

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

---

Stephan Rabanser

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



September 21, 2023

# Warmup: (Supervised) Machine Learning Pipeline

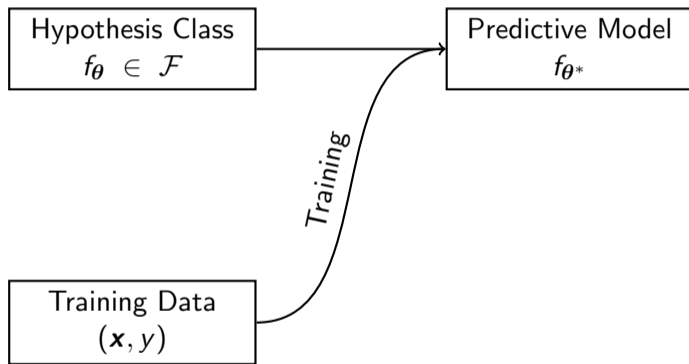
Hypothesis Class

$$f_{\theta} \in \mathcal{F}$$

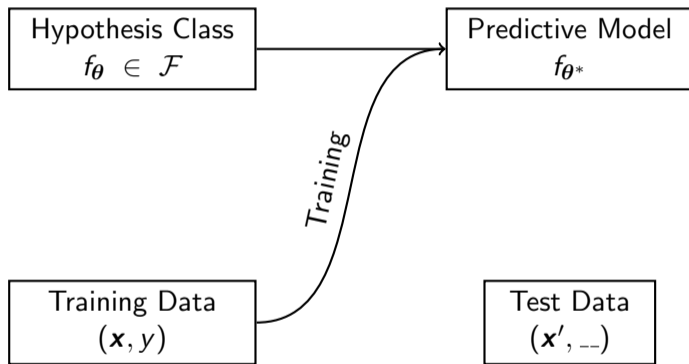
Training Data

$$(\mathbf{x}, y)$$

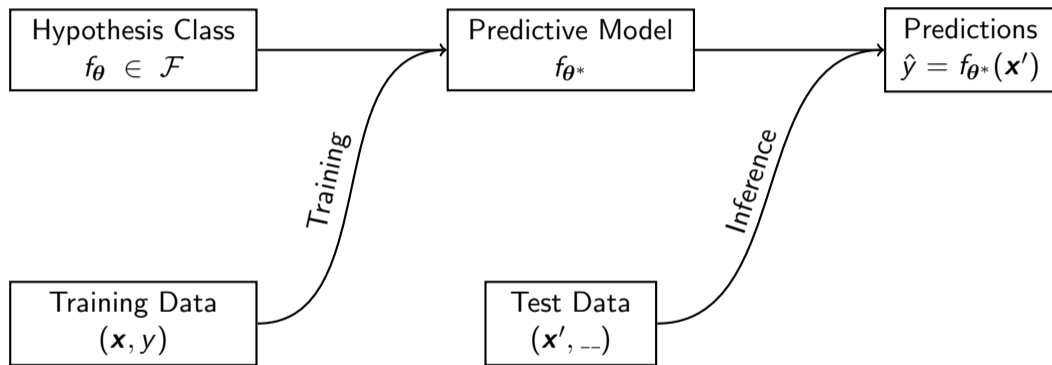
# Warmup: (Supervised) Machine Learning Pipeline



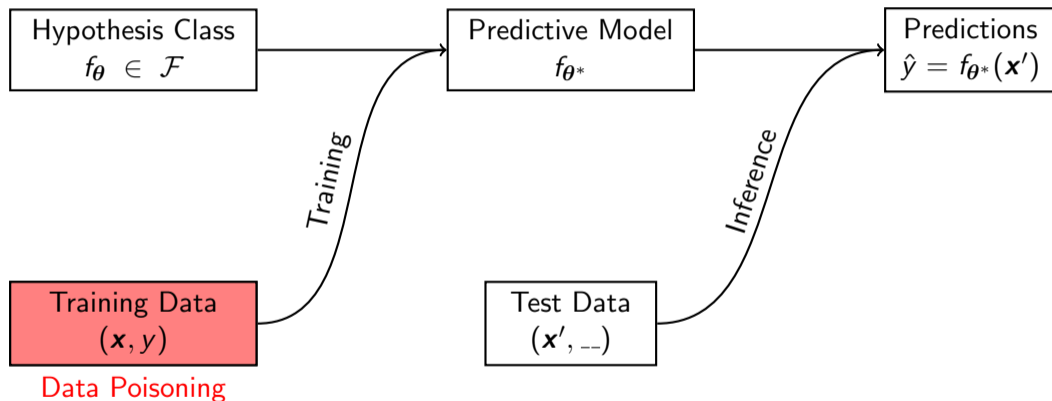
## Warmup: (Supervised) Machine Learning Pipeline



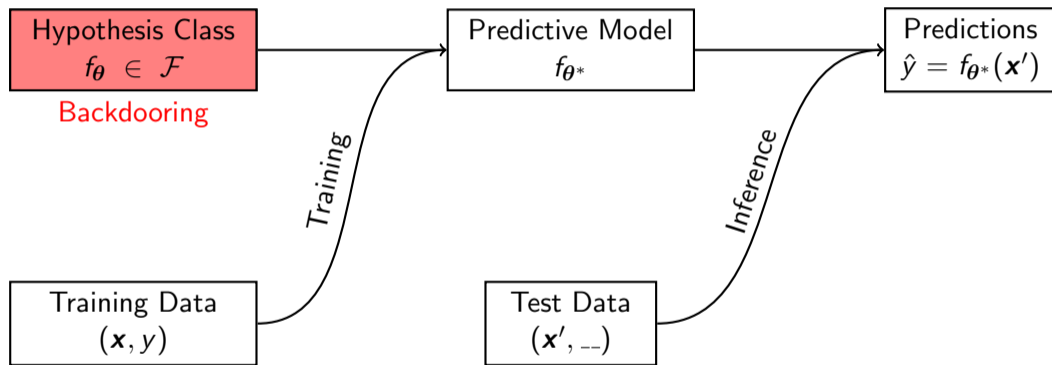
## Warmup: (Supervised) Machine Learning Pipeline



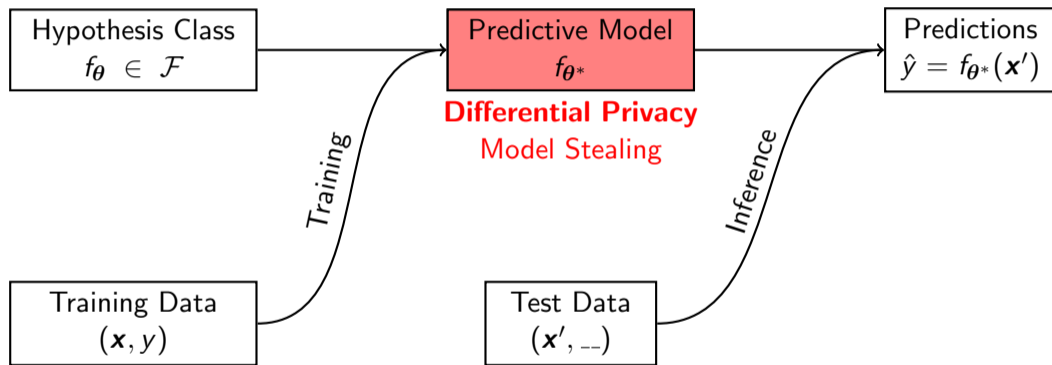
# Warmup: (Supervised) Machine Learning Pipeline



