# Weekly Study Report

Wang Ma

2025-02-18

Electrical, Computer, and Systems Engineering Department
Rensselaer Polytechnic Institute

# Outline

# (ICLR 2025) Enhancing Uncertainty Estimation and Interpretability with Bayesian Non-negative Decision Layer

# 1. Enhancing Uncertainty Estimation and Interpretability with Bayesian Non-negative Decision Layer

In the multi-classification problem:

$$p(y_i|x) = \frac{\exp(zw_i^T + b_i)}{\sum_i \exp(zw_i^T + b_i)},$$

where $zw_i^T + b_i$ is called *logits*. In Last-Layer Lapalce Approximation, they model the weights matrix $W$ of last layer to be random variables; this paper does more.

The classification process with softmax can be seen as the special case of generative process of label $y$ (when $p(z|x) \sim \delta$ distribution):

$$p(y|x) = . \int_z p(y|z)p(z|x)dz.$$

If we model $z$ as a latent variable $\theta$, which follows some distribution, the above equation can be improved to

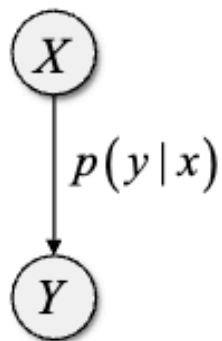$$p(y|x) = \int_\theta p(y|\theta)p(\theta|x)d\theta.$$

Notice, that $\theta$ models the output of previous layers.

To further account for epistemic uncertainty, we also model the last-layer's weight $W$ to follow some distributions $\Phi$. The we can have the following model:
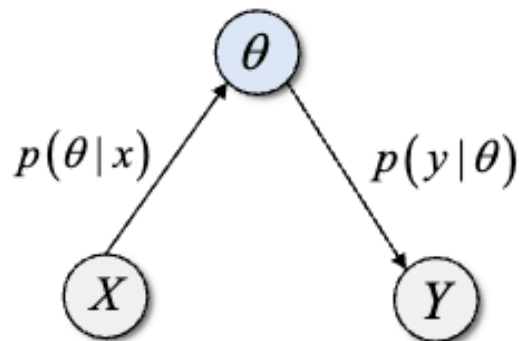
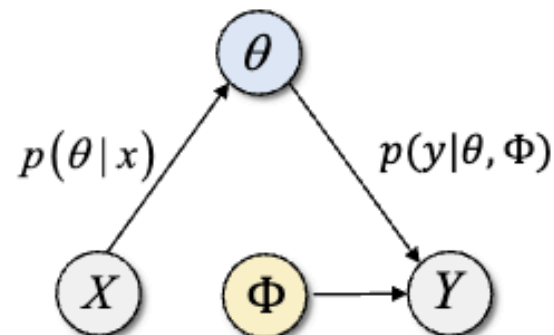$$p(y|x) = \int_{\theta,\Phi} p(y|\theta,\Phi)p(\theta|x)p(\Phi)d\theta d\Phi.$$
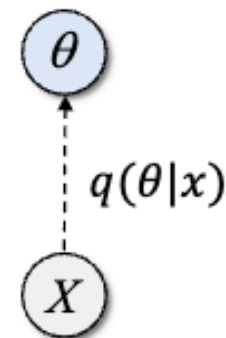
theta: output of previous layers
Phi: the weight matrix of last-layer

(a) DNNs

(b) Extend DNNs with $\theta$

(c) Generative Model of BNDL

(d) Inference Model of BNDL

Gamma distribution: non-linearity and non-negativity (also conjugate).

$$y_j \mid \theta_j \sim \text{Categaory} \left( \theta_j, \Phi \right),$$

$$\theta_j | x_j \sim \Gamma(1,1),$$

$$\Phi \sim \Gamma(1,1).$$

This choice allows for a more comprehensive capture of complex relationships within the model, while also accommodating the characteristics of sparse non-negative variables, thereby enhancing the disentanglement and interpretability of the learned representations $\theta$ and decision layers $\Phi$.

