

Learning to Defer with an Uncertain Rejector via Conformal Prediction

Anonymous authors

Paper under double-blind review

Abstract

Learning to defer (L2D) allows prediction tasks to be allocated to a human or machine decision maker, thus getting the best of both’s abilities. This allocation decision crucially depends on a ‘rejector’ function. In practice, the rejector could be poorly fit or otherwise misspecified. In this work, we perform uncertainty quantification for the rejector sub-component of the L2D framework. We use conformal prediction to allow the rejector to output prediction sets or intervals of a user-defined confidence level (with distribution-free guarantees), instead of just the binary outcome of ‘defer’ or not. On tasks ranging from image to hate speech classification, we demonstrate that the uncertainty in the rejector translates to safer decisions via two forms of selective prediction.

1 Introduction

Learning-to-Defer (L2D) (Madras et al., 2018; Mozannar & Sontag, 2020) is a framework for human-AI collaboration that divides responsibility between machine and human decision makers. For every test instance, a ‘rejector’ function decides if the case should be passed to either a human or model (but not both). The rejector can be seen as a meta-classifier that determines how to assign responsibility based on which decision maker (human or machine) is more likely to make the correct prediction. While L2D systems offer the promise of improved safety and robustness—by having a human available for support—this promise critically depends on the rejector’s performance. Being a predictive model itself, the rejector is susceptible to the usual failure modes, such as distribution shift between training and test data. Yet, unlike with traditional predictive models, there is an extra point of failure in that the distribution of the human’s predictions can also shift (Tailor et al., 2024).

In this paper, we perform principled uncertainty quantification (UQ) for the rejector sub-component of L2D systems. Specifically, we use the framework of *conformal prediction* (CP) (Vovk et al., 2005) to allow the rejector to output sets or interval, instead of just a binary outcome (defer or not). This allows the rejector to express its uncertainty about whether the human or machine should be assigned to make the decision. In turn, this unlocks new abilities for the L2D system. For example, if the rejector is unsure about to whom responsibility should be assigned, the system can simply abstain from making any prediction. This setting would be useful in cases of distribution shift, among others, since we cannot be sure if either can handle the new, shifted data. Alternatively, an uncertain rejector could mean that we should query *both* the human and model for their predictions. If the human and model agree on their prediction, then this is a good indication that that prediction is reliable (i.e. ensembling).

We explore these novel L2D workflows in experiments on tasks ranging from image to hate speech classification. We find that introducing UQ into L2D allows for safe alternative behaviors (e.g. abstention or consensus checking, as mentioned above) that prevent the L2D system from otherwise returning the wrong prediction. We also study various L2D learning objectives, parameterizations, and CP formulations, finding that the one-vs-all parameterization tends to result in better downstream performance (e.g. accuracy) but at the cost of sometimes having too small of uncertainty sets and in turn under-covering the true label. In summary, we make the following contributions:

- **Distribution-Free UQ for the L2D Rejector** (Sections 3 - 3.1): We are the first to formulate a UQ problem for the L2D deferral decision. Moreover, we are the first to apply conformal prediction to the rejector sub-model of L2D, thus providing distribution-free coverage guarantees on expert correctness.¹
- **Novel L2D Workflows** (Section 3.2 - 3.3): We propose four novel, alternative workflows for L2D that operate via (i) abstention, (ii) checking for consensus between expert and model predictions, (iii) preferring to query the model when the human is uncertain (for cost saving), and (iv) preferring to query the human under distribution shift.

2 Background

We first review the necessary background information on L2D and conformal prediction.

2.1 Learning to Defer

Setting, Data, and Model We focus on multiclass L2D (with one expert) (Madras et al., 2018; Mozannar & Sontag, 2020), though the ideas are presented can straightforwardly generalize to L2D-based regression (Zaoui et al., 2020). Let \mathcal{X} denote the feature space and \mathcal{Y} the label space, a categorical encoding of $K \in \mathbb{N}^{\geq 2}$ classes. Let $\mathbf{x}_n \in \mathcal{X}$ denote a feature vector, and $y_n \in \mathcal{Y}$ denotes the associated class index. L2D assumes that we have access to human predictions, denoted $m_n \in \mathcal{Y}$ for the associated feature vector \mathbf{x}_n . The training data then includes the features, the true label, and the human’s prediction: $\mathcal{D} = \{\mathbf{x}_n, y_n, m_n\}_{n=1}^N$. The human is assumed to have some skills at the prediction task but is not an oracle. For example, the feature vector could be a medical image, m_n is the expert’s diagnosis from looking at the image, and y_n is a true label that can only be obtained from a biopsy. L2D also assumes that the human has access to background knowledge that the classifier does not, such as years of medical training in the aforementioned example. The L2D framework requires two sub-models: a classifier and a rejector (Cortes et al., 2016b;a). We denote the *classifier* as $h : \mathcal{X} \rightarrow \mathcal{Y}$ and the *rejector* as $r : \mathcal{X} \rightarrow \{0, 1\}$. When $r(\mathbf{x}) = 0$, the classifier makes the decision, and when $r(\mathbf{x}) = 1$, the classifier abstains and defers the decision to the human. Thus the rejector can be thought of as a ‘meta-classifier,’ predicting which *predictor* would most likely be correct in its prediction.

Learning Learning in L2D requires us to fit both the rejector and classifier. We assume that whoever makes the prediction—model or human—incurs a loss of zero (correct) or one (incorrect). Using the rejector to toggle between the human and model, we have the overall classifier-rejector loss:

$$L_{0-1}(h, r) = \mathbb{E}_{\mathbf{x}, y, m} [(1 - r(\mathbf{x})) + r(\mathbf{x}) \mathbb{I}[m \neq y]] \quad (1)$$

where $\mathbb{I}[h(\mathbf{x}) \neq y]$ denotes an indicator function that checks if the prediction and label are equal. Minimizing this loss results in the Bayes optimal classifier and rejector:

$$h^*(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \mathbb{P}(y = y|\mathbf{x}), r^*(\mathbf{x}) = \mathbb{I} \left[\mathbb{P}(m = y|\mathbf{x}) \geq \max_{y \in \mathcal{Y}} \mathbb{P}(y = y|\mathbf{x}) \right] \quad (2)$$

where $\mathbb{P}(y|\mathbf{x})$ is the probability of the label under the data generating process, and $\mathbb{P}(m = y|\mathbf{x})$ is the probability that the expert is correct. The assumption that the expert has additional knowledge is what allows it to possibly outperform the Bayes optimal classifier.

Surrogate Losses Several consistent surrogate losses have been proposed for Equation 1 (Mozannar & Sontag, 2020; Verma & Nalisnick, 2022; Mao et al., 2024c;b; Cao et al., 2023; Charusaie et al., 2022). For our implementation, we focus on the two surrogates that have demonstrated the ability to learn calibrated predictors in practice—since the more calibrated the predictor, the better the conformal prediction results will be. Specifically, we use Verma & Nalisnick (2022)’s one-vs-all (OvA) parameterization and Cao et al. (2024)’s assymetric softmax (A-SM) parameterization. These parameterizations assume the classifier and

¹An abbreviated version of this paper appeared in the non-archival proceedings of [workshop redacted to preserve anonymity].

rejector are unified via an augmented label space: $\mathcal{Y}^\perp = \mathcal{Y} \cup \{\perp\}$, where \perp denotes the rejection option. Then let $g_k : \mathcal{X} \mapsto \mathbb{R}$ for $k \in [1, K]$ where k denotes the class index, and let $g_{K+1} : \mathcal{X} \mapsto \mathbb{R}$ denote the rejection (\perp) option. The g functions are analogous to the logits of a neural-network-based classifier. The OvA surrogate loss is given as (Verma & Nalisnick, 2022):

$$\begin{aligned} \psi_{\text{OvA}}(g_1, \dots, g_{K+1}; \mathbf{x}, y, m) = & \phi[g_y(\mathbf{x})] + \sum_{y' \in \mathcal{Y}, y' \neq y} \phi[-g_{y'}(\mathbf{x})] + \phi[-g_{K+1}(\mathbf{x})] \\ & + \mathbb{I}[m = y] (\phi[g_{K+1}(\mathbf{x})] - \phi[-g_{K+1}(\mathbf{x})]) \end{aligned} \quad (3)$$

where $\phi : \{\pm 1\} \times \mathbb{R} \mapsto \mathbb{R}_+$ is a binary surrogate loss. For instance, when ϕ is the logistic loss, we have $\phi[f(\mathbf{x})] = \log(1 + \exp\{-f(\mathbf{x})\})$.

