

---

# Improved Online Conformal Prediction via Strongly Adaptive Online Learning

---

Aadyot Bhatnagar<sup>1</sup> Huan Wang<sup>1</sup> Caiming Xiong<sup>1</sup> Yu Bai<sup>1</sup>

## Abstract

We study the problem of uncertainty quantification via prediction sets, in an online setting where the data distribution may vary arbitrarily over time. Recent work develops *online conformal prediction* techniques that leverage regret minimization algorithms from the online learning literature to learn prediction sets with approximately valid coverage and small regret. However, standard regret minimization could be insufficient for handling changing environments, where performance guarantees may be desired not only over the full time horizon but also in all (sub-)intervals of time. We develop new online conformal prediction methods that minimize the *strongly adaptive regret*, which measures the worst-case regret over all intervals of a fixed length. We prove that our methods achieve near-optimal strongly adaptive regret for all interval lengths simultaneously, and approximately valid coverage. Experiments show that our methods consistently obtain better coverage and smaller prediction sets than existing methods on real-world tasks, such as time series forecasting and image classification under distribution shift.

## 1. Introduction

Modern machine learning models make highly accurate predictions in many settings. In high stakes decision-making tasks, it is just as important to estimate the model’s uncertainty by quantifying how much the true label may deviate from the model’s prediction. A common approach for uncertainty quantification is to learn *prediction sets* that associate each input with a set of candidate labels, such as prediction intervals for regression, and label sets for classification. The most important requirement for learned prediction sets is to

achieve valid *coverage*, i.e. they should cover the true label with at least  $1 - \alpha$  (such as 90%) probability.

Conformal prediction (Vovk et al., 2005) is a powerful framework for augmenting any *base predictor* (such as a pretrained model) into prediction sets with valid coverage guarantees (Angelopoulos & Bates, 2021). These guarantees require almost no assumptions on the data distribution, except *exchangeability* (i.i.d. data is a sufficient condition). However, exchangeability fails to hold in many real-world settings such as time series data (Chernozhukov et al., 2018) or data corruption (Hendrycks et al., 2018), where the data may exhibit *distribution shift*. Various approaches have been proposed to handle such distribution shift, such as reweighting (Tibshirani et al., 2019; Barber et al., 2022) or distributionally robust optimization (Cauchois et al., 2022).

A recent line of work develops *online conformal prediction* methods for the setting where the data arrives in a sequential order (Gibbs & Candès, 2021; 2022; Zaffran et al., 2022; Feldman et al., 2022). At each step, their methods output a prediction set parameterized by a single *radius* parameter that controls the size of the set. After receiving the true label, they adjust this parameter adaptively via *regret minimization* techniques—such as Online Gradient Descent (OGD) (Zinkevich, 2003)—on a certain quantile loss over the radius. These methods are shown to achieve empirical coverage frequency close to  $1 - \alpha$ , regardless of the data distribution (Gibbs & Candès, 2021). In addition to coverage, importantly, these methods achieve sublinear regret with respect to the quantile loss (Gibbs & Candès, 2022). Such regret guarantees ensure that the size of the prediction set is reasonable, and rule out “trivial” algorithms that achieve valid coverage by alternating between predicting the empty set and full set (cf. Section 2 for a discussion).

While regret minimization techniques achieve coverage and regret guarantees, they may fall short in more dynamic environments where we desire a strong performance not just over the entire time horizon (as captured by the regret), but also within every *sub-interval* of time. For example, if the data distribution shifts abruptly for a few times, we rather desire strong performance within each contiguous interval between two consecutive shifts, in addition to the entire

---

<sup>1</sup>Salesforce AI Research, Palo Alto, CA, USA. Correspondence to: Aadyot Bhatnagar <aadyotb@gmail.com>, Yu Bai <yu.bai@salesforce.com>.

horizon. Gibbs & Candès (2022) address this issue partially by proposing the Fully Adaptive Conformal Inference (FACI) algorithm, a meta-algorithm that aggregates multiple *experts* (base learners) that are OGD instances with different learning rates. However, their algorithm may not be best suited for achieving such interval-based guarantees, as each expert still runs over the full time horizon and is not really localized. This is also reflected in the fact that FACI achieves a near-optimal  $\tilde{O}(\sqrt{k})$  regret within intervals of a fixed length  $k$ , but is unable to achieve this over all lengths  $k \in [T]$  simultaneously.

In this paper, we design improved online conformal prediction algorithms by leveraging *strogly adaptive regret minimization*, a technique for attaining strong performance on all sub-intervals simultaneously in online learning (Daniely et al., 2015; Jun et al., 2017). Our proposed algorithm, Strongly Adaptive Online Conformal Prediction (SAOCP), is a new meta-algorithm that manages multiple experts, with the key difference that each expert now only operates on its own *active interval*. We summarize our contributions:

- We propose SAOCP, a new algorithm for online conformal prediction. SAOCP is a meta-algorithm that maintains multiple experts each with its own active interval, building on strongly adaptive regret minimization techniques (Section 3). We instantiate the experts as Scale-Free OGD (SF-OGD), an anytime variant of the OGD, which we also study as an independent algorithm.
- We prove that SAOCP achieves a near-optimal strongly adaptive regret of  $\tilde{O}(\sqrt{k})$  regret over all intervals of length  $k$  simultaneously, and that both SAOCP and SF-OGD achieve approximately valid coverage (Section 4).
- We show experimentally that SAOCP and SF-OGD attain better coverage in localized windows and smaller prediction sets than existing methods, on two real-world tasks: time series forecasting and image classification under distribution shift (Section 5).

### 1.1. Related work

**Conformal prediction** The original idea of conformal prediction (utilizing exchangeable data) is developed in the early work of Vovk et al. (1999; 2005); Shafer & Vovk (2008). Learning prediction sets via conformal prediction has since been adopted as a major approach for uncertainty quantification in regression (Papadopoulos, 2008; Vovk, 2012; Lei & Wasserman, 2014; Vovk et al., 2018; Romano et al., 2019; Gupta et al., 2019; Barber et al., 2021; 2022) and classification (Lei et al., 2013; Romano et al., 2020; Cauchois et al., 2020; 2022; Angelopoulos et al., 2021b), with further applications in general risk control (Bates et al., 2021; Angelopoulos et al., 2021a; 2022a), biological imaging (Angelopoulos et al., 2022b), and protein design (Fanjiang et al., 2022), to name a few.

Recent work also proposes to optimize the prediction sets’ *efficiency* (e.g. width or cardinality) in addition to coverage (Pearce et al., 2018; Park et al., 2020; Yang & Kuchibhotla, 2021; Stutz et al., 2022; Angelopoulos et al., 2021a;b; Bai et al., 2022). The regret that we consider can be viewed as a (surrogate) measure for efficiency in the online setting.

**Conformal prediction under distribution shift** For the more challenging case where data may exhibit distribution shift (and thus are no longer exchangeable), several approaches are proposed to achieve approximately valid coverage, such as reweighting (using prior knowledge about the data’s dependency structure) (Tibshirani et al., 2019; Podkopaev & Ramdas, 2021; Candès et al., 2021; Barber et al., 2022), distributionally robust optimization (Cauchois et al., 2020), or doubly robust calibration (Yang et al., 2022).

Our work makes an addition to the online conformal prediction line of work (Gibbs & Candès, 2021; 2022; Zafra et al., 2022; Feldman et al., 2022), which uses regret minimization techniques from the online learning literature (Zinkevich, 2003; Hazan, 2022) to adaptively adjust the size of the prediction set based on recent observations. Closely related to our work is the FACI algorithm of Gibbs & Candès (2022), which is a meta-algorithm that uses multiple experts for handling changing environments. Our meta-algorithm SAOCP differs in style from theirs, in that our experts only operate on their own active intervals, and it achieves a better guarantee on the strongly adaptive regret.

A related line of work studies conformal prediction for time series data. Chernozhukov et al. (2018); Xu & Xie (2021); Sousa et al. (2022) use randomization and ensembles to produce valid prediction sets for time series that are ergodic in a certain sense. Some other works directly apply vanilla conformal prediction to time series either without theoretical guarantees or requiring weaker notions of exchangeability (Dashevskiy & Luo, 2008; Wisniewski et al., 2020; Kath & Ziel, 2021; Stankeviciute et al., 2021; Sun & Yu, 2022).

**Strongly adaptive online learning** Our algorithms adapt techniques from the online learning literature, notably strongly adaptive regret minimization (Daniely et al., 2015; Jun et al., 2017; Zhang et al., 2018) and scale-free algorithms for achieving other kinds of adaptive (e.g. anytime) regret guarantees (Orabona & Pál, 2018).

## 2. Preliminaries

We consider standard learning problems in which we observe examples  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  and wish to predict a label  $y$  from input  $x$ . A *prediction set*  $C : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$  is a set-valued function that maps any input  $x$  to a *set* of predicted labels  $C(x) \subset \mathcal{Y}$ . Two prevalent examples are *prediction intervals*

for regression in which  $\mathcal{Y} = \mathbb{R}$  and  $C(x)$  is an interval, and *label prediction sets* for ( $m$ -class) classification in which  $\mathcal{Y} = [m]$  and  $C(x)$  is a subset of  $[m]$ . Prediction sets are a popular approach to quantify the uncertainty associated with the point prediction  $\hat{y} = f(x)$  of a black box model.

We study the problem of learning prediction sets in the *online* setting, in which the data  $(X_1, Y_1), \dots, (X_T, Y_T)$  arrive sequentially. At each time step  $t \in [T]$ , we output a prediction set  $\hat{C}_t = \hat{C}_t(X_t)$  based on the current input  $X_t$  and past observations  $\{(X_i, Y_i)\}_{i \leq t-1}$ , *before* observing the true label  $Y_t$ . The primary goal of the prediction set is to achieve *valid coverage*:  $\mathbb{P}[Y_t \in \hat{C}_t(X_t)] = 1 - \alpha$ , where  $1 - \alpha \in (0, 1)$  is the target coverage level pre-determined by the user. Standard choices for  $\alpha$  include  $\{0.1, 0.05\}$ , which correspond to  $\{90\%, 95\%\}$  target coverage respectively.

Throughout this paper, we use the standard notation  $\mathcal{O}(\cdot)$  to suppress absolute constants, and  $\tilde{\mathcal{O}}(\cdot)$  to suppress absolute constants and polylogarithmic factors, e.g.  $\log(T)$ .

**Online conformal prediction** We now review the idea of *online conformal prediction*, initiated by Gibbs & Candès (2021; 2022). This framework for learning prediction sets in the online setting achieves coverage guarantees even under distribution shift.

At each time  $t \in [T]$ , online conformal prediction assumes that we have a family of prediction sets  $\mathcal{C}_t = \{\hat{C}_t(x, s)\}_{x \in \mathcal{X}, s \in \mathbb{R}}$  specified by a *radius* parameter  $s \in \mathbb{R}$ , and we need to predict  $\hat{s}_t \in \mathbb{R}$  and output prediction set  $\hat{C}_t = \hat{C}_t(X_t, \hat{s}_t) \subset \mathcal{Y}$ . The family  $(\mathcal{C}_t)_{t \in [T]}$  is typically defined through *base predictors*  $\hat{f}_t$  (for example,  $\hat{f}_t \equiv f$  can be a fixed pretrained model). A standard example in regression is that we have a base predictor  $\hat{f}_t : \mathcal{X} \rightarrow \mathbb{R}$ , and we can choose  $\hat{C}_t(X_t, s) := [\hat{f}_t(X_t) - s, \hat{f}_t(X_t) + s]$  to be a prediction interval around  $\hat{f}_t(X_t)$ , in which case the radius  $s$  is the (half) width of the interval. In general, we allow any  $\mathcal{C}_t$  that are *nested sets* (Gupta et al., 2019) in the sense that  $\hat{C}_t(x, s) \subseteq \hat{C}_t(x, s')$  for all  $x \in \mathcal{X}$  and  $s \leq s'$ , so that a larger radius always yields a larger set.

Online conformal prediction adopts online learning techniques to learn  $\hat{s}_t$  based on past observations. Defining the “true radius”  $S_t := \inf\{s \in \mathbb{R} : Y_t \in \hat{C}_t(X_t, s)\}$  (i.e. the smallest radius  $s$  such that  $\hat{C}_t$  covers  $Y_t$ ), we consider the  $(1 - \alpha)$ -*quantile loss* (aka *pinball loss* (Koenker & Bassett Jr, 1978)) between  $S_t$  and any predicted radius  $\hat{s}$ :

$$\begin{aligned} \ell^{(t)}(\hat{s}) &= \ell_{1-\alpha}(S_t, \hat{s}) \\ &:= \max\{(1 - \alpha)(S_t - \hat{s}), \alpha(\hat{s} - S_t)\}. \end{aligned} \quad (1)$$

Throughout the rest of this paper, we assume that all true radii are bounded:  $S_t \in [0, D]$  almost surely for all  $t \in [T]$ .

After observing  $X_t$ , predicting the radius  $\hat{s}_t$ , and observing

the label  $Y_t$  (and hence  $S_t$ ), the gradient<sup>1</sup>  $\nabla \ell^{(t)}(\hat{s}_t)$  can be evaluated and has the following simple form:

$$\nabla \ell^{(t)}(\hat{s}_t) = \alpha - \mathbb{1}[\hat{s}_t < S_t] = \alpha - \underbrace{\mathbb{1}[Y_t \notin \hat{C}_t]}_{:= \text{err}_t}, \quad (2)$$

where  $\text{err}_t$  is the indicator of miscoverage at time  $t$  ( $\text{err}_t = 1$  if  $\hat{C}_t$  did not cover  $Y_t$ ). Gibbs & Candès (2021) perform an Online Gradient Descent (OGD) step to obtain  $\hat{s}_{t+1}$ :

$$\hat{s}_{t+1} = \hat{s}_t - \eta \nabla \ell^{(t)}(\hat{s}_t) = \hat{s}_t + \eta(\text{err}_t - \alpha), \quad (3)$$

where  $\eta > 0$  is a learning rate, and the algorithm is initialized at some  $\hat{s}_1 \in \mathbb{R}$ . Update (3) increases the predicted radius if  $\hat{C}_t$  did not cover  $Y_t$  ( $\text{err}_t = 1$ ), and decreases the radius otherwise. This makes intuitive sense as an approach for adapting the radius to recent observations.

**Adaptive Conformal Inference (ACI)** The ACI algorithm of Gibbs & Candès (2021) uses update (3) in conjunction with a specific choice of  $\mathcal{C}_t = \mathcal{C}_t^{\text{ACI}}$ , where

$$\hat{C}_t^{\text{ACI}}(X_t, \hat{s}_t) = \left\{ y : \tilde{S}_t(X_t, y) \leq Q_{\hat{s}_t} \left( \left\{ \tilde{S}_\tau \right\}_{\tau=1}^{t-1} \right) \right\}, \quad (4)$$

where  $\tilde{S}_t : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is any function (termed the *score* function),  $\tilde{S}_\tau := \tilde{S}_\tau(X_\tau, y_\tau)$  denotes score of the  $\tau$ -th observation, and  $Q_s(\cdot)$  denotes the  $s \in (0, 1)$ -th empirical quantile of a set (cf. (32)). In words, ACI’s confidence set contains all  $y$  whose score  $\tilde{S}_t(X_t, y)$  is below the  $\hat{s}_t$ -th quantile of the past scores, and they use (3) to learn this quantile level. The framework presented here generalizes the ACI algorithm where we allow any choice of  $\mathcal{C}_t$  that are nested sets, including but not limited to (4). For convenience of discussions, unless explicitly specified, we will also refer to this general version of (3) as ACI throughout this paper.

Empirically, we show in Section 5 that FACI (Gibbs & Candès, 2022) (an extension of ACI) performs better when trained to predict  $S_t$  directly under our more general parameterization, i.e.  $\hat{C}_t(X_t, \hat{s}_t) = \{y : \tilde{S}_t(X_t, y) \leq \hat{s}_t\}$ .

**Coverage and regret guarantees** Gibbs & Candès (2021) show that ACI<sup>2</sup> achieves approximately valid (empirical) coverage in the sense that the empirical miscoverage frequency is close to the target level  $\alpha$ :

$$\text{CovErr}(T) := \left| \frac{1}{T} \sum_{t=1}^T \text{err}_t - \alpha \right| \leq \frac{D + \eta}{\eta T}. \quad (5)$$

In addition to coverage, by standard online learning analyses, ACI achieves a *regret* bound on the quantile losses  $\{\ell^{(t)}\}_{t \in [T]}$ : we have  $\text{Reg}(T) \leq \mathcal{O}(D^2/\eta + \eta T) \leq \mathcal{O}(D\sqrt{T})$  (with optimally chosen  $\eta$ ), where

$$\text{Reg}(T) := \sum_{t=1}^T \ell^{(t)}(\hat{s}_t) - \inf_{s^* \in \mathbb{R}} \sum_{t=1}^T \ell^{(t)}(s^*). \quad (6)$$

<sup>1</sup>More precisely,  $\nabla \ell^{(t)}(\hat{s}_t)$  is a subgradient.

<sup>2</sup>Their results are established on the specific choice of  $\mathcal{C}_t$  in (4), but can be extended directly to any  $\mathcal{C}_t$  that are nested sets.

**Algorithm 1** Strongly Adaptive Online Conformal Prediction (SAOCP), adapted from Jun et al. (2017).

---

**Input:** Target coverage  $1 - \alpha \in (0, 1)$ ; maximum possible radius  $D > 0$

```

1 for  $t = 1, \dots, T$  do
    // Obtain prediction interval by aggregating active experts
2 Initialize new expert  $\mathcal{A}_t = \text{SF-OGD}(\alpha \leftarrow \alpha; \eta \leftarrow D/\sqrt{3}; \hat{s}_1 \leftarrow \hat{s}_{t-1})$  (Algorithm 2), and set weight  $w_{t,t} = 0$ 
3 Compute active set  $\text{Active}(t) = \{i \in [T] : t - L(i) < i \leq t\}$ , where  $L(i)$  is defined in (8)
4 Compute prior probability  $\pi_i \propto i^{-2}(1 + \lceil \log_2 i \rceil)^{-1} \mathbb{1}[i \in \text{Active}(t)]$ 
5 Compute un-normalized probability  $\hat{p}_i = \pi_i [w_{t,i}]_+$  for all  $i \in [t]$ 
6 Normalize  $p = \hat{p} / \|\hat{p}\|_1 \in \Delta^t$  if  $\|\hat{p}\|_1 > 0$ , else  $p = \pi$ 
7 Compute predicted radius  $\hat{s}_t = \sum_{i \in \text{Active}(t)} p_i \hat{s}_{i,t}$  (for  $t \geq 2$ ), and  $\hat{s}_t = 0$  for  $t = 1$ 
8 Observe input  $X_t \in \mathcal{X}$  and return prediction set  $\hat{C}_t(X_t, \hat{s}_t)$ 
    // Use meta loss and per-expert losses to update experts
9 Observe true label  $Y_t \in \mathcal{Y}$ , compute true radius  $S_t = \inf\{s \in \mathbb{R} : Y_t \in \hat{C}_t(X_t, s)\}$  and loss function  $\ell^{(t)}(\cdot) = \ell_{1-\alpha}(S_t, \cdot)$ 
10 for  $i \in \text{Active}(t)$  do
11     Update expert  $\mathcal{A}_i$  with  $(X_t, Y_t)$  and obtain next predicted radius  $\hat{s}_{i,t+1}$ 
12     Compute  $g_{i,t} = \begin{cases} \frac{1}{D} (\ell^{(t)}(\hat{s}_t) - \ell^{(t)}(\hat{s}_{i,t})) & w_{i,t} > 0 \\ \frac{1}{D} [\ell^{(t)}(\hat{s}_t) - \ell^{(t)}(\hat{s}_{i,t})]_+ & w_{i,t} \leq 0 \end{cases}$ 
13     Update expert weight  $w_{i,t+1} = \frac{1}{t-i+1} \left( \sum_{j=i}^t g_{i,j} \right) \left( 1 + \sum_{j=i}^t w_{i,j} g_{i,j} \right)$ 
    
```

---

One advantage of the regret as an additional performance measure alongside coverage is that it rules out certain algorithms that achieve good coverage in a “trivial” fashion and are not useful in practice — for example,  $\hat{C}_t$  may simply alternate between the empty set and the full set  $\{\alpha, 1 - \alpha\}$  proportion of the time, which satisfies the coverage bound (5) on arbitrary data distributions, but suffers from linear regret even on certain simple data distributions (cf. Appendix A.2).

