
Improving Coverage in Combined Prediction Sets with Weighted p-values

Anonymous Author
Anonymous Institution

Abstract

Conformal prediction quantifies the uncertainty of machine learning models by augmenting point predictions with valid prediction sets. For complex scenarios involving multiple trials, models, or data sources, conformal prediction sets can be aggregated to create a prediction set that captures the overall uncertainty, often improving precision. However, aggregating multiple prediction sets with individual $1 - \alpha$ coverage inevitably weakens the overall guarantee, typically resulting in $1 - 2\alpha$ worst-case coverage. In this work, we propose a framework for the *weighted aggregation of prediction sets*, where weights are assigned to each prediction set based on their contribution. Our framework offers flexible control over how the sets are aggregated, achieving tighter coverage bounds that interpolate between the $1 - 2\alpha$ guarantee of the combined models and the $1 - \alpha$ guarantee of an individual model depending on the distribution of weights. Importantly, our framework generalizes to data-dependent weights, as we derive a procedure for weighted aggregation that maintains finite-sample validity even when the weights depend on the data. This extension makes our framework broadly applicable to settings where weights are learned, such as mixture-of-experts (MoE), and we demonstrate through experiments in the MoE setting that our methods achieve adaptive coverage.

1 INTRODUCTION

In recent years, machine learning models have achieved remarkable accuracy across a variety of predictive tasks

(Min et al., 2023; Liang et al., 2024a). Understanding the uncertainty associated with each prediction is essential to decision-making in real-world scenarios, but the black box nature of many machine learning models hinders their deployment in safety-critical applications such as medical diagnosis (Chua et al., 2023; Grote and Berens, 2023), industrial control systems (Kumar et al., 2023; Lawrence et al., 2024), and extreme weather forecasting (Eyring et al., 2024; Lai et al., 2024). Conformal prediction (Vovk et al., 2005) has emerged as a popular wrapper method around machine learning models because it provides a statistically valid quantification of uncertainty. Specifically, it transforms point predictions to prediction sets with finite-sample coverage guarantees, as long as the test data is exchangeable with the training data. Complex scenarios involving multiple predictions—for example, when there are multiple trials, models, or data sources—naturally produce multiple conformal prediction sets (Figure 1).

A number of methods have been proposed to aggregate prediction sets. For conformal prediction, where the popular split conformal variant (Papadopoulos et al., 2002) introduces a one-time random split, there are multiple methods of aggregating predictions to reduce the randomness over multiple splits: popular examples include cross-conformal (Vovk, 2015), CV+ (Barber et al., 2021), and jackknife+ (Barber et al., 2021). These aggregation methods are all *symmetric*, in that individual prediction sets contribute equally to the aggregate set. Although comparatively less studied, it is often useful to aggregate sets *asymmetrically*—that is, to weight sets based on their prior importance to the overall result (Gasparin and Ramdas, 2024). Work in both symmetric and asymmetric aggregation establishes that aggregating individual prediction sets with $1 - \alpha$ coverage guarantees results in an overall coverage guarantee of $1 - 2\alpha$ (Vovk and Wang, 2020).

Our work is based on the observation that if the overall coverage guarantee for asymmetric aggregation reflects the contributions of the individual sets, then we can achieve a tighter guarantee than a constant $1 - 2\alpha$. Consider, as an extreme, the case where almost all

Preliminary work. Under review by AISTATS 2026. Do not distribute.

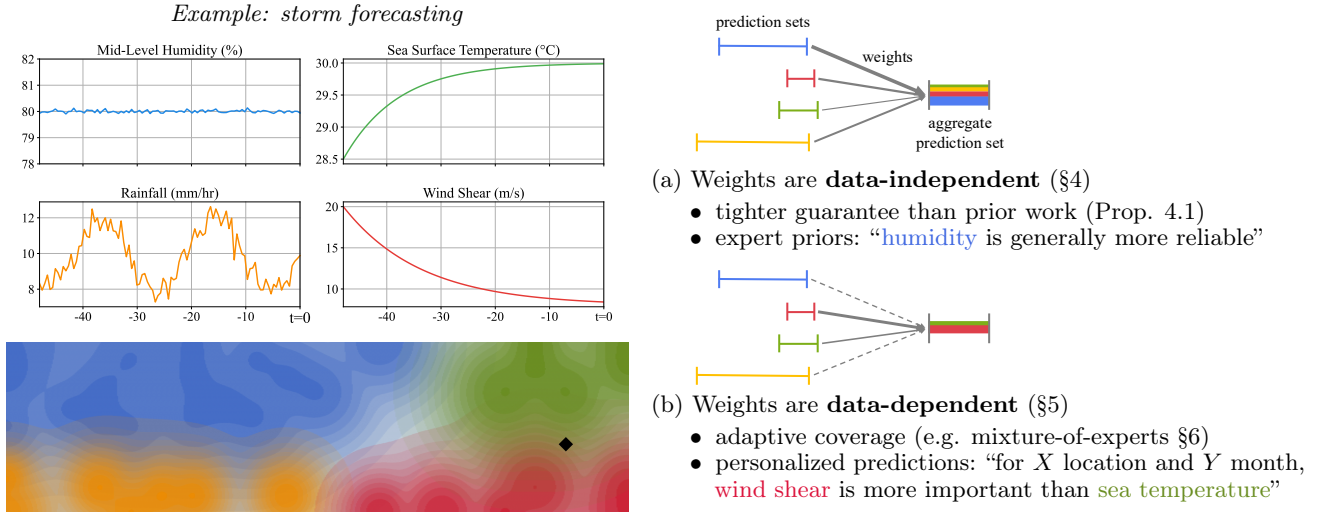


Figure 1: **Left:** Storm forecasting example with different models tracking humidity, sea temperature, rainfall, and wind shear over time. Below, an abstract representation of how models vary in predictive strength across the input space (colored regions). At the given test point (black), the red and green models dominate, so their prediction sets matter most. **Right:** Model prediction sets are combined with weighted aggregation. (a) *Data-independent* weights reflect expert priors (e.g. up-weighting the blue model as it provides the most general coverage). (b) *Data-dependent* weights adapt to context, yielding forecasts better aligned with current conditions (where red and green dominate).

the weight is on one set; then, because the coverage of that set dominates the aggregation, the overall coverage guarantee should be closer to the $1 - \alpha$ guarantee of the dominant individual set. Accordingly, we propose a method for asymmetric aggregation with *data-independent weights*, where the overall coverage guarantee is based on the distribution of weights across the prediction sets. This approach leverages the results of Vovk and Wang (2020) on averaging p-values to achieve tighter guarantees when there exists strong asymmetry in the importance of the prediction sets—i.e. when a single prediction set is significantly more important than the rest—with $1 - 2\alpha$ as the worst-case guarantee when the contributions of the prediction sets are more equal. In this way, we can incorporate expert priors with data-independent weights to potentially tighten coverage (Figure 1a).

Weights from expert priors provide added flexibility and guidance to prediction set aggregation. In practice, however, many applications go beyond fixed weights, and require that weights adapt directly to the data. For example, in the popular mixture-of-experts setting (Jacobs et al., 1991), expert weights are learned from inputs. To deal with this general case, we additionally propose a method for aggregating prediction sets with *data-dependent weights* via a linear transformation on the weighted average of p-values associated with each set. Our method allows us to construct a valid aggregate prediction set that preserves the proportions of individual set weights. Since data-dependent weights adjust the contributions of the individual sets based on the observed input, our method also achieves a form

of conditional/data-adaptive coverage (Figure 1b). We apply our method to the mixture-of-experts setting, and we demonstrate its effectiveness in experiments with real and synthetic data.

To summarize, our main contributions are as follows:

- We propose a framework for the asymmetric aggregation of prediction sets based on weighted p-values (§4). With *data-independent weights*, our framework improves the coverage guarantee beyond the standard $1 - 2\alpha$ (Prop. 4.1).
- We derive a general method for transforming a random variable to a p-value variable, which enables the construction of a prediction set with finite-sample guarantees (§5). This method allows us to extend our framework to *data-dependent weights*, and we demonstrate that incorporating data-dependent weights provides adaptive coverage (Prop. 5.1, Prop. 5.2).
- We apply our method to the mixture-of-experts setting, and we demonstrate that using a weighted aggregation of experts improves local validity (§6).

2 BACKGROUND

We begin by reviewing preliminaries leading up to the basics of conformal prediction, with a focus on the widely-used split conformal method. To provide context for our method based on p-values, we establish the connection between the more typical quantile presentation of conformal prediction and its original p-value presentation (Vovk et al., 2005).

Quantiles For a set of n elements $Z = \{z_1, \dots, z_n\}$, the left α -empirical quantile is given by

$$\hat{Q}_\alpha^-(Z) = \text{the } \lfloor (n+1)\alpha \rfloor\text{-th smallest value of } Z,$$

and the right α -empirical quantile, or the $(1 - \alpha)$ -empirical quantile, by

$$\hat{Q}_\alpha^+(Z) = \text{the } \lceil (n+1)(1 - \alpha) \rceil\text{-th smallest value of } Z,$$

where the $n+1$ term serves as a finite sample correction.

p-values A p -variable (Vovk and Wang, 2020) is a random variable P' such that

$$\mathbb{P}\{P' \leq \alpha\} \leq \alpha \quad \forall \alpha \in (0, 1). \quad (1)$$

The values taken by a p -variable are called p -values. In hypothesis testing, p -values represent the probability of observing results at least as extreme as the ones obtained, assuming the null hypothesis is true.

Conformal prediction Suppose we have training data $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ and a test point (X_{n+1}, Y_{n+1}) with unknown label Y_{n+1} . Let \mathcal{P} denote the joint distribution on $(X, Y) \in \mathcal{X} \times \mathcal{Y}$. Then, assuming that the training and test data are *exchangeable*, or that the distribution of the training and test data is permutation-invariant, conformal prediction can be used to construct a valid prediction set for X_{n+1} with no further assumptions on \mathcal{P} .

Split conformal prediction, a popular variant, starts by splitting the training data indices into disjoint “proper” training and calibration subsets so that $\{1, \dots, n\} = S_{\text{train}} \cup S_{\text{cal}}$. We fit a predictor to the subset S_{train} , so that $\hat{\mu} = \mathcal{A}(\{(X_i, Y_i)\}_{i \in S_{\text{train}}})$, where \mathcal{A} is a model fitting algorithm. Then, given a nonconformity score function $s : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, we can compute the scores $R_i = s(X_i, Y_i)$, or the values of the score function that characterize how nonconformal a label Y_i is from its predicted value $\hat{\mu}(X_i)$. We compute the scores on the subset S_{cal} . At a given significance level α , this procedure allows us to define the prediction set

$$\begin{aligned} \hat{C}_\alpha^{\text{split}}(X_{n+1}) &= \{y \in \mathcal{Y} : \\ &\quad s(X_{n+1}, y) \leq \hat{Q}_\alpha^+(\{s(X_i, Y_i)\}_{i \in S_{\text{cal}}})\}. \end{aligned}$$

By construction, the split conformal prediction set results in the (marginal) coverage guarantee

$$\mathbb{P}\{Y_{n+1} \in \hat{C}_\alpha^{\text{split}}(X_{n+1})\} \geq 1 - \alpha \quad \forall \alpha \in (0, 1).$$

In other words, the prediction set $\hat{C}_\alpha^{\text{split}}(X_{n+1})$ is *valid*, or guaranteed to contain the true label Y_{n+1} with probability at least $1 - \alpha$.

As an alternative perspective, the quantile can be replaced with the p -value function

$$\hat{p}(x, y) = \frac{1 + \sum_{i \in S_{\text{cal}}} \mathbb{1}\{s(x, y) < s(X_i, Y_i)\}}{|S_{\text{cal}}| + 1},$$

which returns the proportion of calibration points that are less conformal than some test point (x, y) , with a finite sample correction. The equivalent prediction set constructed from the p -value function is

$$\hat{C}_\alpha^{\text{split}}(X_{n+1}) = \{y \in \mathcal{Y} : \hat{p}(X_{n+1}, y) > \alpha\}.$$

3 RELATED WORK

We categorize existing work on combining conformal results into two approaches: methods for selecting or combining prediction sets, and methods for combining p -values from multiple runs of the conformal procedure.

Combining prediction sets Combining conformal prediction sets was first introduced by Lei et al. (2018), who propose combining K split conformal prediction sets by taking their intersection. Taking the intersection of these sets reduces the size of the resulting prediction set, but the authors demonstrate that, under general conditions, the intersection is wider than the individual sets with probability tending to 1 with an asymptotic number of samples. Another method proposed by Yang and Kuchibhotla (2021) selects the predictor (from a set of K predictors) that returns the best prediction set. They present two methods: one to select the predictor that gives the most efficient prediction set, but with only approximate validity, and one to select the predictor with the most valid prediction set (with minimal coverage slack), but with a width only close to the minimum. Liang et al. (2024b) point out that the approximate validity of the first method is due to selection bias, and introduce an alternative approach that both selects the most efficient predictor and uses it to construct a valid prediction set.

A recent line of work explores combining conformal prediction sets by majority vote with a $1 - 2\alpha$ coverage guarantee (Cherubin, 2019; Solari and Djordjilović, 2021; Gasparin and Ramdas, 2024). In particular, Gasparin and Ramdas (2024) characterize the width of the majority vote set and introduce many extensions, including a weighted majority vote method that incorporates prior information in the weights. Our work on combining conformal results by weighted p -values is related to this method. However, our formulation allows us to leverage the results of Vovk and Wang (2020) to improve the coverage guarantee beyond the standard $1 - 2\alpha$, and can also be generalized to data-dependent weights to enable adaptive coverage. (For a more detailed comparison, see §C.1.)