The A-SM surrogate loss is defined as follows (Cao et al., 2024):

$$\begin{aligned} \psi_{\text{A-SM}}(g_1, \dots, g_{K+1}; \mathbf{x}, y, m) = & -\log \phi_{\text{A-SM}}(g(\mathbf{x}), y) - \mathbb{I}[m \neq y] \cdot \log(1 - \phi_{\text{A-SM}}(g(\mathbf{x}), K+1)) \\ & - \mathbb{I}[m = y] \cdot \log \phi_{\text{A-SM}}(g(\mathbf{x}), K+1) \end{aligned} \quad (4)$$

$$\text{where } \phi_{\text{A-SM}}(g(\mathbf{x}), y) = \begin{cases} \frac{\exp(g_y(\mathbf{x}))}{\sum_{y'=1}^K \exp(g_{y'}(\mathbf{x}))} & \text{if } y < K+1, \\ \frac{\exp(g_{K+1}(\mathbf{x}))}{\sum_{j=1}^{K+1} \exp(g_j(\mathbf{x})) - \max_{y' \in \mathcal{Y}} \exp(g_{y'}(\mathbf{x}))} & \text{otherwise.} \end{cases}$$

Here, the ‘asymmetry’ is due to the softmax having different terms in the denominator for the class and rejector terms. The symmetric softmax parameterization (Mozannar & Sontag, 2020) has the same denominator for both terms, which leads to issues for estimating the expert’s correctness probability in practice (Verma & Nalisnick, 2022; Cao et al., 2024). For both parameterizations, at test time, the classifier is obtained by taking the maximum over g functions: $\hat{y} = h(\mathbf{x}) = \arg \max_{k \in [1, K]} g_k(\mathbf{x})$. The rejection function is given as:

$$r(\mathbf{x}) = \mathbb{I}[g_{K+1}(\mathbf{x}) \geq \max_k g_k(\mathbf{x})].$$

Expert Correctness Both OvA and A-SM parameterization compute the probability that the expert is correct: For the OvA parameterization, this probability is directly parameterized by the $(K+1)$ th binary classifier:

$$\hat{p}(m = y | \mathbf{x}) = \phi_{\text{OvA}}[g_{K+1}(\mathbf{x})] = \frac{1}{1 + \exp\{-g_{K+1}(\mathbf{x})\}}. \quad (5)$$

A-SM similarly uses the deferral score, but here the parameterization requires evaluating all $K+1$ functions:

$$\hat{p}(m = y | \mathbf{x}) = \phi_{\text{A-SM}}(g(\mathbf{x}), K+1) = \frac{\exp(g_{K+1}(\mathbf{x}))}{\sum_{j=1}^{K+1} \exp(g_j(\mathbf{x})) - \max_{y' \in \mathcal{Y}} \exp(g_{y'}(\mathbf{x}))} \quad (6)$$

Both estimators have been shown to be competitively calibrated when trained by empirical risk minimization and without relying upon post-hoc procedures such as temperature scaling (though they could be employed as well) (Cao et al., 2024).

2.2 Conformal Prediction

Conformal prediction (CP) is a model-agnostic, distribution-free approach to uncertainty quantification with finite-sample guarantees (Shafer & Vovk, 2008). Given a test-time feature vector \mathbf{x}_{N+1} , CP seeks to construct a prediction set $C(\mathbf{x}_{N+1}; \tau) \subseteq \mathcal{Y}$ such that the true label y_{N+1} is included with probability $1 - \alpha$: $\mathbb{P}(y_{N+1} \in C(\mathbf{x}_{N+1}; \tau)) \geq 1 - \alpha$, for $\alpha \in [0, 1]$. τ is a parameter that controls the set size, as will be described below. This statement is a *marginal* guarantee, meaning that it will hold, on average, over test samples but will not necessarily hold for any particular sample. CP’s aforementioned guarantee is built off the crucial assumption that the test data is drawn exchangeably with a calibration set.

To compute the parameter τ that controls the prediction sets, the *split*-CP (a.k.a. *inductive* CP) algorithm (Papadopoulos et al., 2002) is a popular choice due to its computational and sample efficiency (Fang & Bellotti, 2024) and resemblance to the traditional workflow of hyperparameter tuning. Split-CP requires τ be fit to a held-out calibration set \mathcal{D}_2 , which must be drawn exchangeably with the test set for the CP coverage guarantee to hold. Given a classifier already trained on the training set \mathcal{D}_1 , its softmax outputs are denoted $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_K(\mathbf{x})]$. CP then requires a score function be chosen that quantifies how well the model’s prediction conforms to the true label. Using the softmax confidence associated with the true label is a reasonable choice: $s(\mathbf{x}, y; \mathbf{f}) = 1 - f_y(\mathbf{x})$, where $f_y(\mathbf{x})$ is the softmax score for the true label. Others exist that incorporate all dimensions that have higher confidence than the true label (Romano et al., 2020). Split-CP then proceeds by evaluating $s(\mathbf{x}, y; \mathbf{f})$ on all points in the held-out set and setting $\hat{\tau}$ to be the $(1 - \alpha)$ quantile (with a finite-sample correction) of the empirical distribution of scores. For a test time point \mathbf{x}_{N+1} , the prediction set is constructed as:

$$C(\mathbf{x}_{N+1}) = \{j | f_j(\mathbf{x}_{N+1}) \geq 1 - \hat{\tau}\},$$

which represents the softmax dimensions that outscore the threshold $1 - \hat{\tau}$. CP is commonly evaluated by checking that the desired coverage level is achieved in practice while also having efficient set sizes. The latter is crucial since the CP guarantee is trivially met by choosing $C(\mathbf{x}_{N+1}; \tau) = \mathcal{Y}$ for $(1 - \alpha)\%$ of cases.

3 Uncertain Deferral via Conformal Prediction

We will now apply the CP framework to quantify the uncertainty in the rejector sub-component of an L2D system. Concretely, instead of just outputting 0 (model) or 1 (human), we want the CP-based rejector to output a set $C_r(\mathbf{x}; \tau)$, which is an element of the superset $\{\{0\}, \{1\}, \{0, 1\}\}$. $C_r(\mathbf{x}; \tau) = \{0, 1\}$ means that the rejector is unsure if the decision should be allocated to the human or model. Thus, instead of *prediction* sets, we call the uncertainty set of the rejector a *deferral set*. In Section 3.2, we will discuss how to incorporate these sets into downstream decision making.

Ideal Construction Recalling the Bayes optimal decision rule for the rejector (Equation 2), it would be ideal if $C_r(\mathbf{x}; \tau)$ could satisfy the guarantee:

$$\mathbb{P}(r^*(\mathbf{x}_{N+1}) \in C_r(\mathbf{x}_{N+1}; \tau)) \geq 1 - \alpha,$$

which means that, marginally, the probability that the output of the Bayes optimal rejector is in the set is at least $1 - \alpha$. Constructing an adaptive set via validation statistics, unfortunately, requires that we be able to compare $\mathbb{P}(m = y | \mathbf{x})$ vs $\mathbb{P}(y | \mathbf{x})$ to compute a non-conformity score. This comparison requires high-fidelity estimates of two conditional probabilities, and obtaining estimates of such one-off events is known to be impossible (Roth et al., 2023). The only work-around is if we observe multiple samples of both the label y and expert prediction m (Johnson et al., 2024), which would only drastically increase the already high supervision burden of L2D. Thus, we leave this construction as an open problem and turn to a more practical alternative below.

Practical Construction We instead consider constructing the set to capture an alternative quantity: $\mathbb{I}[m_{N+1} = y_{N+1}]$, an indicator function representing if the human will make the correct prediction. Similarly, we wish to construct prediction sets such that this binary variable will have a coverage guarantee:

$$\mathbb{P}(\mathbb{I}[m_{N+1} = y_{N+1}] \in C_r(\mathbf{x}_{N+1}; \tau)) \geq 1 - \alpha. \quad (7)$$

This statement is not equivalent to the one above since the expert could be correct (i.e. $\mathbb{I}[m_{N+1} = y_{N+1}] = 1$) but $\mathbb{P}(y | \mathbf{x})$ still be a better predictive model (i.e. $r^*(\mathbf{x}) = 0$). In other words, this formulation is considering the expert’s performance in isolation of the classifier’s. However, the high-level semantics are retained since $C_r(\mathbf{x}_{N+1}; \tau) = \{0\}$ means that the expert will likely be wrong and so using the classifier is either a good decision or not an inferior one (if the model would also be wrong). Conversely, $C_r(\mathbf{x}_{N+1}; \tau) = \{1\}$ means that the expert will likely make the correction prediction. If $C_r(\mathbf{x}_{N+1}; \tau) = \{0, 1\}$, then the prediction set is unsure if the expert will be correct and still suggests uncertainty in the deferral decision. This relaxation, importantly, allows us to define a practical conformity statistic.

