# Conformal Set-based Human-AI Complementarity with Multiple Experts

Helbert Paat
The Hong Kong Polytechnic University
Hong Kong SAR, China
hapaat@polyu.edu.hk

Guohao Shen
The Hong Kong Polytechnic University
Hong Kong SAR, China
guohao.shen@polyu.edu.hk

## ABSTRACT

Decision support systems are designed to assist human experts in classification tasks by providing conformal prediction sets derived from a pre-trained model. This human-AI collaboration has demonstrated enhanced classification performance compared to using either the model or the expert independently. In this study, we focus on the selection of instance-specific experts from a pool of multiple human experts, contrasting it with existing research that typically focuses on single-expert scenarios. We characterize the conditions under which multiple experts can benefit from the conformal sets. With the insight that only certain experts may be relevant for each instance, we explore the problem of subset selection and introduce a greedy algorithm that utilizes conformal sets to identify the subset of expert predictions that will be used in classifying an instance. This approach is shown to yield better performance compared to naive methods for human subset selection. Based on real expert predictions from the CIFAR-10H and ImageNet-16H datasets, our simulation study indicates that our proposed greedy algorithm achieves near-optimal subsets, resulting in improved classification performance among multiple experts.

## KEYWORDS

Prediction Sets; Conformal Prediction Sets; Human-AI Team; Multiple Experts; Human-AI Interaction; Confusion Matrix; Multiclass Classification; Subset Selection

## 1 INTRODUCTION

In recent years, human experts have increasingly relied on AI-based decision support systems to make informed choices in high-risk fields such as medicine, drug discovery, finance, law, and science [10, 13, 21, 26, 27]. Although much existing research focuses on developing sophisticated algorithms, it is essential to advance the paradigm of AI-assisted decision making, where humans and AI collaborate to improve accuracy. Given the impressive performance of modern machine learning models, it is crucial for humans to learn how to effectively leverage these tools for important real-world

tasks. This collaboration, known as *human-AI complementarity*, can lead to better outcomes than when humans and machines operate independently.

Previous studies aim to achieve human-AI complementarity by suggesting that decision support systems should help individuals identify situations where AI offers substantial advantages and provide explanations that clarify the reasoning behind model predictions [20, 22, 37, 47, 49]. This assistance often involves examining the factors that influence trust [31, 44, 50, 52]. However, these investigations have yielded inconclusive results, leading to uncertainty about how experts can mitigate the risk of developing misplaced trust in AI systems. A study conducted by Straitouri et al. [43] developed a system designed to operate without requiring experts to discern when and how to trust AI, thus shifting the focus of the human-AI model from factors such as calibration and explanation. This system generates a set of label predictions, referred to as a prediction set, from which a human expert selects the most appropriate label. In their proposed framework, the construction of these prediction sets is based on conformal predictors [2]. Furthermore, Toni et al. [45] suggest that the conformal predictor may not be the optimal set-valued predictor in this kind of system. They propose a framework for constructing optimal prediction sets that enable human experts to achieve the highest possible accuracy. However, both of these studies are limited in scope, focusing solely on a single expert. We argue that this limitation restricts the potential accuracy and practicality of combined human-AI models. In practice, decision-making typically comes from multiple experts who engage with decision support systems. Our goal is to establish a framework for human-AI complementarity that involves multiple human experts working with a decision support system. For each instance, we propose a human subset selection algorithm and design a framework that incorporates predictions from this subset of multiple experts to generate final predictions in multiclass classification tasks, thereby enhancing decision-making outcomes in complex scenarios.

**Our contributions**. We explore the scenario of multiple human experts collaborating with a decision support system that offers a set of label predictions and requires each expert to select from this set during inference. We begin by establishing a lower bound on the accuracy of a system that incorporates the experts' predictions that are chosen from the conformal sets. Then we identify the conditions that allow these experts to effectively utilize conformal sets instead of relying on the entire label space. We then proceed with the perspective of subset selection of human experts, asking which subset of experts should be chosen to classify each data sample. Inspired by these findings, we propose a greedy algorithm for selecting a subset of human experts for each instance, leveraging the
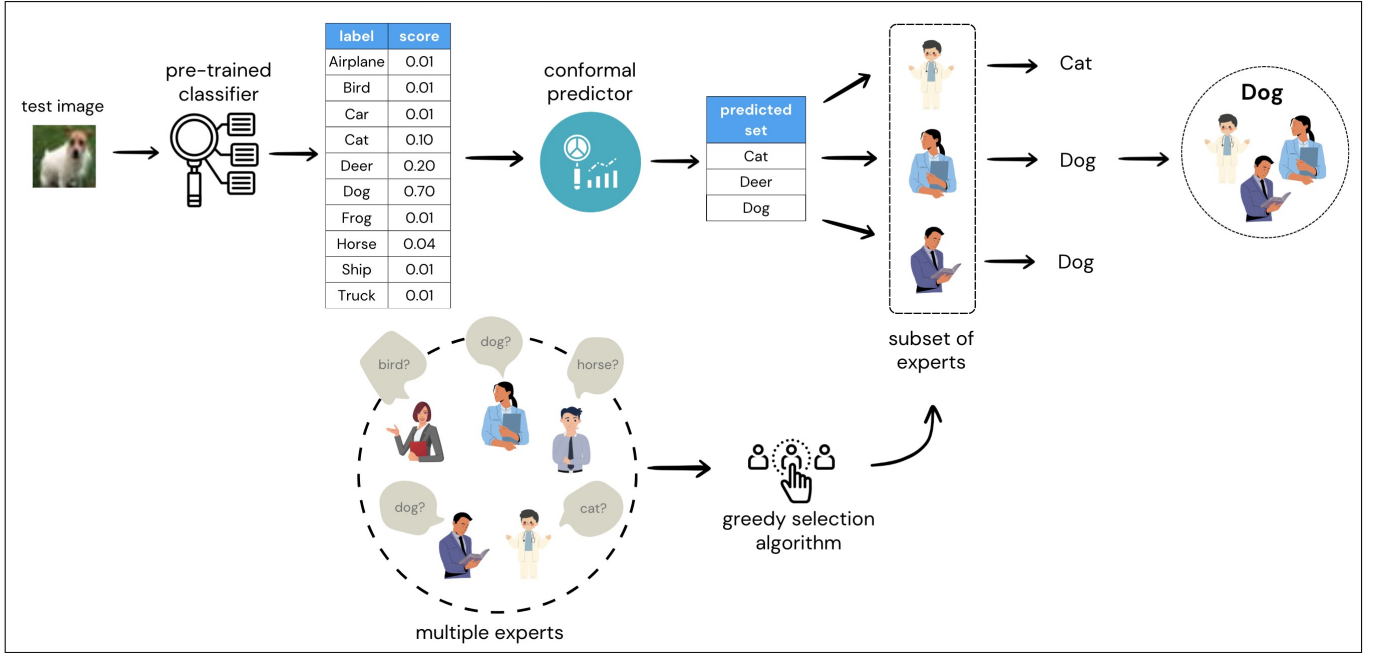
**Figure 1: Illustration of the proposed greedy subset selection algorithm that operates within a multiple expert framework utilizing conformal sets during inference. Initially, a pre-trained classifier computes scores for the conformal predictor. For a given test image, the greedy selection algorithm identifies a subset of human experts. Each selected expert then makes their decision from a narrowed set of options. The final prediction is determined through a combination policy, typically employing a majority decision rule to merge the predictions from multiple experts.**

conformal prediction set to attain near-optimal classification performance. Intuitively, the goal is to favor human experts who, even without prior knowledge of the conformal set, are more inclined to select elements from within that set. This, however, assumes that we understand the likelihood of each expert choosing specific classes based on the true class[1].

Finally, we validate the effectiveness of our approach through simulation studies using real expert predictions on the CIFAR-10H and Imagenet-16H datasets on several multiclass classification tasks. We demonstrate that, for varying numbers of human experts, our proposed greedy selection algorithm consistently outperforms naive approaches of choosing human subsets, such as selecting random subset of humans or choosing the subsets that favor the majority's prediction. This suggests that a strategic human expert selection process can improve classification accuracy. Moreover, we also demonstrate that our approach outperforms subset selection based on top-$k$ set predictors where the top-$k$ sets are the basis for the selection of human subsets. We also show that even as we increase the number of experts, the proposed algorithm is still useful for selecting the appropriate human experts for each instance. Figure 1 shows the general illustration of the proposed framework.

---

[1]We achieve this through the confusion matrix for an expert, which illustrates how an instance is categorized by human experts into different classes given the true label. The conformal set narrows the possible labels for a human, and the confusion matrix is important to assess the relative importance of each label within the conformal set, excluding irrelevant labels.