# Warmup: (Supervised) Machine Learning Pipeline

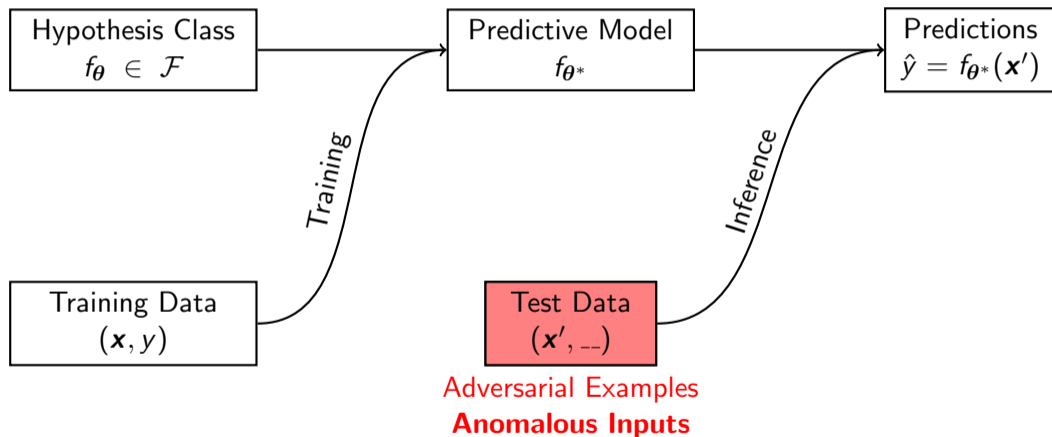


# Warmup: (Supervised) Machine Learning Pipeline

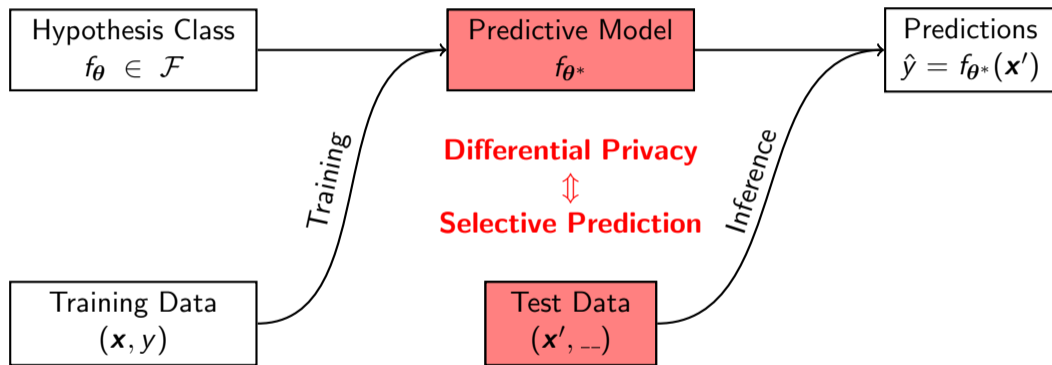




# Warmup: (Supervised) Machine Learning Pipeline

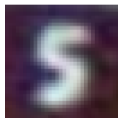
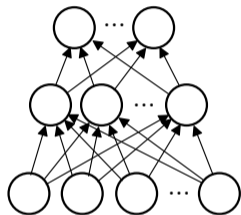


# Warmup: (Supervised) Machine Learning Pipeline



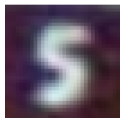
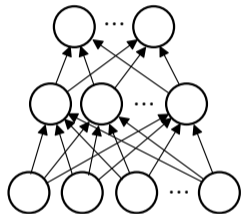
# Motivation: Input Sample Rejection

5

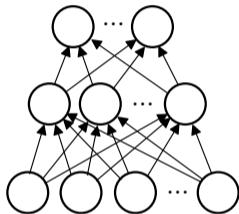


# Motivation: Input Sample Rejection

5

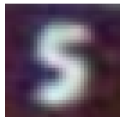
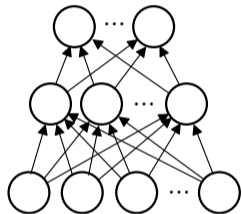


2

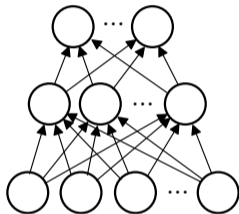


# Motivation: Input Sample Rejection

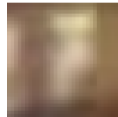
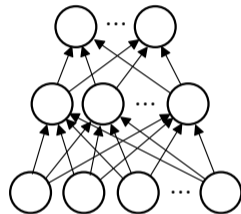
5



2



7? 1? 9?



## Selective Classification (SC)

**Selective classification adds a rejection class  $\perp$  via a gating mechanism.**

# Selective Classification (SC)

Selective classification adds a rejection class  $\perp$  via a gating mechanism.

**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  $\mathbf{x}$ .

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

# Selective Classification (SC)

Selective classification adds a rejection class  $\perp$  via a gating mechanism.

**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  $\mathbf{x}$ .

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

The performance of a selective classifier  $(f, g)$  on a dataset  $D$  is assessed based on

- the *coverage* of  $(f, g)$ , i.e. what fraction of points we predict on; and
- the selective *accuracy* of  $(f, g)$  on the points it accepts.

$$\text{cov}_\tau(f, g) = \frac{|\{\mathbf{x} : g(\mathbf{x}) \leq \tau\}|}{|D|}$$

$$\text{acc}_\tau(f, g) = \frac{|\{\mathbf{x} : f(\mathbf{x}) = y, g(\mathbf{x}) \leq \tau\}|}{|\{\mathbf{x} : g(\mathbf{x}) \leq \tau\}|}$$



