

1	<b>Contents</b>	
2	<b>A Theoretical proof</b>	<b>2</b>
3	A.1 Proof of Theorem 1 . . . . .	2
4	A.2 Proof of Theorem 2 . . . . .	5
5	A.3 Tail-up model . . . . .	6
6	<b>B Taxonomy for Related Works</b>	<b>7</b>
7	<b>C Adaptive Uncertainty Set Construction</b>	<b>8</b>
8	<b>D Additional Experiment Details and Results</b>	<b>9</b>
9	D.1 Computational Resource . . . . .	9
10	D.2 Datasets and Codes . . . . .	9
11	D.3 Simulation Data Generation Details . . . . .	10
12	D.4 Real-world Data Details . . . . .	10
13	D.5 Running Time . . . . .	11
14	D.6 Method Details . . . . .	11
15	D.7 Hyperparameter Study . . . . .	11
16	D.8 Additional Ablation Study: Offline Experiment . . . . .	11

## 17 A Theoretical proof

### 18 A.1 Proof of Theorem 1

19 **Theorem 1** (Validity). *Under the assumptions stated above, the proposed method satisfies the*  
 20 *following conditional coverage guarantee:*

$$|\mathbb{P}(Y_{T+1} \in \mathcal{C}_{T+1}(\alpha) | X_{T+1} = x) - (1 - \alpha)| \leq (4L + 2L\sqrt{\omega} + 2)\sqrt{\omega} + 6\sqrt{\frac{\log(16n)}{n}} + \frac{\log(16n)}{n},$$

21 where

$$\omega = \nu_n^2 r + 2r\nu_n \sqrt{(\kappa_1 + \sqrt{\kappa_2})I} + o(g(n))(\kappa_1 + \sqrt{\kappa_2})I.$$

22 For easy notation, denote  $\hat{A} = A_n$ ,  $\Delta_t = \hat{\epsilon}_t - \epsilon_t$  and sometimes we drop subscript  $t$ .

23 **Lemma 1.** *For any test conformity score  $\hat{s}_t = \hat{\epsilon}_t^T \hat{A} \hat{\epsilon}_t$  and the true conformity score  $s_t = \epsilon_t^T A \epsilon_t$ ,*  
 24 *with probability at least  $1 - \delta$ ,*

$$\sum_{t=T-n+1}^T |\hat{s}_t - s_t| \leq \omega n, \quad (1)$$

25 where

$$\omega = \nu_n^2 r + 2r\nu_n \sqrt{(\kappa_1 + \sqrt{\kappa_2})I} + o(g(n))(\kappa_1 + \sqrt{\kappa_2})I.$$

26 *Proof.* We have:

$$\begin{aligned} |\hat{s}_t - s_t| &= |\epsilon_t^T A \epsilon_t - \hat{\epsilon}_t^T \hat{A} \hat{\epsilon}_t| \leq |\epsilon_t^T A \epsilon_t - \epsilon_t^T \hat{A} \epsilon_t| + |\epsilon_t^T \hat{A} \epsilon_t - \hat{\epsilon}_t^T \hat{A} \hat{\epsilon}_t| \\ &\leq |\Delta^T \hat{A} \Delta| + 2|\Delta^T \hat{A} \epsilon_t| + |\epsilon_t^T (A - \hat{A}) \epsilon_t| \\ &\leq \|\hat{A}\| \|\Delta\|^2 + 2\|\hat{A}\| \|\Delta\| \|\epsilon_t\| + \|\epsilon_t\|^2 \|A - \hat{A}\| \\ &\leq r \|\Delta\|^2 + 2r \|\Delta\| \|\epsilon_t\| + o(g(n)) \|\epsilon_t\|^2. \end{aligned} \quad (i)$$

27 The inequality (i) exists because of the Cauchy-schwartz inequality and Assumption 2. Hence, by  
 28 Assumption 1, we have

$$\begin{aligned} \sum_{t=T-n+1}^T |\hat{s}_t - s_t| &\leq r n \nu_n^2 + 2r \sum_{t=T-n+1}^T \|\Delta_t\| \|\epsilon_t\| + o(g(n)) \sum_{t=T-n+1}^T \|\epsilon_t\|^2 \\ &\leq r n \nu_n^2 + 2r \sqrt{\left( \sum_{t=T-n+1}^T \|\Delta_t\|^2 \right) \left( \sum_{t=T-n+1}^T \|\epsilon_t\|^2 \right)} + o(g(n)) \sum_{t=T-n+1}^T \|\epsilon_t\|^2 \\ &\leq r n \nu_n^2 + 2r \sqrt{n \nu_n^2 \left( \sum_{t=T-n+1}^T \|\epsilon_t\|^2 \right)} + o(g(n)) \sum_{t=T-n+1}^T \|\epsilon_t\|^2. \end{aligned} \quad (3)$$

