

AI, Meet Human: Learning Paradigms for Hybrid Decision Making Systems

CLARA PUNZI, Scuola Normale Superiore, Italy
 ROBERTO PELLUNGRINI, Scuola Normale Superiore, Italy
 MATTIA SETZU, University of Pisa, Italy
 FOSCA GIANNOTTI, Scuola Normale Superiore, Italy
 DINO PEDRESCHI, University of Pisa, Italy

Everyday we increasingly rely on machine learning models to automate and support high-stake tasks and decisions. This growing presence means that humans are now constantly interacting with machine learning-based systems, training and using models everyday. Several different techniques in computer science literature account for the human interaction with machine learning systems, but their classification is sparse and the goals varied. This survey proposes a taxonomy of Hybrid Decision Making Systems, providing both a conceptual and technical framework for understanding how current computer science literature models interaction between humans and machines.

CCS Concepts: • **Computing methodologies** → **Learning paradigms**; • **Human-centered computing** → *HCI theory, concepts and models*.

Additional Key Words and Phrases: hybrid systems, hybrid decision making, cooperative AI

ACM Reference Format:

Clara Punzi, Roberto Pellungrini, Mattia Setzu, Fosca Giannotti, and Dino Pedreschi. 2024. AI, Meet Human: Learning Paradigms for Hybrid Decision Making Systems. *J. ACM* 56, 1, Article 26 (January 2024), 38 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

Advances in Artificial Intelligence (AI) systems have reached a widespread use in a variety of domains and applications. Their effectiveness has grown exponentially over the past decades, and now AI systems are able to achieve surprising performances on tasks previously thought to be out of reach for artificial systems. Often acting as standalone systems with little to no human control, such models are particularly susceptible to distrust [37], misuse and disuse [81] by human actors, with the risk of incurring in undesired patterns of user behavior, such as algorithmic aversion [37] and overreliance [88].

At the same time, AI models have found success in several high-stakes domains, such as medicine [117], finance [49] and law [131]. Here, they often act as decision-support systems

Authors' addresses: Clara Punzi, clara.punzi@sns.it, Scuola Normale Superiore, Piazza dei Cavalieri, 7, Pisa, Italy, 56126; Roberto Pellungrini, roberto.pellungrini@sns.it, Scuola Normale Superiore, Piazza dei Cavalieri, 7, Pisa, Italy, 56126; Mattia Setzu, mattia.setzu@unipi.it, University of Pisa, Largo Bruno Pontecorvo, 3, Pisa, Italy, 56127; Fosca Giannotti, fosca.giannotti@sns.it, Scuola Normale Superiore, Piazza dei Cavalieri, 7, Pisa, Italy, 56126; Dino Pedreschi, dino.pedreschi@unipi.it, University of Pisa, Largo Bruno Pontecorvo, 3, Pisa, Italy, 56127.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Association for Computing Machinery.

0004-5411/2024/1-ART26 \$15.00

<https://doi.org/XXXXXXXX.XXXXXXX>

providing human domain experts with advice, rather than fully replacing them. These “collaborative” decision-making systems are plagued by the same trust, misuse, and disuse issues of stand-alone AI systems, on top of novel and unique issues. Collaborative systems have to define an effective communication channel between the AI system and its human decision makers, and the capability to properly integrate the two. Integrating coarse-grained and fine-grained control of AI systems has proven to be a non-trivial task in itself, indicating that steering AI systems to desired complex behaviors [122], preventing them from learning spurious and undesired correlations [71], and ultimately aligning them with human values [128] is far from a solved problem.

Whether AI systems are standalone or collaborative, we find an inherent mismatch between the two parties: humans and AI systems have fundamentally different ways of understanding and representing the world. Therefore, human and AI reasoning are fundamentally different, and leveraging both their complementary strengths could greatly benefit decision-making processes. Due to the heterogeneity of the interacting agents in these systems, some being human, others being artificial, these systems are often referred to as *Hybrid Decision Making Systems* or *Hybrid Systems* (HS) for short. In this survey, we are going to provide an overview of different learning paradigms for hybrid systems, their evolution over time, a taxonomy to properly categorize existing hybrid systems, and their strengths and weaknesses.

Within hybrid systems we identify two archetypal agents of different nature: a *human* agent, which is endowed with peculiar reasoning abilities, domain, and commonsense knowledge; and a *machine* agent, which is instead endowed with high computing power. The two agents¹, which we simply refer to as *human* and *machine*, are best suited to solve different tasks. Unlike humans, machines often fail in tasks requiring (potentially) complex reasoning, domain knowledge [25], commonsense [23], contextualization, and general human touch [41]. On the other hand, machines excel in computation-heavy tasks and are uniquely capable of performing such tasks at scale. It is no surprise that human and machine agents often disagree [133, 72], and even when they do not, they may be agreeing for different, possibly conflicting, reasons [120]. Further exacerbating this, traditional machine evaluation metrics are often blind to this disagreement, since they are designed to simply target the performance of the machine, rather than its internal reasoning.

Whenever an agent attempts to solve a task, they are performing a *computational step*, also known as *computational turn* [63], which is a *human step* if the agent is human, and a *machine step* if the agent is a machine. The definition of the task of interest, of the agents involved, of their computational steps, and of their joint behavior, defines a *Hybrid System* (HS). Generally, we use the term *system* to identify the whole hybrid system, and *machine (sub)system* and *human (sub)system* to refer to the machine and human component of the system, respectively.

The system is specifically designed to solve a given task on which are defined metrics of performance, e.g., accuracy or mean error. The system looks to maximize performance by carefully leveraging its human and machine subsystems. Optimizing this pairing can yield additional benefits by its own design. Successive observations of the behavior of the machine allow humans to best understand the abilities and limitations of a machine [11], allowing the human agent to develop a *mental model* of the machine, that is, an internal representation of how they perceive the machine. If the machine is also able to observe the human in return and gather feedback and corrections from them, then we close the feedback loop and have a hybrid system with bidirectional communication: the human understanding the machine, the machine understanding the human.

A bidirectional system is also at the basis of several uniquely human benefits: humans who perceive to understand a machine through a good mental model are more likely to trust it and

¹We refer to “agents” simply as entities, i.e., humans and AI systems, within a system, rather than “entities with agency”, that is, we do not attribute intentions to agents.

rely on it [10], thus boosting future use of the machine itself; humans that perceive to be able to “steer” a machine, are more likely to do so [38], and thus ease alignment between themselves and the machine. Machines instead, are able to better integrate uniquely human abilities to which they are not preview to, or that are particularly difficult for them to learn [135]. The ultimate goal of a hybrid system is to promote and facilitate human and machine collaboration, thus improving the overall decision quality, reliability, transparency, fairness, etc..

Example of a Generic Hybrid Decision Making System

For clarity, we will use an example to guide us during the survey. An online forum, e.g., a social network, is faced with a large number of messages, some of which may require moderation or further consideration. The platform may be interested in detecting hate speech or individuals at high-risk of self-harm, or harm to others. Given the sheer number of potential messages of interests, systems managing such forums, e.g., through content moderation, are largely automated. For the sake of simplicity, we can think of moderation as a binary task: either flag a message for moderation or not. This deceptively simple task hides several complexities. Content moderation relies on policy guidelines that are frequently updated and are often open to interpretation due to their inherent vagueness. Messages rely heavily on context, which however is not always understandable from the message alone (e.g., when a message mentions an event that occurred outside of the platform). The context is constantly changing, so messages may or may not be considered suitable for moderation at different times. Irony, hyperbole, and metaphors further increase the complexity of such task. AI systems have become increasingly better at tackling this task, but are still far from being reliable enough for a platform to blindly trust them, hence online forums heavily rely on hybrid systems where human moderators collaborate with said AI systems.

In this paper, we focus on hybrid decision making systems with some degree of interaction between human and machines, and as such we do not specifically address systems where humans are not active agents nor they are pure teaching agents within an environment, such as knowledge injection systems, and reinforcement learning systems. Given their strong focus on interaction and little to no focus on integration of the heterogeneous agents into a cohesive system, we do not focus on conversational agents, and interactive XAI algorithms. We refer to works in the literature on each of the aforementioned topics [149, 132, 104].

1.1 Paradigms of Hybrid Systems

Hybrid systems aim at synergically integrate humans and AI systems, and as such can be categorized according to the depth of such integration. This allows us to identify three distinct families or paradigms of hybrid systems with significant differences: *Human overseeing*, *Learn to Abstain*, and *Learn Together* – we visually depict the paradigms in Figure 1. Each of these systems furthers the control that the human can exert on the AI system, and as a consequence also the level of integration. Unsurprisingly, they also roughly follow in chronological order, with systems becoming increasingly integrated as research progresses.

Human overseeing, the simplest of the three, implements a pipeline with no integration between human and machine, the former overseeing the predictions of the latter. Here, the two steps are one after the other. Since the human is aware of the results of the machine step, and is also able to overwrite it, compliance is maximal. Learn to Abstain instead only allows for one of the two steps to take place, either the human step or the machine step. Here, the step of choice is given by an additional machine step that guides the choice. Unlike machine oversight, agents perform

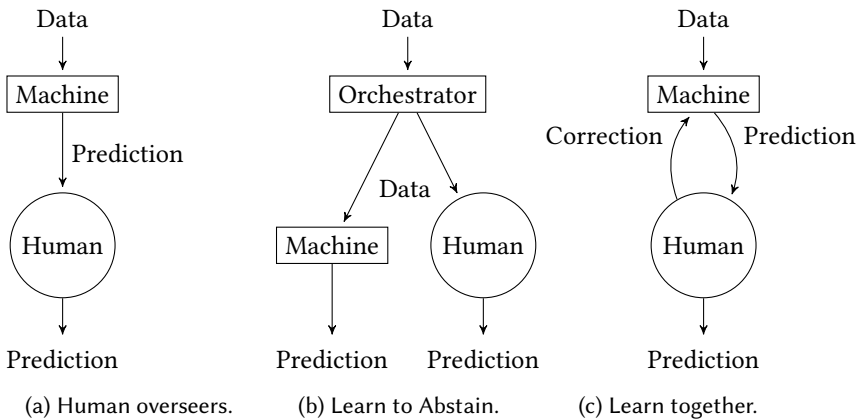


Fig. 1. Paradigms of hybrid systems, where human (circle) and machine (rectangle) steps alternate to form a cohesive system. In human overseers (1a), the machine performs a prediction, and the human accepts it or rejects it in favor of their own. In Learn to Abstain (1b), an orchestrator assigns the prediction task to either of the two, which makes the prediction on their own. In Learn Together (1c), the two agents repeatedly interact by showing and correcting each other's reasoning, finally coming to a prediction.

their step independently of other agents, that is, each prediction is exclusively taken by one of the agents, and an orchestrator decides which agent predicts what. Learning Together systems instead enable a more sophisticated alternation of the two steps, which create a true loop in which human and machine steps are repeated in sequence an arbitrary number of times. Unlike the previous paradigms, the Learning Together paradigm features a communication loop between human and machine agents.

The survey is organized as follows: we provided a general formulation of hybrid systems in Section 2. We then detail the three paradigms, namely, Human overseeing, Learn to Abstain, and Learn Together, in Sections 3, 4, and 5, respectively. We provide an overview of the open problems in Section 6, and finally conclude in Section 7.

1.2 Paper selection criteria

To precisely outline three paradigms, we have explored the scientific literature adapting the paper selection criteria to the significance of the specific paradigm in relation to the scope of this study.

As for the **Paradigm 1**, we have examined papers that have investigated the impact of human oversight on ML models through qualitative or quantitative analysis, but without completing a formal survey. On the contrary, in **Paradigm 2** we conducted a comprehensive review based on the keyword "*learning to reject*", "*learning to defer*", "*learning with a reject option*", "*selective classification*", "*deferral policy*", "*deferral function*" and "*defer to expert*". After a preliminary screening of top conference/journal papers and highly cited papers, where we gathered 150 papers focusing mainly on endowing algorithms with abstention mechanisms, we selected 37 papers representing the wide spectrum of solutions of Learning to Reject, and 53 papers encompassing the most significant findings in Learning to Defer. Finally, in **Paradigm 3**, we examined studies that suggest incorporating human input directly into the learning process of an ML model. The chosen papers represent a promising future path for research on Hybrid Decision Making Systems.

All papers were searched using Google Scholars and DBLP computer science bibliography, selecting papers with high citation counts and/or published in top journals or conferences, such as AAAI, IEEE, ACM, NEURIPS etc..

	Symbol	Definition	Description
Agents	H	\mathcal{H}	Human agent.
	M	\mathcal{M}	Machine agent.
R. Variables	X	$\Omega \rightarrow \mathcal{X}$	Feature matrix. \mathcal{X} includes the empty element \emptyset .
	$Y_{(\cdot)}$	$\Omega \rightarrow \mathcal{Y}_{(\cdot)}$	Label vector: $\mathcal{Y}_{(\cdot)} \subset \mathbb{R}$ (regression), $\mathcal{Y}_{(\cdot)} \subset \mathbb{N}$ (classification) provided by a given agent $(\cdot) \in \{M, H\}$.
	Y^*	$\Omega \rightarrow \mathcal{Y}$	Ground truth label vector.
Machines	f, f_θ	$\mathcal{X} \rightarrow \mathcal{Y}_{(\cdot)}$	Predictor function implemented by a machine agent (\cdot) , belongs to the family of functions \mathcal{F}_M . The parameter $\theta \in \Theta$ is omitted if not relevant.
	f^*	$\mathcal{X} \rightarrow \mathcal{Y}$	Ground truth function. Belongs to the family of functions \mathcal{F}_M
	$\mathcal{L}_{(\cdot)}$	$\mathcal{Y} \times \mathcal{Y}_{(\cdot)} \rightarrow \mathbb{R}_{>0}$	Real-valued loss function of an agent $(\cdot) \in \{M, H\}$.
	$\tilde{\mathcal{L}}_M$	$\mathcal{Y} \times \mathcal{Y}_M \rightarrow \mathbb{R}_{>0}$	Real-valued surrogate loss function of a machine agent M .
	$\rho_{(\cdot)}$	$\mathcal{P}_{(\cdot)} \rightarrow \{0, 1\}$	Deferral or rejection policy of agent $(\cdot) \in \{M, H\}$.
Humans	$Z_{(\cdot)}$	\mathcal{Z}	Set of artifact(s) only available to the human agent (\cdot) .
	$h_{(\cdot)}$	$\mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}_H$	Predictor function implemented by the human agent (\cdot) .
	A	\mathcal{A}	Hybrid artifact, enables communication among agents.
	B	$\mathcal{P}(\mathcal{A})$	Artifact bank, memory storing artifacts.

Table 1. Table of symbols. We provide a definition for each symbol, depending on the context. We will use lowercase letters for elements of a space, e.g., elements x of \mathcal{X} .

2 GENERAL FORMULATION OF HYBRID SYSTEMS

A Hybrid System is composed of two types of agents: a machine agent M and a human agent H . Consider a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and let \mathcal{X} , \mathcal{Z} , \mathcal{Y} , \mathcal{Y}_M , and \mathcal{Y}_H be measurable spaces such that \mathcal{X} represents the machine feature space, \mathcal{Z} the human expertise, that can be modeled as features, decision rules, etc., and \mathcal{Y} , \mathcal{Y}_M , and \mathcal{Y}_H the ground truth, machine and human label spaces for a certain task T , respectively. Full list of symbols can be found in Table 1.

Machine model. Let $X : \Omega \rightarrow \mathcal{X}$ and $Y_M : \Omega \rightarrow \mathcal{Y}_M$ be random variables representing the input and output of a machine M , and let $Y^* : \Omega \rightarrow \mathcal{Y}$ be the random variable representing the true labels. In the classical setting of supervised learning, given independent and identically distributed pairs $\{(X_i, Y_i^*)\}_{i=1}^n \stackrel{\text{iid}}{\sim} (X, Y^*)$ drawn from the same unknown joint distribution over $\mathcal{X} \times \mathcal{Y}^*$, we aim to learn a predictor $f : \mathcal{X} \rightarrow \mathcal{Y}_M$ that approximates, as accurately as possible, the unknown function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ representing the true relationship between the features and the target. Note that f is chosen from a space of hypothesis \mathcal{F}_M parameterized by $\theta \in \Theta$, however, in our work we often omit the parameter θ to simplify our notation. The ultimate goal of the learning algorithm is to find an hypothesis $f \in \mathcal{F}_M$ that minimizes the *empirical expected risk*, which is defined on the training set as follows:

$$\widehat{\mathcal{R}}_n[f] = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_M(Y_{M,i}, f(X_i)) \quad (1)$$

The function $\mathcal{L}_M : \mathcal{Y} \times \mathcal{Y}_M \rightarrow \mathbb{R}_{>0}$ represents the machine loss, which quantifies how far the predictions of a hypothesis f are from the true outcome. In many cases, the formulation of the loss may have a level of complexity that renders direct computation infeasible due to the curse of dimensionality or potential model mis-specification [103]; equation (1) is often not continuous nor differentiable, hence extremely hard to optimize and computationally intractable for many nontrivial classes of functions [50]. The approach typically employed to overcome this issue is to

define and optimize a *surrogate loss function* $\tilde{\mathcal{L}}_M : \mathcal{Y} \times \mathcal{Y}_M \rightarrow \mathbb{R}_{>0}$, that is, a function with good computational guarantees (e.g., differentiability and convexity) that can be easily optimized and whose optimal values approximate well the minimizer of the original computationally hard loss function [12, 86]. The exact formulation of a surrogate loss function is not straightforward, as it depends on the particular task at hand and the desired properties one seeks to guarantee. Refer to Appendix A for an extensive discussion about the fundamental mathematical properties that characterize surrogate losses.