## 2 RELATED WORKS

In literature, we have seen the emergence of methods that both utilize AI and human expert decisions in prediction tasks to leverage their respective strengths – one such area is the development of classifiers that perform predictions to some samples and rely on the human experts for the remaining ones through a triage policy [4, 5, 9, 12, 34, 36, 39, 46]. One common approach involves learning a binary classifier that defers instances with low model confidence to humans [7, 18]. A class of algorithms known as "learning to defer" (L2D) focuses on training models that adapt to human experts by learning when to either make a prediction themselves or defer the decision to the experts. L2D approaches have been extensively researched [11, 17, 33] and have been adapted for scenarios involving multiple human experts [16, 24, 29, 30]. However, a notable limitation of this body of work is that AI models tend to excel only in handling instances with high confidence, which restricts their generalization capabilities across all instances. More relevant to our study, Babbar et al. [3] propose decision support systems that utilize prediction sets. Unlike their approach, however, our pre-trained classifier assists human experts in solving classification tasks by suggesting prediction sets from which the selected experts can select labels for each data sample during inference.

There is a substantial body of research focused on *set-valued predictors*, which aim to develop models that produce a set of label values known as a prediction set [6]. In several studies on cautious or reliable classification [28, 32, 35], the model-generated

set-valued predictions provide a way to understand the uncertainty in the model's predictions. Recent works have explored decision support systems that utilize prediction sets to assist human experts [3, 8, 43, 45, 51], demonstrating that set-valued predictors can enhance the performance of human experts in prediction tasks. Notably, conformal predictors—a type of set-valued predictor with distribution-free guarantees—have shown promise in this context. While the studies by [43, 45] are closely related to our work, we specifically focus on scenarios involving multiple experts, contrasting their single-expert setting.

Another line of research [23, 41] enhances classification performance by directly integrating human predictions with the output probabilities of a pre-trained classifier during inference. In this framework, data samples are not deferred to either the human or the model; instead, predictions are made through a mathematical combination of the human label and the classifier's probabilities. Kerrigan et al. [23] combine the probabilistic outputs of an AI model with human class-level outputs to improve the accuracy of the human-AI collaboration. More closely related to our work, Singh et al. [41] propose a greedy algorithm that incorporates class-level outputs from multiple human experts alongside the probabilistic outputs of a pre-trained classifier. Their approach identifies the optimal subset of workers for a task, enhancing the combined human-AI decision model. Similarly, we also view multiple expert collaboration as a subset selection problem; however, our focus is on selecting human subsets that will perform classification based on knowledge of conformal prediction sets.

## 3 PRELIMINARY

### 3.1 Problem Formulation

In our setting for human-AI multiclass classification tasks, a set of $h$ human experts $\mathcal{H}$ aims to predict the corresponding label $Y \in \mathcal{Y} = \{1, ..., n\}$ from the feature vector $X \in \mathcal{X}$, given access to an automated decision support system $C : \mathcal{X} \to 2^{\mathcal{Y}}$ that predicts a set of potential labels $C(X) \subseteq \mathcal{Y}$. Here, human experts' predictions for the label of $X$ are $H(X) = \{H_i(X)\}_{i=1}^{h}$ where $H_i(X) \in \mathcal{Y} = \{1, ..., n\}$ and the system $C$ is usually a conformal predictor [1, 48], which helps human experts by narrowing the scope of prediction. The system $C$ established based on a pre-trained classifier $\hat{f} : \mathcal{X} \to [0, 1]^{|\mathcal{Y}|}$, where $\hat{f}(X)$ is the normalized probability vector of the prediction on $X$ and $C(X)$ is determined based on the outputs scores $\hat{f}(X)$. Given the system $C$, our objective is to develop a framework to select a subset $\mathcal{S}(X)$ from human experts $\mathcal{H}$.

Similar to Straitouri et al. [43], the system requires that the final prediction of each expert $H_i(X) \in \mathcal{Y}$ be an element of the narrowed set $C(X)$. Multiple human predictions $\{H_i(X)\}_{i \in \mathcal{S}(X)}$ form the basis for generating the final experts prediction $\hat{Y}$.

Ideally, we expect the designed framework to benefit from the collaborative predictions of experts in subset $\mathcal{S}(X)$ in such a way that:

$$\mathbb{P}[\hat{Y} = Y; C | H, Y \in C(X)] \geq \mathbb{P}[H_i = Y; C | Y \in C(X)] \geq \mathbb{P}[H_i = Y; \mathcal{Y}]$$

for any $i \in \{1, ..., h\}$ where $\mathbb{P}[\hat{Y} = Y; C | H, Y \in C(X)]$ indicates the success probability of the subset $\mathcal{S}(X)$ of multiple experts who predicts a class from the narrowed subset $C(X)$, and $\mathbb{P}[H_i = Y; C | Y \in C(X)]$ indicates the success probability of the $i$th single expert

who chooses from the narrowed options $C(X)$. Subset selection of human experts arises from the intuition that within a team of experts, diverse perspectives may exist. Hence, it becomes necessary to choose a specific subset $\mathcal{S}(X)$ for classification in any given instance. One might instinctively favor selecting the subset $\mathcal{S}(X) = \mathcal{H}$ and employing a majority decision rule to get a final prediction. However, this approach may not be the most effective option, as we will show later.

### 3.2 Conformal Prediction

Given a calibration dataset $D_{\text{cal}} = (x_i, y_i)_{i=1}^{l}$ and a test sample $(x_{test}, y_{test})$, conformal prediction aims to construct a prediction set $C$ based on $D_{\text{cal}}$ such that the marginal coverage at a user-specified tolerance level $\alpha \in [0, 1]$ is satisfied, i.e.

$$\mathbb{P}[y_{test} \in C(x_{test})] \geq 1 - \alpha. \tag{1}$$

Constructing a conformal predictor usually requires calculating the conformal score $s : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ by $s(x_i, y_i) = 1 - \hat{f}_{y_i}(x_i)$ on each sample, where $\hat{f}$ denotes a pre-trained classifier and $\hat{f}_{y_i}(x_i)$ denotes the predicted probability for label $y_i$ given $x_i$. A lower conformal score $s(x_i, y_i)$ means better agreement between the input $x_i$ and the label $y_i$. While a conformal score is close to one implies the pre-trained classifier is significantly incorrect on $x_i$. Let $\hat{q}_\alpha$ be the empirical $\lceil (l + 1)(1 - \alpha) \rceil$-th quantile of the conformal scores $s(x_1, y_1), \ldots, s(x_l, y_l)$. Then the conformal prediction set $C(x_{test})$ is constructed by $C(x_{test}) = \{y : s(x_{test}, y) \leq \hat{q}_\alpha\}$. Note that the constructed set $C(x_{test})$ can be established based on any pre-trained classifier $\hat{f}$ without relying on any distributional assumptions about the data [48]. One can show that, conditioned on the calibration set $\mathcal{D}_{\text{cal}}$, the probability of the true label $y_{test}$ belonging to the subset $C(x_{test})$ is guaranteed at level $(1 - \alpha)$ [1].

## 4 UTILIZING CONFORMAL PREDICTION SETS FOR THE SELECTION OF HUMAN SUBSETS

The multiple experts subset selection problem mathematically involves identifying a specific subset of human experts from a larger set. Let $\mathcal{D}_{test} = \{(x_i, y_i)\}_{i=1}^{t}$ denote the testing set. Given any test sample $(x, y)$ in the test set $\mathcal{D}_{test}$, the goal is to find a subset $\mathcal{S}(x) \subseteq \mathcal{H}$ with final expert predictions $\{H_i(x)\}_{i=1,...,|\mathcal{S}(x)|}$ that maximizes the experts' conditional success probability $\mathbb{P}[\hat{Y} = y; C | H, y \in C(x)]$ where $\hat{Y} = \pi(\{H_i(x)\}_{i=1,...,|\mathcal{S}(x)|})$ is the combination of the selected experts ' prediction based on a given combination method or decision rule $\pi$. This can be expressed as follows:

$$\max_{\mathcal{S}(x) \subseteq \mathcal{H}} \mathbb{P}[\pi(\{H_i(x)\}_{i=1,...,|\mathcal{S}(x)|}) = y; C | H, y \in C(x)].$$

Note that several naive selection methods for forming subsets exist in this case. We can choose the subset $\mathcal{S}(x)$ to be a random subset of $\mathcal{H}$ of fixed size $\gamma < h$ where $h$ is the size of $\mathcal{H}$. The subset $\mathcal{S}(x)$ can also be chosen to be of size 1, containing only the best-performing expert in the team $\mathcal{H}$. Furthermore, the subset $\mathcal{S}(x)$ can be chosen as the whole expert team $\mathcal{H}$. However, as we will show in the simulation study, these naive approaches to select subsets of humans are not guaranteed to perform the best during inference.