We remark that the regret has a connection to the coverage error in that  $\text{CovErr}(T) = |\frac{1}{T} \sum_{t=1}^T \nabla \ell^{(t)}(\hat{s}_t)|$ , i.e. the coverage error is equal to the average gradient (derivative) of the losses. However, without further distributional assumptions, regret bounds and coverage bounds do not imply each other in a black-box fashion (see e.g. Gibbs & Candès (2021, Appendix A)) and need to be established separately for each algorithm.

### 3. Strongly Adaptive Online Conformal Prediction

Our approach is motivated from the observation that regret minimization is in a sense limited, as the regret measures performance over the *entire time horizon*  $[T]$ , which may be insufficient when the algorithm encounters changing environments. For example, if  $S_t = 1$  for  $1 \leq t \leq T/2$  and  $S_t = 100$  for  $T/2 < t \leq T$ , then achieving small regret on all (*sub*)-intervals of size  $T/2$  is a much stronger guarantee than achieving small regret over  $[T]$ . For this reason, we seek *localized* guarantees over all intervals simultaneously, to prevent worst-case scenarios such as significant miscoverage or high radius within a specific interval.

The *Strongly Adaptive Regret* (SARegret) (Daniely et al.,

2015; Zhang et al., 2018) has been proposed in the online learning literature as a generalization of the regret that captures the performance of online learning algorithms over all intervals simultaneously. Concretely, for any  $k \in [T]$ , the SARegret of length  $k$  of any algorithm is defined as

$$\text{SAReg}(T, k) := \max_{[\tau, \tau+k-1] \subseteq [T]} \left( \sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(\hat{s}_t) - \inf_{s^*} \sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(s^*) \right) \quad (7)$$

$\text{SAReg}(T, k)$  measures the maximum regret over all intervals of length  $k$ , which reduces to the usual regret at  $k = T$ , but may in addition be smaller for smaller  $k$ .

**Algorithm: SAOCP** We leverage techniques for minimizing the strongly adaptive regret to perform online conformal prediction. Our main algorithm, Strongly Adaptive Online Conformal Prediction (SAOCP, described in Algorithm 1), adapts the work of Jun et al. (2017) to the online conformal prediction setting. At a high level, SAOCP is a meta-algorithm that manages multiple *experts*, where each expert is itself an arbitrary online learning algorithm taking charge of its own *active interval* that has a finite *lifetime*. At each  $t \in [T]$ , Algorithm 1 instantiates a new expert  $\mathcal{A}_t$  with active interval  $[t, t + L(t) - 1]$ , where  $L(t)$  is its lifetime:

$$L(t) := g \cdot \max_{n \in \mathbb{Z}} \{2^n : t \equiv 0 \pmod{2^n}\}, \quad (8)$$

and  $g \in \mathbb{Z}_{\geq 1}$  is a multiplier for the lifetime of each expert. It is straightforward to see that at most  $g \lceil \log_2 t \rceil$  experts are active at any time  $t$  under choice (8), granting Algorithm 1 a total runtime of  $\mathcal{O}(T \log T)$  for any  $g = \Theta(1)$ . Then, at any time  $t$ , the predicted radius  $\hat{s}_t$  is obtained by aggregating

**Algorithm 2** Scale-Free Online Gradient Descent (SF-OGD), adapted from Orabona & Pál (2018).

**Input:**  $\alpha \in (0, 1)$ , learning rate  $\eta > 0$ , init.  $\hat{s}_1 \in \mathbb{R}$

14 **for**  $t \geq 1$  **do**

15     Observe input  $X_t \in \mathcal{X}$

16     **Return** prediction set  $\widehat{C}_t(X_t, \hat{s}_t)$

17     Observe true label  $Y_t \in \mathcal{Y}$  and compute true radius  $S_t = \inf\{s \in \mathbb{R} : Y_t \in \widehat{C}_t(X_t, s)\}$ .

18     Compute loss  $\ell^{(t)}(\cdot) = \ell_{1-\alpha}(S_t, \cdot)$

19     Update predicted radius

$$\hat{s}_{t+1} = \hat{s}_t - \eta \frac{\nabla \ell^{(t)}(\hat{s}_t)}{\sqrt{\sum_{i=1}^t \|\nabla \ell^{(i)}(\hat{s}_i)\|_2^2}} \quad (10)$$

the predictions of active experts (Line 7):

$$\hat{s}_t = \sum_{i \in \text{Active}(t)} p_{i,t} \hat{s}_{i,t},$$

where the weights  $\{p_{i,t}\}_{i \in [t]}$  (Lines 4-6) rely on the  $\{w_{i,t}\}_{i \in [t]}$  computed by the *coin betting* scheme (Orabona & Pál, 2016; Jun et al., 2017) in Lines 12-13.

**Choice of expert** In principle, SAOCP allows any choice of the expert that is a good regret minimization algorithm over its own active interval satisfying *anytime* regret guarantees. We choose the experts to be Scale-Free OGD (SF-OGD; Algorithm 2) (Orabona & Pál, 2018), a variant of OGD that decays its effective learning rate based on cumulative past gradient norms (cf. (10)).

On the quantile loss (1) (executed over the full horizon  $[T]$  with learning rate  $\eta = \Theta(D)$ ;  $\eta = D/\sqrt{3}$  is optimal), SF-OGD enjoys an anytime regret guarantee (Proposition A.2)

$$\text{Reg}(t) \leq \mathcal{O}(D\sqrt{t}) \quad \text{for all } t \in [T], \quad (9)$$

which follows directly by applying Orabona & Pál (2018, Theorem 2). Plugging SF-OGD into Line 2 of Algorithm 1 gives our full SAOCP algorithm.

**SF-OGD as an independent algorithm** As a strong regret minimization algorithm itself, SF-OGD can also be run independently (over  $[T]$ ) as an algorithm for online conformal prediction (described in Algorithm 2). We find empirically that SF-OGD itself already achieves strong performances in many scenarios (Section 5).

## 4. Theory

### 4.1. Strongly Adaptive Regret

We begin by showing the SARegret guarantee of SAOCP. As we instantiate SAOCP with SF-OGD as the experts, the proof follows directly by plugging the regret bound for SF-OGD (9) into the SARegret guarantee for SAOCP (Jun et al., 2017), and can be found in Appendix B.1.

**Proposition 4.1** (SARegret bound for SAOCP). *Algorithm 1 achieves the following SARegret bound simultaneously for all lengths  $k \in [T]$ :*

$$\text{SAReg}(T, k) \leq 15D\sqrt{k(\log T + 1)} \leq \widetilde{\mathcal{O}}(D\sqrt{k}). \quad (11)$$

The  $\widetilde{\mathcal{O}}(D\sqrt{k})$  rate achieved by SAOCP is near-optimal for general online convex optimization problems, due to the standard regret lower bound  $\Omega(D\sqrt{k})$  over any fixed interval of length  $k$  (Orabona, 2019, Theorem 5.1).

**Comparison with FACI** The SARegret guarantee of SAOCP in (11) improves substantially over the FACI (Fully Adaptive Conformal Inference) algorithm (Gibbs & Candès, 2022), an extension of ACI. Concretely, (11) holds *simultaneously for all lengths  $k$* . By contrast, FACI achieves  $\text{SAReg}(T, k) \leq \widetilde{\mathcal{O}}(D^2/\eta + \eta k)$  in our setting (cf. their Theorem 3.2), where  $\eta > 0$  is their meta-algorithm learning rate. This can imply the same rate  $\widetilde{\mathcal{O}}(D\sqrt{k})$  for a *single  $k$*  by optimizing  $\eta$ , but not multiple values of  $k$  simultaneously. As an example, for small interval sizes of order  $k = \mathcal{O}(T^c) = o(T)$  with  $c < 1$ , FACI with the optimized  $\eta \gtrsim T^{-c/2}$  can achieve the optimal rate  $\text{SAReg}(T, k) \leq \widetilde{\mathcal{O}}(DT^{c/2})$  on *this particular  $c$* , but will perform suboptimally on other values of  $c$ . By contrast, SAOCP achieves  $\text{SAReg}(T, k) \leq \widetilde{\mathcal{O}}(DT^{c/2})$  simultaneously for all  $c \in (0, 1)$ .

Also, in terms of algorithm styles, while both SAOCP and FACI are meta-algorithms that maintain multiple experts (base algorithms), a main difference between them is that all experts in FACI differ in their learning rates and are all active over  $[T]$ , whereas experts in SAOCP differ in their active intervals (cf. (8)).

**Dynamic regret** The (worst-case) dynamic regret—which measures the performance of an online learning algorithm against the worst-case (strongest) comparator sequence—is another generalization of the regret for capturing the performance in changing environments (Zinkevich, 2003; Besbes et al., 2015). Building on the reduction from dynamic regret to strongly adaptive regret (Zhang et al., 2018), we show that SAOCP also achieves a worst-case dynamic regret bound on any interval within  $[T]$ , with rate depending on a certain *path length* of the true radii  $\{S_t\}_{t \geq 1}$ ; see Proposition B.1 and the discussions thereafter.<sup>3</sup>

<sup>3</sup>We remark that a more general version of dynamic regret which measures the regret against an *arbitrary* comparator sequence has also been considered in the literature. However, the relationship between the SARegret and the general dynamic regret remains unclear; see e.g. Zhao et al. (2022) for a discussion.

## 4.2. Coverage

Recall the empirical coverage error defined in (5):

$$\text{CovErr}(T) = \left| \frac{1}{T} \sum_{t=1}^T \text{err}_t - \alpha \right|.$$

Without any distributional assumptions, we show that SF-OGD achieves  $\text{CovErr}(t) \leq \mathcal{O}(t^{-1/4} \log t)$  for any  $t \in [T]$ . So its empirical coverage converges to the target  $1 - \alpha$  as  $T \rightarrow \infty$ , similar to ACI (though with a slightly slower rate). The proof (Appendix B.3) builds on a grouping argument and the fact that the effective learning rate  $\eta / \sqrt{\sum_{\tau=1}^t \|\nabla \ell^{(\tau)}(\hat{s}_\tau)\|_2^2}$  of SF-OGD changes slowly in  $t$ .

**Theorem 4.2** (Coverage bound for SF-OGD). *Algorithm 2 with any learning rate  $\eta = \Theta(D)$  and any initialization  $\hat{s}_1 \in [0, D]$  achieves  $\text{CovErr}(T) \leq \mathcal{O}(\alpha^{-2} T^{-1/4} \log T)$  for any  $T \geq 1$ .*

We now provide a distribution-free coverage bound for SAOCP, building on a similar grouping argument as in Theorem 4.2.

**Theorem 4.3** (Coverage bound for SAOCP; Informal version of Theorem B.3). *For any  $T \geq 1$ , a randomized variant of Algorithm 1 where Line 7 is replaced by sampling an expert  $i \sim p_t$  and predicting  $\hat{s}_t := \hat{s}_{t,i}$  achieves*

$$\text{CovErr}(T) \leq \mathcal{O}(\inf_{\beta} (T^{1/2-\beta} + T^{\beta-1} S_{\beta}(T))). \quad (12)$$

Theorem 4.3 considers a randomized variant of SAOCP, and its coverage bound depends on a quantity  $S_{\beta}(T)$  (full definition in Theorem B.3) that measures the smoothness of the expert weights and the cumulative gradient norms for each individual expert. Both are expected for technical reasons and also appear in coverage bounds for other expert-style meta-algorithms such as FACI (Gibbs & Candès, 2022). For instance, if there exists  $\beta \in (1/2, 1)$  so that  $S_{\beta}(T) \leq \tilde{\mathcal{O}}(T^{\gamma})$  for some  $\gamma < 1 - \beta$ , then (12) implies a coverage bound  $\text{CovErr}(T) \leq \tilde{\mathcal{O}}(T^{-\min\{1/2-\beta, \beta-1+\gamma\}}) = o_T(1)$ .

We remark that Theorem 4.3 also holds more generically for other choices of the expert weights  $\{p_t\}_{t \in [T]}$  (Line 5-6) and active intervals, not just those specified in Algorithm 1. In particular, SF-OGD is the special case where there is only a single active expert over  $[T]$ . In this case, we can recover the  $\tilde{\mathcal{O}}(\alpha^{-2} T^{-1/4})$  bound of Theorem 4.2 (see Appendix B.4.1 for a detailed discussion).

**Additional coverage guarantee under distributional assumptions** Under some mild regularity assumptions on the distributions of  $S_1, \dots, S_T$ , we show in Theorem C.3 that SAOCP achieves approximately valid coverage on every sub-interval of time. Its coverage error on any interval  $I = [\tau, \tau + k - 1] \subseteq [1, T]$  is  $\tilde{\mathcal{O}}(k^{-1/(2q)} + (\text{Var}_I/k)^{1/q})$  for a certain  $q \geq 2$  that quantifies the regularity of the distribution, and  $\text{Var}_I$  is a certain notion of variation between

the true quantiles of  $S_t$  over  $t \in I$  (cf. (29)). In particular, we obtain an approximately valid coverage on any interval  $I$  for which  $\text{Var}_I = o(k)$ .

## 5. Experiments

We test SF-OGD (Algorithm 2) and SAOCP (Algorithm 1) empirically on two representative real-world online uncertainty quantification tasks: time series forecasting (Section 5.1) and image classification under distribution shift (Section 5.2). Choices of the prediction sets  $\{\hat{C}_t(x, s)\}_{x, s}$  will be described within each experiment. In both experiments, we compare against the following methods:

1. SCP: standard Split Conformal Prediction (Vovk et al., 2005) adapted to the online setting, which simply predicts the  $(1 - \alpha)$ -quantile of the past radii. SCP does not admit a valid coverage guarantee in our settings as the data may not be exchangeable in general;
2. NExCP: Non-Exchangeable SCP (Barber et al., 2022), a variant of SCP that handles non-exchangeable data by reweighting. We follow their recommendations and use an exponential weighting scheme that upweights more recent observations;
3. FACI (Gibbs & Candès, 2022) with their specific ‘‘quantile parametrization’’ (4), and score function  $\tilde{S}_t$  corresponding to our choice of  $\hat{C}_t$ ;
4. FACI-S: Generalized version of FACI applied to predicting the radius  $\hat{s}_t$ 's on our choice of  $\hat{C}_t$  directly.

Additional details about all methods can be found in Appendix E. Throughout this section we choose the target coverage level to be the standard  $1 - \alpha = 90\%$ .

### 5.1. Time Series Forecasting

**Setup** We consider multi-horizon time series forecasting problems with real-valued observations  $\{y_t\}_{t \geq 1} \in \mathbb{R}$ , where the base predictor  $\hat{f}$  uses the history  $X_t := y_{1:t}$  to predict  $H$  steps into the future, i.e.  $\hat{f}(X_t) = \{\hat{f}^{(h)}(X_t)\}_{h \in [H]} = \{\hat{y}_{t+h}^{(h)}\}_{h \in [H]}$ , where  $\hat{y}_{t+h}^{(h)}$  is a prediction for  $y_{t+h}$ . Using  $\hat{f}(X_t)$ , we produce fixed-width prediction intervals

$$\hat{C}_t^{(h)}(X_t, \hat{s}_t^{(h)}) := [\hat{y}_{t+h}^{(h)} - \hat{s}_t^{(h)}, \hat{y}_{t+h}^{(h)} + \hat{s}_t^{(h)}], \quad (13)$$

where  $\hat{s}_t^{(h)}$  is predicted by an independent copy of the online conformal prediction algorithm for each  $h \in [H]$  (so that there are  $H$  such algorithms in parallel). We form our online setting using a standard rolling window evaluation loop, wherein each *batch* consists of predicting all  $H$  intervals  $\{\hat{C}_t^{(h)}\}_{h \in [H]}$ , observing all  $H$  true values  $\{y_{t+h}\}_{h \in [H]}$ , and moving to the next batch by setting  $t \rightarrow t + H$ . For each  $h \in [H]$ , we only evaluate  $y_{t+h}$  against one interval  $\hat{C}_t^{(h)}(X_t, \hat{s}_t^{(h)})$ . After the evaluation is done, we use all pairs  $\{(y_{t+k}, \hat{y}_{t+k}^{(h)})\}_{k \in [H]}$  to update  $\hat{s}_t^{(h)} \rightarrow \hat{s}_{t+H}^{(h)}$ . To set

Improved Online Conformal Prediction via Strongly Adaptive Online Learning

Method	LGBM (MAE = 0.06)				ARIMA (MAE = 0.18)				Prophet (MAE = 0.12)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
SCP	<u>.844</u> <sub>.004</sub>	.127 <sub>.004</sub>	.252 <sub>.007</sub>	.017 <sub>.001</sub>	<u>.871</u> <sub>.004</sub>	.245 <sub>.029</sub>	.237 <sub>.008</sub>	.039 <sub>.021</sub>	<u>.783</u> <sub>.008</sub>	.178 <sub>.007</sub>	.355 <sub>.012</sub>	.019 <sub>.001</sub>
NExCP	<u>.875</u> <sub>.002</sub>	<u>.134</u> <sub>.004</sub>	.197 <sub>.006</sub>	.013 <sub>.001</sub>	<u>.871</u> <sub>.004</sub>	.245 <sub>.034</sub>	.227 <sub>.007</sub>	.040 <sub>.024</sub>	<u>.856</u> <sub>.003</sub>	.187 <sub>.007</sub>	.231 <sub>.007</sub>	.010 <sub>.001</sub>
FACI	<u>.866</u> <sub>.002</sub>	<b>.113</b> <sub>.003</sub>	.180 <sub>.005</sub>	<b>.009</b> <sub>.001</sub>	<u>.866</u> <sub>.004</sub>	<u>.232</u> <sub>.034</sub>	.214 <sub>.006</sub>	<u>.034</u> <sub>.022</sub>	<u>.867</u> <sub>.002</sub>	<u>.175</u> <sub>.007</sub>	.184 <sub>.005</sub>	<u>.006</u> <sub>.000</sub>
SF-OGD	<u>.889</u> <sub>.002</sub>	.138 <sub>.004</sub>	<u>.154</u> <sub>.004</sub>	.011 <sub>.001</sub>	<u>.877</u> <sub>.003</sub>	.250 <sub>.036</sub>	<u>.195</u> <sub>.006</sub>	.037 <sub>.022</sub>	<u>.888</u> <sub>.001</sub>	.186 <sub>.007</sub>	<b>.138</b> <sub>.003</sub>	.007 <sub>.001</sub>
FACI-S	<u>.883</u> <sub>.002</sub>	.128 <sub>.003</sub>	<u>.163</u> <sub>.004</sub>	.010 <sub>.001</sub>	<u>.872</u> <sub>.004</sub>	.238 <sub>.034</sub>	<u>.201</u> <sub>.006</sub>	.035 <sub>.021</sub>	<u>.885</u> <sub>.001</sub>	.180 <sub>.007</sub>	.144 <sub>.003</sub>	<u>.006</u> <sub>.000</sub>
SAOCP	<u>.882</u> <sub>.002</sub>	<u>.121</u> <sub>.003</sub>	<b>.143</b> <sub>.004</sub>	<b>.009</b> <sub>.000</sub>	<u>.864</u> <sub>.003</sub>	<b>.221</b> <sub>.031</sub>	<b>.190</b> <sub>.005</sub>	<b>.033</b> <sub>.022</sub>	<u>.872</u> <sub>.001</sub>	<b>.173</b> <sub>.007</sub>	<u>.143</u> <sub>.003</sub>	<b>.005</b> <sub>.000</sub>

Table 1. Results on M4 Hourly (414 time series) with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). For all base predictors, SAOCP achieves the best or second-best width, local coverage error, and strongly adaptive regret.

