

# Training Private Models That Know What They Don't Know

---

**Stephan Rabanser**<sup>1,2</sup> Anvith Thudi<sup>1,2</sup> Abhradeep Thakurta<sup>3</sup>  
Krishnamurthy (Dj) Dvijotham<sup>3</sup> Nicolas Papernot<sup>1,2</sup>

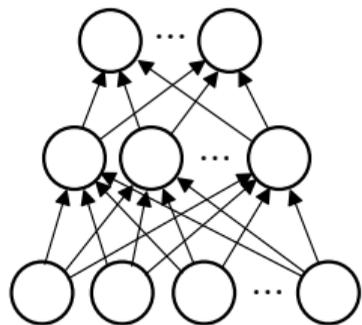
<sup>1</sup>University of Toronto <sup>2</sup>Vector Institute <sup>3</sup>Google DeepMind

Correspondence to: [stephan@cs.toronto.edu](mailto:stephan@cs.toronto.edu)

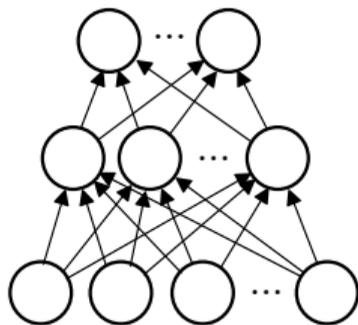


# Motivation: Input Sample Rejection

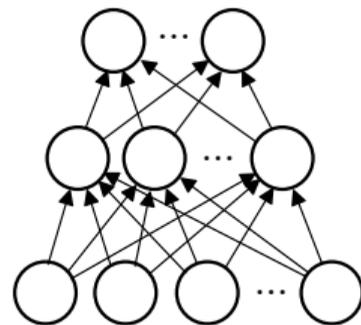
5



2



7? 1? 9?

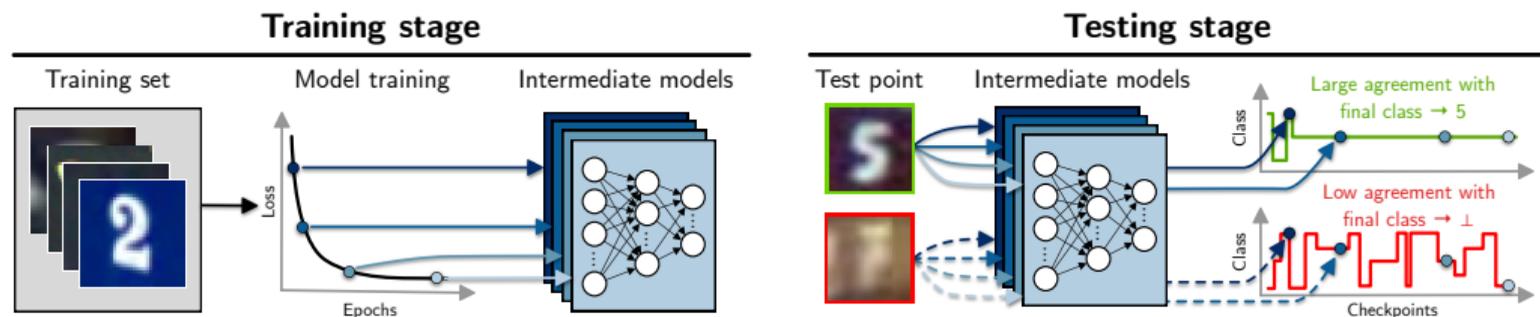


# Selective Classification (SC) with Training Dynamics

**Goal:** Derive a selection function  $g : \mathcal{X} \rightarrow \mathbb{R}$  which, given an acceptance threshold  $\tau$ , determines whether a model  $f : \mathcal{X} \rightarrow \mathcal{Y}$  should predict on a data point  $\mathbf{x}$ .

$$(f, g)(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & g(\mathbf{x}) \leq \tau \\ \perp & \text{otherwise.} \end{cases}$$

New approach: **Selective Classification Training Dynamics (SCTD)**



Rabanser et al. "Selective Classification via Neural Network Training Dynamics." 2022.

