
Learning to Defer under Expert Drift

Ziyao Mou
Johns Hopkins University
zmou1@jhu.edu

Andrea Wynn
Johns Hopkins University
awynn13@jhu.edu

Eric Nalisnick
Johns Hopkins University
nalisnick@jhu.edu

Abstract

The *learning to defer* (L2D) framework enables safety in machine learning and AI systems by allocating difficult or critical decisions to a human expert. However, prior work on L2D assumes fixed expert reliability, overlooking temporal fluctuations in human performance due to factors such as fatigue and cognitive biases. Humans commonly exhibit such substantial, systematic performance shifts in high-stakes domains such as radiology, aviation, and driving. We propose *Expert Drift Adapted Learning-to-Defer* (EDA-L2D), a novel expert-aware L2D framework that explicitly models these temporal shifts by conditioning the deferral head on recent histories of expert outcomes via a sequence model. This history-aware approach enables adaptive deferral policies that anticipate reliability fluctuations and better allocate decisions between human and machine over time. We provide comprehensive experiments across 3 diverse tasks - image classification, medical diagnosis, and hate speech detection - showing that for a time-dependent expert, our approach consistently achieves higher performance and better deferral rates than prior L2D approaches.

1 INTRODUCTION

Hybrid intelligent (HI) systems combine human and machine intelligence to achieve goals unattainable by either alone, emphasizing complementarity over replacement (Kamar, 2016; Dellermann et al., 2019; Akata et al., 2020). One popular HI paradigm is *learning to defer* (L2D): the system learns not only to predict a label but also to decide when to hand off the decision to a human expert. The original formulation of L2D shows that accounting for the expert’s strengths and biases can improve overall accuracy and fairness relative to a static rejection framework (Madras et al.,

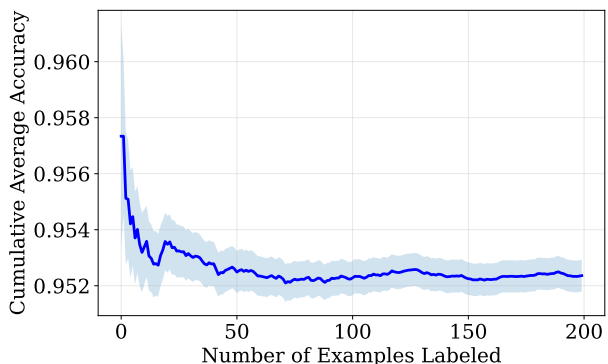


Figure 1: *Annotator Fatigue in CIFAR-10-H*. After removing attention checks, we compute the average cumulative accuracy across all 2,555 human annotators at each time-step. We observe a slight but clear downward drift in performance over time.

2018).

However, existing L2D systems assume that expert behavior is fixed or stationary. In practice, human performance often varies significantly over time due to fatigue (Figalová et al., 2024; Peukert et al., 2023; Pan et al., 2022; Bruni et al., 2012) and various cognitive biases (Legler et al., 2025; Urai et al., 2019; van de Wouw et al., 2024). We show real-world evidence of this effect on a widely studied computer vision task, CIFAR10, by analyzing real human label data from over 2,500 human annotators in the CIFAR10H dataset (Joshua et al., 2019). We show in Fig. 1 that there exists a slight but consistent decline in annotator accuracy over time: the more examples the annotators label, the less accurate they are at producing correct labels. This is a notable observation: CIFAR10 is a relatively simple task (image classification over only 10 classes of common objects), yet even on such a simple task, we observe annotator fatigue. This motivates developing a L2D approach that can account for these structured temporal shifts, in turn allowing it to naturally model and adapt to real human experts.

We address this challenge by proposing a novel frame-

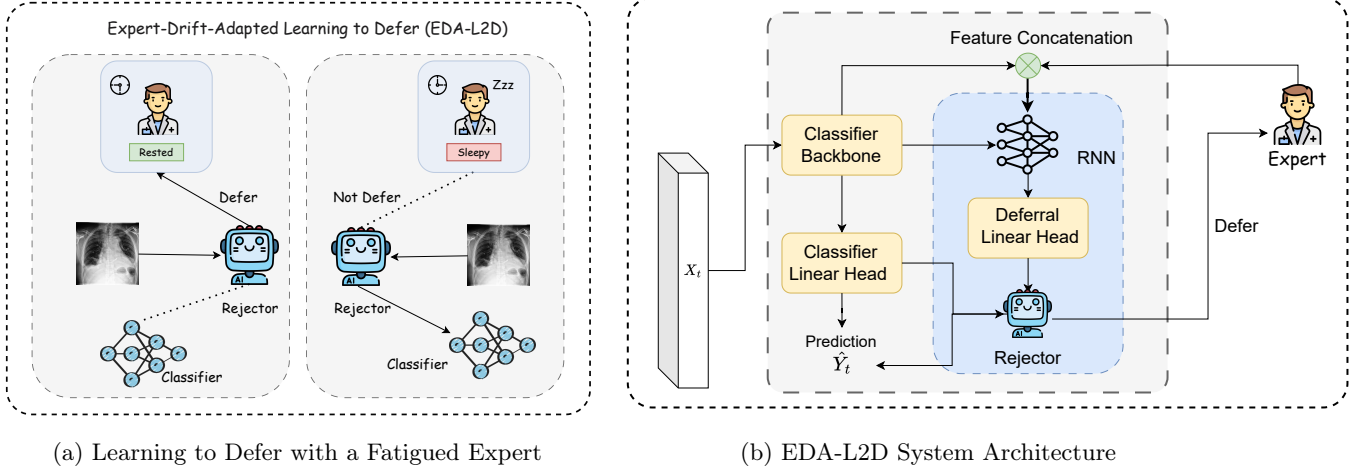


Figure 2: (a) An example use case of our EDA-L2D framework in a setting with expert performance drift: a radiologist whose performance degrades over a workday. Our model monitors past outcomes to infer current expert reliability in real time, deferring to the expert when reliable and assigning cases to the model when the expert’s performance declines. (b) Overview of our expert drift adapted L2D framework. A streaming sequence of inputs is encoded by the backbone into features. At each timestep, the feature vector is concatenated with the expert’s previous prediction and passed to an RNN module. The classifier head outputs class logits and the defer head outputs a deferral score; a deferral gate then assigns the instance to either the expert or the classifier. Because the expert’s competence drifts over time, this design makes the L2D system time-aware and adaptive to temporal shifts.

work for L2D: *Expert-Drift-Adapted Learning to Defer* (EDA-L2D). Our framework explicitly models non-stationarity in expert performance. We modify L2D’s rejector sub-component to condition on short-horizon histories of expert decisions. This allows the system to be aware of the expert’s current behavior, and modify its deferral policy accordingly. We illustrate our approach Fig. 2a. Concretely, we integrate a sequence model into the deferral head to explicitly model expert reliability over time. We validate our approach across diverse tasks including medical image recognition, hate speech detection, and image classification. We demonstrate consistent gains in system accuracy over existing, time-agnostic L2D approaches. In summary, our contributions are as follows:

- We propose a novel *Expert-Drift-Adapted Learning to Defer* (EDA-L2D) that explicitly models variance in human expert performance over time, allowing the system to flexibly adapt to realistic human experts. We propose our approach in §3, including a new time-aware learning objective (§3.2) and model formulation (§3.3).
- We validate the approach with simulated, time-dependent experts on real-world tasks that require HI solutions (§5.1, §5.2, §5.3), and show that this time-variance is present in real human-annotated data (Fig.1).

2 BACKGROUND

We first describe the traditional L2D setting (Mozannar and Sontag, 2021), in which the expert’s abilities are assumed to be static and time-invariant.

2.1 Data & Models

We will exclusively consider the problem of multi-class classification with a feature space \mathcal{X} and the label space $\mathcal{Y} = \{1, \dots, K\}$. Given some input $x \in \mathcal{X}$, we define $\mathbb{P}(y | x)$ as the unknown label-generating distribution and $\mathbb{P}(m | x)$ as the unknown distribution over expert predictions $m \in \mathcal{Y}$. L2D is comprised of two sub-components: a *classifier* $h : \mathcal{X} \rightarrow \mathcal{Y}$ and a *rejector* $r : \mathcal{X} \rightarrow \{0, 1\}$, which decides whether to make a prediction using the classifier ($r(x) = 0$) or by deferring to the expert ($r(x) = 1$).

