
Conformal Prediction for Time-series Forecasting with Change Points

Sophia Sun, Rose Yu
Computer Science and Engineering
University of California, San Diego

Abstract

Conformal prediction has been explored as a general and efficient way to provide uncertainty quantification for time series. However, current methods struggle to handle time series data with change points — sudden shifts in the underlying data-generating process. In this paper, we propose a novel Conformal Prediction for Time-series with Change points (CPTC) algorithm, addressing this gap by integrating a model to predict the underlying state with online conformal prediction to model uncertainties in non-stationary time series. We prove CPTC’s validity and improved adaptivity in the time series setting under minimum assumptions, and demonstrate CPTC’s practical effectiveness on 6 synthetic and real-world datasets, showing improved validity and adaptivity compared to state-of-the-art baselines.

1 Introduction

Uncertainty Quantification (UQ) is a key building block of reliable machine learning systems. Conformal prediction (CP) is a popular distribution-free UQ method that provides finite sample coverage guarantees without placing assumptions on the underlying data generating process [47, 29]. Because of its simplicity and generality, conformal prediction has seen wide adoption in many tasks, ranging from healthcare [7], robotics [31], to validating the factuality of generative models [15].

Existing CP algorithms for time series forecasting mostly adopt an online learning framework [22, 4], where the prediction interval adapts to distribution changes in data *reactively*. The online CP algorithms achieve a asymptotic marginal coverage guarantee (the miscoverage rate converges to a specified fraction as the time horizon goes to infinity). As a result, the algorithms inevitably exhibit miscoverage when distribution shifts occur, and then have to over- or under-cover in later timesteps to compensate. This behavior is evident in the ACI and CP baselines in Figure 1. This is undesirable in practice as these under-covered periods may incur a lot of risk.

This paper builds upon the observation that in real-world scenarios, distribution shifts in time-series are often *predictable*. Take for example the task of forecasting electricity demands: we know that the hidden dynamics may differ between day and night, or during weekdays and weekends. The same logic applies to traffic forecasting and product demand forecasting, where there are patterns in the shifts between “surge times” and “normal times”. In this work, we study time series data that exhibit this type of abrupt shifts, or change points, in the underlying generative process.

State Space Models (SSMs) [25, 9] are a powerful tool for modeling time series data, as they provide a structured framework to capture the underlying temporal dependencies through latent states. In particular, Switching Dynamical Systems (SDS) [1] is a type of SSM with an additional set of latent variables (known as *switches*) that represent the operating mode active at the current timestep. SDS introduces the inductive biases that allow SSMs to switch between a discrete set of dynamics, which can be learned from training data. It is shown to be a flexible, robust, and interpretable representation of time series with varying dynamics exhibiting change points [8, 14]. SDS is deployed not only for

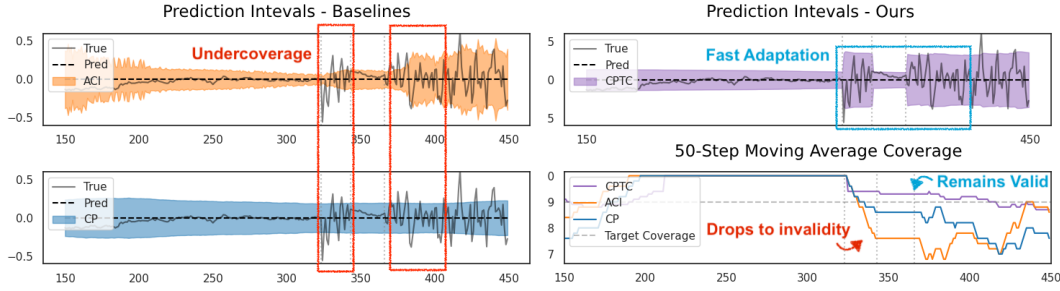


Figure 1: Comparison of the prediction intervals obtained by our algorithm CPTC (purple) against online conformal prediction baselines on synthetic data. The vertical dashed line marks the distribution shifts; ideal behavior is consistent coverage at the horizontal dashed line in the final panel. Bottom right panel shows that CPTC achieves fast adaptation and remains valid when change points occur.

energy [26, 40, 18] and traffic [55, 13] forecasting as examples given above, but also on a wide range of applications including neuroscience [23, 38, 27], engineering [42, 50], and sports analytics [54], all would benefit from robust and calibrated uncertainty quantification.

Our contributions are:

- We propose a new algorithm **Conformal Prediction for Time series with Change points** (CPTC) that utilizes state transition predictions to improve uncertainty quantification for time-series forecasts. Leveraging properties of a SDS model, CPTC consolidates multiple future forecasts to adaptively adjust its prediction intervals when underlying dynamics shift.
- We prove that CPTC achieves asymptotic valid coverage, *without assumptions* on the data generation process or state transition model accuracy. When predicted state transitions align well with distribution shifts, CPTC can anticipate uncertainty and adapt faster.
- We show strong empirical results on 3 synthetic and 3 real-world datasets. Compared to on-line conformal prediction baselines, CPTC achieves more robust coverage with comparable prediction intervals sharpness (example in Figure 1), and is computationally light.

2 Related Work

Probabilistic Forecasting for Time-Series with Change Points. Probabilistic forecasting has become central to modern time-series analysis by modeling a distribution over future outcomes rather than a single-point forecast [12]. Classic approaches often adopt Bayesian or approximate Bayesian strategies for uncertainty quantification: for instance, Bayesian Neural Networks (BNNs) [49, 51, 19] or neural process models [20, 43] use neural networks to parameterize the underlying stochastic process. Other methods, such as DeepAR [39], PatchTST [36], or Temporal Fusion Transformer [30], generate probabilistic forecasts directly by minimizing a pinball loss (i.e., quantile loss), thereby learning predictive distributions.

A key challenge arises when *change points* occur - i.e., abrupt distribution shifts in the time series [35]. Change point detection and segmentation have been studied extensively (beyond our scope, see [3] for a recent survey). Modeling time series with change points typically features state-space models (SSM), and in particular switching dynamical systems (SDS) [24, 8, 34]. These SDS approaches use variational inference to fit neural parameterizations of transitions and emission distributions. However, existing probabilistic forecasters often lack strict calibration guarantees and can be overconfident in practice [28], and many require specialized architecture tuning for different tasks, limiting their applicability to real-world UQ challenges.

Conformal Prediction for time-series. Conformal prediction (CP) [47] has emerged as a leading framework for UQ due to its algorithmic simplicity, broad applicability, and finite-sample coverage guarantees; see [6] for a comprehensive introduction. For time-series forecasting, some works focus on joint coverage over fixed-length horizons [41, 44, 57], while others assume stable underlying dynamics and adapt only to shifts in residual distributions [52]. [45] uses known covariate shifts (by propensity weighting) to achieve theoretical coverage, but can fail in settings where the shift mechanism is unknown, which is typical for time series forecasting. Conformal prediction has also been extended for change-point detection in time series via conformal martingales [48, 46].

Recent works on online conformal prediction develop online strategies that guarantee marginal coverage for adversarial or nonstationary data streams. ACI [21], AgACI [56], DtACI [22], and Conformal PID control [4] leverage online optimization and have a "reaction" period when distribution shifts. Multivalid Prediction (MVP) [11] provides group coverage guarantees and learns a calibration threshold online, although it can require longer calibration windows. Closer to our work, SPCI [53] learns an autocorrelative estimation of the residual's quantiles through Quantile Random Forest; whereas HopCPT [10] condition their conformal interval on similar parts of the time series using a Modern Hopfield Network. The drawback of these two methods is that the regression models are trained during inference time, which (1) requires a long cold start window to learn meaningful associations and (2) is computationally expensive. Our contribution is a CP algorithm that is computationally light during inference time, and can leverage state prediction capabilities within the SDS model to anticipate and adapt quickly to distribution changes.

3 Background

3.1 Conformal Prediction

We briefly review the algorithm and guarantees for conformal prediction, and refer readers to [6] for a thorough introduction. In this paper CP refers to *split conformal prediction* (we consider full conformal prediction out of scope, because its computational complexity renders it nonviable for deep learning settings). The goal of conformal prediction is to produce a *valid* prediction interval (Def. 3.1) for any underlying prediction model.

Definition 3.1 (Validity). Given a new data point (x, y) , where y is the prediction target and x is the covariates, and a desired confidence $1 - \alpha \in (0, 1)$, the prediction interval $\Gamma^{1-\alpha}(x)$ is a subset of \mathcal{Y} containing probable outputs given x . The region $\Gamma^{1-\alpha}$ is valid if

$$P(y \in \Gamma^{1-\alpha}(x)) \geq 1 - \alpha \quad (1)$$

Let $\mathcal{D} = \{(x^i, y^i)\}_{i=1}^n$ be a dataset whose data points (x^i, y^i) are sampled from a distribution on $\mathcal{X} \times \mathcal{Y}$. Conformal prediction splits the dataset into a proper training set \mathcal{D}_{train} and a calibration set \mathcal{D}_{cal} . A prediction model $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$ is trained on \mathcal{D}_{train} . We use a *nonconformity score* function $A : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ to quantify how well a data sample from calibration *conforms* to the training dataset. Typically, we choose a metric of disagreement between the prediction and the true label as the non-conformity score, such as the Euclidean distance:

$$A(x, y) \stackrel{\text{e.g.}}{=} d(y, \hat{f}(x)) \stackrel{\text{e.g.}}{=} \|y - \hat{f}(x)\|_2 \quad (2)$$

Let $\mathcal{S} = \{A(x^i, y^i)\}_{(x^i, y^i) \in \mathcal{D}_{cal}}$ denote the set of nonconformity scores of all samples in the calibration set \mathcal{D}_{cal} . During inference time given a new data x^{n+1} , the conformal prediction interval is constructed as in Eqn 3, where Q is the empirical quantile function.

$$\Gamma^{1-\alpha}(x^{n+1}) := \{y : A(x^{n+1}, y) \leq Q^{1-\alpha}(\mathcal{S} \cup \{\infty\})\} \quad (3)$$

We say a sample is *covered* if the true value lies in the prediction interval $y^{n+1} \in \Gamma^{1-\alpha}(x^{n+1})$. Conformal prediction is guaranteed to produce valid prediction intervals [47] if the calibration data and test data are exchangeable (in a dataset $\{(x^i, y^i)\}_{i=1}^n$ of size n , any of its $n!$ permutations are equally probable).

