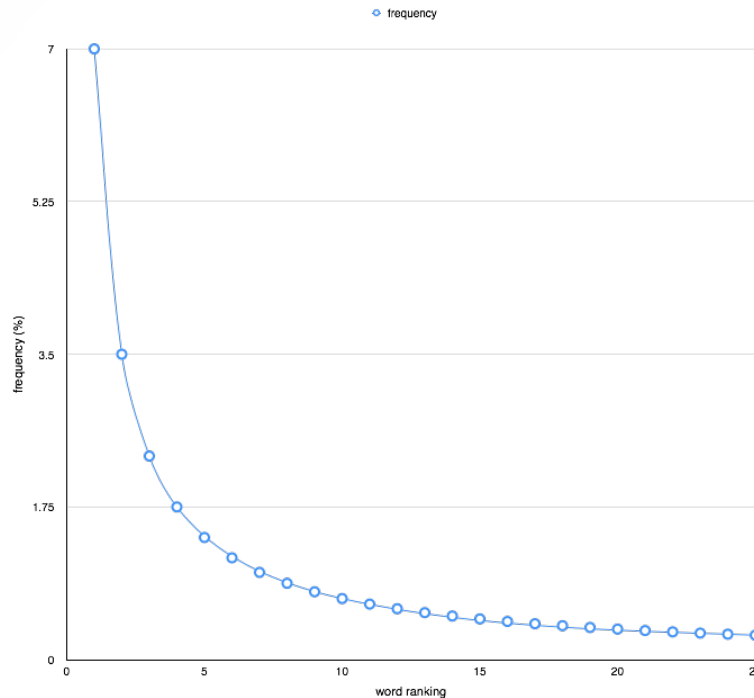
Additional annotations



[Zhu, Anguelov, Ramanan '14]



# The tail rears its head



For  $n \leq t$

$$\text{weight}_{\pi_t} \left[ \frac{1}{2n}, \frac{2}{n} \right] = \Omega \left( \frac{1}{\log t} \right)$$

Zipf distribution

$$\Pr[\alpha = k] \propto 1/k$$

Long-tail:  $\Theta \left( \frac{1}{n} \right)$  frequencies can make up a **significant fraction** of the whole population

# Basic model

- Discrete domain  $X$  of size  $t$  and set of labels  $Y$
- Learning problem: meta-distribution  $\Phi$  over data distributions on  $X \times Y$
- Goal: minimize expected generalization error of algorithm  $A$

$$\overline{\text{err}}_n(\Phi, A) := \mathbf{E}_{P \sim \Phi, S \sim P^n, \hat{h} = A(S)} [\text{err}_P(\hat{h})]$$

- Prior over frequencies  $\pi = (\pi_1, \dots, \pi_t)$
- For each  $x \in X$ , pick randomly and independently a frequency  $p_x$  from  $\pi$  and normalize

$$D(x) = \frac{p_x}{\sum_{x' \in X} p_{x'}}$$

- Arbitrary distribution  $F$  over functions  $X \rightarrow Y$
- $\Phi(\pi, F)$ : pick  $D \sim \pi^t$  and  $f \sim F$ . Define  $P$  as  $(x, f(x))$  for  $x \sim D$

# Benefits of fitting

**Thm:** For every algorithm  $A$ :

$$\mathbf{E}_{P \sim \Phi(\pi, F), S \sim P^n, \hat{h} = A(S)} [\text{err}_P(\hat{h})] \geq \text{opt}_n(\pi, F) + \tau \cdot \mathbf{E}_{P \sim \Phi(\pi, F), S \sim P^n, \hat{h} = A(S)} [\text{err}_S(\hat{h}, 1)]$$

$$\min_A \overline{\text{err}}_n(\Phi(\pi, F), A)$$

$$\tau := \frac{\mathbf{E}_{\alpha \sim \pi^*} [\alpha^2 (1 - \alpha)^{n-1}]}{\mathbf{E}_{\alpha \sim \pi^*} [\alpha (1 - \alpha)^{n-1}]}$$

