052

053

054

000 001

002 003

A Two-Stage Learning-to-Defer Approach for Multi-Task Learning

Anonymous Authors¹

Abstract

The Two-Stage Learning-to-Defer framework has been extensively studied for classification and, more recently, regression tasks. However, many contemporary applications involve both classification and regression in an interdependent manner. In this work, we introduce a novel Two-Stage Learning-to-Defer framework for multi-task learning that jointly addresses these tasks. Our approach leverages a two-stage surrogate loss family, which we prove to be both (\mathcal{G}, \mathcal{R})-consistent and Bayes-consistent, providing strong theoretical guarantees of convergence to the Bayes-optimal rejector. We establish consistency bounds explicitly linked to the cross-entropy surrogate family and the L_1 -norm of the agents' costs, extending the theoretical minimizability gap analysis to the two-stage setting with multiple experts. We validate our framework on two challenging tasks: object detection, where classification and regression are tightly coupled, and existing methods fail, and electronic health record analysis, in which we highlight the suboptimality of current learning-todefer approaches.

1. Introduction

Learning-to-Defer (L2D) integrates predictive models with human experts—or, more broadly, decision-makers—to optimize systems requiring high reliability (Madras et al., 2018). This approach benefits from the scalability of machine learning models and leverages expert knowledge to address complex queries (Hemmer et al., 2021). The Learningto-Defer approach defers decisions to experts when the learning-based model has lower confidence than the most confident expert. This deference mechanism enhances safety, which is particularly crucial in high-stakes scenarios (Mozannar & Sontag, 2020; Mozannar et al., 2023). For example, in medical diagnostics, the system utilizes patientacquired data to deliver an initial diagnosis (Johnson et al., 2023; 2016). If the model is sufficiently confident, its diagnosis is accepted; otherwise, the decision is deferred to a medical expert who provides the final diagnosis. Such tasks, which can directly impact human lives, underscore the need to develop reliable systems (Balagurunathan et al., 2021).

Learning-to-Defer has been extensively studied in classification problems (Madras et al., 2018; Verma et al., 2022; Mozannar & Sontag, 2020; Mozannar et al., 2023; Mao et al., 2023a) and, more recently, in regression scenarios (Mao et al., 2024e). However, many modern complex tasks involve both regression and classification components, requiring deferral to be applied to both components simultaneously, as they cannot be treated independently. For instance, in object detection, a model predicts both the class of an object and its location using a regressor, with these outputs being inherently interdependent (Girshick, 2015; Redmon et al., 2016; Buch et al., 2017). In practice, deferring only localization or classification is not meaningful, as decisionmakers will treat these two tasks simultaneously. A failure in either component-such as misclassifying the object or inaccurately estimating its position-can undermine the entire problem, emphasizing the importance of coordinated deferral strategies that address both components jointly.

This potential for failure underscores the need for a Learning-to-Defer approach tailored to multi-task problems involving both classification and regression. We propose a novel framework for multi-task environments, incorporating expertise from multiple experts and the predictor-regressor model. We focus our work on the two-stage scenario, where the model is already trained offline. This setting is relevant when retraining from scratch the predictor-regressor model is either too costly or not feasible due to diverse constraints such as non-open models (Mao et al., 2023a; 2024e). We approximate the true deferral loss using a surrogate deferral loss family, based on cross-entropy, and tailored for the twostage setting, ensuring that the loss effectively approximates the original discontinuous loss function. Our theoretical analysis establishes that our surrogate loss is both $(\mathcal{G}, \mathcal{R})$ consistent and Bayes-consistent. Furthermore, we study and generalize results on the minimizability gap for deferral loss based on cross-entropy, providing deeper insights into its optimization properties.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

055 Our contributions are as follows:

(i) Novelty: We introduce two-stage Learning-to-Defer for multi-task learning with multiple experts. Unlike previous L2D methods that focus solely on classification or regression, our approach addresses both tasks in a unified framework.

062 (ii) Theoretical Foundation: We prove that our surrogate 063 family is both Bayes-consistent and $(\mathcal{G}, \mathcal{R})$ -consistent for 064 any cross-entropy-based surrogate. We derive tight consis-065 tency bounds that depend on the choice of the surrogate and 066 the L_1 -norm of the cost, extending minimizability gap anal-067 ysis to the two-stage, multi-expert setting. Additionally, we 068 establish learning bounds for the true deferral loss, show-069 ing that generalization improves as agents become more 070 accurate.

(iii) **Empirical Validation**: We evaluate our approach on two challenging tasks. In object detection, our method effectively captures the intrinsic interdependence between classification and regression, overcoming the limitations of existing L2D approaches. In EHR analysis, we show that current L2D methods struggle when agents have varying expertise across classification and regression—whereas our method achieves superior performance.

2. Related Works

071

074

075

076

078

079

081

082

083

085

086

087

088

089

090

091

Learning-to-Defer builds on the foundational ideas of Learning with Abstention (Chow, 1970; Bartlett & Wegkamp, 2008; Cortes et al., 2016; Geifman & El-Yaniv, 2017; Ramaswamy et al., 2018; Cao et al., 2022; Mao et al., 2024a), where the primary goal is to reject inputs when the model lacks sufficient confidence. L2D extends this framework by incorporating a comparison between the model's confidence and the confidence of experts.

092 One-stage L2D. Learning-to-Defer was first introduced 093 by Madras et al. (2018), who proposed a pass function 094 for binary classification, inspired by the predictor-rejector 095 framework of Cortes et al. (2016). Extending this concept to 096 the multi-class setting, Mozannar & Sontag (2020) proposed 097 a score-based approach and demonstrated that employing a 098 log-softmax multi-classification surrogate ensures a Bayes-099 consistent loss. Several subsequent works have further ad-100 vanced or applied this methodology in classification tasks (Verma et al., 2022; Cao et al., 2024; 2022; Keswani et al., 2021; Kerrigan et al., 2021; Hemmer et al., 2022; Benz & Rodriguez, 2022; Tailor et al., 2024; Liu et al., 2024; 104 Montreuil et al., 2024). A seminal contribution by Mozan-105 nar et al. (2023) identified limitations in prior approaches 106 (Mozannar & Sontag, 2020; Verma et al., 2022), highlighting their suboptimality under realizable distributions. The 108 authors argued that Bayes-consistency, while foundational, 109

may not be the most reliable criterion in settings with a *restricted* hypothesis set. To address this, they proposed *hypothesis-consistency* as a more appropriate criterion in such scenarios. Building upon this paradigm, subsequent advances in hypothesis-consistency theory (Long & Servedio, 2013; Zhang & Agarwal, 2020; Awasthi et al., 2022; Mao et al., 2023b) have further refined the theoretical underpinnings of L2D. Notably, Mao et al. (2024c) established that the general score-based formulation of classification L2D is \mathcal{H} -consistent. Mao et al. (2024d) introduced a loss function achieving *realizable-consistency*, ensuring optimal performance under realizable distributions. Moreover, L2D has been successfully extended to regression tasks, with Mao et al. (2024e) demonstrating its applicability in multi-expert deferral settings.

Two-stage L2D. The increasing prominence of large pretrained models has spurred interest in applying L2D to settings where agents (model and experts) are trained offline, reflecting the practical reality that most users lack the resources to train such models from scratch (Mao et al., 2023a; Montreuil et al., 2024). Charusaie et al. (2022) compared one-stage (online predictor) and two-stage (offline predictor) L2D approaches, identifying trade-offs between the two strategies. Mao et al. (2023a) proposed a predictor-rejector framework for two-stage L2D that guarantees both Bayesconsistency and hypothesis-consistency. Beyond classification, Mao et al. (2024e) extended the two-stage L2D framework to regression problems.

Despite significant progress, current two-stage L2D research largely addresses classification and regression independently. However, many contemporary tasks involve both regression and classification components, necessitating their joint optimization. In this work, we extend two-stage L2D to joint classifier-regressor models, addressing this critical gap.

3. Preliminaries

Multi-task scenario. We consider a multi-task setting encompassing both classification and regression problems. Let \mathcal{X} denote the input space, $\mathcal{Y} = \{1, \ldots, n\}$ represent the set of n distinct classes, and $\mathcal{T} \subseteq \mathbb{R}$ denote the space of real-valued targets for regression. For compactness, each data point is represented as a triplet $z = (x, y, t) \in \mathcal{Z}$, where $\mathcal{Z} = \mathcal{X} \times \mathcal{Y} \times \mathcal{T}$. We assume the data is sampled independently and identically distributed (i.i.d.) from a distribution \mathcal{D} over \mathcal{Z} (Girshick, 2015; Redmon et al., 2016; Carion et al., 2020).

We define a *backbone* $w \in W$, or shared feature extractor, such that $w : \mathcal{X} \to \mathcal{Q}$. For example, w can be a deep network that takes an input $x \in \mathcal{X}$ and produces a latent feature vector $q = w(x) \in \mathcal{Q}$. Next, we define a *classifier* $h \in \mathcal{H}$,

- representing all possible classification heads operating on \mathcal{Q} .
- 111 Formally, h is a score function defined as $h : \mathcal{Q} \times \mathcal{Y} \to \mathbb{R}$, 112 where the predicted class is $h(x) = \arg \max_{y \in \mathcal{Y}} h(q, y)$.
- 113 Likewise, we define a *regressor* $f \in \mathcal{F}$, representing all
- 114 regression heads, where $f : \mathcal{Q} \to \mathcal{T}$. These components are
- 115 combined into a single multi-head network $q \in \mathcal{G}$, where 116 $\mathcal{G} = \{ g : g(x) = (h \circ w(x), f \circ w(x)) \mid w \in \mathcal{W}, h \in \mathcal{W} \}$ 117 $\mathcal{H}, f \in \mathcal{F}$. Hence, g jointly produces classification and 118 regression outputs, h(q) and f(q), from the same latent
- 119 representation q = w(x).
- 120
- 121

Consistency in classification. In classification, the pri-122 mary objective is to find a classifier $h \in \mathcal{H}$ that min-123 imizes the true error $\mathcal{E}_{\ell_{01}}(h)$, defined as $\mathcal{E}_{\ell_{01}}(h) =$ 124 $\mathbb{E}_{(x,y)}[\ell_{01}(h(x),y)].$ The Bayes-optimal error is given by $\mathcal{E}^{B}_{\ell_{01}}(\mathcal{H}) = \inf_{h \in \mathcal{H}} \mathcal{E}_{\ell_{01}}(h).$ However, directly minimizing 125 126 $\mathcal{E}_{\ell_{01}}(h)$ is challenging due to the non-differentiability of 127 the true multiclass 0-1 loss (Zhang, 2002; Steinwart, 2007; 128 Awasthi et al., 2022). This motivates the introduction of 129 the cross-entropy multiclass surrogate family, denoted by 130 $\Phi_{01}^{\nu}: \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \to \mathbb{R}^+$, which provides a convex up-131 per bound to the *true multiclass loss* ℓ_{01} . This family is 132 parameterized by $\nu \geq 0$ and encompasses standard surro-133 gate functions widely adopted in the community such as the 134 MAE (Ghosh et al., 2017) or the log-softmax (Mohri et al., 135 2012). 136

$$\begin{array}{l} 137\\ 138\\ 139\\ 140 \end{array} \quad \Phi_{01}^{\nu} = \begin{cases} \frac{1}{1-\nu} \left(\left[\sum_{y' \in \mathcal{Y}} e^{h(x,y') - h(x,y)} \right]^{1-\nu} - 1 \right) & \nu \neq 1\\ \log \left(\sum_{y' \in \mathcal{Y}} e^{h(x,y') - h(x,y)} \right) & \nu = 1. \end{cases}$$

141 The corresponding surrogate error is defined as $\mathcal{E}_{\Phi_{01}^{\nu}}(h) =$ 142 $\mathbb{E}_{(x,y)}[\Phi_{01}^{\nu}(h(x),y)]$, with its optimal value given by $\mathcal{E}^*_{\Phi_{01}^{\nu}}(\mathcal{H}) = \inf_{h \in \mathcal{H}} \mathcal{E}_{\Phi_{01}^{\nu}}(h)$. A crucial property of a surro-143 144 gate loss is *Bayes-consistency*, which guarantees that min-145 imizing the surrogate generalization error also minimizes 146 the true generalization error (Zhang, 2002; Steinwart, 2007; 147 Bartlett et al., 2006; Tewari & Bartlett, 2007). Formally, Φ_{01}^{ν} 148 is Bayes-consistent with respect to ℓ_{01} if, for any sequence 149 ${h_k}_{k\in\mathbb{N}}\subset\mathcal{H}$, the following implication holds:

152

153

$$\begin{aligned}
\mathcal{E}_{\Phi_{01}^{\nu}}(h_k) - \mathcal{E}_{\Phi_{01}^{\nu}}^*(\mathcal{H}) \xrightarrow{k \to \infty} 0 \\
\implies \mathcal{E}_{\ell_{01}}(h_k) - \mathcal{E}_{\ell_{01}}^B(\mathcal{H}) \xrightarrow{k \to \infty} 0.
\end{aligned}$$
(2)

154 This property assumes that $\mathcal{H} = \mathcal{H}_{all}$, a condition that does 155 not necessarily hold for restricted hypothesis classes such 156 as \mathcal{H}_{lin} or \mathcal{H}_{ReLU} (Long & Servedio, 2013; Awasthi et al., 157 2022). To address this limitation, Awasthi et al. (2022) proposed \mathcal{H} -consistency bounds. These bounds depend on 159 a non-decreasing function $\Gamma : \mathbb{R}^+ \to \mathbb{R}^+$ and are expressed 160 as: 161

$$\mathcal{E}_{\Phi_{01}^{\nu}}(h) - \mathcal{E}_{\Phi_{01}^{\nu}}^{*}(\mathcal{H}) + \mathcal{U}_{\Phi_{01}^{\nu}}(\mathcal{H}) \geq$$

$$(3)$$

163
164
$$\Gamma\left(\mathcal{E}_{\ell_{01}}(h) - \mathcal{E}^B_{\ell_{01}}(\mathcal{H}) + \mathcal{U}_{\ell_{01}}(\mathcal{H})\right),$$

where the minimizability gap $\mathcal{U}_{\ell_{01}}(\mathcal{H})$ measures the disparity between the best-in-class generalization error and the expected pointwise minimum error: $\mathcal{U}_{\ell_{01}}(\mathcal{H}) = \mathcal{E}^B_{\ell_{01}}(\mathcal{H}) \mathbb{E}_x \left[\inf_{h \in \mathcal{H}} \mathbb{E}_{y|x}[\ell_{01}(h(x), y)] \right]$. Notably, the minimizability gap vanishes when $\mathcal{H} = \mathcal{H}_{all}$ (Steinwart, 2007; Awasthi et al., 2022). In the asymptotic limit, inequality (3) guarantees the recovery of Bayes-consistency, aligning with the condition in (2).

