

# MULTI-DIMENSIONAL CONFORMAL PREDICTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Conformal prediction has attracted significant attention as a distribution-free method for uncertainty quantification in black-box models, providing prediction sets with guaranteed coverage. However, its practical utility is often limited when these prediction sets become excessively large, reducing its overall effectiveness. In this paper, we introduce a novel approach to conformal prediction for classification problems, which leverages a multi-dimensional nonconformity score. By extending standard conformal prediction to higher dimensions, we achieve better separation between correct and incorrect labels. Utilizing this we can focus on regions with low concentrations of incorrect labels, leading to smaller, more informative prediction sets. To efficiently generate the multi-dimensional score, we employ a self-ensembling technique that trains multiple diverse classification heads on top of a backbone model. We demonstrate the advantage of our approach compared to baselines across different benchmarks.

## 1 INTRODUCTION

Deep learning models become increasingly dominant in almost every domain, ranging from computer vision and natural language processing to speech recognition. However, as deep learning models are deployed in safety-critical applications, such as healthcare Lambert et al. (2024) and autonomous driving Muhammad et al. (2020), it is important to certify their reliability and safety. This highlights the need for robust uncertainty quantification methods that can determine when models are uncertain about their predictions and suggest alternative estimates. Conformal prediction offers a powerful, distribution-free, and model-agnostic framework for uncertainty quantification, providing finite-sample guarantees Vovk et al. (2015); Lei et al. (2013); Barber et al. (2021); Angelopoulos et al. (2020). Its core principle is to transform point predictions from any model into prediction sets or intervals that contain the true value with high probability.

Conformal prediction relies on a nonconformity score that quantifies how unusual or atypical a new input-label pair is relative to the given data. While conformal prediction guarantees valid coverage for any model and data distribution, its practical effectiveness is influenced by both the performance of the model and the choice of the nonconformity score Romano et al. (2020); Angelopoulos et al. (2020). The efficiency of the resulting prediction sets is typically measured by their size, with smaller sets leading to more informative and precise predictions. Consequently, there is active research aimed at developing methods that produce the most efficient and informative prediction sets. This includes proposing new conformity scores Sadinle et al. (2019); Romano et al. (2020); Angelopoulos et al. (2020); Huang et al. (2024); Luo & Zhou (2024a), training models using loss functions that promote efficiency Stutz et al. (2021); Einbinder et al. (2022), and combining different models, scores, or data augmentations Bai et al. (2021); Luo & Zhou (2024b); Lu (2023).

The common approach across existing methods is to optimize a single nonconformity score, while the calibration process—generally involving the computation of a threshold based on a quantile of the calibration scores—remains unchanged. In this paper, we present a novel approach by extending the standard conformal prediction framework from a one-dimensional score to a higher-dimensional space defined by multiple nonconformity scores. The intuition behind this is that higher-dimensional spaces can better separate correct from incorrect labels, potentially leading to more efficient prediction sets with fewer false labels. However, selecting a region in this higher-dimensional space that guarantees exact coverage while optimizing efficiency is non-trivial, as there are infinitely many ways to partition the space. To tackle this challenge, we propose a simple yet effective method that splits the multi-dimensional score space into cells, with cell centers defined by the calibration sam-

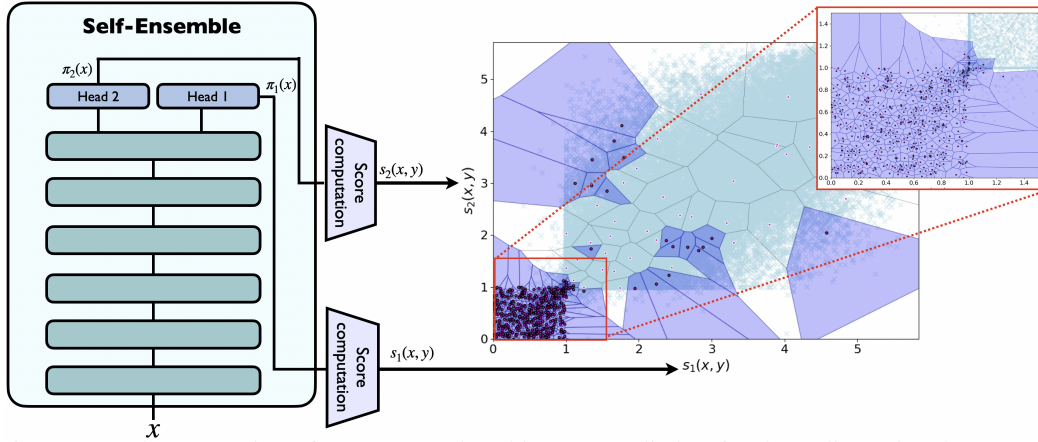


Figure 1: A demonstration of our proposed multi-score prediction for the 2-dimensional case. On the left, our proposed self-ensemble model with 2 classification heads. We compute a nonconformity score for each head. On the right, the selected regions, where circles represent scores of true labels, and light-blue x-marks represent scores of false labels. Selected cells are colored in blue and their centers have black edge color. We see that cells with low number of false labels are chosen.

ples, as illustrated in Fig. 1. The cells are ranked based on increasing ratios of incorrect to correct labels they contain, and we select the top-ranked cells that meet the desired coverage. At test time, the prediction set consists of all labels with scores falling within the selected region. Additionally, we introduce a flexible self-ensembling technique that constructs a multi-dimensional score by combining predictions from multiple diverse classification heads built on top of a single backbone model. We provide theoretical guarantees that this high-dimensional region selection maintains valid coverage, and demonstrate that it offers superior efficiency compared to baseline methods across various settings.

Our main contributions can be summarized as follows:

1. Propose a new multi-dimensional conformal prediction framework that can better identify regions with a large amount of true labels and a small amount of false labels. Our approach is parameter-free, requires no optimization procedures, and is not restricted to a specific coverage level.
2. Theoretically show that this high-dimensional selection procedure maintains the desired coverage with finite-sample guarantees.
3. Present a flexible and cheap self-ensembling approach to obtain multi-dimensional nonconformity scores by training multiple classification heads, while encouraging diversity.
4. Our experimental results demonstrate the superiority of the proposed method over competing baselines, consistently producing smaller and more efficient prediction sets.

## 2 RELATED WORK

**Enhanced nonconformity scores.** Improving the efficiency of conformal prediction has been a central focus of recent research. Several studies have proposed enhanced nonconformity scores aimed at reducing the size of the prediction sets or improving conditional coverage (Sadinle et al., 2019; Romano et al., 2020; Angelopoulos et al., 2020; Huang et al., 2024; Luo & Zhou, 2024a). Another approach involves performing conformal prediction in the feature space, then mapping the intervals from the embedding space back to the output space (Teng et al., 2022). Our approach addresses the use of multiple nonconformity scores, as previous research has demonstrated that combining multiple scores can be more efficient than relying on a single score (Yang et al., 2023b). We introduce a practical method for generating multiple scores without the need for training multiple models or performing additional inference steps. This approach can be applied to any base score and is orthogonal to the advancements in the developing improved nonconformity scores.

**Improving conoformal prediction via training.** Although conformal prediction is usually considered as a wrapper around black-box models, recent approaches suggested to directly train

models to improve conformal prediction efficiency. Stutz et al. (2021) introduced a differentiable conformal prediction pipeline that optimizes the size of prediction sets. In (Einbinder et al., 2022), a regularization term was added to the training loss, encouraging the distribution of the nonconformity scores to match a uniform distribution. Additionally, Bai et al. (2021) explored optimizing conformal prediction within broader function classes. These methods rely on differentiable approximations that may not fully align with the actual target objective, and are often designed for a specific coverage level, requiring retraining for each new level. In contrast, our approach avoids optimization altogether and offers greater flexibility by being independent of the coverage level and the base nonconformity score.

**Combining nonconformity scores.** Several works explore conformal prediction in the context of model fusion and combining nonconformity scores. Ensemble learning methods, which train multiple models on different subsets of the data, have been widely applied to conformal prediction (Linusson et al., 2020). Notable approaches include cross-conformal predictors (Vovk, 2015), bootstrap conformal predictors (Vovk, 2015), and out-of-bag calibrated conformal predictors (Devetyarov & Nourtdinov, 2010). Other methods for combining nonconformity scores have also been investigated. For example, Luo & Zhou (2024b) proposed using a weighted average of multiple scores derived from the same output, with weights learned through optimization. Similarly, Lu (2023) suggested combining nonconformity scores obtained via test-time augmentations of the same image. A more recent work explored the aggregation of multiple prediction sets, assuming no direct access to the underlying scores (Gasparin & Ramdas, 2024). In contrast to existing approaches that focus on combining nonconformity scores via weighted aggregation or majority voting, we propose a general framework that identifies promising regions with low concentrations of false labels in the multi-score space.

**Ensemble methods.** Model ensembles have been shown to enhance various metrics for uncertainty quantification beyond conformal prediction efficiency, such as calibration error (Hansen & Salamon, 1990; Lakshminarayanan et al., 2017). However, ensembles are often considered computationally expensive, as they require training and deploying multiple independent models. To mitigate this, Qendro et al. (2021) proposed an early-exit ensemble, which leverages multiple prediction heads from intermediate layers to improve uncertainty quantification while maintaining a single model. This approach has been shown to enhance computational efficiency (Cai et al., 2020) and improve adversarial robustness (Qendro & Mascolo, 2022). Building on these insights, we propose a self-ensemble model with multiple classification heads to generate a multi-dimensional nonconformity score without significant additional costs.

### 3 BACKGROUND - CONFORMAL PREDICTION

Let  $X \in \mathcal{X}$  represent an input, associated with a label  $Y \in \mathcal{Y}$ , where  $\mathcal{Y} = \{1, \dots, Q\}$ . Consider a classifier  $\pi(x) \in [0, 1]^Q$  that outputs the probability distribution over  $Q$  classes (e.g., a neural network with a softmax layer producing probabilities for each class). The conformal prediction framework starts by selecting a *nonconformity score*  $s : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{S} \subseteq \mathbb{R}$ , which quantifies the uncertainty of the classifier’s prediction for the pair  $(X, Y)$  with respect to existing data. Given a set of i.i.d. calibration data  $\{(X_i, Y_i)\}_{i=1}^m$ , we can form a prediction set for a new test point  $X_{m+1}$  with a guaranteed coverage level of at least  $1 - \alpha$ , where  $\alpha \in (0, 1)$  is set by the user. The prediction set is defined as:

$$\Gamma_\lambda(X_{m+1}) = \{y \in \mathcal{Y} : s(X_{m+1}, y) \leq \lambda\}, \quad (1)$$

where  $\lambda$  is a threshold determined as a quantile of the calibration nonconformity scores:

$$\lambda := \text{Quantile} \left( \left\lceil \frac{(m+1)(1-\alpha)}{m} \right\rceil; \{s(X_i, Y_i)\}_{i=1}^m \right). \quad (2)$$

This procedure ensures that the probability of the true label  $Y_{m+1}$  being included in the prediction set is at least  $1 - \alpha$ , providing a reliable uncertainty estimate based on a finite-size calibration data, as stated in the following theorem.

**Theorem 1.** (*Conformal calibration coverage guarantee*) Let  $\{(X_i, Y_i)\}_{i=1}^{m+1}$  be exchangeable. For any score function  $s : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{S}$  and any significance level  $\alpha \in (0, 1)$ , define a quantile  $\lambda$  by Eq. (2) and prediction set  $\Gamma_\lambda(X_{m+1})$  by Eq. (1). We have:

$$\mathbb{P}(Y_{m+1} \in \Gamma_\lambda(X_{m+1})) \geq 1 - \alpha \quad (3)$$

A commonly used nonconformity score is based on the confidence of the predicted probabilities, defined as  $s_{\text{Thr}}(x, y) = 1 - \pi(x)_y$  (Sadinle et al., 2019), but can sometimes undercover hard examples and overcover trivial ones. Hence, a popular alternative is the adaptive prediction sets (APS) method (Romano et al., 2020), which is based on the cumulative probability  $s_{\text{APS}}(x, y) := \sum_{q=1}^Q \pi(x)_q \mathbf{1}[\pi(x)_q > \pi(x)_y] + u \cdot \pi(x)_y$  where  $u$  a uniform random value that breaks potential ties between different scores. However, the APS method often results in large prediction sets, which is undesirable. To address this, regularized adaptive prediction sets (RAPS) score was introduced (Angelopoulos et al., 2020), which encourages smaller prediction sets by penalizing less likely labels. The RAPS score is defined as:

$$s_{\text{RAPS}}(x, y) := \sum_{q=1}^Q \pi(x)_q \mathbf{1}[\pi(x)_q > \pi(x)_y] + u \cdot \pi(x)_y + \nu \cdot \max(o(x, y) - \kappa, 0), \quad (4)$$

where  $o(x, y)$  denotes the rank of class  $y$ , and  $\nu$  and  $\kappa$  are hyperparameters that control the size of the penalty. More recently, alternative scoring methods have been proposed that rely on the relative rank of the prediction (Huang et al., 2024; Luo & Zhou, 2024a). For example, the sorted adaptive prediction sets (SAPS) score is defined as (Huang et al., 2024):

$$s_{\text{SAPS}}(x, y) := \begin{cases} u \cdot \pi_{\max}(x), & \text{if } o(x, y) = 1, \\ \pi_{\max}(x) + (o(x, y) - 2 + u) \cdot \xi, & \text{otherwise,} \end{cases} \quad (5)$$

where  $\xi$  is a hyperparameter that controls the weight of the ranking information, and  $\pi_{\max}(x)$  is the maximum softmax probability.

## 4 PROPOSED METHOD

### 4.1 MULTI-SCORE CALIBRATION

We now consider a multi-dimensional nonconformity score  $\mathbf{s} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{S}$ , constructed by concatenating  $n$  individual nonconformity scores as  $\mathbf{s}(x, y) = [s_1(x, y), \dots, s_n(x, y)]^T$ , where  $\mathcal{S} := \mathcal{S}_1 \times \dots \times \mathcal{S}_n \subseteq \mathbb{R}^n$ . Before discussing different approaches to handling multi-dimensional nonconformity scores and presenting our proposed method, we first introduce a toy example to highlight the advantages of using a multi-dimensional score over a single-dimensional one.