3.1 Constructing Deferral Sets

We can construct deferral sets that follow the guarantee in Equation 7 by treating the deferral decision as a binary classification problem: whether the expert will make the correct prediction or not. Following CP as it is usually applied to binary classification, we construct the conformity score using the binary probabilities given in Equation 5 for OvA and Equation 6 for A-SM. To obtain the threshold $\hat{\tau}$, we follow the standard procedure of split-CP by computing these non-conformity scores on a held-out calibration set (Angelopoulos et al., 2023), obtaining the $(1 - \alpha)$ empirical quantile, and then applying the threshold at test time as follows:

$$C_r(\mathbf{x}; \hat{\tau}) = \begin{cases} \{0\} & \text{if } 1 - \hat{p}(m = y|\mathbf{x}) \geq 1 - \hat{\tau} \\ \{1\} & \text{if } \hat{p}(m = y|\mathbf{x}) \geq 1 - \hat{\tau} \\ \{0, 1\} & \text{otherwise.} \end{cases} \quad (8)$$

The set $C_r(\mathbf{x}; \hat{\tau})$ should satisfy the coverage guarantee given in Equation 7, assuming the usual assumptions of CP hold, such as exchangeability between calibration and test data.

3.2 Using Deferral Sets in Decision Making

Now that we have detailed how to construct CP deferral sets, we next address how to use them to improve decision making within the L2D framework. While there are surely alternative uses, below we detail three that we believe will be practical and useful in a variety of applications.

Abstention The use that likely first comes to mind is prediction with the option to abstain (Chow, 1957; Cordella et al., 1995; Herbei & Wegkamp, 2006; Hellman, 2007; Geifman & El-Yaniv, 2017). In the traditional case, the classifier only makes a prediction if it is confident; otherwise, it abstains since the consequences of being wrong outweigh the consequences of providing no prediction. This is often appropriate for applications in healthcare: it is better to wait and perform more tests, seek out more opinions, etc. than to give a patient a wrong diagnosis. Our CP deferral sets allow for a similar workflow, but instead of abstaining because the prediction is uncertain, the L2D system will abstain because it is uncertain about to whom to allocate responsibility, the machine or human. Specifically, if $C_r(\mathbf{x}_{N+1}; \hat{\tau}) = \{0, 1\}$, then the L2D system will abstain. Otherwise, the system will defer if $r^*(\mathbf{x}) = 1$. A visualization of this workflow is shown in Figure 1a. As is typically the case with abstention methods, we expect this workflow to improve the system accuracy at the cost of reducing coverage.

Consensus Prediction We next consider how to make a prediction even if $C_r(\mathbf{x}_{N+1}; \hat{\tau}) = \{0, 1\}$. If the rejector is uncertain to defer or not, we propose querying both the model and human for their predictions. If they agree, then that consensus prediction is output as the L2D system’s final prediction. If they do not agree, then the system abstains from making any prediction. This workflow has the same appeal to safety as the abstention-only option, but it will have higher coverage since it will make predictions when the abstention-only workflow would not. This workflow is diagrammed in Figure 1b. We expect this workflow to perform similarly as abstention but with increased coverage, since it can still make predictions even when the deferral set is of maximum size.

Human-Preferred Prediction under Distribution Shift Lastly, we consider cases in which intuition leads us to believe the human is a more robust predictor than the model. Consider the task of image classifier under covariate shift caused by corruption noise (Ovadia et al., 2019). Even low-levels of noise can cause modern classifiers to start to fail, but a human can often be robust to similar noise. Another example of such human superiority is the case of adversarial examples, which—by definition—do not fool humans but fool a classifier. To cover such cases in which safety concerns dictate that the human should have priority, we call our third workflow ‘human-preferred prediction.’ This means that if $C_r(\mathbf{x}_{N+1}; \hat{\tau}) = \{0, 1\}$, the L2D system will still query and return the human’s prediction as the final output, despite the uncertainty. In this case, we are not using the uncertainty set in the same way as the previous two workflows since covariate shift is assumed to be happening. This shift violates the core assumptions of CP, invalidating the coverage guarantee. However, as we will see in the experiments, there is still reason to believe the rejector will be

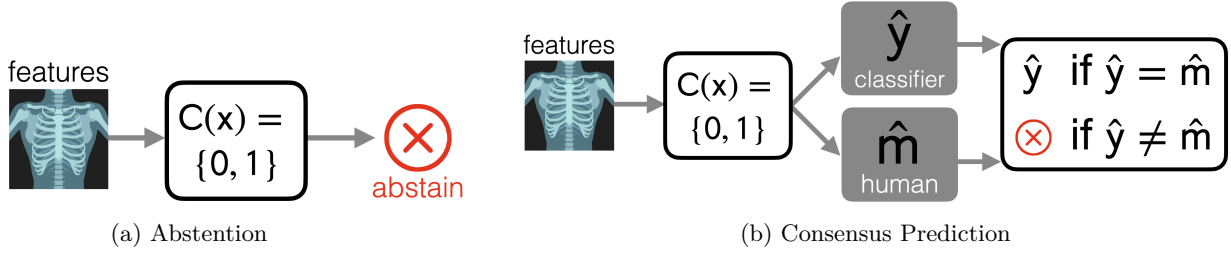


Figure 1: *L2D Decision Making Workflows*. Subfigure (a) shows that the L2D system will *abstain* if the rejector is uncertain. Subfigure (b) shows the alternative workflow in which, if the rejector is uncertain, the system will check for *consensus* between the expert and model predictions and abstain otherwise.

responsive to the shift and reduce its output confidence toward $\hat{p}(m = y|\mathbf{x}) = 0.5$. In turn, the threshold $1 - \hat{\tau}$ will not be met to have a CP set of size one, making the set size of two an indicator of covariate shift.

3.3 Constructing and Using Deferral Intervals

Rather than producing a prediction set for the binary classifier (Section 2.2), we can instead consider an uncertainty interval for the rejector’s confidence itself. This would result in the uncertainty interval of the form $[b_l(\mathbf{x}), b_r(\mathbf{x})]$, such that the conditional endpoints satisfy $b_l(\mathbf{x}) \in [0, 1)$, $b_r(\mathbf{x}) \in (0, 1]$, and $b_l(\mathbf{x}) \leq b_r(\mathbf{x})$. The coverage guarantee would then be $\mathbb{P}(\mathbb{P}(m = y|\mathbf{x}) \in [b_l(\mathbf{x}), b_r(\mathbf{x})]) \geq 1 - \alpha$, where $\alpha \in [0, 1]$ again controls the nominal coverage. Barber (2020) provides an algorithm to compute this interval in practice. Their algorithm can be applied straightforwardly to this case.

Moving on to decision-making, the deferral interval can be used the same way as the deferral set, as outlined in Section 3.2. The only change is that instead of, for example, abstaining when the set is of size two, here we need to specify a certain width that, when exceeded, triggers abstention. One possibility is to simply abstain if the interval is on both sides of 50% (i.e. $0.5 \notin [b_l(\mathbf{x}), b_r(\mathbf{x})]$) and to not abstain when the interval is contained on one side (i.e. $0.5 < b_l(\mathbf{x})$ or $0.5 > b_r(\mathbf{x})$).

Yet unlike with deferral sets, this interval allows for a deferral decision that is close to traditional L2D but can be made more robust. Instead of just comparing the point estimates of the rejector and classifier confidences, we can use the interval and only defer when we are very sure that the human is more likely to be correct than the classifier: **defer** if $b_l(\mathbf{x}) > \max_{y \in \mathcal{Y}} p(y = y|\mathbf{x})$. Doing so will ostensibly save in expert queries since the system will only call the expert when they clearly improves upon the classifier.

While the abstention and consensus prediction workflows in Section 3.2 apply likewise to deferral intervals, we form an additional workflow:

Model-Preferred Prediction for Cost Saving Another reasonable workflow that our UQ procedure allows is for model-preferred prediction. When the uncertainty interval is $[0, 1]$, we are unsure if the expert will be correct, and since experts usually require some expense to query, then we may want to query the expert only when we are sure they will be correct. Otherwise, the L2D system may have just as an acceptable an outcome querying the model, which often requires a negligible cost to query. We expect this workflow to increase the classifier’s coverage while not substantially decreasing overall system accuracy.

Although deferral intervals improved allocation calibration ($\text{ECE } 5.86 \rightarrow 5.17$), the resulting conformal intervals were often too wide to support informative routing decisions at the desired coverage. This reflects a known calibration–efficiency trade-off in conformal inference: marginal coverage is enforced by inflating interval size (Angelopoulos et al., 2023; Barber et al., 2023), especially with split calibration and limited signal, which can yield intervals that are decision-agnostic (e.g., near $[0, 1]$) and thus uninformative for cost-aware thresholds. We therefore report these results in Appendix A.

4 Related Work

The L2D framework (Madras et al., 2018)—along with its precursors (Chow, 1957; Bartlett & Wegkamp, 2008; Yuan & Wegkamp, 2010; Cortes et al., 2016b)—have received much attention of late due to their potential to improve safety via semi-automation (Raghu et al., 2019). The majority of such attention has focused on L2D’s learning objective (Mozannar & Sontag, 2020; Verma & Nalisnick, 2022; Mao et al., 2024c;b; Cao et al., 2023; Charusaie et al., 2022) and its extension to multiple experts (Verma et al., 2023; Tailor et al., 2023; Mao et al., 2024a; Keswani et al., 2021; Hemmer et al., 2022). Only two works have previously considered the uncertainty estimation abilities of the rejector sub-component, with Verma & Nalisnick (2022) first observing the aforementioned pathologies of the symmetric softmax parameterization and Cao et al. (2024) proposing the asymmetric softmax as a remedy. Liu et al. (2022) employed ensembling to estimate the classifier’s uncertainty and used this to inform the deferral decision, but their approach does not model the expert’s abilities nor represent the expert’s uncertainty.