## Training stage

---

Training set



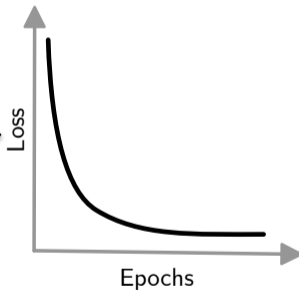
## Training stage

---

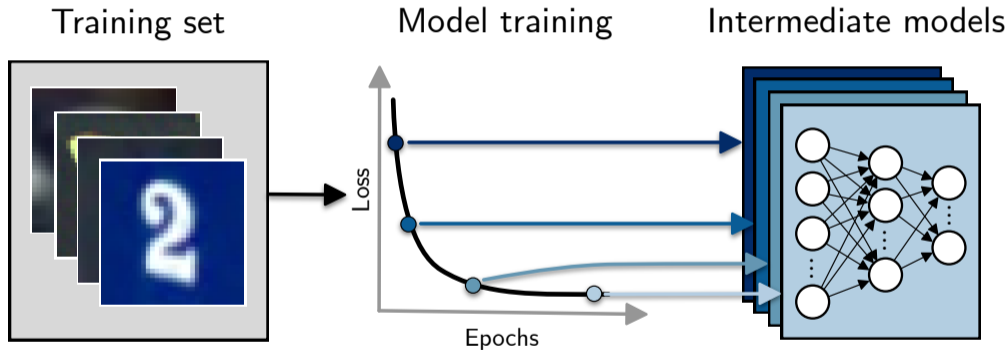
Training set



Model training



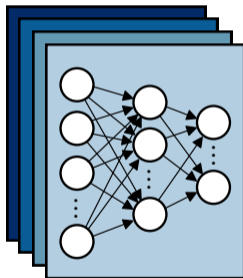
## Training stage



## Testing stage

---

Intermediate models

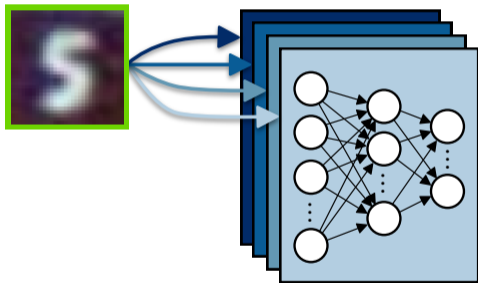


## Testing stage

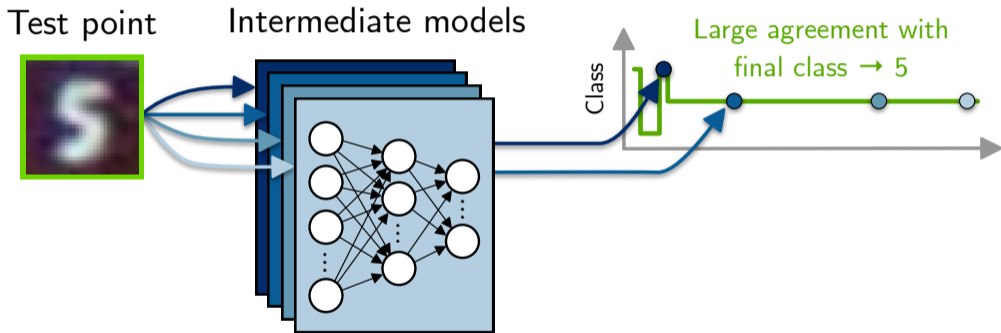
---

Test point

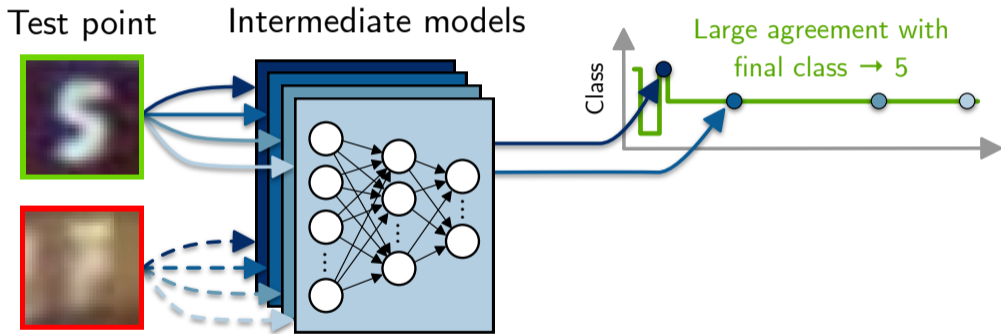
Intermediate models



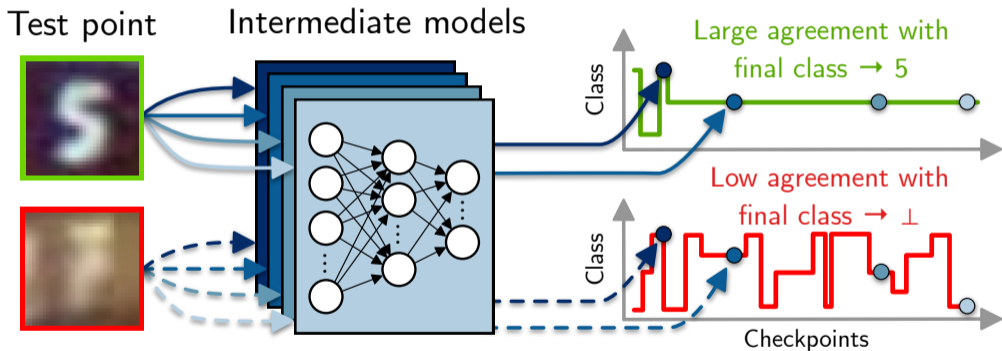
## Testing stage



## Testing stage



## Testing stage





## 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:*

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

The above DP bound is governed by two parameters:

- $\epsilon \in \mathbb{R}_+$  which specifies the privacy level; and
- $\delta \in [0, 1]$  which allows for a small violation of the bound.

The most widely used implementation for ensuring DP in deep neural nets is DP-SGD.

# Impacts of SC on DP Guarantees

## Post-Processing

If a function  $\phi(x)$  satisfies  $(\epsilon, \delta)$ -DP, then for any deterministic or randomized function  $\psi(\cdot)$ ,  $\psi \circ \phi(x)$  continues to satisfy  $(\epsilon, \delta)$ -DP.

**Applicable to:** Softmax Response (SR), Monte-Carlo Dropout (MCDO), Deep Gamblers (DG), Self-Adaptive Training (SAT), Selective Classification Training Dynamics (SCTD)

# Impacts of SC on DP Guarantees

## Post-Processing

If a function  $\phi(x)$  satisfies  $(\epsilon, \delta)$ -DP, then for any deterministic or randomized function  $\psi(\cdot), \psi \circ \phi(x)$  continues to satisfy  $(\epsilon, \delta)$ -DP.

**Applicable to:** Softmax Response (SR), Monte-Carlo Dropout (MCDO), Deep Gamblers (DG), Self-Adaptive Training (SAT), Selective Classification Training Dynamics (SCTD)

## Advanced Sequential Composition

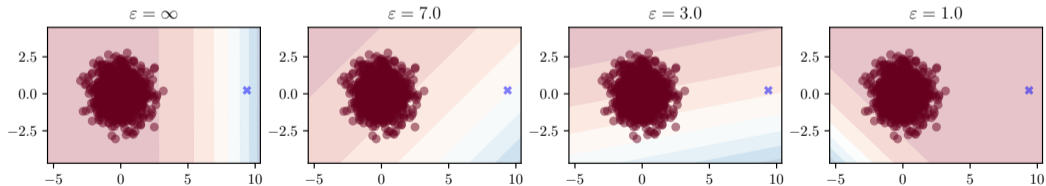
If for a set  $\{\phi_1(x), \dots, \phi_M(x)\}$  each  $\phi_i(x)$  satisfies  $(\epsilon, \delta)$ -DP, then releasing  $\psi(x) = (\phi_1(x), \dots, \phi_M(x))$  satisfies  $\approx (\sqrt{M}\epsilon, M\delta)$ -DP.

If the original  $(\epsilon, \delta)$ -DP constraint should be maintained, each function needs to satisfy  $\approx (\frac{\epsilon}{\sqrt{M}}, \frac{\delta}{M})$ -DP.

**Applicable to:** Deep Ensembles (DE), SelectiveNet (SN)

# Impacts of DP on SC Performance

- We expect DP to impact SC beyond a loss in utility.
- Sample points from a majority class and an outlier point  $\mathbf{x}^*$  from a minority class.
- Train multiple differentially private models with  $\varepsilon \in \{\infty, 7, 3, 1\}$ .
- Non-private model has best accuracy (and uncertainty) but is influenced by  $\mathbf{x}^*$ .
- All models with  $\varepsilon \in \{7, 3, 1\}$  misclassify the outlier and the changing decision boundary increases wrongful overconfidence as  $\varepsilon$  decreases.



## Evaluating SC under DP

- Default approach to quantify SC performance without accuracy bias is to align different SC approaches/models at the same accuracy and evaluate

$$s_{\text{AUC}}(f, g) = \int_0^1 \text{acc}_c(f, g) dc \quad \text{acc}_c(f, g) = \text{acc}_\tau(f, g) \quad \text{for } \tau \text{ s.t. } \text{cov}_\tau(f, g) = c$$

- Accuracy-aligning can have unintended consequences on SC performance.
- Early-stopping is the de-facto way of ensuring accuracy-alignment.
- **But:** Training for less leads to expending less privacy budget.
- Early-stopping yields a DP model with greater privacy than the targeted  $\epsilon$ .

**How do we quantify performance across SC methods and  $\epsilon$ -levels where accuracy-alignment is not possible?**

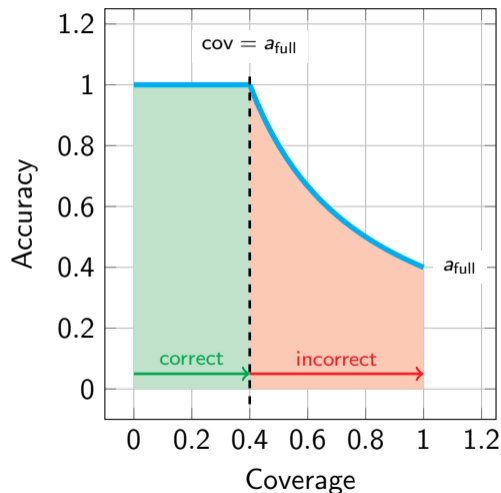
# Upper Bound On Selective Classification Performance

## Definition: Upper SC Perf. Bound

The upper bound for selective classification performance for a fixed full-coverage accuracy  $a_{full} \in [0, 1]$  and a variable coverage level  $c \in [0, 1]$  is given by

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

Optimal SC methods accept all correct points first and incorrect points afterwards.