29 From Assumption 3, we have

$$\mathbb{E} \left[ \frac{1}{n} \sum_{t=T-n+1}^T \|\epsilon_t\|^2 \right] = \frac{1}{n} \sum_{t=T-n+1}^T \mathbb{E}[\|\epsilon_t\|^2] \leq \kappa_1 I. \quad (4)$$

30 Using Chebyshev's inequality, we have

$$\mathbb{P} \left( \frac{1}{n} \sum_{t=T-n+1}^T \|\epsilon_t\|^2 - \mathbb{E}[\|\epsilon_t\|^2] \geq \sqrt{\frac{\text{Var}[\|\epsilon_t\|^2]}{n\delta}} \right) \leq \frac{\text{Var}[\|\epsilon_t\|^2]}{n \cdot \frac{\text{Var}[\|\epsilon_t\|^2]}{n\delta}} = \delta, \quad (5)$$

31 which means that with probability higher than  $1 - \delta$ ,

$$\begin{aligned}
\frac{1}{n} \sum_{t=T-n+1}^T \|\epsilon_t\|^2 &\leq \mathbb{E}[\|\epsilon_t\|^2] + \sqrt{\frac{\text{Var}[\|\epsilon_t\|^2]}{n\delta}} \\
&\leq \kappa_1 I + \sqrt{\frac{\kappa_2 I}{n\delta}} \\
&\leq (\kappa_1 + \sqrt{\frac{\kappa_2}{n\delta I}}) I \leq (\kappa_1 + \sqrt{\kappa_2}) I.
\end{aligned} \tag{6}$$

32 The last inequality is because we can set  $\delta$  such that  $\delta n I < 1$ . Plug into eq. (3), we have with  
 33 probability higher than  $1 - \delta$ , we obtain eq. (1) and the lemma follows.

34  $\square$

35 Denote the empirical CDF:  $\hat{F}_{n+1}(x) = \frac{1}{n} \sum_{i=T-n+1}^T 1_{\hat{s}_i \leq x}$ ,  $\tilde{F}_{n+1}(x) = \frac{1}{n} \sum_{i=T-n+1}^T 1_{s_i \leq x}$  and  
 36 true CDF of score function  $F_s(x) = P(s \leq x)$ .

37 **Lemma 2.** *Under Assumption 3, for any  $n$ , there exists an event  $A_n$  which occurs with probability at  
 38 least  $1 - \sqrt{\frac{\log(16n)}{n}}$ , such that, conditioning on  $A_n$ ,*

$$\sup_x \left| \tilde{F}_{n+1}(x) - F(x) \right| \leq \sqrt{\frac{\log(16n)}{n}}.$$

39 *Proof.* The proof follows Lemma 1 in [16] that utilizes Dvoretzky-Kiefer-Wolfowitz inequality in  
 40 [7].  $\square$

41 **Lemma 3.** *Under Assumption 1, Assumption 3, with high probability,*

$$\sup_x \left| \hat{F}_{n+1}(x) - \tilde{F}_{n+1}(x) \right| \leq (2L + 1)\sqrt{\omega} + 2 \sup_x \left| \tilde{F}_{n+1}(x) - F_e(x) \right|.$$

42 *Proof.* The proof is similar to Lemma B.6 in [15], and is written here for completeness.