2.2 Classifier-Rejector Loss Function

To train our L2D model, we need to fit both the rejector and classifier models. The classifier incurs a loss of zero (correct) or one (incorrect) when it makes a prediction. Similarly, the human expert incurs the same 0-1 loss when making a prediction (i.e. $r(x) = 1$). Combining these two 0-1 losses using the rejector model results in the combined classifier-rejector loss:

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

where \mathbb{I} denotes an indicator function checking the prediction against the ground-truth label. When minimizing this loss, the resulting Bayes optimal classifier and rejector satisfy:

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

where $\mathbb{P}(y | \mathbf{x})$ represents the label probability under the data generating process, and $\mathbb{P}(m = y | \mathbf{x})$ is the probability of the expert making the correct prediction. The expert may have knowledge not available to the classifier, so they could outperform the Bayes optimal classifier.

2.3 Softmax Surrogate Loss

A consistent surrogate loss for the above L_{0-1} loss can be derived following the formulation from [Mozannar and Sontag \(2021\)](#). First, we consider an augmented label space that includes both the label space \mathcal{Y} and the rejection option \perp : $\mathcal{Y}^\perp := \mathcal{Y} \cup \{\perp\}$. Secondly, for a class $k \in [1, K]$, let $g_k : \mathcal{X} \mapsto \mathbb{R}$, and let $g_\perp : \mathcal{X} \mapsto \mathbb{R}$ denote the rejection option. We can combine these $K + 1$ with a loss resembling the cross-entropy loss for a softmax parameterization:

$$\begin{aligned} \phi_{\text{SM}}(g_1, \dots, g_K, g_\perp; \mathbf{x}, y, m) &= \\ &= -\log \left(\frac{\exp\{g_y(\mathbf{x})\}}{\sum_{y' \in \mathcal{Y}^\perp} \exp\{g_{y'}(\mathbf{x})\}} \right) \\ &\quad - \mathbb{I}[m = y] \log \left(\frac{\exp\{g_\perp(\mathbf{x})\}}{\sum_{y' \in \mathcal{Y}^\perp} \exp\{g_{y'}(\mathbf{x})\}} \right). \end{aligned} \quad (3)$$

Here, the first term maximizes g_k for the true label k . The second term maximizes the rejection function g_\perp , but only when the expert’s prediction is correct. At test time, the classifier takes the maximum over the classes: $\hat{y} = h(\mathbf{x}) = \arg \max_{k \in [1, K]} g_k(\mathbf{x})$. Similarly, we formulate the rejection function as $r(\mathbf{x}) = \mathbb{I}[g_\perp(\mathbf{x}) \geq \max_k g_k(\mathbf{x})]$. The minimizers $g_1^*, \dots, g_K^*, g_\perp^*$ of ϕ_{SM} also uniquely minimize the 0-1 loss from Equation 1, $L_{0-1}(h, r)$.

3 EXPERT-DRIFT-ADAPTED LEARNING TO DEFER

This section introduces our *Expert-Drift-Adapted Learning to Defer* (EDA-L2D) framework, which extends L2D to account for non-stationarity in expert performance. Accordingly, we derive the corresponding softmax surrogate loss for this setting. We also propose multiple parameterizations, including one that uses a recurrent neural network to incorporate short-horizon temporal context into the deferral policy.

3.1 Setting: Non-Stationary Expert

We consider the same setting as Section 2—L2D for multi-class classification—except now we assume the expert is non-stationary. Thus, instead of a fixed distribution $\mathbb{P}(m | x)$, the expert’s predictions are generated by $\mathbb{P}_t(m_t | x)$, with the subscript $t \in \mathcal{T}$ denoting a time index. Note that we are *assuming the classification task itself is stationary*, i.e. $y \sim \mathbb{P}(y | x)$, which still has no time dependence.

To ground the notation in an example, consider an L2D system for radiology. The problem of forming a diagnosis or other prediction from the medical image itself is a static task, since we assume the distribution of patients and the mechanism that govern their health is not changing over time (or is at least changing so slowly as to be negligible, e.g. drift in hospital patient demographics). Yet the radiologist who is paired with the L2D system will change, at the very least becoming more tired and distracted over the course of a 10+ hour shift. The expert drift may also be non-monotonic: perhaps her performance degrades over the course of a morning, but after lunch and a mid-day walk, the expert feels rejuvenated for a few more hours.

Returning to the technical details, our EDA-L2D system will have a classifier that is defined just as above: $h : \mathcal{X} \rightarrow \mathcal{Y}$. Again, the classifier is time-invariant since the prediction task itself is not time dependent. The difference, however, arises in the rejector: $r : \mathcal{X} \times \mathcal{T} \mapsto \{0, 1\}$, meaning the rejector is a function of both the feature space \mathcal{X} and time \mathcal{T} . While in the simplest case \mathcal{T} would be a scalar time index, we could also formulate \mathcal{T} as a feature space that describes the current state of the expert (e.g. tabular biomarkers).

3.2 Learning Objective and Bayes Solutions

The learning objective is similar to that in Section 2 (equation 1), except now it takes the form of a summation over time:

$$\begin{aligned} L_{0-1}^{\mathcal{T}}(h, r) &= \sum_{t \in \mathcal{T}} \mathbb{E}_{\mathbf{x}, y, m_t} \left[(1 - r(\mathbf{x}, t)) \mathbb{I}[h(\mathbf{x}) \neq y] \right. \\ &\quad \left. + r(\mathbf{x}, t) \mathbb{I}[m_t \neq y] \right]. \end{aligned} \quad (4)$$

As the distribution of y does not depend on time, the classifier’s Bayes solution is just as before: $h^*(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \mathbb{P}(y = y | \mathbf{x})$. Yet the Bayes optimal rejector does change, becoming time-dependent like so:

$$r^*(\mathbf{x}, t) = \mathbb{I} \left[\mathbb{P}_t(m_t = y | \mathbf{x}) \geq \max_{y \in \mathcal{Y}} \mathbb{P}(y = y | \mathbf{x}) \right].$$

The intuition is the same as above—compare the chance of expert correctness vs classifier confidence for the modal label—yet now the expert correctness term is time dependent, $\mathbb{P}_t(m_t = y | \mathbf{x})$.

3.3 Implementations of the Softmax Surrogate Loss

The derivation for the corresponding softmax surrogate loss follows fairly directly from Mozannar and Sontag (2021)’s original derivation. Yet we consider two distinct parameterizations of the underlying model. The first is a ‘per time step’ version that treats the L2D problem independently at each time step. We consider this the most naive extension of traditional L2D that yet still satisfies our motivating setting of expert drift. We then go on to describe how to use a recurrent neural network (RNN) to parameterize the rejector, which allows the expert’s previous predictions to help inform the deferral policy. For all implementations, we assume we have access to a training set $\mathcal{D} = \left\{ \{(\mathbf{x}_{t,n}, y_{t,n}, m_{t,n})\}_{n=1}^N \right\}_{t \in \mathcal{T}}$, which consists of N feature-label-expert-prediction triplets at each of $|\mathcal{T}|$ time steps. Again, we emphasize that only the underlying distribution of $m_{t,n}$ is time-varying, and the time indices on $\mathbf{x}_{t,n}, y_{t,n}$ are there merely to ‘synchronize’ the feature-label pairs with the expert predictions.