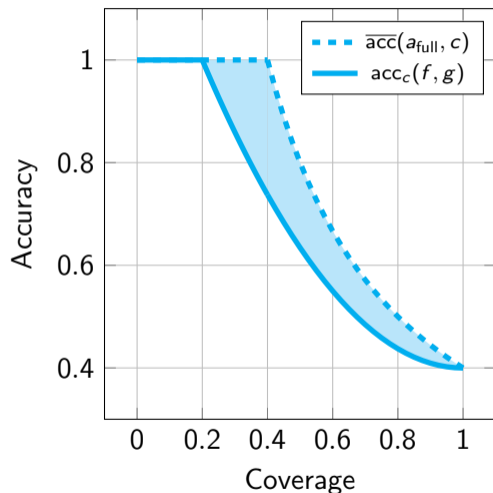
# Accuracy-Normalized Score For Selective Classification

## Definition: Acc-normalized SC Score

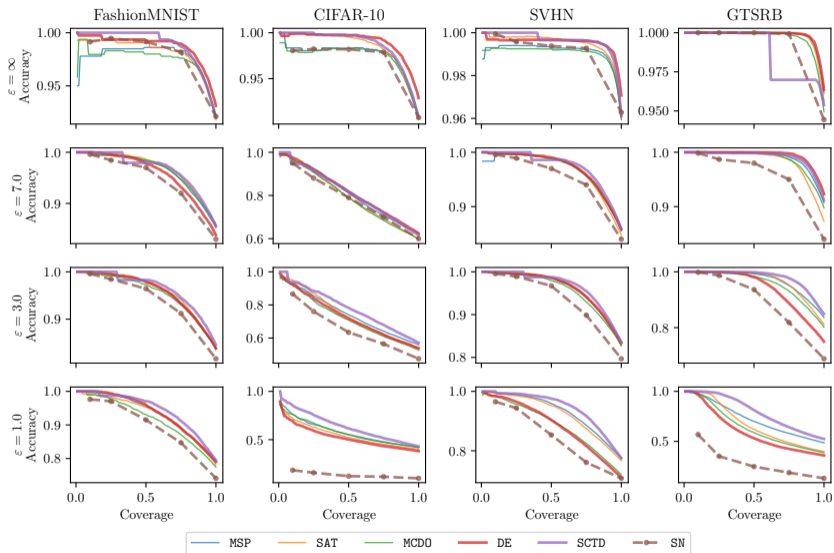
The accuracy-normalized selective classification score  $s_{a_{full}}(f, g)$  for a selective classifier  $(f, g)$  with full-coverage accuracy  $a_{full}$  is given by

$$s_{a_{full}}(f, g) = \int_0^1 (\overline{acc}(a_{full}, c) - acc_c(f, g)) dc$$
$$\approx \sum_c (\overline{acc}(a_{full}, c) - acc_c(f, g))$$

A good selective classifier should achieve a low score ( $s_{a_{full}}(f, g) \approx 0$ ), indicating closeness to the optimal bound  $\overline{acc}(a_{full}, c)$ .