4. Two-stage Multi-Task L2D: Theoretical Analysis

4.1. Formulating the Deferral Loss

We extend the two-stage predictor-rejector framework, originally proposed by Mao et al. (2023a), to a multi-task context. We define an offline-trained model represented by the multi-task model $g \in \mathcal{G}$ defined in Section 3. We consider J offline-trained experts, denoted as M_i for $j \in \{1, \ldots, J\}$, where each expert produces predictions $m_j(x) = (m_j^h(x), m_j^f(x))$ with $m_j^h(x) \in \mathcal{Y}$ and $m_i^f(x) \in \mathcal{T}$. The predictions $m_j(x)$ belong to the corresponding prediction space \mathcal{M}_j , i.e., $m_j(x) \in \mathcal{M}_j$. The combined predictions of all J experts are represented as $m(x) = (m_1(x), \ldots, m_J(x))$, which lies in the joint prediction space \mathcal{M} . We use the notation $[J] = \{1, \ldots, J\}$ exclusively to denote the set of experts and define the agent space $\mathcal{A} = \{0\} \cup [J]$ with $|\mathcal{A}| = J + 1$, the number of agents (model and experts) that we have in our system.

To allocate the decision, we define a rejector function $r \in \mathcal{R}$, where $r : \mathcal{X} \times \mathcal{A} \to \mathbb{R}$. The rejector determines which agent should be assigned the decision by following the rule $r(x) = \arg \max_{j \in \mathcal{A}} r(x, j)$. This formulation naturally leads to the deferral loss $\ell_{def} : \mathcal{R} \times \mathcal{G} \times \mathcal{Z} \times \mathcal{M} \to \mathbb{R}^+$:

Definition 4.1 (True deferral loss). Let an input $x \in \mathcal{X}$, for any $r \in \mathcal{R}$, we have the *true deferral loss*:

$$\ell_{\rm def}(r,g,m,z) = \sum_{j=0}^{J} c_j(g(x),m_j(x),z) \mathbf{1}_{r(x)=j},$$

with a bounded cost c_i that quantifies the cost of allocating a decision to the agent $j \in \mathcal{A}$. When the rejector $r \in \mathcal{R}$ predicts r(x) = 0, the decision is allocated to the multitask model q. This allocation incurs a general cost c_0 , which is defined in this context as $c_0(g(x), z) = \rho(g(x), z)$, where $\rho: \mathcal{Y} \times \mathcal{T} \times \mathcal{Z} \to \mathbb{R}^+$ represents a general function capturing the discrepancy between the model's prediction q(x) and the ground truth z. Conversely, when the rejector predicts r(x) = j for some j > 0, the decision is deferred to expert j, resulting in a deferral cost c_j defined as $c_i(m_i(x), z) = \rho(m_i(x), z) + \beta_i$, where $m_i(x)$ denotes the prediction by expert j, and $\beta_j \ge 0$ accounts for the inherent cost associated with querying expert *j*. For example, β_j could reflect the resources, effort, or time required to obtain the expertise from expert j.

Optimal deferral rule: In Definition 4.1, we introduced the *true deferral loss*, ℓ_{def} , which we aim to minimize by identifying the Bayes rejector $r \in \mathcal{R}$ that minimizes the true error. To formalize this, we analyze the pointwise Bayes rejector $r^B(x)$, which minimizes the conditional risk $C_{\ell_{def}}$. The true risk is given by $\mathcal{E}_{\ell_{def}}(g, r) = \mathbb{E}_x[\mathcal{C}_{\ell_{def}}(g, r, x)]$. The following Lemma 4.2 can be established.

Lemma 4.2 (Pointwise Bayes Rejector). Given an input $x \in \mathcal{X}$ and a distribution \mathcal{D} , the optimal rejection rule that minimizes the conditional risk $C_{\ell_{def}}$ associated with the true deferral loss ℓ_{def} is defined as:

$$r^{B}(x) = \begin{cases} 0 & \text{if } \inf_{g \in \mathcal{G}} \mathbb{E}_{y,t|x}[c_{0}] \leq \min_{j \in [J]} \mathbb{E}_{y,t|x}[c_{j}] \\ j & \text{otherwise}, \end{cases}$$

The proof is provided in Appendix B. Lemma 4.2 establishes that the rejector $r \in \mathcal{R}$ optimally determines whether to utilize the multi-task model $g \in \mathcal{G}$ or to defer to the most cost-effective expert. Specifically, the rejector defers to the expert *j* that minimizes the expected deferral cost, $\mathbb{E}_{y,t|x}[c_j(g(x), m_j(x), z)]$, whenever the expected cost of the optimal multi-task model, $\inf_{g \in \mathcal{G}} \mathbb{E}_{y,t|x}[c_0(g(x), z)]$, exceeds the minimum expected cost of the most confident expert.

Although the *true deferral loss*, ℓ_{def} , and its corresponding Bayes rejector were introduced in Lemma 4.2, the practical computation of this rejector is hindered by the nondifferentiability of ℓ_{def} (Zhang, 2002).

4.2. Surrogate Loss for Two-Stage Multi-Task L2D

Introducing the Surrogate: To address challenges analogous to those posed by discontinuous loss functions (Berkson, 1944; Cortes & Vapnik, 1995), we formalize surrogate losses with desirable analytical properties. Specifically, we use the cross-entropy *multiclass surrogate loss* Φ_{01}^{ν} : $\mathcal{R} \times \mathcal{X} \times \mathcal{A} \to \mathbb{R}^+$ that is convex and upper-bounds the *true multiclass loss* ℓ_{01} . This surrogate family is defined in Equation 1.

Mao et al. (2024e) introduced a convex, upper-bound surrogate for the *true deferral loss*. Building on this, we incorporate new costs c_j to capture the interdependence between classification and regression tasks. The resulting surrogate family, Φ_{def}^{ν} : $\mathcal{R} \times \mathcal{G} \times \mathcal{M} \times \mathcal{Z} \rightarrow \mathbb{R}^+$, yields Lemma 4.3:

Lemma 4.3 (Surrogate deferral losses). Let $x \in \mathcal{X}$ be a given input and Φ_{01}^{ν} a multiclass surrogate loss family. The surrogate deferral losses Φ_{def}^{ν} for J + 1 agents are defined as:

$$\Phi_{def}^{\nu}(r,g,m,z) = \sum_{j=0}^{J} \tau_j \left(g(x), m(x), z\right) \Phi_{01}^{\nu}(r,x,j),$$
(4)
with the aggregated costs $\tau_j \left(g(x), m(x), z\right) =$

$$\sum_{i=0}^{J} c_i(g(x), m_i(x), z) \mathbf{1}_{i \neq j}.$$

The surrogate deferral loss family, Φ_{def}^{ν} , aggregates the weighted surrogate losses for each possible decision: deferring to one of the *J* experts or utilizing the model directly. The weights $\tau_{j\in A}$ represent the cumulative costs associated with each deferral path. Specifically, the weight τ_0 captures the total deferral cost across all experts, serving as a baseline for deferral decisions. Conversely, the weights $\tau_{j\in [J]}$ quantify the combined cost of utilizing the model while deferring to all experts except expert *j*.

The proposed surrogate family is broadly applicable and only requires Φ_{01}^{ν} to admit an \mathcal{R} -consistency bound. Furthermore, the formulated costs c_j accommodate any multitask setting.

Consistency of the Surrogate Losses: In Lemma 4.3, we established that the surrogate losses serve as a convex upper bound for the *true deferral loss*. However, it remains to be shown whether this surrogate family provides a faithful approximation of the target true loss, ℓ_{def} . Specifically, the relationship between $r^*(x)$, the pointwise optimal rejector minimizing a surrogate loss from the family, and $r^B(x)$, the pointwise Bayes-optimal rejector for the true loss, is not immediately evident. To address this, we analyze the discrepancy between the *surrogates' excess risk*, $\mathcal{E}_{\Phi_{def}^{\nu}}(g,r) - \mathcal{E}_{\Phi_{def}^{\nu}}^{B}(\mathcal{G},\mathcal{R})$, and the *excess risk under the true loss*, $\mathcal{E}_{\ell_{def}}(g,r) - \mathcal{E}_{\ell_{def}}^{B}(\mathcal{G},\mathcal{R})$. This analysis is critical for understanding the surrogates' consistency, as explored in prior works (Steinwart, 2007; Zhang, 2002; Bartlett et al., 2006; Awasthi et al., 2022).

Using consistency bounds from (Awasthi et al., 2022; Mao et al., 2024b), we present Theorem 4.4, which establishes the $(\mathcal{G}, \mathcal{R})$ -consistency of the *surrogate deferral loss* family.

Theorem 4.4 ((\mathcal{G}, \mathcal{R})-consistency bounds). Let $g \in \mathcal{G}$ be a multi-task model. Suppose there exists a non-decreasing function $\Gamma^{\nu} : \mathcal{R}^+ \to \mathbb{R}^+$ for $\nu \ge 0$, such that the \mathcal{R} consistency bounds hold for any distribution \mathcal{D} :

$$\begin{aligned} \mathcal{E}_{\Phi_{01}^{\nu}}(r) &- \mathcal{E}_{\Phi_{01}^{\nu}}^{*}(\mathcal{R}) + \mathcal{U}_{\Phi_{01}^{\nu}}(\mathcal{R}) \geq \\ \Gamma^{\nu}(\mathcal{E}_{\ell_{01}}(r) - \mathcal{E}_{\ell_{01}}^{B}(\mathcal{R}) + \mathcal{U}_{\ell_{01}}(\mathcal{R})), \end{aligned}$$

then for any $(g,r) \in \mathcal{G} \times \mathcal{R}$, any distribution \mathcal{D} and any

 $\begin{array}{ll} 220 & x \in \mathcal{X}, \\ 221 & \end{array}$

222 223

224

225

226

227

228

229

230

263

264

265

266

267

268

269

270

271

272

273

274

$$egin{aligned} \mathcal{E}_{\ell_{def}}(g,r) &- \mathcal{E}^B_{\ell_{def}}(\mathcal{G},\mathcal{R}) + \mathcal{U}_{\ell_{def}}(\mathcal{G},\mathcal{R}) \leq \ \overline{\Gamma}^
u \left(\mathcal{E}_{\Phi^
u_{def}}(r) - \mathcal{E}^*_{\Phi^
u_{def}}(\mathcal{R}) + \mathcal{U}_{\Phi^
u_{def}}(\mathcal{R})
ight) \ &+ \mathcal{E}_{c_0}(g) - \mathcal{E}^B_{c_0}(\mathcal{G}) + \mathcal{U}_{c_0}(\mathcal{G}), \end{aligned}$$

where $\overline{\Gamma}^{\nu}(u) = \|\overline{\tau}\|_1 \Gamma^{\nu}\left(\frac{u}{\|\overline{\tau}\|_1}\right)$, with $\Gamma^{\nu}(u) = \mathcal{T}^{-1,\nu}(u)$, and for the log-softmax surrogate, $\mathcal{T}^{\nu=1}(u) = \frac{1+u}{2}\log(1+u) + \frac{1-u}{2}\log(1-u)$.

The proof of Theorem 4.4, along with additional bounds 231 for $\nu \ge 0$, is provided in Appendix C. Theorem 4.4 es-232 tablishes sharper bounds than those in Mao et al. (2024e). 233 Our bounds are specifically tailored for the cross-entropy 234 235 surrogate family and explicitly depend on the parameter ν . Furthermore, the theorem shows that the tightness of these 236 bounds is governed by the L_1 -norm of the aggregated costs, 237 238 $\|\overline{\boldsymbol{\tau}}\|_1$

239 Moreover, we show that our surrogate losses are $(\mathcal{G}, \mathcal{R})$ -240 consistent for a *multiclass surrogate* family Φ_{01}^{ν} that is 241 \mathcal{R} -consistent. Assuming $\mathcal{R} = \mathcal{R}_{all}$ and $\mathcal{G} = \mathcal{G}_{all}$, the 242 minimizability gaps vanish, as demonstrated in Steinwart 243 (2007). Therefore, by minimizing the surrogates' deferral 244 excess risk while accounting for the minimizability gap, we 245 establish that $\mathcal{E}_{\Phi_{def}^{\nu}}(r_k) - \mathcal{E}_{\Phi_{def}^{\nu}}^*(\mathcal{R}_{all}) + \mathcal{U}_{\Phi_{def}^{\nu}}(\mathcal{R}_{all}) \xrightarrow{k \to \infty} 0$. Since the multi-task model g is trained offline, it 246 247 is reasonable to assume that the c_0 -excess risk satisfies 248 $\mathcal{E}_{c_0}(g_k) - \mathcal{E}^B_{c_0}(\mathcal{G}_{all}) + \mathcal{U}_{c_0}(\mathcal{G}_{all}) \xrightarrow{k \to \infty} 0$. This result implies that the left-hand side is bounded above by zero, leading to 249 250 $\mathcal{E}_{\ell_{\mathsf{def}}}(g, r_k) - \mathcal{E}^B_{\ell_{\mathsf{def}}}(\mathcal{G}_{\mathsf{all}}, \mathcal{R}_{\mathsf{all}}) + \mathcal{U}_{\ell_{\mathsf{def}}}(\mathcal{G}_{\mathsf{all}}, \mathcal{R}_{\mathsf{all}}) \xrightarrow{k \to \infty} 0,$ by 251 leveraging the properties of $\overline{\Gamma}^{\nu}$. Consequently, the following 252 253 corollary holds: 254

Corollary 4.5 (Bayes-consistency of the deferral surrogate losses). Under the conditions of Theorem 4.4, assuming $(\mathcal{G}, \mathcal{R}) = (\mathcal{G}_{all}, \mathcal{R}_{all})$ and $\mathcal{E}_{c_0}(g_k) - \mathcal{E}_{c_0}^B(\mathcal{G}_{all}) \xrightarrow{k \to \infty} 0$, the surrogate deferral loss family Φ_{def}^{ν} is Bayes-consistent with respect to the true deferral loss ℓ_{def} . Specifically, minimizing the surrogates' deferral excess risk ensures the convergence of the true deferral excess risk. Formally, for $\{r_k\}_{k\in\mathbb{N}} \subset \mathcal{R}$ and $\{g_k\}_{k\in\mathbb{N}} \subset \mathcal{G}$:

$$\begin{aligned}
\mathcal{E}_{\Phi_{def}^{\nu}}(r_k) - \mathcal{E}_{\Phi_{def}^{\nu}}^{*}(\mathcal{R}_{all}) \xrightarrow{k \to \infty} 0 \\
\implies \mathcal{E}_{\ell_{def}}(g_k, r_k) - \mathcal{E}_{\ell_{def}}^{B}(\mathcal{G}_{all}, \mathcal{R}_{all}) \xrightarrow{k \to \infty} 0.
\end{aligned}$$
(5)

This result demonstrates that as $k \to \infty$, the surrogates Φ_{def}^{ν} achieve asymptotic Bayes optimality for both the multitask model g and the rejector r, effectively bridging the theoretical gap between the surrogate losses and the true deferral loss. Moreover, the pointwise optimal rejector $r^*(x)$ converges to a close approximation of the pointwise Bayesoptimal rejector $r^B(x)$, yielding a deferral rule consistent with the structure described in Lemma 4.2 (Bartlett et al., 2006).