About learning the distributions:
1. Weibull approximate posterior instead of direct Gamma distribution
2. Variational Inferecen to estimate poesterir
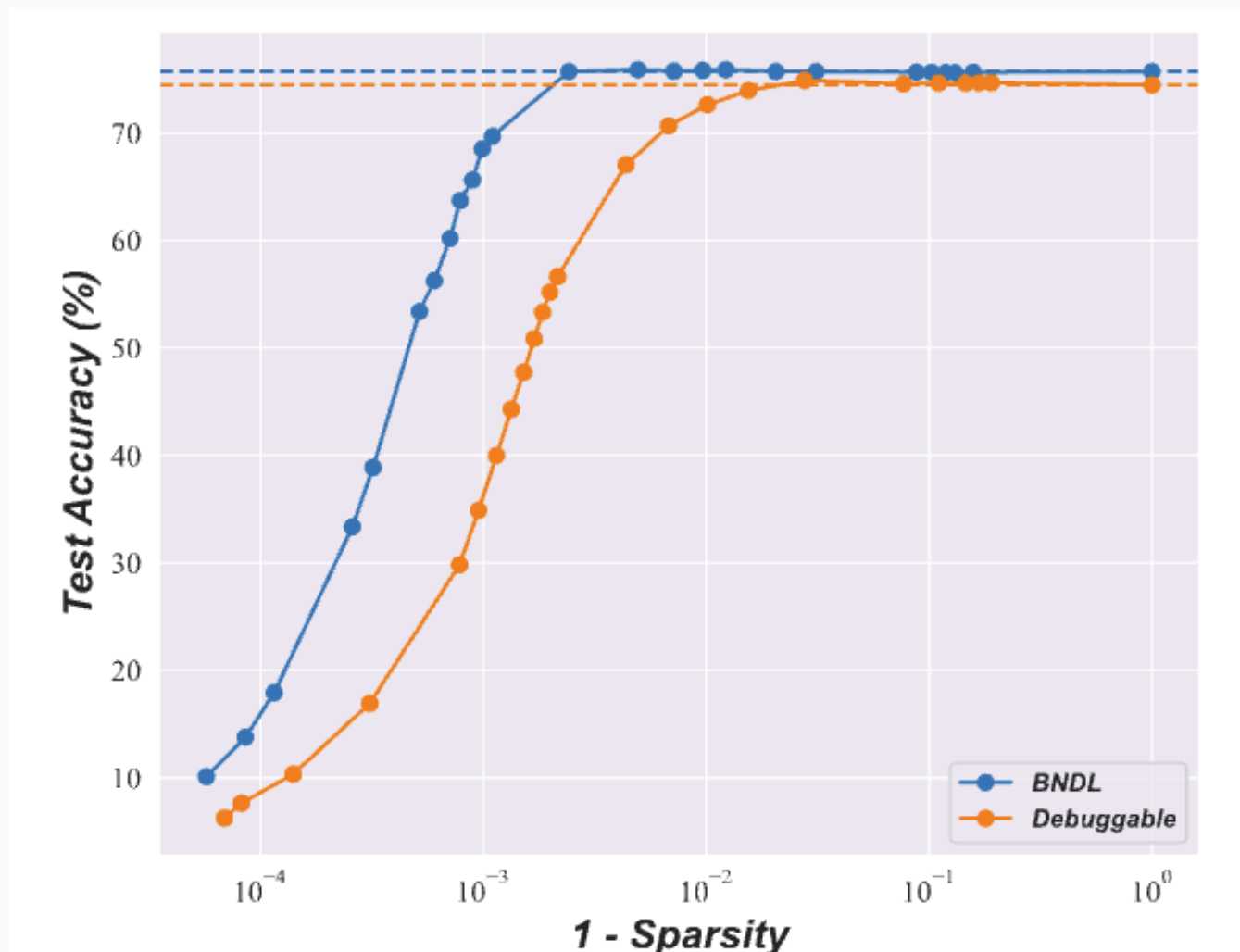
1. Uncertainty Quantification

   Like traditional single-model ways, we just sample many times to get an ensemble of results to computer uncertainty. But this model allows to only pass 1 forward process and get the distribution of $\theta$ and $\Phi$, we just sample from them to generate multiple resutls.