Combining p-values A substantial body of research has been devoted to developing methods for combining p -values. Here, we focus specifically on approaches that accommodate arbitrary dependence, with a particular emphasis on their applications to conformal

prediction. (For a more comprehensive review, see Balasubramanian et al. (2015) and DiCiccio et al. (2020).) An early example of combining p-values with arbitrary dependence is the Bonferroni method, where the minimum of a set of p-values is scaled by the number of p-values. This method has multiple extensions (Rüger, 1978; Hommel, 1983), and was first applied to conformal prediction by Lei et al. (2018). Rüschendorf (1982) finds that twice the average of p-values is a p-value; this was later extended to a more general notion of average by Vovk and Wang (2020). Stutz et al. (2023) use the result of Rüschendorf (1982) to get a $1 - 2\alpha$ guarantee for the average of p-values, but also propose a novel transformation to get a p-value average with improved coverage.

4 COMBINING CONFORMAL PREDICTION SETS BY WEIGHTED P-VALUE

In this section, we extend the method of combining p-values, as proposed by Vovk and Wang (2020), to the setting of conformal prediction. Suppose we have K predictors, denoted by $\hat{\mu}_1, \dots, \hat{\mu}_K$, where each predictor $\hat{\mu}_k$ is associated with a nonconformity score function s_k and calibration set S_k . For each predictor, we define a p-value function $\hat{p}_k(x, y)$ that measures how well a candidate label y conforms to the predicted outcome for a given x . Specifically, for each k , the p-value for a given point is given by

$$\hat{p}_k(x, y) = \frac{1 + \sum_{i \in S_k} \mathbb{1}\{s_k(x, y) < s_k(X_i, Y_i)\}}{|S_k| + 1}. \quad (2)$$

The prediction set for X_{n+1} is then

$$\hat{C}_{\alpha, k}^{\text{split}}(X_{n+1}) = \{y \in \mathcal{Y} : \hat{p}_k(X_{n+1}, y) > \alpha\},$$

with a guaranteed marginal coverage

$$\mathbb{P}\{Y_{n+1} \in \hat{C}_{\alpha, k}^{\text{split}}(X_{n+1})\} \geq 1 - \alpha \quad \forall \alpha \in (0, 1). \quad (3)$$

Consider assigning a weight v_k to each p-value function \hat{p}_k . The weighted average p-value function is then

$$\bar{p}(x, y) = \sum_{k=1}^K v_k \hat{p}_k(x, y) \quad \text{where} \quad \sum_{k=1}^K v_k = 1, \quad (4)$$

with prediction set

$$\hat{C}_{\alpha}^{\text{avg}}(X_{n+1}) = \{y \in \mathcal{Y} : \bar{p}(X_{n+1}, y) > \alpha\}. \quad (5)$$

We now provide a coverage guarantee for this aggregated prediction set.

Proposition 4.1. *Let $\hat{C}_{\alpha, 1}^{\text{split}}(X_{n+1}), \dots, \hat{C}_{\alpha, K}^{\text{split}}(X_{n+1})$ be K prediction sets defined by p-value functions $\hat{p}_1, \dots, \hat{p}_K$ (2) on X_{n+1} , where $1 - \alpha$ coverage (3) holds for each set $k \in [K]$. Then, the prediction set $\hat{C}_{\alpha}^{\text{avg}}(X_{n+1})$ (5) from thresholding the weighted average p-value function (4) gives the coverage guarantee*

$$\mathbb{P}\{Y_{n+1} \in \hat{C}_{\alpha}^{\text{avg}}(X_{n+1})\} \geq 1 - \min\left\{\frac{1}{v}, 2\right\} \alpha$$

for all $\alpha \in (0, 1)$, where $v = \max\{v_1, v_2, \dots, v_K\}$ is the largest weight assigned to any of the p-values.

Weighted aggregation provides more flexible coverage guarantees This result provides a spectrum of coverage guarantees based on the weight distribution among the models. When one predictor dominates (i.e., $v > 1/2$), the guarantee improves beyond the standard $1 - 2\alpha$ of other aggregation methods, approaching $1 - \alpha$ as $v \rightarrow 1$ (recovering split conformal at $v = 1$). Thus, we can interpolate between the standard $1 - 2\alpha$ guarantee of combined models and the $1 - \alpha$ guarantee of individual models, with the weights controlling the trade-off. This improves the guarantee for asymmetric aggregation, and also opens up the utility of asymmetric generalizations of the many established symmetric aggregation methods (e.g. set-weighted versions of cross-conformal, CV+, etc.)

Since the work of Vovk and Wang holds for arbitrarily dependent p-values, this method is robust across a wide range of scenarios. In practice, independent weights allow users to incorporate prior knowledge about the relative quality of different predictors (Vovk and Wang, 2020; Gasparin and Ramdas, 2024). The weighting can reflect, for example, expert insights on which of K models should be prioritized as being more reliable.

From p-values to prediction sets The work of Vovk and Wang is central to our result, allowing us to generalize and improve upon existing conformal guarantees. Still, despite its broad applicability, their work remains underexplored in conformal literature, and the weight-dependent coverage result has not yet been applied to conformal prediction sets.¹ We attribute this oversight to several factors.

Vovk and Wang frame their work in terms of merging functions and p-values, without reference to prediction sets or conformal prediction. As a result, subsequent research has similarly focused on p-values and related test statistics, with limited connection to conformal literature. Meanwhile, conformal literature typically uses a quantile-based construction of prediction sets, rather

¹To the best of our knowledge, Gasparin and Ramdas (2024) is the only work so far to use asymmetric set aggregation, and they derive the standard $1 - 2\alpha$ guarantee of symmetric aggregation methods.

than a p-value construction, and work in conformal prediction set aggregation typically operates directly on the sets, rather than working with the associated p-value functions. These trends may contribute to the pattern where p-value results are not always propagated to the wider conformal prediction community. Our result shows that the p-value presentation offers unique benefits to prediction set aggregation, and we hope that our work may encourage renewed interest in the connection between p-values and conformal prediction.

5 DATA-DEPENDENT WEIGHTS

In many practical settings, model weights are determined by the observed data rather than fixed in advance. Such data-dependent weights naturally adapt the influence of each model to the characteristics of the input, making them central to applications like ensemble learning and mixture-of-experts. This adaptivity, however, creates a technical challenge: once the weights depend on the data, they also depend on the associated p-variables, and the theory of Vovk and Wang (2020) no longer applies. Nevertheless, we can build on their approach to develop a method that allows us to recover a valid coverage guarantee, even when the weights are data-dependent.

The key idea is to directly use the definition of a p-variable (1). The sum of weighted p-variables, where the weights are dependent on the p-variables, is not necessarily a p-variable. However, the weighted sum can *become* a p-variable by a linear transformation that both satisfies the definition of a p-variable and preserves the proportions of the weights.

Let the p-variables of the predictors be $P_1, \dots, P_K \in \mathcal{U}$, where \mathcal{U} is the set of all uniformly distributed random variables. The weights are given by a random vector $\mathbf{W} = (W_1, \dots, W_K)$ in the $(K-1)$ -dimensional simplex $\Delta_{K-1} := \{\mathbf{w} = (w_1, \dots, w_K) \in [0, 1]^K : w_1 + \dots + w_K = 1\}$ depending on the data $\{(X_i, Y_i)\}_{i \in [n]} \cup \{X_{n+1}\}$. We define the weighted average function

$$p_{\text{all}}(x, y; \mathbf{w}) = \sum_k w_k \hat{p}_k(x, y), \quad (6)$$

giving associated random variable $P_{\text{all}} := p_{\text{all}}(X, Y; \mathbf{W})$ with distribution $F_{P_{\text{all}}}$. We propose to learn the scalar

$$\begin{aligned} m^* &= \inf\{m \in \mathbb{R}^+ : \mathbb{P}\{m P_{\text{all}} \leq \alpha\} \leq \alpha \quad \forall \alpha \in (0, 1)\} \\ &= \inf\{m \in \mathbb{R}^+ : F_{P_{\text{all}}}(\alpha/m) \leq \alpha \quad \forall \alpha \in (0, 1)\} \\ &= \sup_{\delta \in (0, \infty)} \frac{F_{P_{\text{all}}}(\delta)}{\delta}. \end{aligned} \quad (7)$$

(See §A.5 for a full derivation.) The scaling defined in (7) transforms the weighted average P_{all} into a valid p-variable, recovering a coverage guarantee.

Algorithm 1 Constructing a weighted aggregate of prediction sets with valid coverage

- 1: **Input:** Data with indices $\{1, \dots, n\} = S_{\text{cal}} \cup S_{\text{merge}}$, K predictors from potentially overlapping datasets, K weights which may depend on data (e.g. learned router weights), test example X_{n+1} .
- 2: **Output:** Prediction set around the test example $\hat{C}_{\alpha}^{\text{scaled}}(X_{n+1})$, with coverage of at least $1 - (\alpha + \epsilon + \delta)$.
- 3: **Step 1.** Derive the p-value function p_{all} .
- 4: 1.1. Using calibration set S_{cal} , derive p-value functions \hat{p}_k for each of the predictors (2).
- 5: 1.2. Compute the aggregated p-value function $p_{\text{all}} = \sum_{k=1}^K w_k \hat{p}_k$, using the (potentially data-dependent) weights w_k .
- 6: **Step 2.** Compute the correction factor \hat{m}^* .
- 7: 2.1. Using the points of merging set S_{merge} with the function p_{all} , get samples of $\hat{F}_{P_{\text{all}}}$.
- 8: 2.2. Derive the correction factor \hat{m}^* from the samples of $\hat{F}_{P_{\text{all}}}$ (8).
- 9: **Step 3.** For test example X_{n+1} , construct the prediction set

$$\hat{C}_{\alpha}^{\text{scaled}}(X_{n+1}) = \{y : \hat{m}^* p_{\text{all}}(X_{n+1}, y) > \alpha\}.$$

Proposition 5.1 (Infinite-sample guarantee). *Let $\hat{C}_{\alpha,1}^{\text{split}}(X_{n+1}), \dots, \hat{C}_{\alpha,K}^{\text{split}}(X_{n+1})$ be K prediction sets defined by p-value functions $\hat{p}_1, \dots, \hat{p}_K$ (2) on X_{n+1} corresponding to p-variables P_1, \dots, P_K . Suppose $1 - \alpha$ coverage (3) holds for each set $k \in [K]$. Let (W_1, \dots, W_K) be a random vector in Δ_{K-1} depending on $\{(X_i, Y_i)\}_{i \in [n]} \cup \{X_{n+1}\}$, and let p_{all} be the weighted average (6), with random variable P_{all} . The prediction set from thresholding the corrected $m^* p_{\text{all}}$ is*

$$\hat{C}_{\alpha}^{\text{scaled}}(X_{n+1}) = \{y : m^* p_{\text{all}}(X_{n+1}, y) > \alpha\},$$

and it satisfies the coverage guarantee

$$\mathbb{P}\{Y_{n+1} \in \hat{C}_{\alpha}^{\text{scaled}}(X_{n+1})\} \geq 1 - \alpha$$

for all $\alpha \in (0, 1)$, where m^* is defined in (7).

From population to finite-sample guarantees By definition, m^* is the minimal scaling factor that makes P_{all} a valid p-variable, and its value is determined by the CDF $F_{P_{\text{all}}}$. Because $F_{P_{\text{all}}}$ is a population quantity that is not accessible in practice, we cannot evaluate m^* exactly. Instead, we define a computable proxy \hat{m}^* using an empirical CDF constructed from a designated *merging set* S_{merge} ².

²So named because S_{merge} is used to learn the correction

The empirical CDF of P_{all} computed from $\mathcal{S}_{\text{merge}}$ is

$$\hat{F}_{P_{\text{all}}}(\alpha) = \frac{1}{|\mathcal{S}_{\text{merge}}|} \sum_{i \in \mathcal{S}_{\text{merge}}} \mathbf{1}\{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) \leq \alpha\}.$$

The empirical correction factor is computed as

$$\hat{m}^* = \max_{i \in \mathcal{S}_{\text{merge}}} \frac{\hat{F}_{P_{\text{all}}}(p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}))}{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)})}, \quad (8)$$

following (7). (See §A.6 for why (8) is equivalent to (7) for an empirical CDF.) The computation of \hat{m}^* yields the following finite-sample guarantee:

Proposition 5.2 (Finite-sample guarantee). *Under the same assumptions as Proposition 5.1, fix a user-chosen failure probability $\delta \in (0, 1)$ and set*

$$\varepsilon = \sqrt{\frac{\log(2/\delta)}{2|\mathcal{S}_{\text{merge}}|}}.$$

The prediction set from thresholding $\hat{m}^* p_{\text{all}}$ is

$$\hat{C}_{\alpha}^{\text{scaled}}(X_{n+1}) = \{y : \hat{m}^* p_{\text{all}}(X_{n+1}, y) > \alpha\},$$

and it satisfies the coverage guarantee

$$\mathbb{P}\{Y_{n+1} \in \hat{C}_{\alpha}^{\text{scaled}}(X_{n+1})\} \geq 1 - (\alpha + \varepsilon + \delta)$$

for all $\alpha \in (0, 1)$, where \hat{m}^* is defined in (8).

(For an alternative formulation where the prediction set is defined directly from p_{all} and the factor m^* appears on the right-hand side of the coverage inequality—in parallel with Proposition 4.1—see §A.4.)

An important practical question is how the quality of the correction \hat{m}^* depends on the merging set $\mathcal{S}_{\text{merge}}$. In §B.1.1, we study this dependence by varying $|\mathcal{S}_{\text{merge}}|$, and we show that even a modest number of samples (< 200) is sufficient for accurate coverage. We summarize our full procedure in Algorithm 1.

Data-dependent weights give a form of conditional coverage Data-dependent weights allow the influence of each model to be adjusted based on how well it performs for a specific data point. This approach enables the construction of aggregated prediction sets that are tailored to the characteristics of the given data, which can be viewed as a form of *locally conditional coverage*. To be clear, true X-conditional coverage—where the coverage guarantee is conditioned on the current input—is impossible without additional distributional assumptions (Lei and Wasserman, 2014; Vovk, 2012).

However, we demonstrate in §6 that data-dependent weights allow us to create aggregated prediction sets

that makes the merging function (Vovk and Wang, 2020) for the weighted average empirically precise.

with data-adaptive coverage, which can greatly improve conditional validity in practice.