43 Using Lemma 1 we have that with probability  $1 - \delta$ ,

$$\sum_{t=T-n+1}^T |s_t - \hat{s}_t| \leq n\omega. \tag{7}$$

44 Let  $S = \{t : |s_t - \hat{s}_t| \geq \sqrt{\omega}\}$ . Then,

$$|S|\sqrt{\omega} \leq \sum_{t=T-n+1}^T |s_t - \hat{s}_t| \leq n\omega. \tag{8}$$

45 So  $|S| \leq n\sqrt{\omega}$ . Then,

$$\begin{aligned}
|\hat{F}_{n+1}(x) - \tilde{F}_{n+1}(x)| &\leq \frac{1}{n} \sum_{t=T-n+1}^T |1\{\hat{s}_t \leq x\} - 1\{s_t \leq x\}| \\
&\leq \frac{1}{n} |S| + \sum_{t \notin S} |1\{\hat{s}_t \leq x\} - 1\{s_t \leq x\}| \\
&\leq \frac{1}{n} |S| + \frac{1}{n} \sum_{t=T-n+1}^T 1\{|s_t - x| \leq \sqrt{\omega}\} \\
&\leq \sqrt{\omega} + P(|s_{T+1} - x| \leq \sqrt{\omega}) \\
&\quad + \sup_x \left| \frac{1}{n} \sum_{t=T-n+1}^T 1\{|s_t - x| \leq \sqrt{\omega}\} - P(|s_{T+1} - x| \leq \sqrt{\omega}) \right| \\
&= \sqrt{\omega} + [F_s(x + \sqrt{\omega}) - F_s(x - \sqrt{\omega})] \\
&\quad + \sup_x [\tilde{F}_{n+1}(x + \sqrt{\omega}) - \tilde{F}_{n+1}(x - \sqrt{\omega}) - (F_s(x + \sqrt{\omega}) - F_s(x - \sqrt{\omega}))] \\
&\leq (2L + 1)\sqrt{\omega} + 2 \sup_x |\tilde{F}_{n+1}(x) - F_s(x)|, \tag{i} \\
&\tag{ii} \\
&\tag{9}
\end{aligned}$$

46 where (i) is because  $|1\{a \leq x\} - 1\{b \leq x\}| \leq 1\{|b - x| \leq |a - b|\}$  for  $a, b \in \mathbb{R}$ , and (ii) is due to  
47 the Lipschitz continuity of  $F_s(x)$ .  $\square$

#### 48 **Proof of Theorem 1**

49 *Proof.* Look at the conditional coverage of  $Y_{T+1}$  given  $X_{T+1}$ :

$$\begin{aligned}
&|\mathbb{P}(Y_{T+1} \in \mathcal{C}_{T+1}^\alpha \mid X_{T+1} = x_{T+1}) - (1 - \alpha)| \\
&= |\mathbb{P}(\hat{s}_{T+1} \leq 1 - \alpha \text{ quantile of } \hat{F}_{n+1} \mid X_{T+1} = x) - (1 - \alpha)| \\
&= |\mathbb{P}(\hat{F}_{n+1}(\hat{s}_{T+1}) \leq 1 - \alpha) - \mathbb{P}(F_s(s_{T+1}) \leq 1 - \alpha)| \\
&= |\mathbb{E}[1\{\hat{F}_{n+1}(\hat{s}_{T+1}) \leq 1 - \alpha\} - 1\{F_s(s_{T+1}) \leq 1 - \alpha\}]| \\
&\leq \mathbb{P}(|F_s(s_{T+1}) - (1 - \alpha)| \leq |\hat{F}_{n+1}(\hat{s}_{T+1}) - F_s(s_{T+1})|). \tag{10} \\
&\tag{11}
\end{aligned}$$

50 Based on Lemma 2, we can define the event  $A_n, \mathbb{P}(A_n) \geq 1 - \frac{\log(16n)}{n}$ , conditional on  $A_n$ , we have:

$$\sup_x |\tilde{F}_{n+1}(x) - F_s(x)| \leq \sqrt{\frac{\log(16n)}{n}}, \tag{12}$$

51 Hence, we can write eq. (11) as

$$\begin{aligned}
&\mathbb{P}(|F_s(s_{T+1}) - (1 - \alpha)| \leq |\hat{F}_{n+1}(\hat{s}_{T+1}) - F_s(s_{T+1})|) \\
&\leq \mathbb{P}(|F_s(s_{T+1}) - (1 - \alpha)| \leq |\hat{F}_{n+1}(\hat{s}_{T+1}) - F_s(s_{T+1})| \mid A_n) + \mathbb{P}(A_n^c) \\
&\leq \mathbb{P}(|F_s(s_{T+1}) - (1 - \alpha)| \leq |\hat{F}_{n+1}(\hat{s}_{T+1}) - F_s(\hat{s}_{T+1})| + |F_s(\hat{s}_{T+1}) - F_s(s_{T+1})| \mid A_n) + \frac{\log(16n)}{n}. \\
&\tag{13}
\end{aligned}$$