2. The behavior of $\theta$ and $\Phi$

   $\theta$: the output of previous layers; $\Phi$ the weight matrix of last layer.

   - **$\Phi$: A learned map: feature $\rightarrow$ Class**. If the model is well-trained, then the variance of $\Phi$'s every column should be small. If the model has a high variance of a column, then it means the model does not have enough knowledge of the corresponding class.
   - **$\theta$: the output of previous layers, $\theta_{f(x)}$.**
     - *in-distribution and high quality data*: more sharp at the mean and variance is small
     - *ood data*: $\theta$ is flat and has a high variance
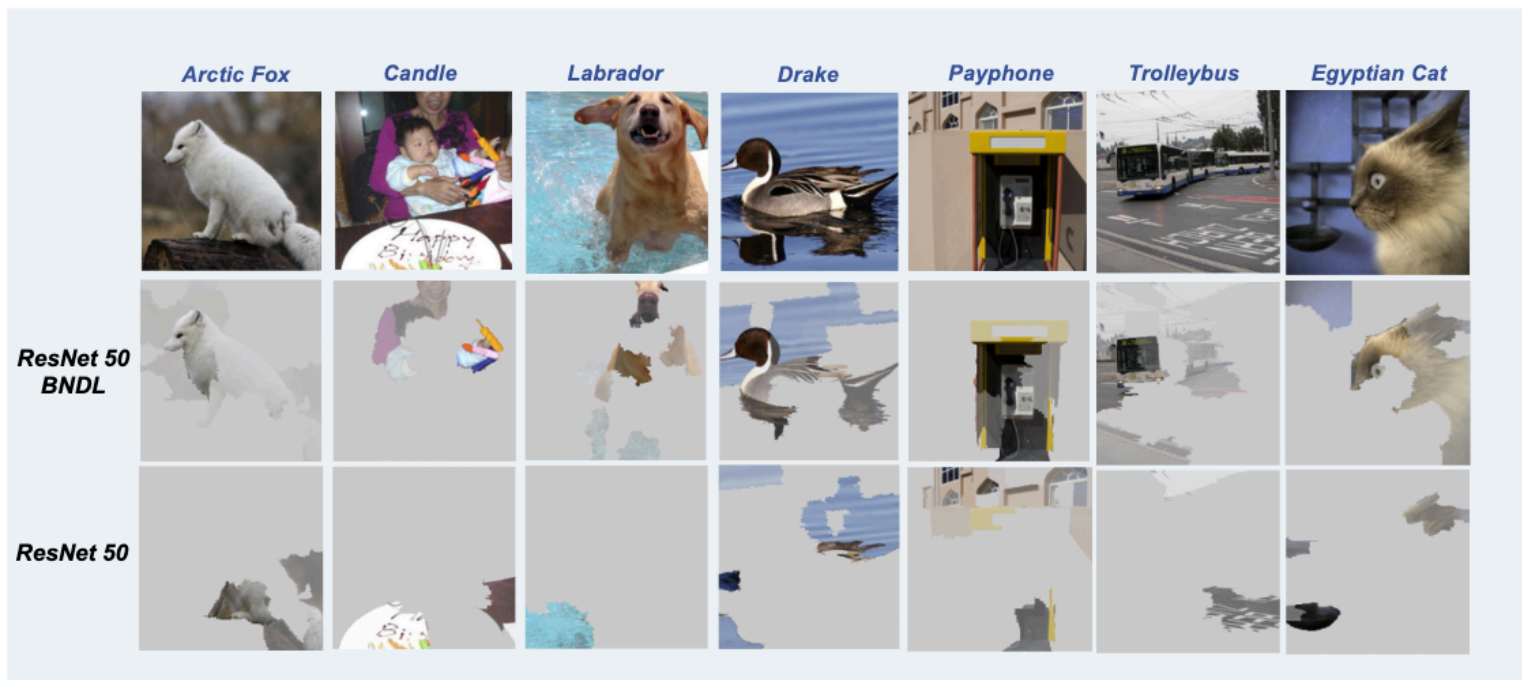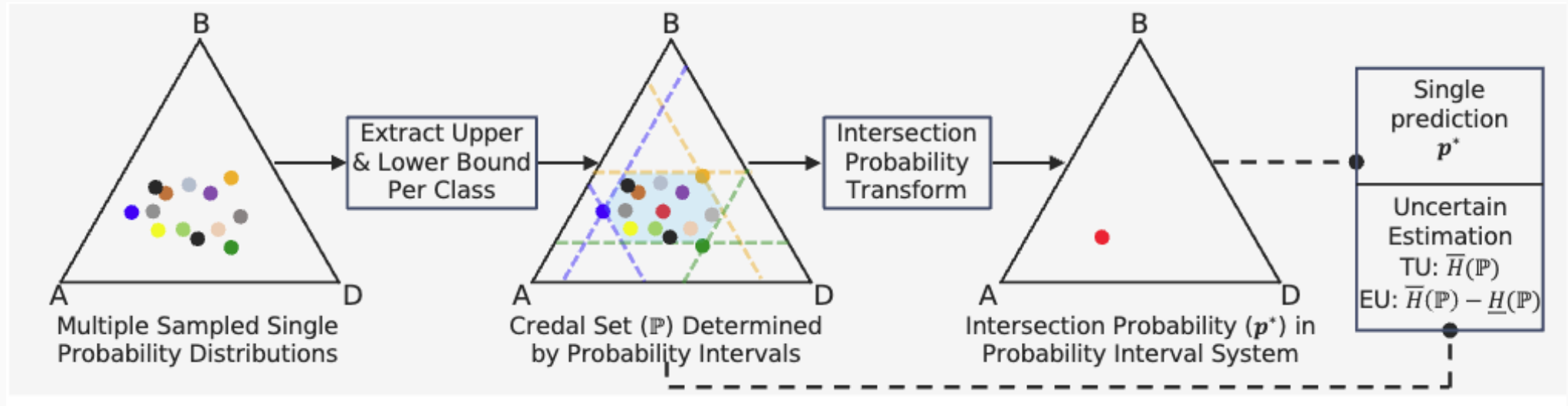     - *low-quality data*: $\theta$ will have higher variance than id data