Conformal prediction has been extended to the online setting. More precisely, the algorithm observes $(x_1, y_1), (x_2, y_2), \dots$ sequentially, and needs to construct a prediction set $\Gamma_t^{1-\alpha}$ for y_t at time step t . Without making any distributional assumptions on the data generation process, online CP algorithms achieve asymptotic guarantee as Eqn 4, where $\mathbb{1}\{\cdot\}$ is the indicator function:

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{1}\{y_t \notin \Gamma_t^{1-\alpha}(x_t)\} = \alpha \quad (4)$$

We refer readers to [5] for a summary of theoretical results in this setting.

3.2 Switching Dynamical Systems for Modeling Time-series with Change Points

Switching dynamics systems (SDS), or mixed dynamics systems, provide a powerful framework for modeling time series with change points, where the underlying dynamics shift between distinct

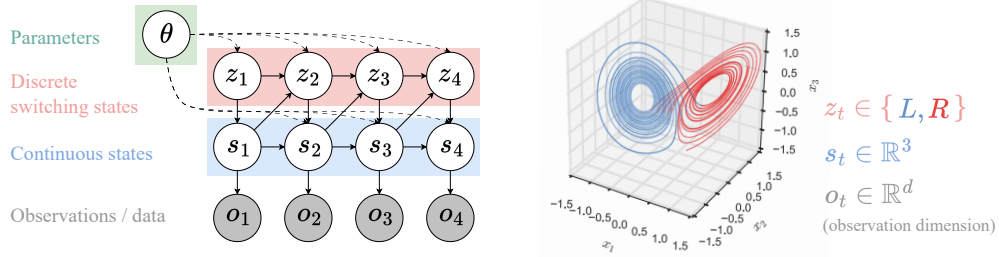


Figure 2: (Left) Generative model of the SDS dynamics as in Eqn 5. Shaded circles represent observed variables; hollow circles are latent variables. (Right) Example to illustrate notation. We show a Lorenz attractor, a canonical nonlinear dynamical system, approximated by a linear SDS [32].

regimes. SDS assumes the system consists of a total of K different base dynamics. It explicitly incorporates discrete switching *states* to represent the abrupt shifts among base dynamics. We denote the discrete switching states at timestep t as $z_t \in \mathcal{Z}$, $\mathcal{Z} = \{1, \dots, K\} \subset \mathbb{Z}^+$, to index one of K base dynamical systems, and represent the observed sequence as o_1, \dots, o_T or $o_{1:T}$ for conciseness.

There have been many different variations and implementations of SDS. Classic SDS have *memory-less* state transitions, i.e. $p(z_t|z_{t-1})$ is parameterized by a stochastic transition matrix $\mathbf{A} \in \mathbb{R}^{K \times K}$. This results in an “open loop” system and the state duration follows a geometric distribution. Recent works have sought to fix this issue, for example by closing the loop using recurrent dependencies on s_t [33], or modeling explicit durations [8]. Our generative model is illustrated in figure 2): for generality, we factor out any non-Markovian covariates to model parameters θ . The joint distribution can be factorized as

$$P(o_{1:T}, s_{1:T}, z_{1:T}) = \prod_{t=1}^T P(o_t|s_t)P(s_t|s_{t-1}, z_t, \theta)P(z_t|s_{t-1}, z_{t-1}, \theta) \quad (5)$$

where $p(s_1|z_1)p(z_1)$ is the initial state prior. The K dynamical systems have continuous state transition $p(z_t|z_{t-1}, s_{t-1}, \theta)$, and the emission $p(o_t|s_t)$ can be linear or nonlinear (e.g. parameterized by a neural network). Moreover, the SDS framework naturally supports forecasting both the next observation o_{t+1} via $P(o_{t+1} | s_{t+1})$ and the next switching state z_{t+1} via $P(z_{t+1} | s_t, z_t, \theta)$.

4 Conformal Prediction for Time series with Change points (CPTC)

We introduce CPTC, a novel conformal prediction algorithm that utilizes the state prediction from SDS to improve uncertainty quantification for time series forecasting.

Problem Setting We are given a sequence of observations $o_t \in \mathbb{R}^d$ for all t potentially with change points. To be consistent with the standard notation in conformal prediction literature, let $\mathcal{X} \subseteq \mathbb{R}^m$ denote the input feature domain and $\mathcal{Y} \subseteq \mathbb{R}^d$ the target observation space. $\mathcal{Z} = \{1, \dots, K\} \subset \mathbb{Z}^+$ indexes one of K underlying states. The input features can include past observations and other environmental covariates. In our implementations, $m = k \cdot d$ with k being the look back window.

We consider the online conformal prediction setting where we have a sequence of data points from $\mathcal{X} \times \mathcal{Y}$: $(x_1, y_1), (x_2, y_2), \dots, (x_T, y_T)$. We have a trained state predictor for $\hat{p}(z_t|x_t)$, a state-conditioned forecaster $\hat{f} : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}$ (a probabilistic SDS model can serve both purposes), and are given a desired confidence level $1 - \alpha$. At each time step t , we have seen $(x_1, y_1), (x_2, y_2), \dots, (x_{t-1}, y_{t-1})$ and the next input feature vector x_t . We want to produce an adaptive, time-dependent prediction interval based on input $\Gamma_t(x_t) \subseteq \mathcal{Y}$ for the unseen true target y_t , which we will learn at the next time step. Let $z_t \in \mathcal{Z}$ be the true (latent switching) state at time t . Our goal is to produce prediction intervals with correct empirical coverage, i.e.

$$P(y_t \in \Gamma_t(x_t)) \geq 1 - \alpha \text{ for any } t \geq 0 \quad (6)$$

4.1 Algorithm

Existing online conformal prediction algorithms often rely on assumptions about how distribution would shift in time series. This is infeasible for data with change points due to unknown shift mechanism. The main contribution our algorithm makes towards addressing this issue is *anticipating* this type of abrupt shifts. Our CPTC algorithm is outlined in pseudo-code in algorithm 1. In this section, we discuss the general procedure and key components of the algorithm.

Algorithm 1: Conformal Prediction for Time series with Change points (CPTC)

Input: nonconformity score function A , probabilistic state model $\hat{p}(z_t|x_t)$ with prior $\hat{p}(z_0)$, forecaster model $\hat{f} : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}$, confidence level $1 - \alpha$, learning rate γ

Output: Prediction intervals $\Gamma(x_t)$ for $t \geq 1$

```
1 for  $z \in \mathcal{Z}$  do
2   | Initialize  $\mathcal{S}_z \leftarrow \{\}$ ,  $\alpha_{z,0} \leftarrow \alpha$ 
3 end for
4 if exists warm start data  $\{(x_t, y_t)\}_{t=-w}^0$  then
5   | for  $t \in [-w, 0]$  do
6     |  $\mathcal{S}_{z_t} \leftarrow \mathcal{S}_{z_t} \cup \{A(x_t, y_t)\}$ , where  $z_t \sim \hat{p}(z_t = z|x_t)$ ; // warm start scores
7   | end for
8 end if
9 for  $t \in [1, T]$  do
10  | for  $z \in \mathcal{Z}$  do
11    |  $\Gamma_{z,t}(x_t) \leftarrow \{y : A(x_t, y) \leq Q^{1-\alpha_{z,t}}(\mathcal{S}_z \cup \{\infty\})\}$ ; // state-specific CP
12  | end for
13  |  $\Gamma_t(x_t) \leftarrow$  Aggregate the  $\Gamma_{z,t}(x_t)$ s by Eqn 10 or Eqn 11;
    | Output:  $\Gamma_t(x_t)$ 
14  | Sample state  $\hat{z}_t \sim \hat{p}(z_t|x_t)$ ; // state-specific coverage target tracking
15  | Update  $\alpha_{\hat{z}_t,t+1} \leftarrow \alpha_{\hat{z}_t,t} + \gamma \cdot (\alpha - err_t)$ , where  $err_t = \mathbb{1}\{y_t \notin \Gamma_t(x_t)\}$ ;
16  | Update scores  $\mathcal{S}_{\hat{z}_t} \leftarrow \mathcal{S}_{\hat{z}_t} \cup \{A(x_t, y_t)\}$ ;
17 end for
```

Switching Behavior and SDS Integration. We explicitly model the abrupt shift in dynamics with the switching dynamical systems (SDS) to anticipate the change. In our formulation, the latent state z_t indicates which dynamical regime (out of K possible regimes) is currently active at time t . Given the SDS formulation (e.g. Eqn 5), we can factor the state transition probability from the conformal prediction process and do calibration for each of the K dynamics separately. The goal in Eqn 6 can therefore be written as:

$$\sum_{z \in \mathcal{Z}} P(y_t \in \Gamma_t(x_t) | z_t = z) P(z_t = z | x_t) \geq 1 - \alpha \quad (7)$$

State Prediction. The SDS model provides us with the state predictor $\hat{p}(z_t | x_t)$. In practice, SDS models [8, 33] have existing z_t forecasts that we can extract. For other model architectures, one can modify them to output an extra covariate prediction for the state, or add an auxiliary model to classify or cluster the input space into different regimes. In practice, z_t doesn't necessarily have to come from a model, and can be any discrete contextual variable such as weekday/weekend/holiday or day/night.

Initialization and Warm Start, For each switching state $z \in \mathcal{Z}$, we initialize (1) a mode-specific set of nonconformity scores \mathcal{S}_z and (2) the mode-specific confidence level $\alpha_{z,0} = \alpha$. If there are w steps of warm-start data $\{(x_t, y_t)\}_{t=-w}^0$, we initialize the score sets \mathcal{S}_z by sampling $z_t \sim \hat{p}(z_t = z|x_t)$ for $t = -w, \dots, 0$, and inserting $A(x_t, y_t)$ into \mathcal{S}_{z_t} . Depending on the size of the warm-start window and desired properties, practitioners can adjust the warm-starting approach. For example, using the entire warm-start window $\mathcal{S}_z = \{A(x_t, y_t)\}_{t=-w}^0$ for all $z \in \mathcal{Z}$ will result in better stability but less adaptability to modes. The approach does not change the algorithm's coverage guarantees.

State-Specific Prediction Intervals. At every time step $t \in [1, T]$, CPTC creates state-specific conformal prediction intervals for all states with nonzero probability. This allows us to generate predictions tailored to the current anticipated dynamics. Specifically, for every state $z \in \mathcal{Z}$ where $\hat{p}(z|x_t) > 0$, we first obtain the point prediction through forecaster model $\hat{f}(x_t, z)$, and then the state-specific prediction interval $\Gamma_{z,t}(x_t)$ by calibrating on nonconformity scores set \mathcal{S}_z to the adaptive confidence level of $1 - \alpha_{z,t}$.