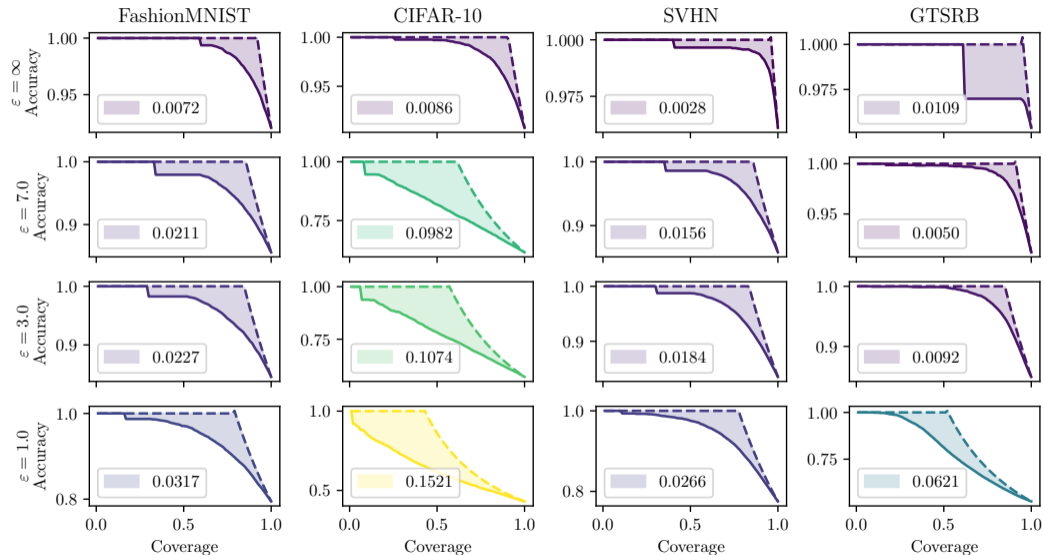


# Accuracy-Coverage Tradeoff Across Datasets & $\epsilon$ Levels





# Upper Bound Closeness for SCTD



# Accuracy-Normalized Selective Classification Performance

	FashionMNIST				CIFAR-10			
	$\epsilon = \infty$	$\epsilon = 7$	$\epsilon = 3$	$\epsilon = 1$	$\epsilon = \infty$	$\epsilon = 7$	$\epsilon = 3$	$\epsilon = 1$
MSP	0.019 ( $\pm 0.000$ )	0.023 ( $\pm 0.000$ )	0.027 ( $\pm 0.002$ )	0.041 ( $\pm 0.001$ )	0.019 ( $\pm 0.000$ )	0.105 ( $\pm 0.002$ )	0.133 ( $\pm 0.002$ )	0.205 ( $\pm 0.001$ )
SAT	0.014 ( $\pm 0.000$ )	<b>0.020 (<math>\pm 0.001</math>)</b>	0.026 ( $\pm 0.002$ )	0.043 ( $\pm 0.002$ )	0.010 ( $\pm 0.000$ )	0.107 ( $\pm 0.000$ )	0.128 ( $\pm 0.000$ )	0.214 ( $\pm 0.002$ )
MCDO	0.020 ( $\pm 0.002$ )	0.023 ( $\pm 0.001$ )	0.030 ( $\pm 0.003$ )	0.053 ( $\pm 0.001$ )	0.021 ( $\pm 0.001$ )	0.110 ( $\pm 0.000$ )	0.142 ( $\pm 0.000$ )	0.201 ( $\pm 0.000$ )
DE	<b>0.010 (<math>\pm 0.003</math>)</b>	0.027 ( $\pm 0.002$ )	0.027 ( $\pm 0.002$ )	0.039 ( $\pm 0.000$ )	<b>0.007 (<math>\pm 0.001</math>)</b>	<b>0.099 (<math>\pm 0.002</math>)</b>	0.138 ( $\pm 0.000$ )	0.222 ( $\pm 0.000$ )
SN	<b>0.008 (<math>\pm 0.002</math>)</b>	0.058 ( $\pm 0.001$ )	0.056 ( $\pm 0.001$ )	0.064 ( $\pm 0.002$ )	0.015 ( $\pm 0.000$ )	0.155 ( $\pm 0.003$ )	0.154 ( $\pm 0.002$ )	0.173 ( $\pm 0.001$ )
SCTD	<b>0.007 (<math>\pm 0.001</math>)</b>	<b>0.021 (<math>\pm 0.001</math>)</b>	<b>0.023 (<math>\pm 0.003</math>)</b>	<b>0.032 (<math>\pm 0.002</math>)</b>	<b>0.009 (<math>\pm 0.002</math>)</b>	<b>0.098 (<math>\pm 0.001</math>)</b>	<b>0.107 (<math>\pm 0.001</math>)</b>	<b>0.152 (<math>\pm 0.001</math>)</b>
	SVHN				GTSRB			
MSP	0.008 ( $\pm 0.001$ )	0.020 ( $\pm 0.001$ )	0.024 ( $\pm 0.001$ )	0.040 ( $\pm 0.001$ )	<b>0.001 (<math>\pm 0.001</math>)</b>	0.006 ( $\pm 0.002$ )	0.017 ( $\pm 0.000$ )	0.109 ( $\pm 0.002$ )
SAT	<b>0.004 (<math>\pm 0.000</math>)</b>	0.019 ( $\pm 0.000$ )	<b>0.021 (<math>\pm 0.002</math>)</b>	0.044 ( $\pm 0.002$ )	<b>0.001 (<math>\pm 0.001</math>)</b>	0.008 ( $\pm 0.001$ )	0.014 ( $\pm 0.000$ )	0.089 ( $\pm 0.000$ )
MCDO	0.009 ( $\pm 0.000$ )	0.019 ( $\pm 0.001$ )	0.027 ( $\pm 0.002$ )	0.069 ( $\pm 0.001$ )	0.002 ( $\pm 0.001$ )	0.007 ( $\pm 0.001$ )	0.023 ( $\pm 0.001$ )	0.110 ( $\pm 0.001$ )
DE	<b>0.004 (<math>\pm 0.001</math>)</b>	0.018 ( $\pm 0.001$ )	<b>0.022 (<math>\pm 0.002</math>)</b>	0.067 ( $\pm 0.003$ )	<b>0.001 (<math>\pm 0.000</math>)</b>	<b>0.003 (<math>\pm 0.002</math>)</b>	0.027 ( $\pm 0.004$ )	0.127 ( $\pm 0.002$ )
SN	<b>0.004 (<math>\pm 0.000</math>)</b>	0.055 ( $\pm 0.001$ )	0.052 ( $\pm 0.000$ )	0.096 ( $\pm 0.000$ )	<b>0.001 (<math>\pm 0.001</math>)</b>	0.050 ( $\pm 0.004$ )	0.044 ( $\pm 0.001$ )	0.091 ( $\pm 0.004$ )
SCTD	<b>0.003 (<math>\pm 0.001</math>)</b>	<b>0.016 (<math>\pm 0.001</math>)</b>	<b>0.018 (<math>\pm 0.002</math>)</b>	<b>0.027 (<math>\pm 0.001</math>)</b>	0.011 ( $\pm 0.002$ )	<b>0.005 (<math>\pm 0.000</math>)</b>	<b>0.009 (<math>\pm 0.000</math>)</b>	<b>0.062 (<math>\pm 0.001</math>)</b>

SCTD has the strongest performance, i.e. the lowest bound distance.

# Coverage Required For Non-Private Full-Coverage Accuracy

	FashionMNIST			CIFAR-10		
	$\epsilon = 7$	$\epsilon = 3$	$\epsilon = 1$	$\epsilon = 7$	$\epsilon = 3$	$\epsilon = 1$
MSP	0.83 ( $\pm 0.01$ )	0.80 ( $\pm 0.01$ )	0.65 ( $\pm 0.03$ )	<b>0.29 (<math>\pm 0.02</math>)</b>	0.14 ( $\pm 0.04$ )	0.00 ( $\pm 0.00$ )
SAT	<b>0.86 (<math>\pm 0.00</math>)</b>	0.81 ( $\pm 0.01$ )	0.67 ( $\pm 0.02$ )	0.25 ( $\pm 0.01$ )	<b>0.19 (<math>\pm 0.02</math>)</b>	0.00 ( $\pm 0.00$ )
MCDO	<b>0.84 (<math>\pm 0.02</math>)</b>	0.79 ( $\pm 0.00$ )	0.56 ( $\pm 0.02$ )	0.25 ( $\pm 0.01$ )	0.12 ( $\pm 0.02$ )	0.00 ( $\pm 0.00$ )
DE	0.75 ( $\pm 0.00$ )	0.75 ( $\pm 0.01$ )	0.61 ( $\pm 0.01$ )	0.22 ( $\pm 0.01$ )	0.09 ( $\pm 0.00$ )	0.00 ( $\pm 0.00$ )
SCTD	<b>0.86 (<math>\pm 0.01</math>)</b>	<b>0.84 (<math>\pm 0.02</math>)</b>	<b>0.73 (<math>\pm 0.01</math>)</b>	<b>0.26 (<math>\pm 0.03</math>)</b>	<b>0.20 (<math>\pm 0.03</math>)</b>	<b>0.04 (<math>\pm 0.04</math>)</b>
	SVHN			GTSRB		
MSP	0.74 ( $\pm 0.00$ )	0.67 ( $\pm 0.01$ )	0.49 ( $\pm 0.02$ )	0.90 ( $\pm 0.01$ )	0.71 ( $\pm 0.03$ )	0.13 ( $\pm 0.00$ )
SAT	0.72 ( $\pm 0.00$ )	0.67 ( $\pm 0.01$ )	0.45 ( $\pm 0.02$ )	0.86 ( $\pm 0.00$ )	0.74 ( $\pm 0.00$ )	0.20 ( $\pm 0.03$ )
MCDO	0.74 ( $\pm 0.00$ )	0.64 ( $\pm 0.00$ )	0.23 ( $\pm 0.03$ )	0.90 ( $\pm 0.01$ )	0.69 ( $\pm 0.01$ )	0.14 ( $\pm 0.01$ )
DE	0.69 ( $\pm 0.01$ )	0.62 ( $\pm 0.01$ )	0.22 ( $\pm 0.00$ )	<b>0.93 (<math>\pm 0.00</math>)</b>	0.57 ( $\pm 0.08$ )	0.10 ( $\pm 0.04$ )
SCTD	<b>0.78 (<math>\pm 0.01</math>)</b>	<b>0.72 (<math>\pm 0.00</math>)</b>	<b>0.59 (<math>\pm 0.02</math>)</b>	<b>0.93 (<math>\pm 0.01</math>)</b>	<b>0.83 (<math>\pm 0.03</math>)</b>	<b>0.30 (<math>\pm 0.02</math>)</b>

**SCTD retains the largest amount of coverage.**

## 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



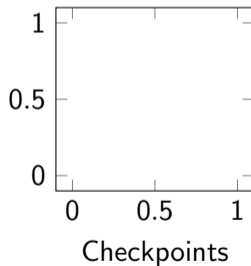
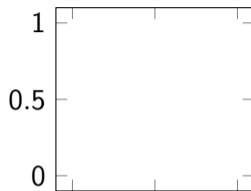
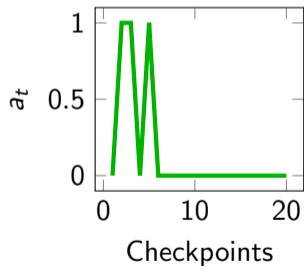
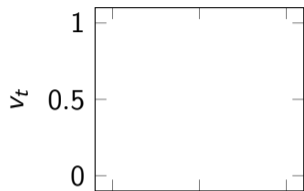
Nicolas

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

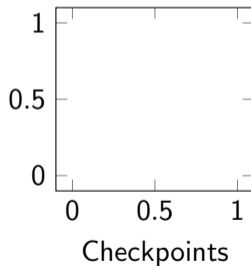
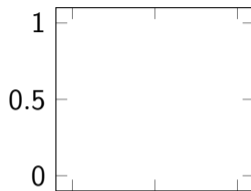
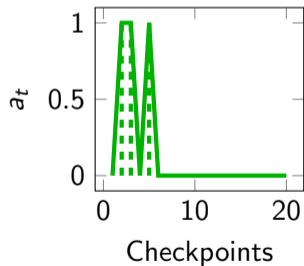
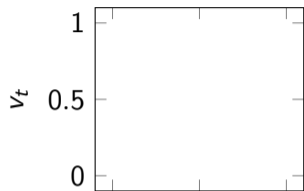
<https://arxiv.org/abs/2305.18393>