In what follows, our goal is to develop a framework to select subsets of experts using the knowledge of the conformal set and

show, theoretically and empirically, that it yields better performance than naive approaches. We first discuss a formulation of the expert's conditional success probability in the context of multiple expert predictions. To establish the theoretical results, we start by considering the entire team of experts. For a feature vector $x$, let $\{p_i\}_{i=1}^h$ be an observation or instantiations of the random variables $H(x) = \{H_i(x)\}_{i=1}^h$. In combining multiple human experts' predictions given the conformal set, we assume independence between human predictions. Using the Bayes rule, we can get the following relationship:

$$\mathbb{P}[\hat{Y} = y | H, y \in C(x)] \propto \mathbb{P}[\hat{Y} = y | y \in C(x)] \times \prod_{i \in [h]} \mathbb{P}[H_i = y | \hat{Y}, y \in C(x)]. \quad (2)$$

Similar to Straitouri et al. [43], the expert's success probability $\mathbb{P}[\hat{Y} = y | H, y \in C(x)]$ is estimated by utilizing the multinomial logit model. For a given sample $(x, y)$ and subset $C(x)$, we assume that the expert's conditional success probability can be estimated as follows:

$$\mathbb{P}[\hat{Y} = y | H, y \in C(x)] = \prod_{i \in [h]} \mathscr{C}_{yp_i}^i = \prod_{i \in [h]} \frac{e^{\mu_{yp_i}^i}}{\sum_{y' \in C(x)} e^{\mu_{y'p_i}^i}} \quad (3)$$

where $\mu_{yp_i}^i$ denotes the preference of the expert $i$ for the label value $p_i \in \mathcal{Y}$, given that the true label is $y$. We set the parameters $\mu_{yp_i}^i = \log C_{yp_i}^i$ where we assume access to a confusion matrix $C$ of the expert predictions in the multiclass classification task. The confusion matrix is estimated based on real expert predictions using maximum likelihood estimation. More specifically, $C = [C_{yy'}]_{y,y' \in \mathcal{Y}}$, where $C_{yy'} = \mathbb{P}[\hat{Y} = y'; \mathcal{Y} | Y = y]$.

To determine the improvement in the classification performance of the framework when we have multiple human expert predictions and knowledge of the conformal sets, we derive a lower bound on the accuracy of the framework (see the proof in Appendix A.1).

LEMMA 1. Given $h$ human expert predictions $H = [p_1, p_2, \ldots, p_h]$, conformal set $C$ with tolerance level $\alpha$, and the $i$th expert's conditional success probability $\mathscr{C}^i$, the lower bound on the accuracy of the combined framework is given as

$$\mathbb{E}[\mathbb{1}(\{\hat{Y} = y\} \cap \{y \in C\})] \geq \mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}_{yp_i}^i}{1 - \mathscr{C}_{yp_i}^i} > 1] \cdot (1 - \alpha). \quad (4)$$

In the following theoretical result, we demonstrate that, in the approach that uses predictions from multiple experts, utilizing a conformal prediction set $C$ within the combination framework results in a tighter lower bound than that derived from considering the entire label space $\mathcal{Y}$ (as established in Lemma 4.1 of Singh et al. [41]), given certain assumptions (see the proof in Appendix A.2).

LEMMA 2. Let $\epsilon = \prod_{i \in [h]} \frac{\mathscr{C}_{yp_i}^i}{1 - \mathscr{C}_{yp_i}^i} - \prod_{i \in [h]} \frac{C_{yp_i}^i}{1 - C_{yp_i}^i}$. Assuming $\epsilon > 1$ and $\alpha = 0$, the combination framework of multiple human expert predictions that utilizes a conformal set achieves a tighter lower bound on accuracy,

$$\mathbb{P}[\prod_{i \in [h]} \frac{C_{yp_i}^i}{1 - C_{yp_i}^i} > \frac{1 - \hat{f}_y}{\hat{f}_y}] \leq \mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}_{yp_i}^i}{1 - \mathscr{C}_{yp_i}^i} > 1]. \quad (5)$$

**Designing a Subset Selection Framework**. Lemma 2 demonstrates the advantage of incorporating the conformal sets $C$ into the combination framework of multiple expert predictions. Motivated by these theoretical results, we describe a framework for selecting a subset of humans $\mathcal{S}(x)$ instead of all humans $\mathcal{H}$ for each test sample. For the task of human subset selection, let $\hat{p}_i{}^2$ denote the initial prediction of the expert $i$. Following the same intuition as Algorithm 1 in Singh et al. [41], the goal is to select a subset such that the derived lower bound (in Equation 7) is maximized. We maximize the left probability term in the lower bound of Lemma 1. Consider that the term $\mathscr{C}_{y\hat{p}_i}^i / (1 - \mathscr{C}_{y\hat{p}_i}^i)$ is greater than 1 if and only if $\mathscr{C}_{y\hat{p}_i}^i > 0.50$. Hence, to maximize the product of these terms, we choose only the corresponding human predictions that satisfy this constraint. However, the ground truth label $y$ is obviously required to maximize this term in Lemma 1. Since we do not have access to $y$ during test, we choose a pseudo label $y^*$ and define it to be the class that maximizes the probability in the derived lower bound (Equation 7). Although similar in intuition as in Singh et al. [41], we note that, instead of choosing from the full class of labels $\mathcal{Y}$, we choose the pseudo label from the conformal set $C(x)$. Thus, given the pseudo label $y^*$, we select the pseudo-optimal human subset as follows:

$$\mathcal{S}(x)^* = \arg \max_{\mathcal{S}(x)} ( \prod_{i \in \mathcal{S}(x)} \frac{\mathscr{C}_{y^*\hat{p}_i}^i}{1 - \mathscr{C}_{y^*\hat{p}_i}^i} ). \quad (6)$$

We emphasize that the pseudo label $y^*$ is in the conformal set, that is, $y^* \in C(x)$. Moreover, we only consider the human predicted labels $\hat{p}_i$'s that are in the conformal set. In addition to excluding humans with corresponding values $\mathscr{C}_{y\hat{p}_i}^i$ of at most 0.50, knowledge of the conformal set also eliminates any human whose initial predicted label is not in the set.

In Algorithm 1, we present the overall greedy humans subset selection method[3] based on the conformal set. Assume that the worst-case (or maximum) size of any conformal set is $c$. Moreover, we assume that the algorithm has access to $\mathscr{C}$, which is calculated based on the confusion matrix $C$ by considering only the probabilities of the classes present in $C(x)$ and normalizing these probabilities row-wise, as in Eq. 3. The algorithm then calculates $\mathscr{C}_{kp_i}^i / (1 - \mathscr{C}_{kp_i}^i)$ for all human predictions $p_i$ and classes $k \in C(x)$, which requires $O(hc)$ steps. Then it calculates the product term $\prod_{i \in \mathcal{H}} \mathscr{C}_{kp_i}^i / (1 - \mathscr{C}_{kp_i}^i)$ but only for each $k \in C$ and when $p_i \in C$. It chooses the pseudo label that maximizes this product term as $y^*$. Calculation of the product term again takes $O(hc)$ steps and pseudo

---

[2]We note that for the subset selection of humans, all experts make initial predictions which are not necessarily in the conformal set. The intuition is to utilize the conformal prediction sets from the decision support system $C$ for additional guidance in selecting the appropriate human experts. However, during inference, the selected human experts choose their final prediction from the conformal set.

[3]The proposed algorithm is a modified version of the one introduced by Singh et al. [41], where we focus on the set $C(x)$ rather than the entire label space $\mathcal{Y}$. In our simulation study, we will also present results using the top-$k$ sets generated by $\hat{f}$ instead of the conformal sets $C(x)$ for human subset selection.

label calculation needs $O(c)$ steps. Finally, given $y^*$, it selects human subsets that satisfies $\mathscr{C}^i_{y^* p_i}/(1-\mathscr{C}^i_{y^* p_i}) > 1$, which needs $O(h)$ steps. Thus, the overall time complexity is $O(hc)$, indicating that the algorithm's performance scales linearly with both the number of humans $h$ and the size of the conformal set $c$.

---

**Algorithm 1** Greedy Selection of Humans based on Conformal Sets $C(x)$

---

**Require:** $h, n, c \in \mathbb{N}$
**Require:** conformal set of indices $C(x)$
**Require:** $n \times n$ matrices $\mathscr{C}^i$, $1 \leq i \leq h$
**Require:** initial predictions $\hat{p}_i(x)$, $1 \leq i \leq h$
1: **for** $i = 1$ to $h$ **do**
2:     **for** $k = 1$ to $c$ **do**
3:         $\mathcal{A} \leftarrow \mathscr{C}^i_{C(x)_k \hat{p}_i(x)}$
4:         $\mathcal{B} \leftarrow 1 - \mathscr{C}^i_{C(x)_k \hat{p}_i(x)}$
5:         $F[i][k] \leftarrow \frac{\mathcal{A}}{\mathcal{B}}$
6:     **end for**
7: **end for**
8: **for** $k = 1$ to $c$ **do**
9:     $S[k] \leftarrow 1$
10:     **for** $i = 1$ to $h$ **do**
11:         **if** $F[i][k] > 1$ **then**
12:             $S[k] \leftarrow S[k] \times F[i][k]$
13:         **end if**
14:     **end for**
15: **end for**
16: $k^* \leftarrow \arg\max_{1 \leq k \leq c} S[k]$
17: $y^* \leftarrow C(x)_{k^*}$
18: $S^* \leftarrow \{1 \leq i \leq h \mid F[i][k^*] > 1\}$

---

**Remarks**. Note that to maximize the lower bound, we can set the user-specified tolerance level $\alpha \approx 0$. However, this may result in prediction sets that cover all classes, rendering them useless and uninformative for practical decision-making. That being noted, in our context—where a classifier is employed in a human-AI collaborative environment—we assume that the classifier chosen for such high-risk tasks has high predictive accuracy and is generally confident in its predictions. Moreover, we assume that the calibration data are well distributed across the classes. In the simulation study where we use accurate pre-trained models and well-distributed calibration data, we achieve a low average set size even though we set a low tolerance level $\alpha$, which results in nearly 100% coverage.

