

Knowing what you don't know: Learning to abstain and beyond

Harikrishna Narasimhan

Google Research

Contributors



Wittawat
Jitkrittum



Ankit S.
Rawat



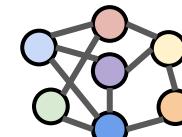
Aditya K.
Menon



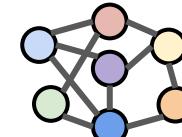
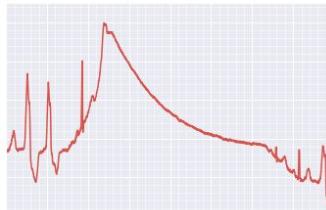
Sanjiv
Kumar

Standard classification paradigm

- Standard classification → **single model** for all samples
- However, it may be challenging to model the entire input space



-1



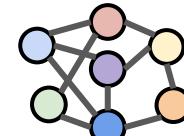
Input sample

Decision maker

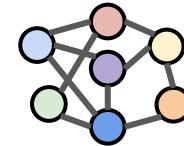
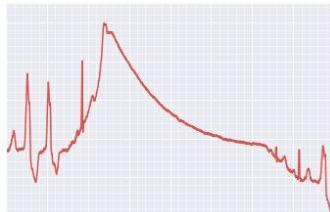
Prediction