Per-Time-Step Implementation and Loss The softmax surrogate in equation 3 can be naturally extended to the drifting-expert setting by instantiating a traditional L2D model at each time step. Specifically, we define a model at each time step and train in locally for that time step via the surrogate:

$$\begin{aligned} \phi_{\text{SM}}^t(g_1, \dots, g_K, g_{\perp}; \mathbf{x}, y, m_t) = & \\ & - \log \frac{\exp\{g_y(\mathbf{x})\}}{\sum_{y' \in \mathcal{Y}^{\perp}} \exp\{g_{y'}(\mathbf{x})\}} \\ & - \mathbb{I}[m_t = y] \log \frac{\exp\{g_{\perp}^t(\mathbf{x})\}}{\sum_{y' \in \mathcal{Y}^{\perp}} \exp\{g_{y'}(\mathbf{x})\}}. \end{aligned} \quad (5)$$

Notice that g_{\perp}^t is specific to the current time step, but g_1, \dots, g_K could be shared across all time steps (since they encode the classifier). As mentioned above, we consider this a naive implementation and will primarily use it to benchmark our second approach, described below. Having $g_{\perp}^t(\mathbf{x})$ only be a function of \mathbf{x} means that it ignores patterns that could be detected by having information about the expert’s behavior at previous time steps. Moreover, this implementation has the clear drawback that it cannot generalize to time steps that were not seen in the training set. For instance, if all training sets had $|\mathcal{T}| = 10$, denoting that the expert worked a 10-hour shift. Then it is not clear how to deploy and use the model when the expert could for a longer period of time than 10 hours.

Loss for RNN-Based Model We next consider a more sophisticated implementation that uses an RNN-based parameterization to incorporate informa-

tion about the expert’s previous predictions. Here we use a surrogate loss defined across all time steps:

$$\begin{aligned} \phi_{\text{SM}}^{\mathcal{T}}(g_1, \dots, g_K, g_{\perp}; \{\mathbf{x}_t, y_t, m_t\}_{t \in \mathcal{T}}) = & \\ \sum_{t \in \mathcal{T}} - \log \frac{\exp\{g_{y_t}(\mathbf{x}_t)\}}{\exp\{g_{\perp}(\mathbf{x}_t, t)\} + \sum_{y' \in \mathcal{Y}} \exp\{g_{y'}(\mathbf{x}_t)\}} & \\ - \mathbb{I}[m_t = y] \log \frac{\exp\{g_{\perp}(\mathbf{x}_t, t)\}}{\exp\{g_{\perp}(\mathbf{x}_t, t)\} + \sum_{y' \in \mathcal{Y}} \exp\{g_{y'}(\mathbf{x}_t)\}}. & \end{aligned} \quad (6)$$

where here $g_{\perp}(\mathbf{x}_t, t)$ is parameterized with an RNN so that it can leverage information from previous time steps. While we have left the time feature t vague, this will contain the expert’s previous prediction m_{t-1} when it is available. We next describe the neural architecture in more detail.

RNN Architecture The RNN first takes as input the feature vector \mathbf{x}_t , producing a hidden representation $\mathbf{f}_t = \psi_0(\mathbf{x}_t)$. For instance, if \mathbf{x}_t is an image, then $\phi(\cdot)$ could be a convolutional NN. The feature representation is then passed into a recurrent module: $\mathbf{z}_t = \psi_1(\mathbf{z}_{t-1}, \mathbf{f}_t, \hat{m}_{t-1})$. During training we use teacher forcing for $\hat{m}_{t-1} \in \{0, 1\}$ (expert correct/incorrect at $t-1$); at test time it is computed from realized routing and the expert’s actual outcome. Lastly, the g -functions found in Equation 6 are produced by a final linear transformation: $g_k(\mathbf{x}_t) = \mathbf{w}_k^{\top} \mathbf{z}_t$ ($k \in 1, \dots, K$) and $g_{\perp}(\mathbf{x}_t) = \mathbf{w}_{\perp}^{\top} \mathbf{z}_t$. The predicted label is chosen via $\arg \max_k g_k(\mathbf{x}_t)$. $g_{\perp}(\mathbf{x}_t)$ quantifies the preference for deferral, with the system deferring if $g_{\perp}(\mathbf{x}_t) \geq \max_k g_k(\mathbf{x}_t)$.

The intuition behind this formulation is that the classifier and rejector should not operate solely on the current input but instead incorporate temporal context that reflects the expert’s evolving reliability. By conditioning the recurrent state \mathbf{z}_t on both the encoded input \mathbf{f}_t and the previous correctness indicator \hat{m}_{t-1} , the system adaptively tracks fluctuations in expert performance. The classifier head then focuses on predicting the label from this time-dependent state, while the deferral head estimates whether the expert should be trusted at the current step.

4 RELATED WORK

Learning to Defer (L2D) Early work on classifiers that can choose to reject or abstain instead of making a prediction (Chow, 1957) inspired more recent work on how to model what happens after the classifier abstains (Madras et al., 2018), with Mozannar and Sontag (2021) deriving the first consistent surrogate loss for the problem. This work has been expanded in recent years to address several limitations, including mis-calibration (Verma and Nalisnick, 2022; Cao et al.,

2023), underfitting (Narasimhan et al., 2022), realizable consistency (Mozannar et al., 2023), sample complexity (Charusaie et al., 2022), data scarcity (Hemmer et al., 2023), and deferral to multiple experts (Verma et al., 2023). However, these approaches do not consider the ‘messiness’ of modeling *human* experts—as we do—and could just as well be applied to a black-box API.

L2D for Sequences and Populations L2D has also been extended to the difficult settings of sequential tasks and encountering novel experts. Regarding the former, *Sequential Learning-to-Defer* (Joshi et al., 2023) incorporates temporal dynamics by framing deferral as a sequential decision-making problem, employing model-based reinforcement learning. This work, however, is primarily concerned with the sequential nature of the prediction task itself and not robustness to expert drift, which is our primary concern. Regarding the latter, *Learning to Defer to a Population* (Tailor et al., 2024) uses meta-learning to adapt to never-before-seen experts at test-time. Their motivation is to be robust to novel experts (with the novelty not being a function of time), whereas we are concerned about modeling the temporal drift of a known expert.

Temporal Shifts in Human Performance In many high-stakes domains, human performance shift over time (due to fatigue or other cognitive factors) are well-documented and widely studied. Fatigue is extremely common in high-stakes domains and long-horizon tasks, such as radiology (Bruni et al., 2012), automated driving (Figalová et al., 2024), air traffic control (Peukert et al., 2023), and flying planes (Pan et al., 2022). Beyond affecting surface-level performance on these tasks, fatigue physically affects the brain, causing reduced alertness and vigilance that can be measured directly, e.g. using EEG signals (Peukert et al., 2023). We extend L2D to be adaptive to these variations in human performance.

5 EXPERIMENTS

We evaluate our proposed EDA-L2D framework through a progression of controlled simulations. We train the two variants of EDA-L2D (per-step and RNN-based) on three benchmarks spanning distinct modalities and annotation regimes: **CIFAR-10** (image classification), **CheXpert** (chest X-ray classification) and **HateSpeech** (text moderation). Our primary baseline is “Native L2D”, meaning that we apply a traditional L2D system that has no awareness of time, collapsing all training data and presenting it to the model without an information about time ordering.

5.1 CIFAR-10

Task and Data. We study standard image classification on CIFAR-10 (Krizhevsky, 2009). The dataset contains 60,000 color images at 32×32 resolution from ten object categories (airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck), with 50,000 training images and 10,000 test images. Given an input image x , the goal is to output either a predicted class label or defer the decision to the expert. For our experiments, we first arrange the 50,000 training images into sequences and perform a 90/10 split into training and validation sets at the sequence level, ensuring that each sequence appears in exactly one set. The official 10,000-image test set is used for final evaluation.

Synthetic Expert In this experiment, we constructed a simple synthetic expert to provide a controlled baseline. Following the setup of Mozannar and Sontag (2021), we assume that the expert initially has full knowledge of all classes, i.e., $K=10$. The expert’s accuracy is then designed to decrease linearly over time from 100% accuracy at the beginning to 50% accuracy at $T=50$. At each timestep, the expert’s prediction is sampled according to a Bernoulli trial with the corresponding time-dependent accuracy, thereby generating stochastic outputs that reflect the intended accuracy trajectory.

Classifier and EDA-L2D architecture. For the classifier, we adopt WRN-28-4 with standard CIFAR-10 normalization. For deferral, we take the features output by the classifier, concatenate them with a one-step indicator of whether the expert was correct at $t-1$ and a normalized timestep indicator, and pass this concatenated vector through a 1-layer GRU with 256 hidden units. The GRU output is then fed to a fully connected deferral head that produces a single defer logit.

