# A  MATHEMATICAL NOTATIONS

| Notation | Meaning |
|---|---|
| $X \in \mathcal{X}$ | Input example |
| $Y \in \mathcal{Y}$ | The ground-truth label |
| $f$ | The soft classifier |
| $C$ | Number of classes |
| $\mathcal{Z} = (\mathcal{X}, \mathcal{Y})$ | Joint space of input-output pairs |
| $\mathcal{P}_{X,Y}$ | Data distribution of $\mathcal{Z}$ |
| $\mathcal{D}_{\text{tr}}$ | Training data |
| $\mathcal{D}_{\text{cal}}$ | Calibration data |
| $\mathcal{D}_{\text{test}}$ | Test data |
| $\epsilon_c$ | The top-$k$ error for class $c$ |
| $r_f(x, y)$ | The rank of $y$ in prediction $f(x)$ |
| $\widehat{\mathcal{C}}(X)$ | Prediction set for input $X$ |
| $V((X, Y) \in \mathcal{Z})$ | Non-conformity scoring function |
| $\alpha$ | Target mis-coverage rate |

Table 2: **Key notations used in this paper.**

# B  TECHNICAL PROOFS OF THEORETICAL RESULTS

## B.1  PROOF OF PROPOSITION 1

**Proposition 2.** *(Proposition 1 restated, class-conditional over- and under-coverage of MCP) Given $\alpha$, assume $|Rob(\alpha)| < |\mathcal{Y}|$. If there exist $\xi, \xi' > 0$ such that for $y \in Rob(\alpha), y' \notin Rob(\alpha)$:*

$$\mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{class}(Y)|Y = y \in Rob(\alpha)\} - \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{MCP}|Y = y \in Rob(\alpha)\} \leq -\xi,$$

$$\mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{class}(Y)|Y = y' \notin Rob(\alpha)\} - \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{MCP}|Y = y' \notin Rob(\alpha)\}$$
$$\geq \frac{1}{n_{y'}} + \xi'.$$

*Then class $y$ and $y'$ are over- and under-covered, respectively:*

$$\mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{MCP}|Y = y\} \geq 1 - \alpha + \xi, \quad \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{MCP}|Y = y'\} \leq 1 - \alpha - \xi'.$$

*Proof.* (of Proposition 1)

(1) Class $y$ is over-covered. We start with the lower bound of the class-conditional coverage using class-conditional quantile for class $y \in \text{Rob}(\alpha)$, i.e., Theorem 1 in Romano *et al.* (2020), as follows.

$$1 - \alpha \leq \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y)|Y = y \in \text{Rob}(\alpha)\}$$
$$= \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = y \in \text{Rob}(\alpha)\}$$
$$+ \mathbb{1}[y \in \text{Rob}(\alpha)] \cdot \left( \mathbb{P}\left\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y)]\Big|Y = y \in \text{Rob}(\alpha)\right\} \right.$$
$$\left. - \mathbb{P}\left\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}\Big|Y = c \in \text{Rob}(\alpha)\right\}\right)$$
$$+ \mathbb{1}[y \notin \text{Rob}(\alpha)] \cdot \left( \mathbb{P}\left\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y)\Big|Y = c \notin \text{Rob}(\alpha)\right\} \right.$$
$$\left. - \mathbb{P}\left\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}\Big|Y = c \notin \text{Rob}(\alpha)\right\}\right). \tag{13}$$

By assumption, for $y \in \text{Rob}(\alpha)$ and $\xi > 0$ such that the following inequality holds:

$$\mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y)|Y = y \in \text{Rob}(\alpha)\} - \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = y \in \text{Rob}(\alpha)\} \leq -\xi, \tag{14}$$

by plugging inequality (14) into inequality (13), we derive the class-conditional over-coverage of CP on class $y$:

$$\mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = y \in Rob(\alpha)\} \geq 1 - \alpha + \xi.$$

(2) Class $c'$ is under-covered. We start with the upper bound of the class-conditional coverage using class-conditional quantile for class $y' \notin \text{Rob}(\alpha)$, i.e., Theorem 1 in Romano *et al.* (2020), as follows.

$$1 - \alpha + \frac{1}{n_{y'}} \geq \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y')|Y = y' \notin \text{Rob}(\alpha)\}$$

$$= \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = y' \notin \text{Rob}(\alpha)\}$$

$$+ \mathbb{1}[y' \in \text{Rob}(\alpha)] \cdot \Bigg( \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y')|Y = y' \notin \text{Rob}(\alpha)\}$$

$$- \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = y' \notin \text{Rob}(\alpha)\} \Bigg)$$

$$+ \mathbb{1}[y' \notin \text{Rob}(\alpha)] \cdot \Bigg( \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y')|Y = y' \notin \text{Rob}(\alpha)\}$$

$$- \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = y' \notin \text{Rob}(\alpha)\} \Bigg). \tag{15}$$

By assumption, for $y' \notin \text{Rob}(\alpha)$ and $\xi' > 0$ such that the following inequality holds:

$$\mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{class}}(y')|Y = y' \notin \text{Rob}(\alpha)\} - \mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = y' \notin \text{Rob}(\alpha)\}$$

$$\geq \frac{1}{n_{y'}} + \xi', \tag{16}$$

by plugging inequality (16) into inequality (15), we derive the class-conditional under-coverage of CP on class $y'$:

$$\mathbb{P}\{V(X,Y) \leq Q_{1-\alpha}^{\text{MCP}}|Y = c' \notin \text{Rob}(\alpha)\} \leq 1 - \alpha - \xi'.$$

The above two results of over-coverage of class $y$ and under-coverage of class $y'$ show that the class-conditional coverage can easily deviate from the marginal coverage as long as there exists a margin (i.e., $\xi, \xi'$) for class-conditional coverage between using marginal quantile (i.e., $Q_{1-\alpha}^{\text{MCP}}$) and class-conditional quantile (i.e., $Q_{1-\alpha}^{\text{class}}(y), Q_{1-\alpha}^{\text{class}}(y')$), as in (14) and (16). $\qquad\square$

### B.2 Proof of Theorem 1

**Theorem 3.** *(Theorem 1 restated, $k$-CCP guarantees class-conditional coverage) Suppose that selecting $\{k(y)\}_{y\in\mathcal{Y}}$ results in class-wise top-$k(y)$ error $\{\epsilon_y\}_{y\in\mathcal{Y}}$. If the nominated mis-coverage probability $\tilde{\alpha}_y$ of CCP for class $y$ is set as*

$$\tilde{\alpha}_y \leq \alpha - \varepsilon_{n_y} - \delta - \epsilon_y, \quad \text{for } 0 < \delta < 1, \ \varepsilon_{n_y} = \sqrt{(3(1-\alpha)\log(2/\delta))/n_y},$$

*then $k$-CCP can achieve the class-conditional coverage as defined in equation (1).*

Before proving Theorem 1, we introduce the following technical lemma.

**Lemma 1.** *(Concentration inequalities for quantiles) Define $\varepsilon_n = \sqrt{3(1-\alpha)\log(2/\delta)/n}$. Let $Q_{1-\alpha} = max\{t : \mathbb{P}_V\{V \leq t\} \geq 1 - \alpha\}$ be the true quantile of a random variable $V$ given $\alpha$, and $\widehat{Q}_{1-\alpha} = V_{(\lceil n(1-\alpha)\rceil)}$ be the empirical quantile estimated by $n$ randomly sampled set $\{V_1, ..., V_n\}_{i=1}^n$. Then with probability at least $1 - \delta$, we have $\widehat{Q}_{1-\alpha-\varepsilon_n-1/n} \leq Q_{1-\alpha} \leq \widehat{Q}_{1-\alpha+\varepsilon_n}$ where $\tilde{O}$ hides the logarithmic factor.*

Lemma 1 has been studied in a previous paper (see Proposition 2(a) in Vovk (2012)). To make the proof and conclusion of the Lemma 1 complete with the Theorem 1, we prove it again at the end of this subsection. Now we begin to prove Theorem 1.

*Proof.* (of Theorem 1)

Let $y \in \mathcal{Y}$ denote any class label. With the lower bound of the coverage on class $y$ (Theorem 1 in Romano *et al.* (2020)), we have

$$1 - (\tilde{\alpha} + \varepsilon_{n_y})$$
$$\leq \mathbb{P}\{Y_{n+1} \in \mathcal{C}_{1-\tilde{\alpha}-\varepsilon_{n_y}}^{\text{CCP}}(X_{n+1})|Y=y\} = \mathbb{P}\{V(X_{n+1}, Y_{n+1}) \leq Q_{1-\tilde{\alpha}-\varepsilon_{n_y}}^{\text{class}}|Y=y\}$$
$$= (\mathbb{P}\{\text{Lemma 1 holds}\} + \mathbb{P}\{\text{Lemma 1 not holds}\}) \cdot \mathbb{P}\{V(X_{n+1}, Y_{n+1}) \leq Q_{1-\tilde{\alpha}-\varepsilon_{n_y}}^{\text{class}}|Y=y\}$$
$$\leq 1 \cdot \mathbb{P}\{V(X_{n+1}, Y_{n+1}) \leq \widehat{Q}_{1-\tilde{\alpha}}^{\text{class}}|Y=y\} + \delta \cdot 1$$
$$= \mathbb{P}\{V(X_{n+1}, Y_{n+1}) \leq \widehat{Q}_{1-\tilde{\alpha}}^{\text{class}}, r_f(X_{n+1}, Y_{n+1}) \leq \widehat{k}(y)|Y=y\}$$
$$\quad + \mathbb{P}\{V(X_{n+1}, Y_{n+1}) \leq \widehat{Q}_{1-\tilde{\alpha}}^{\text{class}}, r_f(X_{n+1}, Y_{n+1}) > \widehat{k}(y)|Y=y\} + \delta$$
$$\leq \mathbb{P}\{V(X_{n+1}, Y_{n+1}) \leq \widehat{Q}_{1-\tilde{\alpha}}^{\text{class}}, r_f(X_{n+1}, Y_{n+1}) \leq \widehat{k}(y)|Y=y\}$$
$$\quad + \underbrace{\mathbb{P}\{r_f(X_{n+1}, Y_{n+1}) > \widehat{k}(y)|Y=y\}}_{\leq \epsilon_y} + \delta$$
$$\leq \mathbb{P}\{Y_{n+1} \in \widehat{\mathcal{C}}_{1-\tilde{\alpha}}^{k\text{-CCP}}(y)|Y=y\} + \epsilon_y + \delta,$$

where the second inequality is due to Lemma 1, and the last inequality is due to the definition of $\epsilon_y$, i.e., $\mathbb{P}\{r_f(X_{n+1}, Y_{n+1}) \leq \widehat{k}(y)|Y=y\} \geq 1 - \epsilon_y$.

Re-arranging the above inequality, we have

$$\mathbb{P}\{Y_{n+1} \in \widehat{\mathcal{C}}_{1-\tilde{\alpha}}^{k\text{-CCP}}(y)|Y=y\} \geq 1 - \tilde{\alpha} - \varepsilon_{n_y} - \delta - \epsilon_y \geq 1 - \alpha,$$

where the last inequality is due to $\tilde{\alpha}_y \leq \alpha - \varepsilon_{n_y} - \delta - \epsilon_y$. This implies that $k$-CCP guarantees the class-conditional coverage on any class $y$. $\square$

After proving Theorem 1, now we show the deferred proof of lemma 1:

*Proof.* (of Lemma 1)

Define $Z_i = \mathbb{1}[V_i \leq Q_{1-\alpha}]$ where $1 \leq i \leq n$ and $\mathbb{1}[\cdot]$ is an indicator function. Then $Z_i$ is a Bernoulli random variable with $\mathbb{P}\{Z_i = 1\} = 1 - \alpha$ and $\mathbb{P}\{Z_i = 0\} = \alpha$ from the definition of $Q(\alpha)$. Let $\widehat{Z} = \frac{1}{n}\sum_{i=1}^{n} Z_i$ and $\mathbb{E}[\widehat{Z}] = 1 - \alpha$.