Method	LGBM (MAE = 0.13)				ARIMA (MAE = 0.06)				Prophet (MAE = 0.32)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
SCP	<u>.769</u> <sub>.004</sub>	.184 <sub>.002</sub>	.466 <sub>.005</sub>	.031 <sub>.001</sub>	<u>.896</u> <sub>.002</sub>	.122 <sub>.002</sub>	.290 <sub>.004</sub>	.018 <sub>.001</sub>	<u>.599</u> <sub>.005</sub>	.349 <sub>.004</sub>	.614 <sub>.005</sub>	.051 <sub>.001</sub>
NExCP	<u>.818</u> <sub>.002</sub>	.183 <sub>.002</sub>	.420 <sub>.004</sub>	.015 <sub>.000</sub>	<u>.891</u> <sub>.001</sub>	.116 <sub>.002</sub>	.296 <sub>.004</sub>	.012 <sub>.001</sub>	<u>.715</u> <sub>.003</sub>	.356 <sub>.004</sub>	.559 <sub>.004</sub>	.019 <sub>.000</sub>
FACI	<u>.846</u> <sub>.002</sub>	.169 <sub>.002</sub>	.308 <sub>.003</sub>	.008 <sub>.000</sub>	<u>.886</u> <sub>.001</sub>	<u>.101</u> <sub>.002</sub>	.259 <sub>.003</sub>	<b>.008</b> <sub>.000</sub>	<u>.767</u> <sub>.003</sub>	.344 <sub>.004</sub>	.397 <sub>.004</sub>	.014 <sub>.001</sub>
SF-OGD	<u>.873</u> <sub>.001</sub>	.173 <sub>.002</sub>	.246 <sub>.003</sub>	.011 <sub>.001</sub>	<u>.892</u> <sub>.001</sub>	.106 <sub>.002</sub>	.245 <sub>.003</sub>	.011 <sub>.001</sub>	<u>.862</u> <sub>.001</sub>	.354 <sub>.004</sub>	.220 <sub>.002</sub>	.008 <sub>.001</sub>
FACI-S	<u>.875</u> <sub>.001</sub>	.169 <sub>.002</sub>	<u>.240</u> <sub>.003</sub>	.010 <sub>.001</sub>	<u>.891</u> <sub>.001</sub>	.103 <sub>.002</sub>	<u>.243</u> <sub>.003</sub>	.010 <sub>.001</sub>	<u>.866</u> <sub>.002</sub>	<b>.349</b> <sub>.004</sub>	<u>.210</u> <sub>.002</sub>	<u>.007</u> <sub>.001</sub>
SAOCP	<u>.869</u> <sub>.001</sub>	<b>.162</b> <sub>.002</sub>	<b>.213</b> <sub>.002</sub>	<b>.007</b> <sub>.000</sub>	<u>.875</u> <sub>.001</sub>	<b>.093</b> <sub>.002</sub>	<b>.238</b> <sub>.002</sub>	<b>.008</b> <sub>.001</sub>	<u>.867</u> <sub>.001</sub>	<u>.352</u> <sub>.004</sub>	<b>.172</b> <sub>.001</sub>	<b>.005</b> <sub>.000</sub>

Table 2. Results on M4 Daily (4227 time series) with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). SAOCP achieves the best width, local coverage error, and strongly adaptive regret for all base predictors. The only methods which achieve global coverage in  $(0.85, 0.95)$  for LGBM and Prophet are the ones that predict  $\hat{s}_{t+1}$  directly, not as a quantile of  $S_1, \dots, S_t$ .

the maximum radius for SF-OGD and SAOCP, we choose  $D/\sqrt{3}$  for each horizon  $h$  to be the largest  $h$ -step residual observed on the calibration split of the training data.

**Base predictors** We consider three diverse types of base predictors (models), and we use their implementations in Merlion v2.0.0 (Bhatnagar et al., 2021):

1. LGBM: A model which uses gradient boosted trees to predict  $\hat{y}_{t+h}^{(h)} = \hat{f}^{(h)}(y_{t-L+1}, \dots, y_t)$ . This approach attains strong performance on many time series benchmarks (Elsayed et al., 2021; Bhatnagar et al., 2021).
2. ARIMA(10,  $d^*$ , 10): The classical AutoRegressive Integrated Moving Average stochastic process model for a time series, where the difference order  $d^*$  is chosen by KPSS stationarity test (Kwiatkowski et al., 1992).
3. Prophet (Taylor & Letham, 2017): A popular Bayesian model which directly predicts the value  $y$  as a function of time, i.e.  $\hat{y}_t = \hat{f}(t)$ .

**Datasets** We evaluate on four datasets totaling 5111 time series: the hourly (414 time series), daily (4227 time series), and weekly (359 time series) subsets of the M4 Competition, a dataset of time series from many domains including industries, demographics, environment, finance, and transportation (Makridakis et al., 2018); and NN5, a dataset of 111 time series of daily banking data (Ben Taieb et al., 2012). We normalize each time series to lie in  $[0, 1]$ .

We use horizons  $H$  of 24, 30, and 26 for hourly, daily, and weekly data, respectively. Each time series of length  $L$  is split into a training set of length  $L - 120$  with 80% for training the base predictor and 20% for initializing the UQ methods, and a test set of length 120 to test the UQ methods.

**Metrics** For each experiment, we average the following statistics across all time series: global coverage, median width, worst-case local coverage error

$$\text{LCE}_k := \max_{[\tau, \tau+k-1] \subseteq [1, T]} \left| \alpha - \frac{1}{k} \sum_{t=\tau}^{\tau+k-1} \text{err}_t \right|, \quad (14)$$

and strongly adaptive regret  $\text{SAReg}(T, k)$  (7), which we abbreviate as  $\text{SAReg}_k$ . In all cases, we use an interval length of  $k = 20$ . We also report the average mean absolute error (MAE) of each base predictor.

**Results** We report results on M4 Hourly and M4 Daily in Tables 1, 2, and on M4 Weekly and NN5 Daily in Tables 4, and 5 (Appendix D).

SAOCP consistently achieves global coverage in  $(0.85, 0.95)$ , and it obtains the best or second-best interval width, local coverage error, and strongly adaptive regret for all base predictors on all 3 M4 datasets. FACI-S generally achieves better  $\text{LCE}_k$  and  $\text{SAReg}_k$  than FACI, showing the benefits of predicting  $\hat{s}_{t+1}$  directly, rather than as a quantile of  $S_1, \dots, S_t$ . The relative performance of FACI-S and SF-OGD varies, though FACI-S is usually a bit better.

Method	LCE <sub>10</sub>	LCE <sub>20</sub>	LCE <sub>30</sub>	LCE <sub>40</sub>	LCE <sub>50</sub>	LCE <sub>60</sub>
SCP	.585	.466	.389	.336	.295	.263
NExCP	.592	.420	.323	.262	.221	.191
FACI	.467	.308	.232	.187	.157	.134
SF-OGD	.394	.246	.183	.146	.123	<u>.104</u>
FACI-S	.385	<u>.240</u>	<u>.179</u>	.144	<u>.122</u>	<u>.104</u>
SAOCP	<b>.337</b>	<b>.213</b>	<b>.159</b>	<b>.128</b>	<b>.107</b>	<b>.091</b>

Table 3. Local coverage error of the LGBM model evaluated on the M4 Daily dataset (4227 time series) with target coverage  $1 - \alpha = 0.9$  and multiple interval sizes  $k$ . Best results are **bold**, while second-best are underlined. SAOCP achieves the best LCE<sub>k</sub> for all  $k$ . Standard errors for all reported values are below 0.005.

However, SAOCP consistently achieves better LCE<sub>k</sub> and SAReg<sub>k</sub> than both FACI-S and SF-OGD.

There are multiple instances where all of SCP/NExCP/FACI fail to attain global coverage in (0.85, 0.95) (Tables 2 and 4). The base predictor’s MAE is at least 0.13 in all these cases, suggesting an advantage of predicting  $\hat{s}_{t+1}$  directly as in SF-OGD/SAOCP when the underlying base predictor is inaccurate.

**Evaluations with different interval lengths** Tables 1, 2, 4, and 5 evaluate LCE<sub>k</sub> of various online conformal prediction methods with an interval size of  $k = 20$  across a wide range of experimental settings. To determine whether these results hold for different interval lengths, we report LCE<sub>k</sub> for the LGBM model on the M4 Daily dataset for  $k \in \{10, 20, 30, 40, 50, 60\}$  (Table 3). Out of the box, SAOCP achieves the best LCE for all  $k$ , further exhibiting its robustness.

**Additional experiments with ensemble models** In Appendix F, we use EnbPI (Xu & Xie, 2021) to train a bootstrapped ensemble, and we compare EnbPI’s results with those obtained by applying NExCP, FACI, SF-OGD, and SAOCP to the residuals produced by that ensemble. The results largely mirror those in the main paper.

## 5.2. Image Classification Under Distribution Shift

**Datasets and Setup** We evaluate the ability of conformal prediction methods to maintain coverage when the underlying distribution shifts in a systematic manner. We use a ResNet-50 classifier (He et al., 2016) pre-trained on ImageNet and implemented in PyTorch (Paszke et al., 2019). Here,  $x \in \mathcal{X}$  is an image, and  $y \in \mathcal{Y} = [m]$  is its class. To construct structured distribution shifts away from the training distribution, we use TinyImageNet-C and ImageNet-C (Hendrycks & Dietterich, 2019), which are corrupted versions of the TinyImageNet ( $m = 200$  classes) (Le & Yang, 2015) and ImageNet ( $m = 1000$  classes) (Deng et al., 2009) test sets designed to evaluate model robustness. These cor-

rupted datasets apply 15 visual corruptions at 5 different severity levels to each image in the original test set.

We consider two regimes: **sudden shifts** where the corruption level alternates between 0 (the base test set) and 5, and **gradual shifts** where the corruption level increases in the order of  $\{0, 1, \dots, 5\}$ . We randomly sample 500 data points for each corruption level before changing to the next level.

**Prediction Sets** We follow Angelopoulos et al. (2021b) to construct our prediction sets. Let  $\hat{f} : \mathbb{R}^d \rightarrow \Delta^m$  be a classifier that outputs a probability distribution on the  $m$ -simplex. At each  $t$ , we sample  $U_t \sim \text{Unif}[0, 1]$  and let

$$S_t(x, y) = \lambda \sqrt{[k_y - k_{reg}]_+} + U_t \hat{f}_y(x) + \sum_{i=1}^{k_y-1} \hat{f}_{\pi(i)}(x) \\ \hat{C}_t(X_t) = \{y : S_t(X_t, y) \leq \hat{s}_t\}, \quad (15)$$

where  $\pi$  is the permutation that ranks  $\hat{f}(x)$  in decreasing order,  $\pi(k_y) = y$ , and  $\lambda$  and  $k_{reg}$  are regularization parameters designed to reduce the size of the prediction set. For TinyImageNet, we use  $\lambda = 0.01$  and  $k_{reg} = 20$ . For ImageNet, we use  $\lambda = 0.01$  and  $k_{reg} = 10$ . Accordingly, we set the maximum radius  $D = \lambda \sqrt{m - k_{reg}}$  for SF-OGD and SAOCP, which we note is the maximum value of  $S_t(\cdot, \cdot)$ .

**Metrics** When evaluating the UQ methods, we plot the local coverage and prediction set size (PSS) of each method using an interval length of  $k = 100$ ,

$$\text{LocalCov}(t) = \frac{1}{100} \sum_{i=t}^{t+99} \mathbb{1}[Y_i \in \hat{C}_i(X_i)] \\ \text{LocalPSS}(t) = \frac{1}{100} \sum_{i=t}^{t+99} |\hat{C}_i(X_i)|.$$

We compare the local coverage to a target of  $1 - \alpha$ , while we compare the local PSS to the  $1 - \alpha$  empirical quantile of the oracle set sizes  $\text{PSS}_t^* = |\{y : S_t(X_t, y) \leq S_t(X_t, Y_t)\}|$ . These targets are the “best fixed” values in each window. We also report the worst-case local coverage error LCE<sub>100</sub> (14).

**Results** We evaluate the UQ methods on TinyImageNet and TinyImageNet-C in Figure 1, and on ImageNet and ImageNet-C in Figure 2 (Appendix G). In both sudden and gradual distribution shift, the local coverage of SAOCP and SF-OGD remains the closest to the target of 0.9. The difference is more notable when the distribution shifts suddenly. When the distribution shifts more gradually, NExCP, FACI, and FACI-S have worse coverage than SAOCP and SF-OGD at the first change point, which is where the largest change in the best set size occurs.

All methods besides SCP predict sets of similar sizes, though FACI’s, FACI-S’s, and NExCP’s prediction set sizes adapt more slowly to changes in the best fixed size (e.g.  $t \in [500, 700]$  for gradual shift in Figure 2). On TinyImageNet, SAOCP obtains slightly better local coverage than

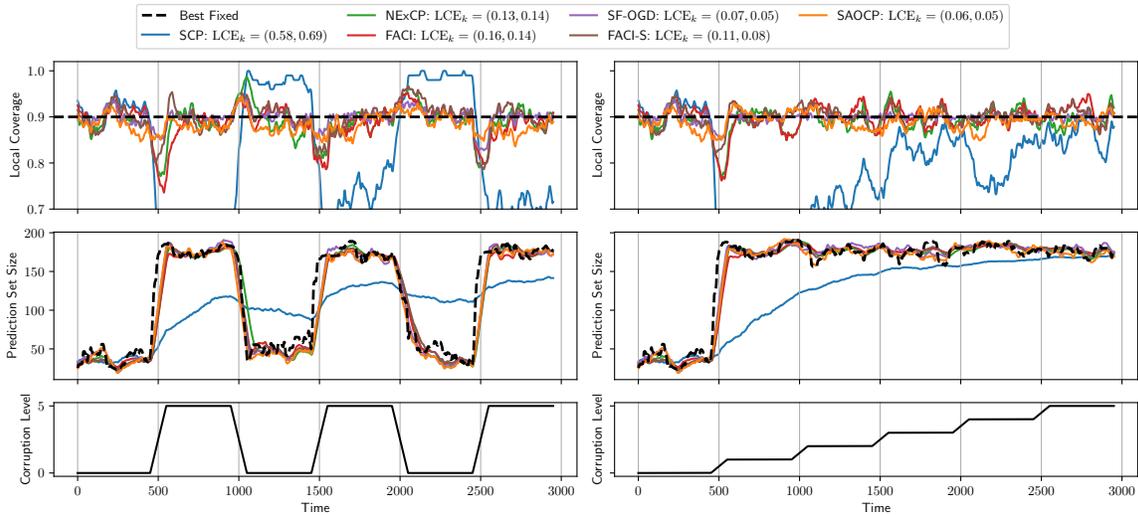


Figure 1. Local coverage (top row) and prediction set size (second row) achieved by various UQ methods when the distribution shifts between TinyImageNet and TinyImageNet-C every 500 steps. We plot moving averages with a window size of 100. Left: sudden shifts between corruption level 0 and 5. Right: gradual shift from level 0 to 5. SAOCP and SF-OGD’s local coverage remain the closest to the target of 0.9, especially at the change points. While their prediction sets have similar size,  $LCE_k$  is lower for SAOCP than SF-OGD.

SF-OGD, and they both have similar prediction set sizes (Figure 1). On ImageNet, SAOCP and SF-OGD attain similar local coverages, but SAOCP tends to attain that coverage with a smaller prediction set (Figure 2).

## 6. Conclusion

This paper develops new algorithms for online conformal prediction under arbitrary distribution shifts. Our algorithms achieve approximately valid coverage and better strongly adaptive regret than existing work. On real-world experiments, our proposed algorithms achieve coverage closer to the target within local windows, and they produce smaller prediction sets than existing methods. Our work opens up many questions for future work, such as obtaining stronger coverage guarantees, or characterizing the optimality of the learned radii under various settings with distribution shift.

## References

- Angelopoulos, A. N. and Bates, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification, 2021. URL <https://arxiv.org/abs/2107.07511>.
- Angelopoulos, A. N., Bates, S., Candès, E. J., Jordan, M. I., and Lei, L. Learn then test: Calibrating predictive algorithms to achieve risk control. *arXiv preprint arXiv:2110.01052*, 2021a.
- Angelopoulos, A. N., Bates, S., Jordan, M., and Malik, J. Uncertainty sets for image classifiers using conformal prediction. In *International Conference on Learning Rep-*
- resentations*, 2021b. URL [https://openreview.net/forum?id=eNdiU\\_DbM9](https://openreview.net/forum?id=eNdiU_DbM9).
- Angelopoulos, A. N., Bates, S., Fisch, A., Lei, L., and Schuster, T. Conformal risk control. *arXiv preprint arXiv:2208.02814*, 2022a.
- Angelopoulos, A. N., Kohli, A. P., Bates, S., Jordan, M., Malik, J., Alshaabi, T., Upadhyayula, S., and Romano, Y. Image-to-image regression with distribution-free uncertainty quantification and applications in imaging. In *International Conference on Machine Learning*, pp. 717–730. PMLR, 2022b.
- Bai, Y., Mei, S., Wang, H., Zhou, Y., and Xiong, C. Efficient and differentiable conformal prediction with general function classes. In *International Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?id=Ht85\\_jyihxp](https://openreview.net/forum?id=Ht85_jyihxp).
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. Predictive inference with the jackknife+. *The Annals of Statistics*, 49(1):486 – 507, 2021. doi: 10.1214/20-AOS1965. URL <https://doi.org/10.1214/20-AOS1965>.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. Conformal prediction beyond exchangeability, 2022. URL <https://arxiv.org/abs/2202.13415>.
- Bates, S., Angelopoulos, A., Lei, L., Malik, J., and Jordan, M. Distribution-free, risk-controlling prediction sets. *Journal of the ACM (JACM)*, 68(6):1–34, 2021.