Backup

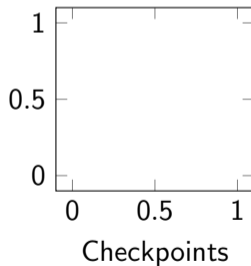
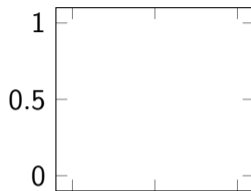
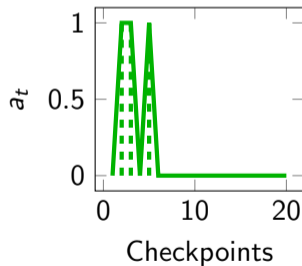
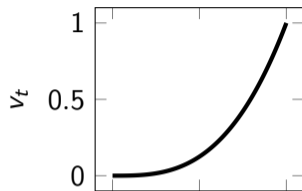
# The Sum Score $s_{\text{sum}}$



# The Sum Score $s_{\text{sum}}$

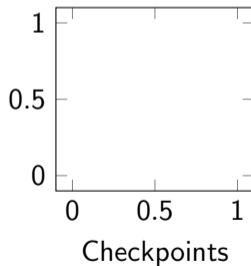
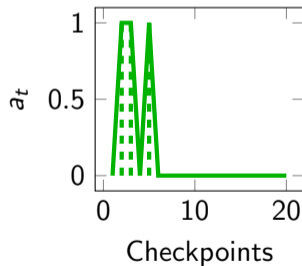
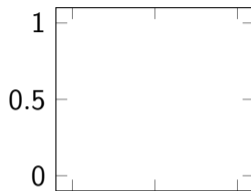
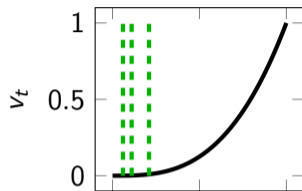


# The Sum Score $s_{\text{sum}}$

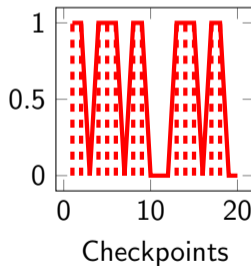
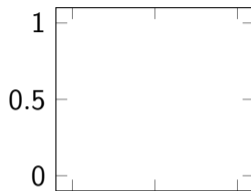
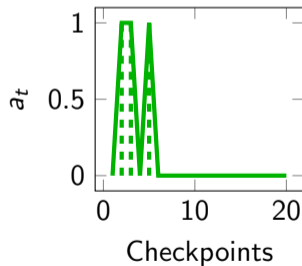
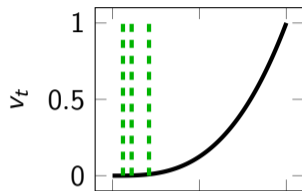




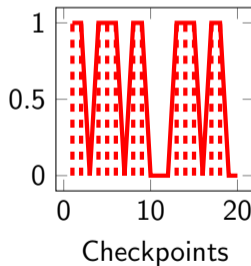
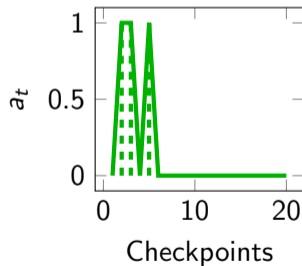
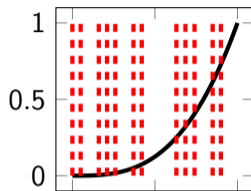
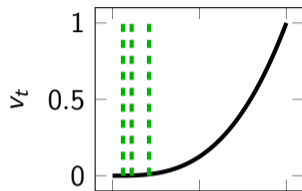
# The Sum Score $s_{\text{sum}}$



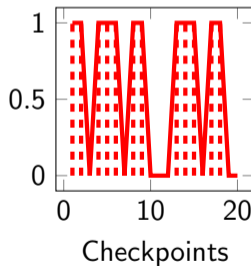
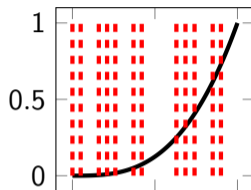
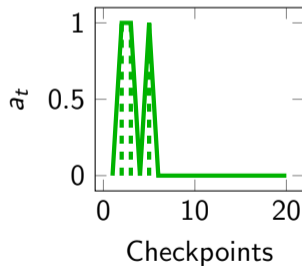
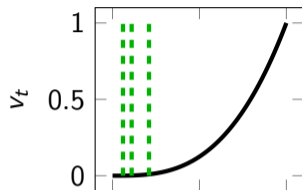
# The Sum Score $s_{\text{sum}}$



# The Sum Score $s_{\text{sum}}$



# The Sum Score $s_{\text{sum}}$



1. Denote  $L = f_T(\mathbf{x})$ , i.e. the label our final model predicts.
2. If  $\exists t$  s.t.  $a_t = 1$ , compute

$$s_{\text{sum}} = \sum v_t a_t$$

else accept  $\mathbf{x}$  with prediction  $L$ .

3. If  $s_{\text{sum}} < \tau$  accept  $\mathbf{x}$  with prediction  $L$ , else reject ( $\perp$ ).

---

## Algorithm 0: SCTD

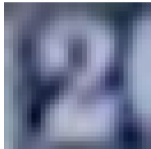
---

**Require:** Checkpointed model sequence  $\{f_1, \dots, f_T\}$ , query point  $\mathbf{x}$ , weighting parameter  $k \in [0, \infty)$ .

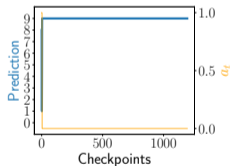
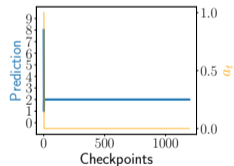
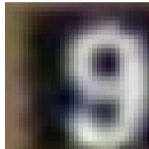
- 1: Compute prediction of last model:  $L \leftarrow f_T(\mathbf{x})$
  - 2: Compute disagreement and weighting of intermediate predictions:
  - 3: **for**  $t \in [T]$  **do**
  - 4:   **if**  $f_t(\mathbf{x}) = L$  **then**  $a_t \leftarrow 0$  **else**  $a_t \leftarrow 1$
  - 5:    $v_t \leftarrow 1 - \left(\frac{t}{T}\right)^k$
  - 6: **end for**
  - 7: Compute sum score:  $s_{\text{sum}} \leftarrow \sum_t a_t v_t$
  - 8: **if**  $s_{\text{sum}} \leq \tau$  **then** accept  $f(\mathbf{x}) = L$  **else** reject with  $f(\mathbf{x}) = \perp$
-

# Individual SVHN Example

Prediction: 2 – Label: 2



Prediction: 9 – Label: 9

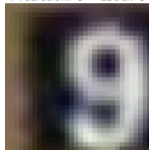


# Individual SVHN Example

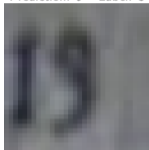
Prediction: 2 – Label: 2



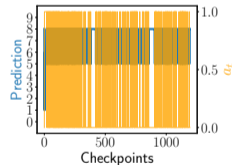
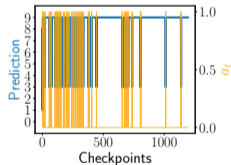
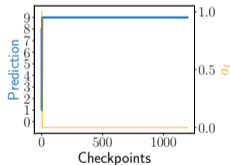
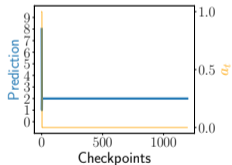
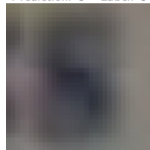
Prediction: 9 – Label: 9