According to Chernoff bound, we know

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{i=1}^{n} Z_i - \mathbb{E}[\widehat{Z}]\right| \geq \varepsilon\mathbb{E}[\widehat{Z}]\right\} \leq 2\exp\left(-n\mathbb{E}[\widehat{Z}]\varepsilon^2/3\right) = 2\exp\left(-n(1-\alpha)\varepsilon^2/3\right).$$

By setting $\delta = 2\exp(-n(1-\alpha)\varepsilon^2/3)$, i.e., $\varepsilon = \frac{\varepsilon_n}{1-\alpha} = \sqrt{(3\log(2/\delta))/((1-\alpha)n)}$, we have with probability at least $1 - \delta$:

$$\left|\frac{1}{n}\sum_{i=1}^{n} \mathbb{1}[V_i \leq Q_{1-\alpha}] - (1-\alpha)\right| \leq (1-\alpha)\varepsilon = \frac{\varepsilon_n}{1-\alpha}(1-\alpha) = \varepsilon_n = \sqrt{(3(1-\alpha)\log(2/\delta))/n}. \tag{17}$$

Recall the definition of the empirical quantile $\widehat{Q}_{1-\alpha}$ given $\alpha$:

$$\widehat{Q}_{1-\alpha} = \max\left\{t : \frac{1}{n}\sum_{i=1}^{n} \mathbb{1}[V_i \leq t] \geq 1 - \alpha\right\}.$$

Then we know the following upper bound and lower bound for $1 - \alpha$:

$$\frac{1}{n}\sum_{i=1}^{n}\mathbb{1}[V_i \le \widehat{Q}_{1-\alpha-1/n}] \le (1-\alpha) \le \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}[V_i \le \widehat{Q}_{1-\alpha}]. \tag{18}$$

Re-arranging (17) and using the above upper/lower bounds, with probability at least $1 - \delta$, we have

$$(1-\alpha)(1-\varepsilon) \le \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}[V_i \le Q_{1-\alpha}] \le (1-\alpha)(1+\varepsilon)$$

$$\Leftrightarrow \quad 1 - \underbrace{(1-(1-\alpha)(1-\varepsilon))}_{\triangleq \alpha'} \le \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}[V_i \le Q_{1-\alpha}] \le 1 - \underbrace{(1-(1-\alpha)(1+\varepsilon))}_{\triangleq \alpha''}$$

$$\stackrel{(18)}{\Rightarrow} \quad \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}[V_i \le \widehat{Q}_{1-\alpha'-1/n}] \le \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}[V_i \le Q_{1-\alpha}] \le \frac{1}{n}\sum_{i=1}^{n}\mathbb{1}[V_i \le \widehat{Q}_{1-\alpha''}]$$

$$\Leftrightarrow \quad \widehat{Q}_{1-\alpha'-1/n} \le Q_{1-\alpha} \le \widehat{Q}_{1-\alpha''}. \tag{19}$$

Finally, we simplify $\alpha'$ and $\alpha''$ as follows

$$\alpha' = 1 - (1-\alpha)(1-\varepsilon) = \alpha + \varepsilon(1-\alpha) = \alpha + \sqrt{3(1-\alpha)\log(2/\delta)/n} = \alpha + \varepsilon_n,$$

$$\alpha'' = 1 - (1-\alpha)(1+\varepsilon) = \alpha - \varepsilon(1-\alpha) = \alpha - \sqrt{3(1-\alpha)\log(2/\delta)/n} = \alpha - \varepsilon_n. \tag{20}$$

Therefore, plugging (20) into (19), we have

$$\widehat{Q}_{1-\alpha-\varepsilon_n-1/n} \le Q_{1-\alpha} \le \widehat{Q}_{1-\alpha+\varepsilon_n}.$$

$\square$

## B.3 PROOF OF THEOREM 2

**Theorem 4.** *(Theorem 2 restated, $k$-CCP produces smaller prediction sets than CCP) Suppose the following inequality holds for any $y \in \mathcal{Y}$:*

$$\sum_{y\in\mathcal{Y}} \sigma_y \cdot \mathbb{P}_{Z_{n+1}}\Big[V(X_{n+1},y) \le \widehat{Q}_{1-\alpha}^{class}(y)\Big] \le \sum_{y\in\mathcal{Y}} \mathbb{P}_{Z_{n+1}}\Big[V(X_{n+1},y) \le \widehat{Q}_{1-\alpha}^{class}(y)\Big].$$

*Then $k$-CCP produces smaller expected prediction sets than CCP, i.e.,*

$$\mathbb{E}_{X_{n+1}}[|\widehat{\mathcal{C}}_{1-\tilde{\alpha}}^{k\text{-}CCP}(X_{n+1})|] \le \mathbb{E}_{X_{n+1}}[|\widehat{\mathcal{C}}_{1-\alpha}^{CCP}(X_{n+1})|].$$

*Proof.* (of Theorem 2)

The proof idea is to reduce the the cardinality of the prediction set made by $k$-CCP to that made by CCP in expectation.

$$\mathbb{E}_{Z_{n+1}}[|\widehat{\mathcal{C}}_{1-\tilde{\alpha}}^{k\text{-}CCP}(X_{n+1})|] = \mathbb{E}_{Z_{n+1}}\Bigg[\sum_{y\in\mathcal{Y}}\mathbb{1}\Big[V(X_{n+1},y) \le \widehat{Q}_{1-\tilde{\alpha}}^{class}(y),\ r_f(X_{n+1},y) \le \widehat{k}(y)\Big]\Bigg]$$

$$= \sum_{y\in\mathcal{Y}}\mathbb{E}_{Z_{n+1}}\Big[\mathbb{1}[V(X_{n+1},y) \le \widehat{Q}_{1-\tilde{\alpha}}^{class}(y),\ r_f(X_{n+1},y) \le \widehat{k}(y)]\Big]$$

$$= \sum_{y\in\mathcal{Y}}\mathbb{P}_{Z_{n+1}}\Big[V(X_{n+1},y) \le \widehat{Q}_{1-\tilde{\alpha}}^{class}(y),\ r_f(X_{n+1},y) \le \widehat{k}(y)\Big]$$

$$\stackrel{(a)}{=} \sum_{y\in\mathcal{Y}}\sigma_y \cdot \mathbb{P}_{Z_{n+1}}\Big[V(X_{n+1},y) \le \widehat{Q}_{1-\alpha}^{class}(y)\Big] \stackrel{(b)}{\le} \sum_{y\in\mathcal{Y}}\mathbb{E}_{Z_{n+1}}\Big[\mathbb{1}[V(X_{n+1},y) \le \widehat{Q}_{1-\alpha}^{class}(y)]\Big]$$

$$= \mathbb{E}_{Z_{n+1}}\Bigg[\sum_{y\in\mathcal{Y}}\mathbb{1}[V(X_{n+1},y) \le \widehat{Q}_{1-\alpha}^{class}(y)]\Bigg] = \mathbb{E}_{Z_{n+1}}[|\widehat{\mathcal{C}}_{1-\alpha}^{CCP}(X_{n+1})|], \tag{21}$$

where the equality $(a)$ is due to the definitions of $\sigma_y$, and inequality $(b)$ is due to the assumption

$$\sum_{y \in \mathcal{Y}} \sigma_y \cdot \mathbb{P}_{Z_{n+1}}\left[V(X_{n+1}, y) \leq \widehat{Q}_{1-\alpha}^{\text{class}}(y)\right] \leq \sum_{y \in \mathcal{Y}} \mathbb{P}_{Z_{n+1}}\left[V(X_{n+1}, y) \leq \widehat{Q}_{1-\alpha}^{\text{class}}(y)\right].$$

This implies that $k$-CCP requires smaller prediction sets to guarantee the class-conditional coverage compared to CCP. □

## C   COMPLETE EXPERIMENTAL RESULTS

### C.1   TRAINING DETAILS

For CIFAR-10 and CIFAR-100, we train ResNet20 using LDAM loss function given in Cao *et al.* (2019) with standard mini-batch stochastic gradient descent (SGD) using learning rate $0.1$, momentum $0.9$, and weight decay $2e - 4$ for 200 epochs. The batch size is $128$. For experiments on mini-ImageNet, we use the same setting. For Food-101, the batch size is $256$ and other parameters are kept the same.

### C.2   CALIBRATION DETAILS

As mentioned in Section 5.1, we balanced split the validation set of CIFAR-10 and CIFAR-100, the number of calibration data is $5000$. For mini-ImageNet, the number of calibration data is $15000$. For Food-101, the total number is $12625$. To compute the mean and standard deviation for the overall performance, we repeat calibration experiments for 10 times. Moreover, we select the $g$ from the interval $[0.1, 1]$ with range $0.05$ to find the minimal $g$ that $k$-CCP and CCP achieves the target class-conditional coverage. We re-emphasize that the we have discussed the assumption in Theorem 2 and Remark 3.

The regularization parameter for RAPS scoring function is from the set $k_{reg} \in \{3, 5, 7\}$ and $\lambda \in \{0.001, 0.01, 0.1\}$ based on the empirical setting in `cluster-CP`. We select the combination of $k_{reg}$ and $\lambda$ for each experiments with same imbalanced type and imbalanced ratio on the same dataset, where the most of $APSS$ values of all methods are minimum. The hyper-parameter $g$ is selected from the interval $[0.1, 1]$ with range $0.05$ to find the minimal $g$ that `CCP`, `cluster-CP`, and $k$-CCP achieve the target class-conditional coverage.

### C.3   ILLUSTRATION OF IMBALANCED DATA



Figure 3: Illustrative examples of the different imbalanced distributions of the number of training examples per class index $c$ on CIFAR-100

### C.4   COMPLETE EXPERIMENT RESULTS

In this subsection, we report complete experimental results over four datasets, three decaying types, five imbalance ratios. Specifically, Table 3, 4, 5 report results on CIFAR-10 with three decaying types. Table 6, 7, 8 report results on CIFAR-100 with three decaying types. Table 9, 10, 11 report results on mini-ImageNet with three decaying types. Table 12, 13, 14 report results on Food-101 with three decaying types. Due to the limited time for rebuttal, we just report the results of

cluster-CP with two imbalanced ratio $\rho = 0.1$ and $\rho = 0.5$. We will add all results in the final paper.

Figure 4 shows overall results on CIFAR-10 and CIFAR-100, including distribution of class-wise quantiles v.s. marginal quantile, histograms of class-conditional coverage and prediction set size achieved by MCP, CCP, cluster-CP, and $k$-CCP, and the histogram of condition numbers $\sigma_y$ in Theorem 2, which are all corresponding to Figure 1 in main text.

Figure 5 shows the sensitivity of CCP, cluster-CP, and $k$-CCP for $g$ on CIFAR-10 and CIFAR-100 with APS scoring function, which are all corresponding to Figure 2 in main text.

Figure 6, Figure 7, Figure 8, Figure 9, Figure 10 and Figure 11 show the class-conditional coverage and the corresponding prediction set sizes on EXP $\rho = 0.1$, EXP $\rho = 0.5$, POLY $\rho = 0.1$, POLY $\rho = 0.5$, MAJ $\rho = 0.1$, MAJ $\rho = 0.5$, respectively. This result on EXP $\rho = 0.1$ is in Figure 1 and Figure 4. Because of the same reason, we lack of visualization results of cluster-CP methods and we will add the complete visualization results in the final paper.