**Example 4.1** (Toy setting). Assume a binary classification problem with  $\mathcal{Y} = \{-1, 1\}$  and prior probabilities  $\mathbb{P}(Y = 1)$  and  $\mathbb{P}(Y = -1)$ . The input  $X$  is generated from a mixture of two Gaussians, with  $\mathbb{P}(X|Y) = \mathcal{N}(Y, \sigma_y^2)$ . In this example, we can compute the posterior  $\mathbb{P}(Y|X)$  using Bayes rule. Let  $p(y|x) = \mathbb{P}(Y = y|X = x)$ , we consider the following three classifiers:

$$\pi_0(x) = p(y|x), \quad \pi_1(x) := \begin{cases} p(y|x), & \text{if } x < 0, \\ (\epsilon, 1 - \epsilon), & \text{if } x > 0 \end{cases}, \quad \pi_2(x) := \begin{cases} (\epsilon, 1 - \epsilon), & \text{if } x < 0, \\ \mathbb{P}(y|x), & \text{if } x > 0 \end{cases} \quad (6)$$

where  $\epsilon \sim \mathcal{N}(0, 1)$ . Here,  $\pi_0(x)$  represents the ideal classifier, while  $\pi_1(x)$  and  $\pi_2(x)$  are ideal classifiers over half the range of  $x$ , but uninformative over the remaining half. Let  $s_i(x, y)$  denote the nonconformity score computed over the  $i$ -th classifier, it is clear that performing conformal prediction over  $s_0(x, y)$  would be the most efficient, however, using either  $s_1(x, y)$  or  $s_2(x, y)$  will lead to suboptimal results in the uninformative region. In this case, performing conformal prediction in the 2-dimensional space defined by  $\mathbf{s}(x, y) = [s_1(x, y), s_2(x, y)]^T$  is advantageous as the individual scores provide complementary information for different ranges of the input  $x$ . This can be seen from Fig. 2 (b), where taking both axes into account can help to identify regions with high density of true points versus false points. Setting  $\alpha = 0.1$ , we obtain the following average set sizes: 1.05 for  $s_0(x, y)$ , 1.48 for  $s_1(x, y)$ , 1.47 for  $s_2(x, y)$  and 1.24 for our proposed method. Figure 2 (c) compares the set sizes per  $x$ -domain. As expected,  $s_0(x, y)$  obtains sets of size 1 for the entire range, except for  $x \in [-1.5, 1.5]$ , where the two Gaussians overlap. In contrast,  $s_1(x, y)$  and  $s_2(x, y)$  produce larger proportions of 2-element sets in the noisy regions, on the right for  $s_1(x, y)$  and on the left for  $s_2(x, y)$ . Our method, utilizes both  $s_1(x, y)$  and  $s_2(x, y)$ , obtaining similar behavior to the ideal case, provided by  $s_0(x, y)$ . More details are provided in Appendix A.

In Section §4.3, we discuss how such multi-dimensional scores can be constructed. In standard conformal prediction, where  $n = 1$ , the threshold  $\lambda$  divides the real line into two regions: scores

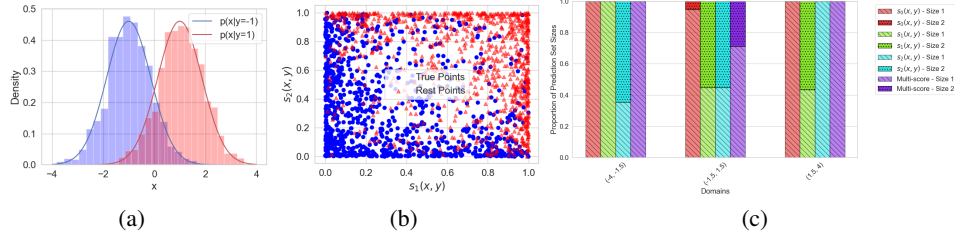


Figure 2: Binary classification example. (a)  $p(x|y)$  for  $y = \{-1, 1\}$ . (b) The 2-dimensional score space defined by the nonconformity scores computed over  $\pi_1(x)$  and  $\pi_2(x)$ . (c) The set-sizes obtained for different domains of  $x$  based on  $s_0(x, y)$ ,  $s_1(x, y)$ ,  $s_2(x, y)$  and our Multi-score method using both  $s_1(x, y)$  and  $s_2(x, y)$ .

less than or equal to  $\lambda$  are included in the prediction set, while scores greater than  $\lambda$  are excluded. This simple thresholding approach is effective because nonconformity scores are expected to be low for true labels, which conform to the patterns in the data, and high for false labels, which deviate from the expected behavior, as explained in §3. When  $n > 1$ , we have a multi-dimensional score. Intuitively, scores that are close to the minimum value in all dimensions correspond to the most conforming labels, while increasing any component of  $s(x, y)$  leads to less conforming scores. However, unlike the one-dimensional case, it is not immediately clear how to optimally partition the score space  $\mathcal{S}$  into regions that should be included or excluded from the prediction set.

A common approach for handling multiple nonconformity scores is to use a weighted sum of the scores:

$$s_w(x, y) = \sum_{i=1}^n w_i s_i(x, y), \quad (7)$$

where the weights  $w_i$  are optimized to minimize the size of the prediction set while maintaining the desired coverage level (Bai et al., 2021; Luo & Zhou, 2024b). Geometrically, this is equivalent to splitting the score space  $\mathcal{S}$  with a hyperplane of the form  $w_1 s_1(x, y) + \dots + w_n s_n(x, y) = \lambda$ , and including only those scores that lie below this hyperplane, i.e. scores satisfying  $w_1 s_1(x, y) + \dots + w_n s_n(x, y) \leq \lambda$ . However, this method imposes a rigid structure on the partitioning, which may not align with the optimal regions that yield the smallest possible prediction sets. Moreover, it requires tuning the weights for a specific coverage level, making it less flexible and potentially unsuitable for different values of  $\alpha$ .

We propose a more flexible approach for managing the multi-dimensional score space that eliminates the need for weight optimization. Our method begins by partitioning the calibration data into two disjoint subsets:  $\mathcal{D}_{\text{cal}} = \mathcal{D}_{\text{cells}} \cup \mathcal{D}_{\text{re-cal}}$ , where  $\mathcal{D}_{\text{cells}} = \{(X_i, Y_i)\}_{i=1}^k$  and  $\mathcal{D}_{\text{re-cal}} = \{(X_i, Y_i)\}_{i=k+1}^m$ , where  $m - k = r$ . The subset  $\mathcal{D}_{\text{cells}}$  is used to partition the score space  $\mathcal{S}$  into distinct regions (cells) and to evaluate the quality of each region. Next, the subset  $\mathcal{D}_{\text{re-cal}}$  is utilized for re-calibration, ensuring that the prediction sets achieve the desired coverage level. Thus, our method consists of three main stages: (i) partitioning, (ii) scoring and ranking, and (iii) calibration, as detailed below.

**(i) Partitioning.** We aim to partition the score space  $\mathcal{S}$  into cells and later decide which cells to include in the prediction sets. Using a uniform grid would be computationally expensive and unscalable as the number of dimensions  $n$  increases. Additionally, score distributions are typically uneven, with some regions being densely populated and others sparse, making uniform partitioning inefficient. Instead, we partition  $\mathcal{S}$  into  $k$  cells, centered around the scores  $s(X_1, Y_1), \dots, s(X_k, Y_k)$  corresponding to the samples in  $\mathcal{D}_{\text{cells}}$ . Specifically, each point in  $\mathcal{S}$  is assigned to the closest center in  $\mathcal{D}_{\text{cells}}$ , formally defined as:

$$\mathcal{C}_i = \left\{ \tilde{s} \in \mathcal{S} \mid i = \arg \min_{1 \leq j \leq k} \|\tilde{s} - s(X_j, Y_j)\| \right\}, \quad i \in \{1, \dots, k\}. \quad (8)$$

This way, the cell resolution adapts to the density of the samples, with smaller cells in high-density regions and larger cells in low-density areas. Note also this approach is analogous to standard conformal prediction, where the calibration scores define segments of varying lengths, and the final selected interval is the union of all segments to the left of the computed quantile (2). Thus, our cell partitioning method can be seen as a generalization of the standard partitioning process to higher dimensions.

(ii) **Scoring and ranking.** In the next stage, we assess which cells should be included in the prediction sets. To achieve this, we compute a ratio that reflects the balance between false labels and true labels in each cell:

$$D_i = \frac{\sum_{j=1}^k \sum_{q=1}^Q \mathbf{1}\{\mathbf{s}(X_j, q) \in \mathcal{C}_i\} \cdot \mathbf{1}\{Y_j \neq q\} + 1}{\sum_{j=1}^k \mathbf{1}\{\mathbf{s}(X_j, Y_j) \in \mathcal{C}_i\}}, \quad i \in \{1, \dots, k\}. \quad (9)$$

This ratio represents the relative amount of false to true labels within each cell  $\mathcal{C}_i$ . Here, we take into account the possibility that several true scores may fall in the same cell (though this becomes increasingly rare as the dimensionality  $n$  increases). Thus, in Eq. (9) we normalize by the number of true labels in each cell, which can be greater than 1. We add 1 to the numerator to distinguish between cells with zero false labels and varying numbers of true labels, as the pure ratio would otherwise be zero regardless of the number of true labels.

In order to improve conformal prediction efficiency and obtain smaller set sizes at test time, we would like to prioritize regions with a low false-to-true label ratio while avoiding regions with high false-to-true ratio. Let  $k'$  denote the number of unique cells, where multiple scores at the same point are treated as a single cell. We then rank the sequence of unique cells  $\mathcal{C}_{(1)}, \mathcal{C}_{(2)}, \dots, \mathcal{C}_{(k')}$  according to  $D_i$ , from lowest to highest, i.e.,  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(k')}$ .

(ii) **Calibration.** The final step is selecting the regions that will ensure exact coverage. In this stage, we utilize the re-calibration set  $\mathcal{D}_{\text{re-cal}}$ . We progressively add cells, starting with the lowest amount of false labels, and continuing until the desired coverage is achieved over  $\mathcal{D}_{\text{re-cal}}$ . Formally, we define the selected region  $\mathcal{C}_{\text{in}}^\eta$  as the union of the top-ranked cells up to index  $\eta$ :

$$\mathcal{C}_{\text{in}}^\eta = \bigcup_{i=1}^{\eta} \mathcal{C}_{(i)}.$$

The required  $\eta$  is determined by:

$$\eta^* = \min \{ \eta \in \{1, \dots, k'\} \mid Y_i \in \mathcal{C}_{\text{in}}^\eta \text{ for at least } \lceil (1 - \alpha)(r + 1) \rceil \text{ samples } (X_i, Y_i) \in \mathcal{D}_{\text{re-cal}} \}$$

and  $\mathcal{C}_{\text{in}}^{\eta^*}$  is the final selected region. Note that it is crucial to use a separate set for calibration, rather than re-using  $\mathcal{D}_{\text{cells}}$ . Since the cell-wise scores are computed using  $\mathcal{D}_{\text{cells}}$ , re-selecting cells based on the same data would introduce bias, leading to undercoverage, as we confirmed empirically during initial experiments (see Appendix A for further justification).

During test time, we receive a new test point  $X_{m+1}$  associated with an unknown label  $Y_{m+1}$ . For each potential  $y \in \mathcal{Y}$ , we compute the multi-dimensional score  $\mathbf{s}(X_{m+1}, y)$  and include in the prediction set only the labels that lie in the selected region  $\mathcal{C}_{\text{in}}^{\eta^*}$ . Accordingly, we obtain the following prediction set:

$$\Gamma_{\eta^*}(X_{m+1}) = \left\{ y \in \{1, \dots, Q\} \mid \mathbf{s}(X_{m+1}, y) \in \mathcal{C}_{\text{in}}^{\eta^*} \right\}. \quad (10)$$

Our method, Multi-score conformal prediction, is summarized in Algorithm 1. Unlike the thresholding mechanism used in the single-score case (1) or the hyperplane splitting for weighted scores (7), our approach defines an unstructured selection region in the multi-dimensional score space. This flexible selection allows us to focus on regions with small amount of false labels, optimizing for smaller prediction sets. It is important to highlight that our method differs from vector quantile regression (VQR), which extends traditional quantile regression to multivariate settings (Carlier et al., 2016; Feldman et al., 2023; Rosenberg et al., 2022). While VQR aims to capture the central  $1 - \alpha$  portion of the output distribution, our focus is on flexible region selection that minimizes set size while maintaining coverage.

Although our prediction set construction differs from standard conformal prediction in its form, it still provides coverage guarantees, as stated in the following theorem.

**Theorem 2.** (Multi-score conformal calibration coverage guarantee). *Let  $\mathcal{D}_{\text{cells}} = \{X_i, Y_i\}_{i=1}^k$  and  $\mathcal{D}_{\text{re-cal}} = \{X_i, Y_i\}_{i=k+1}^{k+m}$  be two disjoint datasets, and  $\{(X_i, Y_i)\}_{i=k+1}^{k+m+1}$  is exchangeable. For any multi-dimensional score function  $\mathbf{s} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{S} \subseteq \mathbb{R}^n$  and any significance level  $\alpha \in (0, 1)$ , the prediction set  $\Gamma_{\eta^*}(X_{m+1})$  defined by Eq. (10) satisfies:*

$$\mathbb{P}(Y_{m+1} \in \Gamma_{\eta^*}(X_{m+1})) \geq 1 - \alpha \quad (11)$$



**Algorithm 1** Multi-Score Conformal Prediction

**Definitions:**  $s(x, y)$  is a multi-dimensional score function.  $\mathcal{D}_{\text{cal}}$  is the calibration data of size  $m$ .  $X_{m+1}$  is a new test sample,  $\alpha$  is the miscoverage level,  $k$  is the number of samples for cell-partitioning and scoring, and  $r = m - k$  is the number of samples for re-calibration.

```

1: function MULTI-SCORE-CP( $s(x, y)$ ,  $\mathcal{D}_{\text{cal}}$ ,  $\alpha$ )
2:   Randomly split  $\mathcal{D}_{\text{cal}}$  to  $\mathcal{D}_{\text{cells}} = \{(X_i, Y_i)\}_{i=1}^k$  and  $\mathcal{D}_{\text{re-cal}} = \{(X_i, Y_i)\}_{i=k+1}^m$ 
3:   Compute the scores  $s(X_i, Y_i)$ ,  $i \in \{1, \dots, m\}$ 
4:   Segment  $\mathcal{S}$  into cells  $\{\mathcal{C}_i\}_{i=1}^k$  centered at  $\{(X_i, Y_i)\}_{i=1}^k$ 
5:   Compute  $D_i$ ,  $i = 1, \dots, k$  according to Eq. (9)
6:   Remove duplicate cells  $\mathcal{C}_1, \mathcal{C}_1, \dots, \mathcal{C}_{k'}$ 
7:   Rank the cells  $\mathcal{C}_{(1)}, \mathcal{C}_{(2)}, \dots, \mathcal{C}_{(k')}$  according to  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(k')}$ 
8:    $\eta^* \leftarrow \min\{\eta \in \{1, \dots, k'\} \mid Y_i \in \mathcal{C}_{\text{in}}^\eta \text{ for at least } \lceil (1 - \alpha)(r + 1) \rceil \text{ samples in } \mathcal{D}_{\text{re-cal}}\}$ 
9:    $\mathcal{C}_{\text{in}}^{\eta^*} \leftarrow \bigcup_{i=1}^{\eta^*} \mathcal{C}_{(i)}$ 
10:  return  $\mathcal{C}_{\text{in}}^{\eta^*}$ 
11: function MULTI-SCORE-EVALUATION( $s(x, y)$ ,  $X_{m+1}$ ,  $\mathcal{C}_{\text{in}}^{\eta^*}$ )
12:  Compute the scores  $s(X_{m+1}, y)$ ,  $y \in \{1, \dots, Q\}$ 
13:  Construct the prediction set  $\Gamma_{\eta^*}(X_{m+1}) = \{y \in \{1, \dots, Q\} \mid s(X_{m+1}, y) \in \mathcal{C}_{\text{in}}^{\eta^*}\}$ 
14:  return  $\Gamma_{\eta^*}(X_{m+1})$ 

```

The proof, detailed in Appendix A, is based on defining a mapping function from the multi-dimensional score to the corresponding cell ratio score, defined in Eq. (9). By formulating the predicted set in Eq. (10) using a thresholding operation over this score, we can align our method with the standard one-dimensional conformal prediction procedure. As a direct result, we establish that valid coverage is guaranteed.