52 Conditional on  $A_n$ :

$$\begin{aligned}
&|\hat{F}_{n+1}(\hat{s}_{T+1}) - F_s(\hat{s}_{T+1})| + |F_s(\hat{s}_{T+1}) - F_s(s_{T+1})| \\
&\leq \sup_x |\hat{F}_{n+1}(x) - F_s(x)| + L|\hat{s}_{T+1} - s_{T+1}| \\
&\leq (2L + 1)\sqrt{\omega} + 3\sqrt{\frac{\log(16n)}{n}} + L\omega. \tag{14}
\end{aligned}$$

53 The last equation exists because of Lemma 3 1 and eq. (12).

54 Note that  $F_s(s_{T+1}) \sim Unif(0, 1)$ , we have

$$\begin{aligned} & \mathbb{P}(|F_s(s_{T+1}) - (1 - \alpha)| \leq |\hat{F}_{n+1}(\hat{s}_{T+1}) - F_s(\hat{s}_{T+1})| + |F_s(\hat{s}_{T+1}) - F_s(s_{T+1})| | A_n) \\ & \leq (4L + 2)\sqrt{\omega} + 6\sqrt{\frac{\log(16n)}{n}} + 2L\omega. \end{aligned} \quad (15)$$

55 Plug into eq. (11), we have

$$\begin{aligned} & \left| \mathbb{P}\left(Y_{T+1} \in \hat{C}_{T+1}^\alpha \mid X_{T+1} = x_{T+1}\right) - (1 - \alpha) \right| \\ & \leq (4L + 2)\sqrt{\omega} + 6\sqrt{\frac{\log(16n)}{n}} + 2L\omega + \frac{\log(16n)}{n}. \end{aligned} \quad (16)$$

56

□

## 57 A.2 Proof of Theorem 2

58 **Theorem 2** (Efficiency). *There exists  $0.1 < \alpha_0 < 1$ , such that when  $\alpha < \alpha_0$ , the optimal solution to*  
 59 *the minimization problem is given by:*

$$A_* := \arg \min_{A \succ 0} V(A, Q_{1-\alpha}(\epsilon^\top A \epsilon)), \quad (17)$$

60 where  $\epsilon \sim \mathcal{N}(0, A_*^{-1})$ .

61 To prove the theorem, we firstly introduce second-order stochastic dominance to compare the  
 62 magnitude of the  $1 - \alpha$  quantile under different  $A$ .

63 **Definition 1** (Second-Order Stochastic Dominance). *Let  $\rho$  and  $\nu$  be probability measures on  $\mathbb{R}$  such*  
 64 *that  $\mathbb{E}_{X \sim \rho}[|X|] < \infty$  and  $\mathbb{E}_{Y \sim \nu}[|Y|] < \infty$ . The following statements are equivalent; whenever*  
 65 *any of them hold we say that  $\rho$  second-order stochastically dominates  $\nu$  and write  $\rho \succeq_{ssd} \nu$ .*

66 • **Utility criterion.** *For every non-decreasing and concave function  $u : \mathbb{R} \rightarrow \mathbb{R}$ ,*

$$\mathbb{E}_{X \sim \rho}[u(X)] \geq \mathbb{E}_{Y \sim \nu}[u(Y)]. \quad (18)$$

67 • **Integrated distribution-function criterion.** *Denoting by  $F_\rho$  and  $F_\nu$  the cumulative distribu-*  
 68 *tion functions of  $\rho$  and  $\nu$ , respectively,*

$$\int_{-\infty}^t F_\rho(x) dx \leq \int_{-\infty}^t F_\nu(x) dx, \quad \forall t \in \mathbb{R}. \quad (19)$$

69 • **Coupling criterion.** *There exist random variables  $X \sim \rho$  and  $Y \sim \nu$  such that*

$$Y = X - \delta + \varepsilon, \quad \delta \geq 0, \quad \mathbb{E}[\varepsilon \mid X - \delta] = 0. \quad (20)$$

70 The equivalence of the definitions is proved in [8, 11]. Second-order stochastic dominance between  
 71  $X$  and  $Y$  implies the the quantile functions of  $X$  and  $Y$  exhibit certain tail behavior, which would be  
 72 further utilized to compare the  $1 - \alpha$  quantile for a given small  $\alpha$ .