- Ben Taieb, S., Bontempi, G., Atiya, A. F., and Sornajama, A. A review and comparison of strategies for multi-step ahead time series forecasting based on the nn5 forecasting competition. *Expert Systems with Applications*, 39(8):7067–7083, 2012. ISSN 0957-4174. doi: <https://doi.org/10.1016/j.eswa.2012.01.039>. URL <https://www.sciencedirect.com/science/article/pii/S0957417412000528>.
- Besbes, O., Gur, Y., and Zeevi, A. Non-stationary stochastic optimization. *Operations research*, 63(5):1227–1244, 2015.
- Bhatnagar, A., Kassianik, P., Liu, C., Lan, T., Yang, W., Cassius, R., Sahoo, D., Arpit, D., Subramanian, S., Woo, G., Saha, A., Jagota, A. K., Gopalakrishnan, G., Singh, M., Krithika, K. C., Maddineni, S., Cho, D., Zong, B., Zhou, Y., Xiong, C., Savarese, S., Hoi, S., and Wang, H. Merlion: A machine learning library for time series. 2021. URL <https://arxiv.org/abs/2109.09265>.
- Candès, E. J., Lei, L., and Ren, Z. Conformalized survival analysis, 2021. URL <https://arxiv.org/abs/2103.09763>.
- Cauchois, M., Gupta, S., Ali, A., and Duchi, J. C. Robust validation: Confident predictions even when distributions shift, 2020. URL <https://arxiv.org/abs/2008.04267>.
- Cauchois, M., Gupta, S., and Duchi, J. C. Knowing what you know: Valid and validated confidence sets in multi-class and multilabel prediction. *J. Mach. Learn. Res.*, 22(1), jul 2022. ISSN 1532-4435.
- Chernozhukov, V., Wüthrich, K., and Yinchu, Z. Exact and robust conformal inference methods for predictive machine learning with dependent data. In Bubeck, S., Perchet, V., and Rigollet, P. (eds.), *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pp. 732–749. PMLR, 06–09 Jul 2018. URL <https://proceedings.mlr.press/v75/chernozhukov18a.html>.
- Daniely, A., Gonen, A., and Shalev-Shwartz, S. Strongly adaptive online learning. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1405–1411, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/daniely15.html>.
- Dashevskiy, M. and Luo, Z. Network traffic demand prediction with confidence. In *IEEE GLOBECOM 2008 - 2008 IEEE Global Telecommunications Conference*, pp. 1–5, 2008. doi: 10.1109/GLOCOM.2008.ECP.284.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Elsayed, S., Thyssens, D., Rashed, A., Jomaa, H. S., and Schmidt-Thieme, L. Do we really need deep learning models for time series forecasting?, 2021. URL <https://arxiv.org/abs/2101.02118>.
- Fannjiang, C., Bates, S., Angelopoulos, A., Listgarten, J., and Jordan, M. I. Conformal prediction for the design problem. *arXiv preprint arXiv:2202.03613*, 2022.
- Feldman, S., Ringel, L., Bates, S., and Romano, Y. Risk control for online learning models, 2022. URL <https://arxiv.org/abs/2205.09095>.
- Gibbs, I. and Candès, E. Adaptive conformal inference under distribution shift. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=6vaActvpcp3>.
- Gibbs, I. and Candès, E. Conformal inference for online prediction with arbitrary distribution shifts, 2022. URL <https://arxiv.org/abs/2208.08401>.
- Gupta, C., Kuchibhotla, A. K., and Ramdas, A. K. Nested conformal prediction and quantile out-of-bag ensemble methods. *arXiv preprint arXiv:1910.10562*, 2019.
- Hazan, E. *Introduction to Online Convex Optimization*. MIT Press, Cambridge, MA, USA, 2022. ISBN 9780262046985.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016. doi: 10.1109/CVPR.2016.90.
- Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HJz6tiCqYm>.
- Hendrycks, D., Mazeika, M., Wilson, D., and Gimpel, K. Using trusted data to train deep networks on labels corrupted by severe noise. *Advances in neural information processing systems*, 31, 2018.
- Jun, K.-S., Orabona, F., Wright, S., and Willett, R. Improved Strongly Adaptive Online Learning using Coin Betting. In Singh, A. and Zhu, J. (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 943–951. PMLR, 20–22 Apr 2017.

- URL <https://proceedings.mlr.press/v54/jun17a.html>.
- Kath, C. and Ziel, F. Conformal prediction interval estimation and applications to day-ahead and intraday power markets. *International Journal of Forecasting*, 37(2):777–799, 2021. ISSN 0169-2070. doi: <https://doi.org/10.1016/j.ijforecast.2020.09.006>. URL <https://www.sciencedirect.com/science/article/pii/S0169207020301473>.
- Koenker, R. and Bassett Jr, G. Regression quantiles. *Econometrica: journal of the Econometric Society*, pp. 33–50, 1978.
- Kwiatkowski, D., Phillips, P. C., Schmidt, P., and Shin, Y. Testing the null hypothesis of stationarity against the alternative of a unit root: How sure are we that economic time series have a unit root? *Journal of Econometrics*, 54(1):159–178, 1992. ISSN 0304-4076. doi: [https://doi.org/10.1016/0304-4076\(92\)90104-Y](https://doi.org/10.1016/0304-4076(92)90104-Y). URL <https://www.sciencedirect.com/science/article/pii/030440769290104Y>.
- Le, Y. and Yang, X. S. Tiny imagenet visual recognition challenge. 2015.
- Lei, J. and Wasserman, L. Distribution-free prediction bands for non-parametric regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):71–96, 2014. doi: <https://doi.org/10.1111/rssb.12021>. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/rssb.12021>.
- Lei, J., Robins, J., and Wasserman, L. Distribution-free prediction sets. *Journal of the American Statistical Association*, 108(501):278–287, 2013. doi: 10.1080/01621459.2012.751873. URL <https://doi.org/10.1080/01621459.2012.751873>. PMID: 25237208.
- Makridakis, S., Spiliotis, E., and Assimakopoulos, V. The m4 competition: Results, findings, conclusion and way forward. *International Journal of Forecasting*, 34(4): 802–808, 2018.
- Orabona, F. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2019.
- Orabona, F. and Pál, D. Coin betting and parameter-free online learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- Orabona, F. and Pál, D. Scale-free online learning. *Theoretical Computer Science*, 716:50–69, 2018. ISSN 0304-3975. doi: <https://doi.org/10.1016/j.tcs.2017.11.021>. URL <https://www.sciencedirect.com/science/article/pii/S0304397517308514>. Special Issue on ALT 2015.
- Papadopoulos, H. Inductive conformal prediction: Theory and application to neural networks. In Fritzsche, P. (ed.), *Tools in Artificial Intelligence*, chapter 18. IntechOpen, Rijeka, 2008. doi: 10.5772/6078. URL <https://doi.org/10.5772/6078>.
- Park, S., Bastani, O., Matni, N., and Lee, I. Pac confidence sets for deep neural networks via calibrated prediction. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=BJxVI04YvB>.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Pearce, T., Brintrup, A., Zaki, M., and Neely, A. High-quality prediction intervals for deep learning: A distribution-free, ensembled approach. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4075–4084. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/pearcel8a.html>.
- Podkopaev, A. and Ramdas, A. Distribution-free uncertainty quantification for classification under label shift. In de Campos, C. and Maathuis, M. H. (eds.), *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161 of *Proceedings of Machine Learning Research*, pp. 844–853. PMLR, 27–30 Jul 2021. URL <https://proceedings.mlr.press/v161/podkopaev21a.html>.
- Romano, Y., Patterson, E., and Candes, E. Conformalized quantile regression. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/5103c3584b063c431bd1268e9b5e76fb-Paper.pdf>.
- Romano, Y., Sesia, M., and Candès, E. J. Classification with valid and adaptive coverage. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

- Shafer, G. and Vovk, V. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008.
- Sousa, M., Tomé, A. M., and Moreira, J. A general framework for multi-step ahead adaptive conformal heteroscedastic time series forecasting, 2022. URL <https://arxiv.org/abs/2207.14219>.
- Stankeviciute, K., M. Alaa, A., and van der Schaar, M. Conformal time-series forecasting. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 6216–6228. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/312f1ba2a72318edaaa995a67835fad5-Paper.pdf>.
- Steinwart, I. and Christmann, A. Estimating conditional quantiles with the help of the pinball loss. *Bernoulli*, 17(1):211 – 225, 2011. doi: 10.3150/10-BEJ267. URL <https://doi.org/10.3150/10-BEJ267>.
- Stutz, D., Dvijotham, K. D., Cemgil, A. T., and Doucet, A. Learning optimal conformal classifiers. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=t80-4LKFVx>.
- Sun, S. and Yu, R. Copula conformal prediction for multi-step time series forecasting, 2022. URL <https://arxiv.org/abs/2212.03281>.
- Taylor, S. J. and Letham, B. Forecasting at scale. *PeerJ Preprints*, 5(e3190v2), Sept 2017. doi: 10.7287/peerj.preprints.3190v2.
- Tibshirani, R. J., Foygel Barber, R., Candès, E., and Ramdas, A. Conformal prediction under covariate shift. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/8fb21ee7a2207526da55a679f0332de2-Paper.pdf>.
- Vovk, V. Conditional validity of inductive conformal predictors. In Hoi, S. C. H. and Buntine, W. (eds.), *Proceedings of the Asian Conference on Machine Learning*, volume 25 of *Proceedings of Machine Learning Research*, pp. 475–490, Singapore Management University, Singapore, 04–06 Nov 2012. PMLR. URL <https://proceedings.mlr.press/v25/vovk12.html>.
- Vovk, V., Gammerman, A., and Saunders, C. Machine-learning applications of algorithmic randomness. In *Proceedings of the Sixteenth International Conference on Machine Learning*, ICML ’99, pp. 444–453, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc. ISBN 1558606122.
- Vovk, V., Gammerman, A., and Shafer, G. *Algorithmic Learning in a Random World*. Springer-Verlag, Berlin, Heidelberg, 2005. ISBN 0387001522.
- Vovk, V., Nouretdinov, I., Manokhin, V., and Gammerman, A. Cross-conformal predictive distributions. In Gammerman, A., Vovk, V., Luo, Z., Smirnov, E., and Peeters, R. (eds.), *Proceedings of the Seventh Workshop on Conformal and Probabilistic Prediction and Applications*, volume 91 of *Proceedings of Machine Learning Research*, pp. 37–51. PMLR, 11–13 Jun 2018. URL <https://proceedings.mlr.press/v91/vovk18a.html>.
- Wisniewski, W., Lindsay, D., and Lindsay, S. Application of conformal prediction interval estimations to market makers’ net positions. In Gammerman, A., Vovk, V., Luo, Z., Smirnov, E., and Cherubin, G. (eds.), *Proceedings of the Ninth Symposium on Conformal and Probabilistic Prediction and Applications*, volume 128 of *Proceedings of Machine Learning Research*, pp. 285–301. PMLR, 09–11 Sep 2020. URL <https://proceedings.mlr.press/v128/wisniewski20a.html>.
- Xu, C. and Xie, Y. Conformal prediction interval for dynamic time-series. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11559–11569. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/xu21h.html>.
- Yang, Y. and Kuchibhotla, A. K. Finite-sample efficient conformal prediction, 2021. URL <https://arxiv.org/abs/2104.13871>.
- Yang, Y., Kuchibhotla, A. K., and Tchetgen, E. T. Doubly robust calibration of prediction sets under covariate shift, 2022. URL <https://arxiv.org/abs/2203.01761>.
- Zaffran, M., Feron, O., Goude, Y., Josse, J., and Dieuleveut, A. Adaptive conformal predictions for time series. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 25834–25866. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/zaffran22a.html>.
- Zhang, L., Yang, T., rong jin, and Zhou, Z.-H. Dynamic regret of strongly adaptive methods. In Dy, J. and Krause,

A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 5882–5891. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/zhang18o.html>.

Zhao, P. and Zhang, L. Improved analysis for dynamic regret of strongly convex and smooth functions. In *Learning for Dynamics and Control*, pp. 48–59. PMLR, 2021.

Zhao, P., Xie, Y.-F., Zhang, L., and Zhou, Z.-H. Efficient methods for non-stationary online learning. *Advances in Neural Information Processing Systems*, 35:11573–11585, 2022.

Zinkevich, M. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning, ICML'03*, pp. 928–935. AAAI Press, 2003. ISBN 1577351894.

## A. Basic Properties of Online Conformal Prediction Algorithms

### A.1. Properties of SF-OGD

We consider the SF-OGD algorithm (Algorithm 2). We first show that the iterates of SF-OGD are bounded within a range slightly larger than the range of the true radii. The proof is similar to Gibbs & Candès (2021, Lemma 4.1).

**Lemma A.1** (Bounded iterates for SF-OGD). *Suppose the true radii are bounded:  $S_t \in [0, D]$  for all  $t \in [T]$ . Then Algorithm 2 with any initialization  $\hat{s}_1 \in [-\eta, D + \eta]$  and learning rate  $\eta > 0$  admits bounded iterates:*

$$\hat{s}_t \in [-\eta, D + \eta] \text{ for all } t \in [T].$$

*Proof.* Recall by (2) that

$$\nabla \ell^{(t)}(\hat{s}_t) = \alpha - \mathbb{1}[\hat{s}_t < S_t] = \alpha - \text{err}_t \in \{-(1 - \alpha), \alpha\} \subset [-1, 1]. \quad (16)$$

for all  $t \in [T]$ . Therefore, Algorithm 2 satisfies for any  $t \geq 1$  that

$$|\hat{s}_{t+1} - \hat{s}_t| = \eta \left| \frac{\alpha - \text{err}_t}{\sqrt{\sum_{\tau=1}^t (\alpha - \text{err}_\tau)^2}} \right| \leq \eta. \quad (17)$$

We prove the lemma by contradiction. Suppose there exists some  $t$  such that  $\hat{s}_t \notin [-\eta, D + \eta]$ . Let  $t \geq 2$  be the smallest such time index (abusing notation slightly). Suppose  $\hat{s}_t > D + \eta$ , then by (17) we must have  $\hat{s}_{t-1} > D$  but  $\hat{s}_{t-1} \leq D + \eta$ . Note that  $\hat{s}_{t-1} > D \geq S_{t-1}$  by our precondition, so that the  $(t - 1)$ -th prediction set must cover and thus  $\text{err}_{t-1} = 0$ . Therefore by the algorithm update (10) we have

$$\hat{s}_t = \hat{s}_{t-1} - \eta \frac{\alpha - \text{err}_{t-1}}{\sqrt{\sum_{\tau=1}^{t-1} (\alpha - \text{err}_\tau)^2}} < \hat{s}_{t-1} \leq D + \eta,$$

contradicting with our assumption that  $\hat{s}_t > D + \eta$ . A similar contradiction can be derived for the other case where  $\hat{s}_t < -\eta$ . This proves the desired result.  $\square$

The following regret bound follows directly by applying the generic regret bound of Scale-Free OGD (Orabona & Pál, 2018, Theorem 2) to the quantile loss (1).

**Proposition A.2** (Anytime regret bound for SF-OGD). *Suppose the true radii are bounded:  $S_t \in [0, D]$  for all  $t \in [T]$ . Then Algorithm 2 with any initialization  $\hat{s}_1 \in [0, D]$  and learning rate  $\eta = D/\sqrt{3}$  achieves the following regret bound for any  $t \in [T]$ :*

$$\text{Reg}(t) \leq (\sqrt{3} + 1)D \sqrt{\sum_{\tau=1}^t \|\nabla \ell^{(\tau)}(\hat{s}_\tau)\|_2^2} \leq \mathcal{O}(D\sqrt{t}).$$

*Proof.* The second inequality follows directly by (16).

To prove the first inequality (the regret bound), we note that Algorithm 2 is a special case of the Scale-Free Mirror Descent algorithm of Orabona & Pál (2018, Section 4) with convex loss  $\ell^{(t)}(\cdot) = \ell_{1-\alpha}(S_t, \cdot)$ , and regularizer  $R(s) := s^2/(2\eta)$  (in their notation) which is  $\lambda = 1/\eta$ -strongly convex with respect to the  $\ell_2$  norm on  $\mathbb{R}$ . Further, by Lemma A.1 we have  $\hat{s}_t \in [-\eta, D + \eta]$  for all  $t \in [T]$ . Therefore, applying Orabona & Pál (2018, Theorem 2) gives that for any  $t \in [T]$ ,

$$\begin{aligned} & \sum_{\tau=1}^t \ell^{(\tau)}(\hat{s}_\tau) - \inf_{s^* \in [0, D]} \sum_{\tau=1}^t \ell^{(\tau)}(s^*) \leq \left( \frac{1}{\lambda} + \sup_{\tau \geq 1} B_R(s^*, \hat{s}_\tau) \right) \cdot \sqrt{\sum_{\tau=1}^t \|\nabla \ell^{(\tau)}(\hat{s}_\tau)\|_2^2} \\ & = \left( \eta + \sup_{s^* \in [0, D], s' \in [-\eta, D + \eta]} \frac{1}{2\eta} (s^* - s')^2 \right) \cdot \sqrt{\sum_{\tau=1}^t \|\nabla \ell^{(\tau)}(\hat{s}_\tau)\|_2^2} \end{aligned}$$

$$= \left( \eta + \frac{(D + \eta)^2}{2\eta} \right) \cdot \sqrt{\sum_{\tau=1}^t \|\nabla \ell^{(\tau)}(\hat{s}_\tau)\|_2^2},$$

where  $B_R(\cdot, \cdot)$  denotes the Bregman divergence associated with  $R$ . Choosing  $\eta = D/\sqrt{3}$ , the leading coefficient is  $(\sqrt{3} + 1)D$ . The desired result follows by noting that

$$\text{Reg}(t) = \sum_{\tau=1}^t \ell^{(\tau)}(\hat{s}_\tau) - \inf_{s^* \in \mathbb{R}} \sum_{\tau=1}^t \ell^{(\tau)}(s^*) = \sum_{\tau=1}^t \ell^{(\tau)}(\hat{s}_\tau) - \inf_{s^* \in [0, D]} \sum_{\tau=1}^t \ell^{(\tau)}(s^*)$$

by our assumption that  $S_t \in [0, D]$  and basic properties of the quantile losses  $\{\ell^{(\tau)}(\cdot)\}_{\tau \geq 1}$ .  $\square$

## A.2. Example of ‘‘Trivial’’ Algorithm with Coverage Guarantee

We consider the following ‘‘trivial’’ online conformal prediction algorithm that does not utilize the data at all: Simply predict the maximum radius  $D$  for  $(1 - \alpha)$  proportion of the time steps, then predict the minimum radius 0 for  $\alpha$  proportion of the time steps:

$$\begin{cases} \hat{s}_t := D & \text{for } t \in \{1, \dots, \lfloor (1 - \alpha)T \rfloor\} =: T_{\text{full}}, \\ \hat{s}_t := 0 & \text{for } t \in \{\lfloor (1 - \alpha)T \rfloor + 1, \dots, T\} =: T_{\text{empty}}. \end{cases} \quad (18)$$

By our assumption that  $S_t \in [0, D]$  almost surely and the nested set structure of  $\hat{C}_t(X_t, \cdot)$ , we have  $\text{err}_t = \mathbb{1}[Y_t \in \hat{C}_t] = \mathbb{1}[\hat{s}_t \geq S_t] = 0$  for all  $t \in T_{\text{full}}$ , and similarly  $\text{err}_t = 1$  for all  $t \in T_{\text{empty}}$ . Therefore, algorithm (18) directly satisfies

$$\text{CovErr}(T) = \left| \frac{1}{T} \sum_{t=1}^T \text{err}_t - \alpha \right| = \left| \frac{|T_{\text{empty}}|}{T} - \alpha \right| = \left| \frac{T - \lfloor (1 - \alpha)T \rfloor}{T} - \alpha \right| \leq \frac{1}{T}, \quad (19)$$

i.e. the algorithm achieves approximately  $(1 - \alpha)$  empirical coverage, with error  $O(1/T)$ . It is also straightforward to see that, by slightly modifying the definition of  $T_{\text{full}}, T_{\text{empty}}$  (making the two index sets alternate), we can make the above coverage bound hold in an anytime sense (for  $t \in [T]$ ).