Learning to reject

- Model can **give up** on a sample, incurring some **cost**



-1



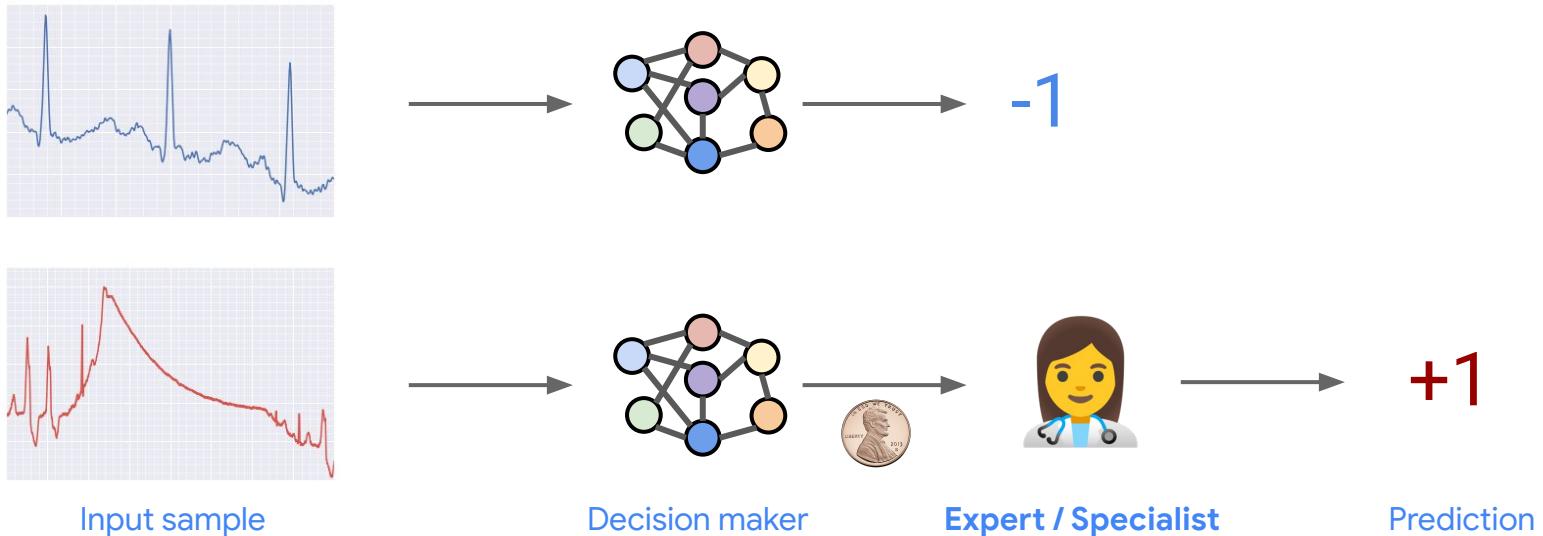
Input sample

Decision maker

Prediction

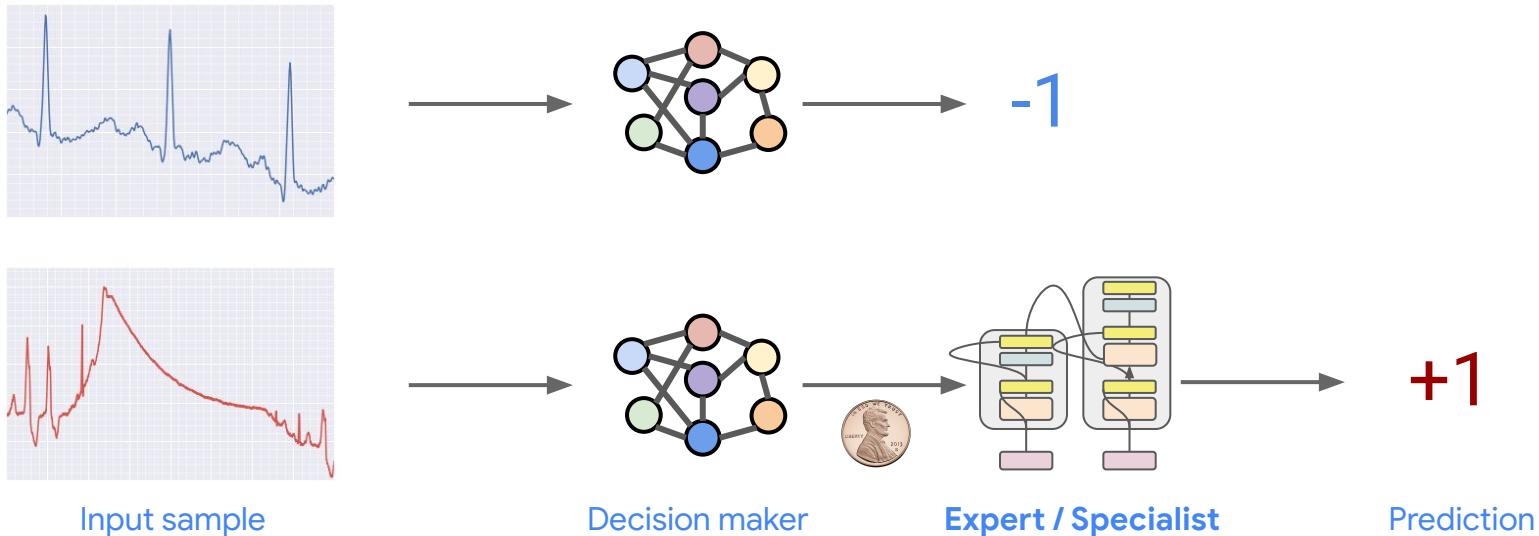
Learning to defer to an expert

- Model can **defer** to an **expert**, incurring some **cost**
 - e.g., human expert



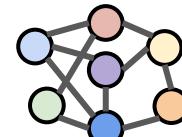
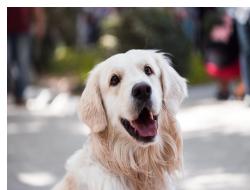
Learning to defer to an expert

- Model can **defer** to an **expert**, incurring some **cost**
 - e.g., human expert, powerful learning model

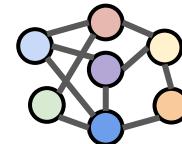


Learning to abstain on outliers

- Model can abstain on samples it deems to be **out-of-distribution (OOD)**



dog



Input sample

Decision maker

Prediction

Goal: learn the base classifier,
and the abstention rule

Cost of abstention: *classical version*

We will denote a joint classifier $h: X \rightarrow [n] \cup \{ \text{🤷} \}$. In the simplest case, one may associate a **constant cost c** to abstaining on a sample

$$1(\hat{y} \neq y, \hat{y} \neq \text{🤷}) + c \cdot 1(\hat{y} = \text{🤷})$$

Usual error when not abstaining

Constant cost when abstaining

Chow's rule: a surprisingly competitive baseline

C. Chow. On optimum recognition error and reject tradeoff.
IEEE Transactions on Information Theory, 16(1):41–46, 1970.

Bayes-optimal rejection rule: abstain on a sample when

$$\max_y \mathbb{P}(y | x) < 1 - c$$

Chow's rule: a surprisingly competitive baseline

C. Chow. On optimum recognition error and reject tradeoff.
IEEE Transactions on Information Theory, 16(1):41–46, 1970.

Bayes-optimal rejection rule: abstain on a sample when

$$\max_y \hat{\mathbb{P}}(y | x) < 1 - c$$

In practice: max softmax probability
from a standard classifier

When Chow's rule fails and ways to remedy it!

- Learning to reject
 - classical Chow's rule is very competitive
- Learning to defer to an expert
 - remedy: expert-aware Chow's rule
- Learning to abstain on outliers
 - remedy: outlier-aware Chow's rule

Cost of abstention: when deferring to an expert

In the learning to defer paradigm, the cost of invoking the expert:

$$\mathbb{1}(\hat{y} \neq y, \hat{y} \neq \text{👤}) + c_{\text{exp}}(x, y) \cdot \mathbb{1}(\hat{y} = \text{👤})$$

Usual error when not abstaining

Cost of invoking the expert

Expert cost: fixed cost + expert's error rate

A natural candidate for the expert cost would include both a **fixed cost** and the **penalty when the expert makes a mistake**