Online Conformal Prediction. The exchangeability assumption does not apply to the online setting of our work. If we naively calibrate with online data, when the data distribution changes, the coverage probability may deviate from the target level $1 - \alpha$. Therefore, we follow online CP methods such as ACI [21] to address this problem by continuously updating $\alpha_{z,t}$ to track a surrogate miscoverage rate $\alpha_{z,t}^*$ as defined in eqn 8 (for each state z in our setting), an internal estimate of the target coverage.

$$\alpha_{z,t}^* = \sup\{\beta \in [0, 1] : M_{z,t}(\beta) \leq \alpha\}, \quad M_{z,t}(\alpha) = p(A(x_t, y_t) > Q^{1-\alpha}(\mathcal{S}_z)) \quad (8)$$

In eqn 8 $M_{z,t}$ measures the probability that the true label y_t falls outside the prediction set $\Gamma_{z,t}$ when the predicted state is z . We track $\alpha_{z,t}^*$ over time by updating the estimate $\alpha_{z,t}$ with online optimization. In our implementation, we use the simple update from ACI (Eqn 9), though more sophisticated approaches such as [4] may be used in its place without affecting the overall algorithm.

$$\alpha_{\hat{z}_t,t+1} \leftarrow \alpha_{\hat{z}_t,t} + \gamma \cdot (\alpha - \text{err}_t), \text{ where } \text{err}_t = \mathbb{1}\{y_t \notin \Gamma_t(x_t)\} \quad (9)$$

Here $\gamma > 0$ is the step size hyperparameter, and $\hat{z}_t \in \mathcal{Z}$, $\hat{z}_t \sim \hat{p}(z_t = z | x_t)$ is the state at time t sampled to be updated. The update rule increases $\alpha_{\hat{z}_t,t}$ if coverage is too conservative and decreases it otherwise, ensuring that for each $z \in \mathcal{Z}$, $\alpha_{z,t}$ converges to a value that maintains long-term coverage.

Aggregation. We aggregate the prediction intervals for each state $\Gamma_{z,t}(x_t)$ into one final set $\Gamma_t(x_t)$ using the weighted level set, as illustrated in Eqn 10.

$$\Gamma_t(x_t) := \left\{ y : \sum_{z \in \mathcal{Z}} \hat{p}(z_t = z | x_t) \mathbb{1}\{y \in \Gamma_{z,t}(x_t)\} \geq 1 - \alpha \right\} \quad (10)$$

Solution to the constraint in Eqn 10 is the minimal set that achieves the marginal coverage. It can be obtained through discretizing \mathcal{Y} into a fine grid and calculating probability mass at each point, which becomes computationally expensive as the grid resolution or dimensionality increases. For faster computation and in implementation, we approximate the weighted level-set in Eqn 10 by taking the union of intervals of the most probable states as in Eqn 11. The two aggregation strategies achieve similar results on all our datasets (Appendix B.4).

$$\Gamma_t(x_t) \approx \cup_{z \in \mathcal{Z}'_t} \Gamma_{z,t}(x_t), \text{ where } \mathcal{Z}'_t = \arg \min_{S \subseteq \mathcal{Z}} |S| \text{ s.t. } \sum_{z \in S} p(z_t = z | x_t) \geq 1 - \alpha. \quad (11)$$

Modularity and the role of the state model. Our algorithm can be viewed as running multiple instances of online conformal inference (lines 14-16 of algorithm 1), one for each underlying state, and adaptively aggregating the confidence intervals when the underlying mode is switching. This differs from the setting where the temporal correlation between residuals needs to be learned online [22, 11, 53, 10]; the state model allows us to leverage training data for such information. Our algorithm is modular - the state model, forecaster model, and the online adaptive conformal prediction algorithms all operate independently of each other. For example, the state-specific online coverage tracking (line 15 of algorithm 1) can be replaced with other online conformal prediction variants for data-specific applications. Our theoretical guarantees are agnostic to the forecaster models used.

4.2 Theoretical analysis

In this section, we discuss theoretical guarantees of the CPTC algorithm. In particular, we establish finite-sample validity under exchangeability, asymptotic coverage in dynamic (potentially nonstationary) settings, robustness to imperfect state classification, and faster adaptation to distribution shifts that align with model-predicted state changes. See Appendix A for proofs and theoretical details.

Finite-sample validity under exchangeability. We start by showing that when data is exchangeable (e.g. independent and identically distributed), CPTC achieves finite sample validity. The proof is standard based on showing the marginal coverage of split conformal prediction.

Proposition 4.1 (Finite-sample validity under exchangeability). *If the data (x_t, y_t) , $t \geq 1$ are exchangeable, prediction intervals obtained via Algorithm 1 satisfy $P(y_t \in \Gamma_t(x_t)) \geq 1 - \alpha$.*

Asymptotic validity without assumptions. Theorem 4.2 ensures that CPTC provides reliable uncertainty quantification in the long run *without assumptions on time-series stationarity or accurate state predictions*. (We do assume that the distribution of states and state predictions exhibit stable long-term behavior in Assumption 1.) We achieve our desiderata of Eqn 6 on average as T grows asymptotically, consistent with the results of other online conformal prediction works.

Theorem 4.2 (Asymptotic validity of CPTC). *For any sample size $T \geq 1$, the CPTC algorithm (with weighted average aggregation) achieves:*

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\text{err}_t] - \alpha \right| \leq \frac{1}{T} \sum_{z=1}^K \frac{1 - (1 - c\gamma)^{|T_z|}}{\gamma} |\alpha - \alpha_z^*|$$

Where:

- $err_t = \mathbb{1}[Y_t \notin \Gamma_t(x_t)]$, the miscoverage rate.
- $\mathcal{T}_z = \{t \in \{1, \dots, T\} : \hat{z}_t = z\}$, the set of time steps the predicted state is z .
- α_z^* is the optimal miscoverage target for timesteps \mathcal{T}_z .
- c is a miscoverage constant in Lemma A.1 [21].

Note the RHS of Eqn 4.2 decays at a rate of $\mathcal{O}(T)$ and satisfies $\lim_{T \rightarrow \infty} \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| = 0$.

Robustness to imperfect state prediction. In real-world applications, the state predictor $\hat{p}(z_t|x_t)$ that plays a large role in our algorithm is often imperfect. The CPTC algorithm accommodates these scenarios by tracking α_z^* for each *predicted* mode z with adaptive online updates. Theorem 4.3 quantifies the effect of imperfect state prediction by showing that the miscoverage rate is bounded by the product of the misclassification rate and state-specific deviations.

Theorem 4.3 (Finite-Sample Miscoverage Bound with Imperfect State Predictions). *For any sample size $T \geq 1$, the CPTC algorithm ensures that:*

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \epsilon \cdot \max_z \delta_{z,T}$$

Where $\epsilon = P(\hat{z}_t \neq z_t)$ is the error rate of the state predictions and $\delta_{z,T}$ is the miscoverage deviation from α for any predicted state z defined in Lemma A.1. For all z , $\delta_{z,T} \rightarrow 0$ as $|T| \rightarrow \infty$.

Theorem 4.3 demonstrates that CPTC can maintain asymptotic valid coverage even when state predictions are imperfect, with the bound tightening as state predictions improve or the number of observations per state increases. In practice, improving the accuracy of the state prediction model directly enhances the coverage performance of CPTC. However, even with imperfect state predictions, the algorithm’s adaptive calibration mechanism ensures that the overall coverage remains close to the desired level $1 - \alpha$ as experiments show in Section 5.2.

Why CPTC? Faster convergence under distribution shifts. The strength of the CPTC algorithm is its fast adaptation in dynamic environments with distributional changes. When a distribution shift aligns with a predicted state transition, the structure of CPTC allows the confidence level $\alpha_{z,t}$ to converge to the new target error rate faster compared to purely online methods.

Consider a data stream with a distribution shift occurring at time t_{shift} . Let α_j^* denote the optimal target error rate for mode j after the shift, and assume that the predicted state \hat{z}_t correctly reflects the shift: t_{shift} and $\alpha_z^* = \alpha_j^*$ for $t > t_{\text{shift}}$. For conciseness of notation in Theorem 4.4, let $\delta_{j,T}$ and $\delta_{ACI,T}$ denote the miscoverage rate *starting from* t_{shift} for CPTC and ACI respectively.

Theorem 4.4 (Miscoverage Ratio under State-Coincident Distribution Shift). *Under a state shift from i to j at time step t_{shift} , coinciding with predicted transition, the CPTC algorithm achieves faster convergence to the new target α_j^* compared to non-state-aware algorithm ACI at a ratio of:*

$$\frac{\delta_{j,T}}{\delta_{ACI,T}} \leq \frac{|\alpha_{j,t_{\text{shift}}-1} - \alpha_j^*|}{|\alpha_{t_{\text{shift}}-1} - \alpha_j^*|}$$

The numerator $|\alpha_{j,t_{\text{shift}}-1} - \alpha_j^*| \rightarrow 0$ as $T \rightarrow \infty$, when in-state adaptation has converged; but the denominator $|\alpha_{t_{\text{shift}}-1} - \alpha_j^*|$ remains greater than zero regardless of the length of the time series. This result highlights the advantage of segmenting nonconformity scores by state: the state-specific $\alpha_{z,t}$ aligns the target miscoverage updates with distribution shifts, allowing CPTC to quickly adjust its prediction intervals. Such accelerated adaptation allows for shorter miscoverage periods and is critical in real-world applications where rapid responses to changing conditions are necessary.

5 Experiments

Baselines. We selected the following baseline methods. **RED-SDS** (Recurrent Explicit Duration Switching Dynamical System) [8] is a state-of-the-art Bayesian model that learns state transitions (implementation details in Appendix B.1). We also use REDSDS as the base predictor for mode switching for our method in the experiments. **CP** is a generalization of conformal prediction to the online setting, also known as Online Sequential Split Conformal Prediction in [56]. We choose **ACI**

[21] to represent online CP algorithms leveraging online optimization. **SPCI** [53] and **HopCPT** [10] are state-of-the-art CP algorithms that learn adaptive predictive intervals by quantile regression and Modern Hopfield Networks learned from the time series respectively.