Training. Training proceeds in two stages. First, the backbone are optimized with cross-entropy, using SGD (momentum 0.9, weight decay 5×10^{-4}), cosine-annealed learning rate, batch size 128, and no dropout. The model minimizes cross-entropy for 200 epochs; this yields 90.27% test accuracy on CIFAR-10. We then fine-tune the entire network end-to-end using the L2D softmax surrogate loss: the classifier backbone continues training with SGD using cosine annealing over 25000 total iterations, while GRU and the two classifier heads use Adam driven by a custom cosine-then-hold schedule implemented via LambdaLR (cosine anneal for the first 25 epochs, then hold for the remaining 25 of a 50-epoch run).

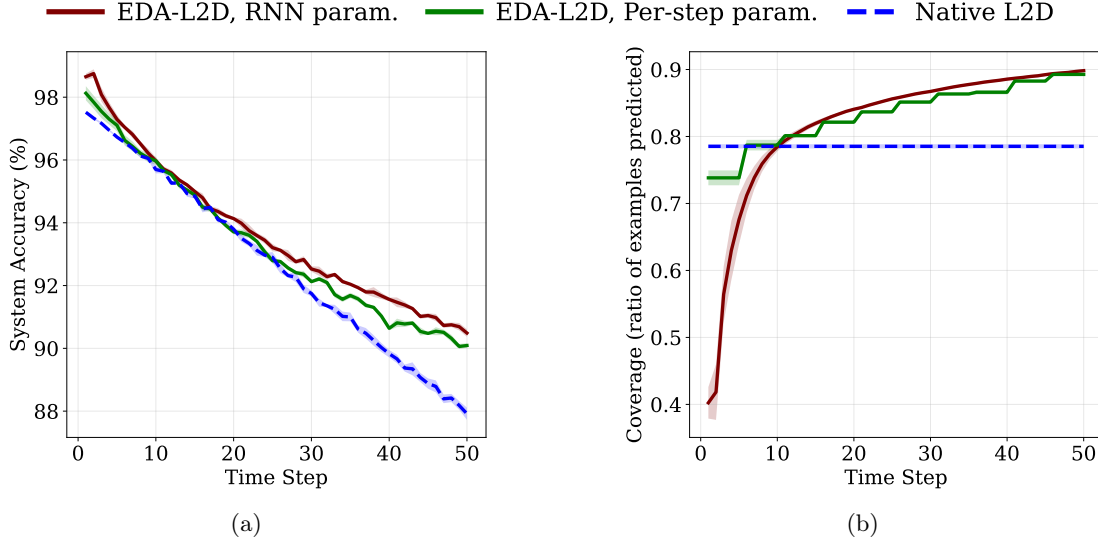


Figure 3: Plot of system accuracy (a) and coverage (proportion of classifier taken examples) (b) comparing our methods with the baseline on the CIFAR-10 test set across 50 time steps. We evaluate using a toy expert whose accuracy linearly decays from 100% to 50%. Error bars represent standard deviations over 10 runs.

Results. In this experiment we considered a simple expert model with linearly decaying accuracy from 100% to 50%. Since the classifier achieves an accuracy of 90.47%, we would expect fewer samples to be deferred to the expert once the expert’s performance falls below the classifier’s. Our results confirm this intuition: as shown in Figure 3b, at time $t=10$ our RNN-based EDA-L2D and the per-step EDA-L2D adapt their deferral rate accordingly. Moreover, Figure 3a demonstrates that EDA-L2D achieves higher overall system accuracy in both settings. These findings indicate that our method makes more appropriate allocations both before and after the crossing point, correctly assigning more samples to the expert when its accuracy is higher and shifting them to the classifier once the expert’s performance declines.

5.2 CheXpert

Task. We study chest X-ray classification on CheXpert (Irvin et al., 2019) with the standard 14 observations, framing the problem as *per-task, per-timestep* binary decisions over image sequences of length T . At each timestep $t \in \{1, \dots, T\}$ and for each task $k \in \{1, \dots, 14\}$, the system must either output a class prediction or defer to an external expert.

Expert. Following Mozannar and Sontag (2021), we simulate a class-dependent, time-varying expert with two accuracy levels, p and q . A designated confounding class receives a lower base accuracy, while all other classes receive a higher one. Over time, both accuracies decay linearly at fixed rates, with a small positive floor

to avoid vanishing performance. In our experiments, we set the base accuracies to $p=1.0$, apply decay rates of $d=0.05$, and $T=10$. The expert label is generated as $\text{Bernoulli}(p_t)$.

Data. We use the downsampled CheXpert release and keep 14 one-vs-rest targets and a binary mask that suppresses tasks marked uncertain or missing under CheXpert’s label encoding. All images are converted to three channels, normalized, and resized to ImageNet-compatible resolution; training augmentation uses random resized crops, horizontal flips, and random rotations up to 15° . Dataset splits are patient-disjoint with an 80/10/10 train/validation/test partition.

Classifier and EDA-L2D architecture. Following Irvin et al. (2019), we use a DenseNet-121 backbone initialized with ImageNet pretraining. We add two heads to the classifier backbone. The first is a per-class classifier that takes the CNN feature at time step t and outputs two logits (negative vs. positive). The second is a one-layer LSTM with 1024 hidden units that takes the CNN feature together with the expert’s previous predictions and produces one deferral logit per class. At inference, for each class, we form a three-way distribution—negative (no disease), positive (diseased), and defer—and we defer to the expert whenever the defer score exceeds the larger of the two class scores; otherwise, we use the classifier’s prediction.

Training. We train the model in two stages over a total of four epochs. In the first stage, we pre-train the CNN and classification head with standard