$$c_{\text{exp}}(x, y) = c_0 + \mathbf{1}(y \neq h_{\text{exp}}(x))$$

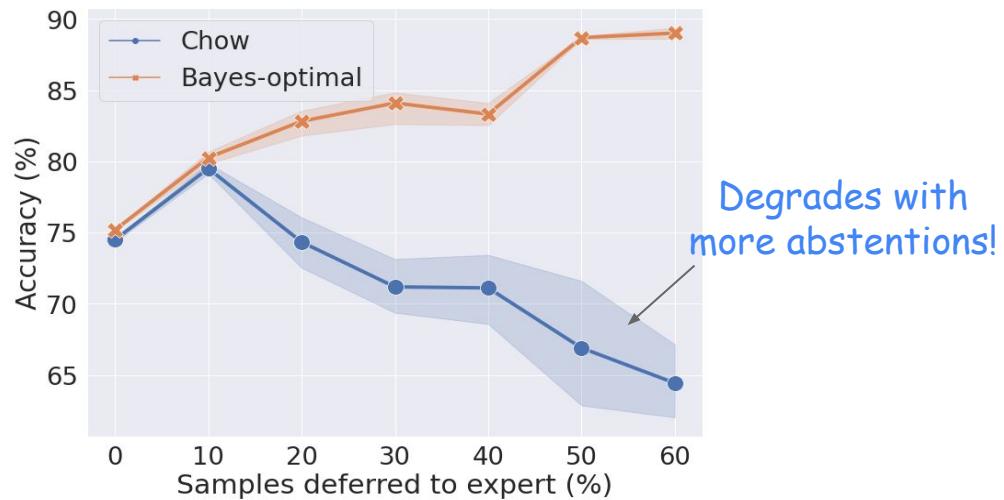
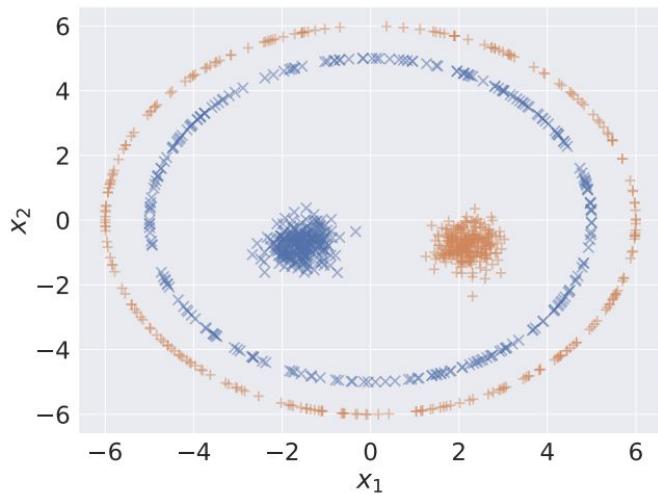
Fixed cost



(e.g. monetary cost)

Expert prediction

Chow's rule can be sub-optimal for this setting



Synthetic dataset
Base model: linear features
Expert model: quadratic features

Expert-aware Chow's rule

Bayes-optimal rule: defer on a sample when

$$\max_y \mathbb{P}(y | x) < \mathbb{E}_{y|x} [\mathbf{1}(y = h_{\text{exp}}(x))] - c_0$$

~ Base classifier's
confidence

~ Expert's confidence

Expert-aware Chow's rule

Bayes-optimal rule: defer on a sample when

When the expert's confidence is highly non-uniform, this is substantially different from Chow's rule

$$\max_y \mathbb{P}(y | x) < \mathbb{E}_{y|x} [\mathbf{1}(y = h_{\text{exp}}(x))] - c_0$$

~ Base classifier's confidence

~ Expert's confidence

Expert-aware Chow's rule

Bayes-optimal rule: defer on a sample when

$$\max_y \hat{\mathbb{P}}(y | x) < \hat{\mathbb{E}}_{y|x} [\mathbf{1}(y = h_{\text{exp}}(x))] - c_0$$


~ Base classifier's confidence ~ Expert's confidence

Unlike classical Chow, we need to estimate the expert's confidence

Separate model for expert's confidence

Raghu et al. '19 suggest training separate model to estimate expert's confidence (using a sample annotated with the expert's predictions)

$$\max_y \hat{\mathbb{P}}(y | x) < \hat{\mathbb{E}}_{y|x}[\mathbf{1}(y = h_{\text{exp}}(x))] - c_0$$

~Softmax probabilities from base classifier

~Separate model to estimate expert's confidence

Separate model for expert's confidence

- This approach has appealing properties:
 - ✓ Simple to compute
 - ✓ Approximates the Bayes deferral rule
 - ! Separate models to estimate base and expert confidence

Cost-sensitive softmax cross-entropy (CSS)

- Mozannar & Sontag '20 suggest training a joint model with an additional label ⊥

Cost-sensitive softmax cross-entropy (CSS)

- Mozannar & Sontag '20 suggest training a joint model with an additional label \perp
- Minimize a **cost-sensitive** softmax cross-entropy (**CSS**) loss

$$\ell_{\text{css}}(x, y, \bar{f}(x)) = -\log(\bar{p}_y(x)) - \mathbf{1}(y = h_{\text{exp}}(x)) \cdot \log(\bar{p}_{\perp}(x)) - c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))$$