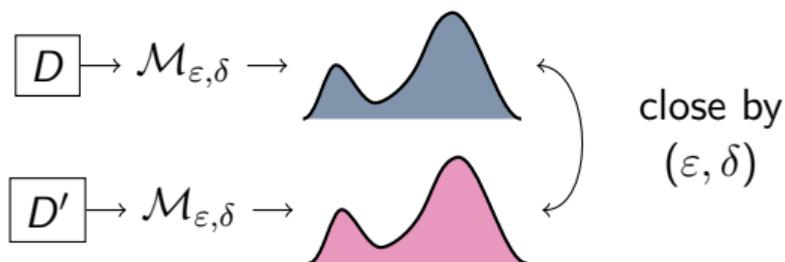
# Differential Privacy

## Definition: Differential Privacy

A randomized algorithm  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$  differential privacy, if for any two datasets  $D, D' \subseteq \mathcal{D}$  that differ in any one record and any set of outputs  $S$  the following inequality holds:

$$\mathbb{P}[\mathcal{M}(D) \in S] \leq e^\epsilon \mathbb{P}[\mathcal{M}(D') \in S] + \delta$$

- $\epsilon \in \mathbb{R}_+$  specifies the privacy level.
- $\delta \in [0, 1]$  allows for a small violation of the bound.
- Most popular implementation for DP in DNNs is DP-SGD.



# Impacts of SC on DP Guarantees: "SC $\rightarrow$ DP"

## No changes in DP guarantees.

### Direct Optimization

Loss function / architecture modifications.

### Post-Processing

Post-hoc modifications and training-time ensembles (SCTD).

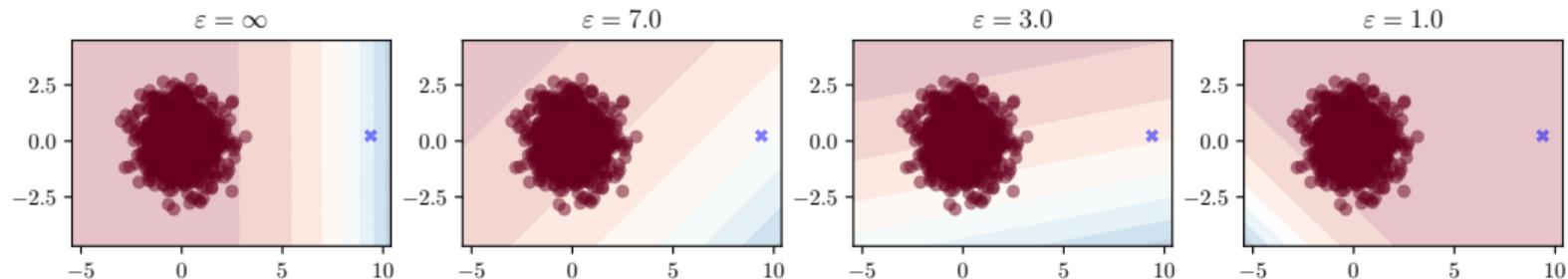
## Worsened DP guarantee.

### Advanced Sequential Composition

Methods iterating over the data multiple times, i.e., classical ensembling methods.

To maintain overall  $(\epsilon, \delta)$ -DP, each model needs to satisfy  $\approx (\frac{\epsilon}{\sqrt{M}}, \frac{\delta}{M})$ -DP.

# Impacts of DP on SC Performance: "DP $\rightarrow$ SC"



Differential privacy degrades selective classification performance beyond a loss in overall utility!

# Consistent SC Evaluations Under DP

## Compare performance across SC approaches

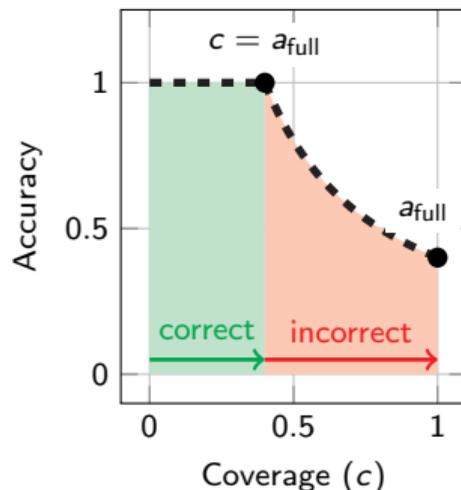
1. align SC methods to the same accuracy.
2. evaluate AUC metric for all SC methods.

**But under DP:**  
accuracy  $\overset{\approx}{\rightleftharpoons} \epsilon$ .

Training for less leads to expending less privacy budget.