Figure 4: The LIME visualizations for BNDL and ResNet-50, focusing on the largest $\theta$ for each image, show that BNDL's features align more closely with the semantic meaning of true labels, suggesting more disentangled representations. As we selected the top-10 super-pixels for visualization, the results may include some less significant super-pixels; this issue is alleviated when we reduce the number of top-$k$ super-pixels.

(ICLR 2025) Credal Wrapper of Model Averaging for Uncertainty Estimation on Out-Of-Distribution Detection

# 2. Credal Wrapper of Model Averaging for Uncertainty Estimation on OOD Detection

**Generation of Credal Sets:**

$$p_{U_k} = \max_{n=1,\ldots,N} p_{n,k} \text{ and } p_{L_k} = \min_{n=1,\ldots,N} p_{n,k}$$

**Credal Sets.**

$$\mathbb{P} = \left\{ \boldsymbol{p} \mid p_k \in \left[ p_{L_k}, p_{U_k} \right], \forall k = 1, 2, ..C \right\} s.t. \sum_{k=1}^{C} p_{L_k} \leq 1 \leq \sum_{k=1}^{C} p_{U_k}.$$

**Uncertainty in Credal Nets.**

- Total Uncertainty: $\overline{H}(\mathbb{P}) = \max \sum_{k}^{C} -p_k \log p_k$

- Aleatoric Uncertainty: $\underline{H}(\mathbb{P}) = \min \sum_{k}^{C} -p_k \log p_k$

- Epistemic Uncertainty: $\text{EU} = \overline{H}(\mathbb{P}) - \underline{H}(\mathbb{P})$

where $\sum_{k}^{C} p_k = 1$.