Classification loss to train base model

Loss to estimate expert's confidence

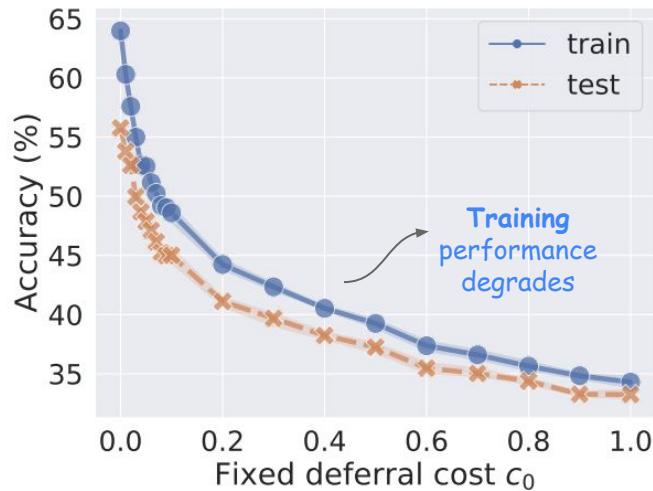
Takes into account fixed cost c_0

The case for CSS

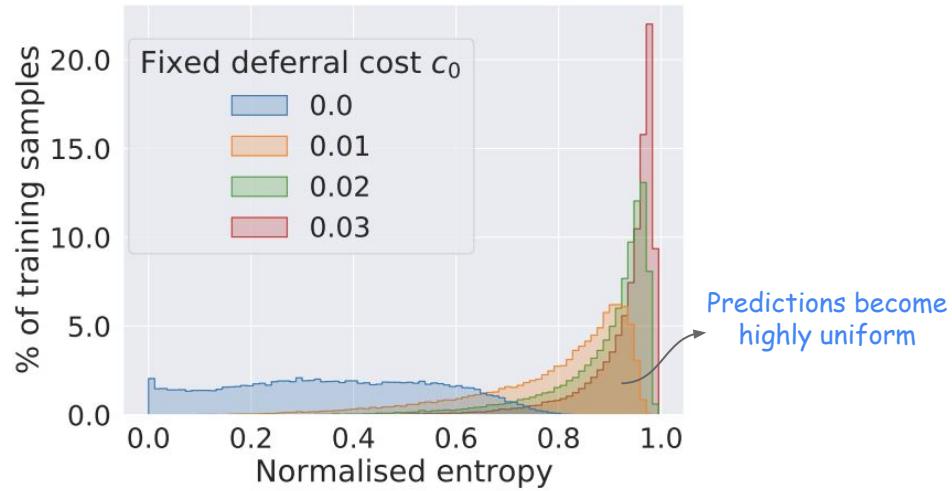
- The CSS loss has a number of appealing characteristics:
 - ✓ **Joint model** for both base classifier and expert's confidence
 - ✓ Optimal solution matches the **Bayes-optimal classifier**
 - ✓ **Empirically effective** on several benchmarks
 - ! when fixed cost $c_0 = 0$...

The case against CSS?

- CSS strongly **underfits** when there is non-zero fixed deferral cost c_0 !



CIFAR 100
ResNet8 base
ResNet32 expert



A label smoothing perspective

- CSS equivalently applies high level of **label smoothing**:

$$\ell_{\text{css}}(x, y, \bar{f}(x)) = -\log(\bar{p}_y(x)) - \mathbf{1}(y = h_{\text{exp}}(x)) \cdot \log(\bar{p}_{\perp}(x)) - c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))$$

- Encourages predictions to become highly uniform
- Low separation between true label and competing labels

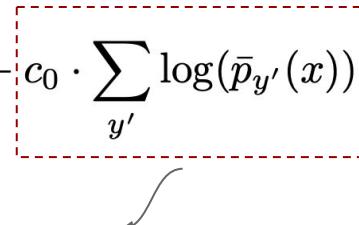
Treat all labels as candidate positive

A label smoothing perspective

- CSS equivalently applies high level of **label smoothing**:

$$\ell_{\text{CSS}}(x, y, \bar{f}(x)) = -\log(\bar{p}_y(x)) - \mathbf{1}(y = h_{\text{exp}}(x)) \cdot \log(\bar{p}_{\perp}(x)) - c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))$$

- Encourages predictions to become highly uniform
- Low separation between true label and competing labels

 Treat all labels as candidate positive

- Not apparent when $c_0 = 0$ (as in prior work)!

- $c_0 > 0$ is crucial in practical settings (e.g. when the expert is a larger model)

Solution: Set $c_0 = 0$ during training; include it in a post-hoc step

- Train base model with $c_0 = 0$, i.e., by minimizing:

$$\ell_{\text{css}}(x, y, \bar{f}(x)) = -\log(\bar{p}_y(x)) - \mathbf{1}(y = h_{\text{exp}}(x)) \cdot \log(\bar{p}_{\perp}(x)) - c_0 \cdot \sum_{y'} \log(\bar{p}_{y'}(x))$$



| | | | |
|-------------|-------------|---------|-------------|
| \bar{p}_1 | \bar{p}_2 | \dots | \bar{p}_L |
|-------------|-------------|---------|-------------|

Class probabilities

| |
|-------------------|
| \bar{p}_{\perp} |
|-------------------|

Probability that the
expert is correct

Solution: Set $c_0 = 0$ during training; include it in a post-hoc step

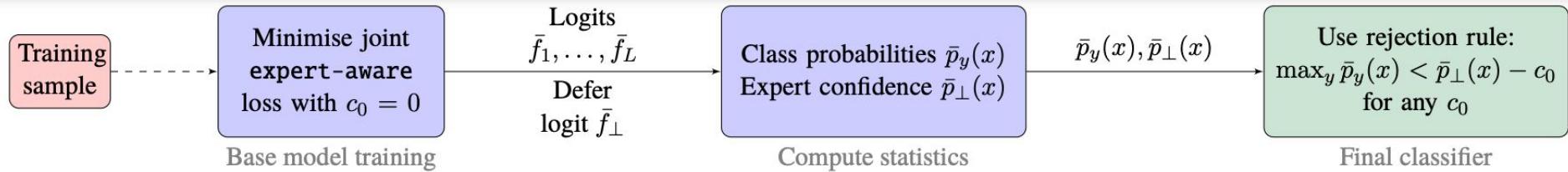
Construct a **post-hoc rejector** to include c_0 (that mimics the Bayes-optimal rule):

$$\max_y \bar{p}_y(x) < \bar{p}_{\perp}(x) - c_0$$

Probability that the expert is correct

Deferral cost

Proposal: two-step plug-in approach [Narasimhan et al '22]