Conformal Prediction for L2D CP has previously been incorporated into L2D and related frameworks. Straitouri et al. (2023) and Babbar et al. (2022) both proposed performing CP for a classifier and then passing the set to a human to choose the label that will be the final prediction. Yet like the aforementioned approached by Liu et al. (2022), the classifier’s uncertainty is being quantified, not the human’s, which is the focus of our methodology. The work of Verma et al. (2023) is more related: they apply CP to multi-expert L2D to quantify the uncertainty in who is the *best* expert of the multiple available. Their coverage guarantee is formulated with the goal of including this best expert in the set. Our approach could be applied to multi-expert L2D, but it would construct a deferral set per expert, not across experts as they do.

5 Experiments

We now experimentally demonstrate that incorporating uncertainty via CP into the deferral decision can have tangible benefits to the safety and robustness of L2D systems. Our experiments follow closely the setup in previous works on L2D (Mozannar et al., 2023; Verma et al., 2023; Cao et al., 2024) while introducing uncertainty quantification for the rejector. We trained L2D models using the OvA and A-SM surrogate losses. Taking this base L2D model, we then apply the CP procedure described in Section 3. Our implementation will be publicly available upon acceptance of the manuscript. See the supplementary materials for additional details about more algorithm detail, training hyperparameters and backbone architectures.

Datasets We utilize three datasets tailored to different tasks: **CIFAR-10** (Krizhevsky et al., 2009) for image classification, **HAM10000** (Tschandl et al., 2018) for skin lesion diagnosis, and **Hate Speech** (Davidson et al., 2017) for hate speech detection. The **CIFAR-10** dataset comprises 60,000 instances, divided into training, calibration, and test sets at 70%, 10%, and 20%, respectively. Similarly, **Hate Speech** contains 25,000 instances, split into the same proportions. The **HAM10000** with 10,015 dermatoscopic images, is divided into 60% training, 20% calibration, and 20% test splits.

Models and Experts We follow previous work’s L2D experimental settings (Mozannar et al., 2023; Verma & Nalisnick, 2022; Verma et al., 2023), including their choice of base model backbones and expert simulations. We apply a three-layer convolutional neural network (CNN) for **CIFAR-10**, a 34-layer residual network (ResNet34) for **HAM10000**, and a linear network and SBERT embedding (Reimers, 2019) for **Hate Speech**. Expert simulations mirror Mozannar et al. (2023): on **CIFAR-10**, an oracle predicts perfectly on the first $k=6$ classes and uniformly at random on the remaining $(10-k)$; on **Hate Speech**, we use a stochastic “random-annotator” baseline; on **HAM10000**, an MLP-Mixer trained on meta-data emulates an expert with contextual information beyond pixels (Tolstikhin et al., 2021).

5.1 Coverage and Efficiency

We first experimentally verify that the target coverage is met, validating CP’s guarantee (Equation 7). In Table 1, we report the empirical coverage and average set size for the three aforementioned datasets. Both parameterizations meet the target coverage level (90%) for all datasets except for OvA on **CIFAR-10** ($\sim 87\%$).

Table 1: *Coverage and Efficiency.* We report mean and standard deviation of the empirical coverage and the average size of the deferral set for a confidence level of $1 - \alpha = 90\%$.

Dataset	Parameterization	Coverage (%)	Average Set Size
CIFAR-10	OvA	86.94 \pm 0.86	1.07 \pm 0.03
	A-SM	90.53 \pm 0.56	1.37 \pm 0.01
HAM100000	OvA	90.65 \pm 0.63	1.25 \pm 0.01
	A-SM	91.13 \pm 0.58	1.28 \pm 0.03
HateSpeech	OvA	90.35 \pm 0.53	1.03 \pm 0.03
	A-SM	90.67 \pm 0.52	1.01 \pm 0.01

In all cases, the sets are quite efficient, with the average set size always being less than 1.3. The exceptionally small set size of 1.07 for OvA on CIFAR-10 leads to its mis-coverage. We suspect the mis-coverage is due to (natural) train-test distribution shift.

5.2 Learning to Defer with Abstention and Consensus Prediction

We next investigate the efficacy of the abstention and consensus decision-making workflows presented in Section 3.2. Table 2 reports the system accuracy, ratio of test points deferred, and the coverage of the system (i.e. the fraction of points for which the system does not abstain) again for CIFAR-10, Hate Speech, and HAM10000. We see that both OvA and A-SM improve upon the accuracy of the base L2D model for CIFAR-10 and HAM10000, with improvements ranging from 2% to 5%. However, the coverage reduction is variable, ranging from modest (-8%) to substantial (-38%), meaning that the accuracy improvement would be practical in some cases (e.g. OvA for CIFAR-10) and not in others (e.g. A-SM for CIFAR-10). On Hate Speech, abstention occurred for very few points, leading to uninteresting accuracy results. We do not see a clear superiority between the parameterizations.

Table 2: *Abstention and Consensus Prediction.* We report mean and standard deviation of system accuracy, fraction of points deferred, and test-set coverage.

	Parameterization	Method	System Accuracy	Fraction Deferred	System Coverage
CIFAR-10	OvA	Base Model	84.71 \pm 0.46	55.26 \pm 1.76	100
		Abstention	86.72 \pm 1.02	56.41 \pm 2.30	92.14 \pm 0.48
		Consensus	86.79 \pm 1.07	56.38 \pm 2.31	93.32 \pm 0.52
	A-SM	Base Model	84.01 \pm 0.45	56.63 \pm 3.73	100
		Abstention	87.05 \pm 0.76	84.13 \pm 4.56	62.53 \pm 0.75
		Consensus	87.58 \pm 0.61	79.62 \pm 4.31	67.57 \pm 0.75
HAM10k	OvA	Base Model	82.1 \pm 0.49	33.71 \pm 2.39	100
		Abstention	87.48 \pm 0.51	35.91 \pm 2.84	75.23 \pm 1.40
		Consensus	85.72 \pm 0.63	34.27 \pm 2.52	88.39 \pm 1.85
	A-SM	Base Model	78.92 \pm 0.29	26.68 \pm 3.07	100
		Abstention	87.05 \pm 0.87	28.11 \pm 3.45	72.82 \pm 1.19
		Consensus	84.76 \pm 0.44	27.49 \pm 3.16	84.48 \pm 0.95
Hate Speech	OvA	Base Model	92.09 \pm 0.07	42.41 \pm 0.99	100
		Abstention	92.28 \pm 0.14	42.48 \pm 0.96	99.38 \pm 0.43
		Consensus	92.25 \pm 0.13	42.42 \pm 0.96	99.78 \pm 0.22
	A-SM	Base Model	91.82 \pm 0.32	67.91 \pm 1.76	100
		Abstention	91.88 \pm 0.15	67.79 \pm 1.74	99.16 \pm 0.75
		Consensus	91.88 \pm 0.12	67.81 \pm 1.73	99.65 \pm 0.28

5.3 Learning to Defer under Covariate Shift

To evaluate three workflows under covariate shift, we induce out-of-distribution (OOD) shift with a six-level severity index (1–6), where levels 1–5 increase smoothly. We then define level 6 as the *extreme-shift* condition across all datasets. The transition from level 5 to 6 represents a markedly larger distributional change than the preceding increments. On CIFAR-10, we utilized the brightness corruption subset of CIFAR-10-C for