Number of examples occurring once in  $S$  but not fit by  $\hat{h}$

For all  $\pi, n$ ,  $\tau = \Omega\left(\frac{\text{weight}_\pi\left[\frac{1}{2n}, \frac{2}{n}\right]}{n}\right)$

For Zipf and  $n \leq t$ ,  $\tau = \Omega\left(\frac{1}{n \log t}\right)$

$F$  can be arbitrarily complex  
Extends to label noise

# Fitting and memorization

**Def:** For dataset  $S = ((x_1, y_1), \dots, (x_n, y_n))$ ,  $i \in [n]$  and learning algorithm  $A$  let

$$\text{mem}(A, S, i) = \Pr_{\hat{h}=A(S)} [\hat{h}(x_i) = y_i] - \Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x_i) = y_i]$$

For sufficiently complex  $F$ : for many  $i \in [n]$

$$\Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x_i) = y_i] \ll 1$$

Thus need to be memorized to fit

# Empirical validation

**Def:** For dataset  $S = ((x_1, y_1), \dots, (x_n, y_n))$ ,  $i \in [n]$  and learning algorithm  $A$  let

$$\text{mem}(A, S, i) = \Pr_{\hat{h}=A(S)} [\hat{h}(x_i) = y_i] - \Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x_i) = y_i]$$

Requires training  $\Omega(n)$  models. Computationally infeasible!

**Def:** For  $m \leq n$ , let  $T$  be a random subset of  $S$  of size  $m$  that includes  $i$

$$\text{mem}_m(A, S, i) = \mathbf{E}_T \left[ \Pr_{\hat{h}=A(T)} [\hat{h}(x_i) = y_i] - \Pr_{\hat{h}=A(T \setminus i)} [\hat{h}(x_i) = y_i] \right]$$

Closely-related and requires training  $O(1)$  models! Use  $m = 0.7n$

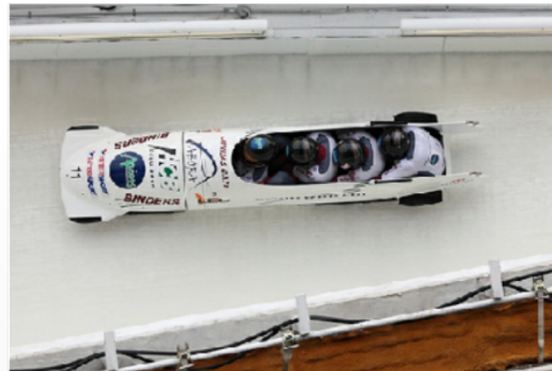
# Examples of memorization estimates

High memorization estimate



ImageNet: **bobsled** class

Low memorization estimate



# Influence of memorized examples

**Def:** For dataset  $S = ((x_1, y_1), \dots, (x_n, y_n))$ ,  $i \in [n]$  and testing example  $(x, y)$

$$\text{infl}(A, S, i, (x, y)) = \Pr_{\hat{h}=A(S)} [\hat{h}(x) = y] - \Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x) = y]$$