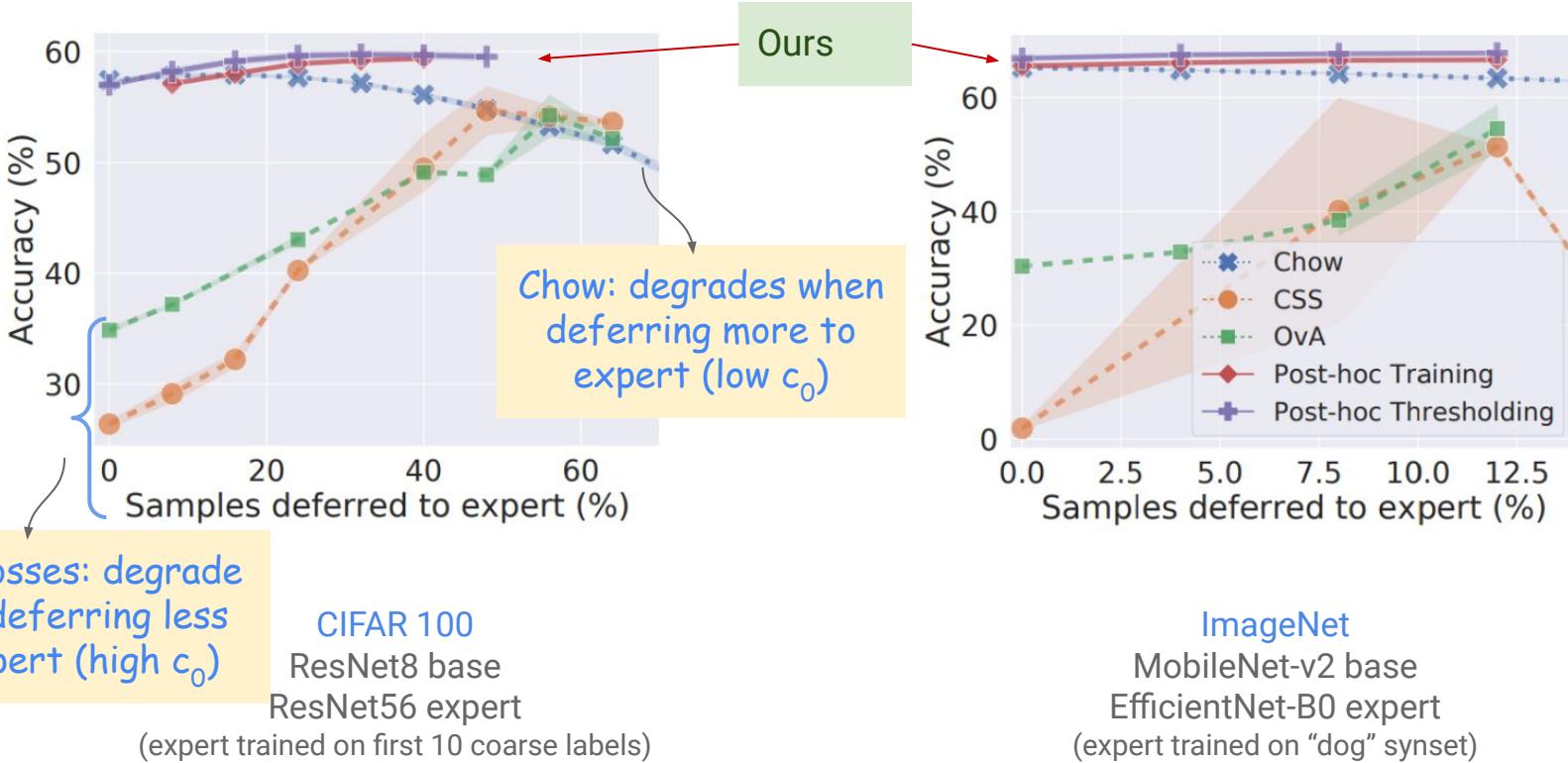
Experimental setup

- **Specialist** expert
 - Model allowed to defer to a “specialist” expert trained on a subset of labels
- Baselines
 - **Chow**: confidence thresholding based only on the deferral cost c_0
 - **CSS**: in-training loss of [Mozannar & Sontag \(2020\)](#) with c_0 included
 - **OvA**: in-training loss of [Verma and Nalisnick \(2022\)](#) with c_0 included

Ignores expert error

Underfits when c_0 is large

Experimental results: expert-aware abstention



When Chow's rule fails and ways to remedy it!

- Learning to reject
 - classical Chow's rule is very competitive
- Learning to defer to an expert
 - remedy: expert-aware Chow's rule
- Learning to abstain on outliers
 - remedy: outlier-aware Chow's rule

Learning to abstain on outliers

Abstain on “**out-of-distribution**” samples that come from distribution different from the one used for training