Analysis of the minimizability gap: As shown by Awasthi et al. (2022), the minimizability gap does not vanish in general. Understanding its conditions, quantifying its magnitude, and identifying mitigation strategies are essential to ensuring that surrogate-based optimization aligns with task-specific objectives.

We provide a strong and novel characterization of the minimizability gap in the two-stage setting with multiple experts, extending the results of Mao et al. (2024f), who analyzed the gap in the context of learning with abstention (constant cost) for a single expert and a specific distribution.

Theorem 4.6 (Characterization Minimizability Gaps). Assume \mathcal{R} symmetric and complete. Then, for the crossentropy multiclass surrogates Φ_{01}^{ν} and any distribution \mathcal{D} , it follows for $\nu \geq 0$:

$$\mathcal{C}_{\Phi_{def}^{\nu,*}}^{\nu,*} = \begin{cases} \|\overline{\tau}\|_1 H\left(\frac{\overline{\tau}}{\|\overline{\tau}\|_1}\right) & \nu = 1\\ \|\overline{\tau}\|_1 - \|\overline{\tau}\|_{\infty} & \nu = 2\\ \frac{1}{\nu - 1} \left[\|\overline{\tau}\|_1 - \|\overline{\tau}\|_{\frac{1}{2-\nu}}\right] & \nu \in (1,2)\\ \frac{1}{1-\nu} \left[\left(\sum_{k=0}^J \overline{\tau}_k^{\frac{1}{2-\nu}}\right)^{2-\nu} - \|\overline{\tau}\|_1\right] & otherwise, \end{cases}$$

then the minimizability gap is,

$$\mathcal{U}_{\Phi_{def}^{\nu}}(\mathcal{R}) = \mathcal{E}_{\Phi_{def}^{\nu}}^{*}(\mathcal{R}) - \mathbb{E}_{x}[\inf_{r \in \mathcal{R}} \mathcal{C}_{\Phi_{def}^{\nu}}^{\nu}(r, x)]$$

with $\overline{\tau} = \{\mathbb{E}_{y,t|x}[\overline{\tau}_0], \dots, \mathbb{E}_{y,t|x}[\overline{\tau}_J]\}$, the aggregated costs $\tau_j = \sum_{k=0}^J c_k \mathbf{1}_{k \neq j}$, and the Shannon Entropy H.

We provide the proof in Appendix D. Theorem 4.6 characterizes the minimizability gap $\mathcal{U}_{\Phi_{def}^{\nu}}(\mathcal{R})$ for cross-entropy multiclass surrogates over symmetric and complete hypothesis sets \mathcal{R} . The gap varies with $\nu > 0$, exhibiting distinct behaviors across different surrogate. For $\nu = 1$, the gap is proportional to the Shannon entropy of the normalized expected cost vector $\frac{\overline{\tau}}{\|\overline{\tau}\|_1}$, increasing with entropy and reflecting higher uncertainty in misclassification distribution. At $\nu = 2$, it simplifies to the difference between the L_1 norm and L_{∞} -norm of $\overline{\tau}$, where a smaller gap indicates concentrated misclassifications, reducing uncertainty. For $\nu \in (1,2)$, the gap balances the entropy-based sensitivity at $\nu = 1$ and the margin-based sensitivity at $\nu = 2$. As $\nu \to 1^+$, it emphasizes agents with higher misclassification counts; as $\nu \to 2^-$, it shifts toward aggregate misclassification counts. For $\nu < 1$, where $p = \frac{1}{2-\nu} \in (0,1)$, the gap is more sensitive to misclassification distribution, increasing when errors are dispersed. For $\nu > 2$, where p < 0, reciprocal weighting reduces sensitivity to dominant errors, potentially decreasing the gap but at the risk of underemphasizing critical misclassifications.

275 In the setting of learning with abstention and a single expert 276 (J = 1), assigning costs $\tau_0 = 1$ and $\tau_J = 1 - c$ recovers the 277 minimizability gap introduced in (Mao et al., 2024f). Thus, 278 our minimizability gap can be seen as a generalization to 279 multiple experts, non-constant costs, and to any distribution 280 \mathcal{D} .

282 **4.3. Generalization Bound**

281

283

284

285

286

287

288

289

290

291

292

306

307

308

327

329

We aim to quantify the generalization capability of our system, considering both the complexity of the hypothesis space and the quality of the participating agents. To this end, we define the empirical optimal rejector \hat{r}^B as the minimizer of the empirical generalization error:

$$\widehat{r}^B = \arg\min_{r\in\mathcal{R}} \frac{1}{K} \sum_{k=1}^{K} \ell_{\mathsf{def}}(g, m, r, z_k), \tag{6}$$

where ℓ_{def} denotes the *true deferral loss* function. To characterize the system's generalization ability, we utilize the Rademacher complexity, which measures the expressive richness of a hypothesis class by evaluating its capacity to fit random noise (Bartlett & Mendelson, 2003; Mohri et al., 2012). The proof of Lemma 4.7 is provided in Appendix E.

300 **Lemma 4.7.** Let \mathcal{L}_1 be a family of functions mapping \mathcal{X} 301 to [0, 1], and let \mathcal{L}_2 be a family of functions mapping \mathcal{X} to 302 $\{0, 1\}$. Define $\mathcal{L} = \{l_1 l_2 : l_1 \in \mathcal{L}_1, l_2 \in \mathcal{L}_2\}$. Then, the 303 empirical Rademacher complexity of \mathcal{L} for any sample S of 304 size K is bounded by:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{L}) \leq \widehat{\mathfrak{R}}_{S}(\mathcal{L}_{1}) + \widehat{\mathfrak{R}}_{S}(\mathcal{L}_{2}).$$
(7)

309 For simplicity, we assume costs $c_0(g(x), z)$ = 310 $\ell_{01}(h(x), y) + \ell_{\text{reg}}(f(x), t)$ and $c_{j>0}(m_j(x), z) = c_0(m_j(x), z)$. We assume the regression loss ℓ_{reg} to be 311 312 non-negative, bounded by L, and Lipschitz. Furthermore, we assume that $m_{k,j}^h$ is drawn from the conditional 313 314 distribution of the random variable M_j^h given parameters 315 $\{X = x_k, Y = y_k\}$, and that $m_{k,j}^f$ is drawn from the 316 conditional distribution of M_i^f given $\{X = x_k, T = t_k\}$. 317 We define the family of deferral loss functions as 318 $\mathcal{L}_{def} = \{\ell_{def} : \mathcal{G} \times \mathcal{R} \times \mathcal{M} \times \mathcal{Z} \to [0,1]\}.$ Under these 319 assumptions, we derive the generalization bounds for the 320 binary setting as follows: 321

Theorem 4.8 (Learning bounds of the deferral loss). For any expert M_j , any distribution \mathcal{D} over \mathcal{Z} , we have with probability $1 - \delta$ for $\delta \in [0, 1/2]$, that the following bound holds at the optimum:

$$\mathcal{E}_{\ell_{def}}(h, f, r) \leq \widehat{\mathcal{E}}_{\ell_{def}}(h, f, r) + 2\mathfrak{R}_{K}(\mathcal{L}_{def}) + \sqrt{rac{\log 1/\delta}{2K}},$$

with

$$\begin{aligned} \mathfrak{R}_{K}(\mathcal{L}_{def}) &\leq \frac{1}{2} \mathfrak{R}_{K}(\mathcal{H}) + \mathfrak{R}_{K}(\mathcal{F}) + \sum_{j=1}^{J} \Omega(m_{j}^{h}, y) \\ &+ \Big(\sum_{j=1}^{J} \max \ell_{reg}(m_{j}^{f}, t) + 2\Big) \mathfrak{R}_{K}(\mathcal{R}), \end{aligned}$$

with $\Omega(m_j^h, y) = \frac{1}{2} \mathcal{D}(m_j^h \neq y) \exp\left(-\frac{K}{8} \mathcal{D}(m_j^h \neq y)\right) + \mathcal{R}_{K\mathcal{D}(m_j^h \neq y)/2}(\mathcal{R}).$

We prove Theorem 4.8 in Appendix F. The terms $\Re_K(\mathcal{H})$ and $\Re_K(\mathcal{F})$ denote the Rademacher complexities of the hypothesis class \mathcal{H} and function class \mathcal{F} , respectively, indicating that the generalization bounds depend on the complexity of the pre-trained model. The term $\Omega(m_j^h, y)$ captures the impact of each expert's classification error on the learning bound. It includes an exponentially decaying factor, $\frac{\mathcal{D}(m_j^h \neq y)}{2} \exp\left(-\frac{K\mathcal{D}(m_j^h \neq y)}{8}\right)$, which decreases rapidly as the sample size K grows or as the expert's error rate $\mathcal{D}(m_j^h \neq y)$ declines (Mozannar & Sontag, 2020). This reflects the intuition that more accurate experts contribute less to the bound, improving overall generalization. Finally, the last term suggests that the generalization properties of our *true deferral loss* depend on the expert's regression performance.

5. Experiments

In this section, we present the performance improvements achieved by the proposed Learning-to-Defer surrogate in a multi-task context. Specifically, we demonstrate that our approach excels in object detection, a task where classification and regression components are inherently intertwined, and where existing L2D methods encounter significant limitations. Furthermore, we evaluate our approach on an Electronic Health Record task, jointly predicting mortality (classification) and length of stay (regression), comparing our results with Mao et al. (2023a; 2024e).

For each experiment, we report the mean and standard deviation across four independent trials to account for variability in the results. All training and evaluation were conducted on an NVIDIA H100 GPU. We give our training algorithm in Appendix A. Additional figures and details are provided in Appendix G. To ensure reproducibility, we have made our implementation publicly available.

5.1. Object Detection Task

We evaluate our approach using the Pascal VOC dataset (Everingham et al., 2010), a multi-object detection benchmark. This is the first time such a multi-task problem has been explored within the L2D framework as previous L2D

approaches require the classification and regression component to be independent (Mao et al., 2023a; 2024e).

333 Dataset and Metrics: The PASCAL Visual Object 334 Classes (VOC) dataset (Everingham et al., 2010) serves 335 as a widely recognized benchmark in computer vision for 336 evaluating object detection models. It consists of annotated images spanning 20 object categories, showcasing diverse 338 scenes with varying scales, occlusions, and lighting condi-339 tions. To assess object detection performance, we report the 340 mean Average Precision (mAP), a standard metric in the 341 field. The mAP quantifies the average precision across all 342 object classes by calculating the area under the precision-343 recall curve for each class. Additionally, in the context of 344 L2D, we report the allocation metric (All.), which represents 345 the ratio of allocated queries per agent.

346

371

372

373

374

375

376

377

378

379

380

381

382

383

384

347 Agents setting: We trained three distinct Faster R-CNN 348 models (Ren et al., 2016) to serve as our agents, differen-349 tiated by their computational complexities. The smallest, 350 characterized by GFLOPS = 12.2, represents our model 351 $g \in \mathcal{G}$ with $\mathcal{G} = \{g : g(x) = (h \circ w(x), f \circ w(x)) \mid w \in$ 352 $\mathcal{W}, h \in \mathcal{H}, f \in \mathcal{F}$. The medium-sized, denoted as Expert 353 1, has a computational cost of GFLOPS = 134.4, while the 354 largest, Expert 2, operates at GFLOPS = 280.3. To account 355 for the difference in complexity between Experts 1 and 2, 356 we define the ratio $R_G = 280.3/134.4$ and set the query 357 cost for Expert 1 as $\beta_1 = \beta_2/R_G$. This parameterization 358 reflects the relative computational costs of querying experts. 359 We define the agent costs as $c_0(g(x), z) = mAP(g(x), z)$ 360 and $c_{i \in [J]}(m_i(x), z) = mAP(m_i(x), z)$. We report the per-361 formance metrics of the agents alongside additional training 362 details in Appendix G.1. 363

Rejector: The rejector is trained using a smaller version of
the Faster R-CNN model (Ren et al., 2016). Training is performed for 200 epochs using the Adam optimizer (Kingma & Ba, 2017) with a learning rate of 0.001 and a batch size
of 64. The checkpoint achieving the lowest empirical risk
on the validation set is selected for evaluation.

Results: In Figure 5.1, we observe that for lower cost values, specifically when $\beta_1 < 0.15$, the system consistently avoids selecting Expert 1. This outcome arises because the cost difference between β_1 and β_2 is negligible, making it more advantageous to defer to Expert 2 (the most accurate expert), where the modest cost increase is offset by superior outcomes. When $\beta_2 = 0.15$, however, it becomes optimal to defer to both experts and model at the same time. In particular, there exist instances $x \in \mathcal{X}$ where both Expert 1 and Expert 2 correctly predict the target (while the model does not). In such cases, Expert 1 is preferred due to its lower cost $\beta_1 < \beta_2$. Conversely, for instances $x \in \mathcal{X}$ where Expert 2 is accurate and Expert 1 (along with the model)



Figure 1. Performance comparison across different cost values β_2 on Pascal VOC (Everingham et al., 2010). The table reports the mean Average Precision (mAP) and the allocation ratio for the model and two experts with mean and variance. We report these results in Appendix Table 3.

is incorrect, the system continues to select Expert 2, as β_2 remains relatively low. For $\beta_2 \ge 0.2$, the increasing cost differential between the experts shifts the balance in favor of Expert 1, enabling the system to achieve strong performance while minimizing overall costs.

This demonstrates that our approach effectively allocates queries among agents, thereby enhancing the overall performance of the system, even when the classification and regression tasks are interdependent.

5.2. EHR Task

We compare our novel approach against existing two-stage L2D methods (Mao et al., 2023a; 2024e). Unlike the first experiment on object detection (Subsection 5.1), where classification and regression tasks are interdependent, this evaluation focuses on a second scenario where the two tasks can be treated independently.