However, it is straightforward to construct examples of data distributions for which the trivial algorithm (18) suffers linear regret on the quantile loss  $\ell^{(t)}$  defined in (1), and such data distributions can be chosen to be fairly simple. For example, suppose all data points admit the same true radius  $D/2$ , i.e.

$$S_t \equiv D/2 \quad \text{for all } t \in [T].$$

Then for  $s^* = D/2$  we have  $\ell^{(t)}(s^*) = \ell_{1-\alpha}(D/2, D/2) = 0$  for all  $t \in [T]$ , which achieves total loss  $\sum_{t=1}^T \ell^{(t)}(s^*) = 0$  (the smallest possible, since  $\ell^{(t)}(\cdot) \geq 0$ ). On the other hand, algorithm (18) achieves loss

$$\ell^{(t)}(\hat{s}_t) = \ell_{1-\alpha}(S_t, \hat{s}_t) = \begin{cases} \ell_{1-\alpha}(D/2, D) = \alpha(D/2) & \text{for } t \in T_{\text{full}}, \\ \ell_{1-\alpha}(D/2, 0) = (1 - \alpha)(D/2) & \text{for } t \in T_{\text{empty}}. \end{cases}$$

Therefore, we have

$$\begin{aligned} \text{Reg}(T) &= \sum_{t=1}^T \ell^{(t)}(\hat{s}_t) - \inf_{s^*} \sum_{t=1}^T \ell^{(t)}(s^*) = \sum_{t=1}^T \ell^{(t)}(\hat{s}_t) = \alpha D/2 \cdot |T_{\text{full}}| + (1 - \alpha)D/2 \cdot |T_{\text{empty}}| \\ &= \alpha D/2 \cdot \lfloor (1 - \alpha)T \rfloor + (1 - \alpha)D/2 \cdot (T - \lfloor (1 - \alpha)T \rfloor) \geq \alpha(1 - \alpha)DT = \Omega(T), \end{aligned}$$

i.e. algorithm (18) suffers from linear regret. This demonstrates sublinear regret as a sensible criterion for ruling out trivial algorithms like (18).

## B. Proofs for Section 4

### B.1. Proof of Proposition 4.1

The proof follows by plugging in the regret bound for SF-OGD (Proposition A.2) into Jun et al. (2017, Theorem 2). Define  $u(t) := \max_n \{2^n : t \equiv 0 \pmod{2^n}\}$ . Fix any  $k \in [T]$  and  $\tau \in [T - k + 1]$ . Their proof starts by splitting the interval  $[\tau, \tau + k - 1]$  into consecutive sub-intervals  $\bar{J}^{(1)}, \dots, \bar{J}^{(n)}$ , where  $\bar{J}^{(i)} = [\tau_i, \max\{\tau + k, \tau_i + u(\tau_i)\} - 1]$  is a prefix of expert  $\mathcal{A}_{\tau_i}$ 's active interval.

We have for any fixed  $s^* \in \mathbb{R}$  that

$$\begin{aligned}
 \text{Regret}_\tau^k(s^*) &:= \sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(\hat{s}_t) - \sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(s^*) \\
 &= \sum_{i=1}^n \sum_{t \in \bar{J}^{(i)}} \left( \ell^{(t)}(\hat{s}_t) - \ell^{(t)}(\hat{s}_{t, \tau_i}) \right) + \sum_{i=1}^n \sum_{t \in \bar{J}^{(i)}} \left( \ell^{(t)}(\hat{s}_{t, \tau_i}) - \ell^{(t)}(s^*) \right) \\
 &\leq D \underbrace{\sum_{i=1}^n \sqrt{|\bar{J}^{(i)}| (7 \log T + 5)}}_{\text{Jun et al. (2017, Lemma 2)}} + D \underbrace{\sum_{i=1}^n \sqrt{|\bar{J}^{(i)}| (1 + \sqrt{3})}}_{\text{Proposition A.2}} \\
 &\leq D \left( \sqrt{7 \log T + 5} + \sqrt{1 + \sqrt{3}} \right) \underbrace{\sum_{j=0}^{\infty} \sqrt{k 2^{-j}}}_{\text{Jun et al. (2017, Lemma 3)}} \\
 &= \frac{D\sqrt{2}}{\sqrt{2}-1} \sqrt{k} \left( \sqrt{7 \log T + 5} + \sqrt{1 + \sqrt{3}} \right) \leq 15D \sqrt{k(\log T + 1)}
 \end{aligned}$$

Taking supremum over all  $s^* \in \mathbb{R}$  and all intervals  $[\tau, \tau + k - 1] \subset [T]$ , we obtain the desired bound on  $\text{SAReg}(T, k)$ .  $\square$

### B.2. Dynamic Regret for SAOCP

**Proposition B.1** (Dynamic regret bound for SAOCP). *Algorithm 1 achieves the following (worst-case) dynamic regret bound: For any interval  $[\tau, \tau + k - 1] \subset [T]$  of length  $k \in [T]$ , we have*

$$\sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(\hat{s}_t) - \min_{s^*_{\tau, \tau+k-1}} \sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(s_t^*) = \sum_{t=\tau}^{\tau+k-1} \left[ \ell^{(t)}(\hat{s}_t) - \ell^{(t)}(S_t) \right] \leq \tilde{\mathcal{O}} \left( D \left[ V_{[\tau, \tau+k-1]}^{1/3} k^{2/3} + \sqrt{k} \right] \right), \quad (20)$$

where

$$V_{[\tau, \tau+k-1]} := \sum_{t=\tau+1}^{\tau+k-1} |S_t - S_{t-1}|$$

is the path length of the true radii within  $[\tau, \tau + k - 1]$ .

The dynamic regret<sup>4</sup>  $\tilde{\mathcal{O}}(V_{[\tau, \tau+k-1]}^{1/3} k^{2/3} + \sqrt{k})$  obtained in Proposition B.1 matches minimax optimal rate (Besbes et al., 2015) for general online convex optimization problems under gradient feedback; though we remark that faster rates are achievable under full-information feedback of the entire loss function (Zhao & Zhang, 2021, Theorem 3).

**Comparison of dynamic regret with FOCI** Dividing (20) by  $k$ , we obtain the following average dynamic regret bound for SAOCP on  $[\tau, \tau + k - 1]$ :

$$\tilde{\mathcal{O}} \left( D \left[ (V_{[\tau, \tau+k-1]}/k)^{1/3} + 1/\sqrt{k} \right] \right),$$

simultaneously for all lengths  $k$  and  $\tau \in [T - k + 1]$ .

<sup>4</sup>More precisely, the intermediate result with  $V_{[\tau, \tau+k-1]}$  replaced by the standard total variation of losses  $\tilde{V}_{[\tau, \tau+k-1]}$  in (21).

In comparison, the FACI algorithm (adapted to our setting) with learning rate  $\eta$  achieves average dynamic regret bound (Gibbs & Candès, 2022, Theorem 3.2)

$$\tilde{\mathcal{O}}\left(D\left[\left(V_{[\tau, \tau+k-1]}/k\right)^{1/2} + \eta/D + D/(\eta k)\right]\right).$$

When the path length  $V_{[\tau, \tau+k-1]} = o(k)$ , FACI achieves a better dependence on the average path length ( $V_{[\tau, \tau+k-1]}/k = o(1)$ ), yet a worse dependence on  $k$  itself due to the inability to choose the optimal  $\eta$  simultaneously for all  $k$ , similar as the comparison of their SARegret bounds (Section 4.1).

**Proof of Proposition B.1** We apply the dynamic regret bound of Zhang et al. (2018, Corollary 5) for the SAOCP algorithm on the interval  $[\tau, \tau + k - 1]$ , and note that our iterates  $\hat{s}_t \in [-\eta, D + \eta] \subset [-D, 2D]$  by Lemma A.1 and our choice  $\eta = D/\sqrt{3}$  in Algorithm 1. Therefore we obtain

$$\sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(\hat{s}_t) - \min_{s_{\tau:\tau+k-1}^*} \sum_{t=\tau}^{\tau+k-1} \ell^{(t)}(s_t^*) \leq \tilde{\mathcal{O}}\left(D\left[\tilde{V}_{[\tau, \tau+k-1]}^{1/3} k^{2/3} + \sqrt{k}\right]\right),$$

where

$$\tilde{V}_{[\tau, \tau+k-1]} = \sum_{t=\tau+1}^{\tau+k-1} \sup_{s' \in [0, D]} \left| \ell^{(t)}(s') - \ell^{(t-1)}(s') \right| \stackrel{(i)}{\leq} \sum_{t=\tau+1}^{\tau+k-1} |S_t - S_{t-1}| = V_{[\tau, \tau+k-1]}, \quad (21)$$

where (i) follows by the fact that  $|\ell^{(t)}(s') - \ell^{(t-1)}(s')| = |\ell_{1-\alpha}(s', S_t) - \ell_{1-\alpha}(s', S_{t-1})| \leq |S_t - S_{t-1}|$  by the 1-Lipschitzness of the quantile loss (1) with respect to the second argument. This proves the desired result.  $\square$

### B.3. Proof of Theorem 4.2

We first note that, by (2) and (10), Algorithm 2 simplifies to the update

$$\hat{s}_{t+1} = \hat{s}_t + \eta \frac{\text{err}_t - \alpha}{\sqrt{\sum_{s=1}^t (\text{err}_s - \alpha)^2}} = \hat{s}_1 + \eta \sum_{s=1}^t \frac{\text{err}_s - \alpha}{\sqrt{\sum_{i=1}^s (\text{err}_i - \alpha)^2}}. \quad (22)$$

Note that we have  $\hat{s}_{t+1} \in [-\eta, D + \eta]$  for all  $t \geq 0$  (Lemma A.1), which implies that

$$\left| \sum_{t=t_0+1}^{t_f} \frac{\text{err}_t - \alpha}{\sqrt{\sum_{s=1}^t (\text{err}_s - \alpha)^2}} \right| = \frac{1}{\eta} |\hat{s}_{t_f+1} - \hat{s}_{t_0+1}| \leq \frac{D + 2\eta}{\eta} \quad \text{for any } 0 \leq t_0 < t_f.$$

Note that  $|\text{err}_t - \alpha| \in [\alpha, 1]$  for all  $t$ . Therefore, we can invoke Lemma B.2 below with  $a_t = \text{err}_t - \alpha$  and  $M = (D + 2\eta)/\eta$  to obtain that for any  $T \geq 1$ ,

$$\left| \frac{1}{T} \sum_{t=1}^T \text{err}_t - \alpha \right| \leq 2 \left( \frac{D + 3\eta}{\eta} + \alpha^{-2} \log T \right) T^{-1/4} \leq \mathcal{O}(\alpha^{-2} T^{-1/4} \log T),$$

where the later bound holds for any  $\eta = \Theta(D)$ . This proves Theorem 4.2.  $\square$

**Lemma B.2.** Suppose the sequence  $\{a_t\}_{t \in [T]} \in \mathbb{R}$  satisfies  $\alpha \leq |a_t| \leq 1$  for some  $\alpha > 0$ , and

$$\left| \sum_{t=t_0+1}^{t_f} \frac{a_t}{\sqrt{\sum_{s=1}^t a_s^2}} \right| \leq M \quad \text{for any } 0 \leq t_0 < t_f \leq T.$$

Then we have

$$\left| \frac{1}{T} \sum_{t=1}^T a_t \right| \leq 2(M + 1 + \alpha^{-2} \log T) T^{-1/4}.$$

*Proof.* The proof builds on a grouping argument. Define integers

$$L = \lceil T^\beta \rceil, \quad K = \lceil T/L \rceil \leq T^{1-\beta} + 1,$$

where  $\beta \in (0, 1)$  is a parameter to be chosen. For any  $k \in [K]$ , define the  $k$ -th group to be

$$G_k = \{t_{k-1} + 1, \dots, t_k\} := \{(k-1)L + 1, \dots, \min\{kL, T\}\}, \quad (23)$$

so that we have  $\bigcup_{k=1}^K G_k = [T]$ ,  $|G_k| = L$  for all  $k \in [K-1]$ , and  $|G_K| \leq L$ .

Next, for any fixed  $k \geq 2$ , define sums

$$S_k := \sum_{t \in G_k} \frac{a_t}{\sqrt{\sum_{s=1}^t a_s^2}}, \quad \tilde{S}_k := \sum_{t \in G_k} \frac{a_t}{\sqrt{\sum_{s=1}^{t_{k-1}} a_s^2}}.$$

By our precondition, we have  $|S_k| \leq M$  for all  $k \in [K]$ . Further, we have

$$\begin{aligned} |S_k - \tilde{S}_k| &\leq \sum_{t \in G_k} |a_t| \cdot \left( \frac{1}{\sqrt{\sum_{s=1}^{t_{k-1}} a_s^2}} - \frac{1}{\sqrt{\sum_{s=1}^t a_s^2}} \right) \leq |G_k| \cdot \left( \frac{1}{\sqrt{\sum_{s=1}^{t_{k-1}} a_s^2}} - \frac{1}{\sqrt{\sum_{s=1}^{t_k} a_s^2}} \right) \\ &\stackrel{(i)}{\leq} L \cdot \frac{\sum_{s=t_{k-1}+1}^{t_k} a_s^2}{2 \left( \sum_{s=1}^{t_{k-1}} a_s^2 \right)^{3/2}} \stackrel{(ii)}{\leq} L \cdot \frac{L}{2(\alpha^2(k-1)L) \cdot \left( \sum_{s=1}^{t_{k-1}} a_s^2 \right)^{1/2}} = \frac{L}{2\alpha^2(k-1) \cdot \sqrt{\sum_{s=1}^{t_{k-1}} a_s^2}}, \end{aligned}$$

where (i) uses the inequality  $\frac{1}{\sqrt{x}} - \frac{1}{\sqrt{x+y}} \leq \frac{y}{2x^{3/2}}$  for  $x, y \geq 0$ , and (ii) uses the bounds  $\sum_{s=t_{k-1}+1}^{t_k} a_s^2 \leq (t_k - t_{k-1}) \leq L$  and  $\sum_{s=1}^{t_{k-1}} a_s^2 \geq \alpha^2 t_{k-1} = \alpha^2(k-1)L$ . By the triangle inequality, this implies that

$$|\tilde{S}_k| \leq |S_k| + |\tilde{S}_k - S_k| \leq M + \frac{L}{2\alpha^2(k-1) \cdot \sqrt{\sum_{s=1}^{t_{k-1}} a_s^2}},$$

and thus for any  $k \geq 2$  that

$$\begin{aligned} \left| \sum_{t \in G_k} a_t \right| &= \underbrace{\left| \sum_{t \in G_k} \frac{a_t}{\sqrt{\sum_{s=1}^{t_{k-1}} a_s^2}} \right|}_{|\tilde{S}_k|} \cdot \sqrt{\sum_{s=1}^{t_{k-1}} a_s^2} \leq \left( M + \frac{L}{2\alpha^2(k-1) \cdot \sqrt{\sum_{s=1}^{t_{k-1}} a_s^2}} \right) \cdot \sqrt{\sum_{s=1}^{t_{k-1}} a_s^2} \\ &\leq M \sqrt{\sum_{s=1}^{t_{k-1}} a_s^2} + \frac{L}{2\alpha^2(k-1)} \leq M \sqrt{(k-1)L} + \frac{L}{2\alpha^2(k-1)}. \end{aligned}$$

For  $k = 1$ , we have trivially  $|\sum_{t \in G_1} a_t| \leq |G_1| \leq L$ . Summing the bounds over  $k \in [K]$  yields

$$\begin{aligned} \left| \sum_{t=1}^T a_t \right| &\leq L + \sum_{k=1}^K \left| \sum_{t \in G_k} a_t \right| \leq L + M\sqrt{L} \cdot \sum_{k=2}^K \sqrt{k-1} + \frac{L}{2\alpha^2} \sum_{k=2}^K \frac{1}{k-1} \\ &\leq L + \frac{2}{3} M\sqrt{L} K^{3/2} + \frac{L}{2\alpha^2} \log_2 K \\ &\leq \lceil T^\beta \rceil + \frac{2}{3} M \sqrt{\lceil T^\beta \rceil} \cdot T^{3(1-\beta)/2} + \frac{1}{2\alpha^2} \lceil T^\beta \rceil \log_2(T^{1-\beta}) \\ &\leq 2T^\beta + 2MT^{3/2-\beta} + \frac{2}{\alpha^2} T^\beta \log T. \end{aligned}$$

Choosing  $\beta = 3/4$ , we obtain

$$\left| \sum_{t=t_0+1}^{t_f} a_t \right| \leq 2(M+1 + \log T/\alpha^2) T^{3/4}.$$

Dividing by  $T$  on both sides yields the desired result.  $\square$

#### B.4. Coverage of SAOCP

**Theorem B.3** (Coverage bound for SAOCP). *Consider a randomized version of Algorithm 1 where Line 7 is changed to sampling an expert  $i \sim p_{t, \cdot} \in \Delta([t])$  and outputting radius  $\widehat{s}_{t,i}$ . Consider the corresponding expected miscoverage error*

$$\widetilde{\text{err}}_t := \sum_{i=1}^t p_{t,i} \underbrace{\mathbb{1}[\widehat{s}_{t,i} < S_t]}_{:= \text{err}_{t,i}}. \quad (24)$$

Then we have for any  $T \geq 1$  that

$$\left| \frac{1}{T} \sum_{t=1}^T \widetilde{\text{err}}_t - \alpha \right| \leq \mathcal{O} \left( \inf_{\beta \in (1/2, 1)} \left\{ T^{1/2-\beta} + T^{\beta-1} \times \underbrace{\left( 1 + \sum_{j=2}^{\lceil T^{1-\beta} \rceil} \max_{t \in G_j} \sum_{i=1}^t \left| p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right| \right)}_{S_\beta(\{p_t\}_{t \in [T]}, \{\sum_{\tau=i}^t \|\nabla \ell^{(\tau)}(\widehat{s}_{\tau,i})\|_2^2\}_{i \leq t}) =: S_\beta(T)} \right\} \right)$$

(understanding  $p_{t_{j-1},i} := 0$  for any  $i > t_{j-1}$ ), where for each  $\beta \in (1/2, 1)$ ,  $\{G_j\}_{j=1}^{\lceil T^{1-\beta} \rceil}$  with  $|G_j| \leq \lceil T^\beta \rceil$ ,  $G_j = \{t_{j-1} + 1, \dots, \min\{t_j, T\}\}$  is the even grouping of  $[T]$  as in (23), and

$$G_{i:t}^i := \sqrt{\sum_{\tau=i}^t \|\nabla \ell^{(\tau)}(\widehat{s}_{\tau,i})\|_2^2} = \sqrt{\sum_{\tau=i}^t (\text{err}_{\tau,i} - \alpha)^2}$$

is the cumulative squared gradients received by expert  $\mathcal{A}_i$  for any  $t > i$  (understanding experts as running until time  $T$  even after they become inactive).