Inlier samples

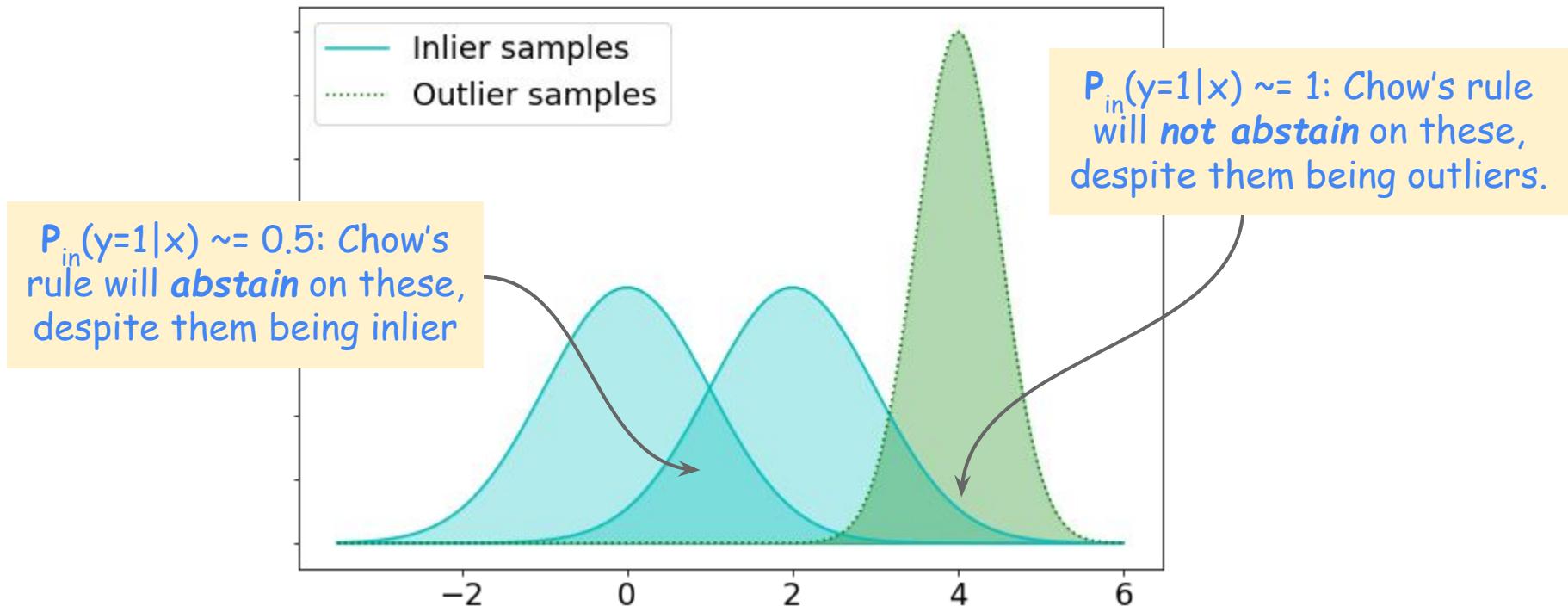
Outlier samples

Chow's rule (or the MSP scorer) is a popular baseline!

Thresholding the maximum softmax probability (MSP) from a standard classifier is a common baseline in this literature [Hendrycks et al. '17; Vaze et al. '22].

$$\max_{y \in [L]} \hat{\mathbb{P}}_{\text{in}}(y | x) < t$$

Chow's rule can fail for outlier detection



Cost of abstention: when abstaining on outliers

We need to account for both inlier **and** outlier abstentions.

$$\mathbb{P}_{\text{in}}(\hat{y} \neq y, \hat{y} \neq \text{😊}) + \alpha \cdot \mathbb{P}_{\text{in}}(\hat{y} = \text{😊}) + \beta \cdot \mathbb{P}_{\text{out}}(\hat{y} \neq \text{😊})$$

↓ ↓ ↓

Error on inlier samples
(when not abstaining)

Cost of abstaining
on inlier samples

Cost of *not* abstaining
on outlier samples

Outlier-aware Chow's rule

Bayes-optimal rule: abstain on a sample when [Narasimhan et al. '23]

$$\max_y \mathbb{P}_{\text{in}}(y | x) < 1 - \alpha + \beta \cdot \frac{\mathbb{P}_{\text{out}}(x)}{\mathbb{P}_{\text{in}}(x)}$$


Inlier class probabilities Outlier-to-inlier density ratio

Outlier-aware Chow's rule

Bayes-optimal rule: abstain on a sample when [Narasimhan et al. '23]

$$\max_y \widehat{\mathbb{P}}_{\text{in}}(y | x) < 1 - \alpha + \beta \cdot \frac{\widehat{\mathbb{P}}_{\text{out}}(x)}{\widehat{\mathbb{P}}_{\text{in}}(x)}$$

A diagram illustrating the components of the Bayes-optimal rule. Two curved arrows point from blue text labels to specific terms in the equation. The first arrow points from "Inlier class probabilities" to the term $\max_y \widehat{\mathbb{P}}_{\text{in}}(y | x)$. The second arrow points from "Outlier-to-inlier density ratio" to the term $\frac{\widehat{\mathbb{P}}_{\text{out}}(x)}{\widehat{\mathbb{P}}_{\text{in}}(x)}$.