Figure 12, Figure 13, and Figure 14 show the distribution of class-wise quantiles with EXP $\rho = 0.5$, POLY $\rho = 0.5$, and MAJ $\rho = 0.5$, respectively.

Figure 15, Figure 16, Figure 17, Figure 18, Figure 19 and Figure 20 verify the condition numbers $\sigma_y$ on EXP $\rho = 0.1$, EXP $\rho = 0.5$, POLY $\rho = 0.1$, POLY $\rho = 0.5$, MAJ $\rho = 0.1$, MAJ $\rho = 0.5$, respectively.



Figure 4: Justification experiments: CIFAR-10 in first row and CIFAR-100 in second row with ResNet20 model. First column: distribution of class-wise quantiles v.s. marginal quantile with imbalance type EXP and imbalance ratio $\rho = 0.5$. Second and third columns: histograms of class-conditional coverage and prediction set size achieved by MCP, CCP, cluster-CP, and $k$-CCP with imbalance type EXP and imbalance ratio $\rho = 0.1$. The final column: the histogram of condition numbers $\sigma_y$ in Theorem 2 with imbalance type EXP and imbalance ratio $\rho = 0.1$.



Figure 5: Results for under coverage ratio and average prediction set size achieved by CCP, cluster-CP, and $k$-CCP methods as a function of $g$ using APS scoring function with imbalance type EXP for imbalance ratio $\rho = 0.1$. $k$-CCP degenerates to CCP in CIFAR-10, so overlopping with CCP (the black line overlaps with the red one).

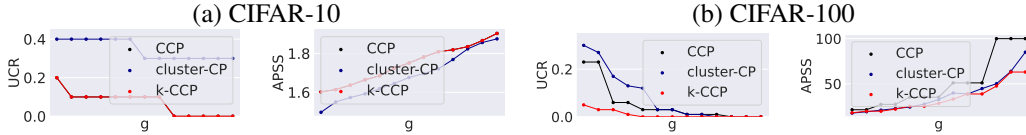| Measure | Methods | EXP | | | | |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
|---|---|---|---|---|---|---|
| UCR | MCP | $0.460 \pm 0.025$ | $0.390 \pm 0.03$ | $0.4 \pm 0.032$ | $0.4 \pm 0.014$ | $0.49 \pm 0.026$ |
| | CCP | $\mathbf{0.110 \pm 0.030}$ | $\mathbf{0.14 \pm 0.029}$ | $\mathbf{0.140 \pm 0.032}$ | $\mathbf{0.08 \pm 0.031}$ | $\mathbf{0.040 \pm 0.020}$ |
| | cluster-CP | $0.160 \pm 0.025$ | — | — | — | $0.080 \pm 0.012$ |
| | $k$-CCP | $\mathbf{0.110 \pm 0.030}$ | $\mathbf{0.14 \pm 0.029}$ | $\mathbf{0.140 \pm 0.032}$ | $\mathbf{0.008 \pm 0.031}$ | $\mathbf{0.040 \pm 0.020}$ |
| APSS | MCP | $1.132 \pm 0.033$ | $1.13 \pm 0.03$ | $1.169 \pm 0.034$ | $1.227 \pm 0.039$ | $1.406 \pm 0.035$ |
| | CCP | $1.481 \pm 0.082$ | $\mathbf{1.508 \pm 0.090}$ | $\mathbf{1.625 \pm 0.083}$ | $\mathbf{1.711 \pm 0.101}$ | $\mathbf{2.032 \pm 0.096}$ |
| | cluster-CP | $\mathbf{1.445 \pm 0.017}$ | — | — | — | $2.323 \pm 0.0115$ |
| | $k$-CCP | $1.481 \pm 0.082$ | $\mathbf{1.508 \pm 0.090}$ | $\mathbf{1.625 \pm 0.083}$ | $\mathbf{1.711 \pm 0.101}$ | $\mathbf{2.032 \pm 0.096}$ |

Table 3: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type EXP and APS scoring function on dataset CIFAR-10. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

| Measure | Methods | POLY | | | | |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
|---|---|---|---|---|---|---|
| UCR | MCP | $0.51 \pm 0.033$ | $0.54 \pm 0.021$ | $0.41 \pm 0.03$ | $0.47 \pm 0.028$ | $0.32 \pm 0.031$ |
| | CCP | $0.14 \pm 0.032$ | $\mathbf{0.17 \pm 0.028}$ | $\mathbf{0.12 \pm 0.031}$ | $\mathbf{0.07 \pm 0.025}$ | $\mathbf{0.01 \pm 0.009}$ |
| | cluster-CP | $\mathbf{0.14 \pm 0.021}$ | — | — | — | $\mathbf{0.01 \pm 0.009}$ |
| | $k$-CCP | $0.14 \pm 0.032$ | $\mathbf{0.17 \pm 0.028}$ | $\mathbf{0.12 \pm 0.031}$ | $\mathbf{0.07 \pm 0.025}$ | $\mathbf{0.01 \pm 0.009}$ |
| APSS | MCP | $1.17 \pm 0.028$ | $1.107 \pm 0.028$ | $1.138 \pm 0.032$ | $1.57 \pm 0.033$ | $1.214 \pm 0.038$ |
| | CCP | $\mathbf{1.487 \pm 0.09}$ | $\mathbf{1.465 \pm 0.090}$ | $\mathbf{1.571 \pm 0.086}$ | $\mathbf{1.652 \pm 0.084}$ | $\mathbf{1.945 \pm 0.087}$ |
| | cluster-CP | $\mathbf{1.612 \pm 0.013}$ | — | — | — | $\mathbf{2.102 \pm 0.015}$ |
| | $k$-CCP | $\mathbf{1.487 \pm 0.09}$ | $\mathbf{1.465 \pm 0.090}$ | $\mathbf{1.571 \pm 0.086}$ | $\mathbf{1.652 \pm 0.084}$ | $\mathbf{1.945 \pm 0.087}$ |

Table 4: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type POLY and APS scoring function on dataset CIFAR-10. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

| Measure | Methods | MAJ | | | | |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
|---|---|---|---|---|---|---|
| UCR | MCP | $0.38 \pm 0.019$ | $0.33 \pm 0.014$ | $0.45 \pm 0.029$ | $0.51 \pm 0.03$ | $0.5 \pm 0.014$ |
| | CCP | $\mathbf{0.12 \pm 0.024}$ | $\mathbf{0.13 \pm 0.025}$ | $\mathbf{0.11 \pm 0.03}$ | $\mathbf{0.06 \pm 0.021}$ | $0.01 \pm 0.009$ |
| | cluster-CP | $0.15 \pm 0.025$ | — | — | — | $\mathbf{0.009 \pm 0.013}$ |
| | $k$-CCP | $\mathbf{0.12 \pm 0.024}$ | $\mathbf{0.13 \pm 0.025}$ | $\mathbf{0.11 \pm 0.03}$ | $\mathbf{0.06 \pm 0.021}$ | $0.01 \pm 0.009$ |
| APSS | MCP | $1.17 \pm 0.028$ | $1.107 \pm 0.028$ | $1.138 \pm 0.032$ | $1.57 \pm 0.033$ | $1.406 \pm 0.035$ |
| | CCP | $\mathbf{1.487 \pm 0.09}$ | $\mathbf{1.465 \pm 0.090}$ | $\mathbf{1.571 \pm 0.086}$ | $\mathbf{1.652 \pm 0.084}$ | $\mathbf{2.032 \pm 0.096}$ |
| | cluster-CP | $1.787 \pm 0.019$ | — | — | — | $2.969 \pm 0.025$ |
| | $k$-CCP | $\mathbf{1.481 \pm 0.082}$ | $\mathbf{1.508 \pm 0.090}$ | $\mathbf{1.625 \pm 0.083}$ | $\mathbf{1.711 \pm 0.101}$ | $\mathbf{2.032 \pm 0.096}$ |

Table 5: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type MAJ and APS scoring function on dataset CIFAR-10. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

| Measure | Methods | EXP | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | 0.386±(0.009) | 0.394±(0.010) | 0.361±(0.012) | 0.382±(0.011) | 0.384±(0.018) |
| | CCP | 0.009±(0.003) | 0.015±(0.004) | 0.011±(0.004) | 0.016±(0.003) | 0.011±(0.002) |
| | cluster-CP | 0.004±(0.002) | — | — | — | 0.004±(0.002) |
| | $k$-**CCP** | **0.0±(0.0)** | **0.0±(0.0)** | **0.0±(0.0)** | **0.0±(0.0)** | **0.001±(0.001)** |
| APSS | MCP | 10.303±(0.111) | 10.848±(0.104) | 12.480±(0.113) | 12.909±(0.115) | 14.544±(0.119) |
| | CCP | 44.194±(0.514) | 44.447±(0.566) | 47.688±(0.569) | 46.955±(0.500) | 50.963±(0.482) |
| | cluster-CP | 30.922±(0.454) | — | — | — | 43.883±(1.070) |
| | $k$-**CCP** | **20.355±(0.357)** | **20.540±(0.356)** | **22.550±(0.306)** | **23.163±(0.265)** | **25.185±(0.279)** |

Table 6: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type EXP and APS scoring function on dataset CIFAR-100. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

| Measure | Methods | POLY | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | 0.395±(0.010) | 0.382±(0.014) | 0.409±(0.013) | 0.383±(0.015) | 0.410±(0.010) |
| | CCP | 0.011±(0.003) | 0.008±(0.002) | 0.016±(0.003) | 0.011±(0.004) | 0.015±(0.003) |
| | cluster-CP | 0.001±(0.001) | 0.008±(0.002) | 0.016±(0.003) | 0.011±(0.004) | 0.015±(0.003) |
| | $k$-**CCP** | **0.0±(0.0)** | **0.0±(0.0)** | **0.0±(0.0)** | **0.0±(0.0)** | **0.0±(0.0)** |
| APSS | MCP | 15.730±(0.126) | 16.738±(0.170) | 17.670±(0.165) | 20.422±(0.211) | 25.888±(0.197) |
| | CCP | 49.896±(0.490) | 54.011±(0.572) | 56.018±(0.529) | 59.893±(0.438) | 64.366±(0.390) |
| | cluster-CP | 56.696±(0.393) | — | — | — | 63.208±(0.364) |
| | $k$-**CCP** | **25.843±(0.300)** | **26.851±(0.222)** | **29.655±(0.286)** | **31.933±(0.218)** | **37.035±(0.245)** |

Table 7: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type POLY and APS scoring function on dataset CIFAR-100 that will be finished in the final version. The dash symbol (—) means missing results. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

| Measure | Methods | MAJ | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | 0.352±(0.010) | 0.412±(0.014) | 0.361±(0.009) | 0.374±(0.010) | 0.401±(0.008) |
| | CCP | 0.017±(0.003) | 0.008±(0.002) | 0.018±(0.004) | 0.011±(0.004) | 0.008±(0.004) |
| | cluster-CP | 0.005±(0.002) | — | — | — | 0.019±(0.005) |
| | $k$-**CCP** | **0.002±(0.001)** | **0.001±(0.001)** | **0.000±(0.000)** | **0.000±(0.000)** | **0.001±(0.001)** |
| APSS | MCP | 11.680±(0.117) | 14.034±(0.128) | 15.171±(0.121) | 18.516±(0.152) | 23.796±(0.159) |
| | CCP | 48.323±(0.548) | 49.193±(0.403) | 53.688±(0.537) | 55.024±(0.402) | 64.640±(0.621) |
| | cluster-CP | 33.623±(0.395) | — | — | — | 50.382±(0.711) |
| | $k$-**CCP** | **21.196±(0.320)** | **24.058±(0.282)** | **25.501±(0.249)** | **29.344±(0.315)** | **35.630±(0.232)** |