## 4.2 ADDITIONAL VARIANTS OF MULTI-SCORE CONFORMAL PREDICTION

In the following, we present two additional variants of our proposed method.

**Jackknife+ Multi-Score Conformal Prediction.** A limitation of our approach is that it uses only part of the data to perform the calibration, which may impact efficiency. This is especially critical if the sample size  $m$  is small. An alternative solution is to adopt a jackknife+ approach (Romano et al., 2020; Barber et al., 2021), which is computationally more expensive but often provides tighter prediction sets. This approach is summarized in Algorithm B.1. The core idea is to consider the entire calibration data but each time exclude the  $i$ th point from the set of centers and from the score computation in Eq. 9. We sweep over all possible labels  $y \in \mathcal{Y}$  and assign it to the closest center while removing the  $i$ th center. We perform a similar assignment for the  $i$ th score  $s(X_i, Y_i)$ . Then, we compare the rank of the chosen cells, and include  $y$  in the final prediction if its rank is smaller than  $(1 - \alpha)(m + 1)$  hold-out ranks evaluated on the true labeled data for every  $i \in \{1, \dots, m\}$ . Note that this establishes that the coverage is  $1 - 2\alpha$ . In practice, it was shown that the obtained coverage is around  $1 - \alpha$  even without compensating by  $\alpha' = \alpha/2$ .

**Soft Multi-Score Conformal Prediction.** We present a more general variant of our approach, where instead of considering a hard assignment of each point into the closest center, we consider a softer assignment to  $b$  nearest neighbors. This variant is summarized in Algorithm B.2. Here we detect the  $b$  nearest neighbors, and include the point in the prediction set if at least half of its neighbors are in the chosen region.

## 4.3 MULTI-SCORE CONSTRUCTION

Multi-dimensional nonconformity scores can be constructed in several ways. One approach involves using an ensemble of different models, where a nonconformity score is computed for each model. However, as is generally the case with ensembling, this requires training and managing multiple models, which can be resource-intensive. An alternative method, proposed by Lu (2023), uses test-time augmentations to generate multiple scores for different augmented versions of the input image. Similarly to ensembling, test-time augmentation requires running multiple inferences, increasing

computational costs. It is also more appropriate for image data, and is not easily adaptable to other types of data. A more efficient option is to compute various types of scores from the output of a single model (Luo & Zhou, 2024b), such as the ones described in §3. While this requires only a single forward pass, its drawback is that all scores are computed based on a single output, thus does not reflect varied viewpoints on the predictive uncertainty.

In principle, our multi-score conformal prediction approach can be applied to any of the multi-dimensional scores mentioned above. However, we propose a new method for obtaining diverse nonconformity scores that represent varied perspectives without increasing computational complexity. To achieve this, we attach multiple classification heads,  $\{\pi_i(x)\}_{i=1}^H$ , to the penultimate layer (the second-to-last layer) of the model. Training these heads using only the cross-entropy (CE) loss may result in highly similar outputs, providing little additional insight into uncertainty estimation. To address this, we follow (Qendro et al., 2021) and introduce a regularization term that promotes diversity among heads by minimizing their similarity. **This is further motivated by the scenario illustrated in Example 4.1, suggesting that we should encourage the classification heads to produce complementary predictions by specializing in specific input domains. Then, combining their scores in the multi-dimensional space, we can offer an improved uncertainty quantification for the entire input space.** Specifically, the heads are trained using the following loss function:

$$\mathcal{L} = \frac{1}{H} \sum_{i=1}^H L_{\text{CE}}(\pi_i(x), y) - \frac{\beta}{H(H-1)} \sum_{i=1}^H \sum_{i \neq j} \text{sim}(\pi_i(x), \pi_j(x)), \quad (12)$$

where  $L_{\text{CE}}(\cdot, \cdot)$  denotes CE loss, and  $\text{sim}(\cdot, \cdot)$  denotes cosine similarity. Here,  $\beta$  is a hyperparameter that controls the relative weight of the diversity regularization term. We set  $\beta = 1$  in our experiments. While (Qendro et al., 2021) proposed generating classification heads from various layers and depths of the network, our empirical findings revealed that, as expected, heads from shallower layers tend to be weak, leading to lower accuracies. As a result, their contribution in the multi-score setting is minimal. By attaching the classification heads to the penultimate layer and applying diversity regularization (12), we were able to achieve robust classification heads with minimal correlation between them.

## 5 EXPERIMENTS

### 5.1 DATASETS AND MODELS

We test our method over three image classification datasets, with varying number of classes and difficulty levels: CIFAR100 (Krizhevsky et al., 2009), Tiny ImageNet (Le & Yang, 2015), and PathM-NIST (Yang et al., 2023a). For all datasets we use a ResNet50 backbone model pretrained on ImageNet. We first attach a single classification head and fine-tune the full model with CE loss. Next, we add additional 6 classification heads ( $H = 7$  heads in total) and train only the classification heads using the loss defined in Eq. (12). Based on the heads’ output probabilities, we compute the nonconformity scores, where we use either RAPS (4) or SAPS (5) as base scores. We refer to Appendix C for further details on datasets, models other nonconformity base scores, and the training procedure.

### 5.2 EVALUATION

We compare the proposed multi-score method to the following baselines:

- **Best single head (Best Head)** - Conformal prediction is applied separately to each classification head using the entire  $\mathcal{D}_{\text{cal}}$  dataset, **and not only over the  $\mathcal{D}_{\text{re-cal}}$  dataset as in our method.** We report the results for the head that achieves the smallest set size among the evaluated heads.
- **Uniform average (Uniform)** - We average the scores obtained for each head and perform standard conformal prediction on the entire  $\mathcal{D}_{\text{cal}}$  dataset.
- **Optimized weights (Optimized)** - We use the weighted score  $s_W$  defined in Eq. (7). We perform constrained optimization, minimizing the mean set size with a constraint that the empirical mis-coverage does not exceed  $\alpha$ . The weights are optimized using Optuna (Akiba et al., 2019) over  $\mathcal{D}_{\text{cells}}$  with 100 optimization steps. Then, we perform standard conformal prediction over  $\mathcal{D}_{\text{re-cal}}$ .



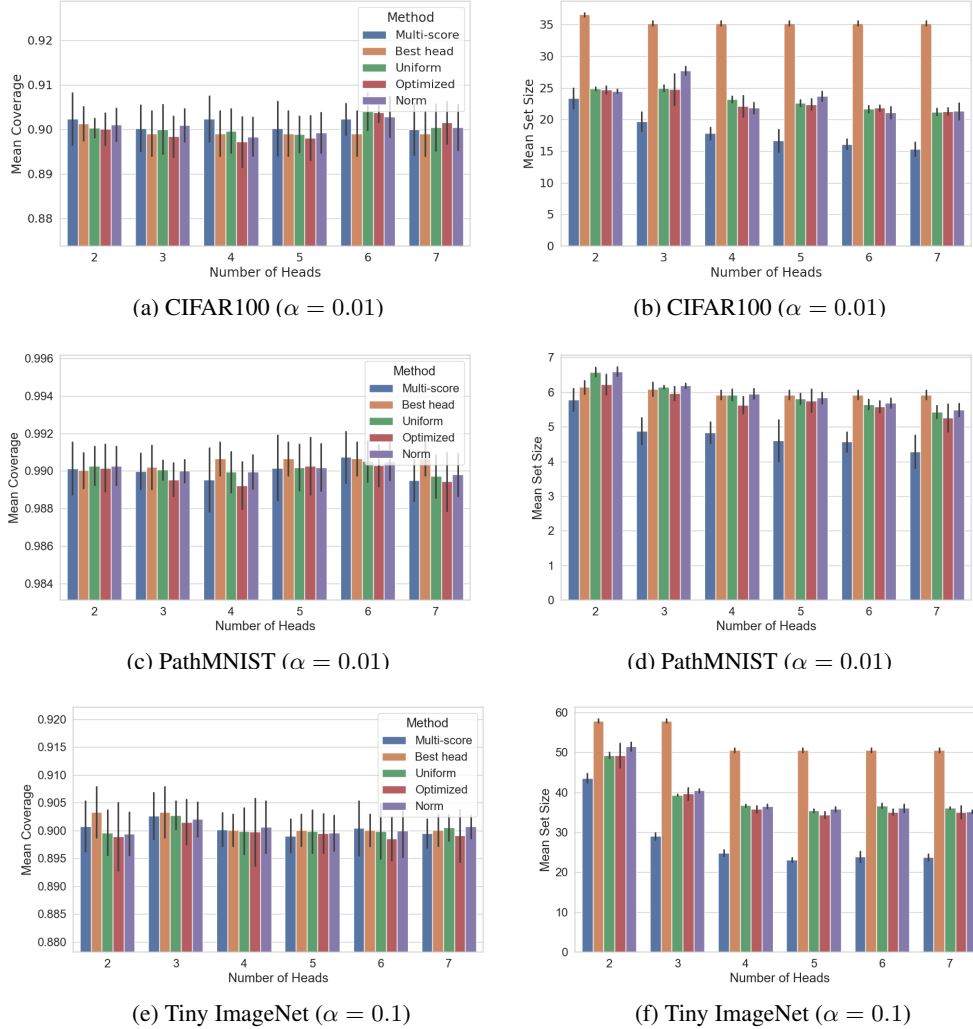


Figure 3: Conformal prediction with RAPS as a function of the number of classification heads. Results compare multi-score conformal prediction and the baselines (Best head, Optimized, Uniform, and Norm) across two metrics: empirical coverage (left column), and mean set size (right column), over: CIFAR100, Tiny ImageNet and PathMNIST.

- **Norm-based score (Norm)** - Conformal prediction is performed over  $\mathcal{D}_{\text{cal}}$  with a norm-based score, defined as  $s_N(x, y) = \|\mathbf{s}(x, y)\| = \sqrt{\sum_{i=1}^n s_i^2(x, y)}$ .

We evaluate the different methods in terms of the empirical coverage  $\frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(X, Y) \in \mathcal{D}_{\text{test}}} \mathbf{1}\{Y \in \Gamma(X)\}$  and the mean set size  $\frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(X, Y) \in \mathcal{D}_{\text{test}}} |\Gamma(X)|$ , computed over the test data  $\mathcal{D}_{\text{test}}$ . We report the average results and the standard deviation over 10 random splits to calibration and test.

### 5.3 RESULTS

**Performance with varying number of heads.** Results as a function of the number of heads are shown in Figs. 3 and D.2, with RAPS and SAPS as base scores, respectively. Results for additional  $\alpha$  levels are provided in Figs. D.1, and D.3. As expected all methods obtain the required coverage. We observe that the the proposed multi-score calibration leads to smaller prediction sets, with decreased sizes as the number of heads increases. Similar trends are observed for both RAPS and SAPS scores. Fig. D.7 illustrates how the set size distribution changes while the number of heads increases, showing the benefit of the proposed method in producing smaller set sizes compared to the baselines, especially as the number of heads increases.

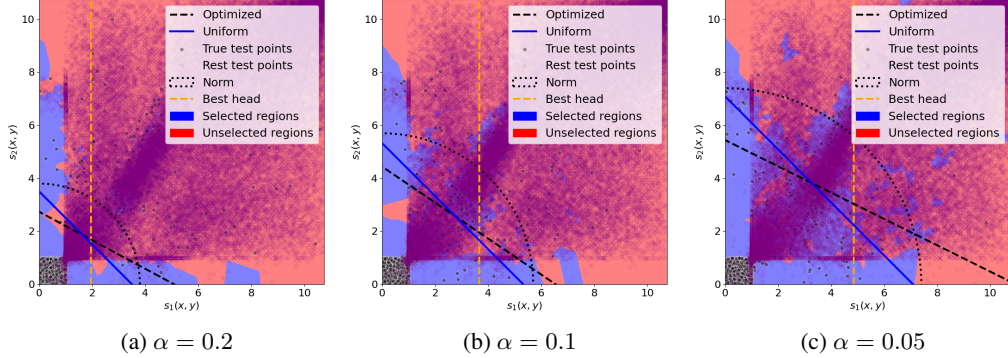


Figure 4: We present the selection regions for a 2-dimensional score ( $n = 2$ ) using RAPS on the Tiny ImageNet dataset, across different  $\alpha$  levels. The region selected by our method is shaded in blue, while the unselected region is shaded in orange. The decision boundaries of the baseline methods are shown with dashed lines. For the baselines, the selected region lies to the left of the “Best head” boundary and below the boundaries for Norm, Uniform, and Optimized methods. True test points are depicted as green circles, while test points with incorrect labels are marked by purple x-marks.

**Selection region.** Figure 4 demonstrates the results obtained for the 2-dimensional case ( $n = 2$ ). We present the region selected by the proposed method (blue area), and the decision boundaries for the baseline methods. In addition, the test scores for true and false labels are presented. We observe that the test scores are concentrated in two square regions: on the left bottom corner there are mostly true labels, whereas on the right upper corner there are mostly false labels. We see that our method focuses on regions that are less populated by false labels, while the baselines have a fixed structure and thus unavoidably include areas with a large density of false labels. Figure D.4 illustrates the cell selection order defined by Eq. (9), where the cells are colored according to their normalized rank in  $[0, 1]$ . Here too, we observe that cells near the bottom left corner are preferable, as well as cells that have low score in either dimension.