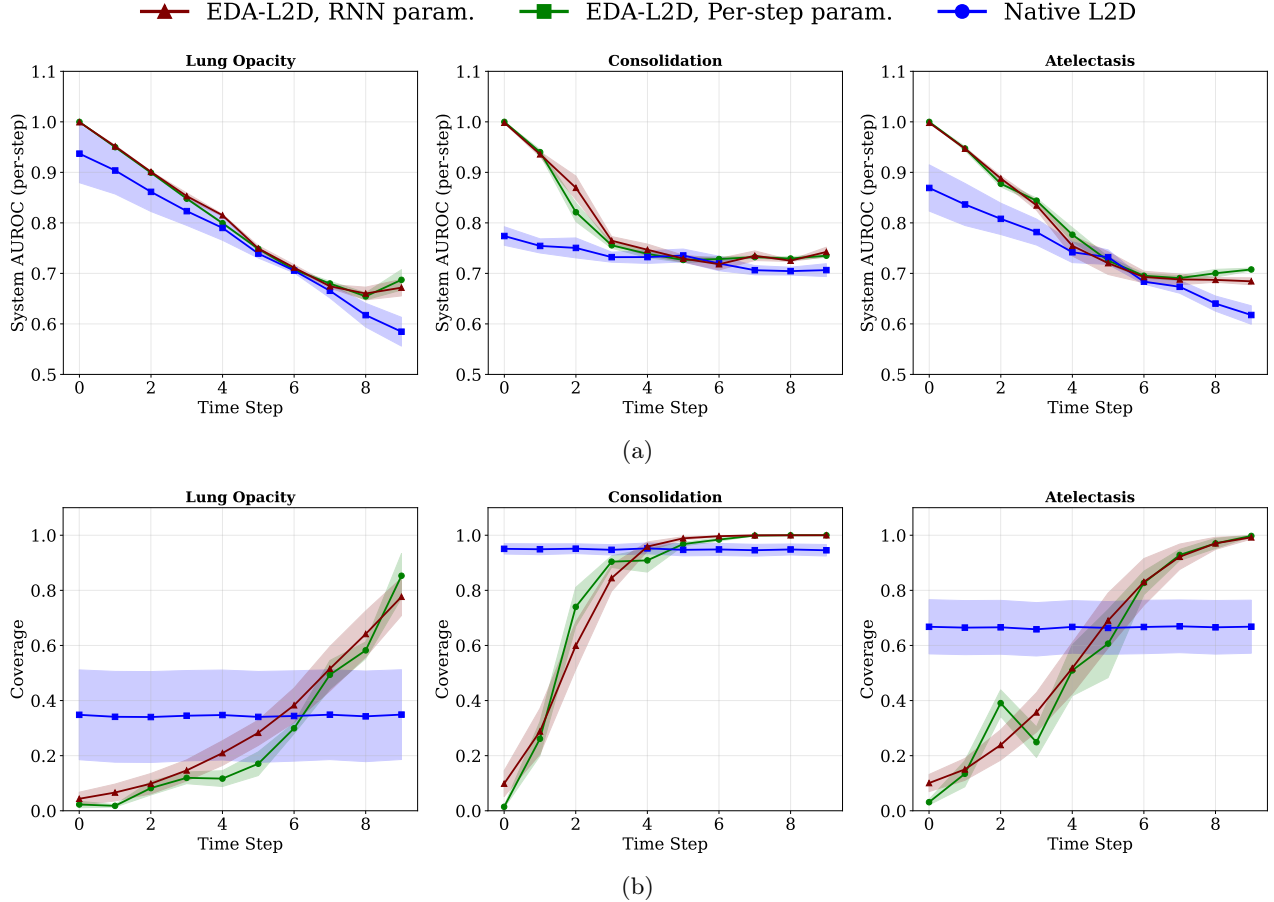


Figure 4: Plot of system accuracy (a) and coverage (ratio of classifier-taken examples) (b) comparing our methods with the baseline on the CheXpert dataset. For each timestep t , we evaluate on sequences built from a 10% patient-ID subset, using a fair expert whose accuracy decays linearly from 100% to 50%.

cross-entropy for three epochs, averaging the loss over samples. In the second stage, we fine-tune the entire network end-to-end for one epoch with L_{CE} , accumulating losses across each sequence and normalizing by its length T . All four epochs use Adam (learning rate 1×10^{-4} , weight decay 1×10^{-5}) with a Reduce-on-Plateau scheduler.

Results In Fig. 4, we compare our method with the baseline L2D approach and the per-step EDA-L2D on the CheXpert dataset with our synthetic expert. On the Lung Opacity and Atelectasis tasks, our approach clearly outperforms native L2D and performs on par with the per-step EDA-L2D; on the Consolidation task, RNN-based EDA-L2D outperforms both the per-step EDA-L2D and native L2D.

5.3 Hate Speech

We study detection of abusive content on Twitter posts using the Davidson et al. (2017) corpus of 24,783 English tweets annotated into three mutually exclusive

classes: *hate speech*, *offensive but not hate*, and *neither*.

Expert. We follow prior work that uses the Twitter-AAE lexical model to probabilistically identify tweets written in African-American English (AAE) and binarize group membership with a 0.5 threshold (Blodgett et al., 2016). We instantiate a synthetic fair expert whose accuracy is identical across demographics. At timestep t , the expert’s correctness probability decays linearly from 1.0 to 0.5 over the sequence. For each tweet, we sample a Bernoulli variable to decide whether the expert outputs the ground-truth label; if not, the expert predicts uniformly at random among the remaining classes. We construct sequences of length $T=10$ for evaluation.

Classifier architecture. Our hate-speech classifier models are lightweight TextCNNs. Tweets are tokenized with spaCy and mapped to a 25k vocabulary; embeddings are initialized with GloVe-6B (100d) and fine-tuned. The encoder applies three parallel con-

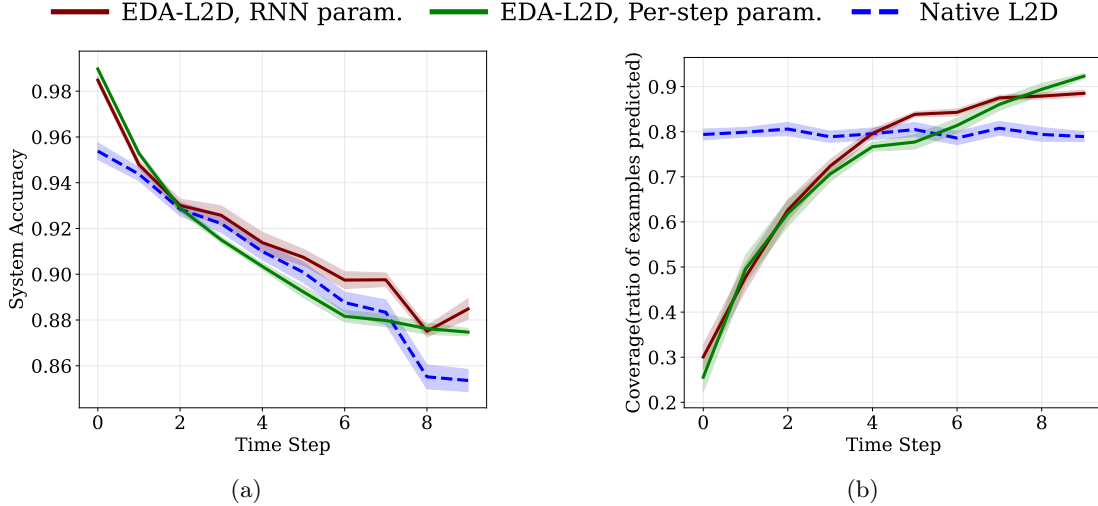


Figure 5: Plot of system accuracy (a) and coverage (ratio of classifier-taken examples) (b) comparing our methods with the baseline on the HateSpeech dataset across 10 time steps. We simulate using a toy expert whose accuracy decays linearly from 100% to 50% over 10 timesteps. Results are averaged over 10 runs, with error bars showing standard deviations.

volutional filters of widths 3, 4, and 5 (300 channels each) over the embedding matrix, followed by ReLU and max-over-time pooling. The pooled features are concatenated, passed through dropout (0.5), and fed to a linear classifier.

EDA-L2D architecture. We employ a TextCNN encoder to obtain a per-timestep feature representation for each tweet in a length- T sequence. To model temporal dependencies and expert reliability, we concatenate a binary indicator of the expert’s correctness at the previous timestep to the CNN feature and feed the resulting vector into a single-layer LSTM with 256 hidden units. The TextCNN encoder feeds a fully connected classification head that outputs three logits (hate, offensive, neither), while the LSTM feeds a fully connected deferral head that outputs a single defer logit. Concatenating these yields a four-way prediction over hate, offensive, neither, defer. At inference, we defer to the expert when the defer logit exceeds the maximum of the three class logits; otherwise, the classifier’s prediction is returned.

Results. From Fig. 5a, we observe that our model trained on the HateSpeech dataset demonstrates a clear advantage. In terms of system accuracy, our method consistently outperforms both the native L2D and per-step EDA-l2d approach, with an average margin of 1.26% and 0.70%, respectively. Moreover, from Fig. 5b, we observe that our method identifies more favorable timesteps than the per-step model at which the system coverage surpasses that of the general baseline, resulting in higher system accuracy at those points.

6 CONCLUSION & FUTURE WORK

We propose *Expert-Drift-Adapted Learning to Defer (EDA-L2D)*, a framework that enables improved generalization to real human experts by explicitly modeling temporal drift in their performance—variance that is common and caused by a wide variety of factors, such as fatigue and cognitive biases. We achieve this by explicitly incorporating time and prediction history into the L2D model, allowing the model to learn representations of the expert that vary with each time-step and dynamically adapt to the expert’s behavior to achieve improved deferral decisions and greater human-AI team performance. In future work, we aim to extend the framework to capture complex cognitive biases, such as the availability heuristic (Tversky and Kahneman, 1973), gambler’s fallacy (Kovic and Kristiansen, 2019), delay discounting (Kurth-Nelson et al., 2012), and the recency bias (Turvey and Freeman, 2012).