Table 8: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type MAJ and APS scoring function on dataset CIFAR-100 that will be finished in the final version. The dash symbol (—) means missing results. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

| Measure | Methods | EXP | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | $0.408 \pm 0.008$ | $0.424 \pm 0.008$ | $0.406 \pm 0.009$ | $0.405 \pm 0.01$ | $0.414 \pm 0.012$ |
| | CCP | $0.007 \pm 0.003$ | $0.001 \pm 0.001$ | $0.0 \pm 0.0$ | $0.002 \pm 0.002$ | $0.001 \pm 0.001$ |
| | cluster-CP | $0.009 \pm 0.003$ | — | — | — | $0.002 \pm 0.001$ |
| | $k$-**CCP** | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ |
| APSS | MCP | $9.705 \pm 0.101$ | $9.498 \pm 0.102$ | $9.438 \pm 0.098$ | $9.361 \pm 0.092$ | $8.931 \pm 0.093$ |
| | CCP | $26.666 \pm 0.415$ | $30.437 \pm 0.352$ | $29.777 \pm 0.409$ | $30.007 \pm 0.33$ | $34.867 \pm 0.445$ |
| | cluster-CP | $27.786 \pm 0.307$ | — | — | — | $33.114 \pm 0.418$ |
| | $k$-**CCP** | $\mathbf{18.129 \pm 0.454}$ | $\mathbf{17.546 \pm 0.453}$ | $\mathbf{18.944 \pm 0.381}$ | $\mathbf{18.81 \pm 0.368}$ | $\mathbf{17.769 \pm 0.463}$ |

Table 9: Results comparing `MCP`, `CCP`, `cluster-CP`, and $k$-`CCP` with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type EXP and APS scoring function on dataset mini-ImageNet. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-`CCP` the same as or better than that of `CCP` and `cluster-CP` for a fair comparison of prediction set size.

| Measure | Methods | POLY | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | $0.412 \pm 0.013$ | $0.4 \pm 0.009$ | $0.427 \pm 0.011$ | $0.407 \pm 0.01$ | $0.41 \pm 0.009$ |
| | CCP | $0.005 \pm 0.002$ | $0.002 \pm 0.001$ | $0.003 \pm 0.001$ | $0.001 \pm 0.001$ | $0.001 \pm 0.001$ |
| | cluster-CP | $0.028 \pm 0.005$ | — | — | — | $0.015 \pm 0.003$ |
| | $k$-**CCP** | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ |
| APSS | MCP | $9.81 \pm 0.102$ | $9.838 \pm 0.091$ | $9.801 \pm 0.107$ | $9.528 \pm 0.099$ | $9.665 \pm 0.101$ |
| | CCP | $26.620 \pm 0.369$ | $30.236 \pm 0.331$ | $30.912 \pm 0.401$ | $31.639 \pm 0.422$ | $29.852 \pm 0.36$ |
| | cluster-CP | $21.273 \pm 0.369$ | — | — | — | $25.550 \pm 0.279$ |
| | $k$-**CCP** | $\mathbf{17.784 \pm 0.438}$ | $\mathbf{17.751 \pm 0.466}$ | $\mathbf{19.388 \pm 0.441}$ | $\mathbf{19.342 \pm 0.443}$ | $\mathbf{19.153 \pm 0.412}$ |

Table 10: Results comparing `MCP`, `CCP`, `cluster-CP`, and $k$-`CCP` with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type POLY and APS scoring function on dataset mini-ImageNet. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-`CCP` the same as or better than that of `CCP` and `cluster-CP` for a fair comparison of prediction set size.

| Measure | Methods | MAJ | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | $0.408 \pm 0.009$ | $0.405 \pm 0.012$ | $0.424 \pm 0.01$ | $0.43 \pm 0.01$ | $0.411 \pm 0.007$ |
| | CCP | $0.011 \pm 0.004$ | $0.005 \pm 0.002$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ |
| | cluster-CP | $0.015 \pm 0.005$ | — | — | — | $0.011 \pm 0.003$ |
| | $k$-**CCP** | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ |
| APSS | MCP | $9.84 \pm 0.091$ | $9.929 \pm 0.112$ | $9.817 \pm 0.092$ | $9.499 \pm 0.094$ | $9.123 \pm 0.086$ |
| | CCP | $27.306 \pm 0.377$ | $31.114 \pm 0.456$ | $30.741 \pm 0.345$ | $30.608 \pm 0.433$ | $34.186 \pm 0.32$ |
| | cluster-CP | $25.288 \pm 0.226$ | — | — | — | $25.229 \pm 0.352$ |
| | $k$-**CCP** | $\mathbf{18.11 \pm 0.414}$ | $\mathbf{17.874 \pm 0.511}$ | $\mathbf{19.711 \pm 0.439}$ | $\mathbf{19.592 \pm 0.376}$ | $\mathbf{18.594 \pm 0.439}$ |

Table 11: Results comparing `MCP`, `CCP`, `cluster-CP`, and $k$-`CCP` with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type MAJ and APS scoring function on dataset mini-ImageNet. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-`CCP` the same as or better than that of `CCP` and `cluster-CP` for a fair comparison of prediction set size.

| Measure | Methods | EXP | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | $0.444 \pm 0.007$ | $0.452 \pm 0.007$ | $0.439 \pm 0.005$ | $0.427 \pm 0.006$ | $0.364 \pm 0.01$ |
| | CCP | $0.001 \pm 0.001$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ |
| | cluster-CP | $0.007 \pm 0.003$ | — | — | — | $0.003 \pm 0.002$ |
| | $k$-**CCP** | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ |
| APSS | MCP | $9.57 \pm 0.076$ | $9.687 \pm 0.075$ | $10.437 \pm 0.064$ | $11.404 \pm 0.076$ | $13.998 \pm 0.089$ |
| | CCP | $40.408 \pm 0.378$ | $41.156 \pm 0.405$ | $40.881 \pm 0.398$ | $42.207 \pm 0.356$ | $60.762 \pm 0.531$ |
| | cluster-CP | $28.828 \pm 0.294$ | — | — | — | $44.885 \pm 0.589$ |
| | $k$-**CCP** | $\mathbf{17.281 \pm 0.225}$ | $\mathbf{17.294 \pm 0.202}$ | $\mathbf{17.928 \pm 0.21}$ | $\mathbf{18.63 \pm 0.241}$ | $\mathbf{20.61 \pm 0.222}$ |

Table 12: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type EXP and APS scoring function on dataset Food-101. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

| Measure | Methods | POLY | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | $0.466 \pm 0.009$ | $0.446 \pm 0.011$ | $0.456 \pm 0.006$ | $0.465 \pm 0.007$ | $0.451 \pm 0.008$ |
| | CCP | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.001 \pm 0.001}$ | $0.002 \pm 0.001$ | $\mathbf{0.0 \pm 0.0}$ | $0.001 \pm 0.001$ |
| | cluster-CP | $0.007 \pm 0.002$ | — | — | — | $0.010 \pm 0.003$ |
| | $k$-**CCP** | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.001 \pm 0.001}$ | $0.001 \pm 0.001$ | $\mathbf{0.001 \pm 0.001}$ |
| APSS | MCP | $12.267 \pm 0.079$ | $12.349 \pm 0.085$ | $13.533 \pm 0.09$ | $14.357 \pm 0.08$ | $16.468 \pm 0.095$ |
| | CCP | $45.148 \pm 0.342$ | $45.572 \pm 0.355$ | $46.134 \pm 0.347$ | $47.788 \pm 0.407$ | $65.672 \pm 0.515$ |
| | cluster-CP | $32.873 \pm 0.307$ | — | — | — | $38.326 \pm 0.248$ |
| | $k$-**CCP** | $\mathbf{20.452 \pm 0.209}$ | $\mathbf{20.503 \pm 0.192}$ | $\mathbf{21.606 \pm 0.206}$ | $\mathbf{22.62 \pm 0.187}$ | $\mathbf{24.771 \pm 0.192}$ |

Table 13: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type POLY and APS scoring function on dataset Food-101. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.
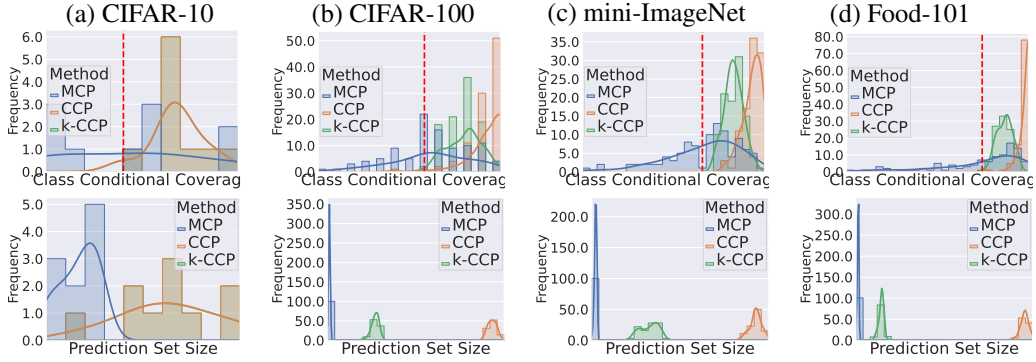
| Measure | Methods | MAJ | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | $\rho = 0.5$ | $\rho = 0.4$ | $\rho = 0.3$ | $\rho = 0.2$ | $\rho = 0.1$ |
| UCR | MCP | $0.462 \pm 0.008$ | $0.469 \pm 0.009$ | $0.47 \pm 0.01$ | $0.459 \pm 0.007$ | $0.467 \pm 0.006$ |
| | CCP | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.001 \pm 0.001}$ | $0.003 \pm 0.001$ | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.0 \pm 0.0}$ |
| | cluster-CP | $0.007 \pm 0.003$ | — | — | — | $0.005 \pm 0.003$ |
| | $k$-**CCP** | $\mathbf{0.0 \pm 0.0}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.002 \pm 0.001}$ | $\mathbf{0.0 \pm 0.0}$ | $0.002 \pm 0.001$ |
| APSS | MCP | $9.964 \pm 0.078$ | $10.742 \pm 0.075$ | $11.57 \pm 0.088$ | $12.654 \pm 0.085$ | $16.256 \pm 0.088$ |
| | CCP | $41.453 \pm 0.335$ | $43.252 \pm 0.385$ | $44.884 \pm 0.348$ | $45.492 \pm 0.288$ | $66.633 \pm 0.622$ |
| | cluster-CP | $33.258 \pm 0.450$ | — | — | — | $46.430 \pm 0.337$ |
| | $k$-**CCP** | $\mathbf{19.398 \pm 0.223}$ | $\mathbf{18.97 \pm 0.217}$ | $\mathbf{20.375 \pm 0.218}$ | $\mathbf{21.172 \pm 0.242}$ | $\mathbf{25.164 \pm 0.197}$ |