**Additional results.** We briefly highlight a few additional results that are included in Appendix D. We examined the performance in terms of the conditional coverage, defining groups by set size ranges and reporting the maximum coverage violation across these groups. The results in Table D.1 show that all methods exhibit similar behavior regarding the maximum coverage violation. Additionally, we conducted several experiments to demonstrate the versatility of our approach in other settings where multiple scores can be extracted, including: (i) a standard ensemble of separately trained models, (ii) test-time augmentation, and (iii) multiple scores computed from a single head. Moreover, we conducted additional experiments on text classification task, and ImageNet to verify that our method is suitable for different data types and large datasets. In all cases, our method achieves comparable or smaller set sizes with respect to the baselines. We examined the performance of the two additional variants of our method, concluding that the Jackknife+ version achieves smaller set sizes, and the soft version with  $b > 1$  is not preferable over  $b = 1$  in most cases. We also performed an ablation study to assess the impact of the diversity-regularized loss (12) used in training the classification heads. Furthermore, we show that our method is not highly sensitive with respect to the size of  $\mathcal{D}_{\text{cal}}$  and the size of  $\mathcal{D}_{\text{cells}}$ .

## 6 CONCLUSIONS

We propose a multi-score conformal prediction procedure that combines multiple nonconformity scores to select regions in the high-dimensional score space. This selection provides the desired coverage while minimizing the number of false labels that fall in the selected region. Unlike existing approaches, our method requires no optimization (neither black-box nor gradient-based) and is applicable to any coverage level. We propose to construct such a score using a cost-efficient self-ensemble model with multiple classification heads, trained with both CE and diversity losses. Experimental results demonstrate the superiority of our multi-score method over state-of-the-art baselines across several benchmark datasets.

## REFERENCES

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 2623–2631, 2019.
- Anastasios Nikolas Angelopoulos, Stephen Bates, Michael Jordan, and Jitendra Malik. Uncertainty sets for image classifiers using conformal prediction. In *International Conference on Learning Representations*, 2020.
- Yu Bai, Song Mei, Huan Wang, Yingbo Zhou, and Caiming Xiong. Efficient and differentiable conformal prediction with general function classes. In *International Conference on Learning Representations*, 2021.
- Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. Predictive inference with the jackknife+. *The Annals of Statistics*, 49(1):486–507, 2021.
- Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once for all: Train one network and specialize it for efficient deployment. In *International Conference on Learning Representations*, 2020.
- Guillaume Carlier, Victor Chernozhukov, and Alfred Galichon. Vector quantile regression: An optimal transport approach. *The Annals of Statistics*, pp. 1165–1192, 2016.
- Dmitry Devetyarov and Ilia Nouretdinov. Prediction with confidence based on a random forest classifier. In *Artificial Intelligence Applications and Innovations: 6th IFIP WG 12.5 International Conference, AIAI 2010, Larnaca, Cyprus, October 6-7, 2010. Proceedings 6*, pp. 37–44. Springer, 2010.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pp. 4171–4186, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- Bat-Sheva Einbinder, Yaniv Romano, Matteo Sesia, and Yanfei Zhou. Training uncertainty-aware classifiers with conformalized deep learning. *Advances in Neural Information Processing Systems*, 2022.
- Shai Feldman, Stephen Bates, and Yaniv Romano. Calibrated multiple-output quantile regression with representation learning. *Journal of Machine Learning Research*, 24(24):1–48, 2023.
- Matteo Gasparin and Aaditya Ramdas. Merging uncertainty sets via majority vote. *arXiv preprint arXiv:2401.09379*, 2024.
- Lars Kai Hansen and Peter Salamon. Neural network ensembles. *IEEE transactions on pattern analysis and machine intelligence*, 12(10):993–1001, 1990.
- Jianguo Huang, Huajun Xi, Linjun Zhang, Huaxiu Yao, Yue Qiu, and Hongxin Wei. Conformal prediction for deep classifier via label ranking. In *International Conference on Machine Learning (ICML)*, 2024.
- Shayan Kiyani, George J Pappas, and Hamed Hassani. Length optimization in conformal prediction. In *The Annual Conference on Neural Information Processing Systems*, 2024.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.

- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- Benjamin Lambert, Florence Forbes, Senan Doyle, Harmonie Dehaene, and Michel Dojat. Trustworthy clinical ai solutions: a unified review of uncertainty quantification in deep learning models for medical image analysis. *Artificial Intelligence in Medicine*, pp. 102830, 2024.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- Jing Lei and Larry Wasserman. Distribution-free prediction bands for non-parametric regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 76(1):71–96, 2014.
- Jing Lei, James Robins, and Larry Wasserman. Distribution-free prediction sets. *Journal of the American Statistical Association*, 108(501):278–287, 2013.
- Henrik Linusson, Ulf Johansson, and Henrik Boström. Efficient conformal predictor ensembles. *Neurocomputing*, 397:266–278, 2020.
- Helen Lu. *Data Augmentation and Conformal Prediction*. PhD thesis, Massachusetts Institute of Technology, 2023.
- Rui Luo and Zhixin Zhou. Trustworthy classification through rank-based conformal prediction sets. *arXiv preprint arXiv:2407.04407*, 2024a.
- Rui Luo and Zhixin Zhou. Weighted aggregation of conformity scores for classification. *arXiv preprint arXiv:2407.10230*, 2024b.
- Khan Muhammad, Amin Ullah, Jaime Lloret, Javier Del Ser, and Victor Hugo C de Albuquerque. Deep learning for safe autonomous driving: Current challenges and future directions. *IEEE Transactions on Intelligent Transportation Systems*, 22(7):4316–4336, 2020.
- Lorena Qendro and Cecilia Mascolo. Towards adversarial robustness with early exit ensembles. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pp. 313–316. IEEE, 2022.
- Lorena Qendro, Alexander Campbell, Pietro Lio, and Cecilia Mascolo. Early exit ensembles for uncertainty quantification. In *Machine Learning for Health*, pp. 181–195. PMLR, 2021.
- Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. *Advances in Neural Information Processing Systems*, 33:3581–3591, 2020.
- Aviv A Rosenberg, Sanketh Vedula, Yaniv Romano, and Alexander Bronstein. Fast nonlinear vector quantile regression. In *The Eleventh International Conference on Learning Representations*, 2022.
- Mauricio Sadinle, Jing Lei, and Larry Wasserman. Least ambiguous set-valued classifiers with bounded error levels. *Journal of the American Statistical Association*, 114(525):223–234, 2019.
- David Stutz, Krishnamurthy Dj Dvijotham, Ali Taylan Cemgil, and Arnaud Doucet. Learning optimal conformal classifiers. In *International Conference on Learning Representations*, 2021.
- Jiaye Teng, Chuan Wen, Dinghuai Zhang, Yoshua Bengio, Yang Gao, and Yang Yuan. Predictive inference with feature conformal prediction. In *The International Conference on Learning Representations*, 2022.
- Vladimir Vovk. Cross-conformal predictors. *Annals of Mathematics and Artificial Intelligence*, 74: 9–28, 2015.
- Vladimir Vovk, Ivan Petej, and Valentina Fedorova. Large-scale probabilistic predictors with and without guarantees of validity. *Advances in Neural Information Processing Systems*, 28, 2015.
- Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Scientific Data*, 10(1):41, 2023a.
- Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter Abbeel, and Dale Schuurmans. Foundation models for decision making: Problems, methods, and opportunities. *arXiv:2303.04129*, 2023b.

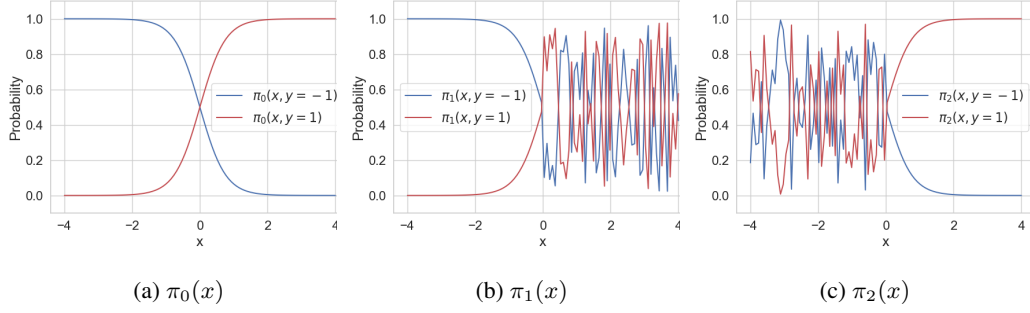


Figure A.1: The three classifiers in the binary classification problem.

## A MATHEMATICAL DETAILS

### A.1 TOY EXAMPLE DETAILS

Assume a binary classification problem with  $Y \in \{-1, 1\}$  and prior probabilities  $\mathbb{P}(Y = 1)$  and  $\mathbb{P}(Y = -1)$ . The input  $x$  is generated from a mixture of two Gaussians, with  $\mathbb{P}(X|Y) = \mathcal{N}(y, \sigma_y^2)$ . In this example, we can compute the posterior  $\mathbb{P}(Y|X)$  using Bayes rule:

$$\mathbb{P}(Y = a|X) = \frac{\mathbb{P}(Y = a)\mathbb{P}(X|Y = a)}{\mathbb{P}(Y = 1)\mathbb{P}(X|Y = 1) + \mathbb{P}(Y = -1)\mathbb{P}(X|Y = -1)}, \quad a = -1, 1 \quad (13)$$

We set  $\mathbb{P}(Y = 1) = \mathbb{P}(Y = -1) = 0.5$ ,  $\sigma_1^2 = \sigma_{-1}^2 = 0.75$ . Consider the following three classifiers:

$$\pi_0(x) = \mathbb{P}(y|x), \quad \pi_1(x) := \begin{cases} \mathbb{P}(y|x), & \text{if } x \leq 0, \\ \epsilon, 1 - \epsilon, & \text{if } x > 0 \end{cases}, \quad \pi_2(x) := \begin{cases} \epsilon, 1 - \epsilon, & \text{if } x \leq 0, \\ \mathbb{P}(y|x), & \text{if } x > 0 \end{cases} \quad (14)$$

where  $\epsilon \sim \mathcal{N}(0, 1)$ . The classifiers are illustrated in Fig. A.1. Here  $\pi_0(x)$  represents the ideal classifier, and  $\pi_1(x)$  and  $\pi_2(x)$  are ideal only in half of the range of  $x$  and uninformative for the other half. We generate 2000 points using  $p(x)$ , and use 1000 for validation and 1000 for calibration. We compute the Thr nonconformity score. Results are averaged over 20 random trials.

When  $\mathbb{P}(Y|X)$  is known, it was shown that the optimal set with minimal size under coverage constraint is given by  $\Gamma^*(x) = \{y \in \mathcal{Y} | p(y|x) > q_\alpha\}$  (Lei & Wasserman, 2014; Sadinle et al., 2019; Kiyani et al., 2024). This implies that the optimal set is a level set of the distribution  $p(y|x)$ . Thus, in our example, thresholding  $s_0(x, y) = 1 - p(y|x)$  results in the optimal set. Using only  $s_1(x, y)$  and  $s_2(x, y)$  the optimal set can be equivalently defined as:

$$\begin{aligned} \Gamma^*(x) &= \{y \in \mathcal{Y} | (\mathbf{1}\{x \leq 0\} \cdot s_1(x, y) + \mathbf{1}\{x > 0\} \cdot s_2(x, y)) < 1 - q_\alpha\} \\ &= \{y \in \mathcal{Y} | \mathbf{s}^T(x, y) \cdot \mathbf{i}_x(x) < 1 - q_\alpha\} \end{aligned}$$

where  $\mathbf{i}_x(x) = [\mathbf{1}\{x \leq 0\}, \mathbf{1}\{x > 0\}]^T$ . Thus, we obtain that the optimal set is a function of the 2-dimensional nonconformity score  $\mathbf{s}(x, y)$ . In this example, each classifier specializes on a different subdomain of the input space  $\mathcal{X}$ . Another practical case is when classifiers specialize on different parts of the output space  $\mathcal{Y}$ . For example, consider  $\mathcal{Y} = \{0, 1, 2\}$ , and the following three classifiers:

$$\pi_a(x) := \begin{cases} \mathbb{P}(y = a|x), & \text{if } y = a, \\ \epsilon, & \text{if } y = (a + 1) \bmod 3, \\ 1 - \mathbb{P}(y = a|x) - \epsilon, & \text{if } y = (a + 2) \bmod 3 \end{cases}, \quad a \in \{0, 1, 2\}, \quad (15)$$

where  $\epsilon \sim \mathcal{N}(0, 1)$ . In this case, the optimal set is given by:

$$\Gamma^*(x) = \{y \in \mathcal{Y} | \mathbf{s}^T(x, y) \cdot \mathbf{i}_y(y) < 1 - q_\alpha\} \quad (16)$$

where  $\mathbf{i}_y(y) = [\mathbf{1}\{y = 1\}, \mathbf{1}\{y = 2\}, \mathbf{1}\{y = 3\}]^T$ . We conclude that whenever  $p(y|x) = \phi(\mathbf{s}(x, y); x, y)$ , where  $\phi : \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$  is a non-degenerate function of the multi-score vector, the optimal set relies on  $\mathbf{s}(x, y)$ , i.e.:

$$\Gamma^*(x) = \{y \in \mathcal{Y} | \phi(\mathbf{s}(x, y); x, y) < 1 - q_\alpha\} \quad (17)$$

In contrast, relying solely on a single score will result in a suboptimal solution. Note that, according to Eq. (17), the ideal set corresponds to a level set of  $\phi(\mathbf{s}(x, y); x, y)$ , rather than  $\mathbf{s}(x, y)$  itself. This implies that, in general, the decision boundaries in the multi-dimensional score space can be arbitrarily complex, depending on the properties of  $\phi$ .