Human model. Let $Z : \Omega \rightarrow \mathcal{Z}$ and $Y_H : \Omega \rightarrow \mathcal{Y}_H$ be random variables representing, respectively, the human expertise and the predictions of a human agent H for a certain task T , where H is modeled as a predictor $h : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}_H$. Generally, we consider $|\mathcal{X} \cap \mathcal{Z}| \geq 0$, implying that there may exist shared information between the machine and the human, or there may not.

It is worth noting that there may be a divergence between the predictions made by human agents and the ground truth labels. Indeed, different levels of background knowledge, experience, or personal biases can lead to distinct decision-making outcomes, resulting in both correct and incorrect predictions across various domains within the input space [78]. Therefore, we additionally take into account a loss function $\mathcal{L}_H : \mathcal{Y} \times \mathcal{Y}_H \rightarrow \mathbb{R}_{>0}$ as an indicator of the quality of human predictions. The observation that various types of errors are not only made by distinct human agents but also occur between humans and AI models provides support for the transition towards paradigms that incorporate hybrid combinations of human and machine predictions [73].

3 PEERING INTO THE MACHINE: HUMAN OVERSEERS

Human oversight [76] is probably the simplest and most straightforward form of hybrid system. In this first paradigm, machine and human agents are independent of each other, the former performing a task, and the latter *verifying* its predictions. Informally, the human agents perform a straightforward task: given the machine computation and/or the input data, either accept or reject the computation. More formally, a *human oversight policy* ρ_H is a binary function that, given a machine M implementing a function $f \in \mathcal{F}_M$, an overseeing human H with additional expertise $Z \in \mathcal{Z}$, and some input data $x \in \mathcal{X}$, either *accepts* or *rejects* the prediction given by M . In its most general formulation, ρ_H is formulated as follows:

$$\rho_H : \mathcal{X} \times \mathcal{Y}_M \times \mathcal{Z} \times \mathcal{F}_M \rightarrow \{\text{accept, reject}\}.$$

In many cases, the human agent either does not have access or does not take into consideration the machine itself, thus the above formulation is often reduced to:

$$\rho_H : \mathcal{X} \times \mathcal{Y}_M \times \mathcal{Z} \rightarrow \{\text{accept, reject}\}.$$

Example of online content moderation: Human Oversight

In the paradigm of Human Oversight, the human moderator is presented with a message X , the prediction Y_M given by the machine, and has to decide whether to agree ($\rho_H(X, Y_M, Z) = 1$) or not ($\rho_H(X, Y_M, Z) = 0$) with the prediction of the machine Y_M , e.g., whether the message is hate speech, or is indicative of an individual at risk of self-harm.

Other than simply leveraging their own expertise Z , overseers often leverage external factors in their decision. In the simplest of cases, rejected patterns of misbehavior are limited to **machine-specific** failures in which the underlying context is of little or no impact. In these context-independent scenarios, the overseers aim to identify machine failures induced by machine-specific causes by tracking a set of *subjects of monitoring* around which the overseeing policy will

be centered. Machine-specific failures may be induced by wildly different factors in each machine, hence we abstract over the underlying causes, since the goal of machine oversight is to *identify*, rather than *diagnose*, undesired behavior. We tackle behavior correction later in Section 5. Common subjects of monitoring are:

- *Data shift*. Data shift is generally intended as a change in the data distribution [113], and may be due to a change in feature distribution, i.e., *covariate shift* or knowledge drift, or label distribution, i.e., *prior shift*; these shifts may occur frequently in real-world scenarios, namely due to data seasonality, sampling bias, or naturally-occurring distribution changes. Unlike spontaneous *outlier* instances, dataset shifts are responsible for consistent and predictable failures in the model, thus they have to be accounted for.
- *(Partial) Model performance*. Partially stemming from data shift, model performance is another primary subject of interest. Here, we distinguish between two classes of performance metrics: “global metrics”, where the model is evaluated *wholly*, and “partial metrics”, where the model is evaluated *partially* on a suitable subset of data. Data shift may indeed only affect a subset of data, hence global metrics may easily deceive the overseers.
- *Model uncertainty*. Classical learning algorithms are trained to predict, setting *confidence calibration* aside. Simply put, they operate under a *closed world assumption* where the only available option is to give a prediction, that is, the model has no notion of uncertainty nor of the unknown. This projects a false high confidence, seldom tricking the algorithm user into overestimating its competency [20].
- *Decision complexity*. Even with highly sophisticated and precise models, some decisions are inherently suitable for human rather than algorithmic reasoning [107]. In this context, it is crucial to anticipate which decision is which.

Operating as self-contained agents, machines often lack the decision context wherein their predictions are evaluated and applied. Failures in this domain are said to be **context-dependent**. Context can be of primary importance, and humans are far better suited than machines in understanding it and integrating it in their decision-making process. For a human, concerns such as fairness, legality, and explainability of the decision are strong contextual motivations that a machine does not necessarily take into proper account. Unsurprisingly, most of them are already being encoded in several legislatures, which strongly discourage or punish discriminative or otherwise illegal [7], unexplainable [94] decisions and behaviors. Much of this stems from the current use of machine agents in ethically-charged contexts. For instance, machines are leveraged in monitoring and discouraging [7] illegal activities, where they often yield unfair or biased predictions [41]; furthermore, they are a critical component of speech regulation, an extremely dynamic use case where human scrutiny and decision autonomy are essential, yet they often regulate marginally- or fully-free [48] speech; they aid hiring in public and private companies, yet they are biased [114].

What’s more, context is often dynamic and loosely defined [102], and thus integrating it into the machine is an open challenge in and of itself. Jointly, machine-specific and context-specific failures offer a strong motivation for machine oversight. Yet, even though the *why* is clear, *how* machine oversight is to be implemented is still an open problem. Even worse, machine oversight poses a set of inherently human problems to face.

3.1 Monitoring pitfalls

While technical solutions to machine-specific failures have already been developed, context-dependent failures cause a plethora of additional and more complex problems. Given that humans are usually the time bottleneck when it comes to decision-making, one cannot let the overseers monitor every prediction, thus one needs to understand *when* to let the overseer monitor the

machine. A conventional solution, which we explore in Section 4, is to let the machine itself call the human into action. Here, we focus instead on the overseers and their inherently human fallacies, which lead to some natural pitfalls of the whole monitoring process.

Human Executors and Skeptics. Overseers need to be aware of two possible cognitive biases: algorithmic aversion [37] and overreliance [88]. Algorithmic aversion pushes the overseers towards excessively doubting the machine, thus introducing unnecessary monitoring in the decision-making process. Algorithmic aversion manifest itself independently of the performance of the machine [37, 92], and more strongly when the machine fails. In other words, every single perceived mistake of the machine compounds in increasing the rejection rate of the overseer.

Unlike algorithm aversion, algorithm overreliance occurs when human agents under-monitor a machine system, and thus act as mere executors. The two biases are well-documented in the literature, and regardless of the machine system they human agents interface with, they are to be accounted for.

Biased monitoring. Automation bias is particularly strong when human agents oversee fairness-related tasks where the task directly involves other humans. When monitoring decision on pre-held stereotypes, say on vulnerable groups, overseers either avoid monitoring in the first place [2] or further confirm the stereotypes [8, 41]. On similar reasoning, particularly in cases of fairness evaluation, human agents tend to under-monitor when they perceive affinity towards the monitor case at hand [53], or simply when they deem their reasoning more “human” than the machine’s [82]. On an even more biological level, intrinsic demographic traits are also likely to play a role in the decision-making of the human agents [144, 108]. To further increase the complexity of setting up a set of overseers, it is often the case that human agents, in part for the aforementioned reasons, have a low level of agreement on the correct task solution [52].

Failure to oversee and trust calibration. Even more worrying than biased monitoring is the failure to reject obvious machine failures. In a pilot study with legal experts, [40] showed that, when assisted by a faulty machine agent, domain experts incorporate into their decision-making machine recommendations based on irrelevant or random factors, with extreme cases in which the domain experts were knowingly introducing random factors themselves – a clear case of placebo effect where the mere presence of a machine prediction, regardless of its correctness, induces an almost blind trust in the human agent. Unsurprisingly, overseers have repeatedly shown to be unable to properly assess their ability to assess the performance of a given machine [130, 1].

3.2 Enhanced monitoring: the case for Explainable AI

Overseeing a machine simply through its predictions and uncertainty provides minimal tools to a human agent, which can easily fall into one or more of the aforementioned pitfalls. To enhance their overseeing power, humans are often accompanied by *explanation algorithms*, that is, algorithms able to further explain the predictions given by a machine. Explainable AI (XAI) [55] is a recent field of research aiming to shed light on the prediction process of AI models by extracting human-understandable explanations. Explanations allow an overseer to peer into the machine and get a grasp of what features a machine is relying upon [89] and what rule-like logic it is following [54] to make its predictions, what training instances have had a particular influence on the learning process [75], and how one could change the input instance to achieve a different prediction [146]. Explanations have shown to empower the human agent into better understanding the machine, thus improving their ability to monitor it.

4 LEARNING TO ABSTAIN

To mitigate human failures in the monitoring process of the aforementioned HS paradigm, a potential approach involves developing an enhanced machine architecture that enables the machine learning model to refrain from predicting on certain instances. In broad terms, whenever the machine prediction can't be relied upon, it may be better to pay an extra costs and defer the prediction to a human agent [31]. The ultimate goal of a Learning to Abstain system is thus to maximize the overall performance of the hybrid system by effectively identifying, for each single instance, the agent best suited for a given prediction. This particular setting does not involve any formally defined interaction between human and machine agents. Instead, it can be better characterized as a policy in which a machine autonomously directs input data towards an agent, whether it be a human or an AI, based on an estimation of superior performance in each individual case. In other words, the algorithms in this paradigm produce models optimized to leverage the option to abstain w.r.t. their primary predictive task. The human interaction in this paradigm may come afterwards, e.g., operating on those instances rejected by the Learning to Abstain system.

At the most general level, the machine learning systems that fall within this framework can be categorized into **Learning to Reject** (L2R) [24, 31] and **Learning to Defer** (L2D) [90, 99] systems, depending on whether the model assumes a pre-specified cost for abstaining or instead is designed to work adaptively with the human agent. Both classes of algorithms share the same set-up, namely an HS where the machine (or a system of machines) should learn not only a standard classification function but also a *rejection function* (also called *deferral function* in the L2D domain) under the optimization objective of maximizing the performance of the human-AI system as a whole. The deferral action incurs an additional cost which is offset by the payoff of an expected gain in performance resulting from querying a human expert, who can use information that the machine cannot rely on (e.g., professional experience and common sense). Although L2R and L2D have been tackled as separate and independent problems, a seminal paper by Madras et al. [90] demonstrated that L2R can be regarded as the specific case of L2D in which a fixed cost is allocated to each deferred instance. Consequently, they proposed a revised definition of L2D as an *adaptive* L2R approach. Therefore, we devote most of our discussion to L2D (Section 4.2) while giving a more general overview of L2R (Section 4.1).

4.1 Learning to Reject

Learning to Reject (L2R) was first introduced by Chow [24], and its general formulation constitutes a base for more advanced learning to defer algorithms (Section 4.2.3).

Example of online content moderation: Learning to Reject

In the L2R paradigm, the algorithm is trained to consider a third alternative with respect to the original binary classification task, namely the rejection option. Specifically, the model decides to abstain from making a prediction in those cases where it is more likely to make a wrong prediction. For example, a model could be trained to recognize those messages that contain sarcasm or humor, typical cases where a classification model for moderation is likely to make a wrong prediction, and abstain on those messages. In this way, a person responsible for monitoring the system can focus solely on the instances that are rejected by the model, rather than having to deal with all of them.

In the literature, this area of research is referred to with several names: *learning to reject* [150], *selective classification* [148], or *machine learning with a reject option* [62]. The general idea is to

provide machine learning algorithms with the means to learn when to abstain from actually making a certain prediction. Here, the focus is on optimizing machine learning algorithms to include a reject option, thus improving the final performance by refraining from predicting on critical instances. Therefore, humans are not an active agent in this framework, as the focus is on finding the best combination of machine prediction and machine abstention. However, in the context of Hybrid Systems, one could envision real life scenarios where L2R models can be used to focus the attention of a human overseer on the subset of instances that are rejected by the model. Learning to Reject is a widely studied problem with a long history of research, already touched upon in several previous surveys. Therefore, we report a similar general definition to the work of Hendrickx et al. [62] and Zhang et al. [150]. The goal of L2R algorithms is to learn a model f_{ρ_M} composed of two parts: a prediction function $f : \mathcal{X} \rightarrow \mathcal{Y}_M$ and a rejection policy $\rho_M : \mathcal{X} \rightarrow \{0, 1\}$. The composed system can be defined as $f_{\rho_M} : \mathcal{X} \rightarrow \mathcal{Y}_M \cup \{\emptyset\}$ such that:

$$f_{\rho_M}(x) = \begin{cases} \emptyset & \text{if } \rho_M(x) = 1 \\ f(x) & \text{otherwise} \end{cases} \quad (2)$$

That is, if the rejection policy ρ_M rejects x , then no prediction is made. If instead ρ_M accepts x , then the prediction function f is applied to x and the result $f(x)$ is obtained. The function f is generally assumed to be a classifier. Ideally, ρ_M should be able to prevent misprediction of f while conversely accepting examples for which a good prediction is more likely. The policy ρ_M that actually performs the rejection operation is called *rejection function* or *rejector*. The focus of a rejector, as mentioned, is to refuse those examples for which the classifier f is expected to output an incorrect prediction. More generally, the focus of research in this field is to find the optimal balance between accuracy on accepted instances and number of rejected instances. We can generally divide the types of rejection into two macro-categories, depending on the type of data that are being rejected:

- *Novelties rejection*, i.e., rejections of examples that differ significantly from the data points in the training set X of f and therefore are likely to cause a misprediction.
- *Ambiguities rejection*, i.e., rejections of examples in proximity of the decision boundary of f .

Novelties are examples that are, in general, very different from the ones that f was trained on. Novelties are usually data that may be considered out of distribution (OOD), or simply particularly unusual instances. Ambiguities are typically examples that can either be classified as any by f as their probability of being a member of one is almost equal to the probability of belonging to the others [57]. Depending on the kind of mistake that needs to be addressed, different methods for learning the rejector can be used. The rejector can be based only on feature observation or on model evaluation. Therefore, ρ_M may have access to only the feature space \mathcal{X} or to both \mathcal{X} and the output of the predictor f , or even to the architecture and parametrization of f . In essence, a rejector can either be:

- *Independent*, i.e., the rejection function is learned regardless of the predictor, by observing solely the feature space.
- *Dependent*, i.e., the rejection function is built by either querying the predictor or by relying on particular characteristics of the predictor.

4.1.1 Independent rejectors. Independent rejectors are mostly designed to perform *Novelties rejections*. Most of the works in this particular area frame the problem as *Outlier detection* [64] or *Anomaly detection* [136]. Another widely-used term is *Open-set recognition* [91], where the focus is on optimizing the accuracy of a predictor both on in-distribution and out-of-distribution data. The basic idea is to recognize and handle those examples that present out-of-distribution characteristics. This is done independently from the predictor function f , i.e., the rejection mechanism can be applied directly to the data in an agnostic way w.r.t. the trained model f .

Some of the earliest work in this area leverage statistical methods to find anomalies in the data. Seo et al. use the estimated posterior variance of Gaussian processes as a testing point to understand whether to reject a data point or not [125]. Another similar technique was proposed by Coles [30] which relies on statistical models studied in Extreme Value Theory to understand which examples may actually be considered outside expected distributions. This technique has been employed in many different tasks, e.g. facial recognition [124] and network analysis [95]. Recently, [121] proposed the Extreme Value Machine, a model for kernel-free variable bandwidth incremental learning based on Extreme Value Theory. One limitation of this model is that it relies on the distances between the instances belonging to different classes in the training set, so that classes that are unreasonably distant from the majority classes in the training may be rejected in their entirety. Some works address this limitation by relying on generalized Pareto distribution approximation to determine whether a new example is a normal or abnormal data point [141].

Other approaches rely on models trained to recognize the training distribution and perform anomaly detection as a separate task. One such method was proposed by Coenen et al. and it relies on One-Class Support Vector Machines [27] specifically fitted to recognize the training distribution. Alternatively, Gaussian Mixture Models can be employed to estimate the training distribution, and thus detect novelties [79]. Recently [4] proposed a few-shot learning approach where margin-loss w.r.t. the training class is used to train a model to detect anomalies.

While independent rejectors are mainly designed to detect anomalies, Asif and Amir Afsar Minhas [6] propose a framework for generalized rejection based on jointly-trained Neural Networks. These networks are trained with a dual-penalization both on misprediction and incorrect rejection. Since the rejection network is trained independently, it can then be used with any other prediction model, and it can also perform ambiguity rejections.