Dataset and Metrics: The Medical Information Mart for Intensive Care IV (MIMIC-IV) dataset (Johnson et al., 2023) is a comprehensive collection of de-identified health-related data patients admitted to critical care units. For our analysis, we focus on two tasks: mortality prediction and length-ofstay prediction, corresponding to classification and regression tasks, respectively. To evaluate performance, we report accuracy (Acc) for the mortality prediction task, which quantifies classification performance, and Smooth L1 loss (sL1) for the length-of-stay prediction task, which measures the deviation between the predicted and actual values. Additionally, we report the allocation metric (All.) for L2D to capture query allocation behavior.

385 Agents setting: We consider two experts, M_1 and M_2 , 386 acting as specialized agents, aligning with the category al-387 location described in (Mozannar & Sontag, 2020; Verma 388 et al., 2022; Verma & Nalisnick, 2022; Cao et al., 2024). 389 The dataset is partitioned into Z = 6 clusters using the 390 K-means algorithm (Llovd, 1982), where Z is selected via the Elbow method (Thorndike, 1953). The clusters are de-392 noted as $\{C_1, C_2, \ldots, C_Z\}$. Each cluster represents a subset of data instances grouped by feature similarity, enabling features-specific specialization by the experts. The experts 395 are assumed to specialize in distinct subsets of clusters based 396 on the task. For classification, M1 correctly predicts the outcomes for clusters $C_{cla}^{M_1} = \{C_1, C_2, C_4\}$, while M_2 handles clusters $C_{cla}^{M_2} = \{C_1, C_5, C_5\}$. Notably, cluster C_1 is shared between the two experts, reflecting practical scenarios where 397 398 399 400 domain knowledge overlaps. For regression tasks, M1 is accurate on clusters $C_{\text{reg}}^{M_1} = \{C_1, C_3, C_5\}$, while M₂ spe-401 402 cializes in clusters $C_{\text{reg}}^{M_2} = \{C_1, C_4, C_6\}$. Here too, overlap 403 is modeled, with cluster C_1 being common to both experts 404 and classification-regression task. Note that the category 405 assignments do not follow any specific rule.

We assume that each expert produces correct predictions for 407 the clusters they are assigned (Verma et al., 2022; Mozan-408 nar & Sontag, 2020). Conversely, for clusters outside their 409 expertise, predictions are assumed to be incorrect. In such 410 cases, for length-of-stay predictions, the outcomes are mod-411 eled using a uniform probability distribution to reflect uncer-412 tainty. The detailed performance evaluation of these agents 413 is provided in Appendix G.2. 414

406

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

415 The model utilizes two compact transformer architectures 416 (Vaswani et al., 2017) for addressing both classification and 417 regression tasks, formally defined as $\mathcal{G} = \{q : q(x) =$ 418 $(h(x), f(x)) \mid h \in \mathcal{H}, f \in \mathcal{F}$ }. The agent's costs are spec-419 ified as $c_0(g(x), z) = \lambda^{\operatorname{cla}} \ell_{01}(h(x), y) + \lambda^{\operatorname{reg}} \ell_{\operatorname{reg}}(f(x), t)$ 420 and $c_{j \in [J]}(m_j(x), z) = c_0(m_j(x), z) + \beta_j$. Consistent 421 with prior works (Mozannar & Sontag, 2020; Verma et al., 422 2022; Mao et al., 2023a; 2024e), we set $\beta_i = 0$. 423

Rejectors: The two-stage L2D rejectors are trained using a small transformer model (Vaswani et al., 2017) as the encoder, following the approach outlined by Yang et al. (2023), with a classification head for query allocation. Training is performed over 100 epochs with a learning rate of 0.003, a warm-up period of 0.1, a cosine learning rate scheduler, the Adam optimizer (Kingma & Ba, 2017), and a batch size of 1024 for all baselines. The checkpoint with the lowest empirical risk on the validation set is selected for evaluation.

Results: Table 5.2 compares the performance of our proposed Learning-to-Defer (Ours) approach with two existing methods: a classification-focused rejector (Mao et al., 2023a) and a regression-focused rejector (Mao et al., 2024e). The results highlight the limitations of task-specific rejectors

and the advantages of our balanced approach.

Rejector	Acc (%)	sL1	All. Model	All. Expert 1	All. Expert 2
Mao et al. (2023a)	$71.3\pm.1$	$1.45\pm.03$	$.60 \pm .02$	$.01 \pm .01$	$.39 \pm .02$
Mao et al. (2024e)	$50.7 \pm .8$	$1.18\pm.05$	$.38 \pm .01$	$.37 \pm .02$	$.25 \pm .01$
Ours	$70.0\pm.5$	$1.28\pm.02$	$.66\pm.01$	$.12 \pm .02$	$.22 \pm .01$

Table 1. Performance comparison of different two-stage L2D. The table reports accuracy (Acc), smooth L1 loss (sL1), and allocation rates (All.) to the model and experts with mean and variance.

The classification-focused rejector achieves the highest classification accuracy at 71.3% but struggles with regression, as reflected by its high smooth L1 loss of 1.45. On the other hand, the regression-focused rejector achieves the best regression performance with an sL1 loss of 1.18 but performs poorly in classification with an accuracy of 50.7%. In contrast, our method balances performance across tasks, achieving a classification accuracy of 70.0% and an sL1 loss of 1.28. Moreover, it significantly reduces reliance on experts, allocating 66% of queries to the model compared to 60% for Mao et al. (2023a) and 38% for Mao et al. (2024e). Expert involvement is minimized, with only 12% and 22% of queries allocated to Experts 1 and 2, respectively.

Since the experts possess distinct knowledge for the two tasks ($C_{cla}^{M_1}$ and $C_{reg}^{M_1}$ for M₁), independently deferring classification and regression may lead to suboptimal performance. In contrast, our approach models deferral decisions dependently, considering the interplay between the two components to achieve better overall results.

6. Conclusion

We introduced a Two-Stage Learning-to-Defer framework for multi-task problems, extending existing approaches to jointly handle classification and regression. We proposed a two-stage surrogate loss family that is both (\mathcal{G}, \mathcal{R})consistent and Bayes-consistent for any cross-entropy-based surrogate. Additionally, we derived tight consistency bounds linked to cross-entropy losses and the L_1 -norm of aggregated costs. We further established novel minimizability gap for the two-stage setting, generalizing prior results to Learning-to-Defer with multiple experts. Finally, we showed that our learning bounds improve with a richer hypothesis space and more confident experts.

We validated our framework on two challenging tasks: (i) object detection, where classification and regression are inherently interdependent—beyond the scope of existing L2D methods; and (ii) electronic health record analysis, where we demonstrated that current L2D approaches can be suboptimal even when classification and regression tasks are independent.

Impact Statement 440

441

442

443

444

445

446

447

448

457

458

472

473

474

475

This paper advances the theoretical and practical understanding of machine learning, contributing to the development of more effective models and methods. While our research does not present any immediate or significant ethical concerns, we recognize the potential for indirect societal impacts.

449 References

- 450 Awasthi, P., Mao, A., Mohri, M., and Zhong, Y. Multi-class 451 h-consistency bounds. In Proceedings of the 36th Inter-452 national Conference on Neural Information Processing 453 Systems, NIPS '22, Red Hook, NY, USA, 2022. Curran 454 Associates Inc. ISBN 9781713871088. 455
- 456 Balagurunathan, Y., Mitchell, R., and El Naga, I. Requirements and reliability of AI in the medical context. Physica Medica, 83:72-78, 2021.
- 459 Bartlett, P., Jordan, M., and McAuliffe, J. Convexity, 460 classification, and risk bounds. Journal of the Ameri-461 can Statistical Association, 101:138-156, 02 2006. doi: 462 10.1198/01621450500000907. 463
- 464 Bartlett, P. L. and Mendelson, S. Rademacher and Gaussian 465 complexities: Risk bounds and structural results. J. Mach. 466 Learn. Res., 3(null):463-482, March 2003. ISSN 1532-467 4435. 468
- 469 Bartlett, P. L. and Wegkamp, M. H. Classification with a 470 reject option using a hinge loss. The Journal of Machine 471 Learning Research, 9:1823–1840, June 2008.
 - Benz, N. L. C. and Rodriguez, M. G. Counterfactual inference of second opinions. In Uncertainty in Artificial Intelligence, pp. 453-463. PMLR, 2022.
- 476 Berkson, J. Application of the logistic function to 477 Journal of the American Statistical bio-assay. 478 Association, 39:357-365, 1944. URL https: 479 //api.semanticscholar.org/CorpusID: 480 122893121. 481
- Buch, S., Escorcia, V., Shen, C., Ghanem, B., and Car-482 los Niebles, J. SST: Single-stream temporal action pro-483 posals. In Proceedings of the IEEE Conference on Com-484 puter Vision and Pattern Recognition, pp. 2911-2920, 485 2017. 486
- 487 Cao, Y., Cai, T., Feng, L., Gu, L., Gu, J., An, B., Niu, 488 G., and Sugiyama, M. Generalizing consistent multi-489 class classification with rejection to be compatible with 490 arbitrary losses. In Proceedings of the 36th International 491 Conference on Neural Information Processing Systems, 492 NIPS '22, Red Hook, NY, USA, 2022. Curran Associates 493 Inc. ISBN 9781713871088. 494

- Cao, Y., Mozannar, H., Feng, L., Wei, H., and An, B. In defense of softmax parametrization for calibrated and consistent learning to defer. In Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. In Computer Vision - ECCV 2020: 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part I, pp. 213-229, Berlin, Heidelberg, 2020. Springer-Verlag. ISBN 978-3-030-58451-1. doi: 10.1007/978-3-030-58452-8_13. URL https://doi. org/10.1007/978-3-030-58452-8_13.
- Charusaie, M.-A., Mozannar, H., Sontag, D., and Samadi, S. Sample efficient learning of predictors that complement humans. In International Conference on Machine Learning, pp. 2972–3005. PMLR, 2022.
- Chow, C. On optimum recognition error and reject tradeoff. IEEE Transactions on Information Theory, 16(1):41-46, January 1970. doi: 10.1109/TIT.1970.1054406.
- Cortes, C. and Vapnik, V. N. Support-vector networks. Machine Learning, 20:273–297, 1995. URL https: //api.semanticscholar.org/CorpusID: 52874011.
- Cortes, C., DeSalvo, G., and Mohri, M. Learning with rejection. In Algorithmic Learning Theory: 27th International Conference, ALT 2016, Bari, Italy, October 19-21, 2016, Proceedings 27, pp. 67-82. Springer, 2016.
- Everingham, M., Gool, L., Williams, C. K., Winn, J., and Zisserman, A. The pascal visual object classes (voc) challenge. Int. J. Comput. Vision, 88(2): 303-338, June 2010. ISSN 0920-5691. doi: 10.1007/ s11263-009-0275-4. URL https://doi.org/10. 1007/s11263-009-0275-4.
- Geifman, Y. and El-Yaniv, R. Selective classification for deep neural networks. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips. cc/paper_files/paper/2017/file/ 4a8423d5e91fda00bb7e46540e2b0cf1-Paper. pdf.
- Ghosh, A., Kumar, H., and Sastry, P. S. Robust loss functions under label noise for deep neural networks. ArXiv, abs/1712.09482, 2017. URL https://api. semanticscholar.org/CorpusID:6546734.

- 495 Girshick, R. Fast R-CNN. *arXiv preprint arXiv:1504.08083*,
 496 2015.
- Hemmer, P., Schemmer, M., Vössing, M., and Kühl, N.
 Human-AI complementarity in hybrid intelligence systems: A structured literature review. *PACIS*, pp. 78, 2021.
- 502 Hemmer, P., Schellhammer, S., Vössing, M., Jakubik, J., and 503 Satzger, G. Forming effective human-AI teams: Building 504 machine learning models that complement the capabilities 505 of multiple experts. In Raedt, L. D. (ed.), Proceedings of 506 the Thirty-First International Joint Conference on Artifi-507 cial Intelligence, IJCAI-22, pp. 2478-2484. International 508 Joint Conferences on Artificial Intelligence Organization, 509 7 2022. doi: 10.24963/ijcai.2022/344. URL https: 510 //doi.org/10.24963/ijcai.2022/344. Main 511 Track.
- Johnson, A. E., Pollard, T. J., Shen, L., Lehman, L.-w. H.,
 Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Anthony Celi, L., and Mark, R. G. MIMIC-III, a freely
 accessible critical care database. *Scientific data*, 3(1):1–9,
 2016.
- Johnson, A. E., Bulgarelli, L., Shen, L., Gayles, A., Shammout, A., Horng, S., Pollard, T. J., Hao, S., Moody, B.,
 Gow, B., et al. MIMIC-IV, a freely accessible electronic health record dataset. *Scientific data*, 10(1):1, 2023.
- Kerrigan, G., Smyth, P., and Steyvers, M. Combining 523 human predictions with model probabilities via confusion 524 matrices and calibration. In Ranzato, M., Beygelzimer, 525 A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), 526 Advances in Neural Information Processing Systems, 527 volume 34, pp. 4421-4434. Curran Associates, Inc., 528 2021. URL https://proceedings.neurips. 529 cc/paper_files/paper/2021/file/ 530 234b941e88b755b7a72a1c1dd5022f30-Paper. 531 pdf. 532
- 533 Keswani, V., Lease, M., and Kenthapadi, K. Towards 534 unbiased and accurate deferral to multiple experts. In 535 Proceedings of the 2021 AAAI/ACM Conference on 536 AI, Ethics, and Society, AIES '21, pp. 154-165, New 537 York, NY, USA, 2021. Association for Computing 538 Machinery. ISBN 9781450384735. doi: 10.1145/ 539 3461702.3462516. URL https://doi.org/10. 540 1145/3461702.3462516. 541
- 542 Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization, 2017. URL https://arxiv.org/abs/ 1412.6980.
- Liu, S., Cao, Y., Zhang, Q., Feng, L., and An, B. Mitigating underfitting in learning to defer with consistent losses. In *International Conference on Artificial Intelligence and Statistics*, pp. 4816–4824. PMLR, 2024.