References

Zeynep Akata, Dan Balliet, Maarten de Rijke, Frank Dignum, Virginia Dignum, A. E. Eiben, Antske Fokkens, Davide Grossi, Koen V. Hindriks, Holger Hoos, Haley Hung, Catholijn Jonker, Christof Monz, Mark Neerincx, Frans Oliehoek, Henry Prakken, Stefan Schlobach, Linda C. van der Gaag, Frank van Harmelen, Herke van Hoof, M. Birna van Riemsdijk, Aimee van Wynsberghe, Rineke Verbrugge, Bart Verheij, Paul Vossen, and Max Welling. A research agenda for hybrid intelligence: Augmenting human intellect with collaborative, adaptive, responsible,

- and explainable artificial intelligence. *Computer*, 53(8):18–28, 2020. doi: 10.1109/MC.2020.2996587. URL <https://www.computer.org/csdl/magazine/co/2020/08/09153877/11UB5gL2CnS>.
- Su Lin Blodgett, Lisa Green, and Brendan O’Connor. Demographic dialectal variation in social media: A case study of african-american english, 2016. URL <https://arxiv.org/abs/1608.08868>.
- Silvio G Bruni, Eric Bartlett, and Eugene Yu. Factors involved in discrepant preliminary radiology resident interpretations of neuroradiological imaging studies: a retrospective analysis. *American Journal of Roentgenology*, 198(6):1367–1374, 2012.
- Yuzhou Cao, Hussein Mozannar, Lei Feng, Hongxin Wei, and Bo An. In defense of softmax parametrization for calibrated and consistent learning to defer, 2023. URL <https://arxiv.org/abs/2311.01106>.
- Mohammad-Amin Charusaie, Hussein Mozannar, David Sontag, and Samira Samadi. Sample efficient learning of predictors that complement humans, 2022. URL <https://arxiv.org/abs/2207.09584>.
- C. K. Chow. An optimum character recognition system using decision functions. *IRE Transactions on Electronic Computers*, EC-6(4):247–254, 1957. doi: 10.1109/TEC.1957.5222035.
- Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. Automated hate speech detection and the problem of offensive language, 2017. URL <https://arxiv.org/abs/1703.04009>.
- Dominik Dellermann, Philipp Ebel, Matthias Söllner, and Jan Marco Leimeister. Hybrid intelligence. *Business & Information Systems Engineering*, 61(5):637–643, 2019. doi: 10.1007/s12599-019-00595-2. URL <https://link.springer.com/article/10.1007/s12599-019-00595-2>.
- Nikol Figalová, Hans Joachim Bieg, Michael Schulz, Jürgen Pichen, Martin Baumann, Lewis Chuang, and Olga Pollatos. Fatigue and mental underload further pronounced in l3 conditionally automated driving: Results from an eeg experiment on a test track. *ArXiv preprint arXiv:2405.18114*, 2024. L3 supervisor sleepiness increases with prolonged monitoring; EEG measures show drift in alertness.
- Patrick Hemmer, Lukas Thede, Michael Vössing, Johannes Jakubik, and Niklas Köhl. Learning to defer with limited expert predictions, 2023. URL <https://arxiv.org/abs/2304.07306>.
- Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison, 2019. URL <https://arxiv.org/abs/1901.07031>.
- Shalmali Joshi, Sonali Parbhoo, and Finale Doshi-Velez. Learning-to-defer for sequential medical decision-making under uncertainty. *Transactions on Machine Learning Research*, 2023. URL <https://arxiv.org/abs/2109.06312>.
- C. Peterson Joshua M. Battleday Ruairidh L. Griffiths Thomas and Olga Russakovsky. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9617–9626, 2019. doi: 10.1109/ICCV.2019.00971. URL <https://arxiv.org/abs/1908.07086>.
- Ece Kamar. Directions in hybrid intelligence: Complementing ai systems with human intelligence. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI-16)*, 2016. URL <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/HybridIntelligence.pdf>. Position paper outlining research directions for human-AI systems.
- Marko Kovic and Silje Kristiansen. The gambler’s fallacy fallacy (fallacy). *Journal of Risk Research*, 22(3):291–302, 2019. doi: 10.1080/13669877.2017.1378248. URL <https://doi.org/10.1080/13669877.2017.1378248>.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>. Technical Report.
- Zeb Kurth-Nelson, Warren Bickel, and A. David Redish. A theoretical account of cognitive effects in delay discounting. *European Journal of Neuroscience*, 35(7):1052–1064, 2012. doi: <https://doi.org/10.1111/j.1460-9568.2012.08058.x>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1460-9568.2012.08058.x>.
- Eric Legler, Darío Cuevas Rivera, Sarah Schwöbel, Ben J. Wagner, and Stefan Kiebel. Cognitive computational model reveals repetition bias in a sequential decision-making task. *Communications Psychology*, 3(1):92, 2025. doi: 10.1038/s44271-025-00271-0.
- David Madras, Toniann Pitassi, and Richard Zemel. Predict responsibly: Improving fairness and accuracy by learning to defer, 2018. URL <https://arxiv.org/abs/1711.06664>.
- Hussein Mozannar and David Sontag. Consistent esti-

- mators for learning to defer to an expert, 2021. URL <https://arxiv.org/abs/2006.01862>.
- Hussein Mozannar, Hunter Lang, Dennis Wei, Prasanna Sattigeri, Subhro Das, and David Sonntag. Who should predict? exact algorithms for learning to defer to humans, 2023. URL <https://arxiv.org/abs/2301.06197>.
- Harikrishna Narasimhan, Wittawat Jitkrittum, Aditya K Menon, Ankit Rawat, and Sanjiv Kumar. Post-hoc estimators for learning to defer to an expert. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 29292–29304. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/bc8f76d9caadd48f77025b1c889d2e2d-Paper-Conference.pdf.
- T. Pan et al. Research on the identification of pilots’ fatigue status based on fnirs signals and deep models. *Aerospace*, 9(3):173, 2022. doi: 10.3390/aerospace9030173. Identifies pilot fatigue using brain signals; performance drift with fatigue.
- M. Peukert et al. Subjective and objective fatigue dynamics in air traffic control. *Journal of Sleep Research*, 2023. Showing increase in subjective fatigue and drop in vigilance toward end of evening shifts for air traffic controllers.
- Dharmesh Tailor, Aditya Patra, Rajeev Verma, Putra Manggala, and Eric Nalisnick. Learning to defer to a population: A meta-learning approach, 2024. URL <https://arxiv.org/abs/2403.02683>.
- B.E. Turvey and J.L. Freeman. Jury psychology. In V.S. Ramachandran, editor, *Encyclopedia of Human Behavior (Second Edition)*, pages 495–502. Academic Press, San Diego, second edition edition, 2012. ISBN 978-0-08-096180-4. doi: <https://doi.org/10.1016/B978-0-12-375000-6.00216-0>. URL <https://www.sciencedirect.com/science/article/pii/B9780123750006002160>.
- Amos Tversky and Daniel Kahneman. Availability: A heuristic for judging frequency and probability. *Cognitive Psychology*, 5(2):207–232, 1973. ISSN 0010-0285. doi: [https://doi.org/10.1016/0010-0285\(73\)90033-9](https://doi.org/10.1016/0010-0285(73)90033-9). URL <https://www.sciencedirect.com/science/article/pii/0010028573900339>.
- Anne E. Urai, Jan Willem de Gee, Konstantinos Tsetos, and Tobias H. Donner. Choice history biases subsequent evidence accumulation. *eLife*, 8:e46331, 2019. doi: 10.7554/eLife.46331.
- Didrika S. van de Wouw, Ryan T. McKay, and Nicholas Furl. Biased expectations about future choice options predict sequential economic decisions. *Communications Psychology*, 2(1):119, 2024. doi: 10.1038/s44271-024-00172-8.
- Rajeev Verma and Eric Nalisnick. Calibrated learning to defer with one-vs-all classifiers, 2022. URL <https://arxiv.org/abs/2202.03673>.
- 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.

Checklist

- For all models and algorithms presented, check if you include:
 - A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes]
 - An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes]
 - (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes]
- For any theoretical claim, check if you include:
 - Statements of the full set of assumptions of all theoretical results. [Not Applicable]
 - Complete proofs of all theoretical results. [Not Applicable]
 - Clear explanations of any assumptions. [Not Applicable]
- For all figures and tables that present empirical results, check if you include:
 - The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes]
 - All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes]
 - A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes]
 - A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Yes]
- If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:

- (a) Citations of the creator If your work uses existing assets. [Yes]
 - (b) The license information of the assets, if applicable. [Not Applicable]
 - (c) New assets either in the supplemental material or as a URL, if applicable. [Yes]
 - (d) Information about consent from data providers/curators. [Not Applicable]
 - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
- (a) The full text of instructions given to participants and screenshots. [Not Applicable]
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable]
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable]

Appendix

A THEORETICAL ANALYSIS

We study deferral policies that adapt to systematic patterns in human decisions rather than assuming static or i.i.d. performance. Human experts show predictable temporal and contextual effects such as fatigue, recency, group bias, and, crucially, expert drift that alters accuracy over a sequence. We analyze simple settings that isolate these dynamics and measure their impact on joint performance. When accuracy varies in structured ways—through monotone drift or group bias—deferral rules that model this variation outperform policies based on a single average accuracy. This motivates adaptive, human-aware deferral that updates in real time with the collaborator’s reliability and forms the basis for our three L2D formulations: Native L2D, which assumes static expert behavior; EDA-L2D with Per-Step parameterization, which models timestep-specific variations; and EDA-L2D with RNN parameterization, which captures temporal dependencies and evolving expert drift.