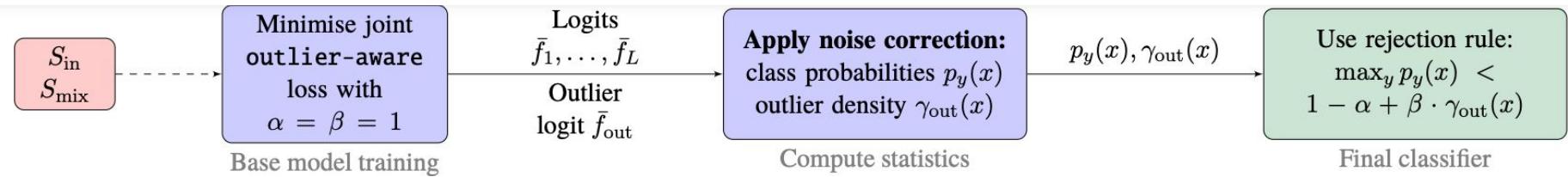
Inlier class probabilities

Outlier-to-inlier density ratio

We need to estimate both the inlier probabilities and the density ratio

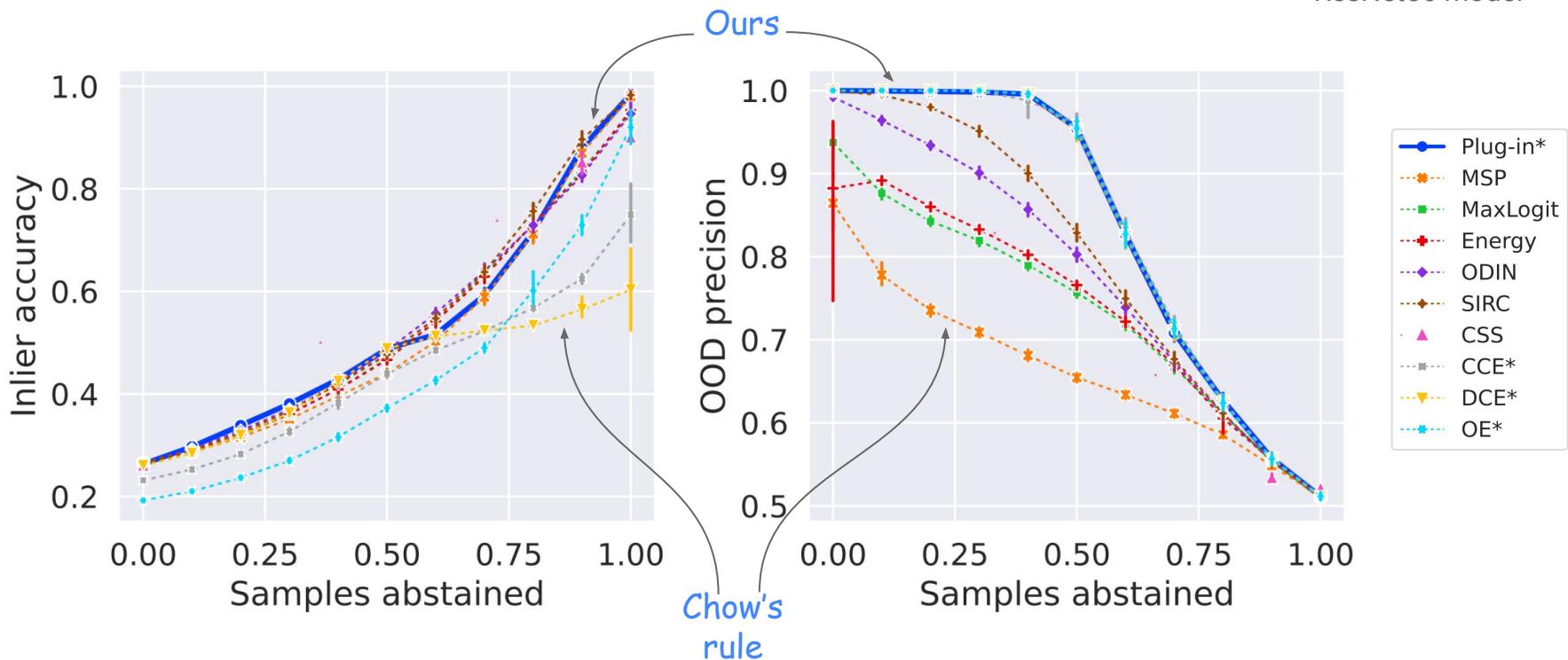
Our proposal: Two-step plug-in approach [Narasimhan et al '23]

Given: labeled **inlier sample** S_{in} , and an unlabeled **mix of inlier and outlier samples** S_{mix}



Experimental results: outlier-aware abstention

Inlier: CIFAR 100
Outlier: OpenImages
ResNet56 model



When Chow's rule fails and ways to remedy it!

- Learning to reject (L2R)
 - Classical Chow's rule is very competitive
- Learning to defer to an expert (L2D)
 - Chow may fail; use *expert-aware* Chow
- Learning to abstain on outliers (OOD)
 - Chow may fail; use *outlier-aware* Chow

| | |
|-----|--|
| L2R | $\max_y \mathbb{P}(y x) < 1 - c$ |
| L2D | $\max_y \mathbb{P}(y x) < \mathbb{E}_{y x}[\mathbf{1}(y = h_{\text{exp}}(x))] - c_0$ |
| OOD | $\max_y \mathbb{P}_{\text{in}}(y x) < 1 - \alpha + \beta \cdot \frac{\mathbb{P}_{\text{out}}(x)}{\mathbb{P}_{\text{in}}(x)}$ |