- Lloyd, S. Least squares quantization in PCM. *IEEE Trans*actions on Information Theory, 28(2):129–137, 1982. doi: 10.1109/TIT.1982.1056489.
- Long, P. and Servedio, R. Consistency versus realizable h-consistency for multiclass classification. In Dasgupta, S. and McAllester, D. (eds.), *Proceedings of the 30th International Conference on Machine Learning*, number 3 in Proceedings of Machine Learning Research, pp. 801– 809, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL https://proceedings.mlr.press/v28/ long13.html.
- Madras, D., Pitassi, T., and Zemel, R. Predict responsibly: Improving fairness and accuracy by learning to defer. In Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18, pp. 6150–6160, 2018.
- Mao, A., Mohri, C., Mohri, M., and Zhong, Y. Two-stage learning to defer with multiple experts. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. URL https://openreview.net/forum? id=GIlsH0T4b2.
- Mao, A., Mohri, M., and Zhong, Y. Cross-entropy loss functions: Theoretical analysis and applications. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23, 2023b.
- Mao, A., Mohri, M., and Zhong, Y. Predictor-rejector multiclass abstention: Theoretical analysis and algorithms. In *International Conference on Algorithmic Learning Theory*, pp. 822–867. PMLR, 2024a.
- Mao, A., Mohri, M., and Zhong, Y. Enhanced *h*-consistency bounds. arXiv preprint arXiv:2407.13722, 2024b.
- Mao, A., Mohri, M., and Zhong, Y. Principled approaches for learning to defer with multiple experts. In *International Workshop on Combinatorial Image Analysis*, pp. 107–135. Springer, 2024c.
- Mao, A., Mohri, M., and Zhong, Y. Realizable *h*-consistent and bayes-consistent loss functions for learning to defer, 2024d. URL https://arxiv.org/abs/2407. 13732.
- Mao, A., Mohri, M., and Zhong, Y. Regression with multi-expert deferral. arXiv 2403.19494, 2024e. https://arxiv.org/abs/2403.19494.
- Mao, A., Mohri, M., and Zhong, Y. Theoretically grounded loss functions and algorithms for score-based multi-class abstention. In *International Conference on Artificial Intelligence and Statistics*, pp. 4753–4761. PMLR, 2024f.

- Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Founda- tions of Machine Learning*. The MIT Press, 2012. ISBN 026201825X.
- Montreuil, Y., Carlier, A., Ng, L. X., and Ooi, W. T.
 Learning-to-defer for extractive question answering,
 2024. URL https://arxiv.org/abs/2410.
 15761.
- Mozannar, H. and Sontag, D. Consistent estimators for learning to defer to an expert. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20, 2020.
- Mozannar, H., Lang, H., Wei, D., Sattigeri, P., Das, S.,
 and Sontag, D. A. Who should predict? exact algorithms for learning to defer to humans. In *Interna- tional Conference on Artificial Intelligence and Statistics*,
 2023. URL https://api.semanticscholar.
 org/CorpusID:255941521.
- Ramaswamy, H. G., Tewari, A., and Agarwal, S. Consistent algorithms for multiclass classification with an abstain option. *Electronic Journal of Statistics*, 12(1):530 – 554, 2018. doi: 10.1214/17-EJS1388. URL https://doi. org/10.1214/17-EJS1388.
- Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. You
 only look once: Unified, real-time object detection. In
 2016 IEEE Conference on Computer Vision and Pattern
 Recognition (CVPR), pp. 779–788, 2016.
- 581 Ren, S., He, K., Girshick, R., and Sun, J. Faster r-cnn:
 582 Towards real-time object detection with region proposal
 583 networks, 2016. URL https://arxiv.org/abs/
 584 1506.01497.

- Steinwart, I. How to compare different loss functions and their risks. *Constructive Approximation*, 26:225–287,
 2007. URL https://api.semanticscholar. org/CorpusID:16660598.
- 590 Tailor, D., Patra, A., Verma, R., Manggala, P., and Nal-591 isnick. E. Learning to defer to a population: A 592 meta-learning approach. In Dasgupta, S., Mandt, S., 593 and Li, Y. (eds.), Proceedings of The 27th Interna-594 tional Conference on Artificial Intelligence and Statis-595 tics, volume 238 of Proceedings of Machine Learn-596 ing Research, pp. 3475-3483. PMLR, 02-04 May 597 2024. URL https://proceedings.mlr.press/ 598 v238/tailor24a.html. 599
- Tewari, A. and Bartlett, P. L. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8(36):1007–1025, 2007. URL http: //jmlr.org/papers/v8/tewari07a.html.

- Thorndike, R. L. Who belongs in the family? *Psy-chometrika*, 18:267–276, 1953. URL https: //api.semanticscholar.org/CorpusID: 120467216.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips. cc/paper_files/paper/2017/file/ 3f5ee243547dee91fbd053c1c4a845aa-Paper. pdf.
- Verma, R. and Nalisnick, E. Calibrated learning to defer with one-vs-all classifiers. In *International Conference* on Machine Learning, pp. 22184–22202. PMLR, 2022.
- Verma, R., Barrejon, D., and Nalisnick, E. Learning to defer to multiple experts: Consistent surrogate losses, confidence calibration, and conformal ensembles. In *International Conference on Artificial Intelligence and Statistics*, 2022. URL https://api.semanticscholar. org/CorpusID:253237048.
- Yang, C., Wu, Z., Jiang, P., Lin, Z., Gao, J., Danek, B., and Sun, J. PyHealth: A deep learning toolkit for healthcare predictive modeling. In *Proceedings of the 27th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD) 2023*, 2023. URL https: //github.com/sunlabuiuc/PyHealth.
- Zhang, M. and Agarwal, S. Bayes consistency vs. h-consistency: The interplay between surrogate loss functions and the scoring function class. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 16927–16936. Curran Associates, Inc., 2020. URL https://proceedings.neurips. cc/paper_files/paper/2020/file/ c4c28b367e14df88993ad475dedf6b77-Paper. pdf.
- Zhang, T. Statistical behavior and consistency of classification methods based on convex risk minimization. *Annals* of *Statistics*, 32, 12 2002. doi: 10.1214/aos/1079120130.

605 A. Algorithm

Algorithm 1 Two-Stage Learning-to-Defer for Multi-Task Learning Algorithm **Input:** Dataset $\{(x_k, y_k, t_k)\}_{k=1}^K$, multi-task model $g \in \mathcal{G}$, experts $m \in \mathcal{M}$, rejector $r \in \mathcal{R}$, number of epochs EPOCH, batch size B, learning rate η . **Initialization:** Initialize rejector parameters θ . for i = 1 to EPOCH do Shuffle dataset $\{(x_k, y_k, t_k)\}_{k=1}^K$. for each mini-batch $\mathcal{B} \subset \{(x_k, y_k, t_k)\}_{k=1}^K$ of size B do Extract input-output pairs $z = (x, y, t) \in \mathcal{B}$. Query model g(x) and experts m(x). {Agents are pre-trained and fixed} {Compute task-specific costs} Evaluate costs $c_0(g(x), z)$ and $c_{j>0}(m(x), z)$. Compute rejector prediction $r(x) = \arg \max_{j \in \mathcal{A}} r(x, j)$. {Rejector decision} Compute surrogate deferral empirical risk $\widehat{\mathcal{E}}_{\Phi_{def}}$: $\widehat{\mathcal{E}}_{\Phi_{\text{def}}} = \frac{1}{B} \sum_{z \in \mathcal{B}} \left[\Phi_{\text{def}}(g, r, m, z) \right].$ Update parameters θ using gradient descent: {Empirical risk computation} $\theta \leftarrow \theta - \eta \nabla_{\theta} \widehat{\mathcal{E}}_{\Phi_{\text{def}}}.$ {Parameter update} end for end for **Return:** trained rejector model r^* .

We will prove key lemmas and theorems stated in our main paper.

B. Proof of Lemma 4.2

We aim to prove Lemma 4.2, which establishes the optimal deferral decision by minimizing the conditional risk.

By definition, the Bayes-optimal rejector $r^B(x)$ minimizes the conditional risk $C_{\ell_{def}}$, given by:

(

$$\mathcal{C}_{\ell_{\text{def}}}(g, r, x) = \mathbb{E}_{y, t|x}[\ell_{\text{def}}(g, r, m, z)].$$
(8)

Expanding the expectation, we obtain:

$$C_{\ell_{\text{def}}}(g, r, x) = \mathbb{E}_{y, t|x} \left[\sum_{j=0}^{J} c_j(g(x), m_j(x), z) \mathbf{1}_{r(x)=j} \right].$$
(9)

Using the linearity of expectation, this simplifies to:

$$C_{\ell_{\text{def}}}(g, r, x) = \sum_{j=0}^{J} \mathbb{E}_{y,t|x} \left[c_j(g(x), m_j(x), z) \right] \mathbf{1}_{r(x)=j}.$$
(10)

Since we seek the rejector that minimizes the expected loss, the Bayes-conditional risk is given by:

$$\mathcal{C}^{B}_{\ell_{\mathrm{def}}}(g,r,x) = \inf_{g \in \mathcal{G}, r \in \mathcal{R}} \mathbb{E}_{y,t|x}[\ell_{\mathrm{def}}(g,r,m,z)].$$
(11)

Rewriting this expression, we obtain:

$$\mathcal{C}^{B}_{\ell_{\text{def}}}(g,r,x) = \inf_{r \in \mathcal{R}} \mathbb{E}_{y,t|x} \left[\inf_{g \in \mathcal{G}} c_0(g(x),z) \mathbf{1}_{r(x)=0} + \sum_{j=1}^J c_j(m_j(x),z) \mathbf{1}_{r(x)=j} \right].$$
(12)

656 This leads to the following minimization problem:

$$\mathcal{C}^{B}_{\ell_{\text{def}}}(g,r,x) = \min\left\{\inf_{g\in\mathcal{G}} \mathbb{E}_{y,t|x}\left[c_{0}(g(x),z)\right], \min_{j\in[J]} \mathbb{E}_{y,t|x}\left[c_{j}(m_{j}(x),z)\right]\right\}.$$
(13)

To simplify notation, we define:

$$\overline{c}_{j}^{*} = \begin{cases} \inf_{g \in \mathcal{G}} \mathbb{E}_{y,t|x}[c_{0}(g(x), z)], & \text{if } j = 0, \\ \mathbb{E}_{y,t|x}[c_{j}(m_{j}(x), z)], & \text{otherwise.} \end{cases}$$
(14)

Thus, the Bayes-conditional risk simplifies to:

$$\mathcal{C}^B_{\ell_{\text{def}}}(g, r, x) = \min_{j \in \mathcal{A}} \overline{c}^*_j.$$
(15)

Since the rejector selects the decision with the lowest expected cost, the optimal rejector is given by:

$$r^{B}(x) = \begin{cases} 0, & \text{if } \inf_{g \in \mathcal{G}} \mathbb{E}_{y,t|x}[c_{0}(g(x),z)] \leq \min_{j \in [J]} \mathbb{E}_{y,t|x}[c_{j}(m_{j}(x),z)],\\ j, & \text{otherwise.} \end{cases}$$
(16)

This completes the proof.

C. Proof Theorem 4.4

Before proving the desired Theorem 4.4, we will use the following Lemma C.1 (Awasthi et al., 2022; Mao et al., 2024e): **Lemma C.1** (\mathcal{R} -consistency bound). Assume that the following \mathcal{R} -consistency bounds holds for $r \in \mathcal{R}$, and any distribution

$$\mathcal{E}_{\ell_{01}}(r) - \mathcal{E}^*_{\ell_{01}}(\mathcal{R}) + \mathcal{U}_{\ell_{01}}(\mathcal{R}) \leq \Gamma^{\nu}(\mathcal{E}_{\Phi^{\nu}_{01}}(r) - \mathcal{E}^*_{\Phi^{\nu}_{01}}(\mathcal{R}) + \mathcal{U}_{\Phi^{\nu}_{01}}(\mathcal{R}))$$

then for $p \in (p_0 \dots p_J) \in \Delta^{|\mathcal{A}|}$ and $x \in \mathcal{X}$, we get

$$\sum_{j=0}^{J} p_j \mathbf{1}_{r(x)\neq j} - \inf_{r \in \mathcal{R}} \sum_{j=0}^{J} p_j \mathbf{1}_{r(x)\neq j} \le \Gamma^{\nu} \Big(\sum_{j=0}^{J} p_j \Phi_{01}^{\nu}(r, x, j) - \inf_{r \in \mathcal{R}} \sum_{j=0}^{J} p_j \Phi_{01}^{\nu}(r, x, j) \Big)$$

Theorem 4.4 ((\mathcal{G}, \mathcal{R})-consistency bounds). Let $g \in \mathcal{G}$ be a multi-task model. Suppose there exists a non-decreasing function $\Gamma^{\nu}: \mathcal{R}^+ \to \mathbb{R}^+$ for $\nu \ge 0$, such that the \mathcal{R} -consistency bounds hold for any distribution \mathcal{D} :

$$\begin{split} \mathcal{E}_{\Phi_{01}^{\nu}}(r) &- \mathcal{E}_{\Phi_{01}^{\nu}}^{*}(\mathcal{R}) + \mathcal{U}_{\Phi_{01}^{\nu}}(\mathcal{R}) \geq \\ \Gamma^{\nu}(\mathcal{E}_{\ell_{01}}(r) - \mathcal{E}_{\ell_{01}}^{B}(\mathcal{R}) + \mathcal{U}_{\ell_{01}}(\mathcal{R})) \end{split}$$

then for any $(g, r) \in \mathcal{G} \times \mathcal{R}$, any distribution \mathcal{D} and any $x \in \mathcal{X}$,

$$egin{aligned} \mathcal{E}_{\ell_{def}}(g,r) &- \mathcal{E}^B_{\ell_{def}}(\mathcal{G},\mathcal{R}) + \mathcal{U}_{\ell_{def}}(\mathcal{G},\mathcal{R}) \leq \ &\overline{\Gamma}^
u\left(\mathcal{E}_{\Phi^
u_{def}}(r) - \mathcal{E}^*_{\Phi^
u_{def}}(\mathcal{R}) + \mathcal{U}_{\Phi^
u_{def}}(\mathcal{R})
ight) \ &+ \mathcal{E}_{c_0}(g) - \mathcal{E}^B_{c_0}(\mathcal{G}) + \mathcal{U}_{c_0}(\mathcal{G}), \end{aligned}$$

where $\overline{\Gamma}^{\nu}(u) = \|\overline{\tau}\|_1 \Gamma^{\nu}\left(\frac{u}{\|\overline{\tau}\|_1}\right)$, with $\Gamma^{\nu}(u) = \mathcal{T}^{-1,\nu}(u)$, and for the log-softmax surrogate, $\mathcal{T}^{\nu=1}(u) = \frac{1+u}{2}\log(1+u)$ $u) + \frac{1-u}{2}\log(1-u).$

Proof. Let denote a cost for $j \in \mathcal{A} = \{0, \dots, J\}$:

$$\bar{c}_j^* = \begin{cases} \inf_{g \in \mathcal{G}} \mathbb{E}_{y,t|x}[c_0(g(x), z)] & \text{if } j = 0\\ \mathbb{E}_{y,t|x}[c_j(m(x), z)] & \text{otherwise} \end{cases}$$