### 2.2.1 Computational Complexity Reduction Method

Key idea: consider the unimportant classes together

---

**Algorithm 1** Probability Interval Approximation

**Input:** $[p_{L_k}, p_{U_k}]$ for $k=1,..C$; Intersection probability $\boldsymbol{p}^*$; Chosen number of classes $J$
**Output:** Approximated probability intervals $[r_{L_j}, r_{U_j}]$ for $k=1,..J$
\# Index vector for sorting $\boldsymbol{p}^*$ in descending order
$\boldsymbol{m} \leftarrow \mathrm{argsort}(\boldsymbol{p}^*)$
\# Upper and lower probability per selected class
$r_{L_j} \leftarrow p_{L_{m_j}}, r_{U_j} \leftarrow p_{U_{m_j}}$ for $j=1,...,J-1$
\# Upper and lower probability for deselected classes
$r_{L_J} \leftarrow \max(1-\sum_{i=m_J}^{m_C} p_{U_i}, \sum_{j=1}^{J-1} r_{L_j}); r_{U_J} \leftarrow \min(1-\sum_{i=m_J}^{m_C} p_{L_i}, \sum_{j=1}^{J-1} r_{U_j})$

---

### 2.2.2 Intersection Probability

Key idea: to get a overall predictive probability from ensemble models, and do not use the average directly.

Normally, we know the true probability has a following expression:

$$p_k^* = p_{L_k} + \alpha_k \cdot \left( p_{U_k} - p_{L_k} \right) \quad \forall k = 1, ..., C.$$

The author wants to find a unified $\alpha$ which can be applied to every class:

$$p_k^* = p_{L_k} + \alpha \cdot \left( p_{U_k} - p_{L_k} \right) \quad \forall k = 1, ..., C.$$
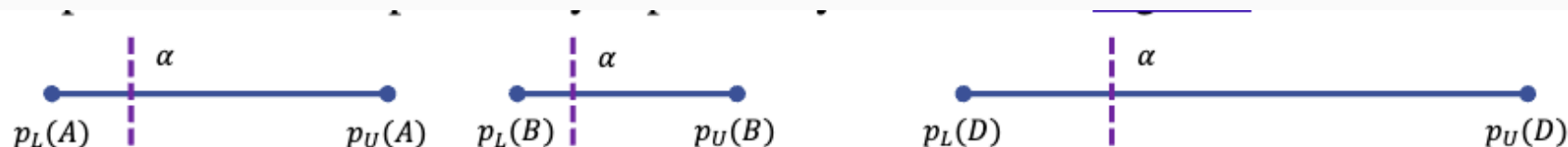
Figure 2: An illustration of the intersection probability for a probability interval system of three classes {A, B, D}[13;12].

$$p_k^* = p_{L_k} + \alpha \cdot \left( p_{U_k} - p_{L_k} \right) \quad \forall k = 1, ..., C.$$

By just confirming that $\sum_{k}^{C} p_k^* = 1$, we can get that $\alpha = \dfrac{1 - \sum_{k}^{C} p_{L_k}}{\sum_{k}^{C} \left( p_{U_k} - p_{L_k} \right)}.$