In practice, we have access to neither the conditional distribution  $\mathbb{P}(Y|X)$  nor the mapping function  $\phi$ . Instead, we aim to solve the problem of minimizing the set size subject to a coverage constraint, similarly to (Stutz et al., 2021; Bai et al., 2021; Kiyani et al., 2024). Since the space of all possible prediction sets is overly complex, the problem must be relaxed. Bai et al. (2021) proposed to optimize an arbitrary class of prediction sets  $\Gamma_\theta$  parametrized by  $\theta$ , while Stutz et al. (2021) optimize a parametrized score  $s_\theta(x, y)$ . In (Kiyani et al., 2024), structured prediction sets of the form  $\Gamma_h^s(x) = \{y \in Y | s(x, y) \leq h(x)\}$  were considered, where  $h : \mathcal{X} \rightarrow \mathbb{R}$  is a learned adaptive threshold. In contrast, we work in the multi-score domain and consider sets defined as a union of a subset  $I \subseteq 2^k$  of cells in  $\mathcal{D}_{\text{cells}}$ , i.e.  $\Gamma_I^s(x) = \{y \in \mathcal{Y} | \mathbf{s}(x, y) \in \cup_{i \in I} \mathcal{C}_i\}$ . Accordingly, the relaxed optimization problem can be written as:

$$\begin{aligned} \arg \min_{I \subseteq 2^k} \mathbb{E} [\text{size}(\Gamma_I^s(X))] \\ \text{s.t. } \mathbb{E} [\mathbf{1}\{Y \in \Gamma_I^s(X)\}] \geq 1 - \alpha \end{aligned} \quad (18)$$

where  $\text{size}(\Gamma_I^s(X)) = \sum_{q=1}^Q \mathbf{1}\{Y \in \Gamma_I^s(X)\} = \sum_{i \in I} \sum_{q=1}^Q \mathbf{1}\{Y \in \mathcal{C}_i(X)\}$ . A finite sample approximation of the expected set size is given by:

$$\mathcal{E} = \frac{1}{k} \sum_{j=1}^k \sum_{i \in I} \sum_{q=1}^Q \mathbf{1}\{\mathbf{s}(X_j, q) \in \mathcal{C}_i\}. \quad (19)$$

Note that this is equivalent to summing the cell scores  $D_i$  defined in Eq. (9) (without removing duplicate cells):

$$\mathcal{E} = \frac{1}{k} \sum_{i \in I} D_i. \quad (20)$$

Therefore, solving the optimization problem in Eq. (18) does not require enumerating all possible sets  $I$ , which is computationally infeasible. Instead, we can rank the cells according to  $D_i$  and take the  $(1 - \alpha)$  proportion of cells with the lowest score values. To obtain exact coverage we perform a recalibration over  $\mathcal{D}_{\text{re-cal}}$ .

We conclude that in practical scenarios, where a single score does not provide the full information on the conditional distribution  $\mathbb{P}(y|x)$ , we benefit from using a multi-dimensional score  $\mathbf{s}(x, y)$ . It may appear that optimizing for set size efficiency in the multi-score space exponentially increases the number of possible prediction sets to be considered, which makes the optimization more challenging compared to the single-dimensional case. However, our cell partitioning and ranking procedure relaxes the problem to a convenient structured prediction with a simple selection rule that does not require any iterative optimization procedures. Note that the number of cell centers and the summation operation over all scores that fall in the chosen region, remain fixed regardless of the dimensionality of  $\mathbf{s}(x, y)$ . However, as  $n$  increases the cells move apart from each other, when the scores are nonidentical and provide complementary information. Moreover, if each dimension contributes information about the actual conditional distribution  $\mathbb{P}(Y|X)$ , we anticipate an improved separation between true and false scores. Consequently, the selected subset of cells is expected to exhibit lower  $D_i$  values, leading to smaller prediction sets. This is demonstrated in Fig. A.2 presenting the distribution of  $D_i$  for the chosen cells. We observe that as  $n$  increases, the values of  $D_i$  become smaller.

## A.2 PROOF OF THEOREM (2)

*Proof.* For a pair of test input  $X_{m+1}$  and a (candidate) label  $y \in \mathcal{Y}$  we compute the multi-dimensional score  $\mathbf{s}(X_{m+1}, y) = [s_1(X_{m+1}, y), \dots, s_n(X_{m+1}, y)]^T$ . We define a combined score function  $s_{\text{multi}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  that maps the multi-dimensional score  $\mathbf{s}(x, y)$  into a single-dimensional score  $s_{\text{multi}}(x, y)$ . The mapping is defined as follows:

$$s_{\text{multi}}(X_{m+1}, y) = \arg \min_{1 \leq j \leq k'} \|\mathbf{s}(X_{m+1}, y) - \mathbf{s}(X_{(j)}, Y_{(j)})\| \quad (21)$$



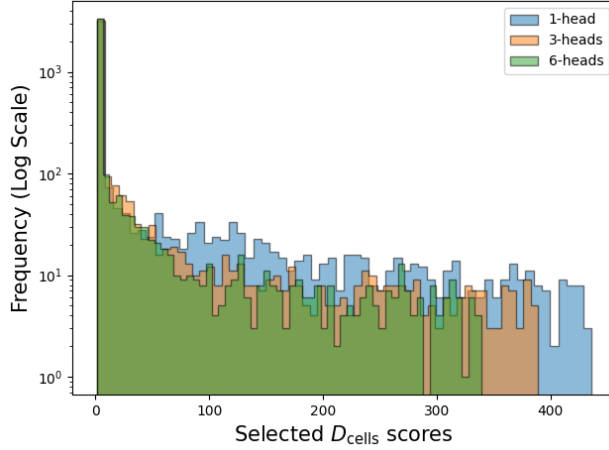


Figure A.2: Histogram of cell scores for the selected number of cells, comparing different number of heads. Values correspond to CIFAR100 dataset, RAPS scores and  $\alpha = 0.1$

where  $(j)$  denotes the index of the samples in  $\mathcal{D}_{\text{cells}}$  after sorting according to the ratio values (9) and eliminating repeating elements, i.e.  $(X_{(1)}, Y_{(1)}), \dots, (X_{(k')}, Y_{(k')})$  are the centers of the cells  $\mathcal{C}_{(1)}, \dots, \mathcal{C}_{(k')}$  where  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(k')}$ . This way each pair  $(X_{m+1}, y)$  is associated with the closet center and the ranking of this cell serves as the one-dimensional score, defined in Eq. (21).

The prediction set  $\Gamma_{\eta^*}(X_{m+1})$  is defined in Eq. (10) by including predictions associated with scores that reside in the selected region. This can be equivalently written in terms of a threshold operation over  $s_{\text{multi}}$ , as we can write:

$$\begin{aligned}
 \Gamma_{\eta^*}(X_{m+1}) &= \left\{ y \in \{1, \dots, Q\} \mid \mathbf{s}(X_{m+1}, y) \in \mathcal{C}_{\text{in}}^{\eta^*} \right\} \\
 &= \left\{ y \in \{1, \dots, Q\} \mid \mathbf{s}(X_{m+1}, y) \in \bigcup_{i=1}^{\eta^*} \mathcal{C}_{(i)} \right\} \\
 &= \left\{ y \in \{1, \dots, Q\} \mid \bigcup_{i=1}^{\eta^*} [\mathbf{s}(X_{m+1}, y) \in \mathcal{C}_{(i)}] \right\} \\
 &\stackrel{(8)}{=} \left\{ y \in \{1, \dots, Q\} \mid \bigcup_{i=1}^{\eta^*} \left[ i = \arg \min_{1 \leq j \leq k'} \|\mathbf{s}(X_{m+1}, y) - \mathbf{s}(X_{(j)}, Y_{(j)})\| \right] \right\} \\
 &\stackrel{(21)}{=} \left\{ y \in \{1, \dots, Q\} \mid \bigcup_{i=1}^{\eta^*} [i = s_{\text{multi}}(X_{m+1}, y)] \right\} \\
 &= \{y \in \{1, \dots, Q\} \mid s_{\text{multi}}(X_{m+1}, y) \leq \eta^*\}.
 \end{aligned} \tag{22}$$

This formulation is the same as a standard one-dimensional conformal prediction procedure, where we define the prediction set by thresholding a one-dimensional score function. The threshold  $\eta^*$  was chosen to obtain exact coverage over the held-out data  $\mathcal{D}_{\text{re-cal}}$ , therefore:

$$\begin{aligned}
 \eta^* &= \min\{\eta \in \{1, \dots, k'\} \mid Y_i \in \mathcal{C}_{\text{in}}^\eta \text{ for at least } \lceil (1 - \alpha)(r + 1) \rceil \text{ samples } (X_i, Y_i) \in \mathcal{D}_{\text{re-cal}}\} \\
 &\stackrel{(22)}{=} \min\{\eta \in \{1, \dots, k'\} \mid s_{\text{multi}}(X_i, Y_i) \leq \eta \text{ for at least } \lceil (1 - \alpha)(r + 1) \rceil \text{ samples in } \mathcal{D}_{\text{re-cal}}\} \\
 &= \text{Quantile} \left( \frac{\lceil (r + 1)(1 - \alpha) \rceil}{r}; \{s_{\text{multi}}(X_i, Y_i)\}_{i=k+1}^m \right).
 \end{aligned}$$

Thus, assuming that  $\mathcal{D}_{\text{re-cal}}$  and  $(X_{m+1}, Y_{m+1})$  are exchangeable, we obtain that:

$$\mathbb{P}(Y_{m+1} \in \Gamma_{\eta^*}(X_{m+1})) \geq 1 - \alpha \tag{23}$$

**Algorithm B.1 Jackknife+ Multi-Score Conformal prediction**

**Definitions:**  $s(x, y)$  is a multi-dimensional score function.  $\mathcal{D}_{\text{cal}}$  is the calibration data of size  $m$ .  $X_{m+1}$  is a new test sample,  $\alpha$  is the miscoverage level,  $k$  is the number of samples for cell-partitioning and scoring, and  $r = m - k$  is the number of samples for re-calibration.

```

1: function MULTI-SCORE-CP( $s(x, y)$ ,  $\mathcal{D}_{\text{cal}}$ ,  $\alpha$ )
2:   Compute the scores  $s(X_i, Y_i)$ ,  $i \in \{1, \dots, m\}$ 
3:   Segment  $\mathcal{S}$  into cells  $\{\mathcal{C}_i\}_{i=1}^k$  centered at  $\{(X_i, Y_i)\}_{i=1}^k$ 
4:   Remove duplicate cells  $\mathcal{C}_1, \mathcal{C}_1, \dots, \mathcal{C}_{k'}$ 
5:   for  $i \in \{1, \dots, k\}$  do
6:     Compute  $D_j^{-i}$ ,  $j = 1, \dots, i-1, i+1, \dots, k'$  using Eq. (9), excluding the  $i$ -th point
7:     Rank the cells  $\mathcal{C}_{(1)}, \mathcal{C}_{(2)}, \dots, \mathcal{C}_{(k'-1)}$  according to  $D_{(1)}^{-i} \leq D_{(2)}^{-i} \leq \dots \leq D_{(k'-1)}^{-i}$ 
8:      $E_i^{-i} \leftarrow \arg \min_{j \in \{1, \dots, k'-1\}} \|s(X_i, Y_i) - s(X_{(j)}, Y_{(j)})\|$ 
9:      $E_{m+1}^{-i}(y) \leftarrow \arg \min_{j \in \{1, \dots, k'-1\}} \|s(X_{m+1}, y) - s(X_{(j)}, Y_{(j)})\|$ ,  $y \in \mathcal{Y}$ 
10:   $\Gamma_{JK+}(X_{m+1}) = \{y \in \{1, \dots, Q\} \mid \sum_{i=1}^m \mathbf{1}\{E_i^{-i} < E_{m+1}^{-i}(y)\} < (1 - \alpha)(m + 1)\}$ 
11:  return  $\Gamma_{JK+}(X_{m+1})$ 

```

**Algorithm B.2 Soft Multi-Score Conformal Prediction**

**Definitions:**  $s(x, y)$  is a multi-dimensional score function.  $\mathcal{D}_{\text{cal}}$  is the calibration data of size  $m$ .  $X_{m+1}$  is a new test sample,  $\alpha$  is the miscoverage level,  $k$  is the number of samples for cell-partitioning and scoring, and  $r = m - k$  is the number of samples for re-calibration and  $b$  is the number of neighbors.

```

1: function SOFT MULTI-SCORE-CP( $s(x, y)$ ,  $\mathcal{D}_{\text{cal}}$ ,  $\alpha, b$ )
2:   Randomly split  $\mathcal{D}_{\text{cal}}$  to  $\mathcal{D}_{\text{cells}} = \{(X_i, Y_i)\}_{i=1}^k$  and  $\mathcal{D}_{\text{re-cal}} = \{(X_i, Y_i)\}_{i=k+1}^m$ 
3:   Compute the scores  $s(X_i, Y_i)$ ,  $i \in \{1, \dots, m\}$ 
4:   Segment  $\mathcal{S}$  into cells  $\{\mathcal{C}_i\}_{i=1}^k$  centered at  $\{(X_i, Y_i)\}_{i=1}^k$ 
5:   Compute  $D_i$ ,  $i = 1, \dots, k$  according to Eq. (9)
6:   Remove duplicate cells  $\mathcal{C}_1, \mathcal{C}_1, \dots, \mathcal{C}_{k'}$ 
7:   Rank the cells  $\mathcal{C}_{(1)}, \mathcal{C}_{(2)}, \dots, \mathcal{C}_{(k')}$  according to  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(k')}$ 
8:    $\mathcal{T}^b(i) \leftarrow b$  nearest cells of  $s(X_i, Y_i)$  in  $\{\mathcal{C}_{(i)}\}_{i=1}^{k'}$ ,  $i \in \{k+1, \dots, m\}$ 
9:    $\eta^* \leftarrow \min \left\{ \eta \mid \sum_{i=k+1}^m \mathbf{1} \left\{ \left( \sum_{t=1}^b \mathbf{1} \{ \mathcal{T}_t^b(i) \subseteq \mathcal{C}_{\text{in}}^\eta \} \right) > \lceil 0.5 \cdot b \rceil \right\} \geq \lceil (1 - \alpha)(r + 1) \rceil \right\}$ 
10:   $\mathcal{C}_{\text{in}}^{\eta^*} \leftarrow \bigcup_{i=1}^{\eta^*} \mathcal{C}_{(i)}$ 
11:  return  $\mathcal{C}_{\text{in}}^{\eta^*}$ 
12: function SOFT MULTI-SCORE-EVALUATION( $s(x, y)$ ,  $X_{m+1}$ ,  $\mathcal{C}_{\text{in}}^{\eta^*}$ )
13:   Compute the scores  $s(X_{m+1}, y)$ ,  $y \in \{1, \dots, Q\}$ 
14:    $\mathcal{T}^b(y) \leftarrow b$  nearest cells of  $s(X_{m+1}, y)$  in  $\{\mathcal{C}_{(i)}\}_{i=1}^{k'}$ ,  $y \in \mathcal{Y}$ 
15:    $\Gamma_{\eta^*}(X_{m+1}) = \{y \in \mathcal{Y} \mid \left( \sum_{t=1}^b \mathbf{1} \{ \mathcal{T}_t^b(y) \subseteq \mathcal{C}_{\text{in}}^{\eta^*} \} \right) > \lceil 0.5 \cdot b \rceil \}$ 
16:  return  $\Gamma_{\eta^*}(X_{m+1})$ 

```

following Theorem (1). □

Note that the split of  $\mathcal{D}_{\text{cal}}$  to two disjoint subsets  $\mathcal{D}_{\text{cells}}$  and  $\mathcal{D}_{\text{re-cal}}$  is critical here, since we cannot claim that the samples in  $\mathcal{D}_{\text{cells}}$  and the test pair  $(X_{m+1}, Y_{m+1})$  are exchangeable as  $\mathcal{D}_{\text{cells}}$  serves for the computation of the score defined in (21) (in a similar way to the fact that the training data cannot be used for calibration in standard conformal prediction).