Metrics. We evaluate calibration and sharpness for each method. For *calibration*, we report the empirical coverage on the test set. Coverage should be as close to the desired confidence level $1 - \alpha$ as possible. Coverage = $\frac{1}{T} \sum_{t=1}^T \mathbb{1}(y_t \in \Gamma_t(x_t))$

For *sharpness*, we report the average width or area of the Prediction Intervals (PI). The measure should be as small as possible while being valid (coverage maintains above the specified confidence level). Width = $\frac{1}{T} \sum_{t=1}^T |\Gamma_t(x_t)|$

Datasets. The three synthetic datasets are designed with increasing randomness in mode changes, challenging the adaptivity of CPTC. **Bouncing Ball** is comprised of univariate time series encoding the height of a ball bouncing between two walls with constant velocity and elastic collisions, following [16, 8]. The two switching states are going up/down, each associated with a different level of Gaussian noise added to observation (Bouncing Ball obs), or the underlying dynamics (Bouncing Ball dyn) which induces uncertainty in the phase as well. **3-mode system** is a switching linear dynamical system with 3 switching states, where each mode samples from a Poisson distribution for duration.

For real-world datasets, the **Electricity** and **Traffic** datasets from [17] have hourly frequency and exhibit seasonality both in terms of the time series itself and volatility. The **honey bee trajectory** dataset [37] is the most complex, composed of 4-dimensional trajectories with length averaging to 900 frames, where the bees’ dance can be decomposed into “left turn”, “right turn” and “waggle”.

Illustrations of the datasets and detailed description can be found in Appendix B. The warm-start window for synthetic and real datasets are 50 and 100 time steps respectively (excluding bees, whose $w = 15$). The datasets represent a diverse set of scenarios to demonstrate the robustness and versatility of CPTC, capturing the complexity and variability of practical time-series forecasting problems.

5.1 Analysis of Coverage and Sharpness Results

Results for uncertainty quantification in table 1 show that we achieve significantly better validity on all datasets compared to baselines. RED-SDS’s intervals fail due to lack of calibration, a common pitfall of probabilistic methods (shown also in [53]). CP’s under-coverage shows that the data are not exchangeable. In the pretense of change points, ACI’s coverage fluctuates dramatically and did not converging within the horizon (more visualizations in Appendix B.5). For SPCI and HopCPT, which also guarantees asymptotic validity, the invalidity is likely because they are designed for large datasets ($T \geq 10,000$ for HopCPT), and did not have enough data to react or learn useful residual patterns for the relatively short prediction horizon of our datasets ($T = 200$ for synthetic datasets, 400 for electricity and traffic, and 60 for bees). Figure 3 provides a qualitative example of how the prediction intervals compare across different methods. Our method (purple shaded regions) does not over-cover during nonvolatile hours, nor under-cover during busy ones, as does ACI.

Table 1: Performance on synthetic and real-world datasets with target confidence $1 - \alpha = 0.9$ (for horizon $T = 200$ for synthetic datasets, 400 for electricity and traffic, and 60 for bee, mean \pm standard deviation of the test samples). Methods that are *invalid* (coverage below 90%) are grayed out. Our method achieves a high level of calibration (coverage is close to 90%) consistently.

		RED-SDS	CP	ACI	SPCI	HopCPT	Ours
Bouncing Ball obs.	Cov	59.14 \pm 23.75	39.27 \pm 6.77	76.97 \pm 1.27	37.90 \pm 9.78	59.45 \pm 12.51	91.03 \pm 1.37
	Width	2.12 \pm 0.05	1.71 \pm 0.93	3.67 \pm 0.91	1.45 \pm 0.85	1.98 \pm 0.90	2.13 \pm 2.15
Bouncing Ball dyn.	Cov	43.46 \pm 27.18	40.41 \pm 6.89	76.63 \pm 0.95	53.62 \pm 7.75	61.12 \pm 4.80	90.44 \pm 1.51
	Width	1.80 \pm 0.03	0.89 \pm 0.45	2.97 \pm 1.07	2.02 \pm 0.23	2.95 \pm 1.12	4.92 \pm 2.79
3-Mode System	Cov	92.98 \pm 3.01	63.97 \pm 5.03	86.25 \pm 0.84	81.72 \pm 2.88	89.85 \pm 2.90	93.35 \pm 3.51
	Width	2.08 \pm 0.16	1.46 \pm 0.51	2.61 \pm 0.53	3.94 \pm 0.85	9.07 \pm 1.23	13.68 \pm 2.49
Traffic	Cov	86.55 \pm 11.20	69.64 \pm 2.62	88.19 \pm 0.33	88.24 \pm 2.13	88.91 \pm 3.55	94.76 \pm 4.93
	Width	2.81 \pm 1.54	0.32 \pm 0.22	56.92 \pm 25.26	4.62 \pm 2.28	7.50 \pm 10.09	8.34 \pm 14.44
Electricity	Cov	75.46 \pm 14.23	68.88 \pm 3.12	87.94 \pm 0.34	86.62 \pm 2.97	86.50 \pm 2.66	92.62 \pm 6.14
	Width	87.47 \pm 119.95	58.62 \pm 456.89	28.28 \pm 169.33	371.69 \pm 305.86	155.71 \pm 130.43	166.88 \pm 129.07
Dancing Bees	Cov	90.79 \pm 5.03	65.58 \pm 7.05	75.16 \pm 3.33	79.2 \pm 0.3	72.11 \pm 1.84	93.43 \pm 11.42
	Width	0.27 \pm 0.05	0.12 \pm 0.01	1.77 \pm 1.38	1.77 \pm 1.38	1.06 \pm 0.51	1.15 \pm 0.56

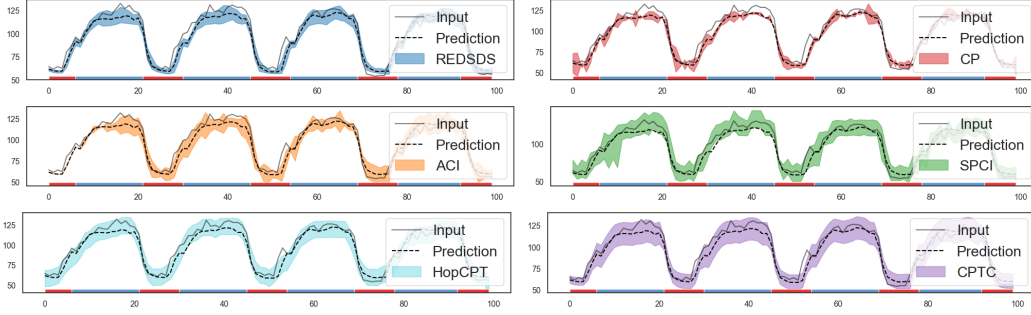


Figure 3: **Visualization of prediction intervals on the Electricity hourly demand dataset.** The red and blue bars in the bottom reflects the underlying switching state of day and night. Our method (purple) adapts to different levels of volatility between day and night, and achieves stabler coverage over time, whereas ACI (yellow) over-covers during the night and under-covers at change points.

5.2 Ablation Studies

Robustness to $p(\hat{z})$ accuracy. We conduct ablation studies on the 3 synthetic datasets of which we have ground truth underlying states. As shown in table 2, it is evident that increasing the error in the state prediction does not affect the coverage performance. This verifies the robustness of our algorithm and theoretical results (Theorem 4.2 and 4.3). Table 3 further shows that the error of state prediction reflects on the width of the prediction intervals: while remaining the same level of coverage, the better the state predictions, the sharper the intervals.

Table 2: Coverage of CPTC for using ground truth (GT) labels and with various levels of injected noise, on synthetic data with $T = 500$. We can see that coverage performance does not change significantly with error in state prediction.

	Ground truth state labels	GT w/ 20% error	GT w/ 50% error
BB obs.	90.00 \pm 0.32	89.94 \pm 0.29	89.92 \pm 0.28
BB dyn.	90.04 \pm 0.34	90.01 \pm 0.35	89.96 \pm 0.30
3-mode	89.34 \pm 0.87	89.40 \pm 0.91	89.64 \pm 0.80

Table 3: Interval width under the same setting. While the coverage guarantee is the same, the more accurate the state prediction is, the sharper the intervals.

	GT	w/ 20% err	w/ 50% err
BB obs.	0.49 \pm 0.04	1.58 \pm 0.24	1.81 \pm 0.23
BB dyn.	0.49 \pm 0.02	1.56 \pm 0.24	2.32 \pm 0.24
3-mode	0.95 \pm 0.04	0.97 \pm 0.04	0.97 \pm 0.04

Additional ablation studies are presented in Appendix B.4. On **Aggregation Methods** (by discretization as Eqn 10 or union as Eqn 11), we found that on our benchmark datasets, union is a close approximation of the true objective, and CPTC achieves similar width and coverage on a sampled subset. For **Long Horizon Forecasting** where $T \geq 10,000$, our method still achieves valid and stable coverage, but have wider PI width compared to SPCI and HopCPT. They achieve this by frequently re-training their models online, which CPTC does not and is less computationally demanding.

6 Conclusion and Discussion

In this paper, we introduced Conformal Prediction for Time-series with Change points (CPTC), a novel conformal prediction algorithm for uncertainty quantification in non-stationary time series. By leveraging state information in switching dynamics systems models, CPTC offers improvements in adaptivity over existing conformal prediction methods. Our theoretical guarantees ensure validity *without assumptions* on time-series stationarity or accurate state predictions. Empirical results corroborate our theory and demonstrate the effectiveness of CPTC across diverse synthetic and real-world datasets, achieving robust coverage under distributional shifts with comparable sharpness compared to state-of-the-art baselines. The adaptivity advantage is most pronounced in shorter time series, when baselines require more data to react or learn temporal correlations.

Limitations of CPTC include (1) slower convergence rate in scenarios with frequent or unpredictable state transitions, as the algorithm requires sufficient observations within each state to achieve optimal calibration and sharpness, (2) wider prediction interval width when applied to very long time series compared methods specialized for such scenarios (e.g. HopCPT), and (3) requirement of a state prediction model, making our method primarily applicable to SDS models and discrete SSMs. Future work could explore extending the framework to continuous states and the theoretical implications thereof, and studying the decision theoretic properties of using conformal prediction for safety-critical applications requiring real-time adaptation.

Acknowledgments and Disclosure of Funding

This work was supported in part by the U.S. Army Research Office under Army-ECASE award W911NF-07-R-0003-03, the U.S. Department Of Energy, Office of Science, IARPA HAYSTAC Program, NSF Grants SCALE MoDL-2134209, CCF-2112665 (TILOS), #2205093, #2146343, and #2134274, as well as CDC-RFA-FT-23-0069 from the CDC’s Center for Forecasting and Outbreak Analytics.