73 **Lemma 4** (Tail-behavior of quantile function). *Let  $X$  and  $Y$  be integrable random variables with*  
 74 *cumulative distribution functions  $F_X, F_Y$  and continuous quantile functions  $Q_X(\alpha) = \inf\{t \in$   
 75  $\mathbb{R}, F_X(t) \geq \alpha\}, Q_Y(\alpha) = \inf\{t \in \mathbb{R}, F_Y(t) \geq \alpha\}$ . If  $X \succeq_{ssd} Y$ , then there exists  $\alpha^* \in$   
 76  $[0, 1]$  such that*

$$Q_X(\alpha) \leq Q_Y(\alpha), \text{ for } \alpha^* < \alpha < 1.$$

77 *Proof.* For integrable  $X$  and  $Y$ , the distribution function criterion for second stochastic dominance  
 78 (19) is equivalent to the quantile-integral condition[13].

$$\Phi(\alpha) := \int_0^\alpha (Q_X(u) - Q_Y(u)) du \geq 0, \quad \forall \alpha \in [0, 1], \quad \Phi(0) = \Phi(1) = 0. \quad (21)$$

79 Set  $\Delta(\alpha) := Q_X(\alpha) - Q_Y(\alpha)$ . Now we prove by contradiction. Assume that there exists a  $b^*$  such  
 80 that  $\int_{b^*}^1 \Delta > 0$ . Then  $\Phi(b^*) = \Phi(1) - \int_{b^*}^1 \Delta < 0$ , contradicting (21). By the continuity of the  
 81 quantile function, there exists an  $\alpha^*$  such that  $\Delta(\alpha) \leq 0$  for  $\alpha > \alpha^*$ . □

82 **Lemma 5.** Let  $Z_1, Z_2, Z_3, \dots, Z_d \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$  and

$$S_{\lambda} = \sum_{i=1}^d \lambda_i Z_i^2, \quad \lambda = (\lambda_1, \lambda_2, \dots, \lambda_d) \in [0, 1]^d, \quad \sum_{i=1}^d \lambda_i = 1.$$

83 Write  $\lambda^* = (\frac{1}{d}, \frac{1}{d}, \dots, \frac{1}{d})$ . Then

$$S_{\lambda} \preceq_{ssd} S_{\lambda^*} \quad \text{for every } \lambda.$$

84 *Proof.* Consider the special case in the coupling criterion (20) when  $\delta = 0$ , which is also known  
 85 as Mean-preserving spread (MPS). Let  $T = \sum_{i=1}^d Z_i^2 \sim \chi_d^2$ . Conditioned on  $T = t$ , the vector  
 86  $(Z_1^2, Z_2^2, \dots, Z_d^2)$  is exchangeable, i.e. we have  $\mathbb{E}[Z_1^2 | T] = \mathbb{E}[Z_2^2 | T] = \dots = \mathbb{E}[Z_d^2 | T]$ . Rewrite  
 87  $S_{\lambda} = S_{\lambda^*} + \varepsilon$ , where  $\varepsilon = \sum_{i=1}^d (\lambda_i - \frac{1}{d}) Z_i^2$ . Then

$$\mathbb{E}(\varepsilon | S_{\lambda^*}) = \sum_{i=1}^d \mathbb{E} \left[ (\lambda_i - \frac{1}{d}) Z_i^2 | T \right] = \sum_{i=1}^d (\lambda_i - \frac{1}{d}) \mathbb{E}[Z_i^2 | T] = 0.$$

88 By the coupling criterion (20),  $S_{\lambda} \preceq_{ssd} S_{\lambda^*}$ . In other words, the uniform weight vector minimises  
 89 risk for every risk-averse preference, completing the proof.  $\square$

## 90 **Proof of Theorem 2**

91 *Proof.* By the definition of ellipsoid volume, we have

$$V(A, Q_{1-\alpha}(\epsilon^{\top} A \epsilon)) = \text{Constant} \times Q_{1-\alpha}(\epsilon^{\top} A \epsilon)^{I/2} [\det(A)]^{-1/2} \quad (22)$$