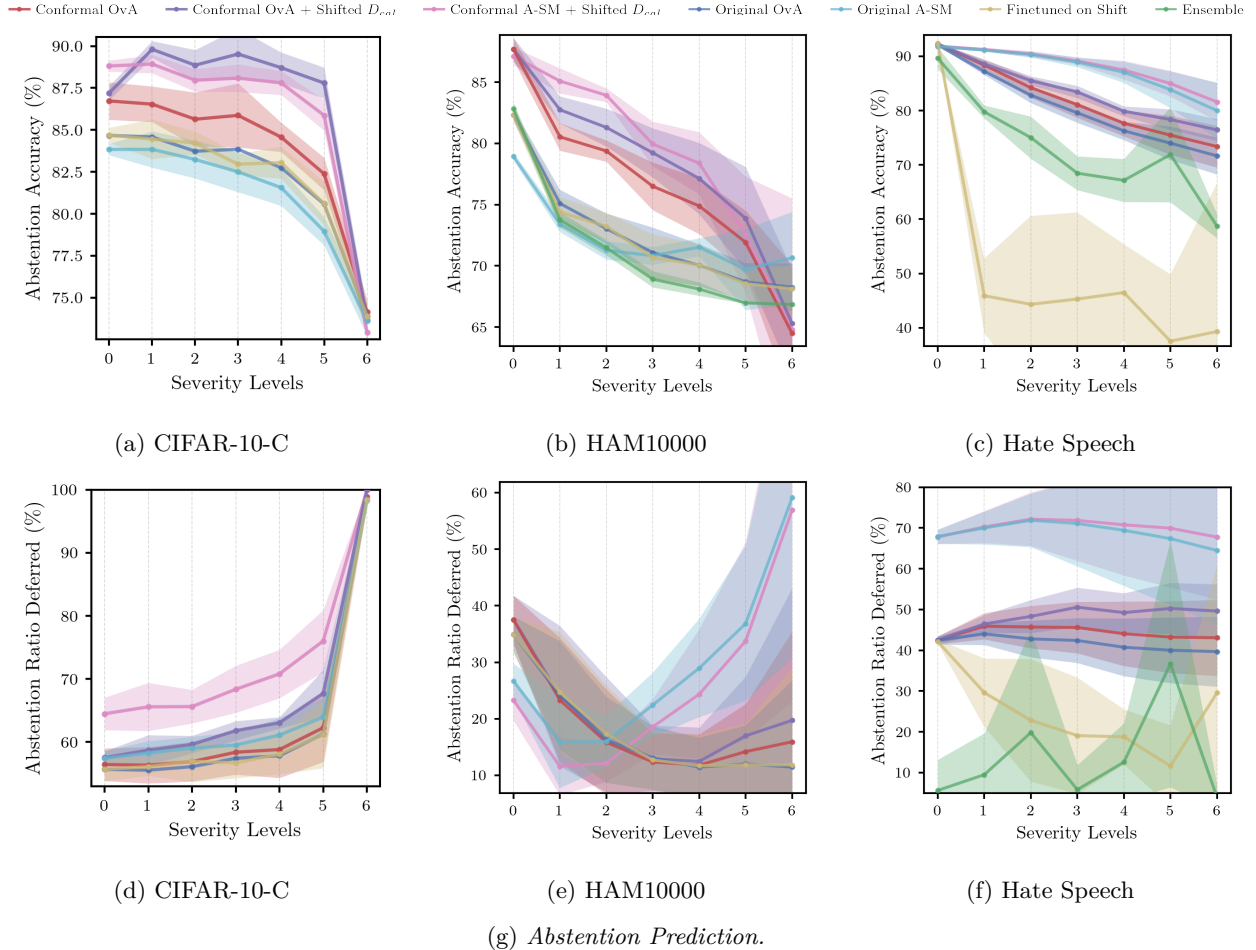
severity levels 1 to 5 and use SVHN at level 6 to induce a semantic shift. For HAM10000, we compose two image corruptions with severity-controlled parameters: (i) masking with fraction p growing from 10% to 50% at level 5 and 80% at level 6, and (ii) Gaussian blur with kernel size increasing from 3×3 to 13×13 and their standard deviation σ from 0 to 4 at level 5 then 6 at level 6. For HateSpeech, we perturb embeddings with adversarial noise injection (Wei & Zou, 2019; Donahue et al., 2016): progressively from level 1 to 5 (i) dimension masking with ratio $r \in [5\%, 25\%]$, (ii) additive jitter with standard deviation increasing from $\approx 2\%$ to $\approx 10\%$, and (iii) simple token-inspired edits with rates up to 25%. Level 6 constitutes an extreme step: $r = 40\%$, jitter $\approx 15\%$, and edit rates $\approx 40\%$. Here, we include additional baselines that test other uncertainty methods (deep ensembling) and if the held-out data we use for split-CP could be better used for finetuning. We detail these additional baselines / method variants below.

Original & Original A-SM We built our original models with OvA surrogate Ψ_{OvA} (Verma et al., 2023) and with A-SM surrogate (Cao et al., 2024) on dataset without shift.

Conformal Following Section 3, we calibrate the rejector with split-CP on a calibration set from the non-shifted dataset to compute the threshold $\hat{\tau}$.

Conformal + Shifted D_{cal} While the above conformal method does not fit to the target distribution, a subset from the target distribution is expected to calibrate the system to target distribution in computing τ in Section 3.

Finetuned on Shift (Baseline) We perform a finetuning to see how L2D framework could take the most advantage of the limited and invaluable data from target distribution. Rejectors were trained on the source distribution will apply on the same subset in the Conformal + Shifted D_{cal} method to finetune.



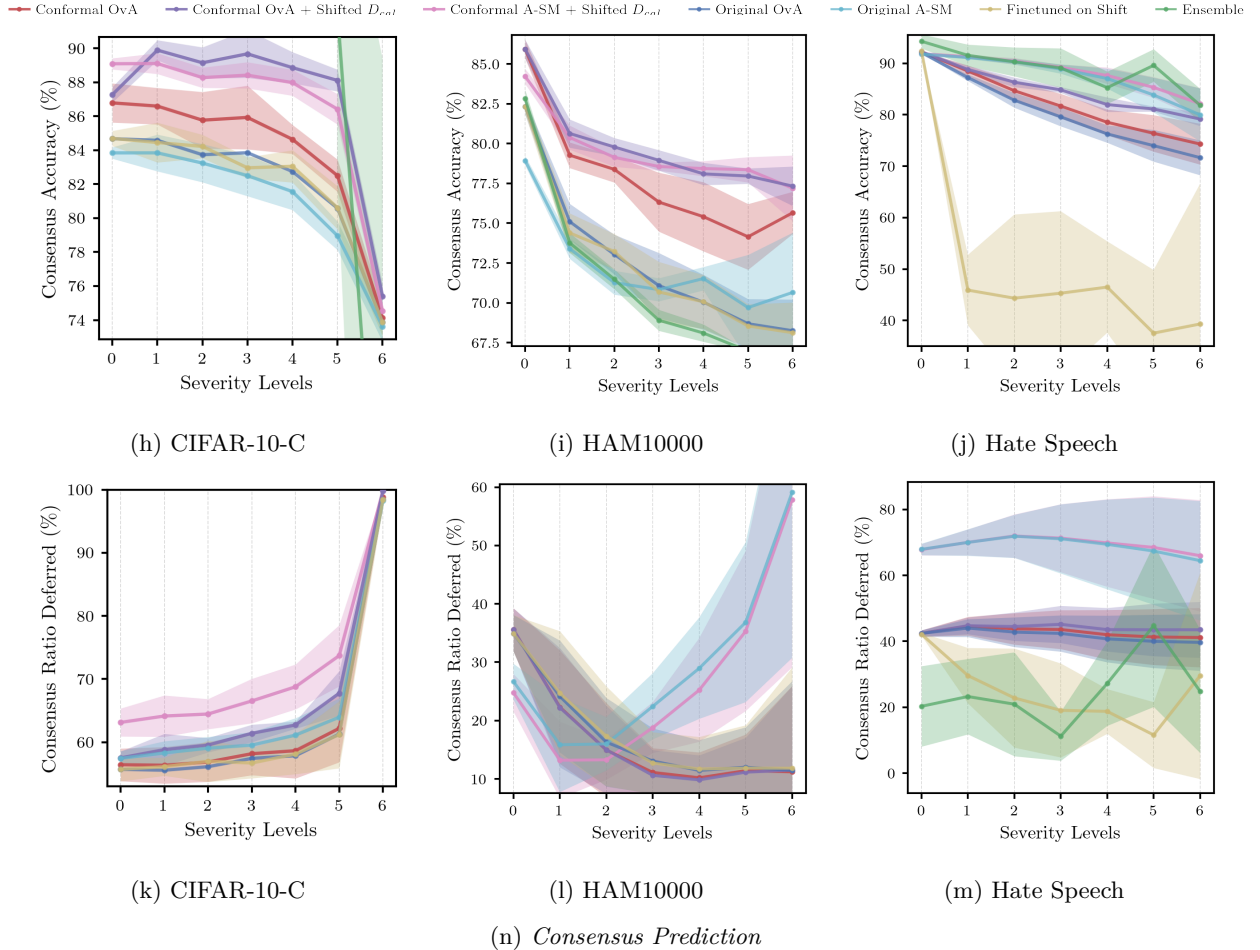
Ensemble (Baseline) We approached the uncertainty ensemble by Liu et al. (2022) in a computationally efficient way by taking advantage of the already trained L2D functions. An uncertainty ensemble can be constructed by reinitializing and trivially retraining m rejector layer functions based on above trained L2D functions. To align with the size of conformal set $C_r(\mathbf{x}; \hat{\tau})$ in this study, we set $m = 2$.