*Proof.* Fix any  $i \in [T]$ . As Algorithm 1 chooses each expert  $\mathcal{A}_i$  to be SF-OGD (Algorithm 2), we have by (10) that for all  $t \geq i$ ,

$$\widehat{s}_{t+1,i} - \widehat{s}_{t,i} = \frac{\eta}{G_{i:t}^i} \cdot (\text{err}_{t,i} - \alpha). \quad (25)$$

Now fix any  $\beta \in (1/2, 1)$ . For any group  $2 \leq j \leq \lceil T^{1-\beta} \rceil$  and  $t \in G_j$ , plugging the above into definition (24) gives that

$$\begin{aligned} \widetilde{\text{err}}_t - \alpha &= \sum_{i=1}^t p_{t,i} (\text{err}_{t,i} - \alpha) = \frac{1}{\eta} \sum_{i=1}^t p_{t,i} G_{i:t}^i (\widehat{s}_{t+1,i} - \widehat{s}_{t,i}) \\ &= \frac{1}{\eta} \sum_{i=1}^{t_{j-1}} p_{t_{j-1},i} G_{i:t_{j-1}}^i (\widehat{s}_{t+1,i} - \widehat{s}_{t,i}) + \frac{1}{\eta} \sum_{i=1}^t \left( p_{t,i} G_{i:t}^i - p_{t_{j-1},i} G_{i:t_{j-1}}^i \right) (\widehat{s}_{t+1,i} - \widehat{s}_{t,i}) \\ &= \frac{1}{\eta} \sum_{i=1}^{t_{j-1}} p_{t_{j-1},i} G_{i:t_{j-1}}^i (\widehat{s}_{t+1,i} - \widehat{s}_{t,i}) + \frac{1}{\eta} \sum_{i=1}^t \left( p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right) \cdot G_{i:t}^i (\widehat{s}_{t+1,i} - \widehat{s}_{t,i}). \end{aligned}$$

Summing this over  $t \in G_j$  and noting that the coefficients  $p_{t_{j-1},i} G_{i:t_{j-1}}^i$  in the first sum does not depend on  $t$ , we get

$$\begin{aligned} &\left| \sum_{t \in G_j} (\widetilde{\text{err}}_t - \alpha) \right| \\ &\leq \left| \frac{1}{\eta} \sum_{i=1}^{t_{j-1}} p_{t_{j-1},i} G_{i:t_{j-1}}^i (\widehat{s}_{t_j+1,i} - \widehat{s}_{t_{j-1}+1,i}) \right| + |G_j| \cdot \max_{t \in G_j} \left| \frac{1}{\eta} \sum_{i=1}^t \left( p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right) \cdot G_{i:t}^i (\widehat{s}_{t+1,i} - \widehat{s}_{t,i}) \right| \\ &\leq \frac{1}{\eta} \max_{i \in [t_{j-1}]} G_{i:t_{j-1}}^i |\widehat{s}_{t_j+1,i} - \widehat{s}_{t_{j-1}+1,i}| + |G_j| \cdot \max_{t \in G_j} \sum_{i=1}^t \left| \frac{1}{\eta} \left( p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right) \cdot G_{i:t}^i (\widehat{s}_{t+1,i} - \widehat{s}_{t,i}) \right| \end{aligned}$$

$$\begin{aligned} &\stackrel{(i)}{\leq} \frac{D + 2\eta}{\eta} \sqrt{T} + |G_j| \cdot \max_{t \in G_j} \sum_{i=1}^t \left| p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right| \\ &\stackrel{(ii)}{\leq} C\sqrt{T} + |G_j| \cdot \max_{t \in G_j} \sum_{i=1}^t \left| p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right| \end{aligned}$$

Above, (i) used  $G_{i:t_{j-1}}^i \leq \sqrt{t_{j-1} - i + 1} \leq \sqrt{T}$  by the definition of  $G_{i:t_{j-1}}^i$ , the bound  $|\widehat{s}_{t_{j+1},i} - \widehat{s}_{t_{j-1}+1,i}| \leq (D + 2\eta)$  which follows by the fact that each expert is initialized within  $[-\eta, D + \eta]$  and applying Lemma A.1, and the bound  $|G_{i:t}^i(\widehat{s}_{t+1,i} - \widehat{s}_{t,i})| \leq \eta$  by (25); (ii) used the fact that  $\eta = D/\sqrt{3}$  in Algorithm 1, so that  $(D + 2\eta)/\eta = 2 + \sqrt{3} =: C$  is an absolute constant. Also, note that for group  $j = 1$ , we directly have

$$\left| \sum_{t \in G_1} (\widetilde{\text{err}}_t - \alpha) \right| \leq |G_1|.$$

Summing all the above bounds over  $j \in \lceil [T^{1-\beta}] \rceil$  gives

$$\begin{aligned} &\left| \sum_{t=1}^T (\widetilde{\text{err}}_t - \alpha) \right| \leq \sum_{j=1}^{\lceil T^{1-\beta} \rceil} \left| \sum_{t \in G_j} (\widetilde{\text{err}}_t - \alpha) \right| \\ &\leq \mathcal{O} \left( T^{3/2-\beta} + |G_1| + \sum_{j=2}^{\lceil T^{1-\beta} \rceil} |G_j| \times \max_{t \in G_j} \sum_{i=1}^t \left| p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right| \right) \\ &\leq \mathcal{O} \left( T^{3/2-\beta} + T^\beta \left( 1 + \sum_{j=2}^{\lceil T^{1-\beta} \rceil} \max_{t \in G_j} \sum_{i=1}^t \left| p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right| \right) \right) \end{aligned}$$

Dividing both sides by  $T$  proves the desired bound for this fixed  $\beta$ . Further taking supremum over  $\beta \in (1/2, 1)$  gives the desired result.  $\square$

#### B.4.1. DISCUSSIONS & SF-OGD AS A SPECIAL CASE

We first note that, the proof of Theorem B.3 does not rely on the specific structure of either the expert weights  $\{p_{t,i}\}_{i < t}$  or the active intervals. Therefore, the result of Theorem B.3 holds generically for any other aggregation scheme over experts with arbitrary active intervals, in addition to that specified in Algorithm 1.

In particular, by setting  $p_{t,1} = 1$  and  $p_{t,i} = 0$  for  $i \geq 2$ , and defining the first expert  $\mathcal{A}_1$  to be active over  $[T]$ , Algorithm 1 (either with or without the randomization, since there is only one active expert) recovers Algorithm 2. In this case, we show that  $S_\beta(T) \leq \tilde{\mathcal{O}}(\alpha^{-2})$  for any  $\beta \in (1/2, 1)$ , so that Theorem B.3 (and its informal version in Theorem 4.3) indeed subsumes Theorem 4.2 as a special case by choosing  $\beta = 3/4$ , as claimed in Section 4.2.

We have

$$S_\beta(T) = 1 + \sum_{j=2}^{\lceil T^{1-\beta} \rceil} \max_{t \in G_j} \sum_{i=1}^t \left| p_{t,i} - p_{t_{j-1},i} \frac{G_{i:t_{j-1}}^i}{G_{i:t}^i} \right| \stackrel{(i)}{=} 1 + \sum_{j=2}^{\lceil T^{1-\beta} \rceil} \max_{t \in G_j} \left| 1 - \frac{G_{1:t_{j-1}}^1}{G_{1:t}^1} \right|, \quad (26)$$

where (i) used the fact that  $p_{t,1} = 1$  and  $p_{t,i} = 0$  for  $i \geq 2$ . For any  $t \in G_j$ , we have

$$\left| 1 - \frac{G_{1:t_{j-1}}^1}{G_{1:t}^1} \right| = 1 - \frac{\sqrt{\sum_{s=1}^{t_{j-1}} (\text{err}_s - \alpha)^2}}{\sqrt{\sum_{s=1}^t (\text{err}_s - \alpha)^2}} \stackrel{(i)}{\leq} \frac{\sum_{s=t_{j-1}+1}^t (\text{err}_s - \alpha)^2}{2 \sum_{s=1}^{t_{j-1}} (\text{err}_s - \alpha)^2} \stackrel{(ii)}{\leq} \frac{t - t_{j-1}}{2\alpha^2 t_{j-1}} \stackrel{(iii)}{\leq} \frac{\lceil T^\beta \rceil}{2\alpha^2 \cdot (j-1) \lceil T^\beta \rceil} = \frac{1}{2\alpha^2(j-1)},$$

where (i) follows from the inequality  $1 - \frac{\sqrt{x}}{\sqrt{x+y}} = \frac{\sqrt{x+y} - \sqrt{x}}{\sqrt{x+y}} \leq \frac{\sqrt{x+y} - \sqrt{x}}{\sqrt{x}} = \sqrt{1 + \frac{y}{x}} - 1 \leq \frac{y}{2x}$  for any  $x, y \geq 0$ ; (ii) follows by the bound  $|\text{err}_s - \alpha| \in [\alpha, 1]$  for any  $s$ ; (iii) follows by definition of the grouping (23). Plugging the above bound

into (26), we obtain

$$S_\beta(T) \leq 1 + \sum_{j=2}^{\lceil T^{1-\beta} \rceil} \frac{1}{2\alpha^2(j-1)} \leq \mathcal{O}(\alpha^{-2} \log T) = \tilde{\mathcal{O}}(\alpha^{-2}),$$

proving the claim.

### C. Distribution-Aware Coverage Guarantees for SAOCP

In this section, we show that under mild density lower bound assumptions on the true radii, a probabilistic variant of the coverage error of SAOCP (Algorithm 1) is bounded by  $\tilde{\mathcal{O}}(k^{-1/(2q)}) + \tilde{\mathcal{O}}((\text{Var}_k/k)^{1/q})$  for every interval of length  $k$ , where  $q \geq 2$  is a parameter of the density lower bound assumption, and  $\text{Var}_k$  measures a certain variance (over intervals of length  $k$ ) in the  $1 - \alpha$  conditional quantiles of the true radii. The proof builds on the strongly adaptive regret guarantee (in the quantile loss) for SAOCP (Proposition 4.1), and bounding parameter estimation errors by excess quantile losses using a *self-calibration inequality* type argument (Steinwart & Christmann, 2011).

**Setting** We consider the online conformal prediction setting described in Section 2. For any  $t \geq 1$ , let  $\mathcal{F}_t := \sigma(\{X_i, \hat{s}_i, S_i\}_{i \in [t-1]}, X_t)$  be the  $\sigma$ -algebra by all observed data  $\{(X_i, \hat{s}_i, S_i)\}_{i \leq t-1}$  as well as  $X_t$ . Note that by definition of the online conformal prediction setting, the predicted radius  $\hat{s}_t$  can only depend on information within  $\mathcal{F}_t$  as well as (possibly) external randomness. Consequently, we have  $S_t \perp\!\!\!\perp \hat{s}_t \mid \mathcal{F}_t$ , i.e.  $S_t$  and  $\hat{s}_t$  are conditionally independent given  $\mathcal{F}_t$ .

We now state our assumptions on the distributions of the true radii.

**Assumption C.1** (Density upper bounds). For all  $t \in [T]$ , there exists a constant  $L > 0$  such that  $S_t \mid \mathcal{F}_t$  is a continuous random variable that is bounded within  $[0, D]$  and has a density  $f_t : [0, D] \rightarrow \mathbb{R}_{\geq 0}$  with  $f_t(s) \leq L/D$  for all  $s \in [0, D]$ .

**Assumption C.2** (Density lower bounds). For all  $t \in [T]$ ,  $S_t \mid \mathcal{F}_t$  is a continuous random variable that is bounded within  $[0, D]$  and has a density  $f_t : [0, D] \rightarrow \mathbb{R}_{\geq 0}$ . With probability one, there exist constants  $b > 0, q \geq 2, \Delta_t > 0$  such that

$$f_t(s) \geq \frac{2b}{D} \left| \frac{2(s - s_t^*)}{D} \right|^{q-2} \quad (27)$$

for all  $s \in [s_t^* - \Delta_t, s_t^* + \Delta_t]$ , where

$$s_t^* := Q_{1-\alpha}(S_t \mid \mathcal{F}_t) \quad (28)$$

is the  $1 - \alpha$  conditional quantile of  $S_t$ .

As examples for Assumption C.2, the case where  $q = 2$  corresponds to a constant lower bound on the conditional density  $f_t(\cdot)$  locally around  $s_t^*$ , which holds e.g. if each  $f_t(\cdot)$  itself has a constant lower bound over  $[0, D]$  (this is the assumption made by Gibbs & Candès (2022)). A larger  $q$  makes the density lower bound (27) easier to satisfy and thus specifies a more relaxed assumption. We also note that  $s_t^*$  is itself a random variable which is measurable on  $\mathcal{F}_t$ .

For any interval  $I = [\tau, \tau + k - 1] \subseteq [T]$ , our coverage result depends on a certain variance between  $s_\tau^*, \dots, s_{\tau+k-1}^*$ . Concretely, define the interval quantile variation

$$\text{Var}_I := \sum_{t=\tau}^{\tau+k-1} \mathbb{E} \left[ \left( \frac{s_t^*}{D} - \frac{1}{kD} \sum_{i=\tau}^{\tau+k-1} \mathbb{E}[s_i^* \mid \mathcal{F}_\tau] \right)^2 \right]. \quad (29)$$

Then, the expected absolute difference between SAOCP's predictions  $\hat{s}_\tau, \dots, \hat{s}_{\tau+k-1}$  and the true conditional quantiles  $s_\tau^*, \dots, s_{\tau+k-1}^*$  is  $\tilde{\mathcal{O}}(k^{-1/(2q)}) + \tilde{\mathcal{O}}(k^{-1/q} \text{Var}_I^{1/q})$ . Due to the Lipschitzness of the CDFs (by Assumption C.1), the coverage error has a similar form. So SAOCP achieves better coverage when the interval quantile variation is lower, and it achieves approximately valid coverage as long as  $\text{Var}_I \leq o(|I|)$ . More formally, we have:

**Theorem C.3.** *Let Assumptions C.1 & C.2 hold. Fix any interval  $I = [\tau, \tau + k - 1] \subseteq [T]$ . Then, letting  $\Delta_I = \min\{\Delta_t : t \in I\}$ , Algorithm 1 achieves quantile estimation error*

$$\frac{1}{|I|} \sum_{t \in I} \mathbb{E} \left[ \left| \frac{\hat{s}_t - s_t^*}{D} \right| \right] \leq \mathcal{O} \left( \frac{1}{b^{1/q}} \frac{D}{\Delta_I} \left( \frac{\log T}{|I|} \right)^{1/(2q)} \right) + \mathcal{O} \left( \frac{L^{1/q}}{b^{1/q}} \frac{D}{\Delta_I} \left( \frac{\text{Var}_I}{|I|} \right)^{1/q} \right)$$

and interval miscoverage error

$$\frac{1}{|I|} \sum_{t \in I} \left| \mathbb{P}[Y_t \in \widehat{C}_t(X_t)] - (1 - \alpha) \right| \leq \mathcal{O} \left( \frac{L}{b^{1/q}} \frac{D}{\Delta_I} \left( \frac{\log T}{|I|} \right)^{1/(2q)} \right) + \mathcal{O} \left( \frac{L^{1+1/q}}{b^{1/q}} \frac{D}{\Delta_I} \left( \frac{\text{Var}_I}{|I|} \right)^{1/q} \right).$$

In Theorem C.3,  $q$  is a parameter quantifying the difficulty of lower bounding the distribution of  $S_t \mid \mathcal{F}_t$  away from its  $1 - \alpha$  conditional quantile  $s_t^*$ . If  $q$  is higher, then closeness to  $s_t^*$  is less correlated with the expected regret on the quantile loss (1). Meanwhile, the term  $D/\Delta_I$  grows larger as  $\alpha$  grows smaller. The inclusion of this term mirrors the inclusion of  $\alpha^{-2}$  in Theorem 4.2, and it indicates that more extreme quantiles are harder to learn.

**Proof of Theorem C.3** The key ingredient is the technical Lemma C.4, which uses the expected dynamic regret of a sequence  $\widehat{s}_\tau, \dots, \widehat{s}_{\tau+k-1}$  to upper bound the expected distance between that sequence and the true  $1 - \alpha$  conditional quantiles of  $S_\tau, \dots, S_{\tau+k-1}$ . We decompose the expected dynamic regret into the interval regret of Algorithm 1 (which we can upper bound by Proposition 4.1) and a term which we can use  $\text{Var}_I$  to upper bound. The desired coverage bound follows by the Lipschitzness of the CDFs of the  $S_t$ 's (implied by the density upper bound Assumption C.1). In this proof, we use  $\mathbb{E}_{\mathcal{F}_t}[X]$  as short-hand for the conditional expectation  $\mathbb{E}[X \mid \mathcal{F}_t]$ .