–

My idea: this way is problematic, because if $\alpha$ is same, then the model tends to have higher prediction on the unsure classes!
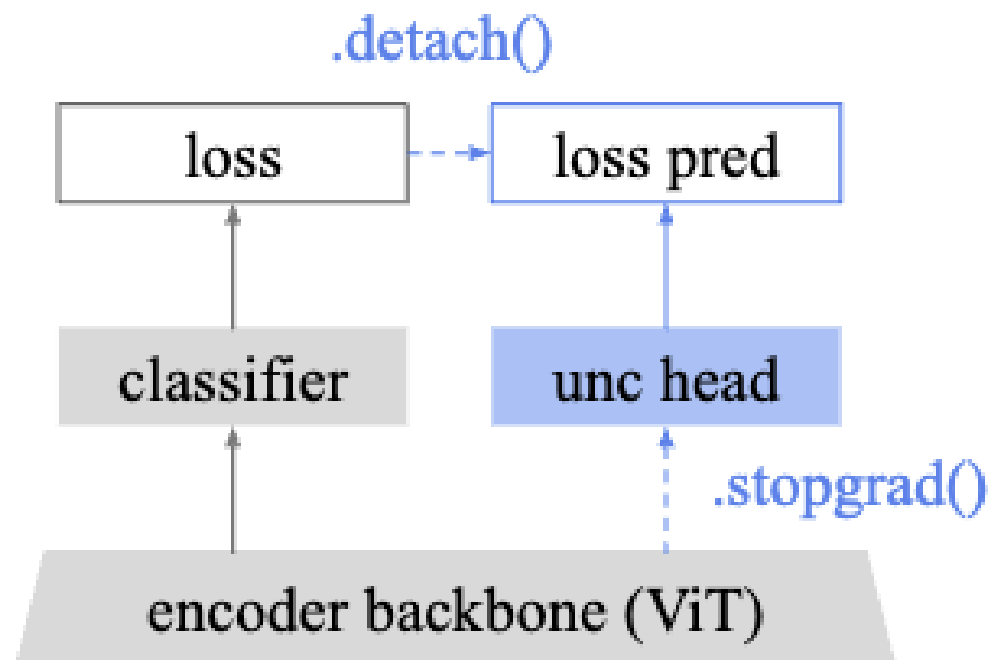
# Pretrained Visual Uncertainties

# 3. Pretrained Visual Uncertainties

- Not real "Uncertainty" — they consider the loss as the "uncertainty"

- Strength: they can train the *uncertainty head* from the trained backbone NN, and the *uncertainty head* can transfer zero-shot.

Figure 2. Pretrained uncertainties are returned by an auxiliary head (blue) that is trained to predict the classification loss of each image.

If trained jointly, the loss function:

$$L = L_{\text{task}}(y, f(x)) + \left(u(e(x)) - L_{\text{task}}^{\text{det}}(y, f(x))\right)^2.$$

The Uncertainty Head module $u$ is implemented as a small MLP, and $e(x)$ is the model representation(feature) outputed by the backbone.

If trained separately , the loss function:

$$L = \left(u(e(x)) - L_{\text{task}}^{\text{det}}(y, f(x))\right)^2,$$

which can be optimized efficiently because we only optimize a 3-layer MLP.

With the current $L_2$ loss, the uncertainty head module can match the scale of the current backbone. However, a downstream user might switch to a different loss on a different scale ∘
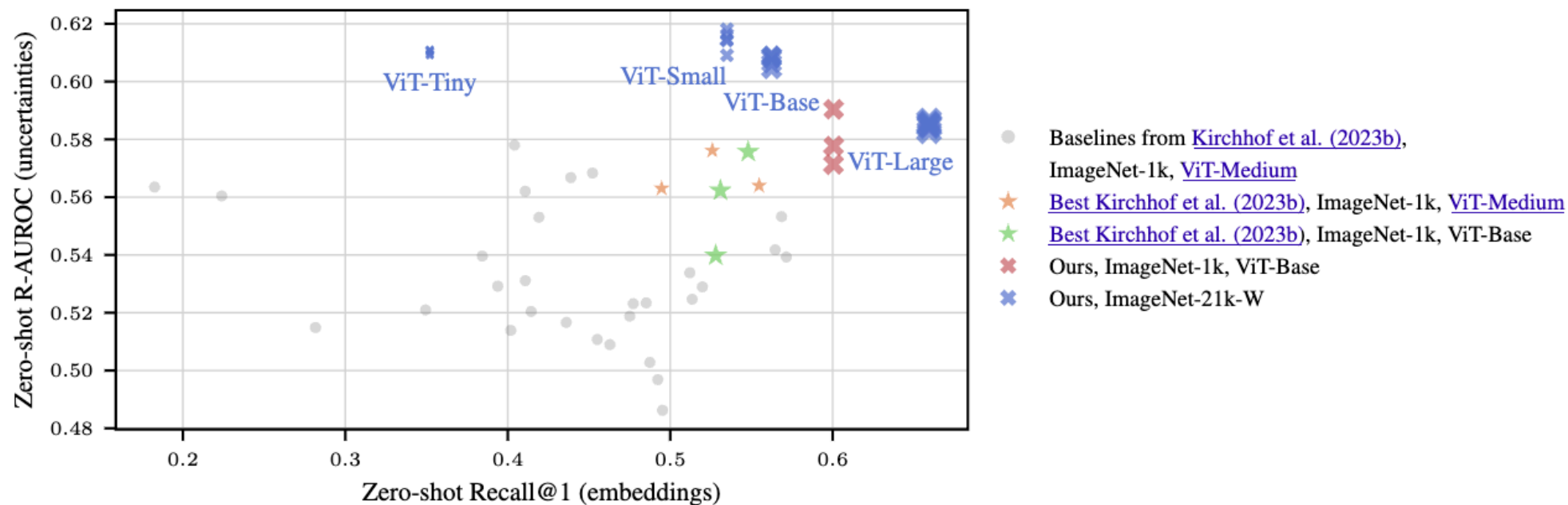
Solution: Ranking-based objective