A.1 Deferral Mechanisms

Static. A single, time-invariant deferral rule that either always defers to the human or never defers, determined from a global average estimate of human accuracy. This matches the spirit of our Native L2D configuration, which is time-agnostic and does not adapt its deferral behavior within a sequence.

Moving Average. A time-indexed rule that estimates human accuracy at step i using only the prefix average of past outcomes and defers when this estimate exceeds the model’s accuracy. This mirrors our EDA-L2D per-step parameterization, where independently trained timestep models provide adaptation but respond to drift with delay and without sharing information across steps.

Dynamic. An oracle policy that knows the human’s instantaneous accuracy at each step and defers exactly when the human currently outperforms the model. This is the behavior our EDA-L2D RNN parameterization seeks to approximate, using shared temporal parameters to track nonstationarity and implement a real-time deferral policy over the sequence.

A.2 Deferring with a Static Accuracy Estimate vs. Perfect Dynamic Model

We first show that a dynamic deferral strategy yields strictly lower expected loss than a static one when the human’s performance degrades monotonically due to decision fatigue.

Setup. Consider a human H and a robot R with constant per-example loss ℓ_R . Let ℓ_{H^t} be the human’s loss at timestep t , evaluated over T timesteps. Assume monotone degradation

$$\ell_{H^1} < \ell_{H^2} < \dots < \ell_{H^T},$$

and a nontrivial crossing

$$\ell_{H^1} < \ell_R < \ell_{H^T}.$$

Let $\overline{\ell_{H^T}} = \frac{1}{T} \sum_{n=1}^T \ell_{H^n}$ denote the human’s average loss. We compare:

- **Static:** always defer if $\overline{\ell_{H^T}} < \ell_R$, else never defer.
- **Dynamic:** defer on timestep i iff $\ell_{H^i} < \ell_R$.

Let $i = \min\{k : \ell_{H^k} > \ell_R\}$. The dynamic policy defers on $1, \dots, i-1$ and assigns i, \dots, T to R , giving

$$\ell_D = \frac{1}{T} \left(\sum_{t=1}^{i-1} \ell_{H^t} + (N-i+1)\ell_R \right).$$

Case 1: $\overline{\ell_{H_T}} < \ell_R$. Static always defers: $\ell_S = \overline{\ell_{H_T}} = \frac{1}{T} \left(\sum_{t=1}^{i-1} \ell_{H^t} + \sum_{t=i}^T \ell_{H^t} \right)$. For $t < i$, both policies incur ℓ_{H^t} . For $t \geq i$, static incurs $\ell_{H^k} \geq \ell_R$ while dynamic incurs ℓ_R . Hence termwise $\ell_D \leq \ell_S$ with strict inequality whenever some $\ell_{H^t} > \ell_R$ for $t \geq i$. Thus $\ell_D < \ell_S$.

Case 2: $\overline{\ell_{H_T}} \geq \ell_R$. Static never defers: $\ell_S = \ell_R$. Dynamic replaces a nonempty prefix $\{1, \dots, i-1\}$ of ℓ_R by strictly smaller $\ell_{H^t} < \ell_R$, so $\ell_D < \ell_S$.

Conclusion. In both cases,

$$\ell_D < \ell_S.$$

Hence, time-aware deferral strictly improves over time-agnostic deferral under monotone fatigue.

A.3 Deferring with a Simple Moving Average vs. Perfect Dynamic Model

We now compare the dynamic policy to a moving-average (MA) policy.

Setup. Define the MA estimate from past observations only:

$$\hat{\ell}_{H^t} = \frac{1}{t-1} \sum_{j=1}^{t-1} \ell_{H^j} \quad (t \geq 2),$$

with any fixed base for $t=1$. The MA policy defers at step t if $\hat{\ell}_{H^t} < \ell_R$. Let

$$j = \min\{t : \ell_{H^t} \geq \ell_R\}$$

be the first step where the human's instantaneous loss reaches the robot's loss.

Switch order ($j < t$). At $t = j$, the MA estimate averages *only* $\{\ell_{H^1}, \dots, \ell_{H^{j-1}}\}$, all strictly $< \ell_R$ by definition of j . Hence $\hat{\ell}_{H^j} < \ell_R$, so MA still defers at j and therefore $t > j$. Thus MA switches strictly *after* the dynamic policy.

Comparison of total losses. Dynamic incurs ℓ_{H^t} for $t < j$ and ℓ_R for $t \geq j$:

$$\ell_D = \frac{1}{N} \left(\sum_{t=1}^{j-1} \ell_{H^t} + (N-j+1)\ell_R \right).$$

MA defers at least until $t > j$, hence assigns ℓ_{H^t} to a strictly longer prefix and uses ℓ_R on a strictly shorter suffix:

$$\ell_{MA} = \frac{1}{N} \left(\sum_{k=1}^{t-1} \ell_{H^k} + (N-i+1)\ell_R \right), \quad i > j.$$

Since $\ell_{H^k} < \ell_R$ for $k < j$ while $\ell_{H^k} \geq \ell_R$ for $k \geq j$, replacing any ℓ_{H^k} with ℓ_R on $k \in \{j, \dots, t-1\}$ strictly decreases loss. Therefore $\ell_D < \ell_{MA}$.

Conclusion.

$$\ell_D < \ell_{MA}.$$

Even an adaptive MA policy that updates from past outcomes lags the true (nonstationary) human performance and is outperformed by a time-aware dynamic policy.

Remark. The above argument also covers fixed-window moving averages. For any window not containing the crossover j , fixed-window MA and the dynamic policy act identically. The only difference arises when the window starts to include post- j outcomes, which can only delay the fixed-window MA switch relative to j , yielding strictly larger total loss.

Discussion. These results substantiate the progression of practical parameterizations: a time-agnostic (Native/Static) policy is dominated by time-aware deferral; per-step policy modeling—akin to moving average policy—offers only coarse adaptivity and still lags, except for degenerate cases where the window happens to straddle the crossover point j ; and shared-parameter temporal models (e.g., RNNs) are natural approximators of the ideal dynamic policy that tracks instantaneous human reliability.

B EXPERIMENT DETAILS

B.1 Training Configuration for CIFAR-10 Experiments

Native L2D. We use a WideResNet- 28×4 with an additional output unit for deferral, yielding $N+1$ logits for an N -class problem. The model is trained with a softmax surrogate loss using SGD (momentum 0.9, Nesterov, weight decay 5×10^{-4}), a cosine-annealed learning rate, batch size 128, and 200 training epochs.

EDA-L2D, Per-Step Param. The training recipe is identical to the native L2D setup; the only change is that we train ten independent models, each responsible for a contiguous 5-timestep window of the full $T=50$ horizon, while keeping the same optimization hyperparameters and preprocessing.

EDA-L2D, RNN Param. We use a two-stage procedure. In the first stage, a WideResNet- 28×4 backbone is trained following the native L2D configuration: SGD with momentum 0.9 and Nesterov, weight decay 5×10^{-4} , cosine-annealed learning rate, batch size 128, for 200 epochs. In the second stage, we attach a single-layer GRU (hidden size 256, dropout 0.1) and classification and deferral heads to the pretrained backbone and train the full model jointly. The backbone is optimized with SGD (learning rate 0.05, momentum 0.9, Nesterov, weight decay 5×10^{-4}) under a CosineAnnealingLR scheduler updated every iteration, while the GRU and linear heads use Adam (learning rate 0.05, $\beta_1=0.9$, $\beta_2=0.999$) with a LambdaLR that follows a cosine-then-hold schedule (cosine decay until 25 epochs remain, then held at a minimum factor of 10^{-3}). This stage runs for 50 epochs. The dual-scheduler design allows the convolutional backbone to keep refining visual features as the GRU learns the temporal deferral policy.