Achieving tighter guarantees Our scaling correction factor m^* provides a coverage guarantee that holds for all significance levels $\alpha \in (0, 1)$. However, guarantees on coverage for all α can lead to overly conservative prediction sets, which may be unnecessarily restrictive in practice when we do not need guarantees for every possible significance level. Thus, we propose two alternatives: for a specific significance level α' , we can learn a correction factor

$$m^{\dagger} = \inf\{m \in \mathbb{R}^+ : \mathbb{P}\{mP_{\text{all}} \leq \alpha\} \leq \alpha \quad \forall \alpha \in (0, \alpha']\}, \quad (9)$$

or the even stricter

$$m^{\ddagger} = \inf\{m \in \mathbb{R}^+ : \mathbb{P}\{mP_{\text{all}} \leq \alpha'\} \leq \alpha'\}. \quad (10)$$

(Note that (9) is not without precedent: for the $1 - 2\alpha$ coverage guarantees of aggregation methods like cross-conformal and jackknife+, $\alpha' = 0.5$ is the highest significance level of interest.)

6 APPLICATION: MIXTURE-OF-EXPERTS

Mixture-of-experts (MoE) is a machine learning framework designed to combine the predictions of multiple specialized models, called experts (Jacobs et al., 1991). Each expert in an MoE focuses on unique aspects of the problem by learning different representations of the input data. A central component of this framework is the routing network, which determines how to combine the experts' outputs. Specifically, for an input x , the routing network of a traditional (soft) MoE assigns weights $W_k(x)$ to each expert output $f_k(x)$, producing the final prediction as a weighted sum

$$f(x) = \sum_k W_k(x) f_k(x). \quad (11)$$

By learning how to route different inputs to the most appropriate experts, the routing network implicitly conditions the model's final prediction on the combination of experts that fits the given input. In this way, the routing network can be viewed as learning a form of conditional relationship between the input features and the expertise of each model.

The routing network's ability to learn data-dependent weights for each expert makes MoE a natural setting for applying our method of aggregating prediction sets by weighted p-values. We collect the routing weights into the simplex-valued vector $W(x) := (W_1(x), \dots, W_K(x))$; then, for a new input X_{n+1} , the routing network outputs the weight vector $\mathbf{W}^{(n+1)} :=$

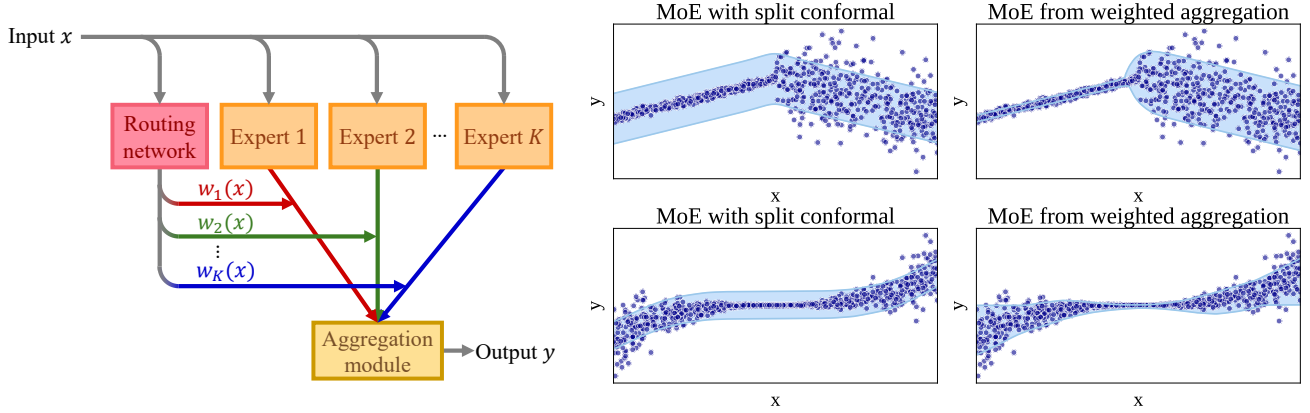


Figure 2: **Left:** Network diagram for MoE. For traditional MoE, the aggregation module takes a weighted sum of the outputs from each expert. To learn weight-dependent prediction sets, we instead propose to combine the prediction sets of each expert by weighted p-value. **Right, top row:** Comparison of split conformal prediction sets with those learned from weighted aggregation. Weighted aggregation allows overall coverage to follow the coverage of the dominant expert, rather than remain purely marginal. **Right, bottom row:** Another comparison of split conformal with weighted aggregation, with the latter showing local coverage with smooth transitions.

$W(X_{n+1})$, whose k th component we denote $W_k^{(n+1)}$. Let \hat{p}_k denote the p-value function of the k th expert. The weighted aggregate p-value function for label y is

$$p_{\text{all}}(X_{n+1}, y; \mathbf{W}^{(n+1)}) = \sum_k W_k^{(n+1)} \hat{p}_k(X_{n+1}, y).$$

At significance level α , we form the MoE prediction set by thresholding the corrected p-variable $\hat{m}^* p_{\text{all}}$, as established in Proposition 5.2:

$$\hat{C}_\alpha^{\text{MoE}}(X_{n+1}) = \{y \in \mathcal{Y} : \hat{m}^* p_{\text{all}}(X_{n+1}, y; \mathbf{W}^{(n+1)}) > \alpha\}.$$

Intuitively, the routing weights adapt the contribution of each expert’s p-value to the input so that experts with higher predictive relevance for X_{n+1} have greater influence. This data-adaptive weighting yields coverage guarantees that are locally more robust than split conformal with the full MoE predictor (Figure 2).

Baselines We refer to our proposed method of aggregating expert p-value functions with learned weights as *weighted aggregation*. As a baseline, we compare against split conformal with the full MoE predictor, and we refer to this simply as *split conformal* for brevity.

To evaluate local validity, we compare against conformal quantile regression (CQR) (Romano et al., 2019), a widely used locally adaptive method. A strength of our framework is that it complements, rather than competes with, other adaptive methods, allowing weighted aggregation to be layered on top of CQR. Accordingly, we also evaluate a hybrid method that combines weighted aggregation with CQR. We assess performance using marginal coverage, worst-slice (WS) coverage (Romano et al., 2020), and prediction set size.

We study two practical weighted aggregation variants:

- *WA targeted* $(0, \alpha']$: Coverage is guaranteed for all $\alpha \in (0, \alpha']$ using the m^\dagger correction (9).

- *WA precise* α' : Coverage is guaranteed for $\alpha = \alpha'$ only, using the m^\ddagger correction (10).

Experiment overview We present our main experimental results on the local validity of our method on real-world data in Figure 3. These results focus on regression tasks, where we can compare the standard absolute residual and CQR nonconformity scores; we evaluate our method on classification tasks in §B.2.1. To complement this analysis, which focuses on WS coverage, we also assess local validity from another perspective, by examining coverage disparities across demographic groups in the Communities and Crimes dataset (Redmond and Baveja, 2002) in §B.2.2.

In the Appendix, we analyze how various factors impact coverage and prediction set size in a series of ablative studies on synthetic data. We investigate: the effects of merging set size $|S_{\text{merge}}|$ (§B.1.1), shared feature information (§B.1.2), and the different weighted aggregation variants (§B.1.3).

The following experiments use an MoE where the routing network and experts are all linear models for simplicity. All results are averaged over 200 trials.

Main regression result In Figure 3, we compare split conformal to weighted aggregation using both absolute residual scores and CQR scores. From the marginal coverage plots (left column), we see that split conformal consistently achieves coverage closest to the nominal $1 - \alpha$ level, while the weighted aggregation methods tend to overcover. However, the WS coverage plots (middle column) expose a well-known limitation of split conformal: it systematically undercovers on the WS slab. While CQR scores improve split conformal WS coverage on some datasets (Airfoil and Communities), they are still insufficient to close the coverage gap.

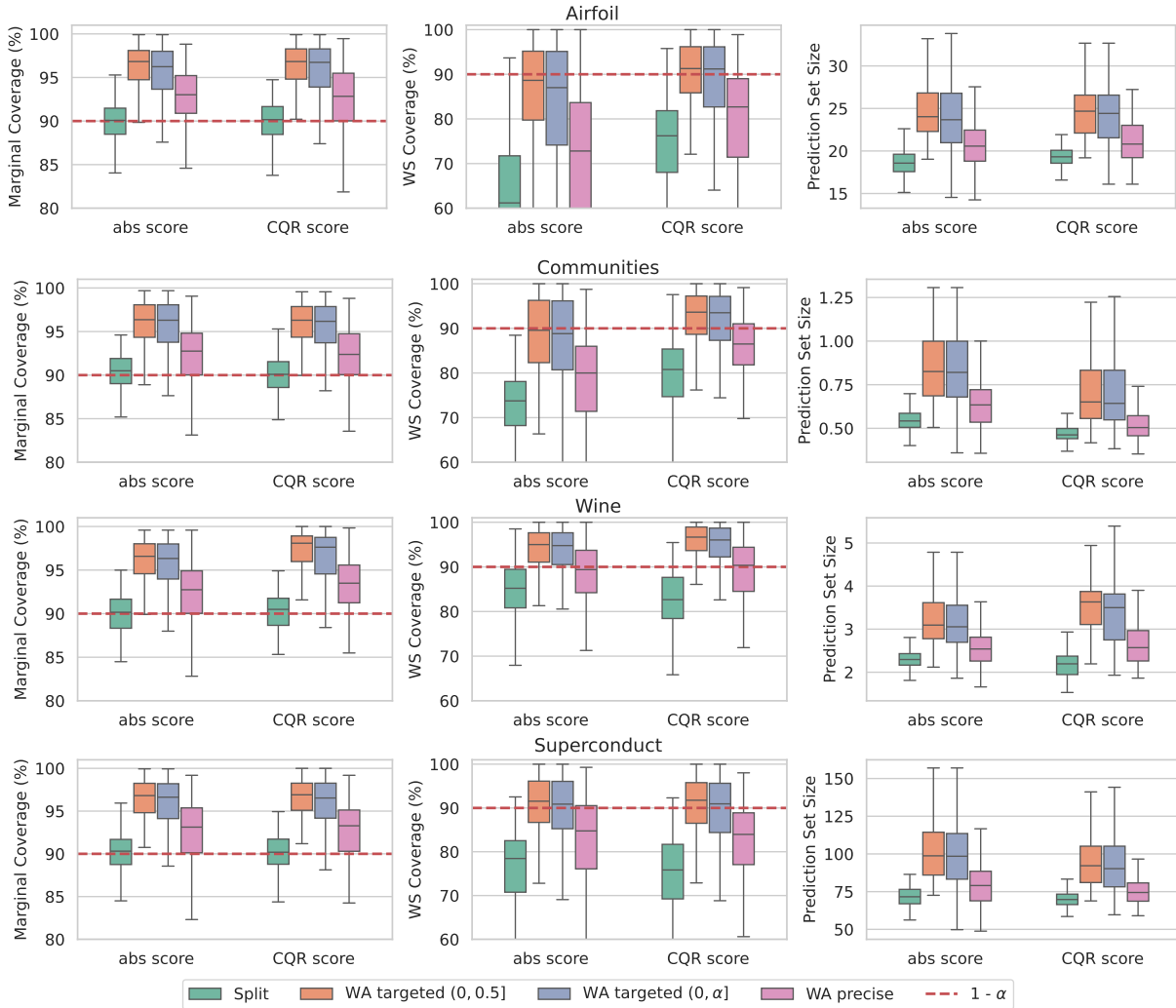


Figure 3: Local validity experiments comparing split conformal to weighted aggregation (WA) using absolute residual scores and CQR scores. Each row corresponds to a dataset, with plots for marginal coverage, WS coverage, and prediction set size from left to right. WA consistently improves WS coverage across datasets, which split conformal undercovers.

Weighted aggregation significantly improves WS coverage, coming close to or meeting the nominal level on all datasets. In addition, the gap between marginal and WS coverage is notably smaller for weighted aggregation than for split conformal, indicating that weighted aggregation provides more uniform coverage even over the challenging regions of the data. These results suggest that weighted aggregation is useful in applications where local validity is a priority.

The prediction set size plots (right column) illustrate the standard trade-off between coverage and efficiency, where higher-coverage methods like WA targeted tend to produce larger prediction sets than lower-coverage methods like split conformal. According to these results, we suggest WA targeted if maintaining coverage in challenging regions is the primary concern. If both efficiency and local validity are important, WA precise provides a reasonable middle ground—offering improved WS coverage over split conformal, while also

maintaining more compact prediction set size.

7 CONCLUSION

The asymmetric (weighted) aggregation of prediction sets is a flexible generalization of standard symmetric aggregation. Intuitively, coverage in this setting should vary with the distribution of weights. The results of Vovk and Wang (2020) allow us to formalize this intuition for conformal prediction, so that for data-independent weights based on expert priors, we obtain improved guarantees when the weights are sufficiently asymmetric. We extend this framework to data-dependent weights (e.g. weights learned from data), enabling adaptive coverage that reflects the observed input. Experiments on WS coverage and demographic subgroups confirm the practical benefits of this extension, showing that weighted aggregation yields more reliable coverage in challenging settings.