Table 14: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-20 model under different imbalance ratio $\rho = 0.5$, $\rho = 0.4$, $\rho = 0.2$, and $\rho = 0.1$ with imbalance type MAJ and APS scoring function on dataset Food-101. The dash symbol (—) means missing results that will be finished in the final version. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.
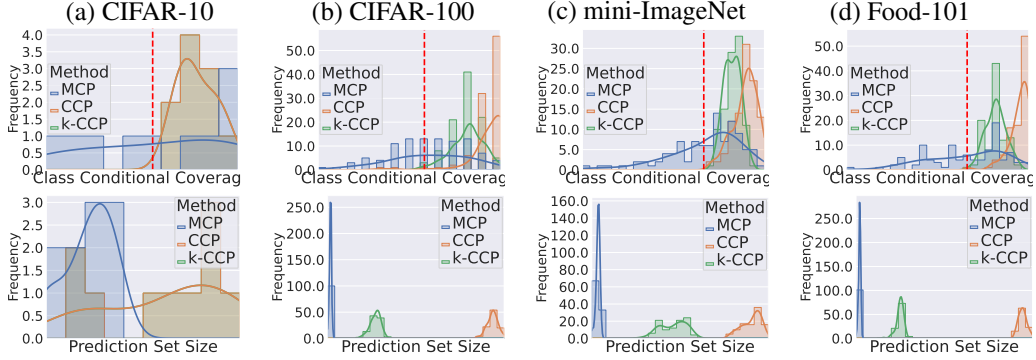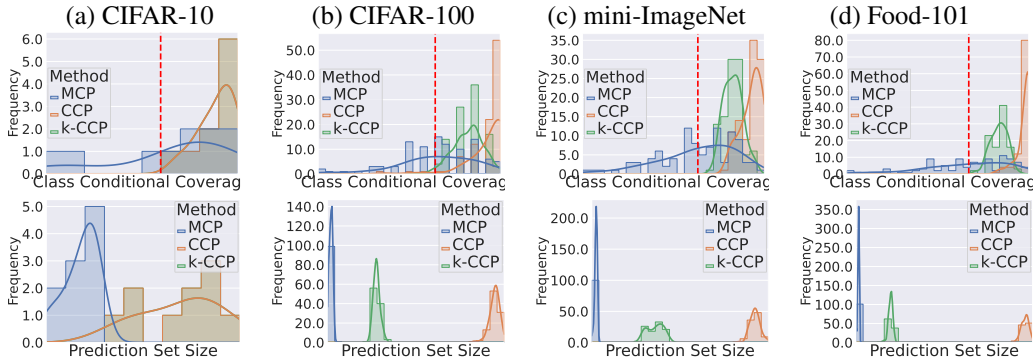
Figure 6: Class-conditional coverage (Top row) and prediction set size (Bottom row) achieved by `MCP`, `CCP`, and $k-$`CCP` methods using ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type EXP for imbalance ratio $\rho = 0.1$. It is clear that $k$-CCP has more densely distributed class-conditional coverage above $0.9$ (the target $1-\alpha$ class-conditional coverage) than CCP with significantly smaller prediction sets on CIFAR-100, mini-ImageNet and Food-101.



Figure 7: Class-conditional coverage (Top row) and prediction set size (Bottom row) achieved by `MCP`, `CCP`, and $k-$`CCP` methods using ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type EXP for imbalance ratio $\rho = 0.5$. It is clear that $k$-CCP has more densely distributed class-conditional coverage above $0.9$ (the target $1-\alpha$ class-conditional coverage) than CCP with significantly smaller prediction sets on CIFAR-100, mini-ImageNet and Food-101.
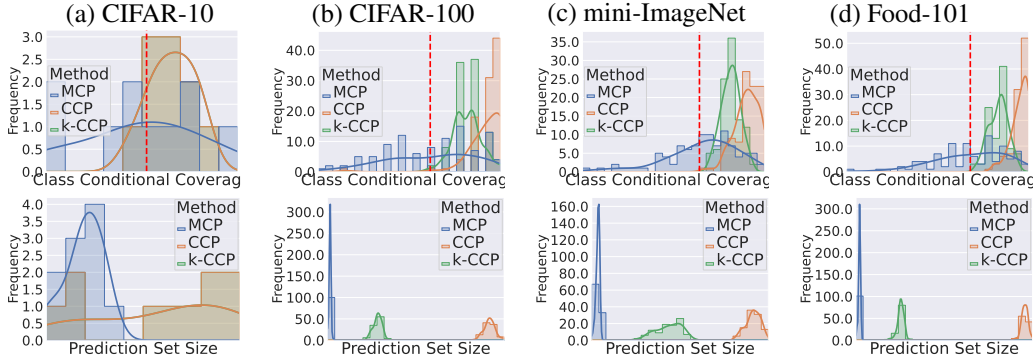


Figure 8: Class-conditional coverage (Top row) and prediction set size (Bottom row) achieved by `MCP`, `CCP`, and $k-$`CCP` methods using ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type POLY for imbalance ratio $\rho = 0.1$. It is clear that $k$-CCP has more densely distributed class-conditional coverage above $0.9$ (the target $1 - \alpha$ class-conditional coverage) than CCP with significantly smaller prediction sets on CIFAR-100, mini-ImageNet and Food-101.
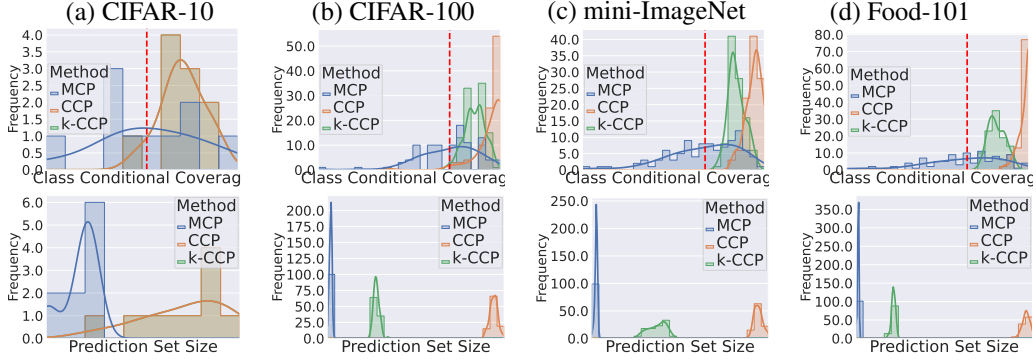
Figure 9: Class-conditional coverage (Top row) and prediction set size (Bottom row) achieved by `MCP`, `CCP`, and $k-$`CCP` methods using ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type POLY for imbalance ratio $\rho = 0.5$. It is clear that $k$-CCP has more densely distributed class-conditional coverage above 0.9 (the target $1 - \alpha$ class-conditional coverage) than CCP with significantly smaller prediction sets on CIFAR-100, mini-ImageNet and Food-101.



Figure 10: Class-conditional coverage (Top row) and prediction set size (Bottom row) achieved by `MCP`, `CCP`, and $k-$`CCP` methods using ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type MAJ for imbalance ratio $\rho = 0.1$. It is clear that $k$-CCP has more densely distributed class-conditional coverage above 0.9 (the target $1 - \alpha$ class-conditional coverage) than CCP with significantly smaller prediction sets on CIFAR-100, mini-ImageNet and Food-101.
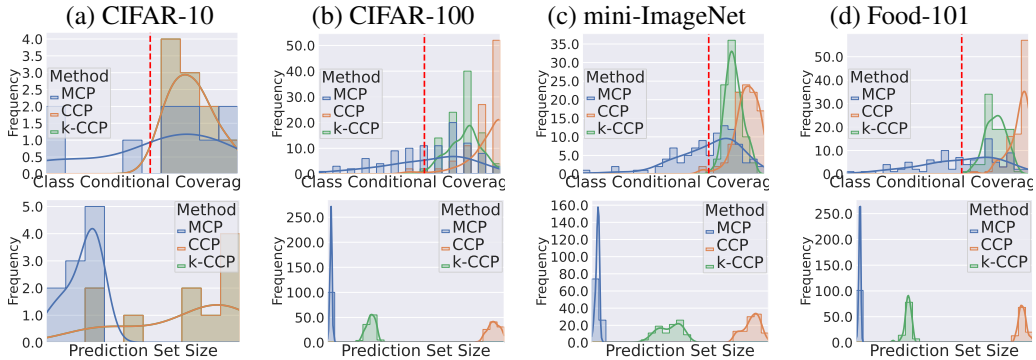


Figure 11: Class-conditional coverage (Top row) and prediction set size (Bottom row) achieved by `MCP`, `CCP`, and $k-$`CCP` methods using ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type MAJ for imbalance ratio $\rho = 0.5$. It is clear that $k$-CCP has more densely distributed class-conditional coverage above 0.9 (the target $1 - \alpha$ class-conditional coverage) than CCP with significantly smaller prediction sets on CIFAR-100, mini-ImageNet and Food-101.
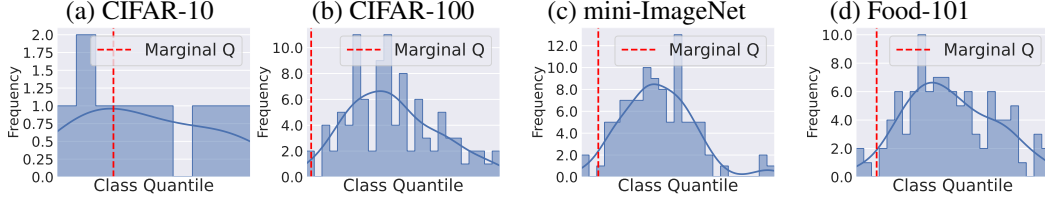
Figure 12: Distribution of class-wise quantiles with ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type EXP and imbalance ratio $\rho = 0.5$. This result verifies that the deviation of class-wise quantiles from the marginal quantile can easily happen, i.e., the assumption of class-conditional under-coverage for MCP in Proposition 1.
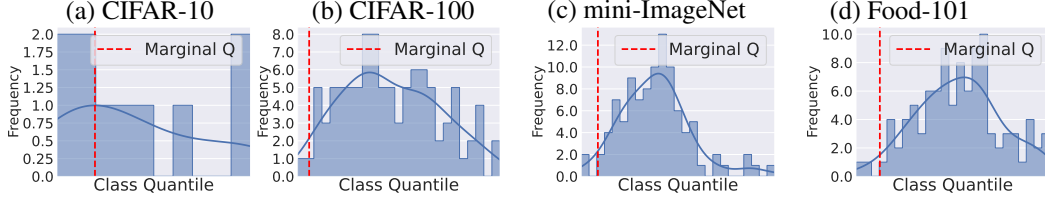


Figure 13: Distribution of class-wise quantiles with ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type POLY and imbalance ratio $\rho = 0.5$. This result verifies that the deviation of class-wise quantiles from the marginal quantile can easily happen, i.e., the assumption of class-conditional under-coverage for MCP in Proposition 1.
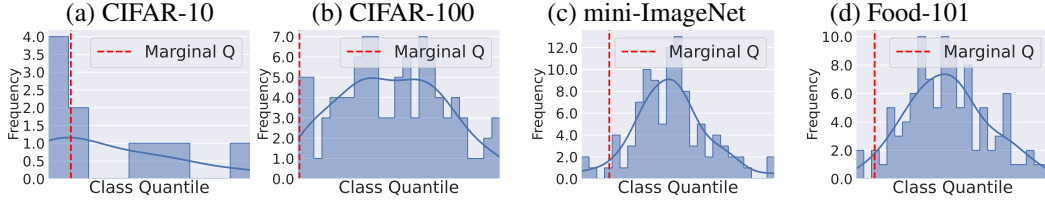


Figure 14: Distribution of class-wise quantiles with ResNet20 model on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type MAJ and imbalance ratio $\rho = 0.5$. This result verifies that the deviation of class-wise quantiles from the marginal quantile can easily happen, i.e., the assumption of class-conditional under-coverage for MCP in Proposition 1.
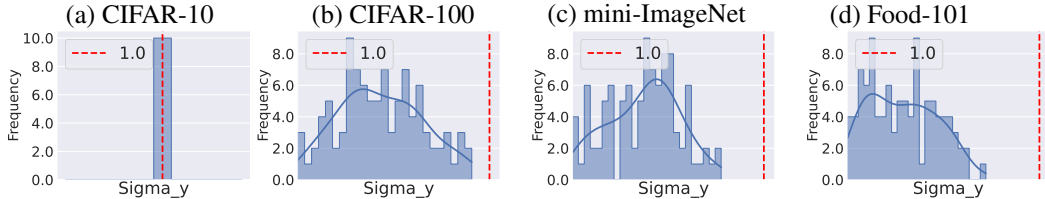


Figure 15: Verification of condition numbers $\{\sigma_y\}_{y=1}^{C}$ in Theorem 2 with $\rho = 0.1$ EXP. Vertical dashed lines represent the value 1, and we observe that all the condition numbers are smaller than 1. This verifies the validity of the condition for Theorem 2, and thus confirms that $k$-CCP produces smaller prediction sets than CCP by the optimized trade-off between calibration on non-conformity scores and calibrated label ranks.