**B ALGORITHMS**

The jackknife+ variant, as described in Section D.6, is provided in Algorithm B.1, while the soft centers variant, as detailed in Section D.6, is presented in Algorithm B.2.

## C IMPLEMENTATION AND DATASET DETAILS

**Datasets and Implementation.** The details of each dataset and the corresponding data splits are summarized in Table C.1, where Tiny ImageNet, CIFAR100 and PathMNIST are used in our main results while [ImageNet and 20NewsGroups](#) are used for the additional experiments, described in § D. The calibration data is split into half for cell computation, and re-calibration.

Table C.1: Datasets Details

Dataset	# Classes	Train	Validation	Calibration	Test	Average Accuracy
Tiny ImageNet	200	71,500	11,000	16,500	11,000	0.58
CIFAR100	100	39,000	6,000	9,000	6,000	0.69
PathMNIST	9	69,667	10,718	16,077	10,718	0.94
<a href="#">20 Newsgroups</a>	20	9,800	2,449	3,298	3,299	0.87
<a href="#">ImageNet</a>	1,000	1,281,184	10,000	20,000	20,000	0.71

For the first three datasets we used ResNet50 model with pretrained weights on ImageNet. Each head is a 3 layer feed-forward neural network with, BatchNorm, ReLU activation and dropout with  $p = 0.1$ .

In the first stage, the full model with a single classification head was fine-tuned on each task with 20, 100 and 200 epochs for Tiny ImageNet, CIFAR100 and PathMNIST, respectively. In the second stage, we freeze the backbone model and train only the classification heads for 20 epochs, using the loss defined in Eq. (12). In both stages, we use Adam optimizer with cosine annealing scheduler, momentum decay of 0.95, weight decay of  $1e - 5$ , and batch size of 16.

For the computation of the RAPS score we used  $\lambda = 0.05$  and  $\kappa = 5$ , and for the SAPS score we set  $\xi = 0.3$ . Performance with respect to other parameter values is reported in Tables D.4 and D.5 for RAPS and SAPS, respectively.

## D ADDITIONAL RESULTS

### D.1 EVALUATING OTHER PERFORMANCE ASPECTS

**Conditional Coverage.** We evaluated the conditional coverage of the proposed method in comparison to the baselines. To do so, we define a set of disjoint strata  $\{S_j\}_{j=1}^J$ , where  $\cup_{j=1}^J S_j = \{1, \dots, Q\}$ . We partition the data into groups with equal numbers of samples. Let  $q_j = \text{Quantile}(\frac{j}{J}; \{|\Gamma(X_i, Y_i)|\}_i)$ ,  $j = 0, \dots, J$ , denote the  $\frac{j}{J}$ -th quantile of the set sizes. The  $j$ th group is then defined as  $\mathcal{G}_j = \{i : q_{j-1} \leq |\Gamma(X_i, Y_i)| \leq q_j\}$  for  $j \in \{1, \dots, J\}$ . Accordingly, the size-stratified coverage violation (SSCV) is defined as (Angelopoulos et al., 2020):

$$\text{SSCV}(\{S_j\}_{j=1}^J) = \sup_j \left| \frac{|\{i : Y_i \in \Gamma(X_i, Y_i), i \in \mathcal{G}_j\}|}{|\mathcal{G}_j|} - (1 - \alpha) \right| \quad (24)$$

For this evaluation, we divided the data into  $J = 10$  groups. The results, summarized in Table D.1, indicate that all methods achieve similar SSCV scores, with no method showing a clear advantage over the others.

**Results with SAPS scores.** To examine the robustness of the multi-score method with respect to the choice of nonconformity scores, we compare the results achieved using the SAPS nonconformity score across different levels of  $\alpha$  and datasets. The results are presented in figs D.2 and D.3. As expected, all methods achieve the required coverage. The proposed multi-score calibration produces smaller prediction sets, with significant improvements for SAPS as the number of heads increases, demonstrating its robustness to the choice of the nonconformity score.

**Results with Thr and APS scores.** Our proposed method can be applied over any type of score. However, we have seen that it is not fully optimal for Thr and APS. The reason is that for these scores the values are more concentrated on specific levels, while for SAPS and RAPS the values

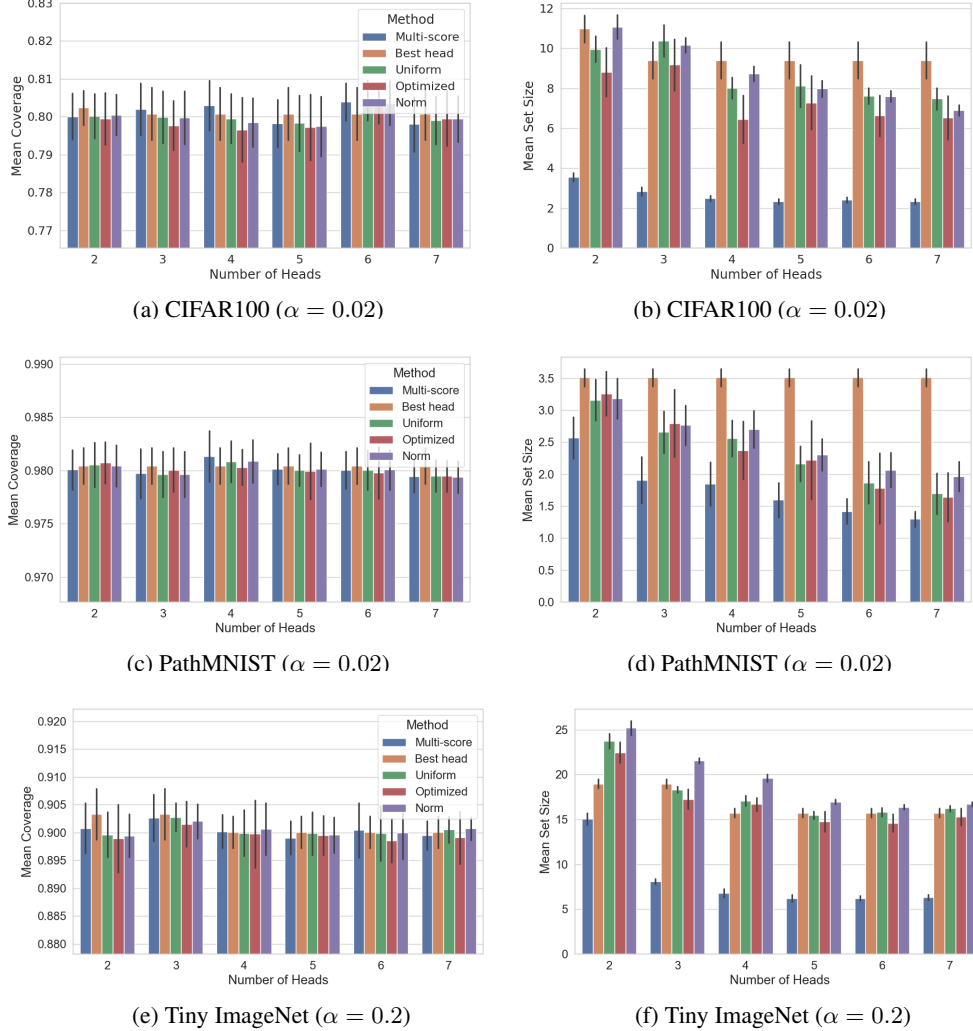


Figure D.1: Conformal prediction with RAPS as a function of the number of classification heads. Results compare multi-score conformal prediction and the baselines (Best head, Optimized, Uniform, and Norm) across two metrics: empirical coverage (left column), and mean set size (right column), over: CIFAR100, Tiny ImageNet and PathMNIST.

are more spread. In order to improve this behavior we use temperature scaling, i.e. we divide the logits by  $T$  before applying Softmax(). As  $T$  increases the entropy of the probabilities increases and they become more spread, as illustrated in Fig. D.5. Figure D.6 presents the set sizes obtained for all methods with respect to different temperatures. We see that the efficiency of the proposed method greatly improves as  $T$  increases for both Thr and APS, while the performance of the baseline methods is less effected by  $T$ . For RAPS and SAPS the proposed method outperforms the baselines regardless of the temperature due to the inherent spread of these scores.

**The distribution of set sizes.** Figure D.7 illustrates how increasing the number of heads (from 1 to 7) shifts the set size distribution toward smaller values, indicating more samples with smaller set sizes. Notably, the multi-score method responds more effectively to the increase in heads compared to the baselines, achieving set size values in a smaller range (2 – 41), while for the other methods the sizes range from 11 – 41. This highlights again that our multi-score method leads to more efficient and precise sets.

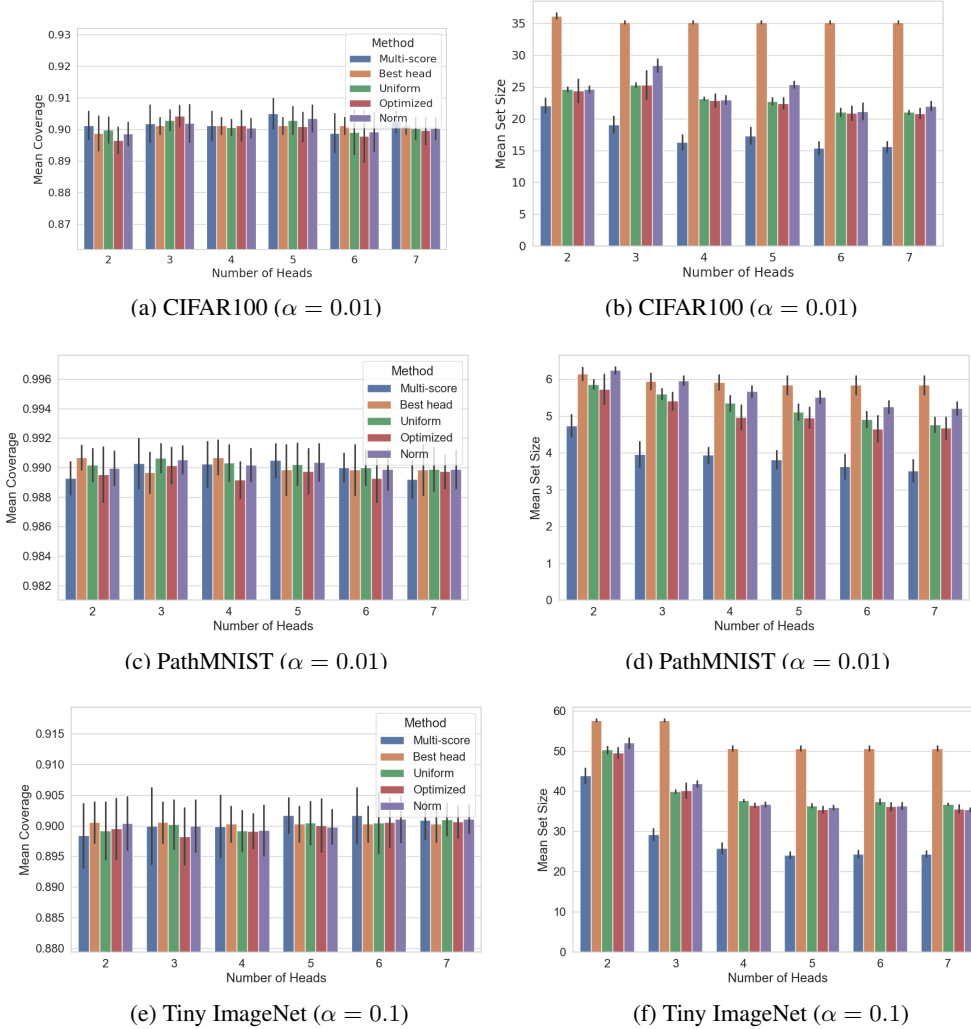


Figure D.2: Conformal prediction with SAPS as a function of the number of classification heads. Results compare multi-score conformal prediction and the baselines (Best head, Optimized, Uniform, and Norm) across two metrics: empirical coverage (left column), and mean set size (right column), over: CIFAR100, Tiny ImageNet and PathMNIST.

## D.2 RESULTS WITH OTHER TYPES OF MULTIPLE SCORES

**Standard Ensemble.** We conduct an experiment with a standard ensemble, consisting of multiple different models that are trained separately on the same dataset. We use ImageNet dataset for evaluation with an ensemble of 5 models pretrained on ImageNet: VGG16, Inception, ResNet50, ResNet152 and DenseNet161. Results are shown in Fig. D.8 for RAPS score and  $\alpha = [0.03, 0.05, 0.1]$ . Similarly to our main results with self-ensemble, here too the proposed method obtains the smallest prediction set sizes. This indicates that our method can be applied to self-ensemble models as well as regular ensembles.

**Test-Time Augmentation.** We evaluate our method on a multi-dimensional score that is formed by test-time augmentations (Lu, 2023). We use ImageNet dataset with Inception-V3 model. We use a simple common test-time augmentation policy (Krizhevsky et al., 2012), which consists of a random crop and a horizontal flip. The random crop pads the original image by 4 pixels and takes a 256x256 crop of the resulting image. We draw five augmentations using this policy. Figure D.9 presents the results for RAPS score and  $\alpha = [0.03, 0.05, 0.1]$ . Our method outperforms the baselines in terms of

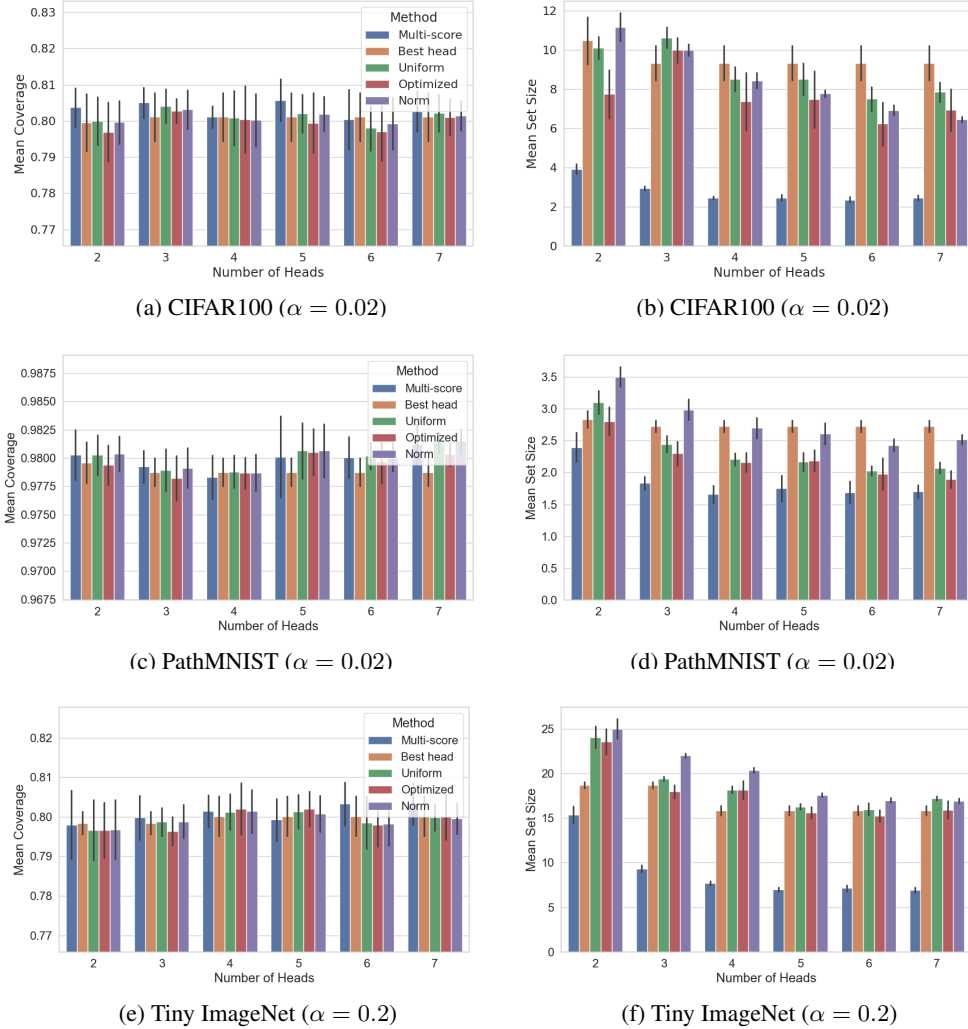


Figure D.3: Conformal prediction with SAPS as a function of the number of classification heads. Results compare multi-score conformal prediction and the baselines (Best head, Optimized, Uniform, and Norm) across two metrics: empirical coverage (left column), and mean set size (right column), over: CIFAR100, Tiny ImageNet and PathMNIST.

prediction set size. We conclude that test-time augmentation can serve as a possible alternative for generating multiple nonconformity scores within our multi-score conformal prediction framework.