Prediction Figure 2 visualized the accuracy and ratios of deferred instances in the population of the abstention, consensus, human-preferred methods shown in panels 2g, 2n and 2u respectively for OOD data. In this scenario, it is notable that under a non-extreme distribution shift, both proposed conformal methods outperform the original L2D framework. In particular, the *Conformal + Shifted D_{cal}* method demonstrates effective deferral behavior by recognizing uncertainty in the test data, thereby maintaining the overall robustness and accuracy of the system. This approach limits performance degradation across shifts from levels 1 to 5 in CIFAR-10 and HAM10000 from level 1 to 5. In the *Hate Speech*, the performance decline was mitigated, showing a slower rate of deterioration.

5.4 Evaluation on Uncertainty of Rejector $r(\mathbf{x})$

In Section 3.3, we incorporate uncertainty interval to defer and abstain conservatively. Specifically, a L2D system should not defer when a human expert is not believed to predict correctly under confidence level $1 - \alpha$. We introduce two metrics that jointly characterize the performance of rejector in deferring: uncertainty rejector accuracy $\widetilde{acc}(m = y)$ and uncertainty classifier accuracy $\widetilde{acc}(y = y)$ as follow:

$$\widetilde{acc}_{m=y} = \frac{\sum \mathbf{1}\{\mathbb{I}[m = y], r(\mathbf{x}) = 1\}}{\sum \mathbf{1}\{r(\mathbf{x}) = 1\}}, \quad \widetilde{acc}_{y=y} = \frac{\sum \mathbf{1}\{\mathbb{I}[y = y], r(\mathbf{x}) = 0\}}{\sum \mathbf{1}\{r(\mathbf{x}) = 0\}},$$



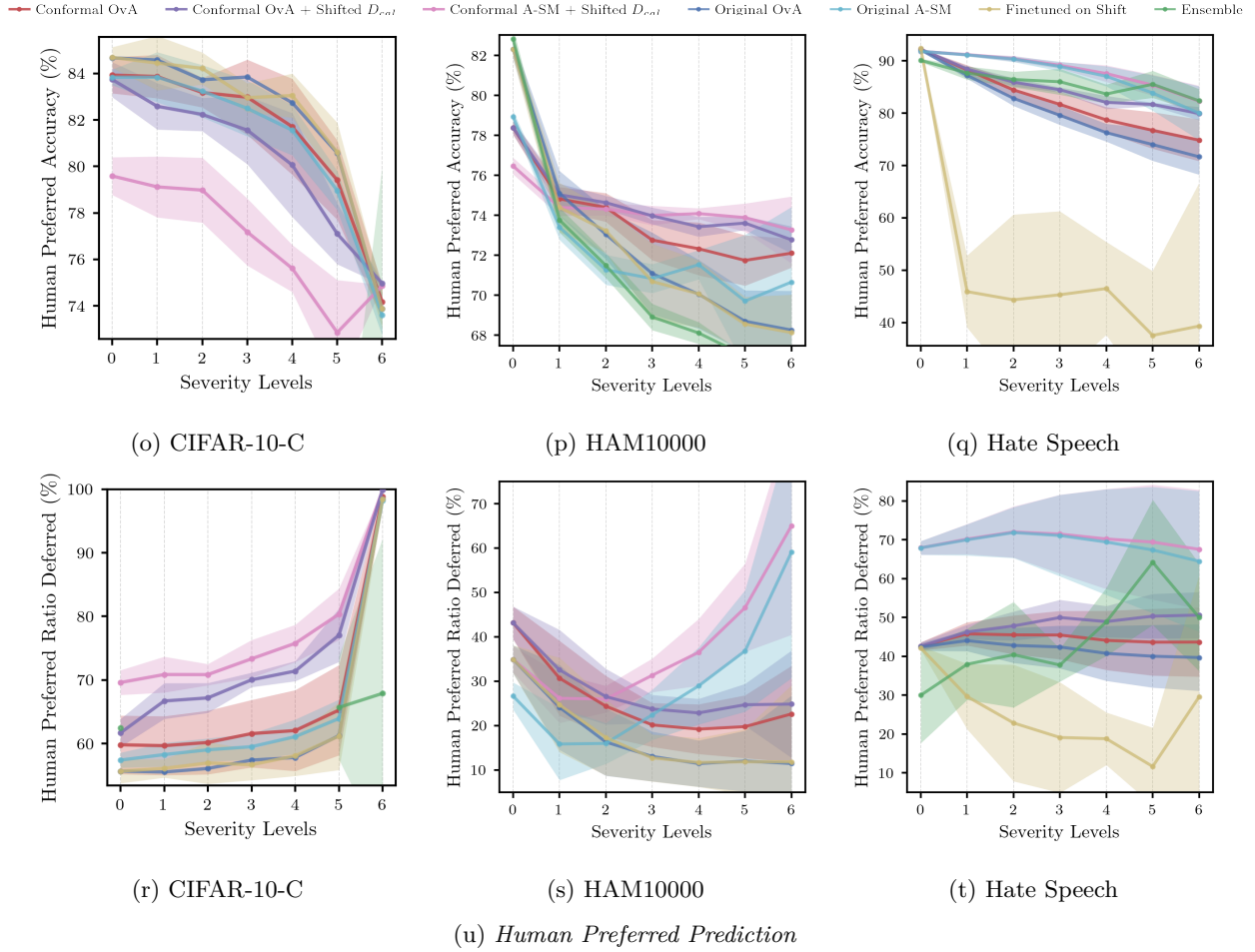


Figure 2: *L2D via Conformal Prediction on OOD*. Figures above report the mean and standard deviation of system accuracy and ratio of deferred instances for methods stated in this section via abstention, consensus, and human-preferred prediction under different levels of distribution shift for **CIFAR-10C**, **HAM10000**, and **Hate Speech**. The proposed L2D framework, employing both OvA and A-SM surrogate losses, exhibits less fluctuation and greater robustness against covariate shift compared to most listed methods, although model deterioration remains evident.

Here, $r(\mathbf{x}) = 1$ indicates deferral, and $r(\mathbf{x}) = 0$ indicates acceptance prediction by the classifier. The metric $\widetilde{\text{acc}}_{m=y}$ measures the correctness of deferred predictions—assuming the original model m approximates the human expert, and $\widetilde{\text{acc}}_{y=y}$ quantifies the accuracy of predictions made autonomously. Together, these metrics succinctly summarize a L2D system’s decision-making performance under uncertainty. We observed improvements in both metrics under L2D with deferral interval, particularly large in gains $\widetilde{\text{acc}}_{y=y}$; detailed results appear in Appendix A. These gains indicate that uncertain instances are predominantly deferred, reducing errors on accepted cases.

6 Conclusions, Limitations, and Future Work

We applied conformal prediction to the rejector component of the learning-to-defer framework with both one-vs-all and asymmetric softmax parameterizations. This approach offers finite-sample, distribution-free guarantees for quantifying uncertainty in the expert’s predictions. Our experiments demonstrate that not only does our method achieve the targeted coverage guarantees with compact prediction sets, but the re-

sulting deferral sets (or intervals) also enable alternative decision-making workflows—such as abstention or expert–model consensus. In particular, we advocate an abstention workflow that empowers the system to respond “I don’t know who knows” when its confidence is insufficient, thereby enhancing the safety and robustness of human–AI collaboration. Such a system alerts users to gather additional information before making a confident decision.

The primary limitation of our work is that our deferral sets are constructed to reflect expert correctness. Thus they reflect but do not perfectly align with the Bayes optimal rejector that compares the expert’s correctness probability with the classifier’s confidence. Thus extending our CP procedure to somehow fuse and compare the uncertainty in the classifier and reject is an exciting and impactful direction for future work.

References

- Anastasios N Angelopoulos, Stephen Bates, et al. Conformal prediction: A gentle introduction. *Foundations and Trends® in Machine Learning*, 16(4):494–591, 2023.
- Varun Babbar, Umang Bhatt, and Adrian Weller. On the utility of prediction sets in human-ai teams. In *International Joint Conference on Artificial Intelligence*, 2022.
- Rina Foygel Barber. Is distribution-free inference possible for binary regression? *Electronic Journal of Statistics*, 14(2):3487 – 3524, 2020. doi: 10.1214/20-EJS1749. URL <https://doi.org/10.1214/20-EJS1749>.
- Rina Foygel Barber, Emmanuel J Candès, Aaditya Ramdas, and Ryan J Tibshirani. Conformal prediction beyond exchangeability. *The Annals of Statistics*, 51(2), 2023.
- Peter L. Bartlett and Marten H. Wegkamp. Classification with a Reject Option Using a Hinge Loss. *Journal of Machine Learning Research*, 2008.
- Yuzhou Cao, Hussein Mozannar, Lei Feng, Hongxin Wei, and Bo An. In defense of softmax parametrization for calibrated and consistent learning to defer. In *Advances in Neural Information Processing Systems*, pp. 38485–38503, 2023.
- Yuzhou Cao, Hussein Mozannar, Lei Feng, Hongxin Wei, and Bo An. In defense of softmax parametrization for calibrated and consistent learning to defer. *Advances in Neural Information Processing Systems*, 36, 2024.
- Mohammad-Amin Charusaie, Hussein Mozannar, David Sontag, and Samira Samadi. Sample Efficient Learning of Predictors that Complement Humans. In *Proceedings of the 39th International Conference on Machine Learning*, 2022.
- C. K. Chow. An Optimum Character Recognition System Using Decision Functions. *IRE Transactions on Electronic Computers*, 1957.
- LP Cordella, C De Stefano, F Tortorella, and M Vento. A method for improving classification reliability of multilayer perceptrons. *IEEE Transactions on Neural Networks*, 6(5):1140–1147, 1995.
- Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Boosting with Abstention. In *Advances in Neural Information Processing Systems*, 2016a.
- Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with Rejection. In *Proceedings of the 27th International Conference on Algorithmic Learning Theory*, 2016b.
- Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. Automated hate speech detection and the problem of offensive language. In *Proceedings of the international AAAI conference on web and social media*, volume 11, pp. 512–515, 2017.
- Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.