References

- Bonan Min, Hayley Ross, Elinor Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Computing Surveys*, 56(2):1–40, 2023.
- Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Foundations & trends in multimodal machine learning: Principles, challenges, and open questions. *ACM Computing Surveys*, 56(10):1–42, 2024a.
- Michelle Chua, Doyun Kim, Jongmun Choi, Nahyoung G Lee, Vikram Deshpande, Joseph Schwab, Michael H Lev, Ramon G Gonzalez, Michael S Gee, and Synho Do. Tackling prediction uncertainty in machine learning for healthcare. *Nature Biomedical Engineering*, 7(6):711–718, 2023.
- Thomas Grote and Philipp Berens. Uncertainty, evidence, and the integration of machine learning into medical practice. In *The Journal of Medicine and Philosophy: A Forum for Bioethics and Philosophy of Medicine*, volume 48, pages 84–97. Oxford University Press US, 2023.
- Sachin Kumar, T Gopi, N Harikeerthana, Munish Kumar Gupta, Vidit Gaur, Grzegorz M Krolczyk, and ChuanSong Wu. Machine learning techniques in additive manufacturing: a state of the art review on design, processes and production control. *Journal of Intelligent Manufacturing*, 34(1):21–55, 2023.
- Nathan P Lawrence, Seshu Kumar Damarla, Jong Woo Kim, Aditya Tulshyan, Faraz Amjad, Kai Wang, Benoit Chachuat, Jong Min Lee, Biao Huang, and R Bhushan Gopaluni. Machine learning for industrial sensing and control: A survey and practical perspective. *Control Engineering Practice*, 145:105841, 2024.
- Veronika Eyring, William D Collins, Pierre Gentine, Elizabeth A Barnes, Marcelo Barreiro, Tom Beucier, Marc Bocquet, Christopher S Bretherton, Hannah M Christensen, Katherine Dagon, et al. Pushing the frontiers in climate modelling and analysis with machine learning. *Nature Climate Change*, pages 1–13, 2024.
- Ching-Yao Lai, Pedram Hassanzadeh, Aditi Sheshadri, Maike Sonnewald, Raffaele Ferrari, and Venkatramani Balaji. Machine learning for climate physics and simulations. *arXiv preprint arXiv:2404.13227*, 2024.
- Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random world*, volume 29. Springer, 2005.
- Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence machines for regression. In *Machine learning: ECML 2002: 13th European conference on machine learning Helsinki, Finland, August 19–23, 2002 proceedings 13*, pages 345–356. Springer, 2002.
- Vladimir Vovk. Cross-conformal predictors. *Annals of Mathematics and Artificial Intelligence*, 74:9–28, 2015.
- Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. Predictive inference with the jackknife+. 2021.
- Matteo Gasparin and Aaditya Ramdas. Merging uncertainty sets via majority vote. *arXiv preprint arXiv:2401.09379*, 2024.
- Vladimir Vovk and Ruodu Wang. Combining p-values via averaging. *Biometrika*, 107(4):791–808, 2020.
- Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.
- Yachong Yang and Arun Kumar Kuchibhotla. Selection and aggregation of conformal prediction sets. *Journal of the American Statistical Association*, pages 1–13, 2021.
- Ruiting Liang, Wanrong Zhu, and Rina Foygel Barber. Conformal prediction after efficiency-oriented model selection. *arXiv preprint arXiv:2408.07066*, 2024b.
- Giovanni Cherubin. Majority vote ensembles of conformal predictors. *Machine Learning*, 108(3):475–488, 2019.
- Aldo Solari and Vera Djordjilović. Multi split conformal prediction. *Statistics & Probability Letters*, 184:109395, 2021.
- Vineeth N Balasubramanian, Shayok Chakraborty, and Sethuraman Panchanathan. Conformal predictions for information fusion: A comparative study of p-value combination methods. *Annals of Mathematics and Artificial Intelligence*, 74(1):45–65, 2015.
- Cyrus J DiCiccio, Thomas J DiCiccio, and Joseph P Romano. Exact tests via multiple data splitting. *Statistics & Probability Letters*, 166:108865, 2020.
- Bernhard Rüger. Das maximale signifikanzniveau des tests: “lehne H_0 ab, wenn k unter n gegebenen tests zur ablehnung führen”. *Metrika*, 25:171–178, 1978.
- Gerhard Hommel. Tests of the overall hypothesis for arbitrary dependence structures. *Biometrical Journal*, 25(5):423–430, 1983.
- Ludger Rüschendorf. Random variables with maximum sums. *Advances in Applied Probability*, 14(3):623–632, 1982.

David Stutz, Abhijit Guha Roy, Tatiana Matejovicova, Patricia Strachan, Ali Taylan Cemgil, and Arnaud Doucet. Conformal prediction under ambiguous ground truth. *arXiv preprint arXiv:2307.09302*, 2023.

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.

Vladimir Vovk. Conditional validity of inductive conformal predictors. In *Asian conference on machine learning*, pages 475–490. PMLR, 2012.

Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression. *Advances in neural information processing systems*, 32, 2019.

Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. *Advances in Neural Information Processing Systems*, 33: 3581–3591, 2020.

Michael Redmond and Alok Baveja. A data-driven software tool for enabling cooperative information sharing among police departments. *European Journal of Operational Research*, 141(3):660–678, 2002.

Isaac Gibbs, John J Cherian, and Emmanuel J Candès. Conformal prediction with conditional guarantees. *arXiv preprint arXiv:2305.12616*, 2023.

Alexandre Decan. portion: Python data structure and operations for intervals. URL <https://github.com/AlexandreDecan/portion>.

Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. The uci machine learning repository. <https://archive.ics.uci.edu>, 2010.

Checklist

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. **[Yes]**
We establish our theoretical contributions in §4 and §5, explicitly stating our setting and assumptions. We also discuss how to implement data-dependent weights in practice in Algorithm 1.
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. **[Yes]**
Please see §D for a discussion on the computational complexity of our method.
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. **[Yes]**
2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. **[Yes]**
All theoretical results include a full list of assumptions at the start of the result.
 - (b) Complete proofs of all theoretical results. **[Yes]**
Please see §A for complete proofs to all theoretical results, with any additional theoretical details that may be of interest.
 - (c) Clear explanations of any assumptions. **[Yes]**
3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). **[Yes]**
The code to reproduce all experiments and figures is in the supplemental material.
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). **[Yes]**
We tried to include all important training details in §6 in “Baselines” and “Experiment overview”. For further details, please see §E.
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). **[Yes]**
 - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). **[Not Applicable]**
All experiments can be run on a consumer-grade CPU.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
 - (a) Citations of the creator if your work uses existing assets. **[Yes]**
 - (b) The license information of the assets, if applicable. **[Yes]**
License information is included in the source code, in the supplementary.
 - (c) New assets either in the supplemental material or as a URL, if applicable. **[Yes]**
The code to reproduce all experiments and figures is in the supplemental material.
 - (d) Information about consent from data providers/curators. **[Not Applicable]**

- (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [**Not Applicable**]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
- (a) The full text of instructions given to participants and screenshots. [**Not Applicable**]
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [**Not Applicable**]
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [**Not Applicable**]

A PROOFS AND ADDITIONAL THEORETICAL DETAILS

A.1 Proof of Proposition 4.1

Following the notation of Vovk and Wang (2020), we define the merging function

$$M_{1,\mathbf{v}}(p_1, \dots, p_K) = (v_1 p_1 + \dots + v_K p_K) \quad \text{where} \quad \sum_{k=1}^K v_k = 1.$$

Recall that we use p-variables to represent our p-value functions applied to data; that is, $P_k = \hat{p}_k(X, Y)$. Then $\bar{P} = M_{1,\mathbf{v}}(P_1, \dots, P_K)$. By Proposition 9 of Vovk and Wang (2020), $A_{\mathbf{v}} M_{1,\mathbf{v}}$ is a precise merging function, where $A_{\mathbf{v}} = \min\{\frac{1}{v}, 2\}$ and $v = \max\{v_1, \dots, v_K\}$. Thus, $A_{\mathbf{v}} \bar{P}$ is a p-variable, and $\mathbb{P}\{A_{\mathbf{v}} \bar{P} \leq \alpha\} \leq \alpha$, or

$$\mathbb{P}\{\bar{P} \leq \alpha\} \leq A_{\mathbf{v}} \alpha = \min\left\{\frac{1}{v}, 2\right\} \alpha.$$

Under exchangeability, the prediction set

$$\widehat{C}_{\alpha}^{\text{avg}}(X_{n+1}) = \{y \in \mathcal{Y} : \bar{p}(X_{n+1}, y) > \alpha\}$$

constructed from \bar{p} has the coverage guarantee

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{\alpha}^{\text{avg}}(X_{n+1})\} \geq 1 - \min\left\{\frac{1}{v}, 2\right\} \alpha.$$

A.2 Proof of Proposition 5.1

We define m^* to be the smallest positive value such that $m^* P_{\text{all}}$ is a p-variable, i.e.

$$\mathbb{P}\{m^* P_{\text{all}} \leq \alpha\} \leq \alpha \quad \forall \alpha \in (0, 1),$$

or, equivalently,

$$\mathbb{P}\{m^* P_{\text{all}} > \alpha\} \geq 1 - \alpha \quad \forall \alpha \in (0, 1).$$

Recall that $P_{\text{all}} := p_{\text{all}}(X, Y, \mathbf{W})$. Then under exchangeability, the prediction set

$$\widehat{C}_{\alpha}^{\text{scaled}}(X_{n+1}) = \{y \in \mathcal{Y} : p_{\text{all}}(X_{n+1}, y; \mathbf{W}^{(n+1)}) > \alpha\}$$

has coverage

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{\alpha}^{\text{scaled}}(X_{n+1})\} \geq 1 - \alpha, \quad \forall \alpha \in (0, 1).$$

This establishes the result.

Note that other transformations of P_{all} can also yield valid p-variables while preserving the relative weighting. For example, if P_{all} is shifted rather than scaled, then the condition

$$\mathbb{P}\{m' + P_{\text{all}} \leq \alpha\} \leq \alpha$$

also leads to coverage at least $1 - \alpha$. We focus on the scale correction factor m^* because it aligns with the framework of Vovk and Wang (2020) and works well in practice.

A.3 Proof of Proposition 5.2

Let \mathcal{G} be the event that the worst-case distance between the true CDF $F_{P_{\text{all}}}$ and the empirical CDF $\widehat{F}_{P_{\text{all}}}$ is bounded by some maximum allowable deviation; that is,

$$\mathcal{G} = \left\{ \sup_{x \in \mathbb{R}} \left| \widehat{F}_{P_{\text{all}}}(x) - F_{P_{\text{all}}}(x) \right| \leq \varepsilon \right\} \quad \text{where} \quad \varepsilon = \sqrt{\frac{\log(2/\delta)}{2|S_{\text{merge}}|}}$$

for some user-chosen failure probability $\delta \in (0, 1)$. By the Dvoretzky–Kiefer–Wolfowitz (DKW) inequality,

$$\mathbb{P}(\mathcal{G}) \geq 1 - \delta.$$

Now, suppose we have a fixed merging set S_{merge} such that \mathcal{G} holds. By definition of a CDF,

$$\mathbb{P}\{\widehat{m}^* P_{\text{all}} \leq \alpha \mid S_{\text{merge}}\} = F_{P_{\text{all}}}\left(\frac{\alpha}{\widehat{m}^*}\right), \quad (12)$$

and under the same fixed S_{merge} ,

$$F_{P_{\text{all}}}(x) \leq \widehat{F}_{P_{\text{all}}}(x) + \varepsilon.$$

We take $x = \alpha/\widehat{m}^*$; this allows us to express the bound as

$$F_{P_{\text{all}}}\left(\frac{\alpha}{\widehat{m}^*}\right) \leq \widehat{F}_{P_{\text{all}}}\left(\frac{\alpha}{\widehat{m}^*}\right) + \varepsilon. \quad (13)$$

By construction of \widehat{m}^* (8), we know that $\widehat{F}_{P_{\text{all}}}(x)/x \leq \widehat{m}^*$ for all x . Substituting again $x = \alpha/\widehat{m}^*$, we have

$$\widehat{F}_{P_{\text{all}}}\left(\frac{\alpha}{\widehat{m}^*}\right) \leq \widehat{m}^* \frac{\alpha}{\widehat{m}^*} = \alpha. \quad (14)$$

Combining (13) and (14) gives

$$F_{P_{\text{all}}}\left(\frac{\alpha}{\widehat{m}^*}\right) \leq \alpha + \varepsilon.$$

Applying this to (12) gives

$$\mathbb{P}\{\widehat{m}^* P_{\text{all}} \leq \alpha \mid S_{\text{merge}}\} \leq \alpha + \varepsilon. \quad (15)$$

Under exchangeability, the prediction set

$$\widehat{C}_{\alpha}^{\text{scaled}}(X_{n+1}) = \left\{y \in \mathcal{Y} : \widehat{m}^* p_{\text{all}}(X_{n+1}, y; \mathbf{W}^{(n+1)}) > \alpha\right\}$$

has a conditional *miscoverage* guarantee of

$$\mathbb{P}\left\{Y_{n+1} \notin \widehat{C}_{\alpha}^{\text{scaled}}(X_{n+1}) \mid S_{\text{merge}}\right\} \leq \alpha + \varepsilon.$$

Let us denote the miscoverage indicator as

$$\text{MC} = \mathbb{1}\left\{Y_{n+1} \notin \widehat{C}_{\alpha}^{\text{scaled}}(X_{n+1})\right\}.$$

The indicator random variable allows us to easily switch between probability and expectation in order to marginalize over S_{merge} :

$$\begin{aligned} \mathbb{P}\{\text{MC} = 1\} &= \mathbb{E}[\text{MC}] \\ &= \mathbb{E}[\mathbb{E}[\text{MC} \mid S_{\text{merge}}]] \\ &= \mathbb{E}[\mathbb{P}\{\text{MC} = 1 \mid S_{\text{merge}}\}] \\ &\leq (\alpha + \varepsilon) \mathbb{P}(\mathcal{G}) + 1 \cdot \mathbb{P}(\mathcal{G}^c) \\ &\leq \alpha + \varepsilon + \delta. \end{aligned}$$

This gives the coverage guarantee

$$\mathbb{P}\left\{Y_{n+1} \in \widehat{C}_{\alpha}^{\text{scaled}}(X_{n+1})\right\} \geq 1 - (\alpha + \varepsilon + \delta).$$

Note that the same argument applies directly if we replace the empirical CDF $\widehat{F}_{P_{\text{all}}}$ by the conservative version $\widehat{F}_{P_{\text{all}}}^{\text{cons}}$ (17) because the two differ by at most $(|S_{\text{merge}}| + 1)^{-1}$, a deterministic offset that can simply be absorbed into the DKW tolerance by replacing ε with $\varepsilon + (|S_{\text{merge}}| + 1)^{-1}$.

A.4 Prediction sets from unscaled p_{all}

For completeness, we also state the guarantees for the *unscaled* prediction sets, where thresholding is applied directly to p_{all} rather than to its corrected version m^*p_{all} . In this formulation, the factor m^* (or \hat{m}^*) appears explicitly in the coverage bound, paralleling Proposition 4.1.