4.1.2 Dependent rejectors. Dependent rejectors are much more studied than independent ones. They are tied to the outputs or properties of the predictor and as such can be learned either after the learning phase of f or simultaneously. Adopting the same naming convention that we will use for Learning to Defer models (Section 4.2) we call the former *staged rejectors* and the latter *joint rejectors*.

Staged rejectors. look at the predictor f output to estimate the *confidence* or *uncertainty* of the model. The works that fall in this category are usually referred to as *confidence based* or *uncertainty based*. These approaches estimate uncertainty by *i)* formulating a good confidence metric in and *ii)* selecting an appropriate threshold for determining the actual rejection [32, 51, 69]. Formally, the rejection of an instance can be seen as:

$$\rho_M(x, f) = \mathbb{1}_{c(x, f) > \tau}$$

where $c(x, f)$ is the confidence metric, which depends both on the instance x and the prediction function f , and τ is the threshold hyperparameter [58, 148]. Many different metrics for confidence have been proposed in the literature. Confidence metrics can be based on:

- *Hard predictions* where the exact class-output of predictor f , i.e., $y_M = f(x)$, is used to assess the confidence of the prediction. These works usually observe multiple predictions and consider the class-wise variance of these predictions as indication of uncertainty [14, 129].
- *Soft predictions* where some scoring output of f is used as an estimation of $P(y_M|x)$ to determine the confidence. This approximation of confidence is more widely studied, and the array of solutions that leverage soft predictions is much varied [36, 18].

There are also other methods that, while still exploiting some output of the predictor, have slightly different mechanisms to approximate the confidence. For example, in the work of Tortorella [137]

the authors exploit directly the score of SVMs as a distance from the decision boundary. In other works, the distance from the k -th nearest prototype is used as a proxy of confidence [17].

The confidence threshold τ can either be global, or local: a global threshold is usually used for those predictors with an equally well calibrated confidence metric over the entire feature space [80, 21, 42]. Multiple local thresholds τ_1, \dots, τ_n [43] are instead best suited for predictors with variable accuracy [109] or for different class-wise variance for the confidence metric.

Joint rejectors. Also sometimes called *integrated rejectors* [62] are directly part of the prediction model, i.e., the rejection option is treated as an additional class to be learned by the predictor. In this setting, predictor f and rejector ρ_M are simultaneously learned with a single algorithm, specifically designed to learn both functions. This means that it is difficult to technically distinguish and separate the predictor from the rejector [31]. Some joint-learning approaches rely on the definition of a specifically formulated objective function able to penalize both incorrect predictions and rejections. Several works resort to surrogate loss functions in place of a more natural discrete loss in order to more easily solve the problem via minimization [13, 119]. Other works aim at specifically leverage some evaluation metric, for example Pugnana and Ruggieri [112] design a model-agnostic approach for probabilistic binary classifiers with a reject option specifically aimed at optimizing Area Under Curve. Other works directly allocate an extra class for rejection to any given predictor and assign a specific penalization cost for predicting such class [152]. Finally, several joint-learning algorithms for rejection have been developed specifically for certain kinds of machine learning models, for example for Support Vector Machines [51, 85] and for Neural Networks [47].

4.1.3 Cost model for rejection. Approaches for L2R need to achieve a balance between predictive performance and rejection rate. In fact, classifier with a reject option can entirely trade performance for coverage of vice-versa. This means that, in theory, to achieve maximum performance, a predictor with a reject option can opt to reject all instances and therefore never actually make a misprediction. Otherwise, the predictor can opt to never reject in order to achieve maximum coverage over the data, thus completely forgoing the benefits of implementing a reject option in the first place. To avoid this, an appropriate cost model must be implemented. One foundational cost model can be found in the work of Cortes et al.[31]:

$$\text{cost}(x) = \begin{cases} 0, & \text{if } \rho_M(x) = 0 \wedge f(x) = y^* \\ 1, & \text{if } \rho_M(x) = 0 \wedge f(x) \neq y^* \\ \mathcal{R} & \text{if } \rho_M(x) = 1. \end{cases} \quad (3)$$

Where y^* is the ground truth for example x , and \mathcal{R} is a fixed, predefined cost for rejection. If we adopt such a cost model, we can express the the learning objective of a L2R algorithm as follows:

$$\min_{f, \rho_M} \sum_{x \in X} [\mathbb{1}_{\rho_M(x)=0} \mathbb{1}_{f(x) \neq y^*} + \mathbb{1}_{\rho_M(x)=1} \mathcal{R}] \quad (4)$$

In essence, we either pay a cost if we mispredict ($\mathbb{1}_{f(x) \neq y^*}$) some instance x that we accepted ($\rho_M(x) = 0$) or we pay a cost \mathcal{R} if we reject that instance ($\rho_M(x) \mathcal{R}$). Although sensible, this cost model has the major drawback of having to determine the rejection cost \mathcal{R} beforehand. Rejection cost can be determined depending on the application domain, but this is not always an option [47].

The work of Geifman and El-Yaniv [46] is among the most important in this field, as it tries to tackle the aforementioned problem. In their proposal, the authors propose a method called Selection with Guaranteed Risk Control that looks at the problem of rejection cost determination from a different angle, that is looking at the coverage of the model. Coverage of a model is defined as the portion of the data that the model accepts for prediction: $\frac{1}{n} \sum_{x \in X} \mathbb{1}_{\rho_M(x)=0}$.

Following Eq. (4) the learning objective can be expressed as a minimization of risk subject to a coverage constraint that forces the model to predict at least a certain portion of instances. Formally:

$$\min_{f, \rho_M} \frac{\sum_{x \in X} \mathbb{1}_{\rho_M(x)=0} \mathbb{1}_{f(x) \neq y^*}}{\sum_{x \in X} (1 - \rho_M(x))} \quad \text{s.t.} \quad \frac{1}{n} \sum_{x \in X} \mathbb{1}_{\rho_M(x)=0} > C$$

where $0 < C < 1$ is a coverage threshold. The problem has thus shifted from having to determine cost \mathcal{R} to having to select threshold C which is easier. It is important to note that if $C = 1$ the model reverts to a standard classifier with no rejection option. This coverage based formulation has been proven to be theoretically equivalent to the original cost model of Eq. (3) [44].

4.1.4 Strengths and limitations of Learning to Reject. In summary, the L2R paradigm allows for the development of ML models with a reject option, which is a foundational starting point for developing models able to interact with humans. Indeed, ML models equipped with the reject option can, in principle, reject exactly those instances that would yield a prediction error, and therefore call for human intervention only when strictly needed. However, the actual benefit of these techniques in a collaborative setting with humans has never been thoroughly investigated. Indeed, if human intervention is called only for those instances for which a decision is difficult, the human expert may find the same difficulties, and thus deem the model as not so useful for solving the task. Moreover, there are studies pointing to possible fairness issues when using classifiers with a reject option [67].

4.2 Learning to Defer

Learning to Defer (L2D) systems embed human knowledge directly into the training process of a ML model. The goal is to equip the model with the ability to call for human intervention on those instances where the human is likely to give accurate prediction and the machine is likely to fail.

Example of online content moderation: Learning to Defer

By comparing predictions made by human workers to the correct labels, an L2D system can be trained to determine which instances can be accurately predicted by AI and which are better handled by humans. For instance, in both cases of hate speech and self-harm content recognition, humans are expected to outperform machines on texts characterized by toxic and satirical content, mostly owing to their greater capacity for comprehending common sense and contextual information, as well as discerning implicit meanings hidden within words and phrases.

In contrast to L2D, an L2R system learns its rejector policy only from the feature set \mathcal{X} (e.g., the text of the message to be flagged) and, possibly, some properties of the predictor used by the AI system. L2D instead actively considers the human expertise in the task domain.

4.2.1 General formulation. Keeping the same notation introduced in Section 2, we consider Hybrid Systems composed of a machine M and a human H . Similarly to L2R, in L2D the machine M is equipped with the possibility of abstaining from making a prediction. In addition, an L2D model also embeds a representation of the human agent H , thereby taking into account their estimated performance when assessing the act of deferral. By doing so, the human expertise \mathcal{Z} is taken into account. Nevertheless, note that the predictor $h : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}_H$ modeling H is *fixed*, meaning that L2D algorithms have no control nor visibility on the function h itself; rather, they only have access to its image, that is, the set $\{Y_{H,i}\}_{i=1}^n := \{h(x_i, z_i)\}_{i=1}^n$ of human predictions about the training data.

Architecture.	<i>Design according to which classifier and deferral policy are integrated.</i>	
Staged learning	Classifier and deferral policy are learnt in two subsequent steps.	[9, 116, 115, 147]
Joint learning	Classifier and deferral policy are learnt jointly.	[22, 45, 59, 74, 87, 90, 99, 98, 111, 118, 140, 139, 147]
Others	Alternative solutions, e.g., iterative models.	[34, 35, 105]
Multiplicity.	<i>Number of human agents in the hybrid system.</i>	
Single	One human agent.	[9, 22, 87, 90, 99, 98, 105, 111, 116, 115, 118, 140, 147]
Multiple (1 predicts)	One human is selected out of many.	[45, 59]
Multiple (j predict)	A subset of agents is selected out of many.	[74, 139]
Theoretical guarantees.	<i>The machine phase in which the interaction occurs.</i>	
Fisher consistency	The optimization objective has the correct target.	[22, 99, 140, 139]
Classification-calibration	Agents are given realistic uncertainty estimates.	[140, 139]
Realizable consistency	The problem is well-defined under specific choices of the classifier and deferral hypothesis function spaces.	[98]
Constraints.	<i>Additional conditions that the hybrid system should satisfy.</i>	
Coverage	Number of instances that can be deferred.	[35, 99, 98, 105]
Budget	Total cost to query human agents.	[115]
Fairness	Metrics to guarantee algorithmic fairness.	[74, 90, 98]
Others	Others, e.g., on the selection of human agents.	[74]

Table 2. Properties of systems in the Learning to Defer paradigm.

In analogy to Eq. (2), a L2D system can be formulated as a function $f_{\rho_M} : \mathcal{X} \rightarrow \mathcal{Y}_M \cup \mathcal{Y}_H$ defined from the classifier $f : \mathcal{X} \rightarrow \mathcal{Y}_M$ and deferral policy $\rho_M : \mathcal{X} \rightarrow \{0, 1\}$:

$$f_{\rho_M}(x) = \begin{cases} h(x) & \text{if } \rho_M(x) = 1 \\ f(x) & \text{otherwise} \end{cases}$$

The goal of L2D is to find the classifier-rejector pair $(\hat{f}, \hat{\rho}_M)$ that minimizes the system loss $\mathcal{L}_{\text{defer}}$. This can be expressed as the summation of the machine loss \mathcal{L}_M and the human loss \mathcal{L}_H :

$$\mathcal{L}_{\text{defer}}(Y^*, Y_M, Y_H, \rho_M) := \underbrace{\mathbb{1}_{\rho_M(X)=0}}_{M \text{ predicts}} \underbrace{\mathcal{L}_M(Y^*, Y_M)}_{\text{machine cost}} + \underbrace{\mathbb{1}_{\rho_M(X)=1}}_{\text{defer to } H} \underbrace{\mathcal{L}_H(Y^*, Y_H)}_{\text{human cost}} \quad (5)$$

where $\mathbb{1}$ denotes the indicator function. Note that the optimization problem has Y_M and ρ_M as unique learnable parameters since both Y_H and Y^* are fixed (i.e., they belong to the training data). In general, the individual losses \mathcal{L}_M and \mathcal{L}_H can take several forms to account for different "costs", such as the misprediction error as in the 0-1 loss, or the cost of querying the human agent.

Notably, when there exists a constant $\mathcal{R} > 0$ such that $\mathcal{L}_H(y^*, y_H) = \mathcal{R}$ for all $(y^*, y_H) \in (Y^*, Y_H)$, then the loss (5) matches the rejection loss described in [31] for the L2R framework [90]:

$$\mathcal{L}_{\text{reject}}(Y^*, Y_M, \rho_M) := \mathbb{1}_{\rho_M(X)=0} \mathcal{L}_M(Y^*, Y_M) + \mathcal{R} \mathbb{1}_{\rho_M(X)=1} \quad (6)$$

Note that Eq. (6) and Eq. (4) are equivalent under the assumption of using the 0/1 cost model for prediction/rejection.

Optimization constraints. Depending on the specific context of use, the application of specific constraints may be necessary for hybrid systems. This objective is commonly accomplished by incorporating regularization terms into the system loss or by imposing specific bounding conditions. Examples of such constraints include:

- Coverage or triage level [35, 99, 98, 105]: the number of instances that can be deferred.
- Fairness metrics [74, 90, 98]: for instance, the Minimax Pareto Fairness criterion [96] or the equalized odds metric with respect to a protected attribute.
- Budget [115]: total cost that can be allocated to query human agents.

Model architectures. L2D systems typically adhere to either of two general designs, referred to as *staged learning* and *joint learning*, which vary in terms of when the classifier and rejector are learned. In the former case, the algorithmic process starts by learning the classifier and only subsequently fits the deferral policy on top of it. On the other hand, in a joint learning setting the classifier and rejector are learnt simultaneously through the direct minimization of the system loss (5). While most of the proposals documented in the literature can be categorized as staged or joint learning models, a few exceptions also exist that do not fit in either category (Section 4.2.4).

Number of human agents. L2D systems can be characterized along another dimension, that is, the size of the pool of human agents to which the decision can be deferred to. Whenever the number of these is greater than one, term *Multiple-Expert L2D* (L2D-ME) is used, as opposed to *Single-Expert L2D* (L2D-SE or L2D), which considers one human only. An example where L2D-ME modeling may be more suitable is in a medical setting, where a critical decision regarding a complex case could be made either by an automated classifier or by one or more doctors chosen from a team of experts with potentially diverse expertise and opinions. As compared to L2D-SE, the formalization of L2D-ME makes the deferral function more complex in nature. Specifically, it should not only determine *when* to defer, but also to *which* human agent(s) [139]. Furthermore, the deferral policy can be designed in a manner that allocates predictions to either *one* single agent or to a *subset* of agents from the available pool.

4.2.2 Staged learning architectures. In L2D models characterized by staged learning architectures, the classifier f and deferral function ρ_M are learnt separately. Specifically, algorithms in this category first fit a classifier on the training dataset, then they learn a second model that predicts the probability that the human makes a mistake on the same dataset, and finally they defer based on which has the lowest error probability instance-wise.

For instance, Raghu et al. [115] developed a basic heuristic for L2D consisting of two independent models trained on the full dataset: a multiclass classifier representing the machine agent, and a binary classifier representing the correctness of the human agent. At inference time, an instance is deferred to the human if the predicted classifier error probability is higher than that of the human. In case of coverage constraints, then the samples whose difference between human and classifier error probability is higher are chosen first. Interestingly, the authors also suggest a reduction of L2D-ME to L2D-SE by modeling the human subsystem in terms of average disagreement between human agents on each single prediction. This approach has been further developed in [116].

Another common baseline for staged learning is the model proposed by Bansal et al. [9], who described a staged learning setting aimed to maximize the expected utility of the system, which is measured in terms of the accuracy of the final decision, the cost of deferring, and the individual accuracy of both the human and machine component. Differently from other L2D models, this method has been claimed to be user-initiated, since the action of deferral is triggered through an (over-simplified) threshold-based policy that represents the humans' mental model of the AI.

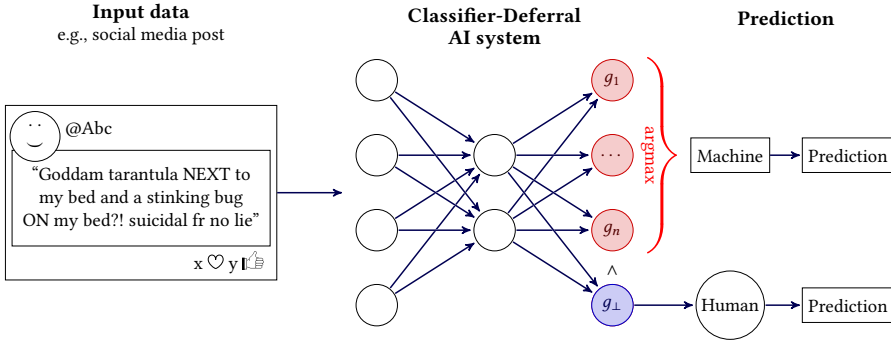


Fig. 2. Overview of the joint learning architecture for the Single-Expert Learning to Defer (L2D-SE) setting, illustrated in the application of flagging online contents for moderation. Adapted from [98].

Finally, a third relevant staged learning method known as *fixed value of information approach* has been proposed in [147]. It consists in training independently three probabilistic models describing, respectively, the distribution of the label given the input data, the human predictions given the input data, and the label given both the input data and human predictions. At inference, the deferral policy evaluates the estimated expected utility of the classifier in two scenarios: when the human is not consulted and when the human is queried, while also taking into account the distribution of human predictions and a constant cost for querying the human.