- Yizirui Fang and Anthony Bellotti. Investigating data usage for inductive conformal predictors. *arXiv preprint arXiv:2406.12262*, 2024.
- Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In *Advances in Neural Information Processing Systems*, volume 30, 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/4a8423d5e91fda00bb7e46540e2b0cf1-Paper.pdf.
- Martin E Hellman. The nearest neighbor classification rule with a reject option. *IEEE Transactions on Systems Science and Cybernetics*, 6(3):179–185, 2007.
- Patrick Hemmer, Sebastian Schellhammer, Michael Vössing, Johannes Jakubik, and Gerhard Satzger. Forming Effective Human-AI Teams: Building Machine Learning Models that Complement the Capabilities of Multiple Experts. In *International Joint Conference on Artificial Intelligence*, 2022.
- Radu Herbei and Marten H Wegkamp. Classification with reject option. *The Canadian Journal of Statistics*, pp. 709–721, 2006.
- Daniel D Johnson, Daniel Tarlow, David Duvenaud, and Chris J Maddison. Experts don’t cheat: learning what you don’t know by predicting pairs. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 22406–22464, 2024.
- Vijay Keswani, Matthew Lease, and Krishnaram Kenthapadi. Towards Unbiased and Accurate Deferral to Multiple Experts. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Jessie Liu, Blanca Gallego, and Sebastiano Barbieri. Incorporating Uncertainty in Learning to Defer Algorithms for Safe Computer-Aided Diagnosis. *Scientific reports*, 2022.
- David Madras, Toniann Pitassi, and Richard Zemel. Predict Responsibly: Improving Fairness and Accuracy by Learning to Defer. In *Advances in Neural Information Processing Systems*, 2018.
- Anqi Mao, Mehryar Mohri, and Yutao Zhong. Principled Approaches for Learning to Defer with Multiple Experts. *International Symposium on Artificial Intelligence and Mathematics*, 2024a.
- Anqi Mao, Mehryar Mohri, and Yutao Zhong. Predictor-rejector multi-class abstention: Theoretical analysis and algorithms. In *Proceedings of The 35th International Conference on Algorithmic Learning Theory*, pp. 822–867, 2024b.
- Anqi Mao, Mehryar Mohri, and Yutao Zhong. Theoretically grounded loss functions and algorithms for score-based multi-class abstention. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, pp. 4753–4761, 2024c.
- Hussein Mozannar and David A. Sontag. Consistent Estimators for Learning to Defer to an Expert. In *In Proceedings of the 37th International Conference on Machine Learning*, 2020.
- Hussein Mozannar, Hunter Lang, Dennis Wei, Prasanna Sattigeri, Subhro Das, and David Sontag. Who Should Predict? Exact Algorithms For Learning to Defer to Humans. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, 2023.
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in Neural Information Processing Systems*, 32, 2019.
- Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence machines for regression. In *13th European Conference on Machine Learning*, pp. 345–356, 2002.
- Maithra Raghu, Katy Blumer, Greg Corrado, Jon Kleinberg, Ziad Obermeyer, and Sendhil Mullainathan. The Algorithmic Automation Problem: Prediction, Triage, and Human Effort. *ArXiv e-Prints*, 2019.

- N Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. In *Advances in Neural Information Processing Systems*, pp. 3581–3591, 2020.
- Aaron Roth, Alexander Tolbert, and Scott Weinstein. Reconciling individual probability forecasts. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pp. 101–110, 2023.
- Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008.
- Eleni Straitouri, Lequn Wang, Nastaran Okati, and Manuel Gomez Rodriguez. Improving expert predictions with conformal prediction. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- Dharmesh Tailor, Mohammad Emtiyaz Khan, and Eric Nalisnick. Exploiting Inferential Structure in Neural Processes. In *Proceedings of the 39th Conference on Uncertainty in Artificial Intelligence*, 2023.
- Dharmesh Tailor, Aditya Patra, Rajeev Verma, Putra Manggala, and Eric Nalisnick. Learning to defer to a population: A meta-learning approach. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, pp. 3475–3483, 2024.
- Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021.
- Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 2018.
- Rajeev Verma and Eric Nalisnick. Calibrated Learning to Defer with One-vs-All Classifiers. In *Proceedings of the 39th International Conference on Machine Learning*, 2022.
- Rajeev Verma, Daniel Barrejón, and Eric Nalisnick. Learning to Defer to Multiple Experts: Consistent Surrogate Losses, Confidence Calibration, and Conformal Ensembles. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, 2023.
- Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World*. Springer, 2005.
- Jason Wei and Kai Zou. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. *arXiv preprint arXiv:1901.11196*, 2019.
- Ming Yuan and Marten Wegkamp. Classification Methods with Reject Option Based on Convex Risk Minimization. *Journal of Machine Learning Research*, 2010.
- Ahmed Zaoui, Christophe Denis, and Mohamed Hebiri. Regression with Reject Option and Application to KNN. *Advances in Neural Information Processing Systems*, 2020.

A Experiment: L2D with Deferral Intervals

This section documents the experiment evaluating L2D workflow with conformal deferral intervals, both without and with **Conformal + Shifted** D_{cal} (i.e., recalibrating on a small target-domain subset under shift).

Reliability Diagram We perform evaluation of calibration and examine the validity of conformal prediction by plotting reliability diagram and computing expected calibration error (ECE). We define the expected accuracy and ECE as

$$\widehat{\text{acc}} = \mathbb{P}(m = y | p_m(x) = c), \quad \text{ECE}(p_m) = \mathbb{E}_{\mathbf{x}} |\mathbb{P}(m = \mathbf{y} | p_m(\mathbf{x}) = c) - c|,$$

where c is the confidence level. From the prediction interval $[b_l(\mathbf{x}), b_r(\mathbf{x})]$, the uncertainty bar of the expected accuracy could be calculated as $\Delta_- = b_l(\mathbf{x}) \mathbb{I}[m = y]$, $\Delta^+ = b_r(\mathbf{x}) \mathbb{I}[m = y]$. Figure 3 presents the reliability diagram and constructs the error bar by distributing the prediction interval across the accuracy.

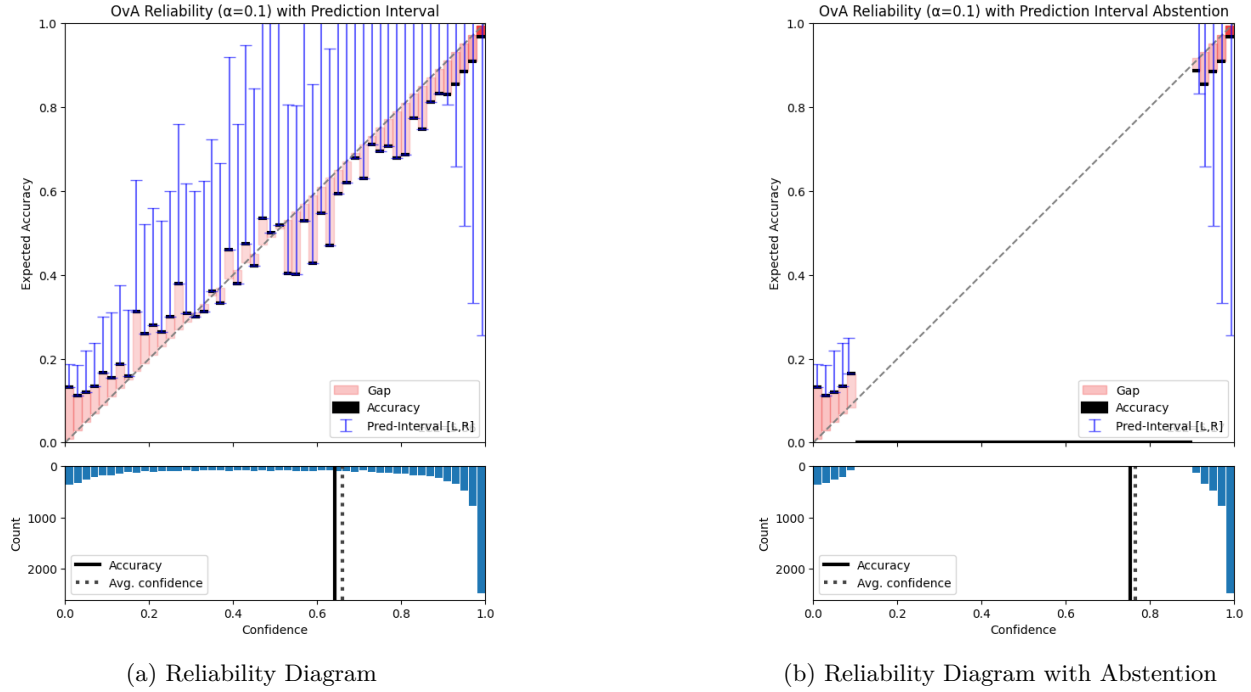


Figure 3: Evaluation of Calibration on CIFAR-10: Subfigure (a) reports a reliability diagram w/o abstention and the expected calibration error (ECE) is 5.86. Subfigure (b) reports a reliability diagram w/ abstention and the ECE is 5.17