**Lemma C.4** (Bounding quantile estimation error by dynamic regret). *Fix any interval  $I = [\tau, \tau + k - 1] \subseteq [T]$ , and let  $\Delta_I = \min\{\Delta_t : t \in I\}$ . Under Assumption C.2, we have*

$$\sum_{t \in I} \mathbb{E} \left[ \left| \frac{\widehat{s}_t - s_t^*}{D} \right|^q \right] \leq \frac{2q(q-1)}{bD} \left( \frac{D}{2\Delta_I} \right)^q \sum_{t \in I} \mathbb{E}[\ell_{1-\alpha}(S_t, \widehat{s}_t) - \ell_{1-\alpha}(S_t, s_t^*)].$$

To prove Theorem C.3, we follow a similar technique to Zhang et al. (2018) and decompose the expected dynamic regret

$$\begin{aligned} \mathbb{E} \left[ \sum_{t \in I} \ell_{1-\alpha}(S_t, \widehat{s}_t) - \ell_{1-\alpha}(S_t, s_t^*) \right] &= \underbrace{\mathbb{E} \left[ \sum_{t \in I} \ell_{1-\alpha}(S_t, \widehat{s}_t) - \inf_s \sum_{t \in I} \ell_{1-\alpha}(S_t, s) \right]}_A \\ &\quad + \underbrace{\mathbb{E} \left[ \inf_s \sum_{t \in I} \ell_{1-\alpha}(S_t, s) - \ell_{1-\alpha}(S_t, s_t^*) \right]}_B. \end{aligned}$$

We first observe that term  $A$  is simply the expected interval regret on  $I$ . Since Proposition 4.1 bounds the strongly adaptive regret with probability one, we can bound  $A \leq 15D\sqrt{|I|(\log T + 1)}$ . Now, we analyze term  $B$ , and note that

$$\mathbb{E} \left[ \inf_s \sum_{t \in I} \ell_{1-\alpha}(S_t, s) - \ell_{1-\alpha}(S_t, s_t^*) \right] \leq \inf_s \mathbb{E} \left[ \sum_{t \in I} \mathbb{E}_{\mathcal{F}_t}[\ell_{1-\alpha}(S_t, s) - \ell_{1-\alpha}(S_t, s_t^*)] \right]$$

by Jensen's inequality and the tower property of conditional expectation. Now, for any  $t \in I$  and  $s \in [0, D]$ ,

$$\begin{aligned} &\mathbb{E}_{\mathcal{F}_t}[\ell_{1-\alpha}(S_t, s) - \ell_{1-\alpha}(S_t, s_t^*)] \\ &= \mathbb{E}_{\mathcal{F}_t}[(1 - \alpha)(S_t - s) \mathbb{1}[S_t > s] + \alpha(s - S_t) \mathbb{1}[S_t \leq s] - (1 - \alpha)(S_t - s_t^*) \mathbb{1}[S_t > s_t^*] - \alpha(s_t^* - S_t) \mathbb{1}[S_t \leq s_t^*]] \\ &= \begin{cases} \mathbb{E}_{\mathcal{F}_t}[(S_t - s) \mathbb{1}[S_t > s] - \alpha(S_t - s) - \alpha(s_t^* - S_t)] & S_t \leq s_t^* \\ \mathbb{E}_{\mathcal{F}_t}[(s - S_t) \mathbb{1}[S_t \leq s] - (1 - \alpha)(s - S_t) - (1 - \alpha)(s_t^* - S_t)] & S_t > s_t^* \end{cases} \\ &= \begin{cases} \mathbb{E}_{\mathcal{F}_t}[(S_t - s) \mathbb{1}[s < S_t \leq s_t^*] + \alpha(s - s_t^*)] & S_t \leq s_t^* \\ \mathbb{E}_{\mathcal{F}_t}[(s - S_t) \mathbb{1}[s_t^* \leq S_t \leq s] - (1 - \alpha)(s - s_t^*)] & S_t > s_t^* \end{cases} \\ &\stackrel{(i)}{=} \mathbb{E}_{\mathcal{F}_t}[(S_t - s) \mathbb{1}[s \leq S_t \leq s_t^*] + (s - S_t) \mathbb{1}[s_t^* \leq S_t \leq s]] + \alpha(s - s_t^*) \mathbb{P}_{\mathcal{F}_t}[S_t \leq s_t^*] - (1 - \alpha)(s - s_t^*) \mathbb{P}_{\mathcal{F}_t}[S_t \leq s_t^*] \\ &\stackrel{(ii)}{=} \mathbb{E}_{\mathcal{F}_t}[(S_t - s) \mathbb{1}[s \leq S_t \leq s_t^*] + (s - S_t) \mathbb{1}[s_t^* \leq S_t \leq s]] \\ &= \left| \int_s^{s_t^*} (x - s) f_t(x) dx \right| \stackrel{(iii)}{\leq} \frac{L(s_t^* - s)^2}{2D} \end{aligned}$$

Above, (i) uses the fact that  $s_t^*$  is  $\mathcal{F}_t$ -measurable, (ii) uses the fact that  $\mathbb{P}_{\mathcal{F}_t}[S_t \leq s_t^*] = 1 - \alpha$  by definition, and (iii) uses the density upper bound  $f_t(x) \leq \frac{L}{D}$  (Assumption C.1). Therefore, taking the infimum over  $s \in \mathbb{R}$  and by definition of  $\text{Var}_I$  (29), we have  $B \leq \frac{LD}{2} \text{Var}_I$ , and

$$\mathbb{E} \left[ \sum_{t \in I} \ell_{1-\alpha}(S_t, \hat{s}_t) - \ell_{1-\alpha}(S_t, s_t^*) \right] \leq 15D \sqrt{|I|(\log T + 1)} + \frac{LD}{2} \text{Var}_I.$$

We combine this result with the power-mean inequality, Lemma C.4, and the facts that  $(q(q-1))^{1/q} = \mathcal{O}(1)$  and  $(x+y)^{1/q} \leq x^{1/q} + y^{1/q}$  to prove the first part of Theorem C.3,

$$\frac{1}{|I|} \sum_{t \in I} \mathbb{E} \left[ \left| \frac{\hat{s}_t - s_t^*}{D} \right| \right] \leq \left( \frac{1}{|I|} \sum_{t \in I} \mathbb{E} \left[ \left| \frac{\hat{s}_t - s_t^*}{D} \right|^q \right] \right)^{1/q} \leq \mathcal{O} \left( \frac{D}{b^{1/q} \Delta_I} \left( \frac{\log T}{|I|} \right)^{1/(2q)} \right) + \mathcal{O} \left( \frac{D}{\Delta_I} \left( \frac{L \text{Var}_I}{b|I|} \right)^{1/q} \right). \quad (30)$$

To prove the desired coverage bound, we note that

$$\frac{1}{|I|} \sum_{t \in I} \left| \mathbb{P}[Y_t \in \hat{C}_t(X_t)] - (1 - \alpha) \right| = \frac{1}{|I|} \sum_{t \in I} \left| \mathbb{E}[\mathbb{E}_{\mathcal{F}_t}[\mathbb{1}[S_t \leq \hat{s}_t] - \mathbb{1}[S_t \leq s_t^*]]] \right| \leq \frac{1}{|I|} \sum_{t \in I} \mathbb{E}[|F_t(\hat{s}_t) - F_t(s_t^*)|],$$

where the final inequality uses Jensen's inequality and the fact that  $\mathbb{E}_{\mathcal{F}_t}[\mathbb{1}[S_t \leq s]] = F_t(s)$  is the CDF of  $S_t \mid \mathcal{F}_t$ . The result follows by combining (30) with Assumption C.1, which implies that the CDF  $F_t$  is  $L/D$ -Lipschitz.  $\square$

**Proof of Lemma C.4** We first consider the following general situation, where  $S$  is any continuous random variable bounded in  $[0, D]$  with a density  $f$ . Define  $s^* = Q_{1-\alpha}(S)$ . Let  $g$  be the density of the normalized random variable  $X = \frac{2S-D}{D} \in [-1, 1]$ , and let  $x^* = Q_{1-\alpha}(X) = \frac{2s^*-D}{D}$ . As in Steinwart & Christmann (2011, Example 2.3), assume that there exist constants  $b > 0, q \geq 2, \Delta > 0$  such that

$$g(x) \geq b|x - x^*|^{q-2} \iff f(s) \geq \frac{2b}{D} \left| \frac{2(s - s^*)}{D} \right|^{q-2} \quad (31)$$

for all  $s \in [s^* - \Delta, s^* + \Delta]$ . Let  $\beta = \frac{b}{q-1}$  and  $\gamma = \beta \left( \frac{2\Delta}{D} \right)^{q-1}$ . By Steinwart & Christmann (2011, Theorem 2.7),

$$\begin{aligned} |x - x^*| &\leq 2^{1-1/q} q^{1/q} \gamma^{-1/q} (\mathbb{E}[\ell_{1-\alpha}(X, x) - \ell_{1-\alpha}(X, x^*)])^{1/q} \\ &= 2 \left( \frac{q}{2\gamma} \mathbb{E}[\ell_{1-\alpha}(X, x) - \ell_{1-\alpha}(X, x^*)] \right)^{1/q} \\ &= 2 \left( \frac{q(q-1)}{b} \left( \frac{D}{2\Delta} \right)^{q-1} \mathbb{E}[\ell_{1-\alpha}(X, x) - \ell_{1-\alpha}(X, x^*)] \right)^{1/q} \end{aligned}$$

Since  $|x - x^*| = \frac{2}{D}|s - s^*|$ ,  $\ell_{1-\alpha}(X, x) = \frac{2}{D}\ell_{1-\alpha}(S, s)$ , and we can obtain

$$\left| \frac{s - s^*}{D} \right|^q \leq \frac{2q(q-1)}{bD} \left( \frac{D}{2\Delta} \right)^q \mathbb{E}[\ell_{1-\alpha}(S, s) - \ell_{1-\alpha}(S, s^*)].$$

Lemma C.4 now follows by fixing any  $t \in I$ , defining  $\mathcal{G}_t = \sigma(\mathcal{F}_t, \hat{s}_t)$  and noticing that  $S_t \mid \mathcal{G}_t \stackrel{\text{dist}}{=} S_t \mid \mathcal{F}_t$  (as  $\mathcal{G}_t$  only involves possibly an additional "external" randomness of the online prediction algorithm over  $\mathcal{F}_t$ ), so  $s_t^* = Q_{1-\alpha}(S_t \mid \mathcal{F}_t) = Q_{1-\alpha}(S_t \mid \mathcal{G}_t)$ . Since  $\hat{s}_t$  and  $s_t^*$  are  $\mathcal{G}_t$ -measurable, we can bound

$$\mathbb{E}_{\mathcal{G}_t} \left[ \left| \frac{\hat{s}_t - s_t^*}{D} \right|^q \right] \leq \frac{2q(q-1)}{bD} \left( \frac{D}{2\Delta_I} \right)^q \mathbb{E}_{\mathcal{G}_t}[\ell_{1-\alpha}(S_t, \hat{s}_t) - \ell_{1-\alpha}(S_t, s_t^*)].$$

The result follows by taking unconditional expectations of both sides, observing that  $\mathbb{E}[|s_t - s_t^*|^q] \leq \mathbb{E}[|s_t - s_t^*|^q]$  by Jensen's inequality, and summing over all  $t \in I$ .  $\square$

## Improved Online Conformal Prediction via Strongly Adaptive Online Learning

Method	LGBM (MAE = 0.19)				ARIMA (MAE = 0.09)				Prophet (MAE = 0.41)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
SCP	<b>.706</b> <sub>.018</sub>	.288 <sub>.008</sub>	.443 <sub>.020</sub>	.040 <sub>.003</sub>	<b>.911</b> <sub>.005</sub>	.203 <sub>.012</sub>	.257 <sub>.013</sub>	.017 <sub>.002</sub>	<b>.555</b> <sub>.019</sub>	.453 <sub>.016</sub>	.571 <sub>.019</sub>	.058 <sub>.004</sub>
NExCP	<b>.764</b> <sub>.013</sub>	.268 <sub>.009</sub>	.411 <sub>.017</sub>	.014 <sub>.001</sub>	<b>.904</b> <sub>.004</sub>	.185 <sub>.012</sub>	.254 <sub>.011</sub>	.009 <sub>.001</sub>	<b>.681</b> <sub>.013</sub>	.462 <sub>.018</sub>	.508 <sub>.017</sub>	.017 <sub>.001</sub>
FACI	<b>.822</b> <sub>.008</sub>	.254 <sub>.009</sub>	.282 <sub>.012</sub>	.008 <sub>.001</sub>	<b>.899</b> <sub>.003</sub>	.169 <sub>.011</sub>	.205 <sub>.008</sub>	<b>.007</b> <sub>.001</sub>	<b>.776</b> <sub>.009</sub>	.458 <sub>.018</sub>	.338 <sub>.013</sub>	.007 <sub>.001</sub>
SF-OGD	<b>.872</b> <sub>.004</sub>	.262 <sub>.009</sub>	.208 <sub>.008</sub>	.011 <sub>.002</sub>	<b>.901</b> <sub>.003</sub>	.170 <sub>.011</sub>	.191 <sub>.008</sub>	.010 <sub>.002</sub>	<b>.871</b> <sub>.003</sub>	.475 <sub>.019</sub>	.209 <sub>.008</sub>	.009 <sub>.002</sub>
FACI-S	<b>.863</b> <sub>.005</sub>	<b>.240</b> <sub>.009</sub>	.207 <sub>.008</sub>	.010 <sub>.002</sub>	<b>.891</b> <sub>.004</sub>	<b>.152</b> <sub>.009</sub>	.197 <sub>.008</sub>	.010 <sub>.002</sub>	<b>.856</b> <sub>.004</sub>	<b>.459</b> <sub>.018</sub>	.197 <sub>.006</sub>	<b>.007</b> <sub>.002</sub>
SAOCP	<b>.872</b> <sub>.004</sub>	.248 <sub>.009</sub>	<b>.170</b> <sub>.006</sub>	<b>.008</b> <sub>.001</sub>	<b>.886</b> <sub>.003</sub>	.155 <sub>.010</sub>	<b>.178</b> <sub>.007</sub>	.008 <sub>.001</sub>	<b>.868</b> <sub>.002</sub>	.473 <sub>.019</sub>	<b>.158</b> <sub>.004</sub>	<b>.007</b> <sub>.001</sub>

Table 4. Results on M4 Weekly (359 time series) with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). For all base predictors, SAOCP achieves the best local coverage error, best strongly adaptive regret, and second-best width. The only methods which achieve global coverage in  $(0.85, 0.95)$  for LGBM and Prophet are the ones that predict  $\hat{s}_{t+1}$  directly, not as a quantile of  $S_1, \dots, S_t$ .

Method	LGBM (MAE = 0.08)				ARIMA (MAE = 0.07)				Prophet (MAE = 0.11)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
SCP	<b>.932</b> <sub>.003</sub>	.205 <sub>.004</sub>	<b>.123</b> <sub>.006</sub>	.012 <sub>.001</sub>	<b>.938</b> <sub>.003</sub>	.199 <sub>.004</sub>	.126 <sub>.007</sub>	.012 <sub>.001</sub>	<b>.912</b> <sub>.006</sub>	.221 <sub>.009</sub>	.154 <sub>.009</sub>	.010 <sub>.000</sub>
NExCP	<b>.922</b> <sub>.002</sub>	.187 <sub>.003</sub>	.133 <sub>.006</sub>	.011 <sub>.000</sub>	<b>.922</b> <sub>.002</sub>	.175 <sub>.003</sub>	.136 <sub>.006</sub>	.012 <sub>.001</sub>	<b>.908</b> <sub>.003</sub>	.208 <sub>.010</sub>	.146 <sub>.006</sub>	.010 <sub>.000</sub>
FACI	<b>.910</b> <sub>.002</sub>	<b>.179</b> <sub>.004</sub>	.130 <sub>.005</sub>	<b>.010</b> <sub>.000</sub>	<b>.906</b> <sub>.002</sub>	<b>.162</b> <sub>.004</sub>	.132 <sub>.005</sub>	<b>.011</b> <sub>.001</sub>	<b>.900</b> <sub>.002</sub>	<b>.200</b> <sub>.010</sub>	.131 <sub>.005</sub>	<b>.009</b> <sub>.000</sub>
SF-OGD	<b>.904</b> <sub>.002</sub>	.190 <sub>.004</sub>	<b>.123</b> <sub>.004</sub>	.011 <sub>.000</sub>	<b>.901</b> <sub>.002</sub>	.176 <sub>.004</sub>	.130 <sub>.005</sub>	.012 <sub>.000</sub>	<b>.898</b> <sub>.002</sub>	.216 <sub>.010</sub>	.128 <sub>.004</sub>	.011 <sub>.000</sub>
FACI-S	<b>.909</b> <sub>.002</sub>	<b>.179</b> <sub>.003</sub>	.127 <sub>.005</sub>	<b>.010</b> <sub>.000</sub>	<b>.910</b> <sub>.002</sub>	<b>.166</b> <sub>.003</sub>	<b>.123</b> <sub>.004</sub>	<b>.011</b> <sub>.000</sub>	<b>.904</b> <sub>.002</sub>	<b>.203</b> <sub>.010</sub>	<b>.125</b> <sub>.004</sub>	<b>.009</b> <sub>.000</sub>
SAOCP	<b>.892</b> <sub>.002</sub>	<b>.179</b> <sub>.003</sub>	.132 <sub>.004</sub>	.012 <sub>.000</sub>	<b>.895</b> <sub>.002</sub>	<b>.166</b> <sub>.003</sub>	.127 <sub>.005</sub>	.012 <sub>.000</sub>	<b>.885</b> <sub>.002</sub>	.207 <sub>.010</sub>	.134 <sub>.004</sub>	.013 <sub>.000</sub>

Table 5. Results on NN5 Daily (111 time series) with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). All methods perform similarly well.

## D. Additional Time Series Experiments

Here, we report the results of our time series experiments (as described in Section 5.1) on M4 Weekly and NN5 Daily (the two smaller datasets) in Tables 4 and 5, respectively. The results on M4 Weekly (Table 4) are quite similar to those on M4 Daily (Table 2). All methods do reasonably well on NN5. Considering that even split conformal attains strong worst-case local coverage error, the residuals likely have a near-exchangeable distribution on NN5 (Barber et al., 2022, Theorems 2a, 3).

## E. Additional Experimental Details

We provide specific implementation details of all methods here. We use  $Q_{1-\alpha}(\cdot)$  to denote the (empirical)  $(1 - \alpha)$ -th quantile of a set of scalars, defined as

$$Q_{1-\alpha}\left(\{S_\tau\}_{\tau=1}^t\right) := \inf \left\{ s \in \mathbb{R} : \frac{1}{t} \sum_{\tau=1}^t \mathbb{1}[S_\tau \leq s] \geq 1 - \alpha \right\}. \quad (32)$$

1. SCP: Split conformal prediction (Vovk et al., 2005) predicts  $\hat{s}_{t+1} = Q_{1-\alpha}\left(\frac{1}{t} \sum_{\tau=1}^t \delta_{S_\tau}\right)$ .
2. NExCP: Non-exchangeable conformal prediction extends SCP by using a *weighted* quantile function  $\hat{s}_{t+1} = Q_{1-\alpha}\left(\frac{1}{t} \sum_{\tau=1}^t w_\tau \delta_{S_\tau}\right)$ . Barber et al. (2022) suggest using geometrically decaying weights to adapt NExCP to situations with distribution shift, so we use  $w_t = (1 - 3\alpha/4)^{1-t} w_1$ .
3. FACI: Fully Adaptive Conformal Inference has 4 hyperparameters: the individual expert learning rates  $\gamma_1, \dots, \gamma_N$ ; a target interval length  $k$ ; and the meta-algorithm learning rate  $\eta$ ; and a smoothing parameter  $\sigma$ . We set  $k = 100$  and follow Gibbs & Candès (2022) to set  $N = 8$ ,  $\sigma = \frac{1}{2k}$ ,  $\gamma = \{0.001, 0.002, 0.004, 0.008, 0.016, 0.032, 0.064, 0.128\}$ , and

$$\eta_t = \sqrt{\frac{\log(Nk) + 2}{\sum_{\tau=t-k}^{t-1} \mathbb{E}[\ell_\alpha(\beta_t, \alpha_t)^2]}}$$

where the expectation is over  $\alpha_t$ . We also tried  $k = 20$  for the time series experiments (to match our evaluation metrics), but the results were worse.