B.2 Training Configuration for CheXpert Experiments

Native L2D. Training follows a two-stage procedure. First, we pretrain the Dense121 networks without deferral, optimizing the standard classification cross-entropy loss to establish a strong task-specific backbone and head. We then finetune with deferral by replacing the head with three logits per class and optimizing softmax surrogate loss L_{CE} , thereby learning to allocate decisions between the model and the expert. Unless otherwise noted, both stages use Adam (learning rate 1×10^{-4} , $\beta_1 = 0.9$, $\beta_2 = 0.999$, weight decay 1×10^{-5}), together with a ReduceLROnPlateau scheduler (factor 0.1, patience 2, monitored on validation loss).

EDA-L2D, Per-Step Param. We train ten independent models, each dedicated to a specific timestep (or fixed-length window). Every model adheres to the same two-stage protocol and loss design as above—standard cross-entropy during pretraining and L_{CE} during finetuning—using the identical Adam configuration and ReduceLROnPlateau schedule.

EDA-L2D, RNN Param. We retain the two-stage protocol but parameterize the policy time-wise over sequences. Pretraining proceeds without deferral using sequence-averaged cross-entropy on the two-class head; finetuning enables deferral with the three-logit head and minimizes the sequence-averaged L_{CE} . Optimization mirrors the configurations above (Adam with the same hyperparameters and ReduceLROnPlateau on validation loss). This time-wise parameterization shares parameters across t and explicitly learns a temporal deferral policy.

B.3 Training Configuration for Hate Speech Experiments

For native L2D, perstep EDA-L2D and RNN-based EDA-L2D, optimization uses the Adam optimizer with its default hyperparameters, trained for 10 epochs with a batch size of 64. Throughout training, we monitor validation performance and retain the best-performing checkpoint based on validation system accuracy.

Table 1: Comparison of three deferral strategies on CIFAR-10H amplified human performance curves ($T=200$, moving-average window size 20, resampled 100 times). We report mean \pm standard error for system accuracy and coverage.

Method	System Accuracy (%)	Coverage (%)	Δ to the Best System Accuracy(%)
Static	89.97 ± 0.06	0.00 ± 0.00	-0.32
Moving Average	90.14 ± 0.08	80.06 ± 5.43	-0.15
Dynamic	90.29 ± 0.06	55.00 ± 0.00	0.00

C ADDITIONAL EXPERIMENTS

C.1 CIFAR-10H

Preprocessing and Accuracy Curve Construction. To stress-test our theoretical claims in a realistic annotation regime, we run an additional study on CIFAR-10H. For each annotator, we rescale trial indices to a common unit interval $t \in [0, 1]$ and partition this interval into B equal bins. Within bin b , we compute the empirical accuracy r_b and use the bin size w_b as a weight. As a noise-robust reference, we fit an isotonic regression with a non-increasing constraint, yielding a monotone accuracy curve over t .

Parametric Families and Model Selection. We then fit simple, interpretable families that capture distinct decline profiles: exponential decay $p(t) = L + Se^{-kt}$, power-law decay $p(t) = L + S(1 + \gamma t)^{-\beta}$, logistic falloff $p(t) = L + \frac{S}{1 + \exp\{+k(t-t_0)\}}$. All parametric fits are constrained to the probability simplex ($0 \leq p(t) \leq 1$ for all t) and regularized to be non-increasing; when needed, we apply a post-fit monotone projection to remove residual violations. Parameters are estimated by weighted least squares using $\{(u_b, r_b, w_b)\}_{b=1}^B$, and model selection is based on Akaike Information Criterion (AIC). Across annotators, the exponential form most often provides the best balance between parsimony and fit quality. We selected power fit as our best model.

Affine Rescaling. To align this curve with our target operating range and the theoretical setup, we begin by assuming the CIFAR-10 classifier attains accuracy 89.5%. To ensure the simulated annotator’s instantaneous accuracy crosses this level—creating the regime where dynamic deferral can improve over static or moving-average policies—we first fit the power-law curve and then apply an affine rescaling that maps the first and last bin values to 100% and 88%. The resulting curve remains monotone nonincreasing, with its maximum above 89.5% and minimum below 89.5%.

Simulation of Deferral Policies. For each resampled human-accuracy trajectory, we run T simulated trials with Bernoulli outcomes and repeat over N resamples. At each step, we evaluate three policies: a dynamic policy that defers when the fitted instantaneous human accuracy from our curve model exceeds the model’s accuracy (the fitted curve also governs the simulated human outcomes); a moving-average policy that decides using the running mean of realized human outcomes up to the current step; and a static policy that compares the model to the trajectory’s global mean human accuracy and therefore either always defers or never defers. As shown in Table 1, the dynamic strategy achieves superior performance on the amplified human performance curves of CIFAR-10H, surpassing MA by 0.15% and Static by 0.32%, which aligns well with our theoretical analysis.

C.2 Synthetic Data

We construct a controlled sequential Gaussian-mixture environment to induce time-varying expert reliability across two sub-populations $A \in \{0, 1\}$. The covariates lie in $\mathcal{X} = \mathbb{R}^d$ with binary targets $Y \in \{0, 1\}$. For each group a and class y , the class-conditional distribution is $X | (Y=y, A=a) \sim \mathcal{N}(\mu_{y,a}, \Sigma_{y,a})$ with $\mu_{y,a} \in [0, 10]^d$ and diagonal $\Sigma_{y,a}$ whose entries are drawn from $\text{Unif}(0, 10)$. Sequences of length T evolve by adding Gaussian noise and a group-specific drift, and labels may flip with small probability to model non-stationarity.

Setup Across 200 trials we sample the group proportion $P(A=1) \sim \text{Unif}(0.02, 0.98)$, fix $d=10$ and $T=10$, and generate $N=6000$ sequences per trial. Within each sequence, features evolve via $X_t = X_{t-1} + \varepsilon_t + \Delta_a$, with $\varepsilon_t \sim \mathcal{N}(0, \sigma^2 I)$ and group drift Δ_a (smaller for $A=0$, larger for $A=1$). Labels follow $Y_t = Y_{t-1}$ with probability $1 - p_{\text{flip}}$ and flip otherwise. Expert predictions E_t are correct with a time-decayed accuracy $p_a(t) = p_{0,a} \exp(-\lambda_a t/T)$,

Table 2: Comparing L2D variants trained with linear models on synthetic data. We report mean \pm standard error of system accuracy, coverage, and classifier accuracy over steps.

Method	System Accuracy(%)	Coverage(%)	Δ to the Best System Accuracy(%)
Native L2D	71.58 ± 0.43	53.69 ± 0.40	-1.38
EDA-L2D, Per-step param.	71.79 ± 0.27	54.21 ± 0.25	-1.17
EDA-L2D, RNN param.	72.96 ± 0.36	53.71 ± 0.77	0.00

yielding a slower decay for $A=0$ and a faster decay for $A=1$. We train with Adam (learning rate 10^{-3} , weight decay 10^{-5}) for 5 epochs per model using a 70/15/15 train/val/test split, and aggregate test metrics over trials.

Models and baselines. We compare three heads that map each time step’s feature vector to three logits over $\{0, 1, \perp\}$. (i) *Native L2D*: a single time-invariant linear head shared across all steps, providing strong parameter sharing and a stationary inductive bias. (ii) *EDA-L2D, Per-step param.*: T independent linear heads for per step, capturing nonstationarity at the cost of T-fold parameters. (iii) *EDA-L2D, RNN param.*: a factorized design with a time-local classifier branch and a sequence-aware deferral branch that feeds $[x_t, e_{t-1}]$ into a GRU followed by a linear head for the defer logit; the three logits are concatenated to form the final output at each step.

Results. Table 2 reports the comparative performance of the three L2D variants in the synthetic sequential environment. The time-agnostic Native L2D serves as a stationary baseline and attains the lowest system accuracy. Per-step EDA-L2D yields a modest gain by capturing coarse nonstationarity, albeit without temporal smoothing. The RNN-based EDA-L2D achieves the best accuracy (+1.38% over Native L2D), confirming the benefit of inserting a sequence-aware module to drive the deferral head in L2D.