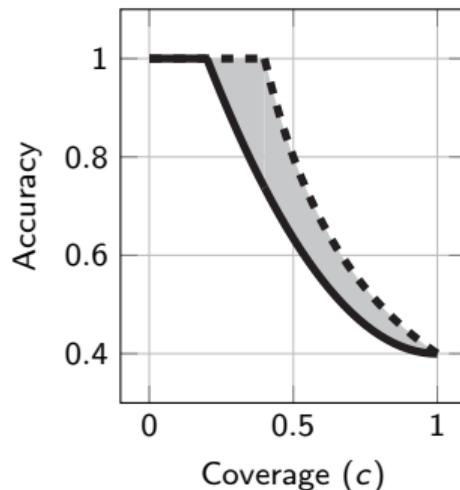
## Upper bound on SC performance

$$\overline{\text{acc}}(a_{\text{full}}, c) = \begin{cases} 1 & 0 < c \leq a_{\text{full}} \\ \frac{a_{\text{full}}}{c} & a_{\text{full}} < c < 1 \end{cases}$$

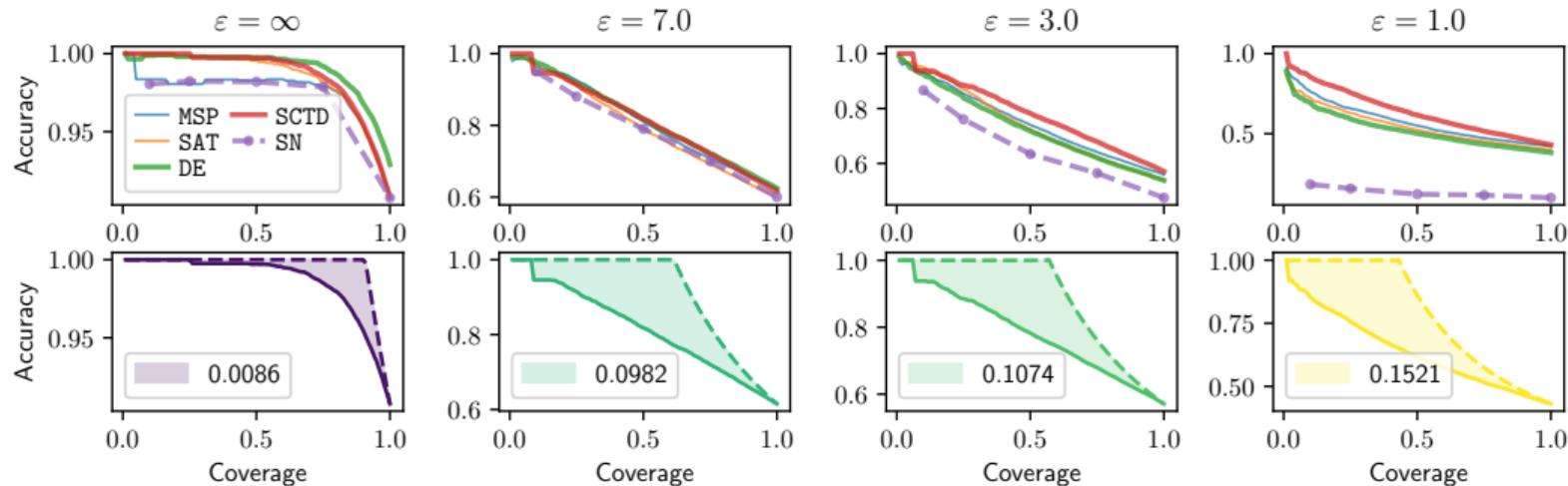


## Our accuracy-normalized SC score

$$s(f, g) = \int_0^1 (\overline{\text{acc}}(a_{\text{full}}, c) - \text{acc}_c(f, g)) dc$$



# Results on CIFAR-10



# Training DP Models That Know What They Don't Know: Conclusion

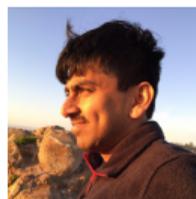
- Analyzed how SC impacts DP guarantees and how DP impacts SC performance.
- Introduced a novel score to disentangle SC performance from baseline utility.
- SC performance degrades with stronger privacy (i.e. as  $\epsilon \rightarrow 0$ ).
- SCTD works best to quantify uncertainty under DP.



Stephan



Anvith



Abhradeep



Dj



Nicolas

**Check out our paper:** <https://openreview.net/forum?id=EgCjf1vjMB>