**Multiple scores computed over a single head.** We examined a setting where instead of considering multiple classification heads, we use a single head and compute different conformity scores: Thr, APS, RAPS and SAPS. To ensure all scores are comparable and fall in the range between  $[0, 1]$ , we apply Softmax() over each score. Table D.2 summarizes the results for all datasets. We see that combining multiple scores improves the results compared to the best single score, and that the proposed method obtains the smallest prediction sets in almost all cases. The advantage of this score fusion is that it does not require any additional modifications to the model or further fine-tuning iterations.



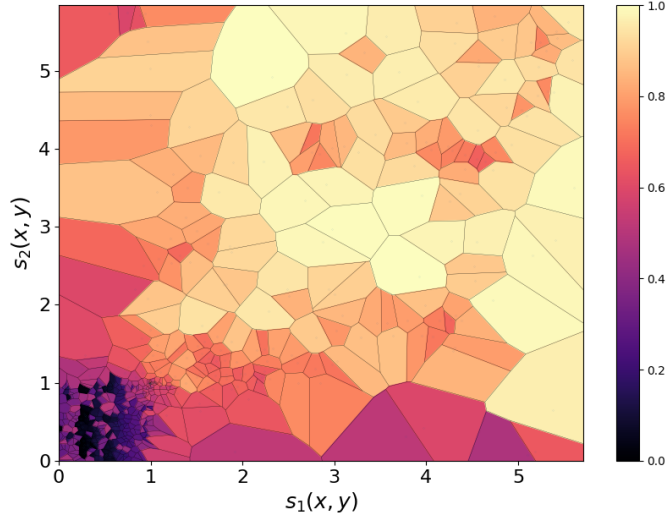


Figure D.4: Visualization of cell selection order over CIFAR100 with RAPS score and  $\alpha=0.1$ . Darker cells were selected earlier in the sequence while lighter cells are selected at a later stage.

Table D.1: SSCV measure for CIFAR100 with RAPS score

$\alpha$	Head	Best head	Multi-score	Norm	Optimized	Uniform
0.1	1	0.174	0.034	0.174	0.159	0.174
	2	0.167	0.108	0.114	0.194	0.141
	3	0.240	0.104	0.066	0.162	0.277
	4	0.298	0.204	0.178	0.167	0.264
	5	0.106	0.155	0.167	0.200	0.252
	6	0.131	0.153	0.172	0.308	0.125
	7	0.056	0.159	0.062	0.191	0.175
0.2	1	0.190	0.037	0.190	0.157	0.190
	2	0.362	0.249	0.292	0.221	0.330
	3	0.342	0.200	0.132	0.187	0.098
	4	0.379	0.206	0.167	0.166	0.279
	5	0.361	0.192	0.159	0.209	0.331
	6	0.220	0.180	0.071	0.274	0.276
	7	0.186	0.203	0.200	0.171	0.245

### D.3 ROBUSTNESS TO HYPERPARAMETERS

**Influence of training with diversity regularization.** We conducted an ablation study to examine the affect of adding diversity regularization to the loss used for training the classification heads, defined in Eq. (12). We use the same fine-tuned model in the first training stage and change only the second stage to optimize the regularized loss in Eq. (12) with different values of  $\lambda$ , controlling the strength of the diversity regularization. Figure D.10 shows the similarity between heads with and without the diversity regularization. As expected, the similarity between heads is smaller for the model trained with the regularized loss. Figure D.11 compares the set sizes obtained for different values of  $\lambda$ , demonstrating that adding the regularization results in smaller sets. The optimal value of  $\lambda$  is around 0.8, after which an increase in set size is observed, apparently due to the fact that increasing head diversity comes at the cost of decreasing the accuracy of each head.

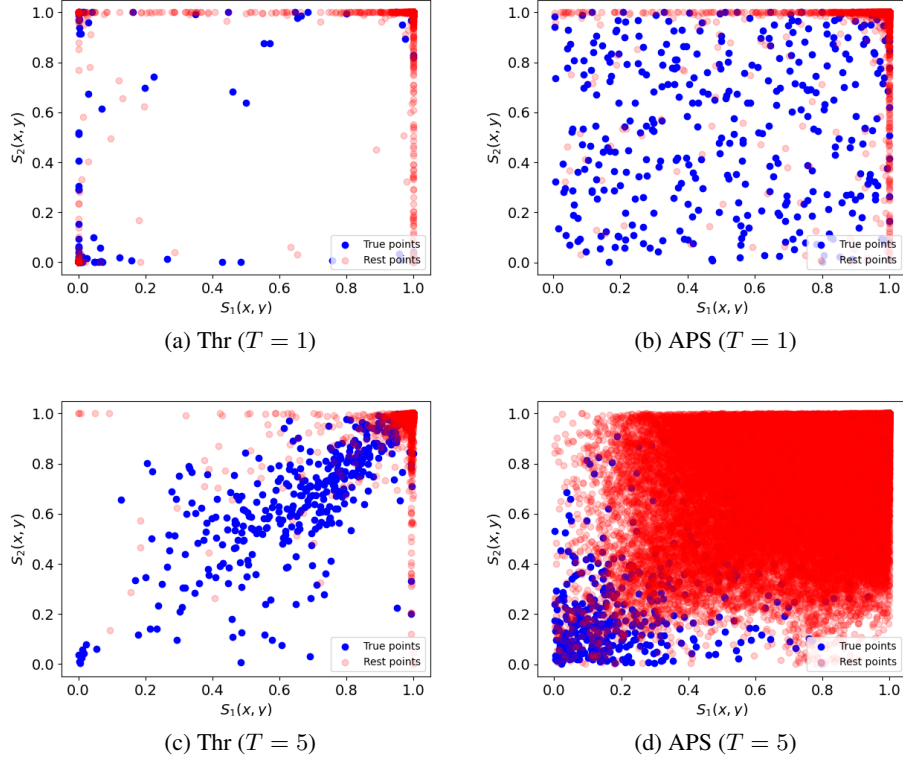


Figure D.5: Demonstration of APS and Thr scores’s distribution with different Temperatures on CIFAR-100 and  $\alpha = 0.1$ .

Table D.2: Results for a single classification head and multiple types of scores.

$\alpha$	Dataset	Methods				
		Multi Score	Best Head	Uniform	Optimized	Norm
<b>0.1</b>	<b>CIFAR100</b>	33.90	35.12	35.32	34.96	35.54
	<b>TinyImageNet</b>	58.08	61.62	62.96	62.70	78.45
<b>0.15</b>	<b>CIFAR100</b>	18.33	27.88	28.12	27.52	28.18
	<b>TinyImageNet</b>	35.84	48.55	49.47	49.48	55.75
<b>0.2</b>	<b>CIFAR100</b>	8.25	9.34	9.19	8.35	13.02
	<b>TinyImageNet</b>	19.99	24.94	27.49	26.94	34.49
<b>0.01</b>	<b>PathMNIST</b>	6.26	6.69	6.67	6.5	6.56
<b>0.02</b>	<b>PathMNIST</b>	3.31	4.56	2.98	2.96	3.07

**Influence of the of size  $\mathcal{D}_{\text{cal}}$ .** We examine how the performance is affected by the size of the calibration data. Here, we vary the proportion  $p$  of samples from  $\mathcal{D}_{\text{cal}}$  that are actually used, i.e., we select a subset  $\mathcal{D}_{\text{cal}}^p \subseteq \mathcal{D}_{\text{cal}}$ , where  $|\mathcal{D}_{\text{cal}}^p| = p \cdot |\mathcal{D}_{\text{cal}}|$ . Then, as before,  $\mathcal{D}_{\text{cal}}^p$  is split into half for cell computation, and re-calibration. Figure D.12 presents the set sizes for different values of  $p$ . It can be seen that the proposed method is almost always preferable, with its advantage becoming more pronounced as the size of the calibration data increases.

**Influence of the of size  $\mathcal{D}_{\text{cells}}$ .** We investigated the effect of  $k$ , the size of  $\mathcal{D}_{\text{cells}}$ , on the results. Recall that  $\mathcal{D}_{\text{cells}}$  is responsible to the definition of the cells in Eq. (8) and the computation of the ratio-based scores in Eq. (9). In this experiment, the dataset is divided, as before, into three fixed

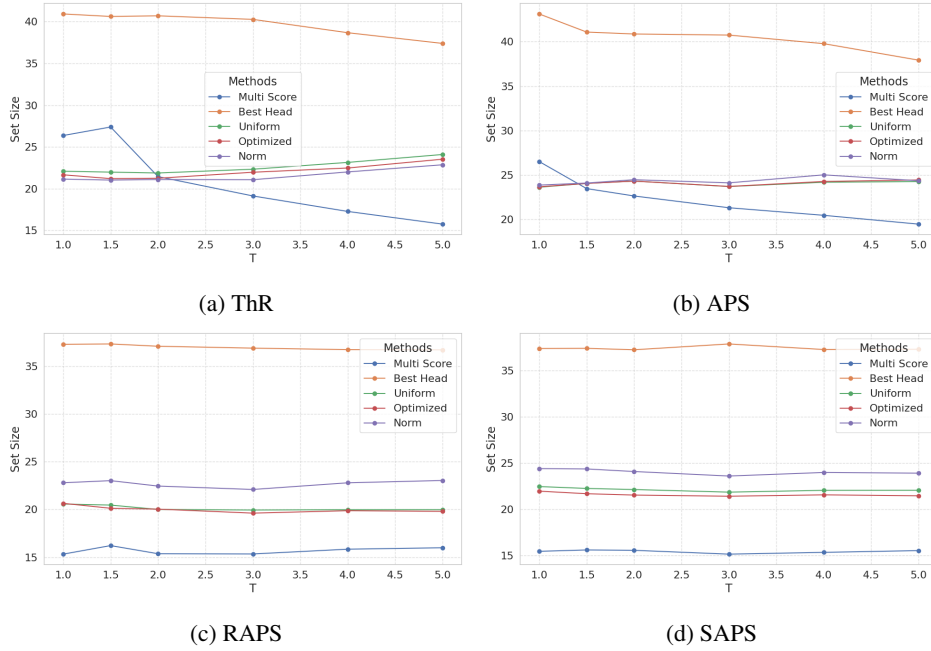


Figure D.6: Set size as a function of the temperature for different nonconformity scores. Results are shown for CIFAR100 with 7 heads and  $\alpha = 0.1$

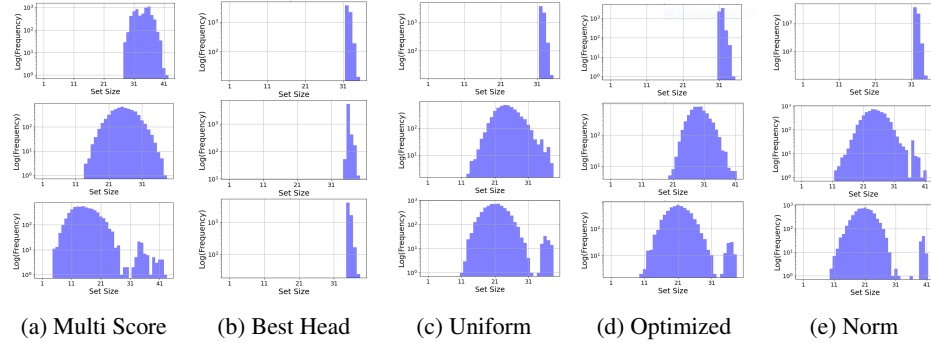


Figure D.7: Histograms of the set size distribution across different numbers of classification heads with RAPS score on Cifar100 and  $\alpha=0.1$ . Rows represent 1, 2, and 7 heads, while columns correspond to different methods.

subsets:  $\mathcal{D}_{\text{cells}}$ ,  $\mathcal{D}_{\text{re-cal}}$ , and  $\mathcal{D}_{\text{test}}$ . We vary the proportion  $p$  of samples from  $\mathcal{D}_{\text{cells}}$  that are actually used, i.e., we select a subset  $\mathcal{D}_{\text{cells}}^p \subseteq \mathcal{D}_{\text{cells}}$ , where  $|\mathcal{D}_{\text{cells}}^p| = p \cdot |\mathcal{D}_{\text{cells}}|$ . The set sizes obtained for different values of  $p$  are presented in Table D.3. As expected, the set size increases as  $p$  decreases. However, the overall behavior remains stable, with only 9.4 – 13.3% increase in set size for 50% reduction in the number of samples in  $\mathcal{D}_{\text{cells}}$ . In addition, we compare to the Optimized baseline, where the set  $\mathcal{D}_{\text{cells}}^p$  is used for optimizing the weights for combining the scores. We see for all  $p$  values Multi-score is advantageous over Optimized.