92 Since the optimization problem is invariant to the rescaling of the scalar (that is,  $A$  and any positive  
 93 scalar multiple  $cA$ , where  $c > 0$ , produce the same mathematical solution), additional constraints,  
 94 such as the bounding of the matrix norm  $A \leq 1$ , can be imposed without loss of generality. Note  
 95 that  $\epsilon \sim N(0, A_*^{-1})$  and  $A \succ 0$  by Cholesky decomposition, we can write  $A_*^{-1} = LL^{\top}$  where  $L$  is a  
 96 lower triangular matrix. Define the matrix  $B = L^{\top} A L$ , we can rewrite the eq. (22) as

$$\min_{B \succ 0} Q_{1-\alpha}^{I/2}(x^{\top} B x) \det(L) [\det(B)]^{-1/2}, \quad (23)$$

97 where  $x \sim N(0, \text{Id})$  and  $\det(L)$  is a constant independent of  $B$ .

98 To further solve the optimization problem, we look at the eigenvalue of  $B$ , suppose  $B =$   
 99  $O \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_I) O^{\top}$  and  $O$  is an orthogonal matrix. Since the optimization value is invariant  
 100 to different scaling of  $B$ , we are imposing additional constraints on  $\lambda_i$ .

$$\min_{\lambda_i > 0, \sum_{i=1}^I \lambda_i = 1} Q_{1-\alpha}^I \left( \sum_{i=1}^I \lambda_i x_i^2 \right) \prod_{i=1}^I \lambda_i^{-\frac{1}{2}}, \quad (24)$$

101 and  $\{x_i\}_{1 \leq i \leq I}$  are i.i.d. random variables  $x_i \sim N(0, 1)$ .

102 Now we would like to prove that the above optimization problem is solved when  $\lambda_1 = \lambda_2 \dots = \lambda_d =$   
 103  $\frac{1}{I}$ . Note that by Cauchy-Schwartz Inequality  $\prod_{i=1}^I \lambda_i^{-\frac{1}{2}}$  is minimized when  $\lambda_1 = \lambda_2 \dots = \lambda_d = \frac{1}{I}$ .  
 104 By Lemma 4 5, there exists an  $\alpha_0$  s.t. when  $\alpha < \alpha_0$ ,  $Q_{1-\alpha}(\sum_{i=1}^I \lambda_i x_i^2)$  is minimized when  
 105  $\lambda_1 = \lambda_2 \dots = \lambda_d = \frac{1}{I}$ . By comparing the quantile function of most majorized(risky) combinations  
 106 of  $\lambda = (1, 0, \dots, 0)$  with the quantile function of the least risky combinations  $\lambda = (\frac{1}{I}, \frac{1}{I}, \dots, \frac{1}{I})$ ,  
 107 empirically we find that  $\alpha_0 > 0.1$ [9].

108  $\square$

## 109 **A.3 Tail-up model**

110 **Lemma 6.** The spatial covariance  $\Sigma$  of any node pair  $(u, v)$  is:

$$\Sigma(u, v) = \int_{\wedge u \cap \wedge v} m(r - u) m(r - v) \frac{w(r)}{\sqrt{w(u)w(v)}} dr. \quad (25)$$

111 *Proof.* Note that

$$\Sigma(u, v) = \text{Cov}\left(\int_{\wedge u} m(s-u) \sqrt{\frac{w(s)}{w(u)}} dB(s), \int_{\wedge v} m(r-v) \sqrt{\frac{w(r)}{w(v)}} dB(r)\right)$$

112 Due to the independence of increments for Brownian motion, only when  $r = s \in \wedge u \cap \wedge v$ , the  
 113 covariance is non-zero. Note that for Brownian motion, we have  $\text{Cov}(dB(r), dB(s)) = 1_{r=s} dr ds$ .  
 114 Hence Lemma 6 follows.  $\square$

115 **Lemma 7.** *If we set the moving function  $m(r-u) = \beta \exp\left(-\frac{d(r,u)}{\phi}\right)$ , with parameters  $\beta > 0$   
 116 (a scale factor) and  $\phi > 0$  (a range or decay parameter), then the covariance matrix between two  
 117 locations  $u, v$  can be expressed as:*

$$\Sigma(u, v) = \sigma^2 \sqrt{\frac{w(u)}{w(v)}} \exp(-d(u, v)/\phi) \mathbf{1}_{\{u \rightarrow v\}}. \quad (26)$$