As noted by Charusaie et al. [22], the staged learning approach presents some important advantages: first of all, it is suited for convenient implementation, since already known appropriate algorithms can be adopted to solve the two stages separately. Secondly, theoretical and experimental results suggest that it outperforms the joint learning approach in realistic scenarios where only a limited portion of data is labeled by the human agent. In these cases, the classifier f can still be optimized over the full dataset; on the other hand, in joint learning f can be learnt from the subset of human labeled data only, thus leading to a reduction in performance dependent on the proportion of unlabeled data. However, [22] also pointed out that staged learning is sub-optimal with respect to joint learning and provide both theoretical and experimental results showing the existence of a performance gap between the two approaches in terms of model complexity.

4.2.3 Joint learning architectures. In L2D systems characterized by a joint learning architecture the classifier f and deferral function ρ_M are learnt simultaneously. In order to implement this design, the task is shaped as a $K + 1$ multiclass problem over an augmented label space $\mathcal{Y}^\theta := \mathcal{Y}_M \cup \{\theta\}$, where $\mathcal{Y}_M = \{1, \dots, K\}$. In particular, we set $\mathbf{g} = (g_1, \dots, g_K, g_\theta)$ to be the set of real-valued scoring functions $g_i : \mathcal{X} \rightarrow \mathbb{R}$, such that g_θ returns the human predictions Y_H , while the classifier and rejector are defined, respectively, as:

$$f(x) = \arg \max_{i \in \mathcal{Y}_M} g_i(x) \quad \rho_M(x) = \begin{cases} 1 & \text{if } \max_{i \in \mathcal{Y}_M} g_i(x) \leq g_\theta(x) \\ 0 & \text{otherwise.} \end{cases}$$

In general, the optimal classifier-rejector pair is found by minimizing the system loss expressed in Eq. (5). However, the function (5) is often computationally hard to optimize. Such a problem is addressed by replacing (5) with a *surrogate loss* $\tilde{\mathcal{L}}_{\text{defer}}$ that is easy to optimize and is chosen to guarantee specific properties with respect to the original loss $\mathcal{L}_{\text{defer}}$ (refer to Appendix A for a formal discussion of desirable properties of surrogate loss functions in the context of L2D). Hence,

in case of appropriate choices of the human and machine surrogate loss functions, the joint learning HS gives theoretical guarantees for optimal performance.

Single-Expert L2D. Most of the literature on joint learning models for L2D-SE focuses on the *cost-sensitive formulation* of the problem over an augmented label space \mathcal{Y}^0 developed by Mozannar and Sontag [99] and illustrated in Figure 2. This setting considers random costs $\mathbf{c} = \{c_1, \dots, c_{K+1}\} \in \mathbb{R}_{>0}^{K+1}$ where each component c_i represents the cost of predicting the label $i \in \mathcal{Y}^0$. In this section, we review and categorize the most relevant proposals according to their statistical properties.

(1) *Fisher Consistency (FC)* has been described as a minimal requirement that surrogate loss functions should satisfy to achieve reasonable performance, since it posits that, if an estimator were computed using the complete population instead of a sample, it would yield the true value of the estimated parameter [86]. To the best of our knowledge, FC approximation of the 0-1 loss in the L2D setting have been implemented (up to adaptations) only by Mozannar and Sontag [99], Charusaie et al. [22], and Verma and Nalisnick [140]. In particular, the surrogate loss \mathcal{L}_{CE}^α [99] consist of a generalization of the cross-entropy loss with the costs corresponding to multiclass misclassification, where the cost of the $K + 1$ class represents the action of deferral, and $\alpha \in \mathbb{R}_{>0}$ is a weighting parameter that modulates deferral. When $\alpha = 1$, \mathcal{L}_{CE}^α has FC and can be expressed as:

$$\mathcal{L}_{CE}^1(\mathbf{g}; x, y^*, y_H) := -\log \left(\frac{\exp(g_{y^*}(x))}{\sum_{y \in \mathcal{Y}^0} \exp(g_y(x))} \right) - \mathbb{1}_{y_H=y^*} \log \left(\frac{\exp(g_\emptyset(x))}{\sum_{y \in \mathcal{Y}^0} \exp(g_y(x))} \right)$$

intuitively, the first term maximizes the scoring function associated with the true label, while the second maximizes the rejection (scoring) function but only if the human's prediction is correct. Notably, [22] describes a family of cost-sensitive surrogate loss functions for the 0-1 loss that generalizes prior work and encompasses \mathcal{L}_{CE} as well.

A few adaptations have been proposed to enhance the L2D algorithms based on the surrogate \mathcal{L}_{CE} , with the aim of better capturing specific properties. These include:

- *Learning to Defer with Uncertainty (LDU)*, where the deferral policy accounts for the epistemic uncertainty of the model (i.e., the uncertainty resulting from limited data availability and lack of knowledge about the system of interest) [87];
- The customization of the model to suit the expertise of particular human agent, which however requires the availability of supplementary data that has been annotated by that particular human [118].

(2) *Confidence calibration* refers to the property of an estimator (e.g., a probabilistic classifier) to produce a predictive distribution that is consistent with the empirical frequencies observed from realized outcomes [33]. Verma and Nalisnick [140] proposed a surrogate loss \mathcal{L}_{OvA} [140] that satisfies both Fisher consistency and classification-calibration. This solution consists in solving the L2D problem via a One-vs-All classification method which breaks down the original $K + 1$ classes into $K + 1$ binary classifier models. The resulting objective function to be optimized is thus a surrogate loss function composed of different logistic loss components, each accounting for the error on one of the $K + 1$ different classes.

$$\mathcal{L}_{OvA}(\mathbf{g}; x, y^*, y_H) := \phi[g_{y^*}(x)] + \sum_{\substack{y \in \mathcal{Y}_M \\ y \neq y^*}} \phi[-g_y(x)] + \phi[-g_\emptyset(x)] + \mathbb{1}_{y_H=y^*} (\phi[g_\emptyset(x)] - \phi[-g_\emptyset(x)])$$

where $\phi : \{0, 1\} \times \mathbb{R} \rightarrow \mathbb{R}_{>0}$ is a binary surrogate loss (e.g., the logistic loss). Experimental findings show that \mathcal{L}_{OvA} results in better calibrated models w.r.t. ones trained with \mathcal{L}_{CE} , with competitive performance w.r.t. other L2D baselines [9, 99, 105, 115].

- (3) (\mathcal{F}_M, P_M) -Realizable consistency is a property that refines the notion of FC by addressing the optimization process over restricted hypothesis classes \mathcal{F}_M and P_M for the predictor and deferral function. For instance, the surrogate \mathcal{L}_{RS} [98] is differentiable, non-convex, and realizable (\mathcal{F}_M, P_M) -consistent for classes \mathcal{F}_M and P_M closed under scaling:

$$\mathcal{L}_{RS}(\mathbf{g}; x, y^*, y_H) := -2 \log \left(\frac{\exp(g_{y^*}(x)) + \mathbb{1}_{y_H=y^*} \exp(g_{\theta}(x))}{\sum_{y \in \mathcal{Y}^0} \exp(g_y(x))} \right)$$

Other formulations of joint learning L2D-SE also exist. This is the case of:

- the seminal paper by Madras et al. [90], which presents a framework for addressing the L2D problem using a Mixture-of-Experts (MoE) approach, with the deferral policy acting as a gating function. The classifier and deferral function are learned together by negative log-likelihood minimization over the augmented label space \mathcal{Y}^0 . An alternative version of the algorithm is also introduced, wherein a regularization component is added to the system loss to account for a fairness. Unfortunately, this method was proven to not have FC [99].
- *Preferential MoE* [111], a variant of [90] where human knowledge is encoded in the form of decision rules that should be followed as much as possible, that is, whenever they are applicable and do not decrease the system performance. The algorithm first checks the applicability of the available rules and, in case of positive response, a deferral function selects whether to rely on the human or machine prediction based on their performance. Notably, the deferral function is chosen to be interpretable (e.g., a linear classifier or decision tree) guaranteeing transparency in the selection of the human agent, and also highlighting reasons for forgoing the human-based rules.
- In the *joint value of information method* [147], the three probabilistic models already introduced in the *fixed value of information method* described in the Section 4.2.2 are trained together through a single neural network which includes a final Platt calibration layer that guarantees the estimation of meaningful expected utilities. Experimental findings show that joint learning yields greater advantages compared to analogous staged learning method.
- The *Mixed Integer Linear Program* (MILP) [98] is a scheme to exactly minimize the misclassification error of the HS. It comes with generalization bounds and allows to provably and easily integrate any linear constraints on the variables. However, it suffers from two limitations: it is computationally expensive and it does not generalize to non-linear predictors.

Multiple-Expert L2D (1 out of J). In this first scenario of L2D-ME, the goal of the multi-expert deferral policy ρ_M^{ME} is to choose either the classifier or *exactly one* human agent from the set of J available ones. Hence, ρ_M^{ME} takes the form of $\rho_M^{ME} : \mathcal{X} \rightarrow \{0, 1, \dots, J\}$, where $\rho_M^{ME}(x) = 0$ means that the classifier decides, while $\rho_M^{ME}(x) = j$ for $j \neq 0$ indicates that the decision is deferred to the j^{th} human agent.

Hemmer et al. [59] adopt a mixture of experts (MoE) approach with the deferral policy serving as a gating function that assigns each instance either to the predictor or one specific human agent. The joint learning of the classifier and deferral function is carried out through a surrogate loss function based on the negative log-likelihood of the system. However, subsequent work [139] has proven this surrogate to be not FC and proposed instead two surrogate loss functions, namely one based on cross-entropy and one on the One-vs-All classification, which are consistent with the 0-1 loss in the L2D-ME setting and extend their single-expert analogue. The experimental findings indicate that the OvA-trained model frequently achieves superior performance compared to both the cross-entropy variant and the MoE baseline [59]. Additionally, it exhibits better calibration in terms of the correctness of agents' decisions.

Finally, Gao et al. [45] study the problem of L2D-ME in a *bandit feedback* setting (i.e., a sequential dynamic allocation problem). Specifically, the deferral policy is specifically learned using a supervised learning model that has been previously trained on historical data that reflect human decisions and corresponding outcomes in order to maximize the complementarity of the machine and human agents. However, by doing so, it is assumed that the human agents who generated the historical data are the same individuals who will be assigned decisions at inference time.

Multiple-Expert L2D (j out of J). In the second scenario of L2D-ME, the multi-expert deferral policy is defined as $\rho_M^{ME} : \mathcal{X} \rightarrow \{0, 1\}^{J+1}$. For each input $x \in \mathcal{X}$, the goal is to choose the committee of agents $C(x) \subseteq \{0, \dots, J\}$, possibly including the classifier, who are likely to make the most accurate decision for x . Hence, the i^{th} vector component of the deferral policy will be defined as $\rho_M^{ME}(x)_{(i)} = 1$ for all $i \in C(x)$, and $\rho_M^{ME}(x)_{(i)} = 0$ for all $i \notin C(x)$. In the event that the designated committee comprises multiple agents, the resulting outcome will be an aggregated decision. This setting has been firstly addressed by Keswani et al. [74], who proposed a joint loss functions obtained by linearly combining the losses associated to the classifier and deferral function via context-dependent hyperparameters. The authors proved that the combined loss is convex with respect to the classifier and deferral function whenever the loss associated to the former is convex; under such assumption, it can be optimized using the projected-gradient descent algorithm. Additionally, the authors outlined a few adaptations of the L2D-ME framework to account for potential real-world constraints and requirements:

- **Fair learning:** this variant takes into account the possibility of performance discrepancies that may occur with respect to individuals belonging to different protected categories.
- **Sparse Committee Selection:** this variant enables the deferral function to exclusively choose a limited number of agents on a per-instance basis.
- **Dropout:** this variant aims to reduce the dependence on a single agent and achieve a more equitable distribution of workload.
- **Regularized versions:** additional constraints can be added to the joint framework as regularizers of the loss function. For instance, this solution can be employed in cases where specific costs associated with individual human agent consultations are provided.

Subsequent work [74] further developed this setting to adapt to *closed* deferral pipeline, wherein the human agents of the HS also provided the training labels. This is achieved through an online framework in which input samples are received in a continuous stream. After each prediction is made, which involves aggregating the outputs of agents in the chosen committee, the sample are utilized to retrain the classifier and deferral function.

Alternatively, Verma et al. [139] suggest to use Conformal Inference [127] to find ensembles of agents $C(x)$ that include the best agent with high marginal probability. The size of $C(x)$ is computed dynamically as a function of the input x , thereby ensuring optimal utilization of agent queries. The authors propose two test statistics for the estimation of $C(x)$: a naive score function that sums up the correctness scores of all agents who correctly predict the given instance, and a regularized statistics that employ conformal risk control [5] to increase the robustness to noise. The experimental findings demonstrate that the latter approach yields a nearly flawless identification of the appropriate number of agents. Moreover, the conformal approach exhibits superior performance in system accuracy compared to a fixed-size ensemble of agents.

4.2.4 Further model architectures. While most of the proposals documented in the literature can be categorized as staged or joint learning models, a few exceptions also exists. A notable example often used as a baseline in the L2D literature is the method proposed by Okati et al. [105], namely, an iterative algorithm that optimizes the classifier and triage policy alternately. At each iteration,

the optimization process is carried out for the classifier on instances where it outperforms the human agents, while the remaining data points are optimized for the triage policy. The authors show that their method converges to a local minimum. Nevertheless, subsequent experimental studies have shown that this method exhibits lower performance in comparison to other L2D algorithms [98]. Additionally, similar algorithms have been implemented to address the issue of L2D for model-specific settings, namely Support Vector Machines [35] and Ridge Regression [34].

4.2.5 L2D with limited human predictions. A significant drawback of L2D is the requirement of human predictions, alongside ground truth labels, for every instance within the training set [83]. Ideally, the L2D system has to be trained on human labels belonging to the same human that will then interact with the system itself. By doing so, the L2D system will learn to complement that specific human [60]. Due to the significant computational and human costs, it is likely that the implementation of the L2D algorithm would be impractical for most real-world scenarios.

A couple proposals have been put forward to address this problem and implement L2D-SE algorithms with only a limited amount of human predictions. Charusaie et al. [22] derived an active learning scheme known as *Disagreement on Disagreements* (DoD) which can learn a classifier-rejector pair by making a minimal number of queries to the human. The DoD algorithm functions in two steps: i) first, a standard active learning algorithm (i.e., CAL [28]) is executed on the hypothesis space to learn the human disagreement with the ground truth: at each round, the disagreement set of the predictors is computed, and then the human is queried on the instances in this set to learn their error boundary; ii) second, a consistent classifier-rejector pair is learnt from the pseudo-labeled data derived from the preceding stage.

Alternatively, Hemmer et al. [60] proposed a three-step methodology that, based on a limited set of human predictions, generates synthetic human labels that simulate their capabilities. The first step consist in training an *embedding model* that maps each instance into feature representations. These, along with the ground truth labels, are used in the second step to learn an *expertise predictor model*, which is designed to approximate the capabilities of a human agent based on their labeling behavior. To harness the potential of both labeled and unlabeled instances, this component incorporates a semi-supervised learning technique. Finally, the expertise predictor model produces synthetic human predictions for instances that have not been labeled by a human and that can thus be used as input of any downstream L2D algorithm. The results of the empirical analysis shows that a small number of human predictions per class is enough to effectively generate synthetic predictions.

4.2.6 Strength and limitations of Learning to Defer. In contrast to algorithms that operate under oversight (see Section 3), Learning to Abstain Hybrid Systems are trained not to predict when their performance is weak. As a result, when using an L2R or L2D algorithm to make decisions, one can expect to receive two kinds of evidence: the machine's prediction concerning the action of deferral, and, if the AI does not abstain, the result of the prediction task. L2D algorithms improve upon L2R by incorporating a representation of human knowledge directly in the training process. In such a way, the deferral policy is trained to adapt to both the AI model and the human decision-maker, ideally the same that will employ the Hybrid System.

Recent empirical investigations involving human subjects yielded evidence for the additional advantages that abstaining systems bring to Hybrid Systems. Hemmer et al. [61] found that such algorithms improve both *human task performance* compared to a human or an AI working alone, and *human task satisfaction* compared to a human working alone. A different study [106] also investigated the effects of employing abstaining Hybrid Systems on the human perception of AI performance and credibility. The results indicate that users are frequently influenced by the system's recommendation also on ambiguous instances, even without conscious awareness, and thus support the adoption of L2D algorithms.

However, L2D also comes with several limitations [83]. Most importantly, these include: data availability issues, which are primarily due to the need of human predictions in addition to the ground truth for all instances within the training set and all human agents involved in the Hybrid System; and fairness concerns that may stem from the introduction of bias by both human and machine agents, as well as from the abstention mechanism itself [68]. Although there have been suggestions to deal with such issues, these proposals still do not offer a straightforward solutions.

In terms of the human's role in the Hybrid System, Learning to Abstain enhances only marginally the paradigm of human oversight over machines (Section 3), as the deferral is exclusively a machine-side operation and there is no direct human-side interaction considered in the design of the algorithms.