## Improved Online Conformal Prediction via Strongly Adaptive Online Learning

Method	LGBM (MAE = 0.05)				ARIMA (MAE = 0.14)				Prophet (MAE = 0.08)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
EnbPI	<b>.803</b> <sub>.006</sub>	.088 <sub>.003</sub>	.331 <sub>.009</sub>	.017 <sub>.001</sub>	<b>.916</b> <sub>.003</sub>	.220 <sub>.021</sub>	.186 <sub>.007</sub>	.032 <sub>.014</sub>	<b>.834</b> <sub>.007</sub>	.141 <sub>.005</sub>	.299 <sub>.011</sub>	.018 <sub>.001</sub>
EnbNEx	<b>.864</b> <sub>.002</sub>	.097 <sub>.003</sub>	.218 <sub>.006</sub>	.010 <sub>.001</sub>	<b>.907</b> <sub>.003</sub>	.198 <sub>.018</sub>	.193 <sub>.007</sub>	<u>.027</u> <sub>.012</sub>	<b>.892</b> <sub>.003</sub>	.151 <sub>.005</sub>	.195 <sub>.006</sub>	.010 <sub>.001</sub>
EnbFACI	<b>.856</b> <sub>.002</sub>	<b>.081</b> <sub>.003</sub>	.200 <sub>.005</sub>	<u>.007</u> <sub>.001</sub>	<b>.900</b> <sub>.003</sub>	<b>.181</b> <sub>.019</sub>	<u>.161</u> <sub>.005</sub>	<b>.021</b> <sub>.007</sub>	<b>.884</b> <sub>.002</sub>	<b>.130</b> <sub>.005</sub>	.164 <sub>.005</sub>	<u>.007</u> <sub>.001</sub>
EnbSF-OGD	<b>.871</b> <sub>.002</sub>	.098 <sub>.004</sub>	<u>.173</u> <sub>.004</sub>	.008 <sub>.001</sub>	<b>.906</b> <sub>.003</sub>	.201 <sub>.020</sub>	.165 <sub>.005</sub>	<u>.027</u> <sub>.009</sub>	<b>.898</b> <sub>.001</sub>	.145 <sub>.005</sub>	<u>.144</u> <sub>.004</sub>	.008 <sub>.001</sub>
EnbSAOCP	<b>.884</b> <sub>.002</sub>	<u>.091</u> <sub>.003</sub>	<b>.134</b> <sub>.002</sub>	<b>.005</b> <sub>.000</sub>	<b>.893</b> <sub>.003</sub>	.192 <sub>.024</sub>	<b>.150</b> <sub>.004</sub>	.053 <sub>.028</sub>	<b>.888</b> <sub>.002</sub>	<b>.130</b> <sub>.005</sub>	<b>.130</b> <sub>.002</sub>	<b>.005</b> <sub>.000</sub>

Table 6. Ensemble results on M4 Hourly with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). EnbSAOCP achieves the best local coverage error and strongly adaptive regret for all models, except ARIMA where its strongly adaptive regret is somewhat high.

Method	LGBM (MAE = 0.11)				ARIMA (MAE = 0.10)				Prophet (MAE = 0.18)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
EnbPI	<b>.512</b> <sub>.005</sub>	.093 <sub>.001</sub>	.716 <sub>.004</sub>	.077 <sub>.001</sub>	<b>.894</b> <sub>.002</sub>	.121 <sub>.002</sub>	.289 <sub>.004</sub>	.056 <sub>.034</sub>	<b>.791</b> <sub>.005</sub>	.288 <sub>.002</sub>	.374 <sub>.005</sub>	.062 <sub>.001</sub>
EnbNEx	<b>.749</b> <sub>.002</sub>	.142 <sub>.002</sub>	.520 <sub>.004</sub>	.020 <sub>.000</sub>	<b>.894</b> <sub>.001</sub>	.116 <sub>.003</sub>	.293 <sub>.004</sub>	<u>.035</u> <sub>.022</sub>	<b>.847</b> <sub>.003</sub>	.343 <sub>.003</sub>	.367 <sub>.005</sub>	.029 <sub>.000</sub>
EnbFACI	<b>.776</b> <sub>.002</sub>	.130 <sub>.002</sub>	.392 <sub>.003</sub>	.015 <sub>.000</sub>	<b>.887</b> <sub>.001</sub>	<u>.100</u> <sub>.007</sub>	.248 <sub>.003</sub>	<b>.015</b> <sub>.005</sub>	<b>.853</b> <sub>.003</sub>	<b>.232</b> <sub>.003</sub>	.230 <sub>.003</sub>	.025 <sub>.001</sub>
EnbSF-OGD	<b>.798</b> <sub>.002</sub>	.143 <sub>.002</sub>	.396 <sub>.004</sub>	.021 <sub>.001</sub>	<b>.898</b> <sub>.001</sub>	<u>.106</u> <sub>.002</sub>	<u>.241</u> <sub>.003</sub>	.047 <sub>.034</sub>	<b>.900</b> <sub>.002</sub>	.292 <sub>.003</sub>	<u>.195</u> <sub>.002</sub>	<u>.023</u> <sub>.001</sub>
EnbSAOCP	<b>.875</b> <sub>.001</sub>	<b>.138</b> <sub>.002</sub>	<b>.203</b> <sub>.002</sub>	<b>.007</b> <sub>.000</sub>	<b>.908</b> <sub>.001</sub>	<b>.096</b> <sub>.002</sub>	<b>.187</b> <sub>.002</sub>	.044 <sub>.033</sub>	<b>.917</b> <sub>.001</sub>	<u>.233</u> <sub>.002</sub>	<b>.139</b> <sub>.001</sub>	<b>.011</b> <sub>.000</sub>

Table 7. Ensemble on M4 Daily with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). EnbSAOCP achieves the best or second best width, local coverage error, and strongly adaptive regret for all models. It is also the only method which achieves valid coverage for LGBM.

4. SF-OGD: Scale-Free Online Gradient Descent. The only hyperparameter is the maximum radius  $D$ . For the time series experiments (Section 5.1, Appendix F), we set  $D/\sqrt{3}$  for each horizon  $h$  equal to the largest  $h$ -step residual observed on the calibration split of the training data. For the  $m$ -way image classification experiments (Section 5.2, Appendix G), we set  $D = 1 + \lambda\sqrt{m - k_{\text{reg}}}$ , where  $\lambda$  and  $k_{\text{reg}}$  are the width regularization parameters in (15).
5. FACI-S: FACI applied to  $S_t$  rather than  $\alpha_t$ . The hyperparameters are the same, except the losses used to compute  $\eta_t$  are  $\ell_{1-\alpha}(S_t, \hat{s}_t)$ , and the learning rates are multiplied by  $D$ . We set  $D$  in the same way as SF-OGD.
6. SAOCP: There are 2 hyperparameters: the maximum radius  $D$  and the lifetime multiplier  $g$  in (8). We set  $D$  in the same way as SF-OGD. We set  $g = 8$  for the time series experiments and  $g = 32$  for the image classification experiments.

## F. Time Series Experiments with Ensemble Models

In this section, we replicate the experiments of Section 5.1 using ensemble models trained with the method of EnbPI (Xu & Xie, 2021). Specifically, we train base learner  $\hat{f}^{(b)}$  on  $(X_t, Y_T)_{t \in I_b}$ , where  $I_b$  is sampled randomly from  $[T]$ . Then, we obtain the residual  $S_t^y = |y - \phi(\hat{f}^{(b)}(X_t) : I_b \not\cong t)|$  by aggregating all models not trained on  $(X_t, Y_t)$ . Finally, the residual of a new observation  $(X_{T+1}, Y_{T+1})$  is  $|Y_{T+1} - \phi(\hat{f}^{(b)}(X_{T+1}) : b \in [B])|$ . We use  $B = 5$  models in the ensemble.

EnbPI predicts the radius  $\hat{s}_{t+1}$  as the  $1 - \alpha$  empirical quantile of the previously observed residuals, as in split conformal prediction. However, these prediction sets can be obtained via an arbitrary function of the scores, including NExCP, FACI, SF-OGD, or SAOCP. Besides EnbPI, we call these hybrid methods EnbNEx, EnbFACI, EnbSF-OGD, and EnbSAOCP respectively. We use the same hyperparameters as described in Appendix E.

This approach puts the other methods on even footing with EnbPI, because EnbPI removes the need for a train/calibration split and changes the underlying model from a single learner to a more accurate ensemble. We consider the contributions of EnbPI orthogonal to our own, and this section shows that their method can successfully be combined with ours.

The results mirror those of Section 5.1. EnbSAOCP generally obtains the best or second-best interval width, worst-case local coverage error, and strongly adaptive regret on all M4 datasets (Tables 6, 7, 8). On NN5 (Table 9), all methods obtain similar strongly adaptive regret. However, EnbSF-OGD and EnbSAOCP obtain the best and second-best worst-case local coverage error at the cost of having slightly wider intervals. Across the board, EnbSAOCP has narrower intervals than EnbSF-OGD.

Improved Online Conformal Prediction via Strongly Adaptive Online Learning

Method	LGBM (MAE = 0.12)				ARIMA (MAE = 0.07)				Prophet (MAE = 0.16)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
EnbPI	<b>.540</b> <sub>.019</sub>	.118 <sub>.006</sub>	.596 <sub>.019</sub>	.079 <sub>.005</sub>	<b>.893</b> <sub>.008</sub>	.174 <sub>.008</sub>	.266 <sub>.015</sub>	.021 <sub>.003</sub>	<b>.785</b> <sub>.016</sub>	.288 <sub>.010</sub>	.369 <sub>.019</sub>	.048 <sub>.004</sub>
EnbNEx	<b>.720</b> <sub>.012</sub>	.168 <sub>.007</sub>	.483 <sub>.016</sub>	.015 <sub>.001</sub>	<b>.910</b> <sub>.004</sub>	.161 <sub>.009</sub>	.252 <sub>.011</sub>	.012 <sub>.001</sub>	<b>.875</b> <sub>.008</sub>	.375 <sub>.013</sub>	.305 <sub>.015</sub>	.027 <sub>.001</sub>
EnbFACI	<b>.784</b> <sub>.009</sub>	.159 <sub>.006</sub>	.336 <sub>.013</sub>	.010 <sub>.001</sub>	<b>.905</b> <sub>.004</sub>	.142 <sub>.008</sub>	.193 <sub>.008</sub>	.009 <sub>.001</sub>	<b>.881</b> <sub>.006</sub>	.235 <sub>.007</sub>	.201 <sub>.010</sub>	.017 <sub>.001</sub>
EnbSF-OGD	<b>.795</b> <sub>.011</sub>	.172 <sub>.007</sub>	.329 <sub>.015</sub>	.016 <sub>.002</sub>	<b>.908</b> <sub>.003</sub>	.152 <sub>.009</sub>	.184 <sub>.008</sub>	.012 <sub>.002</sub>	<b>.915</b> <sub>.004</sub>	.304 <sub>.009</sub>	.164 <sub>.008</sub>	.023 <sub>.002</sub>
EnbSAOCP	<b>.874</b> <sub>.004</sub>	<b>.165</b> <sub>.006</sub>	<b>.173</b> <sub>.006</sub>	<b>.005</b> <sub>.000</sub>	<b>.909</b> <sub>.003</sub>	<b>.131</b> <sub>.006</sub>	<b>.151</b> <sub>.006</sub>	<b>.008</b> <sub>.001</sub>	<b>.907</b> <sub>.003</sub>	<b>.232</b> <sub>.006</sub>	<b>.133</b> <sub>.004</sub>	<b>.012</b> <sub>.001</sub>

Table 8. Ensemble results on M4 Weekly with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). EnbSAOCP achieves the best width, local coverage error, and strongly adaptive regret for all models. It is also the only method which achieves valid coverage for LGBM.

Method	LGBM (MAE = 0.09)				ARIMA (MAE = 0.07)				Prophet (MAE = 0.07)			
	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>	Coverage	Width	LCE <sub>k</sub>	SAReg <sub>k</sub>
EnbPI	<b>.859</b> <sub>.005</sub>	<b>.164</b> <sub>.004</sub>	.225 <sub>.010</sub>	.012 <sub>.001</sub>	<b>.910</b> <sub>.003</sub>	.163 <sub>.004</sub>	.159 <sub>.008</sub>	.010 <sub>.001</sub>	<b>.904</b> <sub>.004</sub>	.163 <sub>.005</sub>	.160 <sub>.007</sub>	<b>.010</b> <sub>.001</sub>
EnbNEx	<b>.882</b> <sub>.003</sub>	<b>.177</b> <sub>.004</sub>	.177 <sub>.007</sub>	.010 <sub>.000</sub>	<b>.920</b> <sub>.002</sub>	.174 <sub>.003</sub>	.138 <sub>.006</sub>	.012 <sub>.001</sub>	<b>.923</b> <sub>.002</sub>	.178 <sub>.006</sub>	.126 <sub>.004</sub>	.012 <sub>.001</sub>
EnbFACI	<b>.883</b> <sub>.002</sub>	<b>.177</b> <sub>.004</sub>	.166 <sub>.006</sub>	<b>.009</b> <sub>.000</sub>	<b>.905</b> <sub>.002</sub>	<b>.156</b> <sub>.003</sub>	.131 <sub>.005</sub>	<b>.009</b> <sub>.001</sub>	<b>.907</b> <sub>.002</sub>	<b>.156</b> <sub>.003</sub>	.123 <sub>.004</sub>	<b>.010</b> <sub>.000</sub>
EnbSF-OGD	<b>.895</b> <sub>.002</sub>	.194 <sub>.004</sub>	<b>.132</b> <sub>.004</sub>	.010 <sub>.000</sub>	<b>.912</b> <sub>.002</sub>	.176 <sub>.004</sub>	<b>.120</b> <sub>.004</sub>	.012 <sub>.001</sub>	<b>.912</b> <sub>.002</sub>	.178 <sub>.004</sub>	<b>.120</b> <sub>.004</sub>	.013 <sub>.001</sub>
EnbSAOCP	<b>.886</b> <sub>.002</sub>	.188 <sub>.004</sub>	<b>.131</b> <sub>.005</sub>	.012 <sub>.000</sub>	<b>.906</b> <sub>.002</sub>	.168 <sub>.003</sub>	.122 <sub>.004</sub>	.012 <sub>.000</sub>	<b>.905</b> <sub>.002</sub>	.170 <sub>.003</sub>	<b>.118</b> <sub>.003</sub>	.013 <sub>.000</sub>

Table 9. Ensemble results on NN5 Daily with target coverage  $1 - \alpha = 0.9$  and interval size  $k = 20$ . Results are formatted as  $\text{mean}_{\text{std}}$ . Best results are **bold**, while second best are underlined, as long as the method’s global coverage is in  $(0.85, 0.95)$  (green). EnbSAOCP and EnbSF-OGD achieve the best and second-best worst-case local coverage error error at the cost of having slightly wider intervals.

### G. Image Classification on ImageNet/ImageNet-C

We replicate the experiments of Section 5.2 using a ResNet-50 classifier on the ImageNet (Deng et al., 2009) base dataset and its corrupted version ImageNet-C (Hendrycks & Dietterich, 2019). We train the model using SGD with learning rate 0.1 (annealed by a factor of 10 every 7 epochs), momentum 0.9, batch size 256, and early stopping if validation accuracy stops improving for 10 epochs. The model achieved a final test accuracy of 52.8%. Figure 2 shows the results.

When performing uncertainty quantification, we use the conformal score (15) with width regularization parameters  $\lambda = 0.01$  and  $k_{\text{reg}} = 20$ . As in Section 5.2, SAOCP and SF-OGD’s local coverages remain the closest to the target of 0.9. The differences are most apparent when the distribution shift is sudden, suggesting that they are able to adapt to these distribution shifts more quickly than other methods. While all methods attain similar prediction set sizes, NExCP and FACI adapt more slowly to the best fixed prediction set size than SAOCP and SF-OGD. SAOCP also has better coverage than SF-OGD.

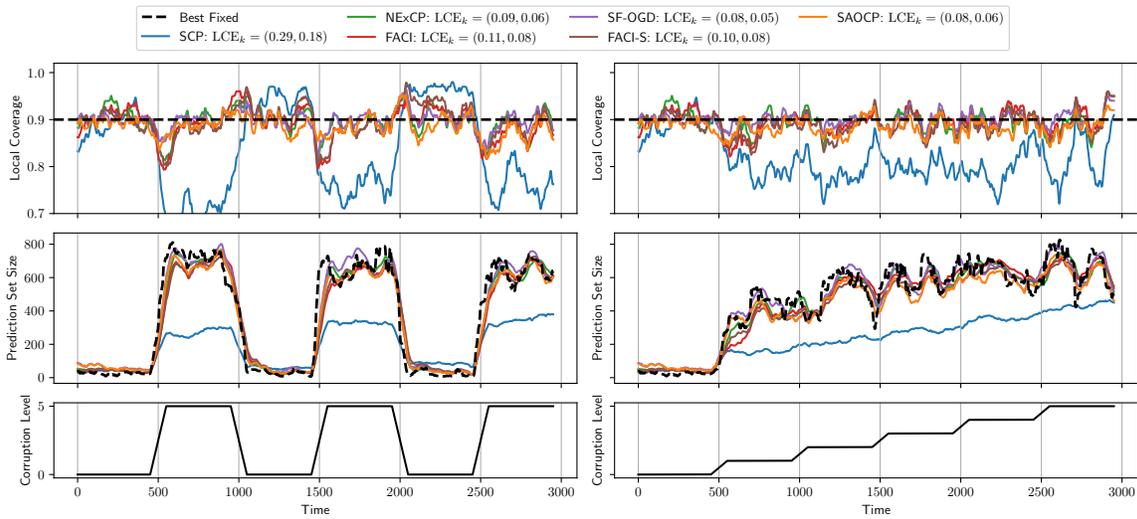


Figure 2. Local coverage (top row) and prediction set size (second row) achieved by various UQ methods when the distribution shifts between ImageNet and ImageNet-C every 500 steps. We plot moving averages with window size  $k = 100$ . Left: sudden shifts between corruption level 0 and 5. Right: gradual shift from level 0 to 5. SAOCP and SF-OGD’s local coverage remain the closest to the target of 0.9, especially at the change points. While the two methods attain similar local coverage, SAOCP returns smaller prediction sets than SF-OGD.