**Naive subset selection methods**. Several naive approaches to human subset selection exist. One such method involves setting $S(x) = \mathcal{H}$, meaning all humans are employed for each data sample, followed by a combination method $\pi$ that applies a majority decision rule[4]. In our simulation study, we refer to this approach as "ALL HUMANS." Another method involves selecting a random subset $S(x)$ from the set of human experts $\mathcal{H}$, ensuring that the average size of these human subsets aligns closely with the average subset size derived from our proposed greedy approach. This approach

also utilizes a majority decision rule and is labeled "RANDOM SUBSET" in our simulations. Additionally, we employ a greedy selection method as described in Algorithm 1, which focuses on the top-$k$ prediction sets for each instance. The top-$k$ classes are determined by the output scores of the pre-trained classifier $\hat{f}$. Instead of using conformal sets $C$ for human subset selection, this method relies on the top-$k$ prediction sets, while still applying the majority decision rule as the combination method $\pi$. In our simulations, we experiment with various values of $k$ and refer to this approach as "TOP-$k$." It is important to note that for both "ALL HUMANS" and "RANDOM SUBSET," the selected humans are required to make their final predictions based on conformal sets. In contrast, for "TOP-$k$," the selected humans must derive their final predictions from the top-$k$ prediction sets rather than from conformal sets.

## 5 EXPERIMENTS ON CIFAR-10H DATASET

In this section, we perform experiments using a dataset of natural images with real expert predictions for a multiclass classification task. We utilize several accurate deep neural network classifiers $\hat{f}$ to compute conformal scores for generating conformal sets, following methodologies similar to those in prior research [43]. Our proposed algorithm is benchmarked against naive subset selection methods and previous human-AI combination approaches that consider both single and multiple experts. Additionally, we compare our results with top-$k$ set-valued predictor baselines. We have included the detailed significance values for each experiment in Appendix B.1.

This simulation study utilizes the (estimated) confusion matrix derived from real expert predictions for the multiclass classification task, along with the multinomial logit model defined in Eq. 3, to evaluate our system's performance—consistent with findings from previous studies [43]. After selecting the subset of humans, we apply a combination method $\pi$, which follows a majority decision rule to determine the final outcome.

### 5.1 Experimental Setup

We conduct experiments on the CIFAR-10H dataset [38], which comprises 10,000 natural images sourced from the CIFAR-10 test set [25]. Each image belongs to one of $n = 10$ classes and includes approximately 50 expert predictions. For this dataset, we utilize three widely recognized deep learning models: DenseNet [19], ResNet [14], and PreResNet [15], as similarly employed in [43]. We split the dataset into three subsets: calibration, estimation, and test set. In our classification tasks, we measure test data accuracy (empirical success probability). Additionally, in frameworks using conformal sets[5], we calculate average set sizes.

### 5.2 Comparison with Naive Approaches

In Figure 2, we compare the empirical success probability $\mathbb{P}[\hat{Y} = Y; C|H]$ on the test data for our proposed greedy selection algorithm with naive methods, such as randomly choosing subsets of humans of average size $\tau$ (the average subset size of our greedy selection approach) and selecting all humans followed by a majority decision rule. We also visualize the results for the best single human

---

[4]Note that we chose the combination method $\pi$ to be the majority decision rule for simplicity. For future work, other combination frameworks can be experimented with.

[5]In the simulation study, the tolerance levels $\alpha$ are set at 0.10%, 0.07%, 0.05%, 0.04%, and 0.03% for calibration data sizes $l$ of 1,000, 1,500, 2,000, 2,500, and 3,000, respectively.

**Table 1: Empirical success probability (in %) achieved by our conformal multi-expert approach during test using the greedy selection algorithm using three accurate neural network-based pre-trained classifiers and five multiple experts ($h = 5$) in comparison with previous baselines on the CIFAR-10H dataset. The single expert's empirical success probability at solving the (original) multiclass task is $\approx 95.24\%$. The calibration and estimation sets each have a size of 1,000. The values are averaged over 10 runs.**

| MODEL NAME | DENSENET | RESNET | PRERESNET |
|---|---|---|---|
| Pre-trained Model Alone | 96.29 ± 0.09 | 92.70 ± 0.20 | 94.41 ± 0.12 |
| SINGLE EXPERT APPROACHES | | | |
| Kerrigan et al. [23] (MAP temp. scal.) | 97.81 ± 0.25 | 97.22 ± 0.40 | 97.56 ± 0.36 |
| Straitouri et al. [43] | 97.70 ± 0.09 | 96.63 ± 0.13 | 97.05 ± 0.18 |
| MULTIPLE EXPERT APPROACHES | | | |
| Singh et al. [41] (Greedy algorithm) | 95.15 ± 0.22 | 95.18 ± 0.21 | 95.25 ± 0.20 |
| Singh et al. [41] (Mode approach) | 97.93 ± 0.13 | 97.23 ± 0.19 | 97.48 ± 0.16 |
| Proposed Greedy Selection Method | **98.48 ± 0.22** | **98.48 ± 0.15** | **98.10 ± 0.38** |

performance and the pre-trained model. While the all-human approach reflects the intuition to favor majority decisions, leading to high empirical success probabilities, it doesn't always yield the best outcomes. In contrast, our greedy algorithm based on the conformal set demonstrates superior effectiveness in selecting human subsets.

### 5.3 Comparison with Top-$k$ Set-Valued Predictors

Figure 3 compares our conformal multi-expert method, which uses a greedy algorithm for selecting human subsets, with a top-$k$ set predictor framework that also employs a greedy selection approach based on the top-$k$ sets. Our method shows a higher empirical success probability.

### 5.4 Comparison with Human-AI Combination Approaches

In Table 1, we compare the empirical success probability of our conformal-based greedy algorithm with existing methods that integrate pre-trained classifier outputs and human predictions. The findings reveal two key insights: multi-expert collaboration yields a higher success probability than relying on a single expert, and our greedy algorithm effectively selects human subsets for classification. Notably, even with a limited set of options, our approach surpasses previous baselines that allowed full label access for humans, demonstrating the conformal predictor's ability to identify meaningful classes for each instance.

## 6 EXPERIMENTS ON IMAGENET-16H DATASET

In this section, we conduct experiments using a different dataset of natural images that includes real expert predictions and features a greater number of classes. For the pretrained models $\hat{f}$, we utilize the VGG19 deep neural network classifier, as provided by Steyvers et al. [42], which has been fine-tuned over 10 epochs. We assess the performance of our proposed greedy subset selection algorithm by

comparing it to naive methods for human subset selection. Additionally, we evaluate our approach against top-$k$ set-valued predictors that implement a modified greedy strategy for selecting human subsets based on the top-$k$ set. Detailed significance values for all experiments can be found in Appendix B.2.

### 6.1 Experimental Setup

We experiment with the ImageNet-16H dataset [42], which comprises 1,200 unique images derived from a subset of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) 2012 database [40]. This dataset includes approximately six predictions made by human experts for each image. Each image can be classified into one of the $n = 16$ categories based on separate human annotations. To increase the difficulty of the annotations for both humans and classifiers, the images are distorted using spatial frequency phase noise with a frequency of $\omega = 80$ Additionally, we randomly divided the images into three subsets: the calibration set, estimation set, and the test set. Typically, these subsets contain 240 images each for the calibration sets and estimation sets, and 720 images for the test set, unless stated otherwise in the experimental results. We use the calibration set[6] to compute the conformal scores necessary for the conformal set predictor. Then we utilize the estimation set (together with the calibration set) to estimate the human confusion matrix. Finally, we utilize the test set to assess the empirical success probability of the experts as applied in various approaches. To illustrate the reduced size of the conformal sets, we also compute the empirical average sizes of these sets.

### 6.2 Comparison with Naive Approaches

Figure 4 illustrates the comparison of empirical success probabilities among the proposed greedy selection method, the random subset selection method, and the all-humans approach when the combination method $\pi$ is set to the majority decision rule. Observably, the proposed greedy method demonstrates superior performance compared to the naive approaches.

### 6.3 Comparison with Top-$k$ Predictors

Figures 5 and 6 compare our conformal multi-expert approach using the greedy algorithm to select human subsets against the framework employing top-$k$ set predictors using mode approach for subset selection. Our framework demonstrates superior empirical success probability, outperforming both single-expert and multi-expert baselines, even when the average prediction set size of the conformal sets is smaller than $k$.