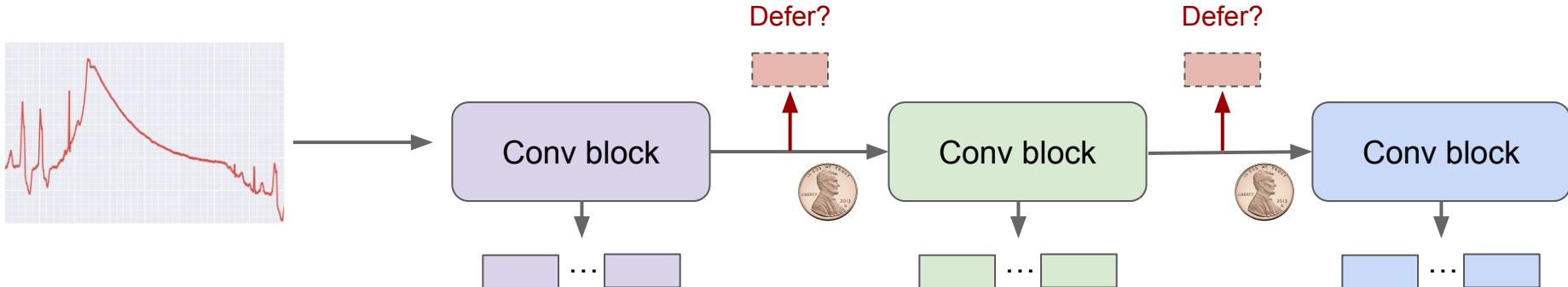
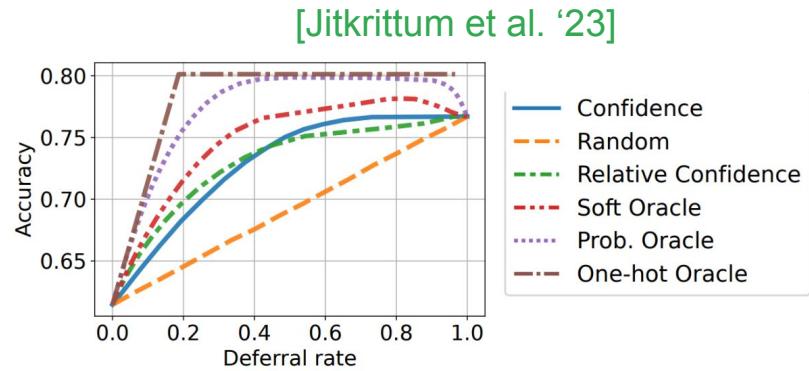
Narasimhan et al. "Post-hoc Estimators for Learning to Defer to an Expert". NeurIPS 2022.

Narasimhan et al. "Learning to Reject Meets OOD Detection: Are All Abstentions Created Equal?". Manuscript, 2023. [arXiv:2301.12386]

Thank you!

Ongoing: application to adaptive inference

- Adaptive inference
 - choose exit from an **early-exit** model
- Subsequent layers seen as “experts”
 - predict at a given exit, or defer downstream



Learning to defer: Comparison to Raghu et al. (2019)

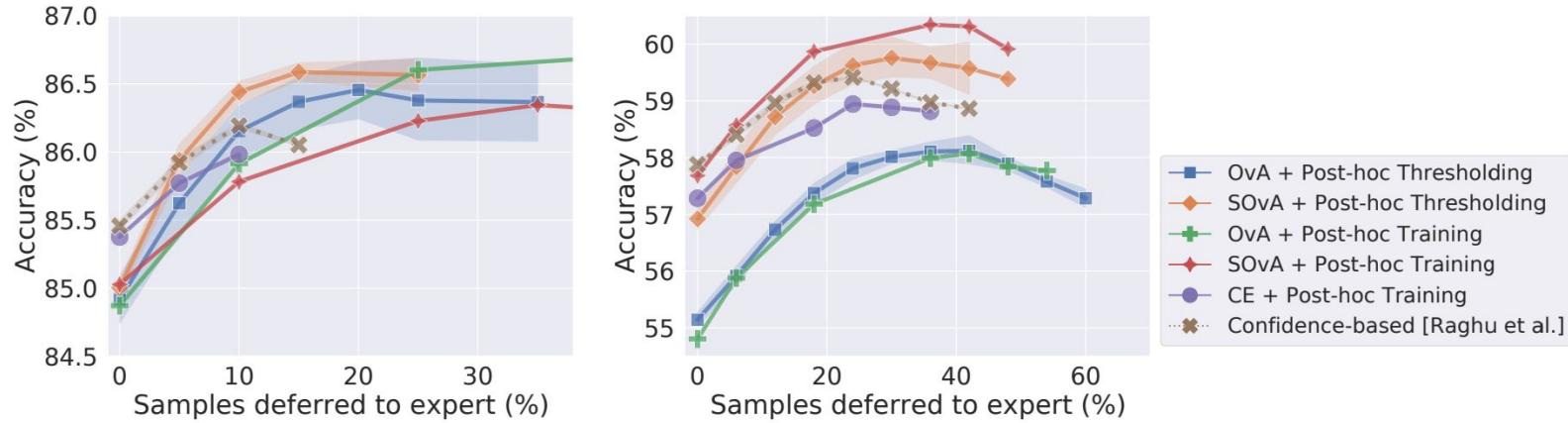


Figure 12: Additional results on CIFAR-10 (left), and CIFAR-100 (middle) in a learning to defer setting where a base model is allowed to defer to a “specialist” expert (same setting as in Figure 5). We compare the post-hoc thresholding and post-hoc training approaches with either the one-vs-all (OvA) loss or the hybrid softmax cross-entropy plus OvA loss (SOvA) loss used to train the base model. We additionally include post-hoc training on a base model trained with the softmax cross-entropy (CE) loss and the confidence-based approach of Raghu et al. [26]. The post-hoc training methods use the rejector parameterisations in Table 1.

Chow's rule can fail for outlier detection

$P_{in}(y|x)$ is the highest for the outlier class