References

- [1] G. Ackerson and K. Fu. On state estimation in switching environments. *IEEE transactions on automatic control*, 15(1):10–17, 1970.
- [2] A. Alexandrov, K. Benidis, M. Bohlke-Schneider, V. Flunkert, J. Gasthaus, T. Januschowski, D. C. Maddix, S. Rangapuram, D. Salinas, J. Schulz, et al. Gluonts: Probabilistic and neural time series modeling in python. *Journal of Machine Learning Research*, 21(116):1–6, 2020.
- [3] S. Aminikhanghahi and D. J. Cook. A survey of methods for time series change point detection. *Knowledge and information systems*, 51(2):339–367, 2017.
- [4] A. Angelopoulos, E. Candes, and R. J. Tibshirani. Conformal pid control for time series prediction. *Advances in neural information processing systems*, 36, 2024.
- [5] A. N. Angelopoulos, R. F. Barber, and S. Bates. Theoretical foundations of conformal prediction. *arXiv preprint arXiv:2411.11824*, 2024.
- [6] A. N. Angelopoulos and S. Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- [7] A. N. Angelopoulos, S. Bates, A. Fisch, L. Lei, and T. Schuster. Conformal risk control. *arXiv preprint arXiv:2208.02814*, 2022.
- [8] A. F. Ansari, K. Benidis, R. Kurle, A. C. Turkmen, H. Soh, A. J. Smola, B. Wang, and T. Januschowski. Deep explicit duration switching models for time series. *Advances in Neural Information Processing Systems*, 34:29949–29961, 2021.
- [9] M. Aoki. *State space modeling of time series*. Springer Science & Business Media, 2013.
- [10] A. Auer, M. Gauch, D. Klotz, and S. Hochreiter. Conformal prediction for time series with modern hopfield networks. *Advances in Neural Information Processing Systems*, 36:56027–56074, 2023.
- [11] O. Bastani, V. Gupta, C. Jung, G. Noarov, R. Ramalingam, and A. Roth. Practical adversarial multivalid conformal prediction. *Advances in Neural Information Processing Systems*, 35:29362–29373, 2022.
- [12] K. Benidis, S. S. Rangapuram, V. Flunkert, Y. Wang, D. Maddix, C. Turkmen, J. Gasthaus, M. Bohlke-Schneider, D. Salinas, L. Stella, et al. Deep learning for time series forecasting: Tutorial and literature survey. *ACM Computing Surveys*, 55(6):1–36, 2022.
- [13] M. Cetin and G. Comert. Short-term traffic flow prediction with regime switching models. *Transportation Research Record*, 1965(1):23–31, 2006.
- [14] Y. Chen and H. V. Poor. Learning mixtures of linear dynamical systems. In *International conference on machine learning*, pages 3507–3557. PMLR, 2022.
- [15] J. J. Cherian, I. Gibbs, and E. J. Candès. Large language model validity via enhanced conformal prediction methods. *arXiv preprint arXiv:2406.09714*, 2024.
- [16] Z. Dong, B. Seybold, K. Murphy, and H. Bui. Collapsed amortized variational inference for switching nonlinear dynamical systems. In *International Conference on Machine Learning*, pages 2638–2647. PMLR, 2020.

- [17] D. Dua and E. Karra Taniskidou. UCI machine learning repository, 2017. Accessed: September 4, 2025.
- [18] S. Fan, L. Chen, and W.-J. Lee. Machine learning based switching model for electricity load forecasting. *Energy Conversion and Management*, 49(6):1331–1344, 2008.
- [19] Y. Gal, J. Hron, and A. Kendall. Concrete dropout. In *NIPS*, pages 3581–3590, 2017.
- [20] M. Garnelo, D. Rosenbaum, C. Maddison, T. Ramalho, D. Saxton, M. Shanahan, Y. W. Teh, D. Rezende, and S. M. A. Eslami. Conditional neural processes. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [21] I. Gibbs and E. Candes. Adaptive conformal inference under distribution shift. *Advances in Neural Information Processing Systems*, 34, 2021.
- [22] I. Gibbs and E. J. Candès. Conformal inference for online prediction with arbitrary distribution shifts. *Journal of Machine Learning Research*, 25(162):1–36, 2024.
- [23] J. Glaser, M. Whiteway, J. P. Cunningham, L. Paninski, and S. Linderman. Recurrent switching dynamical systems models for multiple interacting neural populations. *Advances in neural information processing systems*, 33:14867–14878, 2020.
- [24] M. Johnson, D. Duvenaud, A. B. Wiltschko, R. P. Adams, and S. R. Datta. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in Neural Information Processing Systems*, 2016.
- [25] R. E. Kalman. A new approach to linear filtering and prediction problems. 1960.
- [26] N. V. Karakatsani and D. W. Bunn. Intra-day and regime-switching dynamics in electricity price formation. *Energy Economics*, 30(4):1776–1797, 2008.
- [27] S. Kato, H. S. Kaplan, T. Schrödel, S. Skora, T. H. Lindsay, E. Yemini, S. Lockery, and M. Zimmer. Global brain dynamics embed the motor command sequence of caenorhabditis elegans. *Cell*, 163(3):656–669, 2015.
- [28] B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [29] J. Lei, M. G’Sell, A. Rinaldo, R. J. Tibshirani, and L. Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.
- [30] B. Lim, S. Ö. Arık, N. Loeff, and T. Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. In *International Journal of Forecasting*, volume 37, pages 1748–1764. Elsevier, 2021.
- [31] L. Lindemann, M. Cleaveland, G. Shim, and G. J. Pappas. Safe planning in dynamic environments using conformal prediction. *IEEE Robotics and Automation Letters*, 2023.
- [32] S. Linderman. The nuts and bolts of probabilistic state space models: Part i – foundations. Lecture slides, 2023. <https://research.gatech.edu/sites/default/files/inline-files/2023-03-28%20SSM%20Tutorial%20GATech%20RS.pdf>.
- [33] S. W. Linderman, A. C. Miller, R. P. Adams, D. M. Blei, L. Paninski, and M. J. Johnson. Recurrent switching linear dynamical systems. *arXiv preprint arXiv:1610.08466*, 2016.
- [34] Y. Liu, S. Magliacane, M. Kofinas, and E. Gavves. Graph switching dynamical systems. In *International Conference on Machine Learning*, pages 21867–21883. PMLR, 2023.
- [35] G. Montanez, S. Amizadeh, and N. Laptev. Inertial hidden markov models: Modeling change in multivariate time series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- [36] Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. Patchtst: A time series transformer for forecasting, anomaly detection, and classification. *arXiv preprint arXiv:2306.09364*, 2023.

- [37] S. M. Oh, J. M. Rehg, T. Balch, and F. Dellaert. Learning and inferring motion patterns using parametric segmental switching linear dynamic systems. *International Journal of Computer Vision*, 77:103–124, 2008.
- [38] M. I. Rabinovich, P. Varona, A. I. Selverston, and H. D. Abarbanel. Dynamical principles in neuroscience. *Reviews of modern physics*, 78(4):1213–1265, 2006.
- [39] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.
- [40] Z. Song, Y. Jiang, and Z. Zhang. Short-term wind speed forecasting with markov-switching model. *Applied Energy*, 130:103–112, 2014.
- [41] K. Stankevičiūtė, A. Alaa, and M. van der Schaar. Conformal time-series forecasting. In *Advances in Neural Information Processing Systems*, 2021.
- [42] J. P. Strachan, A. C. Torrezan, G. Medeiros-Ribeiro, and R. S. Williams. Measuring the switching dynamics and energy efficiency of tantalum oxide memristors. *Nanotechnology*, 22(50):505402, 2011.
- [43] S. Sun, W. Chen, Z. Zhou, S. Fereidooni, E. Jortberg, and R. Yu. Data-driven simulator for mechanical circulatory support with domain adversarial neural process. *arXiv preprint arXiv:2405.18536*, 2024.
- [44] S. H. Sun and R. Yu. Copula conformal prediction for multi-step time series prediction. In *The Twelfth International Conference on Learning Representations*, 2023.
- [45] R. J. Tibshirani, R. Foygel Barber, E. Candes, and A. Ramdas. Conformal prediction under covariate shift. *Advances in neural information processing systems*, 32, 2019.
- [46] D. Volkhonskiy, E. Burnaev, I. Nourtdinov, A. Gammerman, and V. Vovk. Inductive conformal martingales for change-point detection. In *Conformal and Probabilistic Prediction and Applications*, pages 132–153. PMLR, 2017.
- [47] V. Vovk, A. Gammerman, and G. Shafer. *Algorithmic learning in a random world*. Springer Science & Business Media, 2005.
- [48] V. Vovk, I. Petej, I. Nourtdinov, E. Ahlberg, L. Carlsson, and A. Gammerman. Retrain or not retrain: Conformal test martingales for change-point detection. In *Conformal and Probabilistic Prediction and Applications*, pages 191–210. PMLR, 2021.
- [49] B. Wang, J. Lu, Z. Yan, H. Luo, T. Li, Y. Zheng, and G. Zhang. Deep uncertainty quantification: A machine learning approach for weather forecasting. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2087–2095, 2019.
- [50] Y. Wang and C.-Q. Xu. Actively q-switched fiber lasers: Switching dynamics and nonlinear processes. *Progress in Quantum Electronics*, 31(3-5):131–216, 2007.
- [51] D. Wu, L. Gao, M. Chinazzi, X. Xiong, A. Vespignani, Y.-A. Ma, and R. Yu. Quantifying uncertainty in deep spatiotemporal forecasting. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 1841–1851, 2021.
- [52] C. Xu and Y. Xie. Conformal prediction interval for dynamic time-series. In *International Conference on Machine Learning*. PMLR, 2021.
- [53] C. Xu and Y. Xie. Sequential predictive conformal inference for time series. In *International Conference on Machine Learning*, pages 38707–38727. PMLR, 2023.
- [54] Y. Yamamoto, A. Kijima, M. Okumura, K. Yokoyama, and K. Gohara. A switching hybrid dynamical system: Toward understanding complex interpersonal behavior. *Applied Sciences*, 9(1):39, 2018.

- [55] G. Yu and C. Zhang. Switching arima model based forecasting for traffic flow. In *2004 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 2, pages ii–429. IEEE, 2004.
- [56] M. Zaffran, O. Féron, Y. Goude, J. Josse, and A. Dieuleveut. Adaptive conformal predictions for time series. In *International Conference on Machine Learning*, pages 25834–25866. PMLR, 2022.
- [57] Y. Zhou, L. Lindemann, and M. Sesia. Conformalized adaptive forecasting of heterogeneous trajectories. *arXiv preprint arXiv:2402.09623*, 2024.

A Theoretical Results and Proofs

A.1 Setup and Notation

Let $z_t \in \mathcal{Z} = \{1, \dots, K\}$ denote the unobserved discrete mode, $x_t \in \mathcal{X}$ denote the continuous state, and $y_t \in \mathcal{Y}$ the observation at time t . We assume we have access to a probabilistic model $\hat{p}(z_t | X_{1:t-1})$ that estimates the mode \hat{z}_t .