5 LEARNING TOGETHER: HUMANS TEACHING MACHINES, AND MACHINES TEACHING HUMANS

The next natural step in hybrid systems is a two-way collaboration in which human agents are not mere executors or overseers, but can directly *interact* with the machine to best infuse their human decision-making abilities directly into the machine. While Learning to Defer aims to identify which agent is best suited for a given prediction, there is little to no effort in integrating the decision-making abilities of one agent to the other. In such a setting, the prediction capability of the overall system is at best but a sum of the prediction abilities of its agents, that is, the system is not synergic. Usually, in machine learning machines learn directly *from data*, while in a Learning Together system they aim to learn *from other agents*. To reach the ultimate goal of learning from other human agents, we need to build on top of a foundational stepping stone that enables machines to learn from other machines, the *Teacher-Student* paradigm.

5.1 Machines Learning from Other Machines: The Teacher-Student Model

At the basis of machine-to-machine learning systems is the *Teacher-Student* paradigm. In this approach, we identify two agents: a *Teacher*, whose goal is to train, and a *Student*, whose goal is to learn from the Teacher. Among Teacher-Student models, four are of particular note: Learning with Privilege [138], Knowledge Distillation [19], Transfer Learning [151], and Active Learning [126].

Learning with Privilege. Initially introduced by Vapnik and Vashist [138], Learning with Privilege is defined as a learning paradigm in which a teaching agent f_T , i.e., a Teacher, provides the learning agent f_S , i.e., the Student, with additional “privileged” information about the input data which is not present in the data itself. The additional data is considered privileged because only available at training time. In our formulation, this is roughly equivalent to integrating a proxy \bar{Z} of the human knowledge Z into the training data. At learning time, a privileged estimator implements a function

$$\mathcal{F}_M : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}_M,$$

which directly integrates the human knowledge into the model. Integration relies on a machine-specific encoding function that allows to map the input data $X \in \mathcal{X}$ and the human knowledge $Z \in \mathcal{Z}$ into the input space \mathcal{X} . At inference time, the privileged information is unavailable, and the encoding function reduces to its projection $(x, \{\emptyset\}) \mapsto x$, thus allowing the machine to seamlessly operate either with or without privileged knowledge.

Knowledge distillation. Knowledge distillation aims to distill a given Teacher into a more suitable Student, usually to reduce the running cost of inference [19]. Unlike Privileged learning, here the privileged information is exclusively derived from a machine, usually involving its internal state, which the Student tries to emulate. Formally, at training time the Student is directly conditioned

on the parameters of the Teacher, thus yielding an estimator f_θ of the family:

$$\mathcal{F}_{M,S} : \mathcal{X} \times \mathcal{F}_{M,T} \rightarrow \mathcal{Y}_M,$$

where we assume without loss of generality that the Teacher and Student are members of the same hypothesis space \mathcal{F}_M . The Student loss $\widetilde{\mathcal{L}}_S$ is thus augmented with an additional component $\widetilde{\mathcal{L}}_{T,S}(f_T, f_S)$ penalizing a distance between Teacher and Student:

$$\widetilde{\mathcal{L}}_S(X, Y) = \frac{1}{n} \sum_{i=1}^n \widetilde{\mathcal{L}}(Y_i, f_S(X_i)) + \lambda_S \widetilde{\mathcal{L}}_{T,S}(f_T, f_S),$$

for a predefined weight λ_S balancing task performance and distance from the Teacher. In neural models, where Knowledge Distillation is prevalent, the Teacher-Student loss $\widetilde{\mathcal{L}}_{T,S}$ is usually implemented as $d(\theta_T, \theta_S)$, e.g., Euclidean or Cosine distance, between the parametrization θ_T of the Teacher and the parametrization θ_S the Student. In both Learning with Privilege and Knowledge Distillation settings, the final goal is to transfer information between Teacher and Student.

Transfer Learning. In this the step where the Students tries to leverage the knowledge of the Teacher, to then adapt it to the specific task at hand [151]. This paradigm is typical of dynamic settings in which tasks or data distributions change rapidly, yet they are all strongly related. In this context, learning from agents which have already learned to solve related task is of great benefit.

Active Learning. The aforementioned approaches all provide a straightforward and direct way to inject learning capabilities from the Teacher to the Student. However, the latter remains a passive agent in this process, resulting in a unidirectional connection between the two. Furthermore, knowledge is injected one-shot at training time, that is, these paradigms are not designed for subsequent Teacher-Student interactions. Active learning [126] builds on this paradigm by empowering the Student to actively seek the most relevant information. Like a student that poses questions to a teacher, an Active Student *queries* the Teacher (in our formulation, the human agent) to maximize its own performance. Formally, the Student has at its disposal a pool of unlabeled data $X_U \in \mathbb{P}(\mathcal{X})$ that can, in principle, be labeled with the predictions of the Teacher. Querying the Teacher is feasible but costly, thus one wishes to maximize Student performance while staying under a given budget. To tackle this problem, we aim to learn a *query policy*

$$\pi : \mathcal{F}_M \times \mathbb{P}(\mathcal{X}) \mapsto \mathbb{P}(\mathcal{X})$$

that, given the Student f_S and the unlabeled data X_U , selects a subset $\hat{X}_U \subset X_U$ for which the Teacher then provides labels $\hat{Y}_U = \langle f_T(x_U) \rangle_{x_U \in \hat{X}_U}$.

Training the Student involves an iterative process of querying and training, and in each iteration the Student seeks to best select the subset of instances $\hat{X}_U \subset X_U$ that will maximize its performance. X_U can vary in its nature: it can be static or dynamic [126], it can be updated with a handful [100] or a large set of novel unlabeled samples [29], and the Student may even populate X_U itself. Initial formulations of Active learning [29] strongly resembles Learn to Reject solutions, as the goal of both is to identify regions in the input space \mathcal{X} where the machine is less accurate.

Teacher-Student Model: Limitations. The Teacher-Student paradigm is limited to one-way learning between agents, which makes it extremely flexible, as it does not have to model the interaction between agents, nor the agents themselves. This introduces several layers of complexity, each compounding on the previous, making for novel approaches to Learning Together, which takes inspiration from the aforementioned paradigms.

Communication language.	<i>The communication language that allows human and machine to interact.</i>	
Hard reasoning	First-Order logic.	[56, 39, 39]
Soft reasoning	Logic-like languages with soft reasoning engines.	[97, 16, 101, 143]
Explanations	Feature relevance, decision rules, and concepts.	[134, 135, 120, 145, 110, 77, 15]
Artifact type.	<i>The type of interaction artifact.</i>	
Intrinsic	Artifact internal to the machine, the human agent cannot directly influence the machine.	[134, 135, 120, 145, 110, 77, 15]
Extrinsic	Artifact external to the machine, the human agent can directly influence the machine.	[56, 39, 97, 101]
Time of interaction.	<i>The machine phase in which the interaction occurs.</i>	
Training	The human agent influences the machine directly in its training process.	[134, 135, 120, 145, 110, 77, 15]
Inference	The human agent influences the machine at inference time.	[56, 39, 39, 97, 16, 101, 143]
Learning cost.	<i>The cost incurred by the machine to integrate the feedback by the human agent.</i>	
Training	The correction prompts an additional training phase of the machine.	[134, 135, 120, 145, 110, 77, 15]
Null	The correction does not cause any significant cost to the machine.	[56, 39, 39, 97, 16, 101, 143]

Table 3. Properties of systems in the Learning Together paradigm.

5.2 Learning Together

Unlike the aforementioned Teacher-Student paradigm, the *Learning Together* paradigm aims to create a two-way street in which human and machines can communicate effectively, one learning from the other: the machine explaining its prediction to the human, and the human explaining its prediction to the machine, improving the overall performance of the system. In this context, improvement can take many forms. For the machine, it can provide more accurate predictions, a better ability to generalize, a shallower learning curve, or higher level of transparency. On the other hand, human agents can derive additional advantages from exerting a tighter grip on *how* the machine solves the task, rather than simply solving the task in place of the machine.

Example of online content moderation: Learning Together

In Learning Together systems, users may interact with a feature-importance based explanation about messages. This would allow them to adjust the importance of different terms based on their own judgement. After receiving corrections from humans, the machine could integrate them to enhance the learning process and better flag future posts.

Generally, a machine implementing this paradigm partially integrates a human agent with expertise Z , thus yielding an estimator of the form: $f_{\theta}(X, Z)$. For the sake of simplicity, here we model a single human agent, but the formulation can be extended to multiple agents.

We characterize systems in the Learning Together paradigm according to four properties that we summarize in Table 3: communication language, which guides the communication between

agents; interaction artifacts, which guide the feedback mechanism; time of interaction, which defines *when* the interaction occurs; and learning cost, which defines what *cost* the interaction creates.

5.2.1 Communication Languages and Interaction Artifacts. The communication language plays a critical role in facilitating successful high-level human-machine interactions. On one hand, its level must be sufficiently low to allow the machine to understand and interface with it; on the other hand, its level should be high enough for the human agent to understand it. Communication languages are comprised of two intertwined components: a *language*, which defines how the agents communicate, and a set of *interaction artifacts*, which comprise the atomic unit of interaction, and enable the feedback of one agent to be integrated into the other.

Artifacts. At the core of a communication language lies a set A of *interaction artifacts* that enable effective communication between humans and machines, and enable the embedding of human component Z into the machine. That is, the human encodes their knowledge by *acting* on the interaction artifacts $A \in \mathcal{A}$ presented by the machine. Then, the machine integrates the human feedback by integrating the action of the human agent on A . Typical artifacts include logic rules, relevant features, and decision rules. The human action can also take several forms, and is inherently dependent on the type of the interaction artifact.

Regardless of their nature and formulation, artifacts are:

- *understandable* by both the human and the machine, since they must enable communication;
- *malleable* by the human agent, as they transfer their expertise Z to the machine by *acting* on the artifact;
- *embeddable* into the machine, since the human action on the artifact is to be integrated into the machine.

The last property puts a strong constraint on machines, which, unlike artifacts, tend to employ subsymbolic, e.g., neural, rather than symbolic models. This often results in systems where the family and architecture of the machine are all but determined by the choice of artifacts, thus communication language and machine end up being tightly coupled.

Interaction artifacts come in one of two forms: *intrinsic* and *extrinsic*, yielding intrinsic and extrinsic machines, respectively. Intrinsic artifacts are machine-oriented, thus after the human acts upon them, the integration step is hidden to the human agent, who has no direct control of how the machine integrates the action. This often allows easier definition of the interaction artifact set at the cost of the effectiveness of the action. The result of an interaction between an intrinsic machine parametrized by θ and a human agent acting on artifact a is in an update of the model parameters θ , which yields an estimator of the form $f_{\theta|a}(X)$, where θ is conditioned on the artifact a acting as proxy feedback for the human agent. Intrinsic artifacts are most often employed for systems using *explanations* as communication language.

Extrinsic artifacts instead, are human-oriented, thus once the human acts upon them, the integration is straightforward for the machine, which does not require additional integration steps. Typically, such artifacts are organized in an artifact bank B that the user can inspect and act upon. With an artifact bank to attend to, the interactive machine M implements an estimator of the form $f_{\theta}(X, B)$. Note that the whole artifact bank is a parameter of the machine, hence the parameters θ of the machine are independent of the artifact bank, Extrinsic artifacts are typically more complex than intrinsic ones, and so are the communication languages built upon them.

Communication languages. Communication languages can be broadly categorized in two families: *reasoning* languages, both *hard* and *soft*, and *explanation* languages.

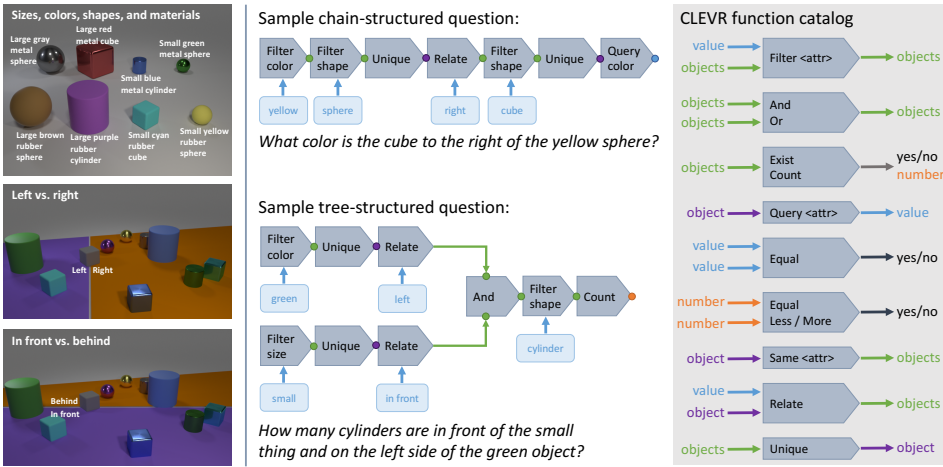


Fig. 3. An example of Question Answering machine employing a hard reasoning language. The agent maps the question to a program by using a set of primitives (right-hand side of the picture), and then executes the program on the input (left-hand side of the picture), providing the human agent with both a prediction, and a malleable program that the user can correct.

Hard reasoning languages: Logic. Turning to languages, logic natively offers languages and artifacts that satisfy all three requirements. Logic programs are naturally readable by humans due to their transparency and similarity to human reasoning, and thus are a suitable candidate for communication language. They are also naturally malleable but it is not straightforward to embed them in the machine due to their symbolic nature, which is in stark contrast with the subsymbolic nature of the vast majority of machines. Learning Together systems using logic as a communication language are strongly coupled and often neuro-symbolic in nature, that is, they combine subsymbolic and symbolic components. The subsymbolic component is typically a neural one, e.g., a neural network, which fully embodies machine reasoning. On the other hand, the symbolic component is geared towards the human, and thus employs a logic language. Neuro-symbolic systems aim to combine and integrate these two components.

Generally, logic and reasoning-like languages define a set of clauses or rules \mathcal{R} , and facts \mathcal{T} that allow a logic engine to reason and present its logical derivations to the human agent as extrinsic artifacts on which they can act. Possible actions on logic artifacts, i.e., logic rules, facts, and compositions of the two, include *addition* and *deletion* of new/existing rules or facts. Hard logic enjoys strong theoretical properties, which makes it suitable for highly compliant systems where providing feedback to the machine is highly likely to result in a successful integration in the machine. Logic as a communication language is expert-driven, rather than data-driven. This means that artifacts in these languages are often hand-crafted, domain-specific, and automating their definition usually involves a significant effort on the human side. Figure 3 provides an example extracted from the CLEVR dataset [66]. Here, the task is to answer questions about the provided image, e.g., “What is the color of the cube to the right of the yellow sphere?” The interaction artifact consists of a program solving the task by first identifying the yellow sphere, then finding objects to its right and finally filtering them by shape. The intermediate result of this program is then used to extract the color of such object. When presented with it, the human agent can inspect and appropriately modify the program to solve the task. We can find another classical example of hard

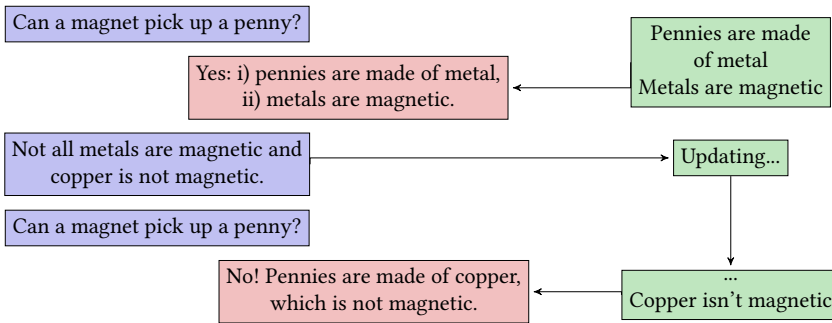


Fig. 4. An example of a Question Answering hybrid system employing a soft reasoning language. Unlike hard reasoning languages, the **human** agent can interact with the **machine** through Natural Language. Here, the user is first provided with a prediction and a rationale by the machine, and is able to correct the rationale, feed it back to the machine, which updates its **artifact bank** (denoted as “Feedback memory” in [97]). On a subsequent interaction, the machine is presented with the same question, but this time it provides the correct prediction and rationale by leveraging the feedback provided by the human agent in the previous step.

reasoning in [56], where the machine is tasked to predict the final outcome of a Tic-Tac-Toe board by leveraging a set of logic rules, which the human can inspect and change at their leisure.

Neurosymbolic models also offer another related subfield, *neural program synthesis*. Like inductive logic programming, neural program synthesis tries to solve tasks by learning interpretable programs, only in this case the program is not expressed in a logic language, but in an actual programming language. DreamCoder [39] is a neurosymbolic model to solve tasks by writing computer programs in a given minimal computer language. Notably, DreamCoder progressively grows a library of functions, small reusable snippets of computer code re-used throughout the program. A human agent can then inspect and change functions in the library that they deem incorrect.

Soft reasoning languages. Soft reasoning communication languages improve the flexibility of Logic languages at the cost of their strong theoretical properties by replacing logic rules and facts with *logic-like* rules, and the symbolic reasoning engine with a subsymbolic approximate one. Like in hard reasoning languages, possible actions on artifacts include *addition* and *deletion* of new/existing rules or facts. These languages are almost exclusively Natural Languages, and are themselves applied to Natural Language models.