The following corollaries restate Propositions 5.1 and 5.2 in the unscaled form.

Corollary A.1 (Infinite-sample guarantee, unscaled). *Under the assumptions of Proposition 5.1, the prediction set*

$$\hat{C}_{\alpha}^{\text{unscaled}}(X_{n+1}) = \{y : p_{\text{all}}(X_{n+1}, y) > \alpha\},$$

satisfies

$$\mathbb{P} \left\{ Y_{n+1} \in \hat{C}_{\alpha}^{\text{unscaled}}(X_{n+1}) \right\} \geq 1 - m^* \alpha \quad \forall \alpha \in (0, 1).$$

Proof. As before, we start with the fact that m^* is the smallest positive value such that m^*P_{all} is a p-variable, so

$$\mathbb{P} \{ m^* P_{\text{all}} \leq \alpha \} \leq \alpha.$$

We can rearrange this condition to be

$$\mathbb{P} \left\{ P_{\text{all}} \leq \frac{\alpha}{m^*} \right\} \leq \alpha;$$

if we define $\alpha' = \alpha/m^*$, this becomes

$$\mathbb{P} \{ P_{\text{all}} \leq \alpha' \} \leq m^* \alpha',$$

or equivalently,

$$\mathbb{P} \{ P_{\text{all}} > \alpha' \} \geq 1 - m^* \alpha'.$$

Under exchangeability, the prediction set

$$\hat{C}_{\alpha}^{\text{unscaled}}(X_{n+1}) = \left\{ y \in \mathcal{Y} : p_{\text{all}}(X_{n+1}, y; \mathbf{W}^{(n+1)}) > \alpha \right\}$$

satisfies

$$\mathbb{P} \left\{ Y_{n+1} \in \hat{C}_{\alpha}^{\text{unscaled}}(X_{n+1}) \right\} \geq 1 - m^* \alpha.$$

If P_{all} is shifted rather than scaled, then the condition

$$\mathbb{P} \{ m' + P_{\text{all}} \leq \alpha \} \leq \alpha$$

leads to a coverage guarantee of $1 - (m' + \alpha)$. \square

Corollary A.2 (Finite-sample guarantee, unscaled). *Under the assumptions of Proposition 5.2, the prediction set*

$$\hat{C}_{\alpha}^{\text{unscaled}}(X_{n+1}) = \{y : p_{\text{all}}(X_{n+1}, y) > \alpha\}$$

satisfies

$$\mathbb{P} \{ Y_{n+1} \in \hat{C}_{\alpha}^{\text{unscaled}}(X_{n+1}) \} \geq 1 - (\alpha \mathbb{E}[\hat{m}^*] + \varepsilon + \delta) \quad \forall \alpha \in (0, 1).$$

Proof. The proof proceeds identically to the proof of Proposition 5.2 (in §A.3) up to inequality (15), which establishes

$$\mathbb{P} \{ \hat{m}^* P_{\text{all}} \leq \alpha \mid S_{\text{merge}} \} \leq \alpha + \varepsilon,$$

or

$$\mathbb{P} \{ P_{\text{all}} \leq \alpha' \mid S_{\text{merge}} \} \leq \hat{m}^* \alpha' + \varepsilon$$

by a change of variable.

The remainder of the proof follows the same conditioning and marginalization argument, with the substitution of $\hat{m}^* \alpha$ for α carried through. To be explicit, the prediction set

$$\hat{C}_{\alpha}^{\text{unscaled}}(X_{n+1}) = \left\{ y \in \mathcal{Y} : p_{\text{all}}(X_{n+1}, y; \mathbf{W}^{(n+1)}) > \alpha \right\}$$

has a conditional miscoverage guarantee of

$$\mathbb{P} \left\{ Y_{n+1} \notin \widehat{C}_\alpha^{\text{unscaled}}(X_{n+1}) \mid S_{\text{merge}} \right\} \leq \widehat{m}^* \alpha + \varepsilon.$$

Then we can use miscoverage indicator $\text{MC} = \mathbf{1} \left\{ Y_{n+1} \notin \widehat{C}_\alpha^{\text{unscaled}}(X_{n+1}) \right\}$, to help us marginalize over S_{merge} :

$$\begin{aligned} \mathbb{P} \{ \text{MC} = 1 \} &= \mathbb{E} [\mathbb{P} \{ \text{MC} = 1 \mid S_{\text{merge}} \}] \\ &\leq (\alpha \mathbb{E}[\widehat{m}^*] + \varepsilon) \mathbb{P}(\mathcal{G}) + 1 \cdot \mathbb{P}(\mathcal{G}^c) \\ &\leq \alpha \mathbb{E}[\widehat{m}^*] + \varepsilon + \delta. \end{aligned}$$

This gives the final bound

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha^{\text{unscaled}}(X_{n+1}) \right\} \geq 1 - (\alpha \mathbb{E}[\widehat{m}^*] + \varepsilon + \delta).$$

□

A.5 More detailed derivation of m^*

For some distribution function $F_{P_{\text{all}}}$, we define

$$m^* = \inf \{ m > 0 : F_{P_{\text{all}}}(\alpha/m) \leq \alpha \text{ for all } \alpha \in (0, 1) \}, \quad c = \sup_{\delta > 0} \frac{F_{P_{\text{all}}}(\delta)}{\delta}.$$

We aim to prove that $m^* = c$ to establish the equivalence in (7). To aid our proof, we define the feasible set

$$S = \{ m > 0 : F_{P_{\text{all}}}(\alpha/m) \leq \alpha \text{ for all } \alpha \in (0, 1) \},$$

where $m^* = \inf(S)$.

We begin by showing that $m^* \leq c$. To this end, consider any $m > c$. By definition of c ,

$$F_{P_{\text{all}}}(\delta) \leq c\delta < m\delta \quad \forall \delta > 0.$$

Pick $\delta = \alpha/m$ for any $\alpha \in (0, 1)$ so that $\delta \in (0, 1/m)$. Substituting yields

$$F_{P_{\text{all}}}(\alpha/m) < m(\alpha/m) = \alpha.$$

This shows that every $m > c$ satisfies the feasibility condition, so every $m > c$ is in S ; it follows that $S \supset (c, \infty)$. Then, as the infimum of S , $m^* \leq c$.

To show that $m^* \geq c$, consider any feasible $m > 0$, i.e. assume

$$F_{P_{\text{all}}}(\alpha/m) \leq \alpha \quad \forall \alpha \in (0, 1).$$

For any $\delta > 0$, there are two cases.

1. If $\delta \geq 1/m$, since $F_{P_{\text{all}}}(\delta) \leq 1$,

$$F_{P_{\text{all}}}(\delta)/\delta \leq 1/\delta \leq m.$$

2. If $\delta \in (0, 1/m)$, let us select $\alpha = m\delta \in (0, 1)$. By feasibility of m ,

$$F_{P_{\text{all}}}(\delta) = F_{P_{\text{all}}}(\alpha/m) \leq \alpha = m\delta,$$

or

$$F_{P_{\text{all}}}(\delta)/\delta \leq m.$$

The two cases above establish $m \geq c$ for every $m \in S$, or that $S \subset (c, \infty)$. Thus, c is a lower bound of S and $m^* \geq c$.

A.6 Computing \hat{m}^* from an empirical CDF

For random variable $P_{\text{all}} = p_{\text{all}}(X, Y; \mathbf{W})$, a merging correction factor λ is a positive scalar that ensures that a λP_{all} is a valid p-value. That is,

$$\lambda \in \{m > 0 : \mathbb{P}\{mP_{\text{all}} \leq \alpha\} \leq \alpha \forall \alpha \in (0, 1)\} = \{m > 0 : F_{P_{\text{all}}}(\alpha/m) \leq \alpha \forall \alpha \in (0, 1)\}.$$

Scaling by a merging correction factor gives us the guarantee

$$\mathbb{P}\{P_{\text{all}} > \alpha\} = 1 - \lambda\alpha.$$

To achieve the tightest guarantee (with application to all $\alpha \in (0, 1)$), we define the minimal merging correction factor m^* to be

$$m^* = \inf\{m > 0 : F_{P_{\text{all}}}(\alpha/m) \leq \alpha \forall \alpha \in (0, 1)\} = \sup_{\delta > 0} \frac{F_{P_{\text{all}}}(\delta)}{\delta}. \quad (16)$$

(Equation (16) is a restatement of (7) from the main body; we replicate it here for easy reference.)

In practice, we compute the minimal merging factor \hat{m}^* by first constructing the empirical CDF of P_{all} over the merging set

$$\hat{F}_{P_{\text{all}}}(\alpha) = \frac{\sum_{i \in S_{\text{merge}}} \mathbb{1}\{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) \leq \alpha\}}{|S_{\text{merge}}|},$$

or its conservative version $\hat{F}_{P_{\text{all}}}^{\text{cons}}$ (17). For the sake of conciseness, let

$$\hat{F}_i := \hat{F}_{P_{\text{all}}}\left(p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)})\right).$$

Since $\hat{F}_{P_{\text{all}}}$ is a right-continuous step function that only jumps at the observed values $\{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) : i \in S_{\text{merge}}\}$, the supremum in (16) is attained at one of these points. In fact, for any δ not equal to one of these values, shifting δ slightly to the right (toward the next jump) increases the denominator without changing the numerator, decreasing the ratio. For this reason, it suffices to take the maximum over the observed p-values, yielding the empirical merging factor

$$\hat{m}^* = \max_{\substack{i \in S_{\text{merge}} \\ p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) > 0}} \frac{\hat{F}_i}{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)})}.$$

In some cases, we are only interested in coverage above a certain significance level α as in (9). Then we only need to find the first point on the empirical CDF that surpasses α , or

$$\bar{\alpha} = \min_{\substack{i \in S_{\text{merge}} \\ \hat{F}_i \geq \alpha}} \hat{F}_i,$$

and then compute the scaling factor to be

$$\hat{m}^\dagger = \max_{\substack{i \in S_{\text{merge}} \\ p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) > 0 \\ \hat{F}_i \leq \bar{\alpha}}} \frac{\hat{F}_i}{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)})}.$$

B ADDITIONAL EXPERIMENTS

B.1 Synthetic data

We generate a simple homoskedastic dataset to simulate a regression task, where each input is a 16-dimensional Gaussian with random parameters, and the label is the sum of the different dimensions with additive noise (see §E.2 for more detail). We use this dataset to investigate how various factors impact the coverage and interval width of our prediction sets.

B.1.1 Coverage and size of prediction sets improve with larger $|S_{\text{merge}}|$

The merging set S_{merge} allows us to construct an empirical CDF for P_{all} . To improve stability in finite samples, we apply a conservative correction to the typical formula for the empirical CDF, and use this conservative CDF (17) to compute the correction factor \hat{m}^* (8). Our conservative CDF tends to be overly conservative when $|S_{\text{merge}}|$ is small, leading to overcoverage. However, as $|S_{\text{merge}}|$ grows and the empirical CDF approaches the true CDF, coverage becomes closer to nominal and the prediction sets become less conservative.

Figure 4 illustrates how the size of the merging set $|S_{\text{merge}}|$ impacts mean coverage. On the left plot, we use the m^\dagger correction (9) for *WA targeted* (that is, $(0, \alpha)$ aggregation), and we compare different ways of assigning features to experts. On the right plot, we use a non-overlapping feature assignment, and we compare the different variants of weighted aggregation. Our results show that larger $|S_{\text{merge}}|$ leads to tighter coverage, with 160 samples being sufficient to achieve $< 3\%$ overcoverage for most feature assignment methods. Interestingly, configurations where experts have no overlapping features tend to overcover the most. We explore the effect of feature assignment in the next section.

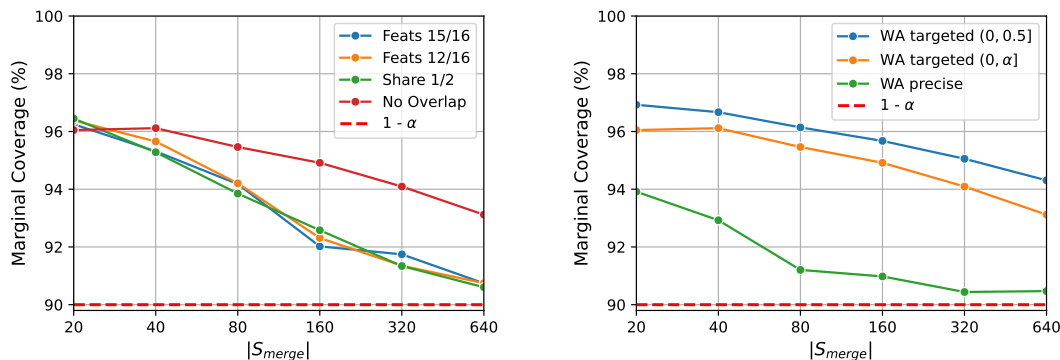


Figure 4: Mean coverage compared to the size of the merging set S_{merge} , with different feature assignments (left) and different weighted aggregation variants (right). Overall, we find that coverage improves as the merging set gets larger for all methods. We note that the merging set does not have to be prohibitively large to produce decent results: for example, most feature assignment methods overcover by only 2% with fewer than 200 samples.

B.1.2 Feature information overlap leads to higher coverage and more efficient prediction sets

To better understand how the allocation of features to experts affects the behavior of MoE weighted aggregation, we define four feature assignment methods (for our MoE of four experts):

- *Features 15/16*: each expert predicts from 15 of the 16 available features.
- *Features 12/16*: each expert predicts from 12 of the 16 available features.
- *Share 1/2*: all experts share 8 of the 16 features and partition the remaining 8 (2 features each).
- *No Overlap*: the experts partition the 16 features (4 features each).

Figure 5 shows how different feature assignment methods affect coverage and prediction set size. Broadly, we observe that greater feature overlap leads to higher coverage (exceeding the nominal level) and more efficient

prediction sets. In the MoE setting, this may be because feature sharing leads to more consistent estimations across experts, improving the reliability of the aggregated p-values and, in turn, reducing the size of the prediction sets. More generally, this may imply that the information redundancy introduced by feature overlap allows for better sample efficiency. This parallels findings in aggregation methods like cross-conformal and jackknife+, which also tend to produce smaller prediction sets than split conformal by reusing data. These results suggest that feature sharing is an important design consideration when aggregating prediction sets from multiple models.