**Influence of the nonconformity score’s parameters .** We examined the influence of the different parameters of the non-conformal RAPS and SAPS scores on the methods performances . The results for the RAPS score are detailed in Table D.4, while those for the SAPS score are presented in Table D.5. Overall, our findings indicate that the Multi-Score method remains consistently advantageous, regardless of the parameter settings.

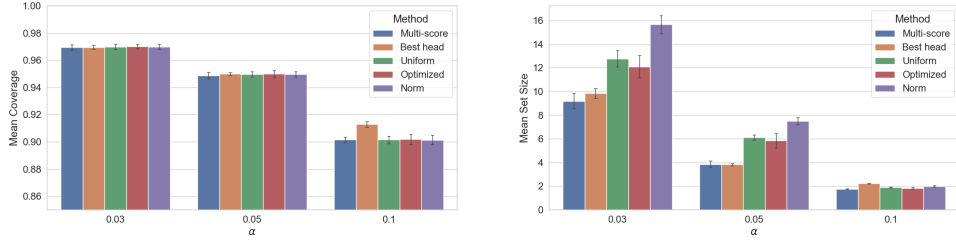


Figure D.8: Results for ImageNet with a standard ensemble of five different models, using RAPS score.

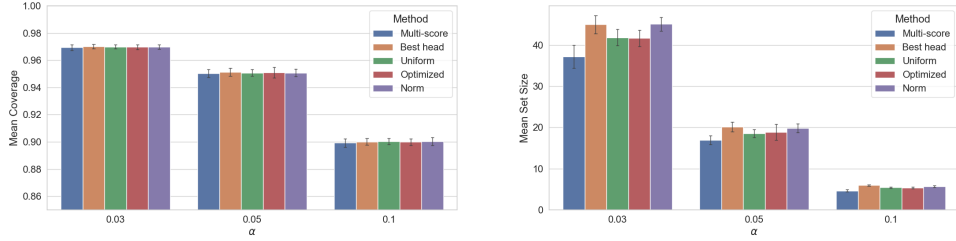


Figure D.9: Results for ImageNet with test-time augmentations, using RAPS score.

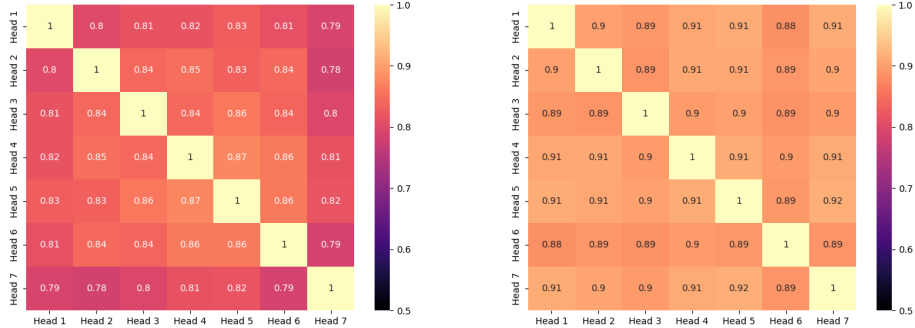


Figure D.10: Similarity matrix between prediction heads for heads trained with (left) or without (right) diversity regularization on CIFAR100 dataset.

Table D.3: Sensitivity of the results to the size of  $\mathcal{D}_{\text{cells}}$ . We select a subset  $\mathcal{D}_{\text{cells}}^p \subseteq \mathcal{D}_{\text{cells}}$ , where  $|\mathcal{D}_{\text{cells}}^p| = p \cdot |\mathcal{D}_{\text{cells}}|$ . The set sizes obtained for different values of  $p$  are reported for CIFAR100 dataset and 7 heads.

		$\alpha$	100%	95%	90%	85%	80%	70%	60%	50%	25%	10%
Multi Score	RAPS	0.1	15.32	15.66	15.68	15.69	15.92	16.26	16.68	17.32	19.04	22.13
		0.2	2.33	2.33	2.34	2.35	2.37	2.49	2.54	2.64	3.09	4.16
	SAPS	0.1	15.46	15.46	15.76	15.72	15.88	16.04	16.32	16.91	19.01	22.17
		0.2	2.36	2.36	2.37	2.41	2.42	2.54	2.6	2.64	3.09	3.61
Optimized	RAPS	0.1	22.47	22.27	22.04	22.09	21.85	21.71	22.47	22.52	23.4	24.2
		0.2	6.53	6.31	6.21	6.3	6.73	6.49	6.14	6.55	6.54	7.21
	SAPS	0.1	21.8	21.83	21.71	21.68	21.51	21.63	21.55	21.66	23.24	25.11
		0.2	6.8	6.45	6.4	6.45	6.31	6.14	5.93	5.93	7.4	8.611

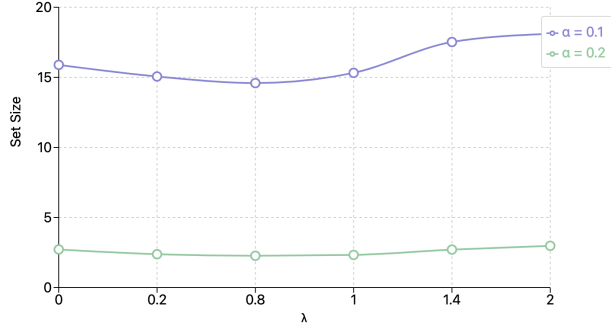


Figure D.11: Set size as a function of the regularization parameter  $\lambda$  for CIFAR100 with RAPS score and 7 heads.

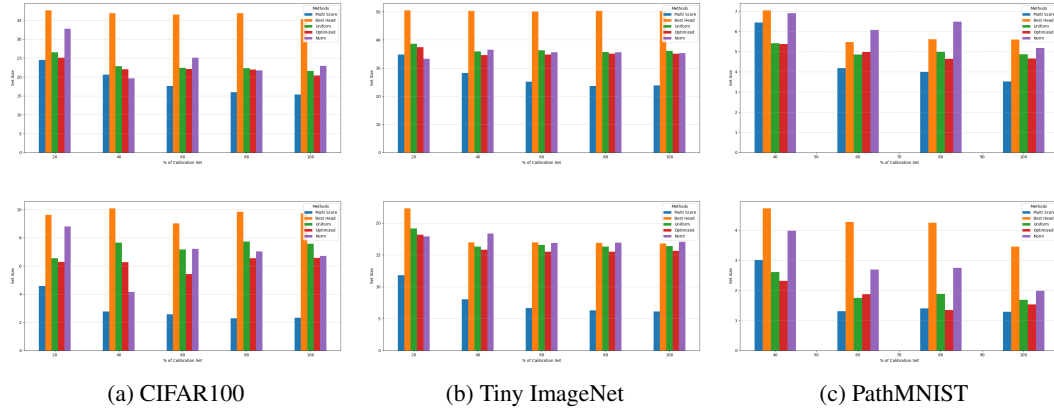


Figure D.12: Sensitivity of the results to the size of  $\mathcal{D}_{\text{cal}}$ . We select a subset  $\mathcal{D}_{\text{cal}}^p \subseteq \mathcal{D}_{\text{cal}}$ , where  $|\mathcal{D}_{\text{cal}}^p| = p \cdot |\mathcal{D}_{\text{cal}}|$ . The set sizes obtained for different values of  $p$  are reported for RAPS score on 7 heads.

**Influence of the underlying model.** To demonstrate that our method improves the efficiency of conformal prediction (CP) regardless of the underlying model, we conducted an experiment using a more powerful ViT model Dosovitskiy et al. (2020). Results are presented on table D.6, showing similar trends to those observed in our main results.

#### D.4 RESULTS ON OTHER DATASETS

**Results on text data classification.** To show that our method can be applied across different types of data, we conducted an experiment with a text classification task. We use the 20 Newsgroups dataset, which comprises newsgroup posts on 20 topics. We use a BERT-base model (Devlin et al., 2019), and attach additional classification heads in a similar way to the other models. The results on Table D.7 show that Multi-Score outperforms all the baselines. Here the norm baseline is the closest in the performance to the proposed method.

**Results for ImageNet.** Table D.8 presents the results obtained for ImageNet. Here the classification heads consist of a single linear layer, and are all initialized by the weights of the pretrained model. Here too we can see the benefit of the proposed method over the baselines.

Table D.4: Set size comparison between the baselines for different values of  $\kappa$  and  $\lambda$ , used in RAPS score, on CIFAR100,  $\alpha = 0.1$  and 7 heads.

$\kappa \backslash \lambda$	0.001	0.01	0.1	1.0
<b>Multi Score</b>				
1	12.79	14.83	15.10	15.45
2	12.84	14.26	15.11	15.53
3	12.88	14.29	15.30	15.54
4	12.81	14.45	15.36	15.53
5	12.80	14.44	15.51	15.70
<b>Best Head</b>				
1	32.04	32.11	32.11	32.11
2	24.41	24.00	24.00	25.00
3	23.01	23.02	23.04	22.79
4	20.77	21.12	21.63	22.70
5	20.32	21.21	23.60	22.60
<b>Uniform</b>				
1	24.38	24.76	24.91	24.96
2	24.41	23.59	24.87	25.05
3	23.10	23.75	25.44	26.14
4	21.20	21.84	23.26	23.33
5	20.40	21.41	22.63	22.75
<b>Optimized</b>				
1	24.45	24.89	25.26	25.84
2	24.52	24.00	24.76	25.10
3	23.09	23.32	24.64	25.07
4	20.82	21.20	22.30	22.70
5	20.34	21.40	22.80	22.60
<b>Norm</b>				
1	18.01	19.62	19.80	22.59
2	18.01	22.99	24.53	24.66
3	22.82	23.05	27.74	28.50
4	20.80	21.21	21.75	23.60
5	20.34	21.21	23.86	25.04

Table D.5: Set size comparison between the baselines for different values of  $\xi$ , used in SAPS score, on CIFAR100,  $\alpha = 0.1$  and 7 heads.

$\xi$	0.01	0.05	0.1	0.5	1.0
<b>Multi Score</b>					
	14.7	16.17	16.28	16.89	17.06
<b>Best Head</b>					
	36.49	37.11	37.21	37.43	37.48
<b>Uniform</b>					
	19.15	21.73	21.91	22.55	22.57
<b>Optimized</b>					
	18.51	21.0	21.60	21.90	22.07
<b>Norm</b>					
	18.61	21.31	23.12	24.56	24.59

## D.5 ADDITIONAL BASELINES

**L1 Norm baseline.** We examined an additional baseline, where we use the L1 norm instead of the L2 norm. Table D.9 compares the two baselines for CIFAR-100 dataset and RAPS score. We observe that both baselines obtain similar performance.



Table D.6: ResNet vs. ViT model set size comparison on CIFAR100 with  $\alpha = 0.1$  and 7 heads.

Model	Acc.	Score	Methods				
			Multi Score	Best Head	Uniform	Optimized	Norm
ViT	0.8	RAPS	4.34	22.11	16.21	15.93	5.51
		SAPS	2.33	12.77	2.43	2.45	2.60
ResNet50	0.69	RAPS	15.32	35.29	24.37	22.47	22.49
		SAPS	15.46	37.41	22.05	21.80	22.79

Table D.7: Comparison of set sizes across different methods on the 20 Newsgroups dataset, with 7 heads, using the BERT-base model fine-tuned for news topic classification.

$\alpha$	Score	Methods				
		Multi Score	Best Head	Uniform	Optimized	Norm
0.08	RAPS	1.46	9.61	9.58	9.49	1.5
	SAPS	1.47	10.66	2.09	2.18	2.35
0.05	RAPS	2.29	10.1	10.01	10.07	2.32
	SAPS	1.97	10.9	2.4	2.46	2.62
0.02	RAPS	4.28	11.73	11.73	11.47	4.36
	SAPS	3.28	11.57	3.45	3.34	3.65

Table D.8: Comparison of set sizes across different methods on the ImageNet dataset, with 7 heads, demonstrating the performance with a large-number-of classes.

$\alpha$	Score	Methods				
		Multi Score	Best Head	Uniform	Optimized	Norm
0.1	RAPS	4.1	4.36	13.47	14.63	13.47
	SAPS	4.94	5.36	13.85	14.99	13.85
0.2	RAPS	1.82	2.4	2.1	2.1	2.1
	SAPS	1.57	1.68	2.04	2.06	2.04

Table D.9: Comparison between L1 and L2 Norm methods on CIFAR100 and 7 heads.

$\alpha$	Score	Methods		
		L1	L2	Multi Score
0.1	RAPS	22.98	22.49	15.32
	SAPS	22.25	22.79	15.46
0.2	RAPS	6.75	7.51	2.33
	SAPS	7.62	7.85	2.36

**Comparison to vanilla baseline.** We examined a vanilla baseline, where the CP procedure is performed over the original classification head, without the addition of multiple classification heads. We observe that the results are similar to the Best Head baseline defined above. Table D.10 summarizes the results

## D.6 VARIANTS OF MULTI-SCORE CONFORMAL PREDICTION

**Jackknife+ Multi-score conformal prediction.** We compare our split version in Algorithm 1 to the jackknife+ version in Algorithm B.1. We obtained that the required  $1-\alpha$  coverage is achieved in both settings, and the set sizes are summarized in Table D.11. As expected jackknife+ obtains

Table D.10: Comparison between the vanilla baseline and the multi-score methods on CIFAR100, with RAPS score and 7 heads.

$\alpha$	Set Size		
	Multi Score	Vanilla RAPS	Best Head
<b>0.1</b>	15.32	39.27	35.29
<b>0.2</b>	2.33	8.09	9.37

Table D.11: Comparison between Jackknife+ and Multi score methods on CIFAR100 and 7 heads.

$\alpha$	Score	Methods	
		Multi Score Jackknife+	Multi Score
<b>0.1</b>	<b>RAPS</b>	14.85	15.32
	<b>SAPS</b>	14.61	15.46
<b>0.2</b>	<b>RAPS</b>	2.26	2.33
	<b>SAPS</b>	2.22	2.36

Table D.12: Comparison of soft centers approach with different number of neighbors  $b$  on CIFAR100 and 7 heads.

$\alpha$	Score	Methods				
		$b = 200$	$b = 100$	$b = 50$	$b = 10$	Multi Score ( $b = 1$ )
<b>0.1</b>	<b>RAPS</b>	23.68	22.73	19.32	16.48	15.32
	<b>SAPS</b>	23.53	22.44	19.12	17.14	15.46
<b>0.2</b>	<b>RAPS</b>	3.01	2.89	2.75	2.26	2.33
	<b>SAPS</b>	3.04	2.87	2.81	2.29	2.36

smaller set sizes. However, the improvement appears to be insignificant in this case and may not justify the additional computational cost.

**Soft Multi-score conformal prediction.** We evaluate the soft version of our proposed approach. The set sizes obtained for different number of neighbors are summarized in Table D.12. We observe that in almost all cases  $b = 1$  (Algorithm 1) is preferable.