Uncertainty Rejector Evaluation we evaluate the performance of the uncertainty rejector $r(\mathbf{x})$ by accuracy metrics $\widehat{\text{acc}}_{m=y}$ and $\widehat{\text{acc}}_{y=y}$ on CIFAR-10 datasets. Table 3 could show a significant increase in L2D with the prediction interval in $\widehat{\text{acc}}_{m=y}$ and $\widehat{\text{acc}}_{y=y}$, while the deferral interval is too large. This confirms that the binary regression would adjust the rejector to make deferral decision and abstention with care of uncertainty. It also reports that with abstention, the system accuracy increases from 83.87% to 93.75%, Table 3 evaluated the $\widehat{\text{acc}}_{m=y}$ and $\widehat{\text{acc}}_{y=y}$ in L2D with deferral interval.

Table 3: Performance under uncertainty rejection across severity levels, reporting standard accuracy, uncertainty-based rejector accuracy ($\widetilde{acc}_{m=y}$), classifier accuracy ($\widetilde{acc}_{y=y}$), non-abstention rejector accuracy, and model preferred accuracy.

Severity	Method	\widetilde{acc} (%)	$\widetilde{acc}_{m=y}$ (%)	$\widetilde{acc}_{y=y}$ (%)
0	Baseline	83.87	84.44	32.75
	Deferral Interval	72.57	92.77	70.06
	Abstention	93.75	92.77	74.14
	Model Preferred	92.77	92.77	70.06
1	Baseline	83.10	84.64	35.27
	Deferral Interval	71.90	92.88	67.65
	Abstention	93.36	92.88	74.95
	Model Preferred	92.88	92.88	67.65
	Deferral Interval + Shifted D_{cal}	69.75	90.76	67.62
	Abstention + Shifted D_{cal}	93.18	90.76	72.83
	Model Preferred + Shifted D_{cal}	90.76	90.76	67.62
2	Baseline	80.65	83.21	35.35
	Deferral Interval	68.80	91.96	65.88
	Abstention	91.95	91.96	73.93
	Model Preferred	91.96	91.96	65.88
	Deferral Interval + Shifted D_{cal}	70.00	89.56	66.59
	Abstention + Shifted D_{cal}	92.27	89.56	72.91
	Model Preferred + Shifted D_{cal}	89.56	89.56	66.59
3	Baseline	77.20	83.62	40.53
	Deferral Interval	62.40	94.55	58.79
	Abstention	91.75	94.55	72.58
	Model Preferred	94.55	94.55	58.79
	Deferral Interval + Shifted D_{cal}	63.00	90.09	59.79
	Abstention + Shifted D_{cal}	92.47	90.09	70.45
	Model Preferred + Shifted D_{cal}	90.09	90.09	59.79
4	Baseline	76.90	83.78	42.53
	Deferral Interval	63.10	94.32	59.06
	Abstention	91.75	94.32	72.75
	Model Preferred	94.32	94.32	59.06
	Deferral Interval + Shifted D_{cal}	64.80	92.19	59.33
	Abstention + Shifted D_{cal}	92.78	92.19	71.33
	Model Preferred + Shifted D_{cal}	92.19	92.19	59.33
5	Baseline	73.75	84.05	46.40
	Deferral Interval	59.80	90.06	56.97
	Abstention	90.62	90.06	72.80
	Model Preferred	90.06	90.06	56.97
	Deferral Interval + Shifted D_{cal}	59.85	90.96	57.03
	Abstention + Shifted D_{cal}	91.83	90.96	72.03
	Model Preferred + Shifted D_{cal}	90.96	90.96	57.03

B Additional Experiment

B.1 Dual CP: Conformal Prediction on Both Classifier and Rejector

While Section 3 parameterizes the expert correctness with $\hat{p}(m = y|\mathbf{x})$. In this section, we compare the uncertainty of both classifier and human to make deferral decision explicitly. To this end, one may simulate the human annotator with a weak classifier $human^*(\mathbf{x}) = m$, which is trained on the same data \mathbf{x} and m as the L2 system and expert prediction. Likewise Section 3 in constructing deferral set $C_r(\mathbf{x}; \hat{\tau})$ for rejector $r^*(\mathbf{x})$, one may explicitly construct classifier prediction sets $C_h(\mathbf{x}; \hat{\tau}_h)$ on $h^*(\mathbf{x})$ and annotator prediction set $C_{human}(\mathbf{x}; \hat{\tau}_{human})$. comparing uncertainty of $h^*(\mathbf{x})$ and $human^*(\mathbf{x})$ to defer. Given these estimators, one may construct the usual conformity score for the classifier $h^*(x)$ prediction:

$$s(\mathbf{x}, y, m; \hat{p}(y|x)) = \sum_{k=1}^k \hat{p}(k|x)$$

where $\hat{p}(k|x)$ is normalized probability of each class estimated by classifier $\hat{p}(k|\mathbf{x})$ is $\max_{k \in [1, K]} \phi(g_k(\mathbf{x}))$.

One may then defer by comparing the size of two prediction set, which means

$$r^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \|C_h(\mathbf{x}; \hat{\tau}_h)\| > \|C_{human}(\mathbf{x}; \hat{\tau}_{human})\|, \\ 0 & \text{otherwise} \end{cases}$$

We reuse the experimental setup and one-vs-all (OvA) parameterization from Section 5. We then apply the dual conformal-prediction (CP) procedure described above to the CIFAR-10, HAM10000, and Hate Speech datasets without shift. The resulting coverage and deferral-rate metrics are reported in Table 4.

Table 4: *Dual CP Experiment*. We report system accuracy on dual CP on the same three datasets but do not notice performance improvement.

Dataset	Method	Sys. Acc.(%)
CIFAR-10	Base Model	84.74
	Dual CP	74.36
HAM10000	Base Model	82.10
	Dual CP	80.72
Hate Speech	Base Model	91.44
	Dual CP	91.24

C Finetuning on Uncertain Instances

To broaden coverage, one may finetune the L2D system on instances that CP flags as uncertain. The initial threshold $\hat{\tau}$ could be again computed over initial calibration set \mathcal{D}_2 as Section 3.1. Building on the dual-CP deferral setup in Section B.1, we now run CP on the training set \mathcal{D}_1 and form prediction set $C_r(\mathbf{x}; \hat{\tau})$. We could then define the uncertain subset by a cardinality (“width”) threshold t_{abs} :

$$\mathcal{D}_{uncertain} = \{\mathbf{x} \in \mathcal{D}_1 : \|C_r(\mathbf{x}; \hat{\tau})\| \geq t_{abs}\}.$$

We reweight the OvA surrogate loss by prediction set width $\|C_r(\mathbf{x}; \hat{\tau})\|$ and perform a fine-tuning to explicitly emphasize these uncertain inputs:

$$\tilde{\psi}_{\text{Re-OvA}}(g_1, \dots, g_{K+1}; \mathbf{x}, y, m) = \phi[g_y(\mathbf{x})] - \frac{\phi[-g_{K+1}(\mathbf{x})]}{\|C_{human}(\mathbf{x}; \hat{\tau}_{human})\|}$$

followed by re-running the dual CP to compute updated thresholds $\hat{\tau}'$ and $\hat{\tau}'_{\text{human}}$ on \mathcal{D}_2 . Resulting coverage and deferral-rate metrics are reported in Table 5. Finetuning on uncertain instances can raise accuracy relative to dual CP without finetuning, but it does *not* surpass the original L2D baseline in our runs: suggesting that the OvA parameterization may capture rejector uncertainty and that post-hoc CP calibration remains a strong, distribution-free control for deferral decisions.

Table 5: *Dual CP Finetuning*. We report system accuracy on top of the previous dual CP experiment where $t_{\text{abs}} = 1$

Dataset	Method	Sys. Acc. (%)
HAM10000	FT Dual CP	81.46
Hate Speech	FT Dual CP	91.42