Figure 16: Verification of condition numbers $\{\sigma_y\}_{y=1}^{C}$ in Theorem 2 with $\rho = 0.5$ EXP. Vertical dashed lines represent the value 1, and we observe that all the condition numbers are smaller than 1. This verifies the validity of the condition for Theorem 2, and thus confirms that $k$-CCP produces smaller prediction sets than CCP by the optimized trade-off between calibration on non-conformity scores and calibrated label ranks.

Figure 17: Verification of condition numbers $\{\sigma_y\}_{y=1}^C$ in Theorem 2 with $\rho = 0.1$ POLY. Vertical dashed lines represent the value 1, and we observe that all the condition numbers are smaller than 1. This verifies the validity of the condition for Theorem 2, and thus confirms that $k$-CCP produces smaller prediction sets than CCP by the optimized trade-off between calibration on non-conformity scores and calibrated label ranks.
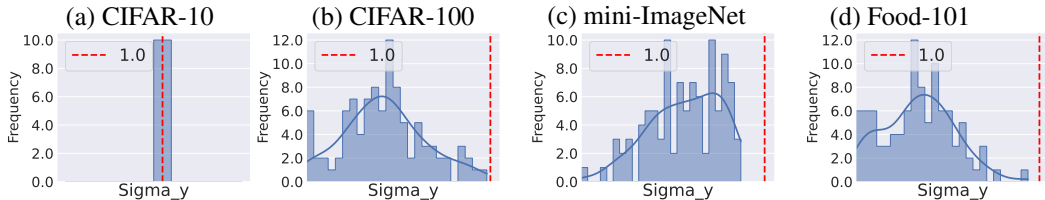


Figure 18: Verification of condition numbers $\{\sigma_y\}_{y=1}^C$ in Theorem 2 with $\rho = 0.5$ POLY. Vertical dashed lines represent the value 1, and we observe that all the condition numbers are smaller than 1. This verifies the validity of the condition for Theorem 2, and thus confirms that $k$-CCP produces smaller prediction sets than CCP by the optimized trade-off between calibration on non-conformity scores and calibrated label ranks.
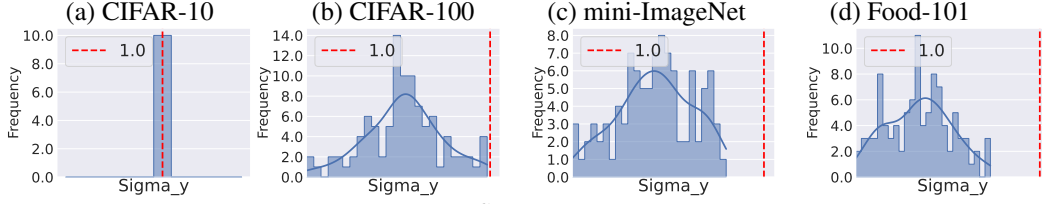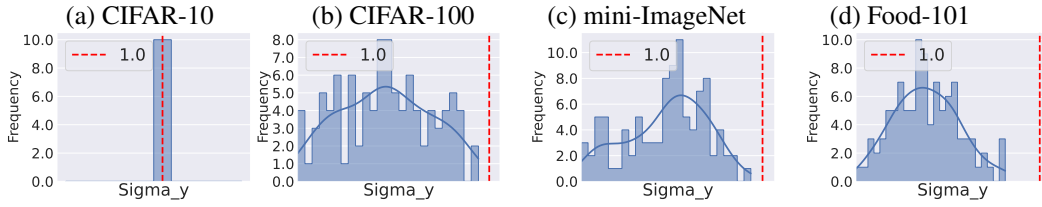


Figure 19: Verification of condition numbers $\{\sigma_y\}_{y=1}^C$ in Theorem 2 with $\rho = 0.1$ MAJ. Vertical dashed lines represent the value 1, and we observe that all the condition numbers are smaller than 1. This verifies the validity of the condition for Theorem 2, and thus confirms that $k$-CCP produces smaller prediction sets than CCP by the optimized trade-off between calibration on non-conformity scores and calibrated label ranks.



Figure 20: Verification of condition numbers $\{\sigma_y\}_{y=1}^C$ in Theorem 2 with $\rho = 0.5$ MAJ. Vertical dashed lines represent the value 1, and we observe that all the condition numbers are smaller than 1. This verifies the validity of the condition for Theorem 2, and thus confirms that $k$-CCP produces smaller prediction sets than CCP by the optimized trade-off between calibration on non-conformity scores and calibrated label ranks.
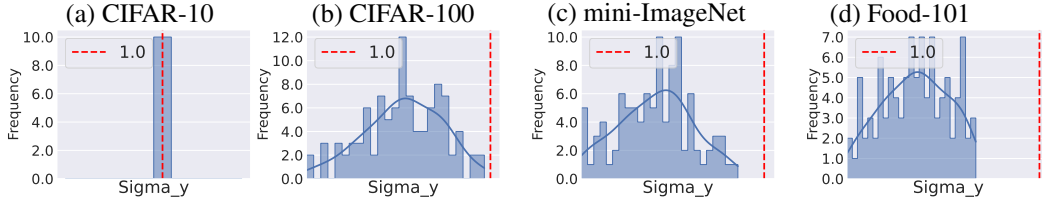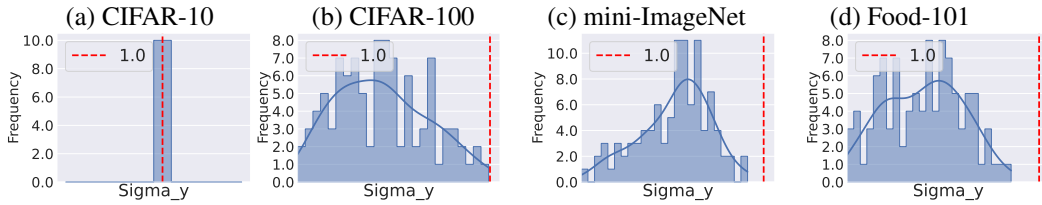
## D Experiments with cluster-CP Using APS Score Function

With the same model, evaluation metrics, and APS score function Romano *et al.* (2020), we add the comparison experiments with cluster-CP Ding *et al.* (2023) [3] on four datasets and summarize the results in Table 15. We highlight that we also select the $g$ from the interval $[0.1, 1]$ with increments of $0.05$ to find the minimal $g$ that $k$-CCP and cluster-CP achieves the target class-conditional coverage.

Based on the results in Table 15, we make the following observations: (i) cluster-CP and $k$-CCP can guarantee the class-conditional coverage; and (ii) $k$-CCP significantly outperforms cluster-CP on three datasets by producing smaller prediction sets.

| Measure | Methods | EXP | | POLY | | MAJ | |
|---|---|---|---|---|---|---|---|
| | | $\rho = 0.5$ | $\rho = 0.1$ | $\rho = 0.5$ | $\rho = 0.1$ | $\rho = 0.5$ | $\rho = 0.1$ |
| | | CIFAR-10 | | | | | |
| UCR | MCP | $0.460 \pm 0.025$ | $0.490 \pm 0.026$ | $0.510 \pm 0.033$ | $0.320 \pm 0.030$ | $0.380 \pm 0.020$ | $0.500 \pm 0.014$ |
| | CCP | $\mathbf{0.110 \pm 0.030}$ | $\mathbf{0.040 \pm 0.020}$ | $0.140 \pm 0.032$ | $\mathbf{0.010 \pm 0.001}$ | $\mathbf{0.120 \pm 0.024}$ | $0.010 \pm 0.001$ |
| | cluster-CP | $0.160 \pm 0.025$ | $0.080 \pm 0.012$ | $\mathbf{0.140 \pm 0.021}$ | $\mathbf{0.010 \pm 0.001}$ | $0.150 \pm 0.025$ | $\mathbf{0.09 \pm 0.013}$ |
| | $k$-CCP | $\mathbf{0.110 \pm 0.030}$ | $\mathbf{0.040 \pm 0.020}$ | $0.140 \pm 0.032$ | $\mathbf{0.010 \pm 0.001}$ | $\mathbf{0.120 \pm 0.024}$ | $0.010 \pm 0.001$ |
| APSS | MCP | $1.132 \pm 0.033$ | $1.406 \pm 0.045$ | $1.117 \pm 0.028$ | $1.214 \pm 0.038$ | $1.196 \pm 0.032$ | $2.039 \pm 0.046$ |
| | CCP | $1.481 \pm 0.082$ | $\mathbf{2.032 \pm 0.096}$ | $\mathbf{1.487 \pm 0.090}$ | $\mathbf{1.945 \pm 0.087}$ | $\mathbf{1.765 \pm 0.093}$ | $\mathbf{2.964 \pm 0.123}$ |
| | cluster-CP | $\mathbf{1.445 \pm 0.017}$ | $2.323 \pm 0.015$ | $1.612 \pm 0.013$ | $2.102 \pm 0.015$ | $1.787 \pm 0.019$ | $2.969 \pm 0.025$ |
| | $k$-CCP | $1.481 \pm 0.082$ | $\mathbf{2.032 \pm 0.096}$ | $\mathbf{1.487 \pm 0.090}$ | $\mathbf{1.945 \pm 0.087}$ | $\mathbf{1.765 \pm 0.093}$ | $\mathbf{2.964 \pm 0.123}$ |
| | | CIFAR-100 | | | | | |
| UCR | MCP | $0.386 \pm 0.009$ | $0.384 \pm 0.018$ | $0.395 \pm 0.010$ | $0.411 \pm 0.010$ | $0.352 \pm 0.010$ | $0.401 \pm 0.008$ |
| | CCP | $0.009 \pm 0.003$ | $0.011 \pm 0.002$ | $0.011 \pm 0.003$ | $0.015 \pm 0.003$ | $0.017 \pm 0.003$ | $0.008 \pm 0.004$ |
| | cluster-CP | $0.004 \pm 0.002$ | $0.004 \pm 0.002$ | $0.001 \pm 0.001$ | $0.004 \pm 0.002$ | $0.005 \pm 0.002$ | $0.019 \pm 0.005$ |
| | $k$-CCP | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.002 \pm 0.001}$ | $\mathbf{0.001 \pm 0.001}$ |
| APSS | MCP | $10.303 \pm 0.111$ | $14.544 \pm 0.119$ | $15.729 \pm 0.126$ | $25.888 \pm 0.197$ | $11.680 \pm 0.117$ | $23.796 \pm 0.159$ |
| | CCP | $44.194 \pm 0.514$ | $50.963 \pm 0.481$ | $49.895 \pm 0.489$ | $64.366 \pm 0.389$ | $48.323 \pm 0.548$ | $64.640 \pm 0.621$ |
| | cluster-CP | $30.922 \pm 0.454$ | $43.883 \pm 1.070$ | $56.696 \pm 0.393$ | $63.208 \pm 0.364$ | $33.623 \pm 0.395$ | $50.382 \pm 0.711$ |
| | $k$-CCP | $\mathbf{20.355 \pm 0.357}$ | $\mathbf{25.185 \pm 0.278}$ | $\mathbf{25.843 \pm 0.300}$ | $\mathbf{37.034 \pm 0.244}$ | $\mathbf{21.196 \pm 0.320}$ | $\mathbf{35.630 \pm 0.232}$ |
| | | mini-ImageNet | | | | | |
| UCR | MCP | $0.408 \pm 0.008$ | $0.414 \pm 0.012$ | $0.412 \pm 0.013$ | $0.410 \pm 0.0018$ | $0.408 \pm 0.010$ | $0.411 \pm 0.007$ |
| | CCP | $0.007 \pm 0.003$ | $0.001 \pm 0.001$ | $0.005 \pm 0.002$ | $0.001 \pm 0.001$ | $0.011 \pm 0.004$ | $0.003 \pm 0.001$ |
| | cluster-CP | $0.009 \pm 0.003$ | $0.002 \pm 0.001$ | $0.028 \pm 0.005$ | $0.015 \pm 0.003$ | $0.015 \pm 0.005$ | $0.011 \pm 0.003$ |
| | $k$-CCP | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ |
| APSS | MCP | $9.705 \pm 0.101$ | $8.930 \pm 0.093$ | $9.810 \pm 0.101$ | $9.665 \pm 0.101$ | $9.840 \pm 0.091$ | $9.123 \pm 0.086$ |
| | CCP | $26.666 \pm 0.415$ | $34.867 \pm 0.445$ | $26.620 \pm 0.369$ | $29.852 \pm 0.360$ | $27.306 \pm 0.377$ | $29.200 \pm 0.379$ |
| | cluster-CP | $27.786 \pm 0.307$ | $33.114 \pm 0.418$ | $21.273 \pm 0.369$ | $25.550 \pm 0.279$ | $25.288 \pm 0.226$ | $25.229 \pm 0.352$ |
| | $k$-CCP | $\mathbf{18.129 \pm 0.453}$ | $\mathbf{17.769 \pm 0.463}$ | $\mathbf{17.784 \pm 0.438}$ | $\mathbf{19.153 \pm 0.412}$ | $\mathbf{18.110 \pm 0.414}$ | $\mathbf{18.594 \pm 0.439}$ |
| | | Food-101 | | | | | |
| UCR | MCP | $0.444 \pm 0.007$ | $0.364 \pm 0.010$ | $0.466 \pm 0.009$ | $0.451 \pm 0.008$ | $0.462 \pm 0.008$ | $0.467 \pm 0.006$ |
| | CCP | $0.001 \pm 0.001$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ |
| | cluster-CP | $0.007 \pm 0.003$ | $0.003 \pm 0.002$ | $0.007 \pm 0.002$ | $0.010 \pm 0.003$ | $0.007 \pm 0.003$ | $0.005 \pm 0.003$ |
| | $k$-CCP | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ |
| APSS | MCP | $9.570 \pm 0.076$ | $13.998 \pm 0.089$ | $12.267 \pm 0.079$ | $16.468 \pm 0.095$ | $9.964 \pm 0.078$ | $23.796 \pm 0.159$ |
| | CCP | $40.408 \pm 0.378$ | $60.762 \pm 0.531$ | $45.148 \pm 0.342$ | $65.6723 \pm 0.515$ | $41.453 \pm 0.335$ | $66.633 \pm 0.622$ |
| | cluster-CP | $28.828 \pm 0.294$ | $44.885 \pm 0.589$ | $32.873 \pm 0.307$ | $38.326 \pm 0.248$ | $33.258 \pm 0.450$ | $46.430 \pm 0.337$ |
| | $k$-CCP | $\mathbf{17.281 \pm 0.225}$ | $\mathbf{20.610 \pm 0.222}$ | $\mathbf{20.452 \pm 0.209}$ | $\mathbf{24.771 \pm 0.192}$ | $\mathbf{19.398 \pm 0.223}$ | $\mathbf{26.584 \pm 0.191}$ |