Prediction: 9 – Label: 3



Prediction: 8 – Label: 8



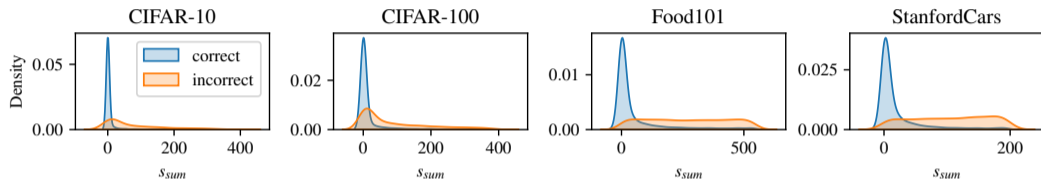
## SC Performance Over Accuracy-Coverage Curve

Dataset	SR	SAT	DE	SCTD	DE+SCTD
CIFAR-10	0.971	0.978	0.980	0.980	<b>0.981</b>
CIFAR-100	0.895	0.900	0.909	0.909	<b>0.912</b>
Food101	0.935	0.939	0.945	0.946	<b>0.947</b>
StanfordCars	0.920	0.927	0.930	0.931	<b>0.934</b>

- SCTD offers comparable performance to DE.
- Combining DE with SCTD (DE+SCTD) delivers new SOTA performance.

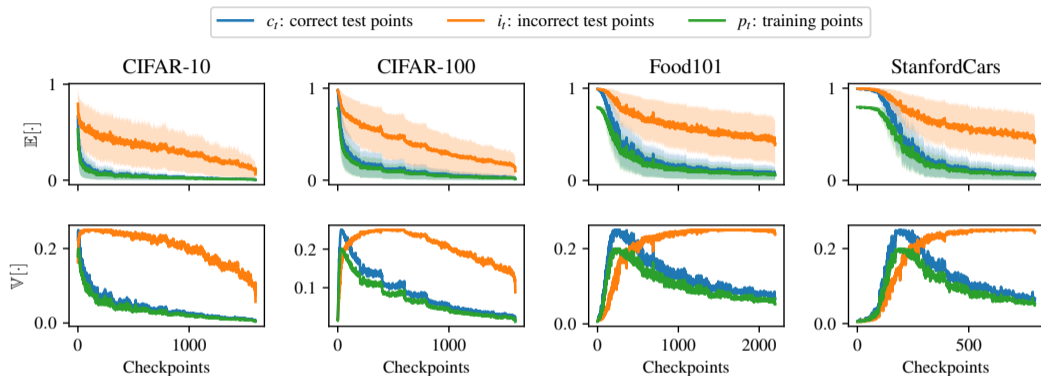


# Score Distributions Of Correct And Incorrect Points



- Correct predictions concentrate at 0 (prediction stability).
- Incorrect predictions spread over a wide score range (prediction instability).

# Monitoring $\mathbb{E}[\cdot]$ and $\mathbb{V}[\cdot]$



- Patterns for optimized points overlap with correctly classified test points.
- Correctly classified points have both  $\mathbb{E}[c_t]$  and  $\mathbb{V}[c_t]$  quickly decreasing to 0.
- Incorrectly classified points exhibit large expectations and variances.

# Performance of $s_{\max}$ vs $s_{\text{sum}}$

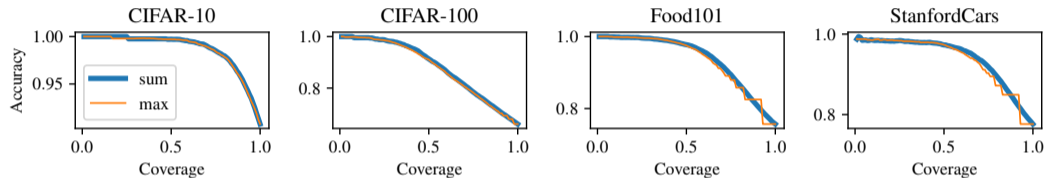
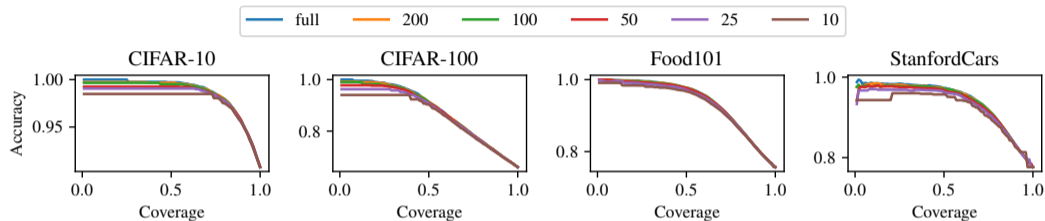


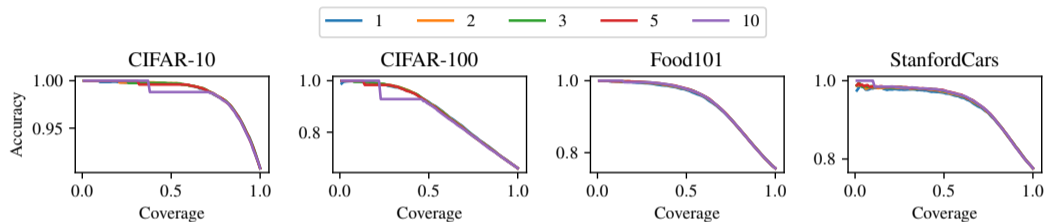
Figure: **Comparing  $s_{\max}$  and  $s_{\text{sum}}$  performance.** It is evident that  $s_{\text{sum}}$  effectively denoises  $s_{\max}$ .

# Ablation On Number of Checkpoints



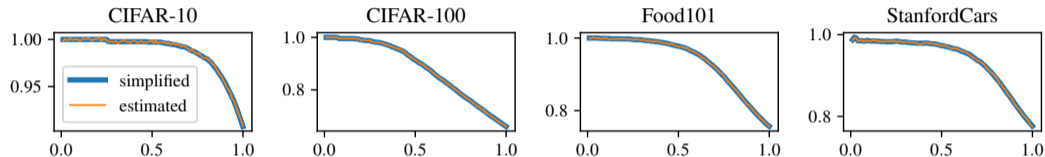
**Figure: Coverage/error trade-off of  $\text{SCTD}(s_{\text{avg}})$  for varying total number of checkpoints.** As the checkpointing resolution decreases, accuracy at low coverage increasingly degrades, thereby showing that a detailed characterization of the training dynamics is helpful to attain high target accuracy.

# Ablation Over $k$



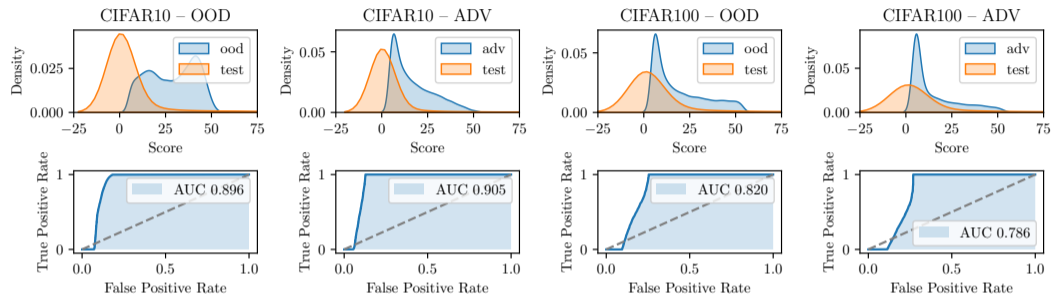
**Figure: Coverage/error trade-off of  $\text{SCTD}(s_{\text{avg}}, k)$  for varying checkpoint weighting  $k$  as used in  $v_t$ . We observe strong performance for  $k \in [2, 5]$  across datasets.**

## Incorporating $e_t$ And $v_t$ Into $s_{\text{sum}}$



**Figure: Coverage/error trade-off when incorporating  $e_t$  and  $v_t$  into  $s_{\text{max}}$  and  $s_{\text{sum}}$ .** We see that our simplifying assumptions match the performance attained from empirical estimation of  $e_t$  and  $v_t$ .