118 *Proof.* By Lemma 6 and substitute  $m(r-u) = \beta e^{-d(r,u)/\phi}$  and  $m(r-v) = \beta e^{-d(r,v)/\phi}$ , we have

$$\Sigma(u, v) = \int_{\wedge u \cap \wedge v} \beta^2 \exp\left(-\frac{d(r,u)}{\phi}\right) \exp\left(-\frac{d(r,v)}{\phi}\right) \frac{w(r)}{\sqrt{w(u)w(v)}} dr.$$

119 The set  $\wedge u \cap \wedge v$  are the segments of networks that flow into *both*  $u$  and  $v$ . Consider the following  
 120 cases:

- 121 • If  $u$  and  $v$  are not flow-connected, then  $\wedge u \cap \wedge v = \emptyset$  and hence  $\Sigma(u, v) = 0$ .
- 122 • If  $u$  and  $v$  are flow-connected, without loss of generality assume  $v$  is downstream of  $u$ . Then  
 123  $\wedge u \cap \wedge v = \wedge u$ , and for each  $r$  in  $\wedge u$ ,  $d(r, v) = d(r, u) + d(u, v)$ . Hence we have

$$\exp\left(-\frac{d(r, u) + d(r, v)}{\phi}\right) = \exp\left(-\frac{d(u, v)}{\phi}\right) \exp\left(-\frac{2d(r, u)}{\phi}\right).$$

124 leads to a remaining integral over  $r \in \wedge u$ . We can write

$$\Sigma(u, v) = \beta^2 \sqrt{\frac{w(u)}{w(v)}} \exp\left(-\frac{d(u, v)}{\phi}\right) \int_{\wedge u} \exp\left(-\frac{2d(r, u)}{\phi}\right) \frac{w(r)}{w(u)} dr,$$

125 Note that  $\int_{\wedge u} \exp\left(-\frac{2d(r, u)}{\phi}\right) \frac{w(r)}{w(u)} dr = \Sigma(u, u)$  is a constant, since the additivity constrain on  
 126  $w(u)$  assures the constant variance of site  $u$ . Thus, the tail-up exponential model yields a covariance  
 127 of the form

$$\Sigma(u, v) = \begin{cases} \left(\text{constant factors}\right) \exp\left(-\frac{d(u, v)}{\phi}\right), & \text{if } u \text{ and } v \text{ are flow-connected,} \\ 0, & \text{otherwise.} \end{cases}$$

128  $\square$

## 129 B Taxonomy for Related Works

130 Figure 1 provides an overview of the conformal prediction (CP) literature, complementing the  
 131 discussion in the related work section. The Venn diagram categorizes existing CP methods based  
 132 on the type of data they are designed for—time series, multivariate data, or both—and the nature of  
 133 the assumptions they make, particularly regarding exchangeability, temporal stationarity, and spatial  
 134 independence.

135 Traditional CP methods, such as split conformal prediction for i.i.d. regression[14], assume full  
 136 exchangeability, and thus cannot be directly applied to time series or spatially structured data  
 137 without modification. In the time series domain, recent works such as [1, 18, 5, 17] relax the

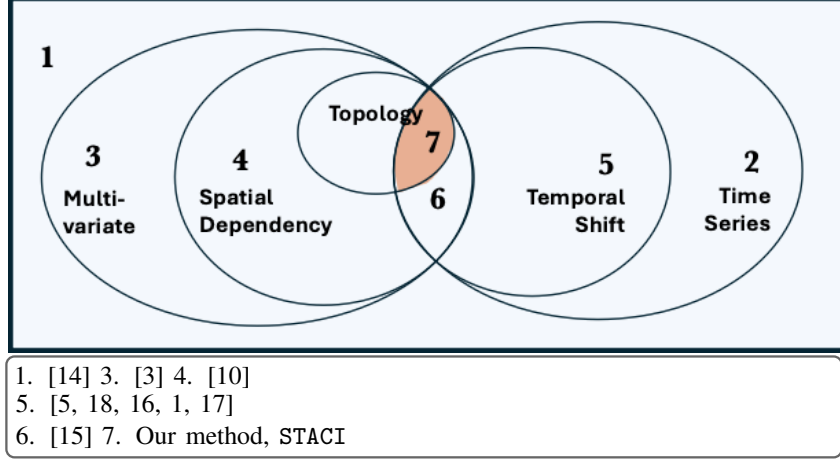


Figure 1: Taxonomy of works in conformal prediction. Among studies that account for both spatial dependency and temporal shift—without assuming spatial and temporal exchangeability—our work is the first to incorporate topology information.