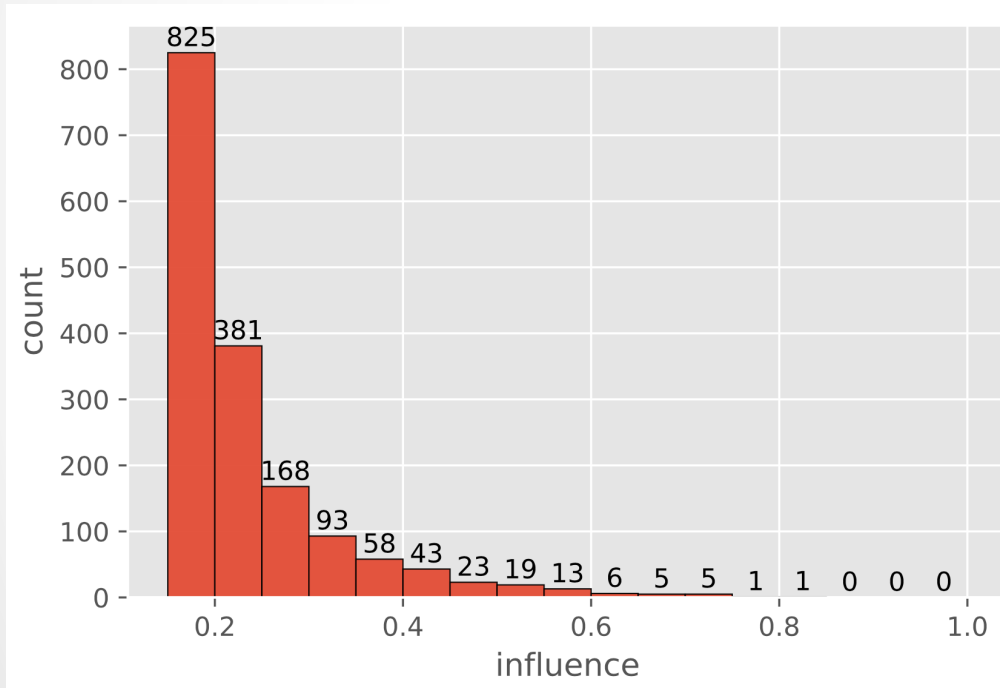
Same trick!

**Def:** For  $m \leq n$ , let  $T$  be a random subset of  $S$  of size  $m$  that includes  $i$

$$\text{infl}_m(A, S, i, (x, y)) = \mathbf{E}_T \left[ \Pr_{\hat{h}=A(T)} [\hat{h}(x) = y] - \Pr_{\hat{h}=A(T \setminus i)} [\hat{h}(x) = y] \right]$$

# High-influence pairs

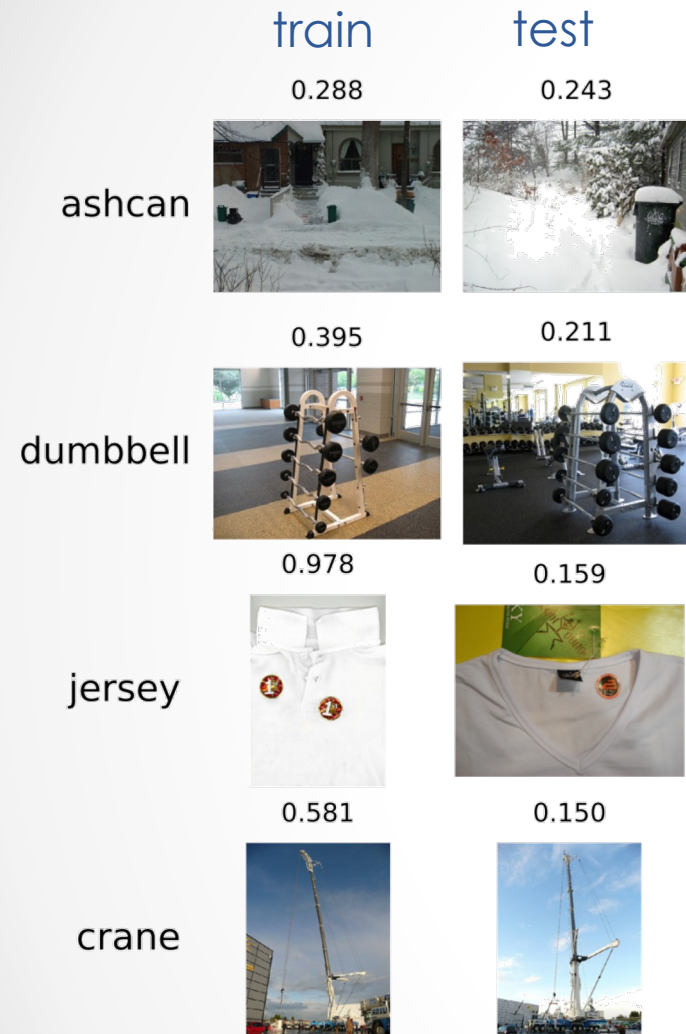
Many memorized training examples each of which significantly influences the accuracy on just one test example