Using the change of variables and the Bayes-conditional risk introduced in the proof of Lemma 4.2 in Appendix B, we have:

709
710
$$\mathcal{C}^B_{\ell_{def}}(\mathcal{G}, \mathcal{R}, x) = \min_{i \in \mathcal{A}} \overline{c}^*_i$$

711
$$j \in \mathcal{A}$$

711
712
713
714

$$\mathcal{C}_{\ell_{def}}(g,r,x) = \sum_{j=0}^{J} \mathbb{E}_{y,t|x} \Big[c_j(g(x), m_j(x), z) \Big] \mathbf{1}_{r(x)=j}$$

(17)

715 We follow suit for our surrogate Φ_{def} and derive its conditional risk and optimal conditional risk.

$$\mathcal{C}_{\Phi_{\text{def}}} = \mathbb{E}_{y,t|x} \Big[\sum_{j=1}^{J} c_j(m(x), z) \Phi_{01}^{\nu}(r, x, 0) + \sum_{j=1}^{J} \Big(c_0(g(x), z) + \sum_{i=1}^{J} c_i(m_i(x), z) \mathbf{1}_{j \neq i} \Big) \Phi_{01}^{\nu}(r, x, j) \\ \mathcal{C}_{\Phi_{\text{def}}}^* = \inf_{r \in \mathcal{R}} \mathbb{E}_{y,t|x} \Big[\sum_{j=1}^{J} c_j(g(x), m(x), z) \Phi_{01}^{\nu}(r, x, 0) + \sum_{j=1}^{J} [c_0(g(x), z) + \sum_{i=1}^{J} c_i(m_i(x), z) \mathbf{1}_{j \neq i}] \Phi_{01}^{\nu}(r, x, j) \Big] \Big]$$

Let us define the function $v(m(x), z) = \min_{j \in [J]} \overline{c}_j(m_j(x), z)$, where $m_j(x)$ denotes the model's output and \overline{c}_j represents the corresponding cost function. Using this definition, the calibration gap is formulated as $\Delta C_{\ell_{def}} := C_{\ell_{def}} - C_{\ell_{def}}^B$, where $C_{\ell_{def}}$ represents the original calibration term and $C_{\ell_{def}}^B$ denotes the baseline calibration term. By construction, the calibration gap satisfies $\Delta C_{\ell_{def}} \ge 0$, leveraging the risks derived in the preceding analysis.

$$\Delta C_{\ell_{def}} = \mathbb{E}_{y,t|x} \Big[\rho(g(x), z) \mathbf{1}_{r(x)=0} + \sum_{j=1}^{J} \Big(\rho(m(x), z) + \beta_j \Big) \mathbf{1}_{r(x)=j} \Big] \\ - v(m(x), z) + \Big(v(m(x), z) - \min_{j \in \mathcal{A}} \overline{c}_j^*(g(x), m(x), z) \Big)$$

Let us consider $\Delta C_{\ell_{def}} = A_1 + A_2$, such that:

$$A_{1} = \mathbb{E}_{y,t|x} \Big[\mathbb{1}_{r(x)=0} \rho(g(x), z) + \sum_{j=1}^{J} \mathbb{1}_{r(x)=j} \Big(\rho(m_{j}(x), z) + \beta_{j} \Big) \Big] - v(m(x), z)$$

$$A_{2} = \Big(v(m(x), z) - \min_{j \in \mathcal{A}} \overline{c}_{j}(g(x), m(x), z) \Big)$$
(18)

By considering the properties of min, we also get the following inequality:

$$v(m(x), z) - \min_{j \in \mathcal{A}} \overline{c}_{j}^{*}(g(x), m(x), z) \le \mathbb{E}_{y,t|x}[c_{0}(g(x), z)] - \inf_{g \in \mathcal{G}} \mathbb{E}_{y,t|x}[c_{0}(g(x), z)]$$
(19)

implying,

$$\Delta \mathcal{C}_{\ell_{\text{def}}} \le A_1 + \overline{c}_0(g(x), z) - \overline{c}_0^*(g(x), z) \tag{20}$$

We now select a distribution for our rejector. We first define $\forall j \in A$,

$$p_0 = \frac{\sum_{j=1}^{J} \bar{c}_j(m_j(x), z)}{J \sum_{j=0}^{J} \bar{c}_j(g(x), m_j(x), z)}$$

757 and

$$p_{j\in[J]} = \frac{\overline{c}_0(g(x), z) + \sum_{j\neq j'}^J \overline{c}'_j(m_j(x), z)}{J \sum_{j=0}^J \overline{c}_j(g(x), m_j(x), z)}$$

which can also be written as:

$$p_j = \frac{\overline{\tau}_j}{\|\overline{\boldsymbol{\tau}}\|_1} \tag{21}$$

765 Injecting the new distribution, we obtain the following:

$$\Delta \mathcal{C}_{\Phi_{\rm def}} = \|\overline{\tau}\|_1 \Big(\sum_{j=0}^J p_j \Phi_{01}^{\nu}(r, x, j) - \inf_{r \in \mathcal{R}} \sum_{j=0}^J p_j \Phi_{01}^{\nu}(r, x, j) \Big)$$
(22)

770	Now consider the first and last term of $\Delta C_{\ell_{def}}$. Following the intermediate step for Lemma 4.3, we have:
771	Ţ
772	$A = \mathbb{E} \left[\left[\left(\left(a(m) \times m \right)^2 \right) \right] + \sum_{n=1}^{J} \mathbb{E} \left[\left[\left(\left(m \times m \right)^2 \right)^2 \right] + \left(\left(\left(m \times m \right)^2 \right)^2 \right) \right] + \left(\left(\left(\left(m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(m \times m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \right)^2 \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \times m \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \times m \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \times m \right)^2 \right)^2 \right) + \left(\left(\left(\left(\left(m \times m \times m \times m \right)^2 \right) + \left(\left(\left(\left(\left(m \times m $
773	$A_1 = \mathbb{E}_{y,t x}[c_0(g(x), z)]1_{r(x)=0} + \sum_{i=1}^{\infty} \mathbb{E}_{y,t x}[c_j(m_j(x), z)]1_{r(x)=j} - v(m(x), z)$
774	j=1
775	$\mathbb{T}_{2} \left[\left(\left(\right) \right) \right]_{1} + \sum_{j=1}^{J} \mathbb{T}_{2} \left[\left(\left(\left(\right) \right) \right]_{1} \right]_{j}$
776	$= \mathbb{E}_{y,t x}[c_0(g(x),z)]1_{r(x)=0} + \sum_{i=1}^{\infty} \mathbb{E}_{y,t x}[c_j(m_j(x),z)]1_{r(x)=j}$
777	j=1
778	
779	$-\inf_{r\in\mathcal{R}}\left[\mathbb{E}_{y,t x}[c_0(g(x),z)]1_{r(x)=0}+\sum \mathbb{E}_{y,t x}[c_j(m_j(x),z)]1_{r(x)=j}\right]$
780	$j{=}1$
781	$\sum_{j=1}^{J} \left(-\frac{1}{2} \right) \left(-\frac{1}{2} \left(-\frac{1}{2} \right) \right) + \sum_{j=1}^{J} \left(-\frac{1}{2} \left(-\frac{1}{2} \right) \right) \left(-\frac{1}{2} \right) \left(-\frac{1}{2} \left(-\frac{1}{2} \right) \right) \right)$
782	$= \sum c_j(z, m_j) 1_{r(x) \neq 0} + \sum \left(c_0(g(x), z) + \sum c_{j'}(m_{j'}(x), z) \right) 1_{r(x) \neq j}$
783	$j=1$ $j=1$ $j\neq j'$
784	$\sum_{i=1}^{J} \left(\sum_{j=1}^{J} \left(\sum_{i=1}^{J} \left(\sum_{j=1}^{J} \left(\sum_{j$
785	$-\inf_{r\in\mathcal{R}}\left[\sum_{c_j}c_j(m_{j'}(x),z)1_{r(x)\neq 0}+\sum_{c_j}\left(c_0(g(x),z)+\sum_{c_{j'}}c_{j'}(m_{j'}(x),z)\right)1_{r(x)\neq j}\right]$
786	$j=1$ $j=1$ $j\neq j'$
787	Then, applying a change of variables to introduce $\ \overline{\tau}\ _1$, we get:
788	
789	
790	$\ \boldsymbol{\tau}\ _{1}p_{0}1_{r(x)\neq 0} + \ \boldsymbol{\tau}\ _{1}\sum_{p,j}p_{j}1_{r(x)\neq j} - \inf_{r\in\mathcal{R}}[\ \boldsymbol{\tau}\ _{1}p_{0}1_{r(x)\neq 0} + \ \boldsymbol{\tau}\ _{1}\sum_{p,j}p_{j}1_{r(x)\neq j}]$
791	j=1 $j=1$
792	$\ -\ \sum_{i=1}^{J} (1 - i) + c \ -\ \sum_{i=1}^{J} (1 -$
793	$= \ oldsymbol{ au}\ _1 \sum p_j \mathbbm{1}_{r(x) eq j} - \inf_{r\in\mathcal{R}} \ oldsymbol{ au}\ _1 \sum p_j \mathbbm{1}_{r(x) eq j}$
794	j=0 $j=0$
795	We now apply Lemma C 1 to introduce Γ

We now apply Lemma C.1 to introduce Γ , 796

$$\sum_{j=0}^{J} p_{j} \mathbf{1}_{r(x)\neq j} - \inf_{r \in \mathcal{R}} \sum_{j=0}^{J} p_{j} \mathbf{1}_{r(x)\neq j} \leq \Gamma \Big(\sum_{j=0}^{J} p_{j} \Phi_{01}^{\nu}(r, x, j) - \inf_{r \in \mathcal{R}} \sum_{j=0}^{J} p_{j} \Phi_{01}^{\nu}(r, x, j) \Big) \\ \frac{1}{\|\overline{\tau}\|_{1}} \Big[\sum_{j=0}^{J} \overline{\tau}_{j} \mathbf{1}_{r(x)\neq j} - \inf_{r \in \mathcal{R}} \sum_{j=0}^{J} \overline{\tau}_{j} \mathbf{1}_{r(x)\neq j} \Big] \leq \Gamma \Big(\frac{1}{\|\overline{\tau}\|_{1}} \Big[\sum_{j=0}^{J} \overline{\tau}_{j} \Phi_{01}^{\nu}(r, x, j) - \inf_{r \in \mathcal{R}} \sum_{j=0}^{J} \overline{\tau}_{j} \Phi_{01}^{\nu}(r, x, j) \Big] \Big)$$

$$\Delta \mathcal{C}_{\ell_{def}} \leq \|\overline{\tau}\|_{1} \Gamma\Big(\frac{\Delta \mathcal{C}_{\Phi_{def}}}{\|\overline{\tau}\|_{1}} \Big)$$

$$(23)$$

We reintroduce the coefficient A_2 such that:

$$\begin{aligned} & & 807 \\ & & 808 \\ & & 809 \\ & & 809 \\ & & 810 \\ & & 811 \end{aligned} \qquad \Delta \mathcal{C}_{\ell_{def}} \leq \|\overline{\tau}\|_1 \Gamma\Big(\frac{\Delta \mathcal{C}_{\Phi_{def}}}{\|\overline{\tau}\|_1}\Big) + \mathbb{E}_{y,t|x}[c_0(g(x),z)] - \inf_{g \in \mathcal{G}} \mathbb{E}_{y,t|x}[c_0(g(x),z)] \quad \text{(upper bounding with Eq 19)} \end{aligned}$$

Mao et al. (2023b) introduced a tight bound for the comp-sum surrogates family. It follows for $\nu \ge 0$ the inverse transformation $\Gamma^{\nu}(u) = \mathcal{T}^{-1,\nu}(u)$:

814 815 816

817

806

$$\left\{\frac{2^{1-\nu}}{1-\nu}\left[1-\left(\frac{(1+\nu)^{\frac{2-\nu}{2}}+(1-\nu)^{\frac{2-\nu}{2}}}{2}\right)^{2-\nu}\right] \qquad \nu \in [0,1)$$

818
819

$$\frac{1+v}{2}\log[1+v] + \frac{1-v}{2}\log[1-v]$$
 $\nu = 1$

$$T^{\nu}(v) = \begin{cases} T^{\nu}(v) = \begin{cases} f(v) + \frac{2-\nu}{2} + f(v) + \frac{2-\nu}{2} \end{cases} \end{cases}$$

$$\begin{bmatrix} \frac{1}{(\nu-1)n^{\nu-1}} \left[\left(\frac{(1+\nu)^{-2} + (1-\nu)^{-2}}{2} \right) & -1 \right] \quad \nu \in (1,2) \end{bmatrix}$$

825 We note $\overline{\Gamma}^{\nu}(u) = \|\overline{\tau}\|_1 \Gamma^{\nu}(\frac{u}{\|\overline{\tau}\|_1})$. By applying Jensen's Inequality and taking expectation on both sides, we get

$$\begin{split} \mathcal{E}_{\ell_{\mathrm{def}}}(g,r) &- \mathcal{E}^{B}_{\ell_{\mathrm{def}}}(\mathcal{G},\mathcal{R}) + \mathcal{U}_{\ell_{\mathrm{def}}}(\mathcal{G},\mathcal{R}) \\ &\leq \overline{\Gamma}^{\nu}(\mathcal{E}_{\Phi_{\mathrm{def}}}(r) - \mathcal{E}^{*}_{\Phi_{\mathrm{def}}}(\mathcal{R}) + \mathcal{U}_{\Phi_{\mathrm{def}}}(\mathcal{R})) + \mathcal{E}_{c_{0}}(g) - \mathcal{E}^{B}_{c_{0}}(\mathcal{G}) + \mathcal{U}_{c_{0}}(\mathcal{G}) \end{split}$$

D. Proof Theorem 4.6

Theorem 4.6 (Characterization Minimizability Gaps). Assume \mathcal{R} symmetric and complete. Then, for the cross-entropy multiclass surrogates Φ_{01}^{ν} and any distribution \mathcal{D} , it follows for $\nu \geq 0$:

$$\mathcal{C}_{\Phi_{def}^{\nu,*}}^{\nu,*} = \begin{cases} \|\overline{\boldsymbol{\tau}}\|_1 H\left(\frac{\overline{\boldsymbol{\tau}}}{\|\overline{\boldsymbol{\tau}}\|_1}\right) & \nu = 1\\ \|\overline{\boldsymbol{\tau}}\|_1 - \|\overline{\boldsymbol{\tau}}\|_{\infty} & \nu = 2\\ \frac{1}{\nu-1} \left[\|\overline{\boldsymbol{\tau}}\|_1 - \|\overline{\boldsymbol{\tau}}\|_{\frac{1}{2-\nu}} \right] & \nu \in (1,2)\\ \frac{1}{1-\nu} \left[\left(\sum_{k=0}^J \overline{\boldsymbol{\tau}}_k^{\frac{1}{2-\nu}}\right)^{2-\nu} - \|\overline{\boldsymbol{\tau}}\|_1 \right] & otherwise, \end{cases}$$

then the minimizability gap is,

$$\mathcal{U}_{\Phi_{def}^{\nu}}(\mathcal{R}) = \mathcal{E}_{\Phi_{def}^{\nu}}^{*}(\mathcal{R}) - \mathbb{E}_{x}[\inf_{r \in \mathcal{R}} \ \mathcal{C}_{\Phi_{def}^{\nu}}^{\nu}(r, x)]$$

with $\overline{\tau} = \{\mathbb{E}_{y,t|x}[\overline{\tau}_0], \dots, \mathbb{E}_{y,t|x}[\overline{\tau}_J]\}$, the aggregated costs $\tau_j = \sum_{k=0}^J c_k \mathbf{1}_{k\neq j}$, and the Shannon Entropy H.