$$L = \max(0, \mathbb{1}_L \cdot (u(e(x_1)) - u(e(x_2)) + m))$$

$$s.t. \mathbb{1}_L := \begin{cases} +1 \ , \ \text{if } L(x_1) > L(x_2) \\ -1 \ , \ \text{else} \end{cases}.$$

For every pair of images $x_1$ and $x_2$, the indicator function compares which images has the higher primary task loss. The the uncertainty $u$ of that sample is forced to be higher than that of the other sample, by a margin of at least $m = 0.1$.
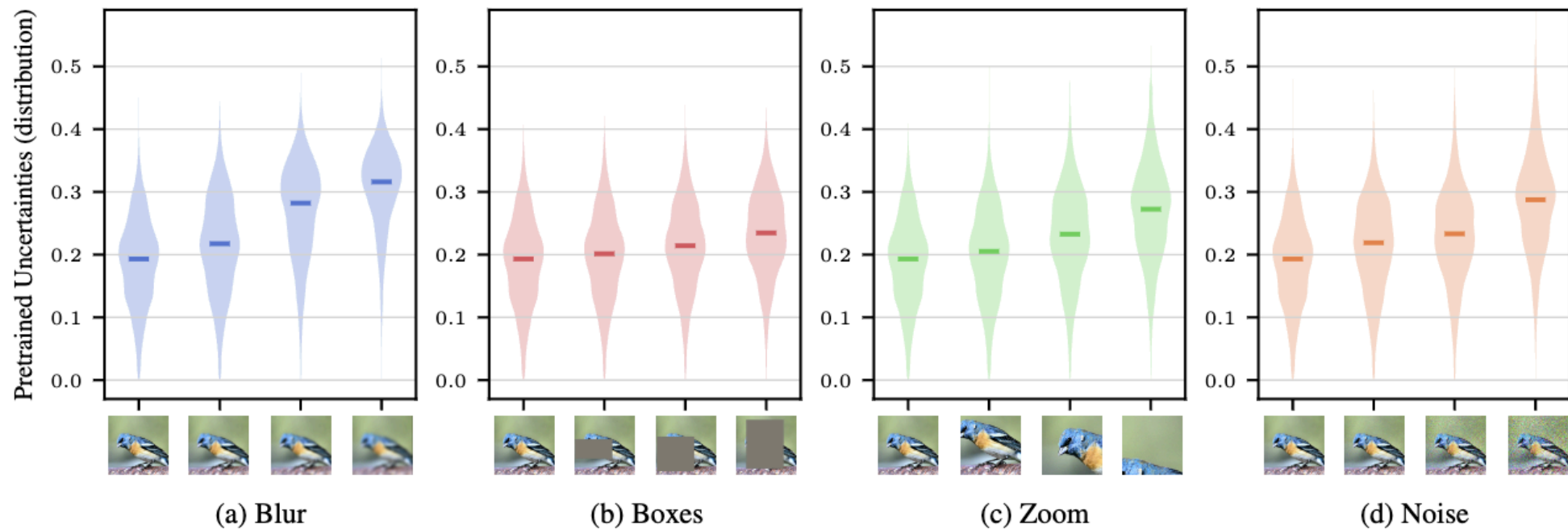
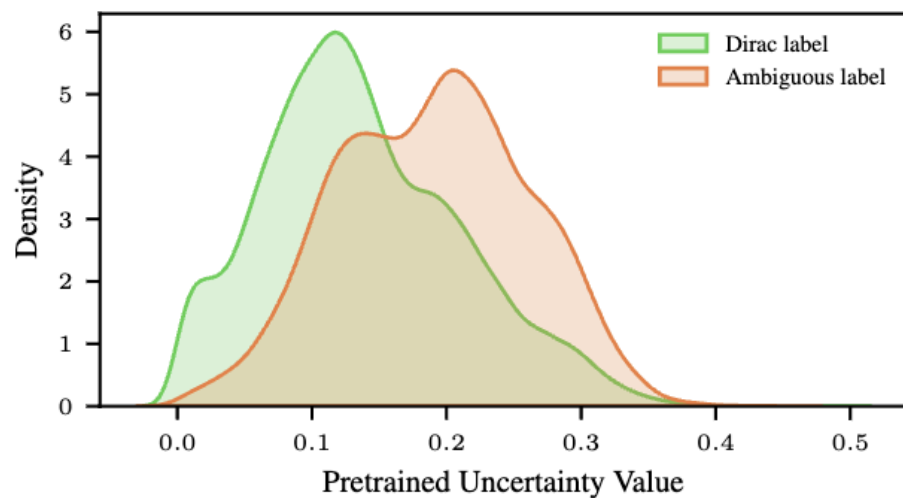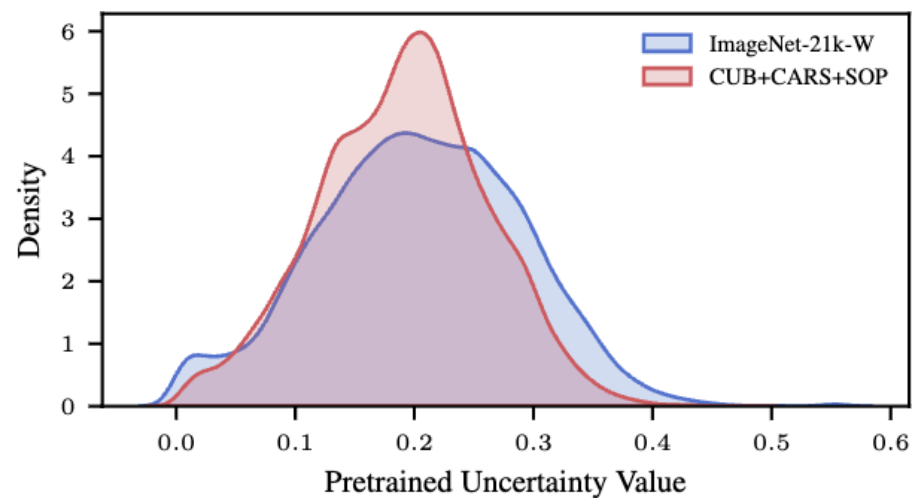Implementation: Sampling pairs from a mini-batch, and summing up all the ranking-based loss

Figure 6. Pretrained uncertainties grow as images are deteriorated. Distributions and medians over unseen datasets (CUB, CARS, SOP). Note that except zooming, the pretrained model is not exposed to these data augmentations during training.

*Figure 7*. Pretrained uncertainties are systematically higher for ImageNet ReaL-H images with multiple possible labels.

*Figure 8*. Pretrained uncertainties are consistent between train and unseen datasets, indicating the absence of epistemic uncertainty.

If we have a model (single model or ensemble) which can do uncertainty quantification, for every training $x_i$, we have the uncertainty $\text{unc}(x_i)$, then we can train a real "uncertainty head":

$$L = (u(x_i) - \text{unc}(x_i))^2,$$

while we can also use the ranking-based objective to confirm the transferability.