# Detectability of OOD and Adversarial Examples



**Figure: Performance of  $SCTD(s_{sum})$  on out-of-distribution (OOD) and adversarial sample detection on CIFAR-10 and CIFAR-100.** The first row shows the score distribution of the in-distribution test set vs the SVHN OOD test set or a set consisting of adversarial samples generated via a PGD attack in the final model. The second row shows the effectiveness of a thresholding mechanism by computing the area under the ROC curve.

---

## Algorithm 1: DP-SGD

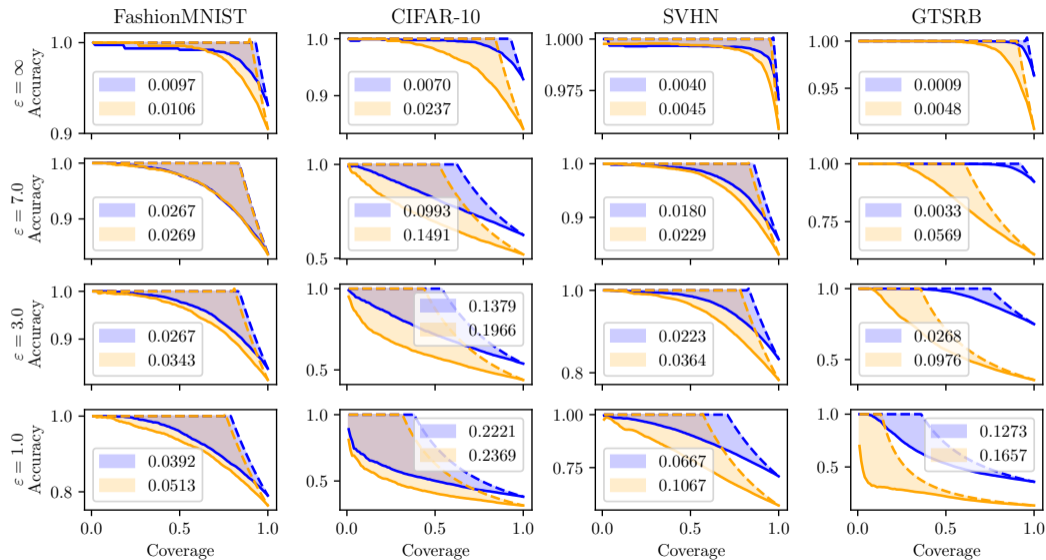
---

**Require:** Training dataset  $D$ , loss function  $\ell$ , learning rate  $\eta$ , noise multiplier  $\sigma$ , sampling rate  $q$ , clipping norm  $c$ , iterations  $T$ .

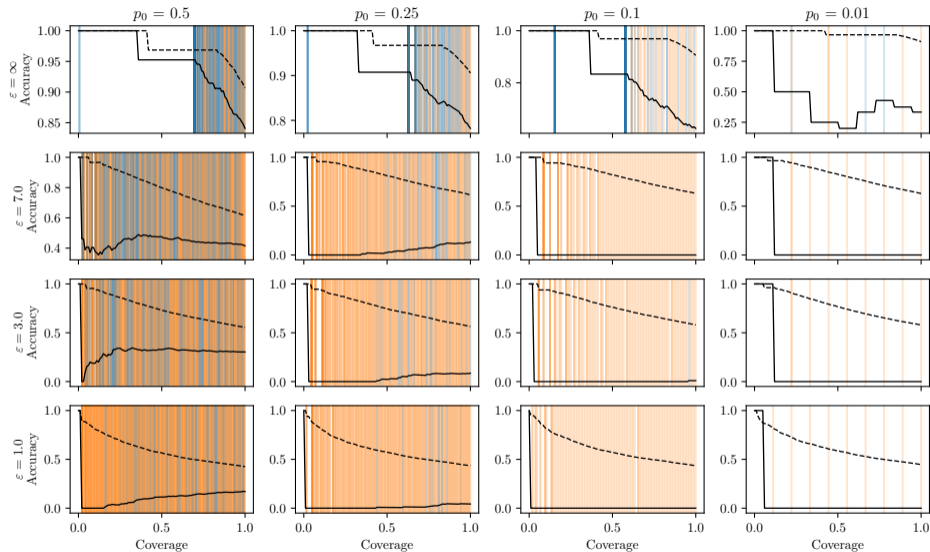
- 1: **Initialize**  $\theta_0$
  - 2: **for**  $t \in [T]$  **do**
  - 3:   **1. Per-Sample Gradient Computation**
  - 4:   Sample  $B_t$  with per-point prob.  $q$  from  $D$
  - 5:   **for**  $i \in B_t$  **do**
  - 6:      $g_t(\mathbf{x}_i) \leftarrow \nabla_{\theta_t} \ell(\theta_t, \mathbf{x}_i)$
  - 7:   **end for**
  - 8:   **2. Gradient Clipping**
  - 9:    $\bar{g}_t(\mathbf{x}_i) \leftarrow g_t(\mathbf{x}_i) / \max\left(1, \frac{\|g_t(\mathbf{x}_i)\|_2}{c}\right)$
  - 10:   **3. Noise Addition**
  - 11:    $\tilde{g}_t \leftarrow \frac{1}{|B_t|} \left(\sum_i \bar{g}_t(\mathbf{x}_i) + \mathcal{N}(0, (\sigma c)^2 \mathbf{I})\right)$
  - 12:    $\theta_{t+1} \leftarrow \theta_t - \eta \tilde{g}_t$
  - 13: **end for**
  - 14: **Output**  $\theta_T$ , privacy cost  $(\epsilon, \delta)$  computed via a privacy accounting procedure
-



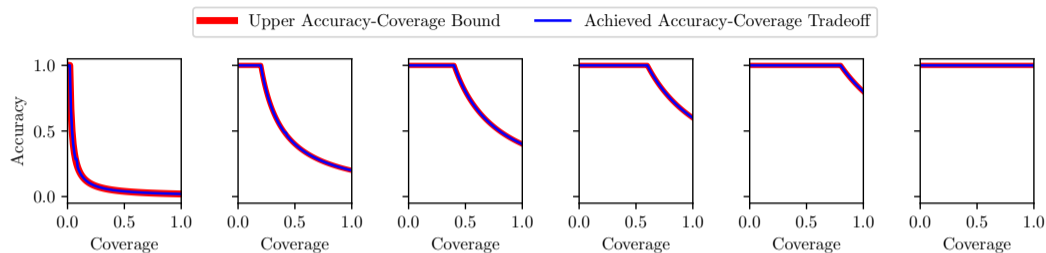
# Composition-Based Deep Ensembles VS Partitioned Deep Ensembles



# Class Imbalance Results on CIFAR-10



# Upper Bound Reachability



- Assume a binary classif. setting with label vector  $\mathbf{y} \in \{0, 1\}^{n_0+n_1}$  and  $n_0 = n_1$ .
- Generate a prediction vector  $\mathbf{p}$  which overlaps with  $\mathbf{y}$  for a fraction of  $a_{\text{full}}$ .
- Sample a scoring vector  $\mathbf{s}$  where each correct prediction is assigned a score  $s_i \sim \mathcal{U}_{0,0.5}$  and each incorrect prediction is assigned a score  $s_i \sim \mathcal{U}_{0.5,1}$ .
- This score is optimal since all  $s_i < 0.5$  correspond to a correct prediction, while all  $s_i \geq 0.5$  correspond to an incorrect prediction.