For a target coverage rate $1 - \alpha$, let $\Gamma_t(X_t)$ be the prediction set at time t . We define the miscoverage indicator M_t as follows, and aim to bound the overall miscoverage rate over T time steps:

$$M_t = \frac{1}{T} \sum_{t=1}^T err_t, \quad \text{where} \quad err_t := \mathbb{1}(Y_t \notin \Gamma_t(X_t)),$$

A.2 Proofs: Lemma

In practical scenarios, the state prediction model may not be perfectly accurate. We now analyze the impact of imperfect state predictions on the coverage guarantees of the CPTC algorithm.

We derive the following lemma from the convergence result of Proposition 4.1 in [21] with two additional assumptions also according to [21]: (1) the existence of a single fixed optimal target $\alpha^* \in [0, 1]$, and (2) $\mathbb{E}[err_t | a_t] = M(a_t)$. These are reasonable assumptions that allow us to quantify distribution shifts and are necessary for analyzing adaptation behavior (whereas proposition 4.1 itself is a worst-case upper bound.) We refer the readers to the original paper for the complete derivation.

Lemma A.1 (Adaptive Conformal Inference Miscoverage Bound Within Predicted States). *Let $\mathcal{T}_z = \{t \in \{1, \dots, T\} : \hat{z}_t = z\}$ denote the set of times when the predicted state is z . For any $T > 1$ and predicted state $\hat{z}_t = z$, the Adaptive Conformal Inference (ACI) algorithm ensures that there exists a constant $\delta_{z,T}$ such that:*

$$\left| \frac{1}{|\mathcal{T}_z|} \sum_{t \in \mathcal{T}_z} \mathbb{E}[err_t] - \alpha \right| \leq \delta_{z,T}$$

Explicitly,

$$\delta_{z,T} = \frac{1 - (1 - c\gamma)^{|\mathcal{T}_z|}}{|\mathcal{T}_z|\gamma} |\alpha - \alpha_z^*|$$

where

- α_z^* is the optimal miscoverage target for timesteps \mathcal{T}_z .
- γ is the step size
- c the miscoverage constant.

Note that $\delta_{z,T}$ decays at a rate of $\mathcal{O}(T)$ and satisfies $\lim_{T \rightarrow \infty} \delta_{z,T} = 0$.

A.3 Proof of Theorem 4.2

Theorem (4.2 Asymptotic validity of CPTC). For any sample size $T \geq 1$, without placing any assumptions, the CPTC algorithm ensures that:

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \frac{1}{T} \sum_{z=1}^K |\mathcal{T}_z| \delta_{z,T} = \frac{1}{T} \sum_{z=1}^K \frac{1 - (1 - c\gamma)^{|\mathcal{T}_z|}}{\gamma} |\alpha - \alpha_z^*|$$

which decays at a rate of $\mathcal{O}(T)$ and satisfies

$$\lim_{T \rightarrow \infty} \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| = 0.$$

To prove theorem 4.2, we assume that both the distribution of the underlying modes z and the predictions \hat{z} is stationary (Assumption 1).

Assumption 1 (Stationary Distribution of States). The sequence of true states $\{z_t\}$ has a stationary distribution $\pi(z)$, and the sequence of predicted states $\{\hat{z}_t\}$ also has a stationary distribution $\hat{\pi}(z)$. Specifically:

$$\forall z, \in \mathcal{Z}, \quad \lim_{t \rightarrow \infty} p(z_t = z) = \pi(z), \quad \lim_{t \rightarrow \infty} p(\hat{z}_t = z) = \hat{\pi}(z),$$

This assumption is well-justified in practical scenarios where the sequence of observations and predictions arises from processes that, while non-stationary over short intervals, exhibit stable long-term behavior. In time-series or dynamic environments, regularities in data generation or prediction models often lead to empirically observed stationarity in distributions. While restrictive, this assumption aligns with prior work in sequential and online prediction frameworks, where stationarity assumptions are standard to derive theoretical guarantees.

Proof of Theorem 4.2. Let $\mathcal{T}_z = \{t \in \{1, \dots, T\} : \hat{z}_t = z\}$ represent the set of timesteps when the predicted state is z . For the total number of timesteps T , the overall miscoverage error can be expressed as:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] &= \frac{1}{T} \sum_{z=1}^K \sum_{t \in \mathcal{T}_z} \mathbb{E}[err_t] \\ &= \frac{1}{T} \sum_{z=1}^K |\mathcal{T}_z| \left(\frac{1}{|\mathcal{T}_z|} \sum_{t \in \mathcal{T}_z} \mathbb{E}[err_t] \right). \end{aligned}$$

From Lemma A.1, for any predicted state z , the miscoverage within \mathcal{T}_z satisfies:

$$\left| \frac{1}{|\mathcal{T}_z|} \sum_{t \in \mathcal{T}_z} \mathbb{E}[err_t] - \alpha \right| \leq \delta_{z,T},$$

where:

$$\delta_{z,T} = \frac{1 - (1 - c\gamma)^{|\mathcal{T}_z|}}{|\mathcal{T}_z|\gamma} |\alpha - \alpha_z^*|.$$

Substituting this bound into the overall miscoverage error, we obtain:

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \frac{1}{T} \sum_{z=1}^K |\mathcal{T}_z| \delta_{z,T}.$$

Expanding $\delta_{z,T}$, the overall bound becomes:

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \frac{1}{T} \sum_{z=1}^K \frac{1 - (1 - c\gamma)^{|\mathcal{T}_z|}}{\gamma} |\alpha - \alpha_z^*|.$$

By assumption 1, as $T \rightarrow \infty$, the size of $|\mathcal{T}_z|$ for each state z grows proportionally to T . The term $1 - (1 - c\gamma)^{|\mathcal{T}_z|}$ approaches 1, while $|\mathcal{T}_z|/T$ converges to the relative proportion of time spent in state z , denoted as p_z . Hence, as $T \rightarrow \infty$, the overall bound $\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \rightarrow 0$.

Thus, the CPTC algorithm ensures that the overall error converges to the target coverage level α at a rate of $\mathcal{O}(T)$, satisfying:

$$\lim_{T \rightarrow \infty} \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| = 0.$$

□

A.4 Proof of Theorem 4.3

Theorem (4.3 Finite-Sample Miscoverage Bound with Imperfect State Predictions). For any sample size $T \geq 1$, the CPTC algorithm ensures that:

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \epsilon \cdot \max_z \delta_{z,T}$$

where:

- $err_t := \mathbb{1}(Y_t \notin \Gamma_t(X_t))$, the miscoverage rate.
- $\mathcal{T}_z = \{t \in \{1, \dots, T\} : \hat{z}_t = z\}$, the set of times when the predicted state is z .
- $\epsilon = \mathbb{P}(\hat{z}_t \neq z_t)$ is the misclassification rate of the state predictions.
- $\delta_{z,T}$ is the deviation from α within any predicted state z as in Lemma A.1.

Proof of Theorem 4.3. First we will partition time into correct vs. incorrect predictions. Define:

$$\mathcal{C} = \{t : \hat{z}_t = z_t\}, \quad \mathcal{I} = \{t : \hat{z}_t \neq z_t\}, \quad \text{and} \quad \frac{|\mathcal{I}|}{T} = \epsilon.$$

Calibrated expected coverage on correctly predicted states \mathcal{C} . Since there are no distribution shifts within each state, we have the same expected error from calibration data. Therefore,

$$\mathbb{E}[err_t] = \alpha \implies \mathbb{E}[err_t] - \alpha = 0.$$

Bounded deviation on incorrectly predicted states \mathcal{I} . On \mathcal{I} , we have $|\mathbb{E}[err_t] - \alpha| \leq \delta_{z,T}$. Therefore,

$$\sum_{t=1}^T |\mathbb{E}[err_t] - \alpha| = \sum_{t \in \mathcal{C}} |\mathbb{E}[err_t] - \alpha| + \sum_{t \in \mathcal{I}} |\mathbb{E}[err_t] - \alpha| \leq 0 + |\mathcal{I}| \delta_{z,T} = |\mathcal{I}| \delta_{z,T}.$$

Since $\frac{|\mathcal{I}|}{T} = \epsilon$, dividing by T yields

$$\frac{1}{T} \sum_{t=1}^T |\mathbb{E}[err_t] - \alpha| \leq \epsilon \delta_{z,T}.$$

Combine to achieve the overall bound. By the triangle inequality, we know that

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \frac{1}{T} \sum_{t=1}^T |\mathbb{E}[err_t] - \alpha|.$$

Hence, from the previous step,

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \epsilon \delta_{z,T}.$$

Since on \mathcal{I} the predicted state \hat{z}_t could vary among different z values, we take $\delta_{z,T} \leq \max_z \delta_{z,T}$. Thus,

$$\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[err_t] - \alpha \right| \leq \epsilon \max_z \delta_{z,T}.$$

□

A.5 Corollaries

Theorem 4.3 indicates that the overall miscoverage rate deviates from the desired level α by at most $\epsilon \cdot \max_z \delta_{z,T}$. The misclassifications rate ϵ contributes directly to the deviation, and the calibration error $\delta_{z,T}$ within each predicted state decreases as more data becomes available. As T increases and $\max_z \delta_{z,T} \rightarrow 0$, the deviation is primarily governed by the misclassification rate ϵ .

Corollary A.2 (Zero miscoverage rate under Accurate State Prediction). *If predicted state probabilities are accurate $\hat{p}(z_t|x_{1:t-1}) = p(z_t|x_{1:t-1})$, then $\epsilon = 0$, therefore $\mathbb{E}[\text{err}_t] = \alpha$ for all T , achieving optimal performance of online conformal prediction.*

Corollary A.3 (Asymptotic calibration under eventually correct State Prediction). *If the predicted state probabilities $\hat{p}(z_t|x_{1:t-1})$ converge in probability to the true state distribution $p(z_t|x_{1:t-1})$ as $T \rightarrow \infty$. With $T \rightarrow 0$, $\epsilon \rightarrow 0$ as well. Hence, the overall miscoverage rate converges asymptotically to 0.*

A.6 Proof of Theorem 4.4

Consider a data stream with a distribution shift occurring at time t_{shift} . Let α_j^* denote the optimal target error rate for mode j after the shift, and assume that the predicted state \hat{z}_t correctly reflects the shift: t_{shift} and $\alpha_z^* = \alpha_j^*$ for $t > t_{\text{shift}}$. For conciseness of notation in Theorem 4.4, let $\delta_{j,T}$ and $\delta_{\text{ACI},T}$ denote the miscoverage rate *starting from* t_{shift} for CPTC and ACI respectively.