Table 15: Results comparing $k$-CCP and cluster-CP with ResNet-20 model and APS score function under different imbalance ratios $\rho = 0.5$ and $\rho = 0.1$. We set UCR of $k$-CCP the same as or better than that of cluster-CP for a fair comparison of prediction set size. The APSS results show that $k$-CCP significantly outperforms cluster-CP in terms of the average prediction set size over all settings on CIFAR-100, mini-ImageNet, and Food-101.

## E Comparison Experiments Using RAPS Score Function

With the same model, evaluation metrics and RAPS score function Angelopoulos *et al.* (2020), we add the comparison experiments with MCP, CCP, and cluster-CP on four datasets with different imbalanced type and imbalance ratio $\rho = 0.5$ and $\rho = 0.1$. The regularization parameter for RAPS score function is from the set $k_{reg} \in \{3, 5, 7\}$ and $\lambda \in \{0.001, 0.01, 0.1\}$. We select the combination of $k_{reg}$ and $\lambda$ for each experiments with same imbalanced type and imbalanced ratio on same

---

[3]https://github.com/tiffanyding/class-conditional-conformal/tree/main

dataset, where the most of $APSS$ values of all methods are minimum. The overall performance is summarized in Table 16. We highlight that we also select the $g$ from the interval $[0.1, 1]$ with increments of $0.05$ to find the minimal $g$ that CCP, cluster-CP, and $k$-CCP achieves the target class conditional coverage.

Based on the results in Table 16, we make the following observations: (i) CCP, cluster-CP, and $k$-CCP can guarantee the class-conditional coverage; and (ii) $k$-CCP significantly outperforms CCP and cluster-CP on three datasets by producing smaller prediction sets.

| Measure | Methods | EXP | | POLY | | MAJ | |
|---|---|---|---|---|---|---|---|
| | | $\rho = 0.5$ | $\rho = 0.1$ | $\rho = 0.5$ | $\rho = 0.1$ | $\rho = 0.5$ | $\rho = 0.1$ |
| | | CIFAR-10 | | | | | |
| UCR | MCP | $0.460 \pm 0.021$ | $0.490 \pm 0.026$ | $0.500 \pm 0.031$ | $0.290 \pm 0.030$ | $0.380 \pm 0.019$ | $0.500 \pm 0.014$ |
| | CCP | $\mathbf{0.010 \pm 0.020}$ | $\mathbf{0.020 \pm 0.013}$ | $\mathbf{0.080 \pm 0.030}$ | $0.050 \pm 0.021$ | $\mathbf{0.090 \pm 0.022}$ | $\mathbf{0.040 \pm 0.015}$ |
| | cluster-CP | $0.160 \pm 0.025$ | $0.080 \pm 0.012$ | $0.140 \pm 0.021$ | $\mathbf{0.040 \pm 0.015}$ | $0.150 \pm 0.025$ | $0.120 \pm 0.013$ |
| | $k$-CCP | $\mathbf{0.010 \pm 0.020}$ | $\mathbf{0.020 \pm 0.013}$ | $\mathbf{0.080 \pm 0.030}$ | $0.050 \pm 0.021$ | $\mathbf{0.090 \pm 0.022}$ | $\mathbf{0.040 \pm 0.015}$ |
| APSS | MCP | $1.143 \pm 0.004$ | $1.419 \pm 0.013$ | $1.118 \pm 0.004$ | $1.233 \pm 0.006$ | $1.196 \pm 0.032$ | $2.043 \pm 0.016$ |
| | CCP | $1.502 \pm 0.007$ | $\mathbf{2.049 \pm 0.013}$ | $\mathbf{1.558 \pm 0.010}$ | $\mathbf{1.776 \pm 0.012}$ | $\mathbf{1.786 \pm 0.020}$ | $\mathbf{2.628 \pm 0.012}$ |
| | cluster-CP | $\mathbf{1.493 \pm 0.017}$ | $2.323 \pm 0.015$ | $1.612 \pm 0.013$ | $1.981 \pm 0.013$ | $1.787 \pm 0.019$ | $2.968 \pm 0.024$ |
| | $k$-CCP | $1.502 \pm 0.007$ | $\mathbf{2.049 \pm 0.013}$ | $\mathbf{1.558 \pm 0.010}$ | $\mathbf{1.776 \pm 0.012}$ | $\mathbf{1.786 \pm 0.020}$ | $\mathbf{2.628 \pm 0.012}$ |
| | | CIFAR-100 | | | | | |
| UCR | MCP | $0.388 \pm 0.007$ | $0.384 \pm 0.018$ | $0.394 \pm 0.010$ | $0.402 \pm 0.010$ | $0.352 \pm 0.010$ | $0.401 \pm 0.007$ |
| | CCP | $0.008 \pm 0.002$ | $0.011 \pm 0.002$ | $0.011 \pm 0.003$ | $0.015 \pm 0.003$ | $0.015 \pm 0.003$ | $0.008 \pm 0.004$ |
| | cluster-CP | $0.004 \pm 0.002$ | $0.004 \pm 0.002$ | $0.001 \pm 0.001$ | $0.003 \pm 0.002$ | $0.005 \pm 0.002$ | $0.018 \pm 0.006$ |
| | $k$-CCP | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.000 \pm 0.000}$ |
| APSS | MCP | $10.300 \pm 0.080$ | $14.554 \pm 0.107$ | $15.755 \pm 0.103$ | $25.850 \pm 0.150$ | $11.684 \pm 0.091$ | $23.708 \pm 0.137$ |
| | CCP | $44.243 \pm 0.340$ | $50.969 \pm 0.345$ | $49.877 \pm 0.354$ | $64.247 \pm 0.234$ | $48.337 \pm 0.355$ | $64.580 \pm 0.536$ |
| | cluster-CP | $30.971 \pm 0.454$ | $43.883 \pm 1.073$ | $56.656 \pm 0.354$ | $63.113 \pm 0.397$ | $33.656 \pm 0.388$ | $50.365 \pm 0.701$ |
| | $k$-CCP | $\mathbf{20.355 \pm 0.005}$ | $\mathbf{25.185 \pm 0.011}$ | $\mathbf{25.843 \pm 0.006}$ | $\mathbf{37.035 \pm 0.005}$ | $\mathbf{21.197 \pm 0.005}$ | $\mathbf{35.631 \pm 0.007}$ |
| | | mini-ImageNet | | | | | |
| UCR | MCP | $0.408 \pm 0.008$ | $0.410 \pm 0.011$ | $0.412 \pm 0.014$ | $0.410 \pm 0.008$ | $0.403 \pm 0.010$ | $0.412 \pm 0.007$ |
| | CCP | $\mathbf{0.007 \pm 0.003}$ | $0.003 \pm 0.001$ | $0.018 \pm 0.002$ | $0.004 \pm 0.002$ | $\mathbf{0.000 \pm 0.000}$ | $0.005 \pm 0.002$ |
| | cluster-CP | $0.008 \pm 0.004$ | $0.003 \pm 0.001$ | $0.031 \pm 0.005$ | $0.002 \pm 0.001$ | $0.004 \pm 0.002$ | $0.011 \pm 0.004$ |
| | $k$-CCP | $0.009 \pm 0.003$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.004 \pm 0.002}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ |
| APSS | MCP | $9.703 \pm 0.076$ | $9.003 \pm 0.067$ | $9.806 \pm 0.079$ | $9.714 \pm 0.075$ | $9.865 \pm 0.060$ | $9.146 \pm 0.063$ |
| | CCP | $26.689 \pm 0.177$ | $29.750 \pm 0.219$ | $21.352 \pm 0.196$ | $26.266 \pm 0.218$ | $36.535 \pm 0.196$ | $25.641 \pm 0.217$ |
| | cluster-CP | $27.466 \pm 0.268$ | $32.991 \pm 0.434$ | $21.212 \pm 0.298$ | $36.061 \pm 0.475$ | $32.085 \pm 0.424$ | $25.269 \pm 0.375$ |
| | $k$-CCP | $\mathbf{15.101 \pm 0.003}$ | $\mathbf{18.418 \pm 0.003}$ | $\mathbf{15.331 \pm 0.003}$ | $\mathbf{17.465 \pm 0.003}$ | $\mathbf{17.388 \pm 0.003}$ | $\mathbf{17.167 \pm 0.004}$ |
| | | Food-101 | | | | | |
| UCR | MCP | $0.445 \pm 0.007$ | $0.363 \pm 0.010$ | $0.466 \pm 0.008$ | $0.450 \pm 0.007$ | $0.457 \pm 0.008$ | $0.465 \pm 0.006$ |
| | CCP | $0.002 \pm 0.001$ | $\mathbf{0.000 \pm 0.000}$ | $0.002 \pm 0.001$ | $\mathbf{0.001 \pm 0.001}$ | $0.001 \pm 0.001$ | $0.008 \pm 0.002$ |
| | cluster-CP | $0.003 \pm 0.002$ | $0.003 \pm 0.002$ | $0.007 \pm 0.002$ | $0.005 \pm 0.002$ | $0.007 \pm 0.002$ | $0.006 \pm 0.003$ |
| | $k$-CCP | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.001 \pm 0.001}$ | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.002 \pm 0.001}$ |
| APSS | MCP | $9.580 \pm 0.037$ | $14.039 \pm 0.055$ | $12.327 \pm 0.046$ | $16.541 \pm 0.060$ | $10.040 \pm 0.051$ | $16.293 \pm 0.047$ |
| | CCP | $40.411 \pm 0.285$ | $60.790 \pm 0.395$ | $36.550 \pm 0.141$ | $41.755 \pm 0.153$ | $32.957 \pm 0.224$ | $36.797 \pm 0.139$ |
| | cluster-CP | $28.919 \pm 0.287$ | $44.583 \pm 0.667$ | $32.928 \pm 0.358$ | $41.785 \pm 0.220$ | $32.983 \pm 0.518$ | $46.078 \pm 0.312$ |
| | $k$-CCP | $\mathbf{17.282 \pm 0.004}$ | $\mathbf{20.610 \pm 0.006}$ | $\mathbf{20.452 \pm 0.002}$ | $\mathbf{24.771 \pm 0.004}$ | $\mathbf{19.398 \pm 0.006}$ | $\mathbf{25.163 \pm 0.002}$ |