Proof. We define the softmax distribution as $s_j = \frac{e^{r(x,j)}}{\sum_{j' \in \mathcal{A}} e^{r(x,j')}}$, where $s_j \in [0,1]$. Let $\overline{\tau}_j = \overline{\tau}_j(g(x), m(x), z)$ with $\tau_j \in \mathbb{R}^+$, and denote the expected value as $\overline{\tau} = \mathbb{E}_{y,t|x}[\tau]$. We now derive the conditional risk for a given $\nu \ge 0$:

$$\mathcal{C}_{\Phi_{def}}^{\nu}(r,x) = \sum_{j=0}^{J} \mathbb{E}_{y,t|x}[\tau_{j}] \Phi_{01}^{\nu}(r,x,j) \\
= \begin{cases} \frac{1}{1-\nu} \sum_{j=0}^{J} \overline{\tau}_{j} \left[\left(\sum_{j' \in \mathcal{A}} e^{r(x,j') - r(x,j)} \right)^{1-\nu} - 1 \right] & \nu \neq 1 \\ \sum_{j=0}^{J} \overline{\tau}_{j} \log \left(\sum_{j' \in \mathcal{A}} e^{r(x,j') - r(x,j)} \right) & \nu = 1 \end{cases} \\
= \begin{cases} \frac{1}{1-\nu} \sum_{j=0}^{J} \overline{\tau}_{j} \left[s_{j}^{\nu-1} - 1 \right] & \nu \neq 1 \\ -\sum_{j=0}^{J} \overline{\tau}_{j} \log(s_{j}) & \nu = 1 \end{cases}$$
(24)

865 For $\nu = 1$: we can write the following conditional risk:

$$\mathcal{C}_{\Phi_{\text{def}}}^{\nu=1}(r,x) = -\sum_{j=0}^{J} \overline{\tau}_j \Big[r(x,j) - \log \sum_{j' \in \mathcal{A}} e^{r(x,j')} \Big]$$
(25)

871 Then,

$$\frac{\partial \mathcal{C}_{\Phi_{\rm def}}^{\nu=1}}{\partial r(x,i)}(r,x) = -\overline{\tau}_i + \Big(\sum_{j=0}^J \overline{\tau}_j\Big)s_i^* \tag{26}$$

At the optimum, we have:

$$s^*(x,i) = \frac{\overline{\tau}_i}{\sum_{j=0} \overline{\tau}_j} \tag{27}$$

880 Then, it follows:

$$\mathcal{C}_{\Phi_{\rm def}}^{*,\nu=1}(\mathcal{R},x) = -\sum_{j=0}^{J} \overline{\tau}_j \log\left(\frac{\overline{\tau}_j}{\sum_{j'=0} \overline{\tau}_{j'}}\right)$$
(28)

As the softmax parametrization is a distribution $s^* \in \Delta^{|\mathcal{A}|}$, we can write this conditional in terms of entropy with $\overline{\tau} = \{\overline{\tau}_j\}_{j \in \mathcal{A}}$:

$$C_{\Phi_{def}}^{*,\nu=1}(\mathcal{R},x) = -\left(\sum_{k=0}^{J} \overline{\tau}_{k}\right) \sum_{j=0} s_{j}^{*} \log(s_{j}^{*})$$

$$= \left(\sum_{k=0}^{J} \overline{\tau}_{k}\right) H\left(\frac{\overline{\tau}}{\sum_{j'=0} \overline{\tau}_{j'}}\right)$$

$$= \|\overline{\tau}\|_{1} H\left(\frac{\overline{\tau}}{\|\overline{\tau}\|_{1}}\right) \quad (\text{as } \tau_{j} \in \mathbb{R}^{+})$$
(29)

For $\nu \neq 1, 2$: The softmax parametrization can be written as a constraint $\sum_{j=0}^{J} s_j = 1$ and $s_j \ge 0$. Consider the objective 897

$$\Phi(\mathbf{s}) = \frac{1}{1-\nu} \sum_{j=0}^{J} \overline{\tau}_j \left[s_j^{\nu-1} - 1 \right].$$
(30)

903 We aim to find $s^* = (s_0^*, \dots, s_J^*)$ that minimizes (30) subject to $\sum_{j=0}^J s_j = 1$. Introduce a Lagrange multiplier λ for the 903 normalization $\sum_{j=0}^J s_j = 1$. The Lagrangian is:

$$\mathcal{L}(\mathbf{s},\lambda) = \frac{1}{1-\nu} \sum_{j=0}^{J} \overline{\tau}_{j} \left[s_{j}^{\nu-1} - 1 \right] + \lambda \left(1 - \sum_{j=0}^{J} s_{j} \right).$$
(31)

908 We take partial derivatives with respect to s_i :

$$\frac{\partial \mathcal{L}}{\partial s_i} = \frac{1}{1-\nu} \,\overline{\tau}_i \left(\nu - 1\right) s_i^{\nu-2} - \lambda = 0. \tag{32}$$

Since $\frac{\nu-1}{1-\nu} = -1$, we get

$$\overline{\tau}_i s_i^{\nu-2} = -\lambda > 0 \implies s_i^{\nu-2} = \frac{\alpha}{\overline{\tau}_i} \text{ for some } \alpha > 0.$$
(33)

Hence

$$s_i = \left(\frac{\alpha}{\overline{\tau}_i}\right)^{\frac{1}{\nu-2}}.$$
(34)

920 Summing s_i over $\{i = 0, ..., J\}$ and setting the total to 1 yields:

$$\sum_{i=0}^{J} \left(\frac{\alpha}{\overline{\tau}_i}\right)^{\frac{1}{\nu-2}} = 1.$$
(35)

(37)

925 Let

$$\alpha^{\frac{1}{\nu-2}} = \frac{1}{\sum_{k=0}^{J} \left(\frac{1}{\overline{\tau}_{k}}\right)^{\frac{1}{\nu-2}}} \implies \alpha = \left[\sum_{k=0}^{J} \left(\frac{1}{\overline{\tau}_{k}}\right)^{\frac{1}{\nu-2}}\right]^{\nu-2}.$$
(36)

929 Therefore, for each i,

 $s_i^* = \left(\frac{\alpha}{\overline{\tau}_i}\right)^{\frac{1}{\nu-2}} = \frac{\overline{\tau}_i^{\frac{1}{2-\nu}}}{\sum\limits_{k=0}^{J} \overline{\tau}_k^{\frac{1}{2-\nu}}}.$

935 This $\{s_i^*\}$ is a valid probability distribution. Let

$$A = \sum_{k=0}^{J} \tau_k^{\frac{1}{2-\nu}}.$$
(38)

940 Then the optimum distribution is

$$s_i^* = \frac{\overline{\tau_i}^{\frac{1}{2-\nu}}}{A}.$$
 (39)

944 Recall

$$\Phi(\mathbf{s}) = \frac{1}{1-\nu} \sum_{j=0}^{J} \overline{\tau}_j \left[s_j^{\nu-1} - 1 \right].$$
(40)

At
$$s_j^*$$
, we have

$$(s_j^*)^{\nu-1} = \left(\frac{\overline{\tau}_j^{\frac{1}{2-\nu}}}{A}\right)^{\nu-1} = \frac{\overline{\tau}_j^{\frac{\nu-1}{2-\nu}}}{A^{\nu-1}}.$$
(41)

952 Hence

$$\sum_{j=0}^{J} \overline{\tau}_{j} \left(s_{j}^{*}\right)^{\nu-1} = \frac{1}{A^{\nu-1}} \sum_{j=0}^{J} \overline{\tau}_{j}^{1+\frac{\nu-1}{2-\nu}} = \frac{1}{A^{\nu-1}} \sum_{j=0}^{J} \overline{\tau}_{j}^{\frac{1}{2-\nu}} = \frac{A}{A^{\nu-1}} = A^{2-\nu}.$$
 (42)

Substituting back,

$$\mathcal{C}_{\Phi_{\rm def}}^{*,\nu\neq1,2}(\mathcal{R},x) = \frac{1}{1-\nu} \left[\left(\sum_{k=0}^{J} \overline{\tau}_{k}^{\frac{1}{2-\nu}} \right)^{2-\nu} - \sum_{j=0}^{J} \overline{\tau}_{j} \right]$$
(43)

We can express this conditional risk with a valid $L^{(\frac{1}{2-\nu})}$ norm as long as $\nu \in (1,2)$.

$$\mathcal{C}_{\Phi_{\mathrm{def}}}^{*,\nu\neq1,2}(\mathcal{R},x) = \frac{1}{\nu-1} \left[\|\overline{\boldsymbol{\tau}}\|_1 - \|\overline{\boldsymbol{\tau}}\|_{\frac{1}{2-\nu}} \right]$$
(44)

For $\nu = 2$: Since $\sum_{j=0}^{J} \overline{\tau}_j = S$, we have

$$\mathcal{C}_{\Phi_{\rm def}}^{\nu=2}(r,x) = \sum_{j=0}^{J} \overline{\tau}_j \left[1 - s_j(r) \right] = \sum_{j=0}^{J} \overline{\tau}_j - \sum_{j=0}^{J} \overline{\tau}_j s_j(r).$$
(45)

Hence

$$\inf_{r \in \mathcal{R}} \mathcal{C}_{\Phi_{def}}^{\nu=2}(r,x) = S - \sup_{r \in \mathcal{R}} \sum_{j=0}^{J} \overline{\tau}_j s_j(r).$$
(46)

975 Therefore, minimizing $C_{\Phi_{def}}^{\nu=2}(r,x)$ is equivalent to maximizing

$$F(r) = \sum_{j=0}^{J} \overline{\tau}_j s_j(r).$$
(47)

Its partial derivative w.r.t. r_i is the standard softmax derivative:

$$\frac{\partial s_j}{\partial r_i} = s_j \left(\delta_{ij} - s_i \right) = \begin{cases} s_i \left(1 - s_i \right), & \text{if } i = j, \\ -s_j s_i, & \text{otherwise.} \end{cases}$$
(48)

Hence, for each
$$i$$

$$\frac{\partial F}{\partial r_i} = \sum_{j=0}^J \overline{\tau}_j \frac{\partial s_j}{\partial r_i} = \overline{\tau}_i s_i (1-s_i) + \sum_{\substack{j=0\\j\neq i}}^J \overline{\tau}_j (-s_j s_i).$$
(49)

985 Her

990 Factor out s_i :

$$\frac{\partial F}{\partial r_i} = s_i \left[\overline{\tau}_i \left(1 - s_i \right) - \sum_{j \neq i} \overline{\tau}_j s_j \right] = s_i \left[\overline{\tau}_i - \left(\sum_{j=0}^J \overline{\tau}_j s_j \right) \right], \tag{50}$$

because $\sum_{j \neq i} \overline{\tau}_j s_j = \sum_{j=0}^J \overline{\tau}_j s_j - \overline{\tau}_i s_i$. Define $F(r) = \sum_{j=0}^J \overline{\tau}_j s_j(r)$. Then:

$$\frac{\partial F}{\partial r_i} = s_i \left[\overline{\tau}_i - F(r) \right]. \tag{51}$$

Setting $\frac{\partial F}{\partial r_i} = 0$ for each *i* implies

$$s_i \left[\overline{\tau}_i - F(r) \right] = 0, \quad \forall i.$$
(52)

Thus, for each index *i*:

$$s_i = 0 \quad \text{or} \quad \overline{\tau}_i = F(r).$$
 (53)

To maximize F(r), notice that:

If \$\overline{\tau_{i^*}\$ is strictly the largest among all \$\overline{\tau_{i}\$, then the maximum is approached by making \$s_{i^*}\$ ≈ 1, so \$F(r)\$ ≈ \$\overline{\tau_{i^*}\$. In the softmax parameterization, this occurs in the limit \$r_{i^*}\$ → +∞ and \$r_k\$ → -∞ for \$k ≠ i^*\$.

• If there is a tie for the largest $\overline{\tau}_i$, we can put mass on those coordinates that share the maximum value. In any case, the supremum is $\max_i \overline{\tau}_i$.

1015 Hence

$$\sup_{r \in \mathcal{R}} F(r) = \max_{0 \le i \le J} \overline{\tau}_i.$$
(54)

1019 Because $C_{\Phi_{def}}^{\nu=2}(r, x) = S - F(r)$,

 $\inf_{r \in \mathcal{R}} \mathcal{C}_{\Phi_{def}}^{\nu=2}(r,x) = S - \sup_{r \in \mathcal{R}} F(r) = \sum_{i=0}^{J} \overline{\tau}_{j} - \max_{i \in \mathcal{A}} \overline{\tau}_{i} = \|\overline{\tau}\|_{1} - \|\overline{\tau}\|_{\infty}$ (55)

Hence the global minimum of $C_{\Phi_{def}}^{\nu=2}$ is $\|\overline{\tau}\|_1 - \|\overline{\tau}\|_{\infty}$. In the "softmax" parameterization, this is only approached in the limit as one coordinate r_{i^*} goes to $+\infty$ and all others go to $-\infty$. No finite r yields an exactly one-hot $s_i(r) = 1$, but the limit is enough to achieve the infimum arbitrarily closely.