### 6.4 What happens as the number of multiple experts $h$ increases?

In Figure 7, we compare the empirical success probabilities of our multi-expert system using the greedy subset selection method against a human-only expert team and a top-$k$ predictor with a top-$k$ subset selection algorithm (where $k = 5$). The expert team reflects majority predictions without relying on prediction sets or models. Our results show that as the number of human experts increases, our conformal set-based greedy selection approach outperforms
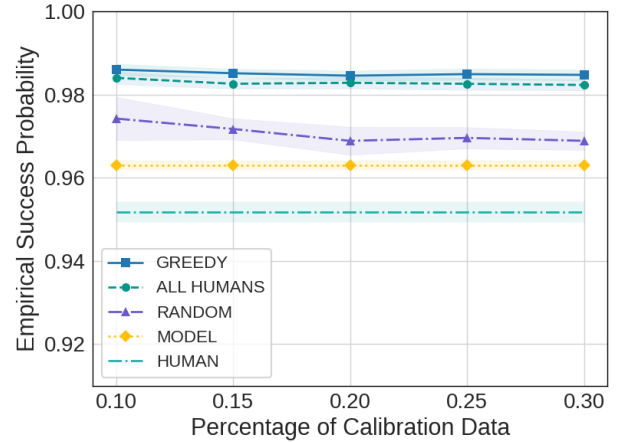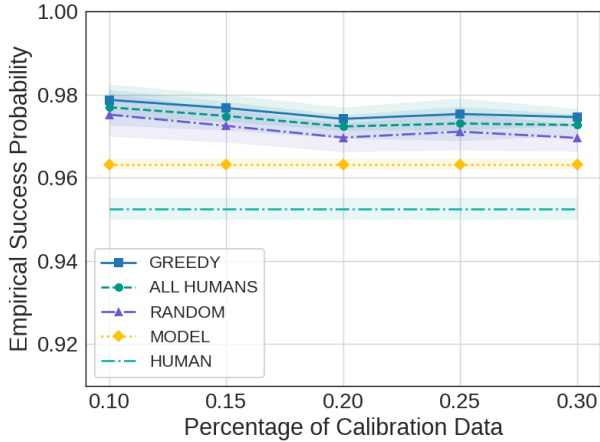
---

Figure 2: Empirical success probability of a multiple expert system utilizing the proposed greedy human subset selection algorithm in comparison with naive approaches across varying percentages of calibration data on the CIFAR-10H dataset. The estimation data matches the calibration data proportionally. The left figure presents results for 3 experts, while the right shows results for 5 experts. "ALL HUMANS" and "RANDOM" refer to all-human and random subset selection methods, respectively, applying a majority decision rule. "MODEL" and "HUMAN" indicate the performance of the classifier and a single expert. The classifier employed is Densenet. The average sizes of the conformal sets are 3.75, 4.79, 5.66, 5.69, and 5.92 for calibration percentages of 10%, 15%, 20%, 25%, and 30%, respectively. Results are averaged over 10 runs, with shaded regions representing one standard deviation from the mean.
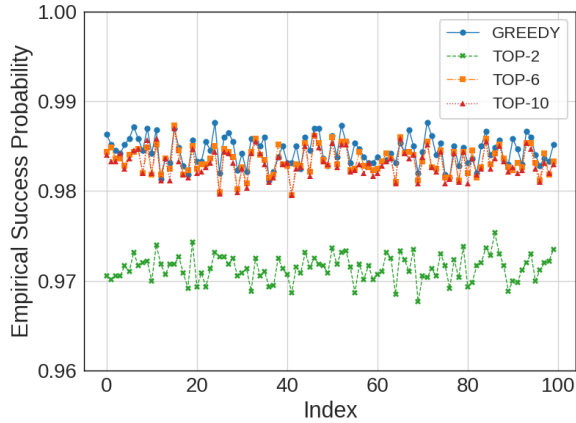


Figure 3: Empirical success probability of our proposed greedy selection algorithm in comparison with top-$k$ set predictor framework for different values of $k$ on the CIFAR-10H dataset using the Resnet classifier for 100 runs (using $h = 5$). The empirical average set size of the conformal sets is 5.2718. The calibration and estimation sets each have a size of 2,000.

both the expert team and the top-5 methods, remaining effective even with a large pool of human experts.

## 7 DISCUSSION

In this section, we explore the assumptions and limitations of the proposed approach.
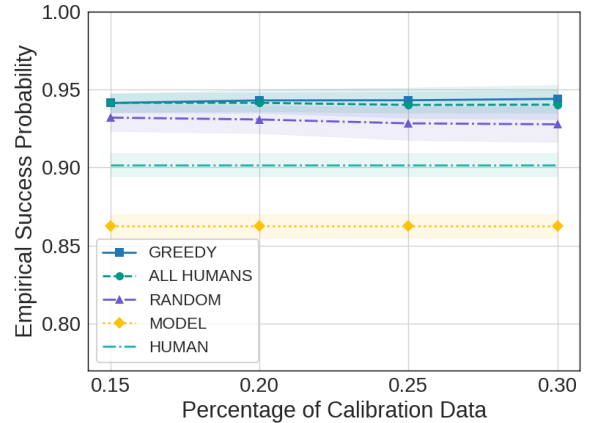


Figure 4: Empirical success probability of the multi-expert system using the proposed greedy human subset selection algorithm compared to naive approaches across various calibration data percentages on the ImageNet-16H dataset, with results displayed for 5 experts. "ALL HUMANS" and "RANDOM" refer to all-human and random subset selection methods, while "MODEL" and "HUMAN" indicate the performance of the classifier and a single expert, respectively. The estimation data is proportionally matched to the calibration data. The average sizes of the conformal sets are 2.67, 2.82, 3.26, and 3.39 for calibration percentages of 15%, 20%, 25%, and 30%. Results are averaged over 100 runs, with shaded regions representing one standard deviation from the mean.
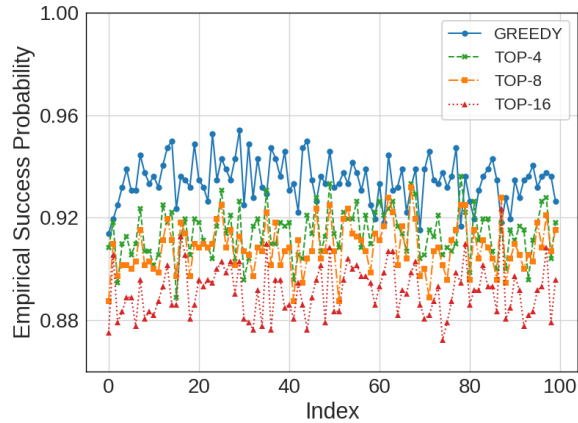
**Figure 5: Empirical success probability achieved by the multi-expert system utilizing the proposed greedy human subset selection algorithm in comparison with top-$k$ set-valued predictors on the ImageNet-16H dataset over 100 runs. The figure illustrates results with 3 experts. Both the calibration and estimation sets comprise 240 images each. The empirical average size of the conformal sets is 2.8247.**
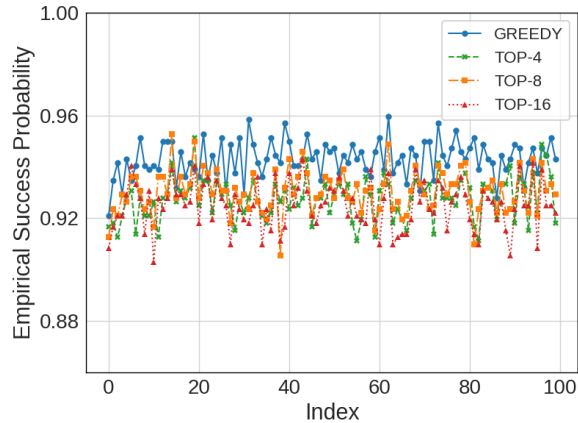


**Figure 6: Empirical success probability of the system using the proposed greedy human subset selection algorithm compared to top-$k$ set-valued predictors on the ImageNet-16H dataset over 100 runs. The figure presents results with 5 experts and maintains the same settings as Figure 5.**

**Framework.** The proposed greedy algorithm relies on an estimated confusion matrix for human performance, based on real expert predictions and following a straightforward Maximum Likelihood Estimation approach similar to that of Kerrigan et al. [23]. Exploring more sophisticated methods for estimating the confusion matrix would be advantageous. Additionally, the assumption of independence among experts may not hold, as their decisions can influence one another. In our framework, we set the tolerance level $\alpha$ close to zero, which minimally impacts average conformal set sizes, but this may not always be practical. Future research
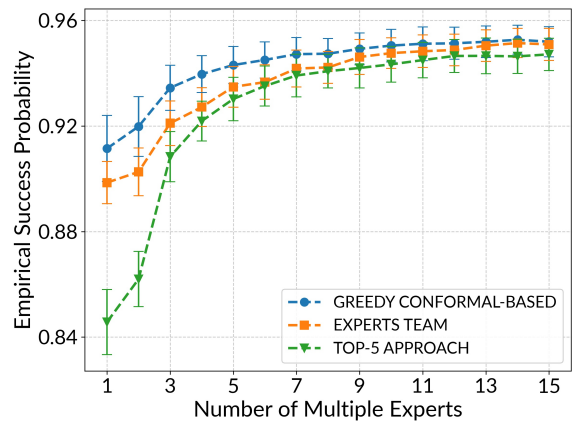


**Figure 7: Empirical success probability of our proposed greedy selection algorithm in comparison with the expert team and top-$5$ set predictor framework for different number of experts $h$ on the ImageNet-16H dataset. The values are averaged over 100 runs. The empirical average set size of the conformal sets for all values of $h$ is 2.8247. The size of the calibration and estimation sets is $l = 240$. The error bars indicate one standard deviation from the mean.**

should examine scenarios where a clear trade-off exists between the tolerance level $\alpha$ and average conformal set sizes.