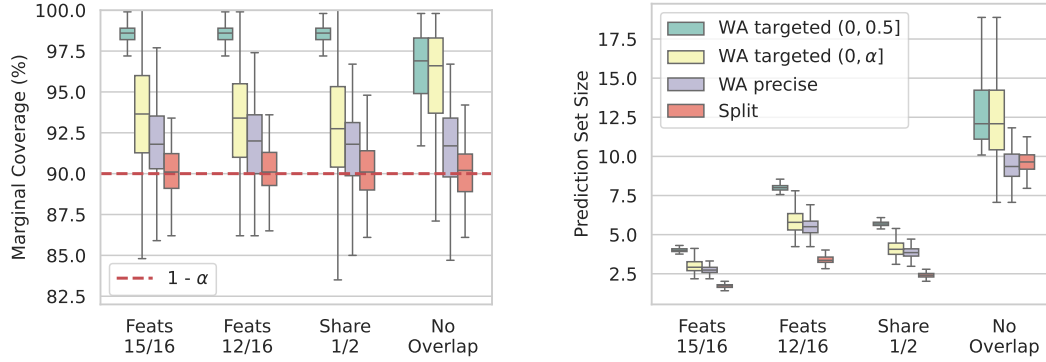


Figure 5: Coverage (left) and prediction set size (right) for different feature assignment methods and weighted aggregation methods. Our results indicate that sharing fewer features leads to tighter coverage, but sharing more features leads to more efficient (smaller) prediction sets. We also see that WA precise tends to have coverage that is closest to nominal and the most efficient prediction sets, albeit with a much looser guarantee.

B.1.3 More general coverage guarantees result in more conservative prediction sets

We now compare the different variants of weighted aggregation described in §6, focusing on how the generality of the guarantee for each variant affects its empirical coverage.

Figure 5 compares the effects of different feature assignment methods and for each weighted aggregation variant (with $\alpha' = 0.1$), and Figure 6 shows how coverage varies across α for each variant.

Unsurprisingly, WA precise—which provides the narrowest guarantee, targeting a single α —achieves coverage closest to the nominal level. In contrast, WA targeted offers more general guarantees over a range of α values, but tends to overcover, reflecting the conservativeness built into the method to accommodate worst-case behavior. We observe the same pattern on UCI data in Figure 3. These results illustrate the trade-off between the conservativeness of a method and the generality of its guarantee, and suggest that more targeted guarantees may be preferable when tighter coverage is important.

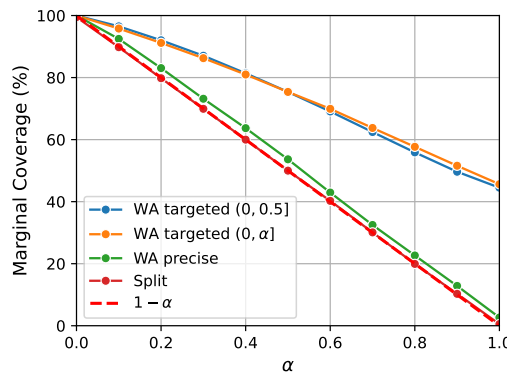


Figure 6: Mean coverage across significance levels α for the different weighted aggregation variants. WA precise achieves close to nominal coverage, while WA targeted tends to overcover (but offers more general guarantees).

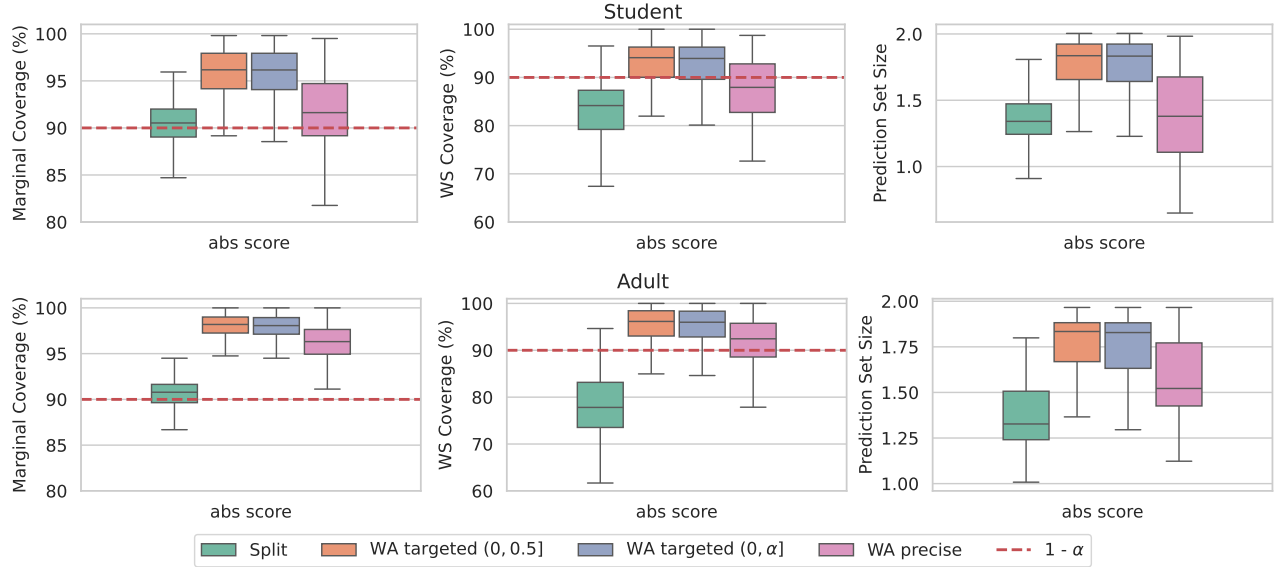


Figure 7: Local validity experiments comparing split conformal (green) to weighted aggregation (orange, purple, and pink) on the classification task. Like in Figure 3, each row corresponds to a dataset, with plots for marginal coverage, WS coverage, and prediction set size from left to right. We find that for classification as well as regression, WA improves WS coverage over split conformal, which undercovers in the WS region.

B.2 Real data

B.2.1 Weighted aggregation improves local validity for classification

The experiments in Figure 7 mirror those in Figure 3, but for classification instead of regression. Like with the other set of experiments, we compare split conformal to weighted aggregation, and we see that split conformal achieves marginal coverage close to nominal but consistently undercovers on the WS slab. In contrast, weighted aggregation maintains much better WS coverage, suggesting it offers better local validity in the classification setting as well.

As before, we note that there is trade-off between coverage and efficiency. WA targeted is the best choice in terms of WS coverage, but WA precise balances the improved WS coverage of weighted aggregation methods with an efficiency that is closer to the split conformal. However, if marginal coverage and prediction set efficiency are the only priorities, then split conformal should be preferred.

B.2.2 Local validity over demographic groups

In this section, we evaluate our methods on the Communities and Crimes dataset (Redmond and Baveja, 2002), where the task is to predict the per capita violent crime rate of a community based on its demographic features, and our primary interest is in understanding how coverage differs across communities with varying racial compositions (Gibbs et al., 2023).

Figures 8 and 9 compare split conformal with weighted aggregation variants in terms of coverage and prediction set size across demographic groups. Unlike Figures 3 and 7, where local validity is assessed via WS coverage, these experiments evaluate local validity in terms of consistency across demographic groups. Thus, we display group-specific performance for each method, which we group further by the two types of nonconformity score, absolute residual and CQR score.

The demographic groups in Figures 8 and 9 represent communities where a particular racial demographic is in the top p -percentile of representation (Gibbs et al., 2023). Figure 8 shows results for $p = 50$, and Figure 9 for $p = 70$.

Marginal coverage and coverage across groups The “All” category in the coverage plots represents coverage across all demographic groups, or marginal coverage. As in prior experiments, we see that split conformal achieves marginal coverage closest to nominal, with WA precise close behind, while WA targeted variants tend

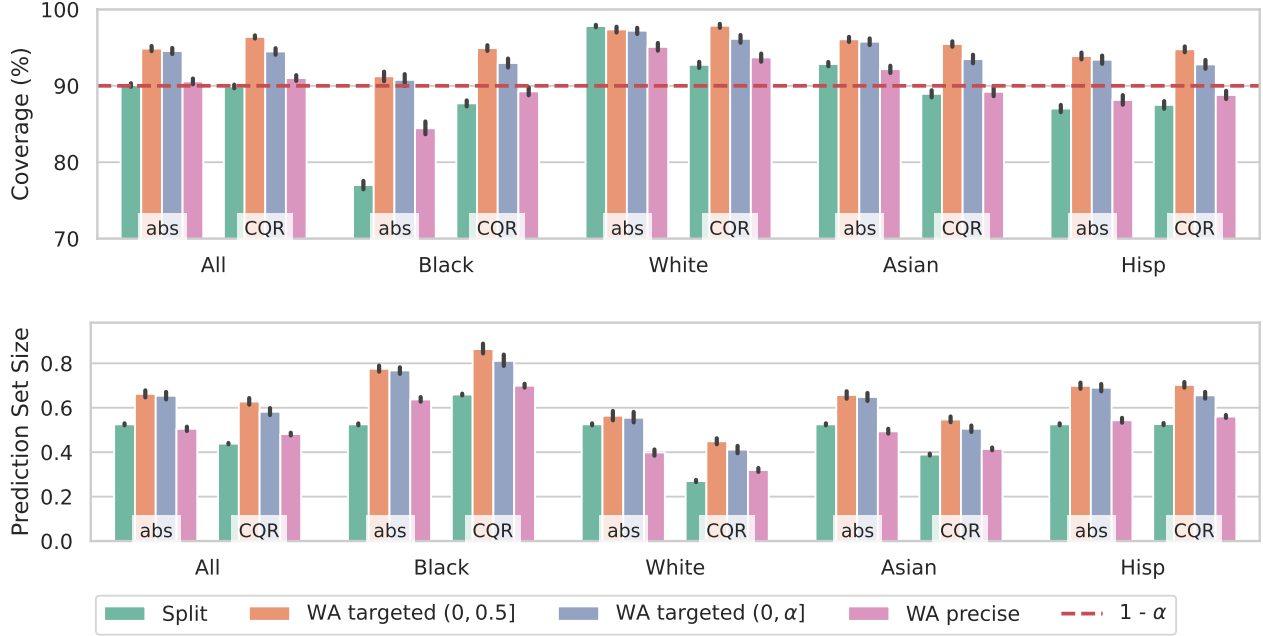


Figure 8: Coverage (top) and prediction set size (bottom) for split conformal (green) and weighted aggregation variants (orange, purple, and pink) across subgroups with top 50th percentile racial representation. Split conformal has precise marginal coverage (“All”), but WA variants have more consistent coverage across subgroups, with WA targeted meeting coverage for all subgroups. Error bars represent 95% confidence intervals.

to overcover. However, although split conformal enjoys precise marginal coverage, it also exhibits substantial disparities in performance across demographic groups, significantly undercovering for Black and Hispanic groups while overcovering for White. In contrast, WA variants display far less demographic variation, with WA targeted achieving coverage for all demographic groups.

Comparing nonconformity scores For all demographic groups except the Asian group, using CQR scores instead of absolute residuals improves coverage for split conformal, with the improvement being most pronounced for the Black group. However, while CQR reduces undercoverage, it is never sufficient to fully close the coverage gap for an undercovered group. Rather, its primary benefit appears to be in reducing the variability in coverage across groups. Across all methods, the variance in CQR-based coverage is lower than that of absolute residuals, suggesting that CQR contributes to more stable group-wise coverage, even if it does not fully mitigate disparities.

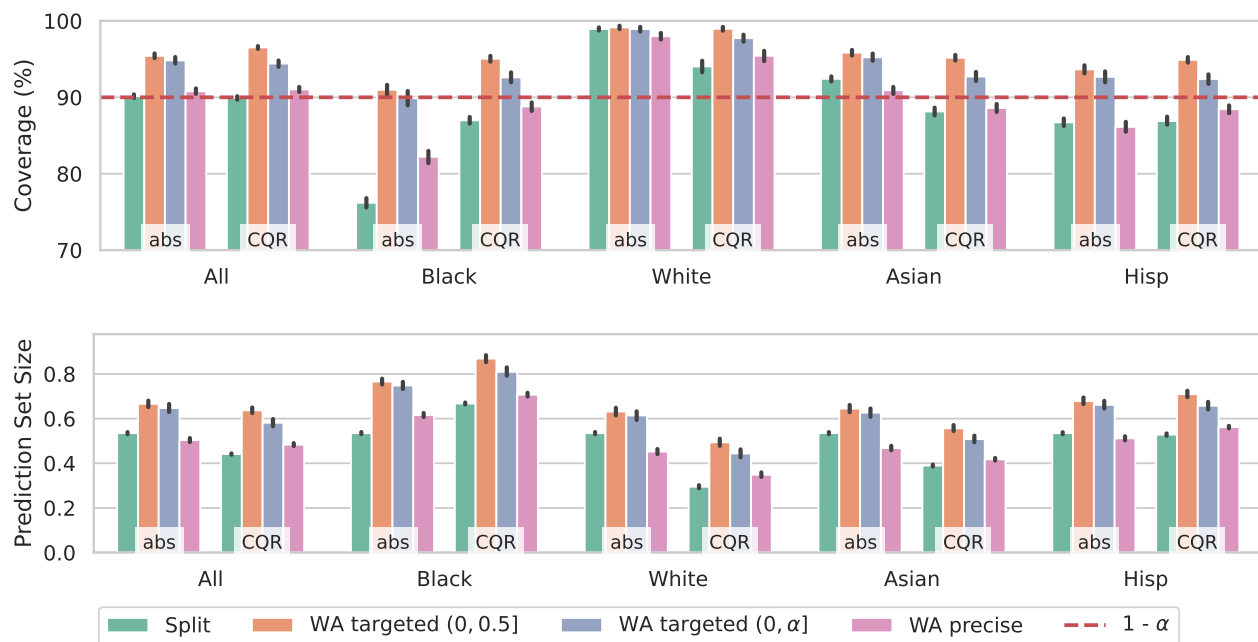


Figure 9: Coverage (top) and prediction set size (bottom) for split conformal and weighted aggregation variants across subgroups with top 70th percentile racial representation. As in Figure 8, WA targeted meets coverage for all subgroups. Error bars represent 95% confidence intervals.