Table 16: Results comparing MCP, CCP, cluster-CP, and $k$-CCP under different imbalance ratios $\rho = 0.5$ and $\rho = 0.1$ using the RAPS score function. The regularization parameter for RAPS score function is selected from the set $[3, 5, 7]$ and $[0.001, 0.01, 0.1]$. We select the best results as each element in the table. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size. The APSS results show that $k$-CCP significantly outperforms CCP and cluster-CP in terms of average prediction set size over all settings on CIFAR-100, mini-ImageNet, and Food-101.

## F    COMPARISON EXPERIMENTS WITH CONFORMAL TRAINING MODEL

To verify the suitability of $k$-CCP, we add the experiments by performing calibration using MCP, CCP, cluster-CP, and $k$-CCP on the conformal training model Stutz *et al.* (2021) [4] in CIFAR-100 dataset. The imbalance type is EXP and the imbalance ratio is $\rho = 0.1$. The regularization parameters of conformal training model are $r \in \{0.01, 0.05, 0.1, 0.5, 1.0\}$. We have selected the best results for each method and summarized results in Table 17. It is clear that $k$-CCP outperforms CCP and cluster-CP, and improvement is not caused by RAPS. Due to the limited time for rebuttal and

---

[4]https://github.com/google-deepmind/conformal_training/tree/main

large computational cost of training conformal classifiers, we plan to add the corresponding results for mini-ImageNet and Food-101 in the final paper.

| Measure | Methods | CIFAR-100 | |
| | | APS | RAPS |
| --- | --- | --- | --- |
| UCR | MCP | $0.355 \pm 0.014$ | $0.377 \pm 0.011$ |
| | CCP | $0.041 \pm 0.004$ | $0.046 \pm 0.004$ |
| | cluster-CP | $0.027 \pm 0.008$ | $0.039 \pm 0.007$ |
| | $k$-**CCP** | $\mathbf{0.000 \pm 0.000}$ | $\mathbf{0.000 \pm 0.000}$ |
| APSS | MCP | $19.786 \pm 0.141$ | $15.175 \pm 0.127$ |
| | CCP | $38.505 \pm 0.301$ | $37.439 \pm 0.593$ |
| | cluster-CP | $36.055 \pm 0.636$ | $36.193 \pm 0.687$ |
| | $k$-**CCP** | $\mathbf{31.281 \pm 0.007}$ | $\mathbf{31.283 \pm 0.006}$ |

Table 17: Results comparing MCP, CCP, cluster-CP, and $k$-CCP with ResNet-18 model by conformal training loss proposed by Stutz *et al.* (2021) and imbalance ratio $\rho = 0.1$ EXP on dataset CIFAR-100. We set UCR of $k$-CCP the same as or better than that of CCP and cluster-CP for a fair comparison of prediction set size.

# G ILLUSTRATION OF GROUP-WISE AVERAGE PREDICTION SIZE

To visualize the average prediction prediction set size gap between group class, we add the experiments with same model, imbalance type and imbalance ratio on four datasets. We select the top $1/4$ classes of largest number of data to the majority group. Similarly, we select the bottom $1/4$ classes of smallest number of data to the minority group, and the remaining $1/2$ classes to the medium group.
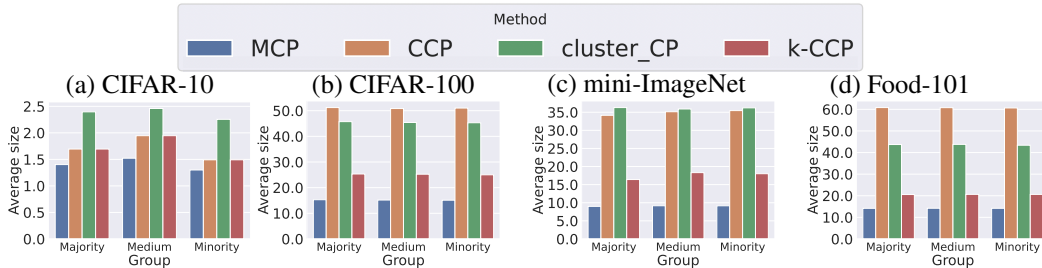


Figure 21: Comparison of average prediction set size with $\rho = 0.1$ EXP and APS score function. These results show that there exists small gap between the average predictions set size on majority, medium, and minority groups.
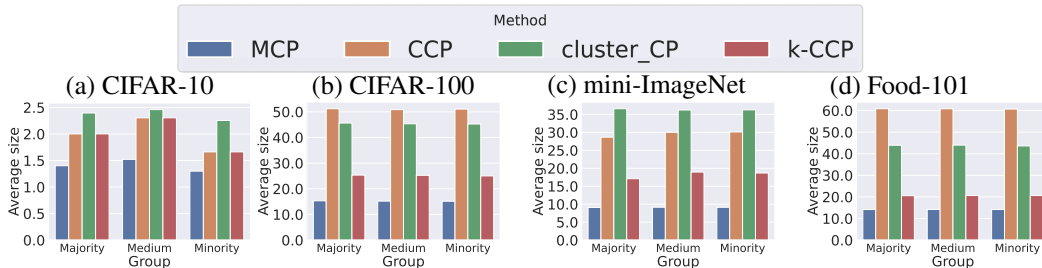


Figure 22: Comparison of average prediction set size with $\rho = 0.1$ EXP and RAPS score function. These results show that there exists small gap between the average predictions set size on majority, medium, and minority groups.

# H ABLATION STUDY FOR HYPER-PARAMETER $g$ WITHIN $k$-CCP

We add the ablation study to verify how the hyper-parameter $g$ affects the performance of CCP, Cluster-CP, and $k$-CCP. We use the under coverage ratio (UCR) and average prediction
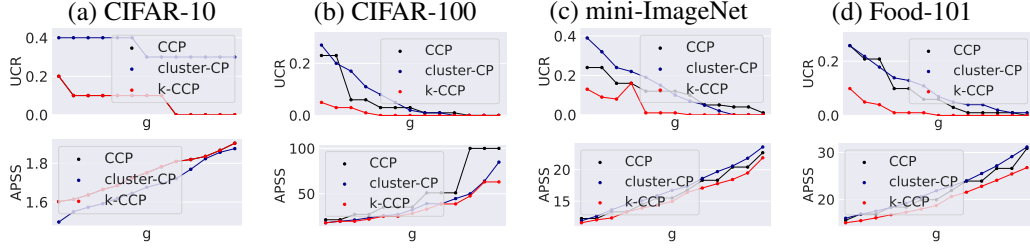
Figure 23: Under coverage ratio (Top row) and average prediction set size (Bottom row) achieved by CCP, cluster-CP, and $k$-CCP methods using ResNet20 model and RAPS score function on CIFAR-10, CIFAR-100, mini-ImageNet, and Food-101 datasets with imbalance type EXP for imbalance ratio $\rho = 0.1$. $k$-CCP degenerates to CCP in CIFAR-10, so there are only two lines (the black line overlaps with the red line).

set size (APSS) as metrics, which are introduced in Section 5.1. We set the range of $g$ from $\{0.1, 0.15, \cdots, 0.7\}$ on four datasets. Based on the results in Figure 2, 5, and 23, the UCR and APSS of $k$-CCP are much smaller than CCP and Cluster-CP with the same $g$ value.

# I    VERIFICATION OF $\sigma_y$ ON CALIBRATION SET

To investigate the validity of $\sigma_y$ on calibration datasets, we add experiments with imbalance ratio $\rho = 0.1$ and imbalance type EXP on four datasets.

Figure 24 verifies the validity of Theorem 2 and confirms that optimized trade-off between the coverage with inflated quantile and the constraint with calibrated rank leads to smaller prediction sets. Experiments even show a stronger condition ($\sigma_y \leq 1$ for all $y$) than the weighted aggregation condition in (12). It also confirms that the condition number $\{\sigma_y\}_{y=1}^{C}$ could be evaluated on calibration datasets without testing datasets and thus decreases the computation cost. We notice that $k$-CCP degenerates to CCP on CIFAR-10, so $\sigma_y = 1$ for all $y$ and there is no trade-off. On the other three datasets, we observe significant conditions for the optimized trade-off in $k$-CCP.
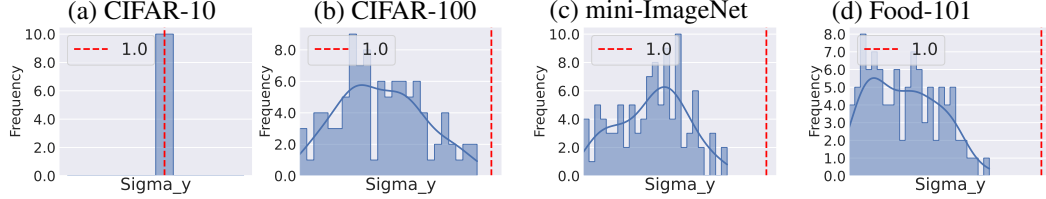


Figure 24: Verification of condition numbers $\{\sigma_y\}_{y=1}^{C}$ in Theorem 2 with $\rho = 0.1$ EXP on calibration datasets. Vertical dashed lines represent the value 1, and we observe that all the condition numbers are smaller than 1. This verifies the validity of the condition for Theorem 2, and thus confirms that the conditional number $\{\sigma_y\}_{y=1}^{C}$ could be evaluated on calibration datasets without testing datasets and thus decrease the computation cost.