**Empirical results.** The experimental results show that the proposed greedy selection of human subsets is non-trivial, highlighting its advantages over majority decisions or random expert choices. However, our experiments assume well-distributed calibration data, indicating a need for further exploration in scenarios with multiple classes and class imbalances. Additionally, we treat all experts as equally important, though some may offer greater value. Future research should examine settings where experts have varying levels of expertise and assess the likelihood of including high-importance experts in the selected human subset.

**Broader impact.** The greedy subset selection of human experts enables a realistic human-AI collaboration, where each expert chooses from a refined set of options. Future work could enhance this framework by accurately modeling experts' class preferences and interdependencies, as well as exploring alternative collaboration methods, such as a scoring system instead of simple subset selection.

## 8 CONCLUSION

We have looked at the challenge of human-AI collaboration among multiple experts through the framework of subset selection. Our theoretical analysis demonstrates the conditions under which selecting multiple experts from conformal sets is more advantageous than choosing from the entire label space. Inspired by these findings, we propose a greedy algorithm for selecting human subsets based on conformal sets to improve classification performance during inference. We have shown that this method outperforms both naive subset selection approaches and greedy strategies based on top-$k$ prediction sets.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Anastasios Nikolas Angelopoulos and Stephen Bates. 2021. A Gentle Introduction to Conformal Prediction and Distribution-Free Uncertainty Quantification. *ArXiv* abs/2107.07511 (2021).

[2] Anastasios Nikolas Angelopoulos and Stephen Bates. 2023. Conformal Prediction: A Gentle Introduction. *Found. Trends Mach. Learn.* 16 (2023), 494–591.

[3] Varun Babbar, Umang Bhatt, and Adrian Weller. 2022. On the Utility of Prediction Sets in Human-AI Teams. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, Lud De Raedt (Ed.). International Joint Conferences on Artificial Intelligence Organization, 2457–2463. https://doi.org/10.24963/ijcai.2022/341 Main Track.

[4] 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*. https://api.semanticscholar.org/CorpusID:231691269

[5] Mohammad-Amin Charusaie, Hussein Mozannar, David A. Sontag, and Samira Samadi. 2022. Sample Efficient Learning of Predictors that Complement Humans. In *International Conference on Machine Learning*. https://api.semanticscholar.org/CorpusID:250340781

[6] Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, and Titouan Lorieul. 2021. Set-valued classification - overview via a unified framework. *ArXiv* abs/2102.12318 (2021). https://api.semanticscholar.org/CorpusID:232035498

[7] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. 2016. Learning with Rejection. In *International Conference on Algorithmic Learning Theory*.

[8] Jesse C. Cresswell, Yi Sui, Bhargava Kumar, and Noël Vouitsis. 2024. Conformal Prediction Sets Improve Human Decision Making. *ArXiv* abs/2401.13744 (2024). https://api.semanticscholar.org/CorpusID:267211902

[9] Abir De, Nastaran Okati, Ali Zarezade, and Manuel Gomez-Rodriguez. 2020. Classification Under Human Assistance. *ArXiv* abs/2006.11845 (2020). https://api.semanticscholar.org/CorpusID:219966931

[10] Davide Dell'Anna, Pradeep K. Murukannaiah, Bernd Dudzik, Davide Grossi, Catholijn M. Jonker, Catharine Oertel, and Pinar Yolum. 2024. Toward a Quality Model for Hybrid Intelligence Teams. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2024, Auckland, New Zealand, May 6-10, 2024*, Mehdi Dastani, Jaime Simão Sichman, Natasha Alechina, and Virginia Dignum (Eds.). International Foundation for Autonomous Agents and Multiagent Systems / ACM, 434–443. https://doi.org/10.5555/3635637.3662893

[11] Krishnamurthy Dvijotham, Jim Winkens, Melih Barsbey, Sumedh Ghaisas, Robert Stanforth, Nick Pawlowski, Patricia Strachan, Zahra Ahmed, Shekoofeh Azizi, Yoram Bachrach, Laura Culp, Mayank Daswani, Jana von Freyberg, Christopher J. Kelly, Atilla P. Kiraly, Timo Kohlberger, Scott Mayer McKinney, Basil Mustafa, Vivek Natarajan, Krzysztof J. Geras, Jan Sylwester Witowski, Zhi Zhen Qin, Jacob Creswell, Shravya Shetty, Marcin Sieniek, Terry Spitz, Greg C. Corrado, Pushmeet Kohli, taylan. cemgil, and Alan Karthikesalingam. 2023. Enhancing the reliability and accuracy of AI-enabled diagnosis via complementarity-driven deferral to clinicians. *Nature Medicine* 29 (2023), 1814–1820.

[12] 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*. https://api.semanticscholar.org/CorpusID:235166375

[13] Nina Grgic-Hlaca, Christoph Engel, and Krishna P. Gummadi. 2019. Human Decision Making with Machine Assistance. *Proceedings of the ACM on Human-Computer Interaction* 3 (2019), 1 – 25. https://api.semanticscholar.org/CorpusID:213944226

[14] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), 770–778. https://api.semanticscholar.org/CorpusID:206594692

[15] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. 2016. Identity Mappings in Deep Residual Networks. In *European Conference on Computer Vision*. https://api.semanticscholar.org/CorpusID:6447277

[16] Patrick Hemmer, Sebastian Schellhammer, Michael Vossing, Johannes Jakubik, and Gerhard Satzger. 2022. Forming Effective Human-AI Teams: Building Machine Learning Models that Complement the Capabilities of Multiple Experts. In *International Joint Conference on Artificial Intelligence*.

[17] Patrick Hemmer, Lukas Thede, Michael Vossing, Johannes Jakubik, and Niklas Kuhl. 2023. Learning to Defer with Limited Expert Predictions. *ArXiv* abs/2304.07306 (2023).

[18] Dan Hendrycks and Kevin Gimpel. 2016. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. *ArXiv* abs/1610.02136 (2016). https://api.semanticscholar.org/CorpusID:13046179

[19] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. 2016. Densely Connected Convolutional Networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 2261–2269. https://api.semanticscholar.org/CorpusID:9433631

[20] Lujain Ibrahim, Mohammad M. Ghassemi, and Tuka Alhanai. 2023. Do Explanations Improve the Quality of AI-assisted Human Decisions? An Algorithm-in-the-Loop Analysis of Factual & Counterfactual Explanations. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2023, London, United Kingdom, 29 May 2023 - 2 June 2023*, Noa Agmon, Bo An, Alessandro Ricci, and William Yeoh (Eds.). ACM, 326–334. https://doi.org/10.5555/3545946.3598654

[21] Wei Jiao, Gurnit Atwal, Paz Polak, Rosa Karlić, Edwin Cuppen, Fatima Gurnit Peter J. Andrew V. Paul C. Peter J. David K Al-Shahrour Atwal Bailey Biankin Boutros Campbell, Fátima Al-Shahrour, Gurnit Atwal, Peter J. Bailey, Andrew V. Biankin, Paul C. Boutros, Peter J. Campbell, David K. Chang, Susanna L. Cooke, Vikram Deshpande, Bishoy Morris Faltas, William C. Faquin, Levi A. Garraway, Gaddy Getz, Sean M. Grimmond, Syeda Rabab Zehra Haider, Katherine A. Hoadley, Wei Jiao, Vera B. Kaiser, Rosa Karlić, Mamoru Kato, Kirsten Kübler, Alexander J. Lazar, Constance H. Li, David N. Louis, Adam A. Margolin, Sancha Martin, Hardeep Nahal-Bose, G. Petur Nielsen, Serena Nik-Zainal, Larsson Omberg, Christine P'ng, Marc D. Perry, Paz Polak, Esther Rheinbay, Mark A Rubin, Colin A. Semple, Dennis C. Sgroi, Tatsuhiro Shibata, Reiner Siebert, Jaclyn Smith, Lincoln D. Stein, Miranda D. Stobbe, Ren X. Sun, Kevin Thai, Derek W. Wright, Chin-Lee Wu, Ke Yuan, Junjun Zhang, Alexandra Danyi, Jeroen de Ridder, Carla van Herpen, Martijn Paul Lolkema, Neeltje Steeghs, Gaddy Getz, Quaid D. Morris, and Lincoln D. Stein. 2020. A deep learning system accurately classifies primary and metastatic cancers using passenger mutation patterns. *Nature Communications* 11 (2020). https://api.semanticscholar.org/CorpusID:211038638