Originating from the work by Kassner et al., this family of machines builds on top of two other families, knowledge injection models [84] and soft reasoners [26], the former providing models able to integrate external possibly human-sourced knowledge in their inference, and the latter providing soft reasoning engines for Natural Language. Due to their strong emphasis on flexibility, soft reasoning languages are usually loosely defined in terms of a set \mathcal{T} of facts and informal rules \mathcal{R} . On inference, the machine presents the human agent with a subset $T \subset \mathcal{T}$ of facts, optionally accompanied by a set of derivation rules $R \subset \mathcal{R}$ it leveraged. According to the complexity of the task, facts may be enriched with intermediate reasoning steps that the machine has derived or generated to solve the task.

Figure 4 shows an example on a Question Answering task where the hybrid system is tasked to answer basic physics questions by the user [97]. The machine is asked whether “A magnet can pick up a penny”, and generated both an answer, “Yes”, and an artifact supporting it in the form of syllogistic reasoning. The human agent notices a flaw in the premises, corrects the artifact by stating that there are non-magnetic materials, one of which is copper. Here, the integration step is



Question: What can the red object on the ground be used for?

Answer: Firefighting.

Support fact: Fire hydrants can be used to fight fires.

Fig. 5. An example of a Question Answering hybrid system employing a soft reasoning language and leveraging Knowledge Graphs. Here, the machine consults its artifact bank to retrieve a supporting fact for its prediction, and provides it to the human agent alongside its prediction.

minimal, and consists in adding the correction by the human agent to the artifacts bank. Then, on a successive iteration, the machine is presented with the same question, and successfully answers by leveraging the claims introduced by the user, without incurring any additional training cost. Other works operate on similar terms by providing full reasoning trees [97], improving on the artifacts by artificial generation [16], or by providing ad-hoc artifacts to correct the machine [101]. The human agent can act upon them by correcting any number of facts and/or rules used in inference, yielding updated facts \mathcal{T}' and \mathcal{R}' .

Another clear advantage of soft reasoning languages is the ability to integrate artifacts from multiple sources to aid reasoning. Knowledge Graphs are a primary source for several reasons: they are easily available, enjoy widespread use, can be mined from Natural Language, cover a wide range of domains, and are immediately understandable to human agents, thus providing high understandability and malleability of the artifacts. Knowledge Graphs can also be easily turned into Natural Language through a process of *verbalization* in which the Knowledge Graph is encoded into a set of Natural Language claims. Improving on Natural Language, they are usually verified by a wide pool of users, hence they make for reliable sources which require minimal to no verification. This last property is particularly of interest because it allows the system to scale to more human agents than the ones who are directly designing the system. To further strengthen the case for Knowledge Graphs, their integration into machines has been a well-studied problem for several years, and already presents several effective solutions [84]. Knowledge Graphs also allow some basic form of reasoning which can be easily extended: they provide a strict reasoning, both commonsense [65] and factual [142], all of which can be integrated in Natural Language. Figure 5 presents an example from [143]. Here, the goal is to understand what the purpose of the red object (the hydrant) in the photo is. We ought to remark the importance of the context, since even though the machine is aware of the concept of “fire hydrant”, it is highly unlikely to have seen either a fire hydrant or a fireman in a forest. Moreover, since the image itself does not provide any information regarding the purpose of the hydrant, one should expect that the machine could not solve the task without the additional information provided by the human agent in the Knowledge Graph.

Explanations. Explanations are designed to *explain* the machine, which makes them highly understandable artifacts out of the box. Common families of explanations include counterfactuals, prototypes, feature relevance, and decision rules, the last two leveraged in Learning Together systems. Feature relevance provides the human agent with the estimated influence that each input feature has on the prediction of the machine: the higher the relevance, the higher the sensitivity of the machine to changes in that feature. Decision rules, instead, provide a descriptive understanding by articulating with logic rules the predictions of the machine. It is worth to mention

that explanations are often extracted *post-hoc* and in a *model-agnostic* fashion, i.e., *after* the machine has been trained and regardless of its form. Jointly, these two characteristics minimize coupling between the interaction artifacts and the machine, thus granting more flexibility in the design of the machine. Some more recent proposals propose conversational agents that are able to provide an interface for the human agent to best query the machine for explanations [93]. On this account, explanations are major artifacts of interest, particularly for systems in which the interaction is the focal point, such as Interactive Machine Learning [3] and eXplainable Interactive Learning (XIL) [135] systems. In particular, the latter allow the human agent to inspect and provide feedback to the machine by correcting its explanation. A XIL machine is based on a simple algorithmic kernel comprised of five steps:

- (1) the machine performs a *learning* step by optimizing its parametrization θ ;
- (2) the machine generates explanations A_M of (a subset of) its predictions;
- (3) the human examines A_M , and provides a (optionally) corrected explanation A_H ;
- (4) the machine performs a learning step by optimizing θ according to the corrected explanation;
- (5) if no stopping criterion is met, the machine returns to step (1).

By iteratively querying the human, the machine parametrization is thus conditioned on the corrected explanations A_H , that is, a XIL machine ends up implementing an estimator of the form: $f_{\theta|A_H}(X)$, where the correction A_H acts as a proxy for the human knowledge Z .

Integration is directly dependent on the family of explanations. By far the most widespread one is feature importance, which assigns a relevance score to each element of the input data X . The interaction then consists in a possible correction of the relevance scores, with the human agent activating or deactivating each feature according to their judgement. Integration follows either a *learning* approach, in which the correction is directly encoded in the machine training objective [120], or a *generative* one, where the correction is implemented via training on additional synthetic data [134, 135]. In a learning approach, the corrections of the human agent are encoded in a correction matrix $C \in \{0, 1\}^{n \times m}$ that states what feature relevance have been corrected for each single instance. In a generative approach, C is instead used to generate synthetic data \tilde{X} to further train the machine. Features with low relevance have randomized or copied values, while features with high relevance are kept as-is [135]. A small subset of machines are designed to explicitly encode feature relevance in their architecture, thus allowing direct manipulation by human agents [145].

The same two approaches are also found in two other explanation types, namely Decision Rules and Concepts. In the former case, we have a generative approach in which the human agent is tasked with directly creating the additional synthetic data themselves [110]. In the latter, we have a learning approach in which the explanation is itself a component of the architecture, thus the correction is again an added component of the machine objective [77].

A third emerging approach, mainly aimed at Concept explanations, is the *structured explanation* approach, where explanations are provided as complex structures that the human can act upon. [15] presents an application on concept hierarchies, where concepts are laid on a tree-like hierarchy such that the concept of a parent node, e.g., “Animal”, is a generalization of the concepts in its children, e.g., “Dog” and “Cat”. The machine, based on a k-NN model, is tasked to solve two tasks: a downstream task, and a *concept drift* task, that is, to identify if the relationships within the structure have changed. Once detected, the machine presents the concepts of interest to the human, who in turn corrects their structure, e.g., by removing or adding concepts, or by acting on the structure itself, i.e., removing or adding parent-child relationships between concepts. The correction is integrated by removal/addition of appropriate instances from the training set of the machine, thus directly impacting the k-NN model.

5.2.2 Time of interaction. Artifacts are integrated at different times in the lifetime of the machine, hence either at training or at inference time. In the former case, the machine integrates the artifact at training time, and cannot be interacted with at inference time. In this case, the human agent is effectively providing feedback only on training. In the latter case, the human agent has more control on the machine, and receives and acts upon interaction artifacts at inference time. Machines based on explanation languages, such as XIL, tend to provide training-time interaction, while more recent approaches based on hard or soft reasoning languages tend to provide inference-time interaction.

5.2.3 Cost of Learning. A critical distinction in Learning Together systems is the cost of learning, that is, the cost the system has to pay to properly integrate the actions of the human within the machine. Systems tend to fall to the two ends of the spectrum. In traditional approaches such as XIL, the interaction triggers a costly training step for the machine. Here, the cost varies with respect to both the magnitude of the human agent correction and the intrinsic features of the machine. Conversely, more recent approaches, such as extrinsic artifact-based systems, have no additional cost due to the nature of the machine itself. Why then rely on traditional approaches if they incur in an inevitable additional cost? Extrinsic artifacts have to be generated in the first place: the cost of populating the artifact bank is as inevitable as the training cost for traditional approaches. As previously mentioned, when such banks already exist, e.g., in knowledge graphs, this cost can be greatly reduced.

5.3 Strengths and limitations

Unlike the Human overseers and Learn to Abstain paradigms, Learn Together systems integrate humans and machine in a fully hybrid system with bidirectional communication. Thus, the two agents are rarely truly decoupled, their design and use being often task, domain, data and user specific. Machines are limited by and tightly coupled with the language of choice, thus *i*) designing such systems requires significant ground-up effort, and *ii*) effective communication in one language does not appear to transfer to other languages. This extreme heterogeneity partially hinders progress, each system often yielding valuable insight only for future systems operating in similar tasks, domains, data type and desired user interaction. In other words, Learn Together systems tend to improve *vertically*, their success being often unpredictable.

At the moment, Learn Together systems are largely static: they are not able to switch language on demand, nor to adapt to different human agents, or defend themselves against possibly incorrect feedback provided by the human agents. As a relatively recent development in Hybrid systems, they also tend to lack proper measures of validation tailored to them. While classical validation measures apply to the system as a whole, as a predictive system, there appears to be little effort in properly assessing the compliance of a machine with the corrections provided by the human agent, and provide guarantees on the effects of such corrections.

6 LITERATURE GAPS

Hybrid Systems are a still novel topic, and there are several problems that are open to new research. Following our taxonomy, we identify three categories of open challenges:

Machine-related problems. Unlike Learning Together systems, Human Oversight and Learning to Abstain systems still lack meaningful means of communication to engage with the human agent. When presented with a prediction to oversee, or with an uncertainty estimate of such prediction, human agents are rarely also presented with suggestions on why the prediction should be accepted or rejected, or why the machine is uncertain of its prediction, let alone how to tackle the uncertainty itself. We have highlighted some first steps in tackling this problem in Subsection 3.2, but this is far from a solved problem.

Human-related problems. Both Learning to Defer and Learning Together systems require a considerable human effort, i.e., a high volume of labels or artifacts, to be implemented. This is a particularly taxing for the human agent, and while current solutions aim to tackle the problem with a given fixed budget, there are no clear solutions as to how to reduce such high cost. When it comes to Learning to Defer, human-AI “collaboration” heavily depends on a global data annotation industry, wherein a vast and often invisible human labor force is engaged with tedious and taxing jobs. Yet, the prevailing labor standards for data annotators are mostly characterized by lax regulations, low wages and few legal protections [123].

Interaction-related problems. Related to the interaction itself, there is still little experimentation on understanding what language is best suited to which task, and how to adapt the language to different users within the same Hybrid System. In particular, current hybrid systems are *monolingual*, and are developed with one language and one type of human agent in mind, thus lack the ability to adapt to different humans. Underlying the development of such systems is the assumption that the human agents in the system and their capabilities and understanding are static, which is rarely the case. As a result, hybrid systems lack flexibility and have little ability to adapt to the heterogeneity of the humans they may interact with.

When Hybrid Systems interact with several users, e.g. in Multiple-Expert Learning to Defer and in Learning Together systems, machine agents are rarely able to seamlessly handle higher multiplicity of agents, and basic concepts such as malicious agents who wish to poison the system, steer it in an undesirable direction, or simply introduce conflicts with existing artifacts are yet to be properly defined and studied.

7 CONCLUSIONS

Hybrid Systems, where humans and machines collaborate in predictive tasks, represent a new paradigm for how AI systems are designed and implemented. Synergetically integrating these two types of heterogeneous agents allows to maximally exploit the strength of both, each also trying to overcome the weaknesses of the other. In this survey, we presented a taxonomy of the fundamental literature for Hybrid Systems, categorizing the research in this field in three paradigms: Human Overseers, Learn to Abstain and Learn Together. Each paradigm represents progressive steps towards greater human interaction with the AI system: while in Human Oversight the human merely verifies the final output of the system in Learn To Abstain, the human’s knowledge is represented in the data and used to fit the AI agent specifically to the human expertise. Finally, in Learn Together, the human actively corrects the reasoning AI agent. We provided a thorough overview of the most important works for understanding each paradigm, and highlighted strengths and weaknesses of each approach. We hope that this survey can serve as a strong foundation for future research on Hybrid Systems.

References

- [1] Alex Albright. 2019. If you give a judge a risk score: evidence from Kentucky bail decisions. *Law, Economics, and Business Fellows’ Discussion Paper Series*, 85.
- [2] Saar Alon-Barkat and Madalina Busuioc. 2023. Human–AI interactions in public sector decision making: “automation bias” and “selective adherence” to algorithmic advice. *Journal of Public Administration Research and Theory*, 33, 1.
- [3] Saleema Amershi, Maya Cakmak, W. Bradley Knox, and Todd Kulesza. 2014. Power to the people: the role of humans in interactive machine learning. *AI Mag.*, 35, 4.
- [4] Shin Ando and Ayaka Yamamoto. [n. d.] Anomaly detection via few-shot learning on normality. In *ECML/PKDD 2023*.
- [5] Anastasios N Angelopoulos, Stephen Bates, Adam Fisch, Lihua Lei, and Tal Schuster. 2022. Conformal risk control. *arXiv preprint arXiv:2208.02814*.
- [6] Amina Asif and Fayyaz ul Amir Afsar Minhas. 2020. Generalized neural framework for learning with rejection. In *2020 International Joint Conference on Neural Networks (IJCNN)*.

- [7] Alexander Babuta and Marion Oswald. 2021. Machine learning predictive algorithms and the policing of future crimes: governance and oversight. In *Predictive Policing and Artificial Intelligence*.
- [8] Lisanne Bainbridge. 1983. Ironies of automation. *Autom.*, 19, 6.
- [9] Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S. Weld. 2021. Is the most accurate AI the best teammate? Optimizing AI for teamwork. In *AAAI Conference on Artificial Intelligence*.
- [10] Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S. Lasecki, Daniel S. Weld, and Eric Horvitz. 2019. Beyond accuracy: the role of mental models in human-ai team performance. In *HCOMP*.
- [11] Gagan Bansal and Daniel S. Weld. 2018. A coverage-based utility model for identifying unknown unknowns. In *AAAI*.
- [12] Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. [n. d.] Convexity, classification, and risk bounds. *Journal of the American Statistical Association*.
- [13] Peter L. Bartlett and Marten H. Wegkamp. [n. d.] Classification with a reject option using a hinge loss. *JMLR*.
- [14] Roberto Battiti and Anna Colla. [n. d.] Democracy in neural nets: voting schemes for classification. *Neural Networks*.
- [15] Andrea Bontempelli, Fausto Giunchiglia, Andrea Passerini, and Stefano Teso. 2022. Human-in-the-loop handling of knowledge drift. *Data Min. Knowl. Discov.*, 36, 5.
- [16] Kaj Bostrom, Xinyu Zhao, Swarat Chaudhuri, and Greg Durrett. 2021. Flexible generation of natural language deductions. In *EMNLP*.
- [17] Johannes Brinkrolf and Barbara Hammer. 2018. Interpretable machine learning with reject option. *at - Automatisierungstechnik*, 66.
- [18] Johannes Brinkrolf and Barbara Hammer. [n. d.] Probabilistic extension and reject options for pairwise lvq. In *International Workshop on Self-Organizing Maps and Learning Vector Quantization, Clustering and Data Visualization*.
- [19] Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. 2006. Model compression. In *SIGKDD 2006*.
- [20] Zana Bućinca, Maja Barbara Malaya, and Krzysztof Z. Gajos. 2021. To trust or to think: cognitive forcing functions can reduce overreliance on AI in ai-assisted decision-making. *Proc. ACM Hum. Comput. Interact.*, 5, CSCW1.
- [21] Hubert Cecotti and Szilárd Vajda. 2013. Rejection schemes in multi-class classification - application to handwritten character recognition. In *ICDAR*. IEEE Computer Society, 445–449.
- [22] Mohammad-Amin Charusaie, Hussein Mozannar, David Sontag, and Samira Samadi. 2022. Sample efficient learning of predictors that complement humans. In *Proc. of the 39th International Conference on Machine Learning*.
- [23] Jiangjie Chen, Wei Shi, Ziquan Fu, Sijie Cheng, Lei Li, and Yanghua Xiao. 2023. Say what you mean! Large language models speak too positively about negative commonsense knowledge. In *ACL (1)*.
- [24] C. Chow. 1970. On optimum recognition error and reject tradeoff. *IEEE Transactions on Information Theory*, 16, 1.
- [25] Manuel R. Ciosici, Joe Cecil, Dong-Ho Lee, Alex Hedges, Marjorie Freedman, and Ralph M. Weischedel. 2021. Perhaps ptlms should go to school - A task to assess open book and closed book QA. In *EMNLP (1)*.
- [26] Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. In *IJCAI*.
- [27] Lize Coenen, Ahmed K. A. Abdullah, and Tias Guns. 2020. Probability of default estimation, with a reject option. *2020 IEEE 7th International Conference on Data Science and Advanced Analytics (DSAA)*.
- [28] David Cohn, Les Atlas, and Richard Ladner. [n. d.] Improving generalization with active learning. *Machine Learning*.
- [29] David A. Cohn, Les E. Atlas, and Richard E. Ladner. 1994. Improving generalization with active learning. *Mach. Learn.*, 15, 2.
- [30] Stuart Coles. 2001. *An introduction to statistical modeling of extreme values*. Springer Series in Statistics. London.
- [31] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. 2016. Learning with rejection. In *International Conference on Algorithmic Learning Theory*. Springer.
- [32] Christoph Dalitz. 2009. Reject options and confidence measures for knn classifiers. *Schriftenreihe des Fachbereichs Elektrotechnik und Informatik Hochschule Niederrhein*, 8, 2009, 16–38.
- [33] A. P. Dawid. [n. d.] The well-calibrated bayesian. *Journal of the American Statistical Association*. <https://doi.org/10.1080/01621459.1982.10477856>.
- [34] Abir De, Paramita Koley, Niloy Ganguly, and Manuel Gomez-Rodriguez. 2020. Regression under human assistance. In *AAAI Conference on Artificial Intelligence*.
- [35] Abir De, Nastaran Okati, Ali Zaregade, and Manuel Gomez-Rodriguez. 2020. Classification under human assistance. In *AAAI Conference on Artificial Intelligence*.
- [36] C. De Stefano, C. Sansone, and M. Vento. 2000. To reject or not to reject: that is the question-an answer in case of neural classifiers. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 30, 1.
- [37] Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. 2015. Algorithm aversion: people erroneously avoid algorithms after seeing them err. *Journal of Experimental Psychology: General*, 144, 1.
- [38] Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. 2018. Overcoming algorithm aversion: people will use imperfect algorithms if they can (even slightly) modify them. *Management science*, 64, 3.