C OTHER METHODS FOR WEIGHTED PREDICTION SET AGGREGATION

C.1 Weighted majority vote

The idea to combine conformal prediction sets by weighted majority vote was introduced by Gasparin and Ramdas (2024) as an extension of the majority vote method first proposed by Cherubin (2019). At first consideration, weighted majority vote appears to differ from our method: weighted majority vote performs weighted aggregation of *prediction sets*, while our method performs weighted aggregation of the *p-values associated with prediction sets*. Despite this distinction, Appendix B of Gasparin and Ramdas (2024) observes that these two methods are, in fact, dual to each other under data-independent weights, the setting considered in their work. Nevertheless, our p-value formulation enables two key extensions that go beyond weighted majority vote.

- Our formulation allows us to apply the result of Vovk and Wang (2020) to strengthen coverage guarantees for *data-independent* weights when they are sufficiently asymmetric.
- Our formulation provides a principled extension to *data-dependent* weights by transforming the weighted average of p-variables to also be a valid p-variable. Not only does this allow us to use weights learned from data, but it also yields a form of local validity, a property not available to existing set aggregation methods.

Because our method is a dual formulation to the weighted majority vote method of Gasparin and Ramdas (2024), we do not include it as a separate baseline to avoid redundancy.

C.2 Extending the p-variable transformation of Stutz et al. (2023) to weighted aggregation

To the best of our knowledge, Gasparin and Ramdas (2024) present the only existing method to address weighted prediction set aggregation, and the method of Stutz et al. (2023) is designed for the different problem of uncertainty in the ground truth labels. To address their problem, Stutz et al. (2023) propose sampling m labels for each calibration point and using the labels to compute m p-values, then taking the unweighted average of these p-values and applying a transformation to obtain a valid p-variable.

Although Stutz et al. (2023) only consider the unweighted average, their transformation is general enough to apply to a weighted average of p-values as well, and can therefore be adapted to our setting.

How the transformations affect the weights Both our method and the method of Stutz et al. (2023) aim to transform a random variable to a p-variable to maintain coverage guarantees. The difference between the two methods lies in the nature of the transformation. Our method applies a *linear* transformation that preserves the proportions of the weights; this can be important when the weights reflect meaningful quantities, like the weights learned by the routing network of an MoE model. In contrast, Stutz et al. (2023) apply a *nonlinear, rank-based* transformation by computing the empirical CDF of the random variable and returning its value at the observed point. That is, given a random variable X , their method estimates its CDF F and uses $F(X)$ as the resulting p-value. While this guarantees validity and preserves ordering, it does not preserve the relative scale between values and therefore discards some of the semantics of the original weights.

Comparing both transformations with our MoE setup In our method for weighted prediction set aggregation with data-dependent weights, we transform the weighted average of the individual prediction set p-values to a valid p-variable using a linear scaling. This transformation allows us to maintain a coverage guarantee for the prediction set defined by the weighted average of the p-values. In the MoE setting, this prediction set corresponds to aggregating the p-values from each expert according to the weights learned by the routing network.

To adapt the method of Stutz et al. (2023) to this context, we substitute their transformation in place of our linear scaling. We now restate their original setting and transformation in more detail to clarify how their method can be extended to our setting.

The transformation proposed by Stutz et al. (2023) was originally developed to address the problem of uncertainty in the ground truth labels. In their setting, each calibration point consists only of an input X_i ; they sample m labels to get m p-values for each calibration point, and then take the unweighted average of these p-values. We denote this unweighted average as P_{avg}^i , with distribution function F . Their method transforms P_{avg}^i into a valid

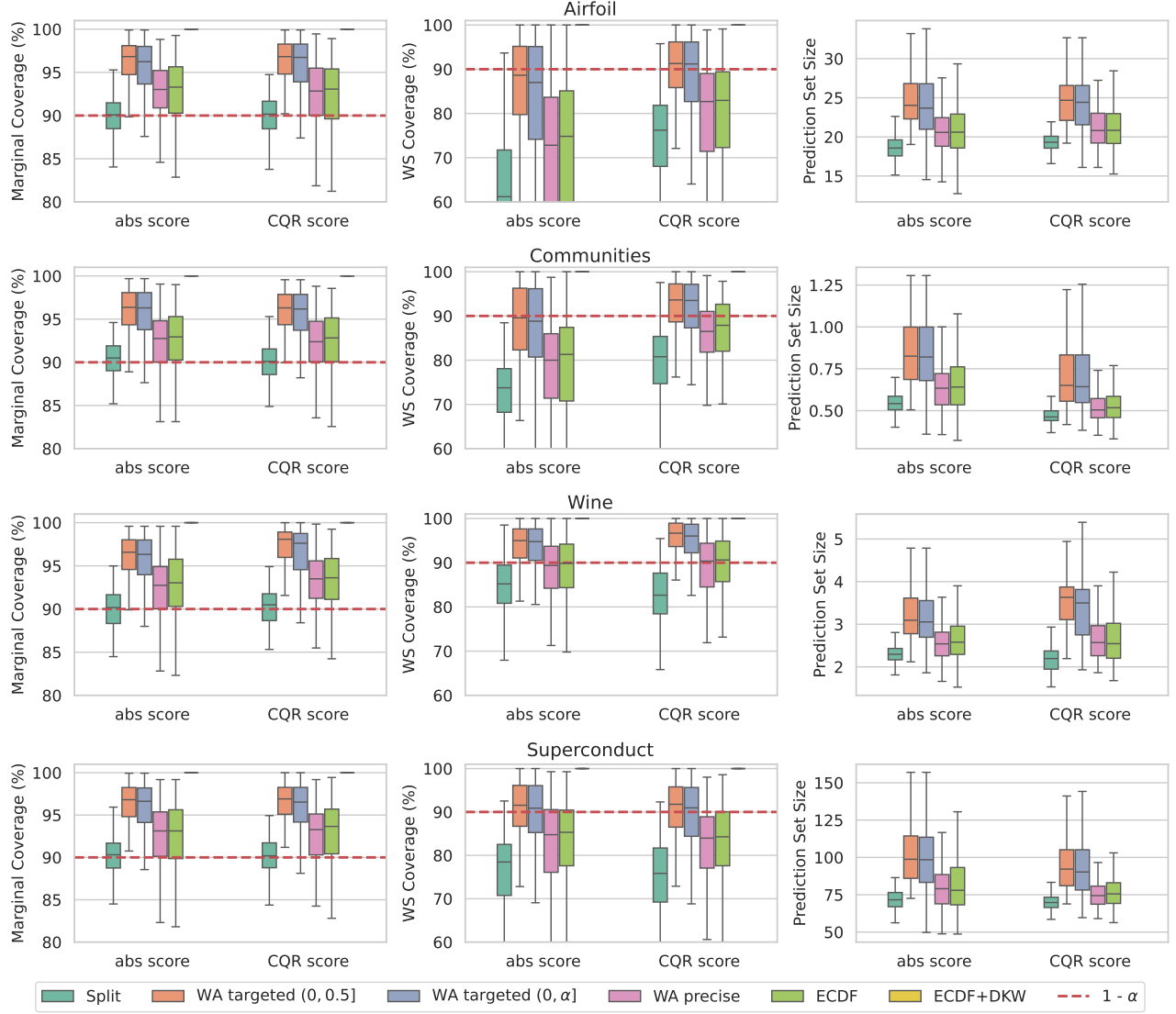


Figure 10: Regression experiments of our linear transformation method with ECDF and ECDF-DKW. Each row corresponds to a dataset, with plots for marginal coverage, WS coverage, and prediction set size from left to right. ECDF performs similarly (sometimes slightly more conservatively) to WA precise in terms of coverage and efficiency. However, ECDF-DKW is so conservative that it covers the entire label space.

p-value via $F(P_{avg}^i)$, and the prediction set is then the set of all labels such that the corrected p-value $F(P_{avg}^i)$ corresponding to each sample is greater than some threshold α .

Stutz et al. (2023) note that the coverage guarantee using the true CDF holds only with an asymptotic number of samples, as the empirical CDF \hat{F} approaches the true CDF F . To establish a finite-sample guarantee, they introduce a DKW-derived correction ϵ and define their prediction sets based on $\hat{F} + \epsilon$. We find that, although their proposed prediction set yields an elegant $(1 - \alpha)(1 - \delta)$ finite-sample guarantee, the ϵ correction is extremely conservative in practice.

Let us refer to the finite-sample empirical CDF method as *ECDF-DKW*, and the variant without the DKW correction as simply *ECDF*. We now present additional experiments where we recreate the main findings of our paper with the ECDF and ECDF-DKW methods.

Figure 10 recreates the regression experiments of Figure 3, with the addition of ECDF and ECDF-DKW. We note that ECDF performs very similarly to WA precise in terms of coverage (left), WS coverage (middle), and prediction set size (right), with ECDF being slightly more conservative in most cases. Like WA precise, ECDF is more conservative than split conformal and less conservative than WA targeted, although it still often falls short

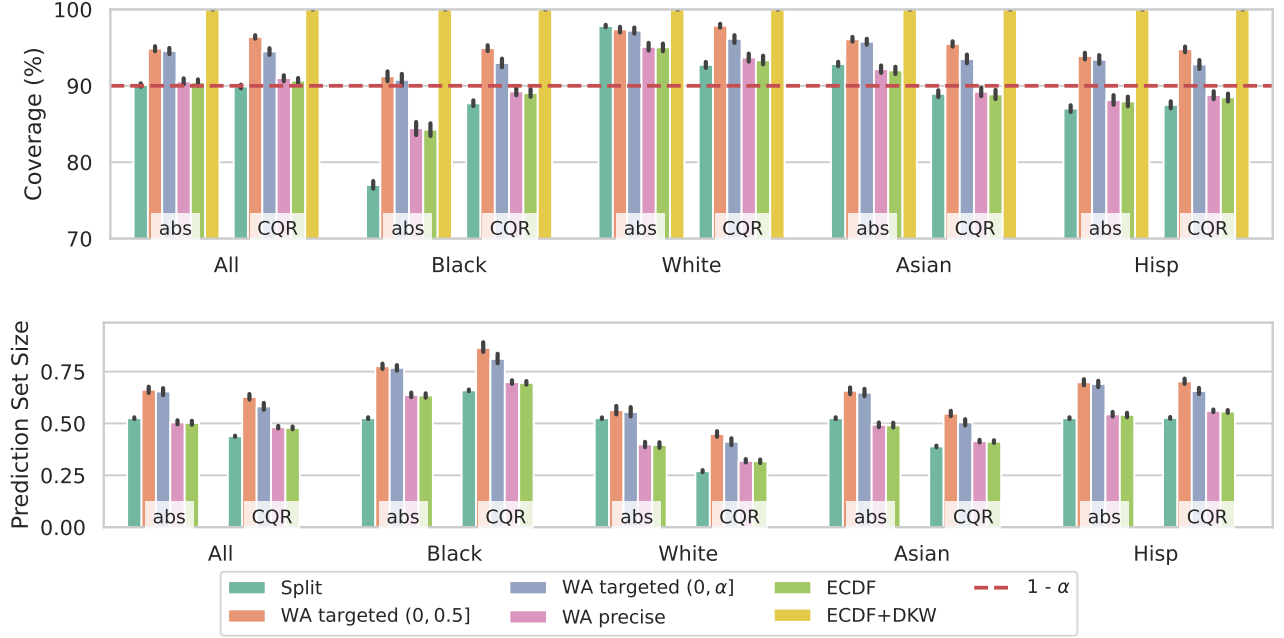


Figure 11: Coverage (top) and prediction set size (bottom) for subgroups with top 50th percentile racial representation in the UCI Communities and Crimes dataset. ECDF performs similarly to WA precise; ECDF-DKW is so conservative that it covers the entire label space. Error bars represent 95% confidence intervals.

of the $1 - \alpha$ guarantee on the WS slab. On the other hand, ECDF-DKW is so conservative that its prediction sets cover the entire label space. We represent this with 100% coverage on both the coverage plots and WS coverage plots, and we omit ECDF-DKW from the prediction set size plots.

Figure 11 recreates the Communities and Crimes experiment of Figure 8 with ECDF and ECDF-DKW. Again, we see that ECDF performs very similarly to WA precise on our demographic-conditioned view of Communities and Crimes, with similar coverage (top) and prediction set size (bottom) across all demographics—with ECDF having slightly lower coverage on most demographics, including demographics where WA precise undercovers. Like before, we also observe that ECDF-DKW is so conservative that it has 100% coverage and unbounded prediction sets.

Why is ECDF-DKW so conservative? For ECDF-DKW, the finite-sample variant of ECDF, Stutz et al. (2023) use DKW to add a finite-sample correction ϵ to the empirical CDF $\hat{F}(P_{\text{all}})$, then compare this sum to the significance level α . The prediction set with finite-sample guarantees is therefore the set of all labels such that $\hat{F}(P_{\text{avg}}) + \epsilon > \alpha$. However, if ϵ is already greater than α , then this condition is always satisfied and the prediction set includes *all* labels, becoming unbounded.

The finite-sample correction ϵ is a function of the number of samples used to compute the empirical CDF. Figure 12 shows that when the merging set size is less than 1000, ϵ is typically large enough to exceed common values of α , making unbounded sets very likely.