Theorem (4.4 Miscoverage Ratio under State-Coincident Distribution Shift). Under a state shift from i to j at time step t_{shift} , coinciding with predicted transition, the CPTC algorithm achieves faster convergence to the new target α_j^* compared to non-state-aware algorithm ACI at a ratio of:

$$\frac{\delta_{j,T}}{\delta_{\text{ACI},T}} \leq \frac{|\alpha_{j,t_{\text{shift}}-1} - \alpha_j^*|}{|\alpha_{t_{\text{shift}}-1} - \alpha_j^*|}$$

Proof of Theorem 4.4. We want to compare the convergence behavior of CPTC and a non-segmented online conformal prediction algorithm (e.g., ACI) when a distribution shift occurs at time t_{shift} and this shift coincides with a predicted state transition from i to j .

We assume that State prediction coincident shift, meaning $\alpha_z^* = \alpha_i^*$ for $t \leq t_{\text{shift}}$ and $\alpha_z^* = \alpha_j^*$ for $t > t_{\text{shift}}$.

From Lemma A.1 and Theorem 4.2, we have bounds on these deviations. Taking the ratio of the bounds, we have:

$$\frac{\delta_{j,T}}{\delta_{\text{ACI},T}} \leq \frac{c \cdot \frac{1}{\gamma} |\alpha_{j,t_{\text{shift}}-1} - \alpha_j^*|}{c \cdot \frac{1}{\gamma} |\alpha_{t_{\text{shift}}-1} - \alpha_j^*|} = \frac{|\alpha_{j,t_{\text{shift}}-1} - \alpha_j^*|}{|\alpha_{t_{\text{shift}}-1} - \alpha_j^*|}$$

□

We elaborate on the implication of this result. The ratio indicates that the CPTC achieves lower miscoverage rate when the initial target $\alpha_{j,t_{\text{shift}}-1}$ for state j in CPTC is closer to the new optimal target α_j^* than the global target $\alpha_{t_{\text{shift}}-1}$ used by ACI. This is because CPTC maintains separate targets for each state, allowing it to better track state-specific optimal targets. The base model takes around 2 hours of training time on the GPU

B Experiment Details and additional results

B.1 Hyperparameters for the base forecasting model

We train our RED-SDS on the synthetic datasets, and use the model checkpoints provided by the authors for the real-world datasets. The model architecture consists of a discrete switching component with $K \in \{2, 3\}$ categories and a continuous state space with dimensionality $d_x \in \{2, 4\}$. For training, we employ the ELBOv2 objective with a learning rate $\eta \in [5 \times 10^{-3}, 7 \times 10^{-3}]$, warmup steps of 1000, and gradient clipping at 10.0. The model uses a batch size $B \in \{32, 50\}$ and is trained for $T \in \{20, 000, 30, 000\}$ steps. The continuous transition and emission models are parameterized by nonlinear MLPs with hidden dimensions $h = 32$, while the inference network uses either a bidirectional RNN or transformer with embedding dimensions $d_e = 4$. For real-world datasets, we apply target transformation and Jacobian correction, while synthetic datasets use raw observations. The model’s capacity is controlled through weight decay $\lambda = 10^{-5}$ and MLP hidden dimensions $h \in \{8, 32, 64\}$ depending on the component. For forecasting, the model is trained to forecast a window of $t = 50$ time steps for synthetic datasets (bouncing ball) and a windows of for real-world datasets (electricity, traffic) as specified by their respective metadata. The model generates $N = 100$ Monte Carlo samples for each forecast, with deterministic latent states (z_t) and stochastic continuous states (x_t) to capture both the switching dynamics and the inherent uncertainty in the predictions. We segment all datasets by a 70/10/20 train/validation/test split. Conforaml prediction results are reported on the test set only.

B.2 Computational Resources

All experiments are done on a server machine with an Nvidia A100 GPU, with some data processing and analytics performed on a Apple Macbook Pro laptop computer with M1 chip. As running time vary greatly depending on implementation and hardware, we provide a rough range of algorithm runtime. Training the RED-SDS base model (1 seed) takes 2-5 hours for our datasets, and inference take around 15-20 minutes for all datasets. Among the conformal prediction baselines, SPCI is the most computationally demanding, taking 8-10 hours for inference on our benchmarks. HopCPT follows, taking around 1 hour. CP, ACI, and our method SPCI do not require extensive inference-time computation, completing all inference within 5 minutes.

B.3 Dataset descriptions, continued

All datasets are visualized in Figure 5 with their underlying states color coded.

Electricity and Traffic. The Electricity and Traffic dataset are from the UCI machine learning repository [17] and we use the train/test split from GluonTS [2]. The electricity dataset consists of hourly electricity consumption data for 370 customers, each with 5,833 observations. The traffic dataset consists of hourly occupancy rates from 963 road sensors, each with 4,000 data points.

Dancing Bees. The honey bee trajectory dataset is from [37]. Honey bees convey information about the location and distance to a food source through a dance carried out inside the hive. This dance consists of three distinct phases: “left turn”, “right turn”, and “waggle”. The duration and orientation during the waggle phase represent the distance and direction to the food source. Experiments are conducted on 6 trajectories, with lengths of 1058, 1125, 1054, 757, 609, and 814 frames, respectively. In this paper, the prediction interval is on the first two dimensions (the coordinates of the bee) only, but all dimensions are used for prediction.

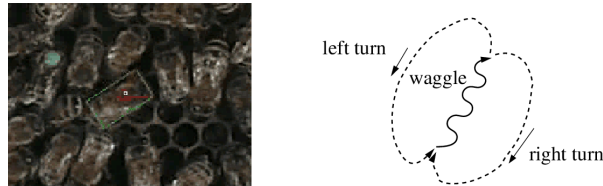


Figure 4: Dancer bees are tracked by a appearance based tracker from video sequences. The tracked bee is shown in green rectangle in the left figure above. The right figure shows a stylized bee dance through which bees talk to the other bees about the orientation and distance to the food sources.

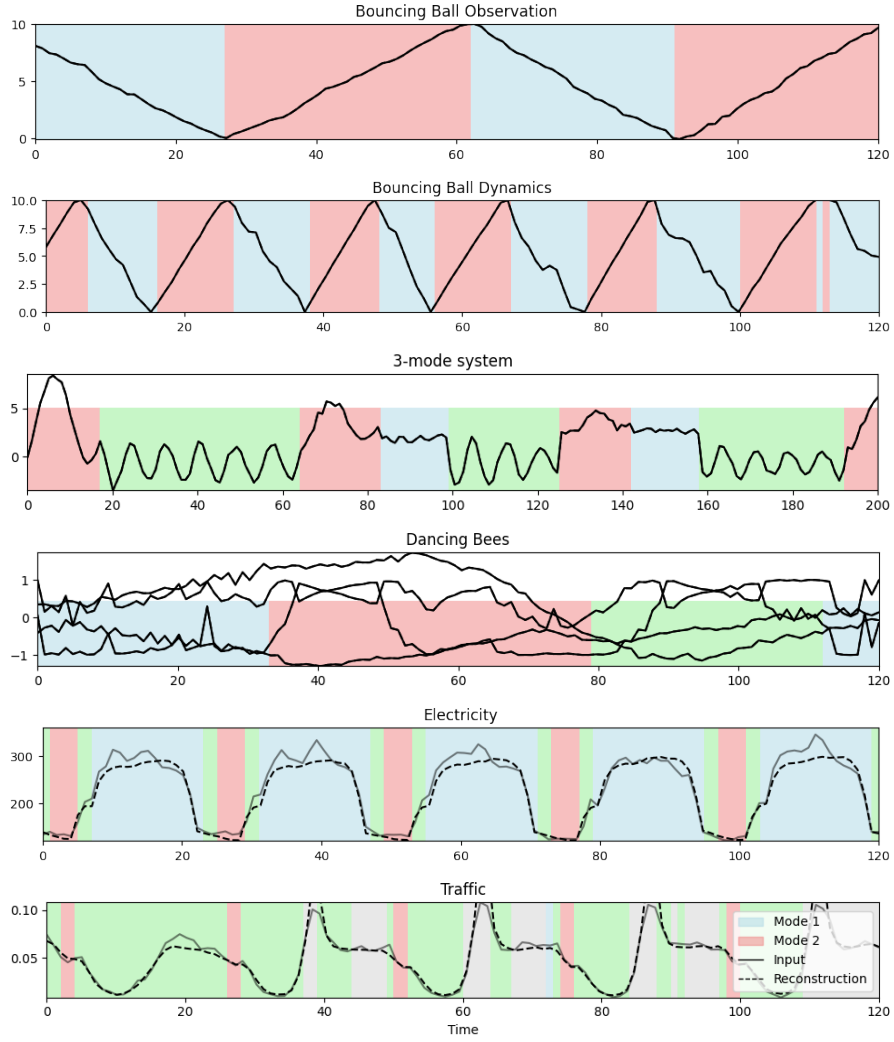


Figure 5: Examples of datasets and prediction results. Dotted line represents prediction whereas solid lines represents ground truth; background colors represent the different operating modes segmented by our base predictor RED-SDS.

B.4 Additional Experiments

Longer horizon Data Result for the same set of experiments performed on longer horizon ($T = 10,000$ for synthetic datasets and $T = 4,000$ for Electricity and Traffic datasets) is presented in table 4. We added the setting where the added noise to the bouncing ball dataset is varying over the course of the time series, which is more challenging than the original settings. In the longer horizon case, both SPCI and HopCPT achieve their state of the art performance. Our method remains valid yet achieves less sharp interval width.

Table 4: Longer time horizon ($T = 10,000$ for synthetic datasets and $T = 4,000$ for Electricity and Traffic datasets). Performance in synthetic and real-world datasets with target confidence $1 - \alpha = 0.9$. Methods that are *invalid* (coverage below 90%) are grayed out. Our method achieves a high level of calibration (coverage is close to 90%) consistently.