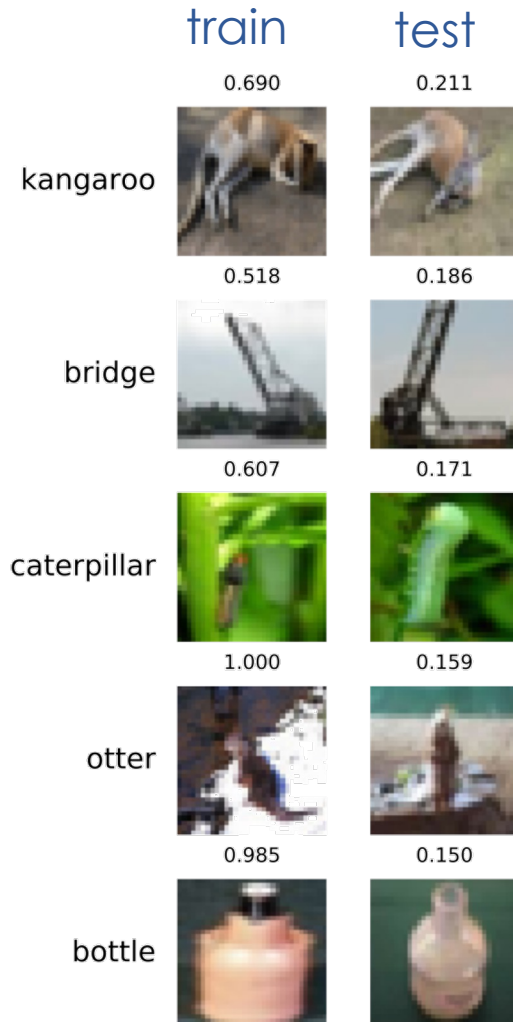
Affects  $\approx 3\%$  of the ImageNet test set



# Examples of high-influence



ImageNet



CIFAR-100

# Implications: differential privacy

**Def:** For dataset  $S = ((x_1, y_1), \dots, (x_n, y_n))$ ,  $i \in [n]$  and learning algorithm  $A$  let

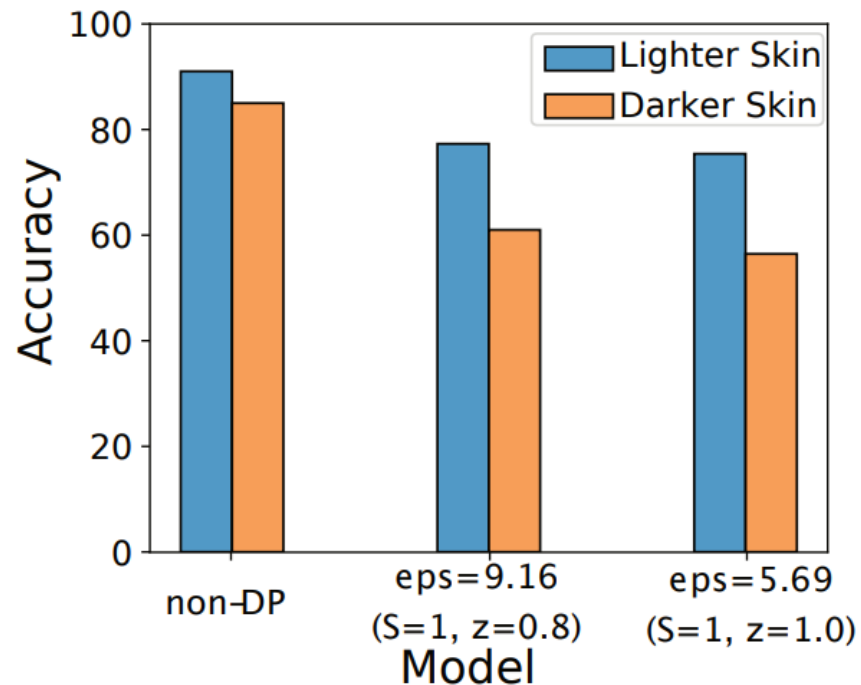
$$\text{mem}(A, S, i) = \Pr_{\hat{h}=A(S)} [\hat{h}(x_i) = y_i] - \Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x_i) = y_i]$$