- [39] Kevin Ellis, Catherine Wong, Maxwell I. Nye, Mathias Sablé-Meyer, Lucas Morales, Luke B. Hewitt, Luc Cary, Armando Solar-Lezama, and Joshua B. Tenenbaum. 2021. Dreamcoder: bootstrapping inductive program synthesis with wake-sleep library learning. In *PLDI*.
- [40] Birte Englich, Thomas Mussweiler, and Fritz Strack. 2006. Playing dice with criminal sentences: the influence of irrelevant anchors on experts' judicial decision making. *Personality and Social Psychology Bulletin*, 32, 2.
- [41] Andrew Guthrie Ferguson. 2020. High-tech surveillance amplifies police bias and overreach. *The Conversation*, 12.
- [42] L. Fischer, Barbara Hammer, and H. Wersing. 2015. Efficient rejection strategies for prototype-based classification. *Neurocomputing*, 169, (Apr. 2015).
- [43] Lydia Fischer, Barbara Hammer, and Heiko Wersing. [n. d.] Local rejection strategies for learning vector quantization. In *ICANN '14*. Springer.
- [44] Vojtech Franc and Daniel Prusa. [n. d.] On discriminative learning of prediction uncertainty. In *Proceedings of ICML '19*.
- [45] Ruijiang Gao, Maytal Saar-Tsechansky, Maria De-Arteaga, Ligong Han, Min Kyung Lee, and Matthew Lease. 2021. Human-ai collaboration with bandit feedback. In *International Joint Conference on Artificial Intelligence*.
- [46] Yonatan Geifman and Ran El-Yaniv. 2017. Selective classification for deep neural networks. In *Advances in Neural Information Processing Systems*. Vol. 30.
- [47] Yonatan Geifman and Ran El-Yaniv. 2019. Selectivenet: a deep neural network with an integrated reject option. In *International Conference on Machine Learning*.
- [48] Tarleton Gillespie. 2018. *Custodians of the Internet: Platforms, content moderation, and the hidden decisions that shape social media*.
- [49] John W Goodell, Satish Kumar, Weng Marc Lim, and Debidutta Pattnaik. 2021. Artificial intelligence and machine learning in finance: identifying foundations, themes, and research clusters from bibliometric analysis. *Journal of Behavioral and Experimental Finance*, 32.
- [50] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. <http://www.deeplearningbook.org>.
- [51] Yves Grandvalet, Alain Rakotomamonjy, Joseph Keshet, and Stéphane Canu. 2008. Support vector machines with a reject option. In *Advances in Neural Information Processing Systems*. Vol. 21.
- [52] Nina Grgić-Hlača, Elissa M. Redmiles, Krishna P. Gummadi, and Adrian Weller. 2018. Human perceptions of fairness in algorithmic decision making: A case study of criminal risk prediction. In *Proc. of the 2018 World Wide Web Conference on World Wide Web*.
- [53] Nina Grgić-Hlača, Gabriel Lima, Adrian Weller, and Elissa M. Redmiles. 2022. Dimensions of diversity in human perceptions of algorithmic fairness. In *Equity and Access in Algorithms, Mechanisms, and Optimization*.
- [54] Riccardo Guidotti, Anna Monreale, Fosca Giannotti, Dino Pedreschi, Salvatore Ruggieri, and Franco Turini. 2019. Factual and counterfactual explanations for black box decision making. *IEEE Intell. Syst.*, 34, 6.
- [55] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2019. A survey of methods for explaining black box models. *ACM Comput. Surv.*, 51, 5.
- [56] Lijie Guo, Elizabeth M. Daly, Ozgur Alkan, Massimiliano Mattetti, Owen Corneic, and Bart P. Knijnenburg. 2022. Building trust in interactive machine learning via user contributed interpretable rules. In *IUI 2022*.
- [57] László Györfi, Z. Györfi, and István Vajda. 1979. Bayesian decision with rejection. *Problems of Control and Information Theory*, 8, (Jan. 1979).
- [58] Lars Kai Hansen, Christian Liisberg, and Peter Salamon. 1997. The Error-Reject Tradeoff. *Open Systems & Information Dynamics*, 4, 2, (Apr. 1997).
- [59] Patrick Hemmer, Sebastian Schellhammer, Michael Vössing, Johannes Jakubik, and Gerhard Satzger. 2022. Forming effective human-AI teams: building machine learning models that complement the capabilities of multiple experts. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*. (July 2022).
- [60] Patrick Hemmer, Lukas Thede, Michael Vössing, Johannes Jakubik, and Niklas Köhl. 2023. Learning to defer with limited expert predictions. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37, 5, (June 2023).
- [61] Patrick Hemmer, Monika Westphal, Max Schemmer, Sebastian Vetter, Michael Vössing, and Gerhard Satzger. 2023. Human-ai collaboration: the effect of ai delegation on human task performance and task satisfaction. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*.
- [62] Kilian Hendrickx, Lorenzo Perini, Dries Van der Plas, Wannes Meert, and Jesse Davis. 2021. Machine learning with a reject option: A survey. *CoRR*, abs/2107.11277. arXiv: 2107.11277.
- [63] Mireille Hildebrandt and Katja de Vries. 2013. *Privacy, due process and the computational turn: The philosophy of law meets the philosophy of technology*.
- [64] Victoria J. Hodge and Jim Austin. [n. d.] A Survey of Outlier Detection Methodologies. *Artificial Intelligence Review*.
- [65] Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. (comet-) atomic 2020: on symbolic and neural commonsense knowledge graphs. In *AAAI*.

- [66] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. 2017. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*.
- [67] Erik Jones, Shiori Sagawa, Pang Wei Koh, Ananya Kumar, and Percy Liang. 2021. Selective classification can magnify disparities across groups. In *ICLR*. OpenReview.net.
- [68] Erik Jones, Shiori Sagawa, Pang Wei Koh, Ananya Kumar, and Percy Liang. 2021. Selective classification can magnify disparities across groups. In *ICLR 2021*.
- [69] Pedro Ribeiro Mendes Júnior, Roberto Medeiros de Souza, Rafael de Oliveira Werneck, Bernardo V. Stein, Daniel V. Pazinato, Waldir R. de Almeida, Otávio Augusto Bizetto Penatti, Ricardo da Silva Torres, and Anderson Rocha. 2017. Nearest neighbors distance ratio open-set classifier. *Machine Learning*, 106.
- [70] Nora Kassner, Oyvind Tafjord, Hinrich Schütze, and Peter Clark. 2021. Beliefbank: adding memory to a pre-trained language model for a systematic notion of belief. In *EMNLP 21*.
- [71] Shachar Kaufman, Saharon Rosset, Claudia Perlich, and Ori Stitelman. 2012. Leakage in data mining: formulation, detection, and avoidance. *ACM Trans. Knowl. Discov. Data*, 6, 4.
- [72] Hendrik Kempt, Jan-Christoph Heilinger, and Saskia K Nagel. 2022. “i’m afraid i can’t let you do that, doctor”: meaningful disagreements with ai in medical contexts. *AI & society*.
- [73] Gavin Kerrigan, Padhraic Smyth, and Mark Steyvers. 2021. Combining human predictions with model probabilities via confusion matrices and calibration. In *Advances in Neural Information Processing Systems*. Vol. 34.
- [74] Vijay Keswani, Matthew Lease, and Krishnaram Kenthapadi. 2021. Towards unbiased and accurate deferral to multiple experts. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*.
- [75] Pang Wei Koh and Percy Liang. [n. d.] Understanding black-box predictions via influence functions. In *ICML*.
- [76] Riikka Koulu. 2020. Proceduralizing control and discretion: human oversight in artificial intelligence policy. *Maas-tricht Journal of European and Comparative Law*, 27, 6.
- [77] Isaac Lage and Finale Doshi-Velez. 2020. Learning interpretable concept-based models with human feedback. *CoRR*, abs/2012.02898. arXiv: 2012.02898.
- [78] Thomas A. Lampert, André Stumpf, and Pierre Gançarski. 2016. An empirical study into annotator agreement, ground truth estimation, and algorithm evaluation. *IEEE Transactions on Image Processing*, 25, 6.
- [79] Thomas Landgrebe, David Tax, Pavel Paclik, and Robert Duin. 2006. The interaction between classification and reject performance for distance-based reject-option classifiers. *Pattern Recognition Letters*, 27, (June 2006).
- [80] Hoel Le Capitaine and C. Frelicot. 2010. An optimum class-rejective decision rule and its evaluation. *International Conference on Pattern Recognition*, (Aug. 2010).
- [81] John D Lee and Katrina A See. 2004. Trust in automation: designing for appropriate reliance. *Human factors*, 46, 1.
- [82] Min Kyung Lee. 2018. Understanding perception of algorithmic decisions: fairness, trust, and emotion in response to algorithmic management. *Big Data & Society*, 5, 1.
- [83] Diogo Leitão, Pedro Saleiro, Mário A. T. Figueiredo, and Pedro Bizarro. 2022. Human-ai collaboration in decision-making: beyond learning to defer. (2022).
- [84] Patrick S. H. Lewis et al. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *NeurIPS 2020*.
- [85] Dongyun Lin, Lei Sun, Kar-Ann Toh, Jing Bo Zhang, and Zhiping Lin. 2018. Twin svm with a reject option through roc curve. *Journal of the Franklin Institute*, 355, 4.
- [86] Yi Lin. [n. d.] A note on margin-based loss functions in classification. *Statistics & Probability Letters*, ().
- [87] Jessie Liu, Blanca Gallego, and Sebastiano Barbieri. 2022. Incorporating uncertainty in learning to defer algorithms for safe computer-aided diagnosis. *Scientific Reports*, 12, 1, (Feb. 2022).
- [88] Jennifer M Logg, Julia A Minson, and Don A Moore. 2019. Algorithm appreciation: people prefer algorithmic to human judgment. *Organizational Behavior and Human Decision Processes*, 151.
- [89] Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *NIPS*.
- [90] David Madras, Toni Pitassi, and Richard Zemel. 2018. Predict responsibly: improving fairness and accuracy by learning to defer. In *Advances in Neural Information Processing Systems*. Vol. 31.
- [91] Atefeh Mahdavi and Marco Carvalho. [n. d.] A survey on open set recognition. In *2021 IEEE Fourth International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*.
- [92] Frank Main. 2016. Cook county judges not following bail recommendations: study. *Chicago Sun-Times*.
- [93] Lorenzo Malandri, Fabio Mercorio, Mario Mezzanzanica, and Navid Nobani. 2023. Convxai: a system for multimodal interaction with any black-box explainer. *Cogn. Comput.*, 15, 2.
- [94] Gianclaudio Malgieri and Giovanni Comandé. 2017. Why a right to legibility of automated decision-making exists in the general data protection regulation. *International Data Privacy Law*, 7, 4.
- [95] Natalia Markovich and Marijus Vaičiulis. [n. d.] Extreme value statistics for evolving random networks. *Mathematics*.
- [96] Natalia Martinez, Martin Bertran, and Guillermo Sapiro. 2020. Minimax pareto fairness: a multi objective perspective. In *Proceedings of the 37th International Conference on Machine Learning (ICML’20)*.

- [97] Bhavana Dalvi Mishra, Oyvind Tafjord, and Peter Clark. 2022. Towards teachable reasoning systems: using a dynamic memory of user feedback for continual system improvement. In *EMNLP '22*.
- [98] Hussein Mozannar, Hunter Lang, Dennis Wei, Prasanna Sattigeri, Subhro Das, and David A. Sontag. 2023. Who should predict? exact algorithms for learning to defer to humans. *ArXiv*, abs/2301.06197.
- [99] Hussein Mozannar and David Sontag. 2020. Consistent estimators for learning to defer to an expert. In *ICML 2020*.
- [100] Thomas Müller, Guillermo Pérez-Torró, Angelo Basile, and Marc Franco-Salvador. 2022. Active few-shot learning with FASL. In *International Conference on Applications of Natural Language to Information Systems*.
- [101] Shikhar Murty, Christopher D. Manning, Scott M. Lundberg, and Marco Tulio Ribeiro. 2022. Fixing model bugs with natural language patches. In *EMNLP*.
- [102] Yifat Nahmias and Maayan Perel. 2021. The oversight of content moderation by ai: impact assessments and their limitations. *Harv. J. on Legis.*, 58.
- [103] Matey Neykov, Jun S. Liu, and Tianxi Cai. 2016. On the characterization of a class of fisher-consistent loss functions and its application to boosting. *Journal of Machine Learning Research*, 17, 70.
- [104] Navid Nobani, Fabio Mercorio, and Mario Mezzanzanica. 2021. Towards an explainer-agnostic conversational XAI. In *IJCAI*.
- [105] Nastaran Okati, Abir De, and Manuel Rodriguez. 2021. Differentiable learning under triage. In *Advances in Neural Information Processing Systems*. Vol. 34.
- [106] Andrea Papenmeier, Daniel Hienert, Yvonne Kammerer, Christin Seifert, and Dagmar Kern. 2023. Know what not to know: users' perception of abstaining classifiers. In *2023 ACM Designing Interactive Systems Conference (DIS 2023)*.
- [107] Jagdish Parikh, Alden Lank, and Friedrich Neubauer. 1994. *Intuition: The new frontier of management*.
- [108] Emma Pierson. 2017. Gender differences in beliefs about algorithmic fairness. *CoRR*, abs/1712.09124.
- [109] Ignazio Pillai, Giorgio Fumera, and Fabio Roli. [n. d.] Multi-label classification with a reject option. *Pattern Recognition*.
- [110] Teodora Popordanoska, Mohit Kumar, and Stefano Teso. 2020. Machine guides, human supervises: interactive learning with global explanations. *CoRR*, abs/2009.09723. arXiv: 2009.09723.
- [111] M. Pradier, J. Zazo, S. Parbhoo, R. Perlis, M. Zazzi, and F. Doshi-Velez. 2021. Preferential mixture-of-experts: interpretable models that rely on human expertise as much as possible. In vol. 1.
- [112] Andrea Pugnana and Salvatore Ruggieri. 2023. Auc-based selective classification. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*. Vol. 206. PMLR, (Apr. 2023).
- [113] Joaquin Quinonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. 2008. *Dataset shift in machine learning*.
- [114] Manish Raghavan, Solon Barocas, Jon M. Kleinberg, and Karen Levy. 2020. Mitigating bias in algorithmic hiring: evaluating claims and practices. In *FAT* '20: Conference on Fairness, Accountability, and Transparency*.
- [115] Maithra Raghu, Katy Blumer, Greg S Corrado, Jon M. Kleinberg, Ziad Obermeyer, and Sendhil Mullainathan. 2019. The algorithmic automation problem: prediction, triage, and human effort. *ArXiv*, abs/1903.12220.
- [116] Maithra Raghu, Katy Blumer, Rory Sayres, Ziad Obermeyer, Bobby Kleinberg, Sendhil Mullainathan, and Jon Kleinberg. 2019. Direct uncertainty prediction for medical second opinions. In *ICML 2019*. Vol. 97. (June 2019).
- [117] Pranav Rajpurkar, Emma Chen, Oishi Banerjee, and Eric J Topol. 2022. Ai in health and medicine. *Nature medicine*.
- [118] Naveen Raman and Michael Yee. 2021. Improving learning-to-defer algorithms through fine-tuning. *ArXiv*, abs/2112.10768.
- [119] H. G. Ramaswamy, Ambuj Tewari, and Shivani Agarwal. 2018. Consistent algorithms for multiclass classification with an abstain option. *Electronic Journal of Statistics*, 12.
- [120] Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. 2017. Right for the right reasons: training differentiable models by constraining their explanations. In *IJCAI*.
- [121] Ethan M. Rudd, Lalit P. Jain, Walter J. Scheirer, and Terrance E. Boulton. 2018. The extreme value machine. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40, 3.
- [122] Abulhair Saparov and He He. 2023. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. In *ICLR*.
- [123] Advait Sarkar. 2023. Enough with "human-AI collaboration". In *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems*. (Apr. 2023).
- [124] Walter Scheirer, Anderson Rocha, Ross Micheals, and Terrance Boulton. 2011. Meta-recognition: the theory and practice of recognition score analysis. *IEEE transactions on pattern analysis and machine intelligence*, 33, (Mar. 2011).
- [125] Sambu Seo, Marko Wallat, Thore Graepel, and Klaus Obermayer. 2000. Gaussian process regression: active data selection and test point rejection. *Informatik aktuell*, (Jan. 2000).
- [126] Burr Settles. 2009. Active learning literature survey.
- [127] Glenn Shafer and Vladimir Vovk. 2008. A tutorial on conformal prediction. *Journal of Machine Learning Research*.
- [128] Irene Solaiman and Christy Dennison. 2021. Process for adapting language models to society (PALMS) with values-targeted datasets. In *NeurIPS*.