 $\inf_{r \in \mathcal{R}} \mathcal{C}^{\nu}_{\Phi_{def}}(r, x) = \begin{cases} \|\overline{\tau}\|_1 H\left(\frac{\overline{\tau}}{\|\overline{\tau}\|_1}\right) & \nu = 1\\ \|\overline{\tau}\|_1 - \|\overline{\tau}\|_{\infty} & \nu = 2\\ \frac{1}{\nu - 1} \left[\|\overline{\tau}\|_1 - \|\overline{\tau}\|_{\frac{1}{2-\nu}}\right] & \nu \in (1, 2)\\ \frac{1}{\nu - 1} \left[\left(\sum_{l=0}^{J} \overline{\tau}_{l}^{\frac{1}{2-\nu}}\right)^{2-\nu} - \|\overline{\tau}\|_1\right] & \text{otherwise} \end{cases}$

1020 It follows for $\overline{\tau} = {\{\overline{\tau}_j\}}_{j \in \mathcal{A}}$ and $\nu \ge 0$:

Building on this, we can infer the minimizability gap:



 $\mathcal{U}_{\Phi_{\rm def}}(\mathcal{R}) = \mathcal{E}^*_{\Phi_{\rm def}}(\mathcal{R}) - \mathbb{E}_x[\inf_{x \in \mathcal{R}} \mathcal{C}^{\nu}_{\Phi_{\rm def}}(r, x)]$

(57)

(56)

1045 E. Proof Lemma 4.7

Lemma 4.7. Let \mathcal{L}_1 be a family of functions mapping \mathcal{X} to [0, 1], and let \mathcal{L}_2 be a family of functions mapping \mathcal{X} to $\{0, 1\}$. Define $\mathcal{L} = \{l_1 l_2 : l_1 \in \mathcal{L}_1, l_2 \in \mathcal{L}_2\}$. Then, the empirical Rademacher complexity of \mathcal{L} for any sample S of size K is bounded by:

$$\widehat{\mathfrak{R}}_{S}(\mathcal{L}) \leq \widehat{\mathfrak{R}}_{S}(\mathcal{L}_{1}) + \widehat{\mathfrak{R}}_{S}(\mathcal{L}_{2}).$$
(7)

1054 *Proof.* We define the function ψ as follows:

$$\psi: \begin{array}{ccc} \mathcal{L}_1 + \mathcal{L}_2 & \longrightarrow & \mathcal{L}_1 \mathcal{L}_2 \\ l_1 + l_2 & \longmapsto & (l_1 + l_2 - 1)_+ \end{array}$$
(58)

Here, $l_1 \in \mathcal{L}_1$ and $l_2 \in \mathcal{L}_2$. The function ψ is 1-Lipschitz as we have $t \mapsto (t-1)_+$ for $t = l_1 + l_2$. Furthermore, given that ψ is surjective and 1-Lipschitz, by Talagrand's lemma (Mohri et al., 2012), we have:

$$\hat{\mathfrak{R}}_{S}(\psi(\mathcal{L}_{1}+\mathcal{L}_{2})) \leq \hat{\mathfrak{R}}_{S}(\mathcal{L}_{1}+\mathcal{L}_{2}) \leq \hat{\mathfrak{R}}_{S}(\mathcal{L}_{1}) + \hat{\mathfrak{R}}_{S}(\mathcal{L}_{2})$$
(59)

This inequality shows that the Rademacher complexity of the sum of the losses is bounded by the sum of their individual complexities. \Box

1067

1069

1073

1062 1063 1064

¹⁰⁶⁸ **F. Proof Theorem 4.8**

Theorem 4.8 (Learning bounds of the deferral loss). For any expert M_j , any distribution \mathcal{D} over \mathcal{Z} , we have with probability $1 - \delta$ for $\delta \in [0, 1/2]$, that the following bound holds at the optimum:

$$\mathcal{E}_{\ell_{def}}(h, f, r) \leq \widehat{\mathcal{E}}_{\ell_{def}}(h, f, r) + 2\mathfrak{R}_{K}(\mathcal{L}_{def}) + \sqrt{\frac{\log 1/\delta}{2K}},$$

 $\Re_K(\mathcal{L}_{def}) \leq \frac{1}{2} \Re_K(\mathcal{H}) + \Re_K(\mathcal{F}) + \sum_{i=1}^J \Omega(m_j^h, y)$

+ $\Big(\sum_{j=1}^{J} \max \ell_{reg}(m_j^f, t) + 2\Big) \Re_K(\mathcal{R}),$

⁰⁷⁶ with

1077

1079

100

1082

1083

100

108

- 1087
- 100

Proof. We are interested in finding the generalization of $u = (g, r) \in \mathcal{L}$:

with $\Omega(m_i^h, y) = \frac{1}{2} \mathcal{D}(m_i^h \neq y) \exp\left(-\frac{K}{8} \mathcal{D}(m_i^h \neq y)\right) + \mathcal{R}_{K\mathcal{D}(m^h \neq y)/2}(\mathcal{R}).$

1100 Let's consider j = 0: 1102 $\frac{1}{K} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathcal{L}} \sum_{k=1}^{K} \sigma_{k} c_{0} \mathbf{1}_{r(x_{k})=0} \Big] = \frac{1}{K} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathcal{L}} \sum_{k=1}^{K} \sigma_{k} [\mathbf{1}_{h(x_{k}) \neq y} + \ell_{\text{reg}}(f(x_{k}), b_{k})] \mathbf{1}_{r(x_{k})=0} \Big]$ 1104 1105 $\leq \frac{1}{K} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathcal{L}} \sum_{k=1}^{K} \sigma_{k} \mathbf{1}_{h(x_{k}) \neq y} \mathbf{1}_{r(x_{k})=0} \Big] + \frac{1}{K} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathcal{L}} \sum_{k=1}^{K} \sigma_{k} \ell_{\text{reg}}(f(x_{k}), b_{k}) \mathbf{1}_{r(x_{k})=0} \Big] \Big]$ 1106 1108 $\leq \Big[\frac{1}{2} \mathfrak{R}_{K}(\mathcal{H}) + \mathfrak{R}_{K}(\mathcal{R}) \Big] + \Big[\mathfrak{R}_{K}(\mathcal{F}) + \mathfrak{R}_{K}(\mathcal{R}) \Big] \quad (\text{using Lemma 4.7})$ 1109 1100 $= \frac{1}{2} \mathfrak{R}_{K}(\mathcal{H}) + \mathfrak{R}_{K}(\mathcal{F}) + 2\mathfrak{R}_{K}(\mathcal{R})$ (60)

1113 Let's consider j > 0:

$$\frac{1}{K} \sum_{j=1}^{J} \mathbb{E}_{\sigma} \left[\sup_{r \in \mathcal{R}} \sum_{k=1}^{K} \sigma_{k} c_{j} \mathbf{1}_{r(x_{k})=j} \right] \leq \frac{1}{K} \sum_{j=1}^{J} \mathbb{E}_{\sigma} \left[\sup_{r \in \mathcal{R}} \sum_{k=1}^{K} \sigma_{k} \mathbf{1}_{m_{k,j}^{h} \neq y} \mathbf{1}_{r(x_{k})=j} \right]
+ \frac{1}{K} \sum_{j=1}^{J} \mathbb{E}_{\sigma} \left[\sup_{r \in \mathcal{R}} \sum_{k=1}^{K} \sigma_{k} \ell_{\text{reg}}(m_{k,j}^{f}, b_{k}) \mathbf{1}_{r(x_{k})=j} \right]$$
(61)

Using learning-bounds for single expert in classification (Mozannar & Sontag, 2020), we have: 1121

$$\frac{1}{K}\mathbb{E}_{\sigma}\left[\sup_{r\in\mathcal{R}}\sum_{k=1}^{K}\sigma_{k}\mathbf{1}_{m_{k}^{h}\neq y}\mathbf{1}_{r(x_{k})=1}\right] \leq \frac{\mathcal{D}(m^{h}\neq y)}{2}\exp\left(-\frac{K\mathcal{D}(m^{h}\neq y)}{8}\right) + \mathcal{R}_{K\mathcal{D}(m^{h}\neq y)/2}(\mathcal{R})$$
(62)

Applying it to our case:

$$\frac{1}{K}\sum_{j=1}^{J}\mathbb{E}_{\sigma}\left[\sup_{r\in\mathcal{R}}\sum_{k=1}^{K}\sigma_{k}\mathbf{1}_{m_{k,j}^{h}\neq y}\mathbf{1}_{r(x_{k})=j}\right] \leq \sum_{j=1}^{J}\left(\frac{\mathcal{D}(m_{j}^{h}\neq y)}{2}\exp\left(-\frac{K\mathcal{D}(m_{j}^{h}\neq y)}{8}\right) + \mathcal{R}_{K\mathcal{D}(m_{j}^{h}\neq y)/2}(\mathcal{R})\right)$$
(63)

For the last term,

$$\frac{1}{K}\sum_{j=1}^{J}\mathbb{E}_{\sigma}\left[\sup_{r\in\mathcal{R}}\sum_{k=1}^{K}\sigma_{k}\ell_{\text{reg}}(m_{k,j}^{f},b_{k})\mathbf{1}_{r(x_{k})=j}\right] \leq \sum_{j=1}^{J}\left(\max\ell_{reg}(m_{j}^{f},t)\mathfrak{R}_{K}(\mathcal{R})\right)$$
(64)

Then, it leads to:

$$\Re_{K}(\mathcal{L}_{def}) \leq \frac{1}{2} \Re_{K}(\mathcal{H}) + \Re_{K}(\mathcal{F}) + \sum_{j=1}^{J} \Omega(m_{j}^{h}, y) + \Big(\sum_{j=1}^{J} \max \ell_{\text{reg}}(m_{j}^{f}, t) + 2\Big) \Re_{K}(\mathcal{R})$$
$$\Omega(m_{j}^{h}, y) = \frac{\mathcal{D}(m_{j}^{h} \neq y)}{2} \exp\left(-\frac{K\mathcal{D}(m_{j}^{h} \neq y)}{8}\right) + \mathcal{R}_{K\mathcal{D}}(m_{j}^{h} \neq y)/2(\mathcal{R})$$

G. Experiments

with

G.1. PascalVOC Experiment

Since an image may contain multiple objects, our deferral rule is applied at the level of the entire image $x \in \mathcal{X}$, ensuring that the approach remains consistent with real-world scenarios.

	Model	M_1	M_2	
mAP	39.5	43.3	52.8	

Table 2. Agent accuracies on the CIFAR-100 validation set. Since the training and validation sets are pre-determined in this dataset, the agents' knowledge remains fixed throughout the evaluation.

A Two-Stage Learning-to-Defer Approach for Multi-Task Learning

1155		$\overline{\operatorname{Cost} \beta_2}$	mAP (%)	Model Allocation (%)	Expert	1 Allocation (%)	Expert 2 Allocation (%)	
1156		0.01	52.8 ± 0.0	0.0 ± 0.0		0.0 ± 0.0	100.0 ± 0.0	
1157		0.05	52.5 ± 0.1	7.3 ± 0.8		0.0 ± 0.0	92.7 ± 0.3	
1158		0.1	49.1 ± 0.6	48.0 ± 0.7		0.0 ± 0.0	52.0 ± 0.2	
1159		0.15	44.2 ± 0.4 42.0 ± 0.2	68.1 ± 0.3 77.5 ± 0.2		19.7 ± 0.4	12.2 ± 0.1	
1160		0.2	42.0 ± 0.2 40.1 ± 0.2	77.5 ± 0.2 98.1 ± 0.0		22.3 ± 0.3 19 ± 01	0.0 ± 0.0 0.0 ± 0.0	
1161		0.5	40.1 ± 0.2 39.5 ± 0.0	100.0 ± 0.0		0.0 ± 0.0	0.0 ± 0.0 0.0 ± 0.0	
1160								
1102	Table 3. Det	ailed resu	ilts across di	fferent cost values β_{i}	2. Errors	s represent the st	andard deviation over n	ultiple runs.
1103				, .				
1164								
1165	G.2. MIMIC-IV Ex	perime	nts					
1166		1	2022) :	1	J. : J.			
1167	MINIC-IV (Johnson	n et al.,	2025) is a	large collection of	de-ide	numed nearth-	related data covering	over forty thousand
1168	patients who stayed	in critica	al care units	s. This dataset inclu	ides a v	vide variety of	information, such as	demographic details,
1169	vital signs, laborator	y test re	sults, medi	cations, and procee	lures. F	for our analysis	s, we focus specificall	y on features related
1170	to <i>procedures</i> , which	h corres	pond to me	dical procedures p	erform	ed during hos	pital visits, and <i>diagn</i>	oses received by the
1171	patients.							
1172	Using these features	waadd	race two n	adiativa taskar (1)	o oloco	fightion tools to	a pradiat whathar a p	tiont will die during
1173	Using these realties	, we add	iless two pi	L'enformention from			(2) a maximum (2) a maximum (3)	allent will die during
1174	their next nospital vi	sit basec	i on clinica	i information from	the cur	rent visit, and	(2) a regression task t	o estimate the length
1175	of stay for the curren	it hospit	al visit base	ed on the same clin	ical inf	ormation.		
1175	A key challenge in t	his task	is the seve	re class imbalance	nartic	ularly in predic	cting mortality. To m	itigate this issue we
1170	sub-sample the nega	tive mor	tality class	retaining a balance	ed data	set with $K =$	5005 samples compr	ising 48.2% positive
11//	sub-sample the nega	1.8% no	antivo mort	, lotanning a Darano	dol is t	roined on 80%	of this detect while	the remaining 20% is
11/8	holtanty cases and 5	1.0% lie	gative mon					the remaining 20% is
1179	neid out for validatio	on. To er	isure consis	stency in the result	s, we m	xed the training	g and validation partit	lons.
1180								
1181				Μ	odel	$M_1 = M_2$		
1182				Accuracy	0.0	39.7 46.2		
1183				Smooth I 1 1	45	2 31 1 02		
1184					.40	2.01 1.02		
1185	Table 1 Performance	of the an	ents on the	MIMIC-IV dataset	evaluat	ed in terms of	accuracy and Smooth 1	1 loss We fixed the
1186	training/validation set	such that	the agents' l	nowledge remains fi	xed thro	ughout the evalu	uation	Li ioss. we fixed the
1187	training, variation set	such that	the ugents 1	liowieuge remains n	Act the	agnout the even	uution.	
1188								
1189								
1190								
1191								
1102								
1102								
1173								
1194								
1 1 1 1 1 1								
1195								
1195 1196								
1195 1196 1197								
1195 1196 1197 1198								
1195 1196 1197 1198 1199								
1195 1196 1197 1198 1199 1200								
1195 1196 1197 1198 1199 1200 1201								
1195 1196 1197 1198 1199 1200 1201 1202								
1195 1196 1197 1198 1199 1200 1201 1202 1203								
1195 1196 1197 1198 1199 1200 1201 1202 1203 1204								
1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205								
1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205								
1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205 1206								
1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205 1206 1207								
1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205 1206 1207 1208								
1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205 1206 1207 1208 1209								