[22] Subbarao Kambhampati. 2019. Synthesizing Explainable Behavior for Human-AI Collaboration. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS '19, Montreal, QC, Canada, May 13-17, 2019*, Edith Elkind, Manuela Veloso, Noa Agmon, and Matthew E. Taylor (Eds.). International Foundation for Autonomous Agents and Multiagent Systems, 1–2. http://dl.acm.org/citation.cfm?id=3331663

[23] Gavin Kerrigan, Padhraic Smyth, and Mark Steyvers. 2021. Combining Human Predictions with Model Probabilities via Confusion Matrices and Calibration. *ArXiv* abs/2109.14591 (2021). https://api.semanticscholar.org/CorpusID:238215419

[24] Vijay Keswani, Matthew Lease, and Krishnaram Kenthapadi. 2021. Towards Unbiased and Accurate Deferral to Multiple Experts. *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society* (2021). https://api.semanticscholar.org/CorpusID:232046150

[25] Alex Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images. https://api.semanticscholar.org/CorpusID:18268744

[26] Jijia Liu, Chao Yu, Jiaxuan Gao, Yuqing Xie, Qingmin Liao, Yi Wu, and Yu Wang. 2024. LLM-Powered Hierarchical Language Agent for Real-time Human-AI Coordination. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2024, Auckland, New Zealand, May 6-10, 2024*, Mehdi Dastani, Jaime Simão Sichman, Natasha Alechina, and Virginia Dignum (Eds.). International Foundation for Autonomous Agents and Multiagent Systems / ACM, 1219–1228. https://doi.org/10.5555/3635637.3662979

[27] Ruishan Liu, Shemra Rizzo, Sam Whipple, Navdeep Pal, Arturo López Pineda, Michael Lu, Brandon Arnieri, Ying Lu, William B. Capra, Ryan Copping, and James Zou. 2021. Evaluating eligibility criteria of oncology trials using real-world data and AI. *Nature* 592 (2021), 629 – 633. https://api.semanticscholar.org/CorpusID:233183554

[28] Liyao Ma and Thierry Denoeux. 2021. Partial classification in the belief function framework. *Knowl. Based Syst.* 214 (2021), 106742. https://api.semanticscholar.org/CorpusID:232023309

[29] Anqi Mao, Christopher Mohri, Mehryar Mohri, and Yutao Zhong. 2023. Two-Stage Learning to Defer with Multiple Experts. In *Neural Information Processing Systems*.

[30] Anqi Mao, Mehryar Mohri, and Yutao Zhong. 2023. Principled Approaches for Learning to Defer with Multiple Experts. *ArXiv* abs/2310.14774 (2023).

[31] Siddharth Mehrotra. 2021. Modelling Trust in Human-AI Interaction. In *AAMAS '21: 20th International Conference on Autonomous Agents and Multiagent Systems, Virtual Event, United Kingdom, May 3-7, 2021*, Frank Dignum, Alessio Lomuscio, Ulle Endriss, and Ann Nowé (Eds.). ACM, 1826–1828. https://doi.org/10.5555/3463952.3464253

[32] Thomas Mortier, Marek Wydmuch, Krzysztof Dembczynski, Eyke Hüllermeier, and Willem Waegeman. 2019. Efficient set-valued prediction in multi-class classification. *Data Mining and Knowledge Discovery* 35 (2019), 1435 – 1469. https://api.semanticscholar.org/CorpusID:219753835

[33] 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. In *International Conference on Artificial Intelligence and Statistics*. https://api.semanticscholar.org/CorpusID:255941521

[34] Hussein Mozannar and David A. Sontag. 2020. Consistent Estimators for Learning to Defer to an Expert. In *International Conference on Machine Learning*.

[35] Vu-Linh Nguyen and Eyke Hüllermeier. 2021. Multilabel Classification with Partial Abstention: Bayes-Optimal Prediction under Label Independence. *J. Artif. Intell. Res.* 72 (2021), 613–665. https://api.semanticscholar.org/CorpusID:242398874

[36] Nastaran Okati, Abir De, and Manuel Gomez-Rodriguez. 2021. Differentiable Learning Under Triage. *ArXiv* abs/2103.08902 (2021). https://api.semanticscholar.org/CorpusID:232240535

[37] Andrea Papenmeier, Gwenn Englebienne, and Christin Seifert. 2019. How model accuracy and explanation fidelity influence user trust. *ArXiv* abs/1907.12652 (2019). https://api.semanticscholar.org/CorpusID:198985791

[38] Joshua C. Peterson, Ruairidh M. Battleday, Thomas L. Griffiths, and Olga Russakovsky. 2019. Human Uncertainty Makes Classification More Robust. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019), 9616–9625. https://api.semanticscholar.org/CorpusID:201103726

[39] 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 (2019). https://api.semanticscholar.org/CorpusID:88524064

[40] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. 2014. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* 115 (2014), 211 – 252. https://api.semanticscholar.org/CorpusID:2930547

[41] Sagalpreet Singh, Shweta Jain, and Shashi Shekhar Jha. 2023. On Subset Selection of Multiple Humans To Improve Human-AI Team Accuracy. In *Adaptive Agents and Multi-Agent Systems*.

[42] Mark Steyvers, Heliodoro Tejeda, Gavin Kerrigan, and Padhraic Smyth. 2022. Bayesian modeling of human–AI complementarity. *Proceedings of the National Academy of Sciences of the United States of America* 119 (2022).

[43] Eleni Straitouri, Lequn Wang, Nastaran Okati, and Manuel Gomez Rodriguez. 2023. Improving Expert Predictions with Conformal Prediction. In *International Conference on Machine Learning*. https://api.semanticscholar.org/CorpusID:259309473

[44] Harini Suresh, Natalie Lao, and Ilaria Liccardi. 2020. Misplaced Trust: Measuring the Interference of Machine Learning in Human Decision-Making. *Proceedings of the 12th ACM Conference on Web Science* (2020). https://api.semanticscholar.org/CorpusID:218863174

[45] G. D. Toni, Nastaran Okati, Suhas Thejaswi, Eleni Straitouri, and Manuel Gomez-Rodriguez. 2024. Towards Human-AI Complementarity with Predictions Sets. *ArXiv* abs/2405.17544 (2024). https://api.semanticscholar.org/CorpusID:270067597

[46] Rajeev Verma and Eric Nalisnick. 2022. Calibrated Learning to Defer with One-vs-All Classifiers. In *International Conference on Machine Learning*.

[47] Kailas Vodrahalli, Tobias Gerstenberg, and James Zou. 2022. Uncalibrated Models Can Improve Human-AI Collaboration. *ArXiv* abs/2202.05983 (2022). https://api.semanticscholar.org/CorpusID:246823481

[48] Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. 2005. Algorithmic Learning in a Random World.

[49] Xinru Wang and Ming Yin. 2021. Are Explanations Helpful? A Comparative Study of the Effects of Explanations in AI-Assisted Decision-Making. *Proceedings of the 26th International Conference on Intelligent User Interfaces* (2021). https://api.semanticscholar.org/CorpusID:233224125

[50] Ming Yin, Jennifer Wortman Vaughan, and Hanna M. Wallach. 2019. Understanding the Effect of Accuracy on Trust in Machine Learning Models. *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (2019). https://api.semanticscholar.org/CorpusID:109927933

[51] Dongping Zhang, Angelos Chatzimparmpas, Negar Kamali, and Jessica R. Hullman. 2024. Evaluating the Utility of Conformal Prediction Sets for AI-Advised Image Labeling. *Proceedings of the CHI Conference on Human Factors in Computing Systems* (2024). https://api.semanticscholar.org/CorpusID:267027836

[52] Yunfeng Zhang, Qingzi Vera Liao, and Rachel K. E. Bellamy. 2020. Effect of confidence and explanation on accuracy and trust calibration in AI-assisted decision making. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (2020). https://api.semanticscholar.org/CorpusID:210023849

# A PROOFS

## A.1 Proof of Lemma 1

LEMMA 1. Given $h$ human expert predictions $H = [p_1, p_2, \ldots, p_h]$, conformal set $C$ with tolerance level $\alpha$, and the $i$th expert's conditional success probability $\mathscr{C}^i$, the lower bound on the accuracy of the combined framework is given as

$$\mathbb{E}[\mathbb{1}(\{\hat{Y} = y\} \cap \{y \in C\})] \geq \mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} > 1] \cdot (1 - \alpha) \tag{7}$$

PROOF:

$$\begin{aligned}
\mathbb{E}[\mathbb{1}(\{\hat{Y} = y\} \cap \{y \in C\})] &= \mathbb{P}[\hat{Y} = y | H, y \in C] \cdot \mathbb{P}[y \in C | H] \\
&\geq \mathbb{P}[\hat{Y} = y | H, y \in C] \cdot (1 - \alpha) \\
&= \mathbb{P}[y = \arg\max_{k \in C} \prod_{i \in [h]} \mathscr{C}^i_{kp_i}] \cdot (1 - \alpha) \\
&= \mathbb{P}[\prod_{i \in [h]} \mathscr{C}^i_{yp_i} > \max_{k \neq y, k \in C} \prod_{i \in [h]} \mathscr{C}^i_{kp_i}] \cdot (1 - \alpha) \\
&\geq \mathbb{P}[\prod_{i \in [h]} \mathscr{C}^i_{yp_i} > \prod_{i \in [h]} \max_{k \neq y, k \in C} \mathscr{C}^i_{kp_i}] \cdot (1 - \alpha) \\
&\geq \mathbb{P}[\prod_{i \in [h]} \mathscr{C}^i_{yp_i} > \prod_{i \in [h]} (1 - \mathscr{C}^i_{yp_i})] \cdot (1 - \alpha) \\
&= \mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} > 1] \cdot (1 - \alpha)
\end{aligned}$$