- [129] Ricardo Sousa, Ajalmar Rocha Neto, Jaime Cardoso, and Guilherme Barreto. 2015. Robust classification with reject option using the self-organizing map. *Neural Computing and Applications*, 26, (Jan. 2015).
- [130] Aaron Springer, Victoria Hollis, and Steve Whittaker. 2018. Dice in the black box: user experiences with an inscrutable algorithm. *CoRR*, abs/1812.03219. arXiv: 1812.03219.
- [131] Harry Surden. 2019. Artificial intelligence and law: an overview. *Georgia State University Law Review*, 35.
- [132] Richard S. Sutton and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction*.
- [133] Sarah Tan, Julius Adebayo, Kori Inkpen, and Ece Kamar. 2018. Investigating human + machine complementarity for recidivism predictions. *CoRR*, abs/1808.09123. arXiv: 1808.09123.
- [134] Stefano Teso. 2019. Toward faithful explanatory active learning with self-explainable neural nets. In *Proceedings of the Workshop on Interactive Adaptive Learning (IAL 2019)*. CEUR Workshop Proceedings.
- [135] Stefano Teso and Kristian Kersting. 2019. Explanatory interactive machine learning. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society, AIES 2019*.
- [136] Srikanth Thudumu, Philip Branch, Jiong Jin, and Jugdutt (Jack) Singh. 2020. A comprehensive survey of anomaly detection techniques for high dimensional big data. *Journal of Big Data*, 7, 1, (July 2020).
- [137] Francesco Tortorella. 2004. Reducing the classification cost of support vector classifiers through an roc-based reject rule. *Pattern Analysis and Applications*, 7, 2.
- [138] Vladimir Vapnik and Akshay Vashist. 2009. A new learning paradigm: learning using privileged information. *Neural Networks*, 22, 5-6.
- [139] Rajeev Verma, Daniel Barrejón, and Eric T. Nalisnick. 2023. Learning to defer to multiple experts: consistent surrogate losses, confidence calibration, and conformal ensembles. In *AISTATS 2023*. Vol. 206.
- [140] Rajeev Verma and Eric Nalisnick. 2022. Calibrated learning to defer with one-vs-all classifiers. In *Proceedings of the 39th International Conference on Machine Learning*. Vol. 162. PMLR.
- [141] Edoardo Vignotto and Sebastian Engelke. [n. d.] Extreme value theory for anomaly detection – the GPD classifier. *Extremes*.
- [142] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57, 10.
- [143] Peng Wang, Qi Wu, Chunhua Shen, Anthony Dick, and Anton Van Den Hengel. 2017. Fvqa: fact-based visual question answering. *IEEE transactions on pattern analysis and machine intelligence*, 40, 10.
- [144] Ruotong Wang, F. Maxwell Harper, and Haiyi Zhu. 2020. Factors influencing perceived fairness in algorithmic decision-making: algorithm outcomes, development procedures, and individual differences. In *CHI '20: CHI Conference on Human Factors in Computing Systems*.
- [145] Zijie J. Wang, Alex Kale, Harsha Nori, Peter Stella, Mark E. Nunnally, Duen Horng Chau, Mihaela Vorvoreanu, Jennifer Wortman Vaughan, and Rich Caruana. 2022. Interpretability, then what? editing machine learning models to reflect human knowledge and values. In *Proc. of KDD '22*.
- [146] Greta Warren, Mark T. Keane, Christophe Guéret, and Eoin Delaney. 2023. Explaining groups of instances counterfactually for XAI: A use case, algorithm and user study for group-counterfactuals. *CoRR*, abs/2303.09297.
- [147] Bryan Wilder, Eric Horvitz, and Ece Kamar. 2020. Learning to complement humans. *ArXiv*, abs/2005.00582.
- [148] Ran El-Yaniv and Yair Wiener. 2010. On the foundations of noise-free selective classification. *Journal of Machine Learning Research*, 11, (May 2010).
- [149] Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu, Qingyun Wang, Heng Ji, and Meng Jiang. 2022. A survey of knowledge-enhanced text generation. *ACM Comput. Surv.*, 54, 11s.
- [150] Xu-Yao Zhang, Guo-Sen Xie, Xiuli Li, Tao Mei, and Cheng-Lin Liu. 2023. A survey on learning to reject. *Proceedings of the IEEE*, 111, 2.
- [151] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. 2021. A comprehensive survey on transfer learning. *Proceedings of IEEE*.
- [152] Liu Ziyin, Zhikang T. Wang, Paul Pu Liang, Ruslan Salakhutdinov, Louis-Philippe Morency, and Masahito Ueda. 2019. Deep gamblers: learning to abstain with portfolio theory. *Advances in Neural Information Processing Systems*.

APPENDICES

A PROPERTIES OF SURROGATE LOSS FUNCTIONS IN CLASSIFICATION TASKS

Shifting the learning objective from the empirical risk minimization of a loss function to that of a *surrogate* loss function has valuable computational gains. However, a number of statistical properties need to be assessed in order to guarantee the theoretical robustness of optimization results. This section provides a succinct overview of the statistical implications associated with the key conditions that are deemed desirable for a surrogate loss function.

A.1 Pointwise Fisher consistency

Fisher Consistency (FC) is a desirable property of statistical estimators, which posits that if an estimator were computed using the complete population instead of a sample, it would yield the true value of the estimated parameter [5]. When the estimation procedure is the minimization of the risk associated to a loss function, FC is a necessary condition for reasonable performance, since it guarantees that the loss represents the correct objective function [8]. Furthermore, consistency results ensure that the optimization of a surrogate function does not hinder the search for a function that attains the Bayes risk. Consequently, it provides a theoretical foundation for the use of computationally efficient techniques that are specifically suited to convex functions, such as surrogate losses, and not to discrete losses [1]. In the context of Learning to Abstain, the formal definition FC with respect to surrogate losses follows the one stated for K -class classification, where the optimization objective is to minimize the average loss function $\mathcal{L}(y, g_1(x), \dots, g_k(x))$ in order to learn the real-valued scoring function g_i , for $i \in \{1, \dots, K\}$. These are then used to classify the samples as $\hat{y} = \arg \max_i g_i(x)$ [10]. Formally:

Definition A.1. [Fisher Consistency in L2D [11, 18]] A surrogate loss $\tilde{\mathcal{L}}$ is a consistent loss function with respect to another loss \mathcal{L} if minimizing the surrogate loss over all measurable scoring function is equivalent to minimizing the original loss.

The existence of classes of loss functions, whether convex or not, that can attain FC has been proven for many learning problems, including binary classification [8] and multiclass classification [9, 15, 19]. To date, in the L2D setting the surrogate loss functions that meet the Fisher consistency criteria are only the two proposed in [12, 18]. For a more complete treatment of Fisher consistency the reader can refer to [5, 6] and for its application in the context of surrogate loss function to [1, 8, 13, 14, 15]. Note that this notion of consistency has been also called *classification-calibration* [1] or simply *consistency* throughout the L2D literature. However, it should not be confused with other concepts commonly referred to with the same name, such as the notion of *asymptotic consistency*, which instead refers to the property that an estimator has if it converges in probability to the true value.

A.2 Confidence calibration

The notion of *confidence calibration* refers to the property of an estimator (e.g., a probabilistic classifier) to generate a predictive distribution that is consistent with the empirical frequencies observed from realized outcomes [3, 18, 16]. That is, a calibrated estimator correctly quantifies the level of uncertainty or confidence associated with its predictions [4]. For instance, in the case of a HS, if a human decision-makers has a 80% probability of being correct, then the system should forecast that the human agent will be correct around 80% of the time. In the context of decision-making, this property is desirable as it guarantees that the outcomes of a model can be interpreted as real-world probabilities (i.e., the true uncertainties of both human and machine), thus promoting the trustworthiness of the system [16]. When considering a HS, such as a Learning to

Abstain model, a properly calibrated machine loss is not enough to ensure optimal performance, as a faulty deferral policy can severely and negatively impact the system performance. Indeed, while the calibration of the underlying classifier can be accomplished by employing post-hoc techniques such as temperature scaling to the classifier sub-component of the optimization objective function [7, 4], it is equally necessary to calibrate the correctness of the human agent. This entails ensuring that the estimated conditional probability, as determined by the deferral function, of the i^{th} human agent being correct on a specific input aligns with their actual probability of being correct on similar inputs. Formally:

Definition A.2 (Confidence calibration [17, 18]). Let $t : \mathcal{X} \rightarrow (0, 1)$ be an estimator of the correctness of a human H , i.e., $t(\cdot) = \mathbb{P}[Y_H = Y \mid X = \cdot]$. Then t is said to be *calibrated* if, for any confidence level $c \in (0, 1)$, the actual proportion of times H is correct is equal to c :

$$\mathbb{P}[Y_{H,i} = Y_i \mid t(X_i) = c] = c \quad \forall c \in (0, 1) \quad \forall i \in \{1, \dots, n\}.$$

Moreover, the level of calibration of t can be computed through the *Expected Calibration Error* (ECE), which is defined as follows:

$$\text{ECE}(t) = \mathbb{E}_X [\mathbb{P}(Y_H = Y \mid t(X) = c) - c].$$

The explicit formulation of t has been derived in [18] for the case of Single-Expert L2D, while it has been derived in [17] for the case of Multiple-Expert L2D.

A.3 Realizable consistency

The notion of consistency of a surrogate loss with respect to another loss introduced in Definition A.1 assumes the optimization to be performed over the entire set of measurable functions. However, in realistic scenarios, the minimization will be very likely carried out only on restricted classes of function, resulting in near-optimal performance whenever the Bayes predictor does not belong to the selected classes. In the context of multiclass classification, this issue has been observed and discussed by Long and Servedio [10], who firstly introduced the terminology of \mathcal{F} -*realizable consistency* when referring to consistency with respect to a restricted class \mathcal{F} of scoring functions. In particular, the empirical findings in [10] indicate that the property of \mathcal{F} -realizable consistency holds greater significance and applicability compared to FC in scenarios where learning algorithms are constrained to a specific class \mathcal{F} of scoring functions, as it is more closely related to classification accuracy. Successively, the same concept has been recalled and adapted to the frameworks of L2R [2] and L2D [12, 11]. Notably, in these settings the hypothesis spaces of both the classifier (i.e., \mathcal{F}_M) and rejector (i.e., \mathcal{P}_M) are taken into consideration.

Definition A.3 (realizable $(\mathcal{F}_M, \mathcal{P}_M)$ -consistency [2, 11, 12]). Let $\tilde{\mathcal{L}}$ be a surrogate loss function of the loss \mathcal{L} . Then $\tilde{\mathcal{L}}$ is said to be realizable $(\mathcal{F}_M, \mathcal{P}_M)$ -consistent if, for all distributions \mathbf{P} over $X \times Y_M$ and any human agent H for which there exist $f^* \in \mathcal{F}_M$ and $\rho_M^* \in \mathcal{P}_M$ such that $\mathcal{L}(f^*, \rho_M^*) = 0$, then for any $\epsilon > 0$ there exists $\delta > 0$ such that if $(\hat{f}, \hat{\rho}_M)$ satisfies:

$$|\tilde{\mathcal{L}}(\hat{f}, \hat{\rho}_M) - \inf_{(f, \rho_M) \in (\mathcal{F}_M, \mathcal{P}_M)} \tilde{\mathcal{L}}(f, \rho_M)| \leq \delta \quad \Rightarrow \quad \mathcal{L}(\hat{f}, \hat{\rho}_M) \leq \epsilon.$$

For instance, Mozannar et al. [11] propose a differentiable and realizable $(\mathcal{F}_M, \mathcal{P}_M)$ -consistent surrogate loss for the L2D setting suitable for hypothesis classes \mathcal{F}_M and \mathcal{P}_M of scoring functions that are close under scaling², e.g., linear functions and feedforward neural networks.

²A class \mathcal{F} of scoring functions is closed under scaling if $f \in \mathcal{F} \Rightarrow \alpha f \in \mathcal{F} \quad \forall \alpha \in \mathbb{R}$.

References

- [1] Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. [n. d.] Convexity, classification, and risk bounds. *Journal of the American Statistical Association*.
- [2] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. 2016. Learning with rejection. In *International Conference on Algorithmic Learning Theory*. Springer.
- [3] A. P. Dawid. [n. d.] The well-calibrated bayesian. *Journal of the American Statistical Association*. <https://doi.org/10.1080/01621459.1982.10477856>.
- [4] Telmo Silva Filho, Hao Song, Miquel Perello-Nieto, Raul Santos-Rodriguez, Meelis Kull, and Peter Flach. 2023. Classifier calibration: a survey on how to assess and improve predicted class probabilities. *Machine Learning*, 112, 9, (May 2023), 3211–3260. doi: 10.1007/s10994-023-06336-7.
- [5] Ronald A. Fisher. 1922. On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 222, 594-604, (Jan. 1922), 309–368. doi: 10.1098/rsta.1922.0009.
- [6] Ken Gerow. 1989. Fisher consistency - the evolution of a concept: it's hard to get it right the first time. *Biometrics Unit Technical Reports*, BU-1022-M, 1–23. <https://ecommons.cornell.edu/handle/1813/31595>.
- [7] 2019. *Beyond temperature scaling: obtaining well-calibrated multiclass probabilities with dirichlet calibration*. *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. Curran Associates Inc., Red Hook, NY, USA.
- [8] Yi Lin. [n. d.] A note on margin-based loss functions in classification. *Statistics & Probability Letters*, ().
- [9] Yufeng Liu. 2007. Fisher consistency of multiclass support vector machines. In *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics* (Proceedings of Machine Learning Research). Marina Meila and Xiaotong Shen, (Eds.) Vol. 2. PMLR, San Juan, Puerto Rico, 291–298. <https://proceedings.mlr.press/v2/liu07b.html>.
- [10] Phil Long and Rocco Servedio. 2013. Consistency versus realizable h-consistency for multiclass classification. In *Proceedings of the 30th International Conference on Machine Learning* (Proceedings of Machine Learning Research) number 3. Sanjoy Dasgupta and David McAllester, (Eds.) Vol. 28. PMLR, Atlanta, Georgia, USA, 801–809. <https://proceedings.mlr.press/v28/long13.html>.
- [11] Hussein Mozannar, Hunter Lang, Dennis Wei, Prasanna Sattigeri, Subhro Das, and David A. Sontag. 2023. Who should predict? exact algorithms for learning to defer to humans. *ArXiv*, abs/2301.06197.
- [12] Hussein Mozannar and David Sontag. 2020. Consistent estimators for learning to defer to an expert. In *ICML 2020*.
- [13] Matey Neykov, Jun S. Liu, and Tianxi Cai. 2016. On the characterization of a class of fisher-consistent loss functions and its application to boosting. *Journal of Machine Learning Research*, 17, 70.
- [14] Ingo Steinwart. 2007. How to compare different loss functions and their risks. *Constructive Approximation*, 26, 2, (Apr. 2007), 225–287. doi: 10.1007/s00365-006-0662-3.
- [15] Ambuj Tewari and Peter L. Bartlett. 2007. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8, 36, 1007–1025. <http://jmlr.org/papers/v8/tewari07a.html>.
- [16] Juozas Vaicenavicius, David Widmann, Carl R. Andersson, Fredrik Lindsten, Jacob Roll, and Thomas Bo Schön. 2019. Evaluating model calibration in classification. In *International Conference on Artificial Intelligence and Statistics*.
- [17] Rajeev Verma, Daniel Barrejón, and Eric T. Nalisnick. 2023. Learning to defer to multiple experts: consistent surrogate losses, confidence calibration, and conformal ensembles. In *AISTATS 2023*. Vol. 206.
- [18] Rajeev Verma and Eric Nalisnick. 2022. Calibrated learning to defer with one-vs-all classifiers. In *Proceedings of the 39th International Conference on Machine Learning*. Vol. 162. PMLR.
- [19] Hui Zou, Ji Zhu, and Trevor Hastie. 2008. New multiclass boosting algorithms based on multiclass fisher-consistent losses. *The Annals of Applied Statistics*, 2, 4, (Dec. 2008). doi: 10.1214/08-aoas198.