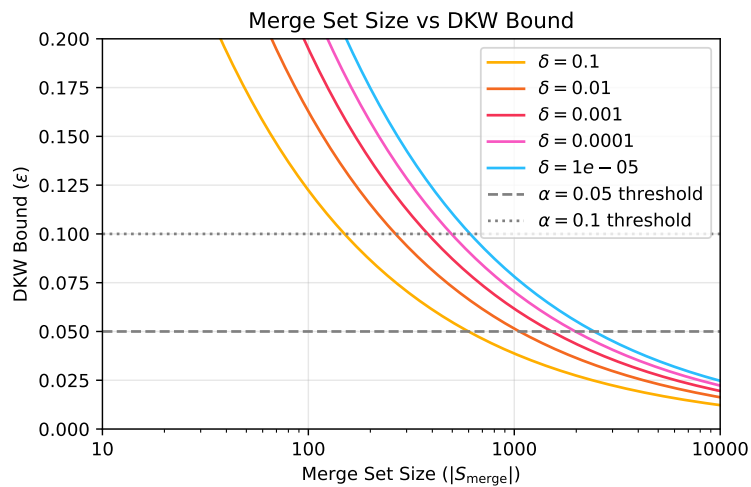


Figure 12: Finite-sample correction ϵ used in ECDF-DKW as a function of the merging set size $|S_{\text{merge}}|$ and user-specified significance level δ . This correction is the offset required to ensure the finite sample guarantee of $(1 - \alpha)(1 - \delta)$ in Stutz et al. (2023). Dashed lines mark α levels of 0.05 and 0.1. When $\epsilon > \alpha$, the prediction set includes all labels and becomes unbounded, yielding 100% coverage.

D A NOTE ON COMPUTATIONAL COMPLEXITY

The computational complexity of our method matches the complexity of existing prediction set aggregation methods for the case of data-independent weights, and includes an additional one-time cost to compute an empirical CDF for the case of data-dependent weights.

The complexity of prediction set aggregation methods was first observed by Cherubin, who compares the computational overhead of their majority vote aggregation method with the overhead of p-value aggregation methods. They note that while majority vote requires simpler operations to determine whether each label is included in the final prediction set, it also requires that predictions be recomputed for each significance level, making it less efficient when sets must be constructed at multiple thresholds. In contrast, p-value methods allow the aggregation to be computed once and then applied to any significance level without additional computation.

These trade-offs in speed and cost may influence which method is better suited to a given application—for example, p-value aggregation may be preferable if prediction sets need to be dynamically thresholded at test time. However, the time complexity of both methods is the same: combining K prediction sets for N test objects with a label space size of $L = |\mathcal{Y}|$ has complexity $\mathcal{O}(KLN)$. This complexity is necessary for all prediction set aggregation methods, as it reflects the cost of evaluating multiple prediction sets across the label space.

With data-independent weights, our method matches this $\mathcal{O}(KLN)$ complexity directly. With data-dependent weights, the only additional computation required is a one-time estimation of a correction factor \hat{m}^* from data split S_{merge} . This step involves computing the empirical CDF of the weighted average p-values on S_{merge} , which has a complexity of $\mathcal{O}(M \log M)$ for a split of size M . Importantly, this correction is computed once and does not require retraining, and the rest of the procedure for data-dependent weights has the same $\mathcal{O}(KLN)$ cost as other prediction set aggregation methods.

E FURTHER IMPLEMENTATION DETAILS

E.1 Conservative empirical CDF

To evaluate the correction factor \hat{m}^* in practice, we approximate the distribution function $F_{P_{\text{all}}}$ by its empirical counterpart on a designated merging set $\mathcal{S}_{\text{merge}}$. The standard empirical CDF of the random variable $P_{\text{all}} = p_{\text{all}}(X, Y; \mathbf{W})$ is

$$\hat{F}_{P_{\text{all}}}(\alpha) = \frac{\sum_{i \in \mathcal{S}_{\text{merge}}} \mathbb{1}\{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) \leq \alpha\}}{|\mathcal{S}_{\text{merge}}|}.$$

When the weights depend on the data, this empirical CDF is used to compute \hat{m}^* . In our experiments, however, we found that the naive estimator can be unstable for small or moderate $|\mathcal{S}_{\text{merge}}|$. To mitigate this, we use a slightly more conservative version,

$$\hat{F}_{P_{\text{all}}}^{\text{cons}}(\alpha) = \frac{\mathbb{1}\{\min_{i \in \mathcal{S}_{\text{merge}}} p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) \leq \alpha\} + \sum_{i \in \mathcal{S}_{\text{merge}}} \mathbb{1}\{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}) \leq \alpha\}}{1 + |\mathcal{S}_{\text{merge}}|}. \quad (17)$$

This modification ensures that the CDF accounts for the minimum observed value and thereby avoids degenerate behavior in finite samples. The correction factor is then computed as

$$\hat{m}^* = \max_{i \in \mathcal{S}_{\text{merge}}} \frac{\hat{F}_{P_{\text{all}}}^{\text{cons}}(p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)}))}{p_{\text{all}}(X_i, Y_i; \mathbf{W}^{(i)})},$$

as in (8). Importantly, using $\hat{F}_{P_{\text{all}}}^{\text{cons}}$ in place of the standard empirical CDF does not affect the theoretical guarantees of our method (see §A.3), but in practice it yields more stable and reliable estimates of \hat{m}^* .

E.2 Synthetic dataset

For our ablative/expository experiments, we generate a simple homoskedastic dataset to simulate a regression task. Each input is a 16-dimensional vector drawn from a standard normal distribution, and the output label is the sum of the feature values with additive Gaussian noise. Specifically, for each sample X_i , we have

$$y_i = \sum_{j=1}^{16} X_{ij} + \varepsilon_i,$$

where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ represents additive noise with standard deviation $\sigma = 0.1$. We generate both training and test datasets by independently drawing samples and computing the corresponding target labels.

E.3 Mixture-of-experts model

In all of our experiments, we use an MoE model of N experts, where each expert is a linear model defined as $f_i(x; \theta_i)$. Here, x represents the input features and θ_i the parameters of the i th expert. Each expert considers a different subset of features depending on the experiment (§E.5 and §B.1.2). The experts are trained independently using L-BFGS to minimize mean square error (MSE) on the training data.

The routing network is also a linear model, and is responsible for generating a set of weights $\{w_i(x)\}_{i=1}^N$ that determine the contribution of each expert to the final prediction. To ensure that the weights are positive and sum to one, softmax is applied to the routing network outputs:

$$w_i(x) = \frac{\exp(g_i(x; \phi))}{\sum_{j=1}^N \exp(g_j(x; \phi))},$$

where $g_i(x; \phi)$ denotes the output of the routing network for the i th expert, and ϕ represents the parameters of the routing network.

The routing network is trained to minimize MSE of the aggregated prediction

$$\ell_{\text{routing}} = \frac{1}{|S_{\text{train}}|} \sum_{i \in S_{\text{train}}} \left(y_i - \sum_{j=1}^N w_j(x_i) f_j(x_i; \theta_j) \right)^2.$$

In our experiments, the trained MoE model is utilized in two ways. For split conformal, we use the weighted sum of experts as a black-box predictive model. For weighted aggregation, we use the weights from the routing network to scale the p-value functions of each expert.

E.4 Computing prediction set length

To compute prediction sets, we leverage what we know about the structure of the p-value function. Observe that for a fixed input x , the function $\hat{p}_k(x, y)$ (2) is piecewise constant with discontinuities only at the values of y that solve $s_k(x, y) = s_k(X_i, Y_i)$. (For an absolute residual score function $s_k(x, y) = |y - \hat{\mu}_k(x)|$, these values are equal to $\hat{\mu}_k(x) \pm s(X_i, Y_i)$ for $i \in S_k$.) For each test point x , \hat{p} is therefore a step function of y with finite discontinuities, and the weighted average p-value function is likewise a step function of y , with its discontinuities as a union of the K separate sets of discontinuities.

We use the `portion` library (Decan) to represent and manipulate our prediction sets; this also allows us to compute their Lebesgue measures without any discretization.

E.5 Expert feature assignment for UCI experiments

For the UCI experiments, the features for each dataset are partitioned into groups of semantically related features, and each expert in the MoE specializes in a single group of features. We list these groups and their features below, where the feature names are provided with the original dataset (Kelly et al., 2010).

Table 1: Airfoil dataset feature groups

Group	Features
Aerodynamics	frequency, free-stream-velocity
Geometry	attack-angle, chord-length, suction-side-displacement-thickness

Table 2: Wine dataset feature groups

Group	Features
Acidity	fixed_acidity, volatile_acidity, citric_acid, pH
Sugar/alcohol	residual_sugar, density, alcohol
Sulfur/salinity	chlorides, free_sulfur_dioxide, total_sulfur_dioxide, sulphates

Table 3: Communities dataset feature groups

Group	Features
Population	population, householdsize, numbUrban, pctUrban, LandArea, PopDens, agePct12t21, agePct12t29, agePct16t24, agePct65up
Race/ethnicity	racepctblack, racePctWhite, racePctAsian, racePctHispanic, PctForeignBorn, PctSpeakEnglOnly, PctNotSpeakEnglWell, PctBornSameState, PctSameHouse85, PctSameCity85
Income/poverty	medIncome, medFamInc, perCapInc, whitePerCap, blackPerCap, indianPerCap, AsianPerCap, hispPerCap, NumUnderPov, PctPopUnderPov
Employment/industry	pctWWage, pctWFarmSelf, pctWInvInc, pctWSocSec, pctWPubAsst, pctWRetire, PctUnemployed, PctEmploy, PctEmplManu, PctEmplProfServ
Occupation/education	PctOccupManu, PctOccupMgmtProf, PctWorkMomYoungKids, PctWorkMom, PctUsePubTrans, PctLess9thGrade, PctNotHSGrad, PctBSorMore, MalePctDivorce, MalePctNevMarr
Family structure	FemalePctDiv, TotalPctDiv, PersPerFam, PctFam2Par, PctKids2Par, PctYoungKids2Par, PctTeen2Par, PctLargHouseFam, PctLargHouseOccup, PctSameState85
Housing characteristics	PersPerOccupHous, PersPerOwnOccHous, PersPerRentOccHous, PctPersOwnOccup, PctHousNoPhone, PctHousLess3BR, MedNumBR, HousVacant, PctHousOccup, PctHousOwnOcc
Housing quality/costs	PctPersDenseHous, PctVacantBoarded, PctVacMore6Mos, MedYrHousBuilt, PctWOFullPlumb, OwnOccLowQuart, OwnOccMedVal, OwnOccHiQuart, RentLowQ, RentMedian
Housing costs/Homelessness	RentHighQ, MedRent, MedRentPctHousInc, MedOwnCostPctInc, MedOwnCostPctIncNoMtg, NumInShelters, NumStreet, NumIlleg, PctIlleg, LemasPctOfficDrugUn
Immigration	NumImmig, PctImmigRecent, PctImmigRec5, PctImmigRec8, PctImmigRec10, PctRecentImmig, PctRecImmig5, PctRecImmig8, PctRecImmig10

Table 4: Superconductivity dataset feature groups

Feature Group	Features
Atomic mass	mean_atomic_mass, wtd_mean_atomic_mass, gmean_atomic_mass, wtd_gmean_atomic_mass, entropy_atomic_mass, wtd_entropy_atomic_mass, range_atomic_mass, wtd_range_atomic_mass, std_atomic_mass, wtd_std_atomic_mass, number_of_element
Atomic radius	mean_atomic_radius, wtd_mean_atomic_radius, gmean_atomic_radius, wtd_gmean_atomic_radius, entropy_atomic_radius, wtd_entropy_atomic_radius, range_atomic_radius, wtd_range_atomic_radius, std_atomic_radius, wtd_std_atomic_radius
Density	mean_Density, wtd_mean_Density, gmean_Density, wtd_gmean_Density, entropy_Density, wtd_entropy_Density, range_Density, wtd_range_Density, std_Density, wtd_std_Density
Electron affinity	mean_ElectronAffinity, wtd_mean_ElectronAffinity, gmean_ElectronAffinity, wtd_gmean_ElectronAffinity, entropy_ElectronAffinity, wtd_entropy_ElectronAffinity, range_ElectronAffinity, wtd_range_ElectronAffinity, std_ElectronAffinity, wtd_std_ElectronAffinity
FIE	mean_fie, wtd_mean_fie, gmean_fie, wtd_gmean_fie, entropy_fie, wtd_entropy_fie, range_fie, wtd_range_fie, std_fie, wtd_std_fie
Fusion heat	mean_FusionHeat, wtd_mean_FusionHeat, gmean_FusionHeat, wtd_gmean_FusionHeat, entropy_FusionHeat, wtd_entropy_FusionHeat, range_FusionHeat, wtd_range_FusionHeat, std_FusionHeat, wtd_std_FusionHeat
Thermal conductivity	mean_ThermalConductivity, wtd_mean_ThermalConductivity, gmean_ThermalConductivity, wtd_gmean_ThermalConductivity, entropy_ThermalConductivity, wtd_entropy_ThermalConductivity, range_ThermalConductivity, wtd_range_ThermalConductivity, std_ThermalConductivity, wtd_std_ThermalConductivity
Valence	mean_Valence, wtd_mean_Valence, gmean_Valence, wtd_gmean_Valence, entropy_Valence, wtd_entropy_Valence, range_Valence, wtd_range_Valence, std_Valence, wtd_std_Valence

Table 5: Adult dataset feature groups

Group	Features
Demographics	age, race, sex
Education	education, education-num
Economic status	fnlwgt, capital-gain, capital-loss, hours-per-week
Family relationship	marital-status, relationship

Table 6: Student dataset feature groups

Group	Features
Personal details	Marital Status, Gender, Age at enrollment, Nationality, International, Mother’s qualification, Father’s qualification, Mother’s occupation, Father’s occupation
Academic details	Application mode, Application order, Previous qualification, Previous qualification (grade), Admission grade, Daytime/evening attendance, Course, Displaced, Educational special needs
Performance	Curricular units 1st sem (credited), Curricular units 1st sem (enrolled), Curricular units 1st sem (evaluations), Curricular units 1st sem (approved), Curricular units 1st sem (grade), Curricular units 1st sem (without evaluations), Curricular units 2nd sem (credited), Curricular units 2nd sem (enrolled), Curricular units 2nd sem (evaluations), Curricular units 2nd sem (approved), Curricular units 2nd sem (grade), Curricular units 2nd sem (without evaluations)
Economic and financial	Scholarship holder, Tuition fees up to date, Debtor, Unemployment rate, Inflation rate, GDP