## A.2 Proof of Lemma 2

LEMMA 2. Let $\epsilon = \prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} - \prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}}$. Assuming $\epsilon > 1$ and $\alpha = 0$, the combination framework of multiple human expert predictions that utilizes a conformal set achieves a tighter lower bound on accuracy,

$$\mathbb{P}[\prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} > \frac{1 - \hat{f}_y}{\hat{f}_y}] \leq \mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} > 1]. \tag{8}$$

PROOF: To prove this result, we start from the right side of the inequality.

$$\begin{aligned}
\mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} > 1] &= \mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} - 1 > 0] \\
&= \mathbb{P}[\prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} - 1 + \prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} - \prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} > 0] \\
&= \mathbb{P}[\prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} + \prod_{i \in [h]} \frac{\mathscr{C}^i_{yp_i}}{1 - \mathscr{C}^i_{yp_i}} - \prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} - 1 > 0] \\
&= \mathbb{P}[\prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} + \epsilon - 1 > 0] \\
&\geq \mathbb{P}[\prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} > 0] \text{ since we assume } \epsilon > 1 \\
&\geq \mathbb{P}[\prod_{i \in [h]} \frac{C^i_{yp_i}}{1 - C^i_{yp_i}} > \frac{1 - \hat{f}_y}{\hat{f}_y}]
\end{aligned}$$

# B    SIGNIFICANCE VALUES FOR THE EXPERIMENTS

## B.1    CIFAR-10H Dataset

We show significance values for the accuracy differences in our experiments, comparing our greedy method to ALL HUMANS method using the CIFAR-10H dataset. We used either paired t-test or Wilcoxon signed-rank test (based on data normality).

Table 2: Significance values for the CIFAR-10H experiments with 3 multiple experts comparing the accuracy of our proposed method with the ALL HUMANS approach in Figure 2.

| Calib. data (in %) | p-value | Conclusion |
|---|---|---|
| 10 | 0.002 | Our method is statistically significantly better |
| 15 | 0.002 | Our method is statistically significantly better |
| 20 | 0.004 | Our method is statistically significantly better |
| 25 | 0.002 | Our method is statistically significantly better |
| 30 | 0.002 | Our method is statistically significantly better |

Table 3: Significance values for the CIFAR-10H experiments with 5 multiple experts comparing the accuracy of our proposed method with the ALL HUMANS approach in Figure 2.

| Calib. data (in %) | p-value | Conclusion |
|---|---|---|
| 10 | 0.002 | Our method is statistically significantly better |
| 15 | 0.002 | Our method is statistically significantly better |
| 20 | 0.008 | Our method is statistically significantly better |
| 25 | 0.002 | Our method is statistically significantly better |
| 30 | <0.001 | Our method is statistically significantly better |

In the table that follows, we show significance values for the differences in accuracy in our CIFAR-10H experiments, comparing our proposed greedy method to the top-*6* and top-*10* approaches. We used either paired t-test or Wilcoxon signed-rank test (depending on data normality).

Table 4: Significance values for the CIFAR-10H experiments with 5 multiple experts comparing the accuracy of our proposed method with the top-*6* and top-*10* approach in Figure 3.

| | p-value | Conclusion |
|---|---|---|
| ours vs top-6 | <0.001 | Our method is statistically significantly better |
| ours vs top-10 | <0.001 | Our method is statistically significantly better |

## B.2    ImageNet-16H Dataset

We now show significance values for the accuracy differences in our experiments, comparing our greedy method to the ALL HUMANS method using the ImageNet-16H dataset. We used either paired t-test or Wilcoxon signed-rank test (based on data normality). We observe that with relatively small calibration data, our proposed approach is not statistically significantly better than the ALL HUMANS approach. This highlights the necessity of having sufficient calibration when using our proposed greedy method.

Table 5: Significance values for the ImageNet-16H experiments with 5 multiple experts comparing the accuracy of our proposed method with the ALL HUMANS approach in Figure 4.

| Calib. data (in %) | p-value | Conclusion |
|---|---|---|
| 15 | 0.926 | No significant difference |
| 20 | <0.001 | Our method is statistically significantly better |
| 25 | <0.001 | Our method is statistically significantly better |
| 30 | <0.001 | Our method is statistically significantly better |

In the following tables, we show significance values for the differences in accuracy in our ImageNet-16H experiments, comparing our proposed greedy method to the top-4, top-8, and top-16 approaches. We used either paired t-test or Wilcoxon signed-rank test (depending on data normality).

**Table 6: Significance values for the ImageNet-16H experiments with 3 multiple experts comparing the accuracy of our proposed method with the top-4, top-8, and top-16 approach in Figure 5.**

|  | p-value | Conclusion |
| --- | --- | --- |
| ours vs top-4 | <0.001 | Our method is statistically significantly better |
| ours vs top-8 | <0.001 | Our method is statistically significantly better |
| ours vs top-16 | <0.001 | Our method is statistically significantly better |

**Table 7: Significance values for the ImageNet-16H experiments with 5 multiple experts comparing the accuracy of our proposed method with the top-4, top-8, and top-16 approach in Figure 5.**

|  | p-value | Conclusion |
| --- | --- | --- |
| ours vs top-4 | <0.001 | Our method is statistically significantly better |
| ours vs top-8 | <0.001 | Our method is statistically significantly better |
| ours vs top-16 | <0.001 | Our method is statistically significantly better |

## C   IMPLEMENTATION DETAILS

To conduct our experiments, we utilize PyTorch 2.1.0, NumPy 1.26.4, and Scikit-learn 1.0.2 in Python 3.9.19. To ensure reproducibility, we kept a fixed random seed across all random processes for each run.

**CIFAR-10H.** For each image in this dataset, there are approximately 50 expert predictions wherein each image can be categorized in one of the 10 classes.

**Imagenet-16H.** This dataset comprises only 1200 samples, with a portion allocated for model training. The classes are chair, oven, knife, bottle, keyboard, clock, boat, bicycle, airplane, truck, car, elephant, bear, dog, cat, and bird. The original labels from the ILSRVR database are used as ground truth labels. To compute the confusion matrix for the humans, we utilize the empirical expert distribution corresponding to the estimation and calibration data, which are not used for testing. We simulate multiple humans for testing by sampling from the expert distribution for test data. This distribution represents a categorical distribution for each class provided by the human experts.

**Execution Time and Memory Consumption** For our experiments, we evaluate our proposed approach on the CIFAR-10H and ImageNet-16H datasets using an NVIDIA A100 GPU accelerator. Our experimental setup utilizes CUDA version 11.8 to accelerate the computationally intensive workloads involved in training and evaluating our models on these datasets. For the CIFAR-10H dataset, each experimental run had modest memory requirements, consuming only a few megabytes, and completed swiftly within seconds. In contrast, the ImageNet-16H dataset posed a more computational challenge, requiring several minutes to complete.

## D   TESTING DIFFERENT VALUES OF THE TOLERANCE LEVEL $\alpha$

We show results for the proposed greedy algorithm for different $\alpha$ values in the CIFAR-10H dataset experiments, using 5 experts and a calibration proportion of 20%. Setting $\alpha$ involves balancing desired coverage with implications on prediction set sizes. Our findings show that increasing $\alpha$ leads to a smaller average set size and a reduced empirical success probability. Since the reduction in the average set size is not drastic (an acceptable trade-off), it is reasonable to select a small $\alpha$ value in this setting.

**Table 8: Mean set size and accuracy of the proposed greedy algorithm for different values of $\alpha$ in the CIFAR-10H dataset.**

| $\alpha$ (in %) | Mean set size | Accuracy (in %) |
| --- | --- | --- |
| 1 | 1.17 | 98.48 |
| 3 | 1.01 | 97.00 |
| 5 | 0.97 | 94.93 |
| 7 | 0.94 | 92.92 |

# E ADDITIONAL REMARKS ON THE SIGNIFICANCE AND LIMITATIONS OF THE STUDY

The accuracy improvement of the proposed framework is crucial for safety-critical applications, such as classification in medical AI or autonomous driving, where tolerance to error is minimal. A limitation of our approach occurs when the calibration data is small. In CIFAR-10H data, when calibration data is 1% (100 samples) for h=5, significance testing using paired t-test shows that our greedy framework and ALL HUMANS do not differ significantly in accuracy (p-value=0.09). In addition, we did not test our method beyond image classification due to the lack of publicly available datasets with multiple expert predictions per sample and more than 3 classes. Datasets like Hatespeech, COMPASS, and NIH Chest X-ray have real human predictions but are limited to 2 or 3 classes, restricting the analysis of the conformal predictor.