For an  $(\epsilon, \delta)$ -DP algorithm  $A$ :

$$\Pr_{\hat{h}=A(S)} [\hat{h}(x_i) = y_i] \leq e^\epsilon \cdot \Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x_i) = y_i] + \delta$$

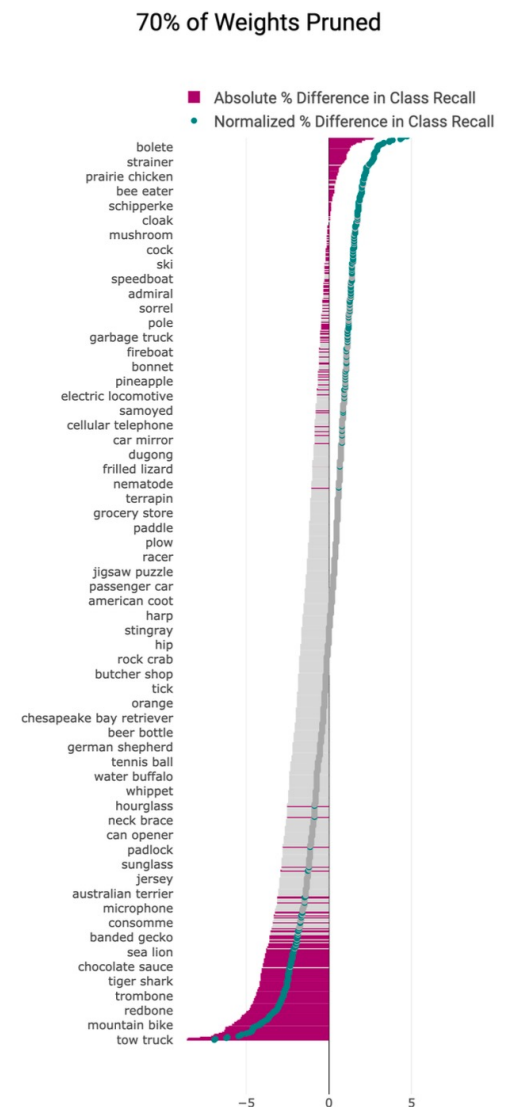
$$\text{mem}(A, S, i) \leq (e^\epsilon - 1) \cdot \Pr_{\hat{h}=A(S \setminus i)} [\hat{h}(x_i) = y_i] + \delta$$

# Implications: disparate effects



[Bagdarasyan, Shmatikov 19]

Cost of non-memorization depends on subgroup size, subpopulation structure, relative hardness of the problem



[Hooker, Courville, Clark, Dauphine, Frome 19]

# Conclusions

- Label memorization is necessary for optimal generalization on long-tailed data distributions
  - Not algorithm-specific
  - Can be modelled theoretically
  - Implications of limiting memorization: effect on the accuracy depends on subgroup frequency and other properties

## Additional materials

- **Theory:** STOC 2020, <https://arxiv.org/abs/1906.05271>
- **Experiments:** NeurIPS 2020, <https://arxiv.org/abs/2008.03703>
- **Results**/additional examples: <https://pluskid.github.io/influence-memorization/>
- **Follow-up:** memorization of data points
  - with Gavin Brown, Mark Bun, Adam Smith and Kunal Talwar
  - STOC 2021: <https://arxiv.org/abs/2012.06421>

## Future work

- Faithful models of learning
- Applications of estimators