		RED-SDS	CP	ACI	SPCI	HopCPT	Ours
Bouncing Ball obs.	Cov	89.31 \pm 5.15	81.55 \pm 7.10	89.25 \pm 1.40	90.05 \pm 0.04	90.02 \pm 0.02	90.53 \pm 0.23
	Width	2.65 \pm 0.11	2.65 \pm 0.85	3.53 \pm 0.81	1.95 \pm 0.12	1.99 \pm 0.08	2.25 \pm 0.54
Bouncing Ball dyn.	Cov	45.12 \pm 25.50	42.88 \pm 6.95	87.95 \pm 1.05	90.04 \pm 0.04	90.01 \pm 0.04	90.05 \pm 0.04
	Width	1.67 \pm 0.13	0.95 \pm 0.51	3.10 \pm 1.04	3.24 \pm 0.14	3.07 \pm 0.10	4.18 \pm 0.91
3-Mode System	Cov	88.52 \pm 3.53	65.24 \pm 5.51	88.75 \pm 0.90	90.05 \pm 0.02	90.02 \pm 0.04	90.08 \pm 0.05
	Width	2.56 \pm 0.21	2.33 \pm 0.63	8.74 \pm 0.61	8.07 \pm 0.94	9.52 \pm 1.31	14.10 \pm 2.63
Bouncing Ball obs. varying	Cov	60.83 \pm 22.53	40.25 \pm 6.52	88.60 \pm 1.25	90.02 \pm 0.04	90.04 \pm 0.00	90.10 \pm 0.03
	Width	2.10 \pm 0.45	1.73 \pm 0.91	3.70 \pm 0.95	2.11 \pm 0.82	2.14 \pm 0.73	2.41 \pm 1.22
Bouncing Ball dyn. varying	Cov	44.25 \pm 26.11	41.73 \pm 7.12	88.05 \pm 1.00	90.05 \pm 0.03	90.02 \pm 0.04	90.08 \pm 0.03
	Width	1.85 \pm 0.12	0.92 \pm 0.55	3.13 \pm 1.14	2.92 \pm 0.21	3.05 \pm 0.19	3.40 \pm 0.32
Traffic	Cov	87.14 \pm 10.51	70.52 \pm 2.83	89.15 \pm 0.38	90.01 \pm 0.04	90.03 \pm 0.03	90.02 \pm 0.03
	Width	2.91 \pm 1.62	0.43 \pm 0.25	55.10 \pm 24.13	4.81 \pm 2.42	7.84 \pm 10.53	8.55 \pm 14.11
Electricity	Cov	76.10 \pm 13.52	69.53 \pm 3.21	89.05 \pm 0.33	90.02 \pm 0.03	90.00 \pm 0.02	90.43 \pm 0.03
	Width	91.30 \pm 125.10	62.10 \pm 451.20	31.40 \pm 171.30	162.02 \pm 11.40	161.70 \pm 15.20	173.38 \pm 31.10

Aggregation methods. In this study we compare the two different aggregation methods. First is the case where we chose the minimum size, i.e. $\Gamma_t(x_t) = \text{Eqn 10}$. Algorithmically, we discretize the output space \mathcal{Y} into intervals of length 0.02, and calculate probability mass for each interval. We refer to this method as *min*. Aggregation by union numbers (the same as in the main text using Eqn 11) are referred to as *union*. Results are present in table 5 where the performance are very similar. When there are two modes (bouncing ball cases) the prediction intervals are almost identical.

Table 5: Ablation Study on two different Union (*union*, eqn 11) and grid-discretization (*min*, eqn 10).

	BB obs.	BB dyn.	3-mode system	traffic	electricity	bees
Coverage (<i>union</i>)	91.03	90.44	93.35	94.76	92.62	93.43
Coverage (<i>min</i>)	91.03	90.35	91.50	90.87	89.82	90.28
Width (<i>union</i>)	2.13	4.92	13.68	8.34	166.88	1.15
Width (<i>min</i>)	2.13	4.85	12.43	7.95	153.57	1.08

B.5 Qualitative Results

We provide additional visualizations to present qualitative comparisons of CPTC compared to baseline models. Figures 6, 7, and 8 show consistent results with the analysis in section 5, where CPTC shows stronger adaptivity and validity for uncertainty quantification in time series with these nonstationary dynamics.

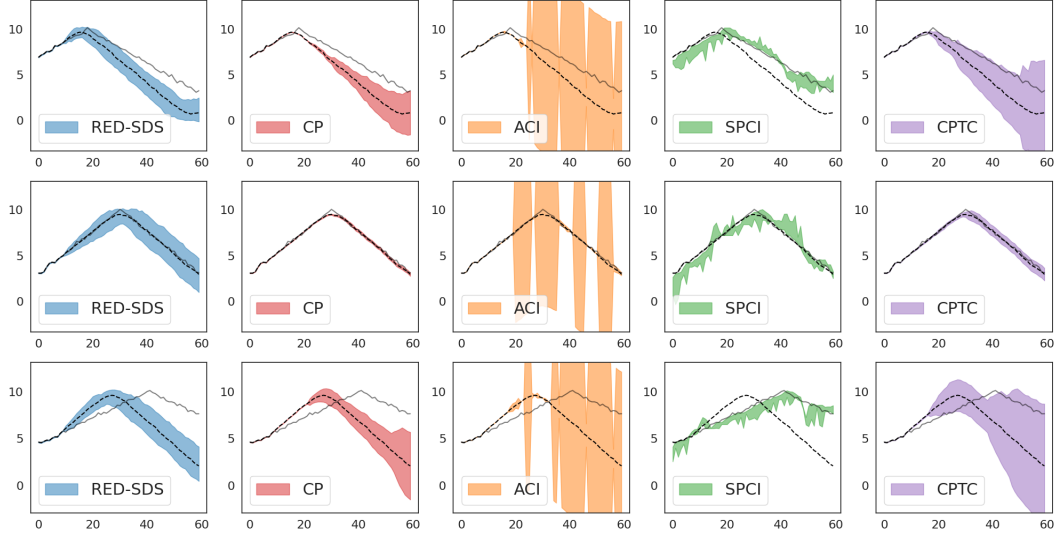


Figure 6: **Qualitative comparison of the prediction interval on the bouncing ball w/ dynamics noise dataset.** The dashed line is the predicted state whereas the solid gray line is the observation. We see that our CPTC algorithm can increase the interval in anticipation of an uncertain state change (top panel), is stabler when the prediction is accuracy compared to ACI and SPCI (middle panel), and is adaptive to inaccurate predictions (botom panel).

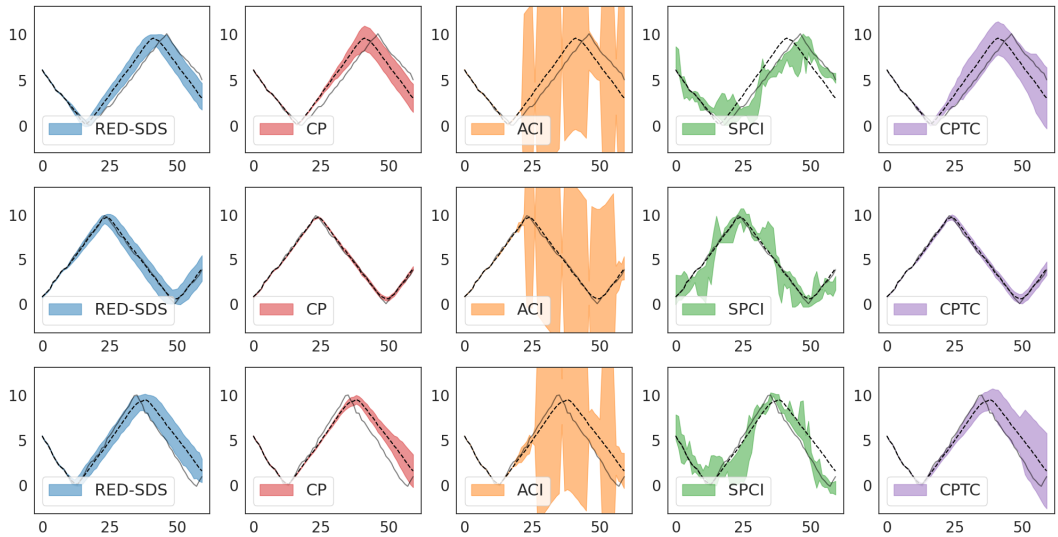


Figure 7: Visualization of prediction intervals on the Bouncing Ball Obs dataset.

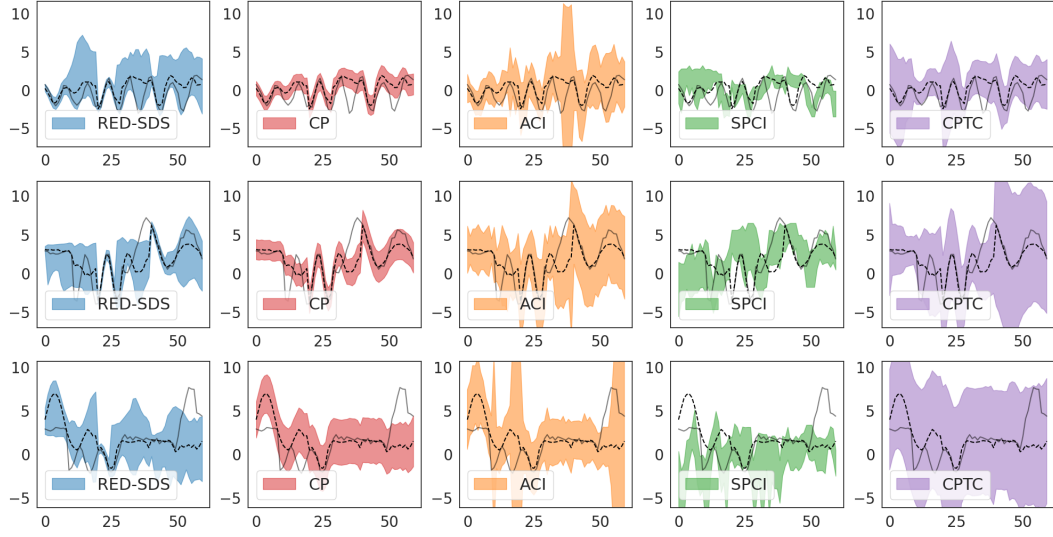


Figure 8: Visualization of prediction intervals on the 3 mode system datasets.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: Main content and contributions are introduced in the abstract and introduction

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: Limitations are discussed in the discussion section.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: I have verified the theoretical claims to be correct as best of my knowledge.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Code to reproduce our results are attached in the supplementary materials. If the paper is published, we will open source the code and code for baseline models.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All data and baselines used are open access. We will open open source our code as well.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Important information discussed in the experiments section, and details are in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: included error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer “Yes” if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Described in detail in appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We followed the code of ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Addressed positive impacts but not negative ones.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: not applicable

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: no proprietary assets

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: N/A

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No crowdsourcing was done.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: no crowdsourcing was done.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: Only used for editing and formatting this paper.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.