exchangeability assumption via online calibration, sliding window methods, scorecasters or dynamic programming. These methods handle distribution shift over time but typically operate in univariate or low-dimensional settings. In contrast, CP methods for multivariate or high-dimensional data [10, 3] often focus on constructing prediction sets that exploit geometry or sparsity, but generally assume no temporal structure.

Our proposed method, *STACI*, is positioned at the intersection of these axes in the Venn diagram. It is designed for high-dimensional time-indexed multivariate data arising from spatio-temporal stream networks. Our method explicitly accounts for both spatial dependencies and temporal shifts, leveraging the underlying topological structure of the network to enhance predictive performance.

## C Adaptive Uncertainty Set Construction

We construct a spatio-temporally adaptive prediction set for a new observed history  $X_{T+1}$  using our proposed nonconformity score as follows:

$$\mathcal{C}(X_{T+1}; \alpha) = \{y : s(X_{T+1}, y) \leq \hat{Q}_{1-\alpha}\},$$

where  $\hat{Q}_{1-\alpha}$  is the  $(1 - \alpha)$ -th quantile of the empirical cumulative distribution function of  $\{s(X_t, Y_t), t \in \mathcal{D}_{\text{cal}}\}$ . The complete *STACI* algorithm is outlined in Algorithm 1.

To account for potential temporal distribution shifts in the predictive error, we adopt the Adaptive Conformal Inference (ACI) framework proposed in [5]. This approach dynamically updates the confidence level  $\alpha_t$  over time, ensuring that the prediction set remains responsive to evolving data distributions.

Specifically, we iteratively update  $\alpha_t$ , and reconstruct the prediction set  $\mathcal{C}(X_t, \alpha_t)$  accordingly. At the initial test time  $T + 1$ , the confidence level is set as  $\alpha_{T+1} = \alpha$ . For subsequent time steps  $t > T + 1$ , we update  $\alpha_t$  with a step size  $\gamma \geq 0$  as follows:

$$\alpha_{t+1} = \alpha_t + \gamma(\alpha - \mathbb{1}\{Y_t \notin \mathcal{C}(X_t; \alpha_t)\}), \quad \forall t \geq T + 1. \quad (27)$$

The rationale behind ACI is that if the prediction set fails to cover the true value at time  $t$ , the effective error level is reduced, leading to a wider prediction interval at time  $t + 1$ , thereby increasing the likelihood of coverage. A larger step size  $\gamma$  makes the method more responsive to observed distribution shifts but also introduces greater fluctuations in  $\alpha_t$ . When  $\gamma = 0$ , the method reduces to standard (non-adaptive) conformal prediction.

We present the analysis of the average coverage guarantee of *STACI* without any assumption about  $\epsilon_t$ . The proof follows from Proposition 4.1 in [5].



---

**Algorithm 1** STACI

---

**Input:** Data  $\mathcal{D}$ ; Network topology  $\mathcal{G}$ ; Model  $f(\cdot)$ ; Hyper-parameters  $\lambda$ ; Confidence level  $\alpha$ .

**Output:** Prediction set  $\mathcal{C}(X_{T+1}; \alpha)$ .

```
1: // Training
2:  $\hat{f} \leftarrow$  Fit  $f$  using  $\mathcal{D} \setminus \mathcal{D}_{\text{cal}}$ ;
3: // Calibration
4:  $\mathcal{E} \leftarrow \{\hat{\epsilon}_t = Y_t - \hat{f}(X_t) - \frac{\sum_{t \in \mathcal{D}_{\text{cal}}} (Y_t - \hat{f}(X_t))}{|\mathcal{D}_{\text{cal}}|}\}_{t \in \mathcal{D}_{\text{cal}}}$ ;
5:  $\hat{\Sigma}_n \leftarrow \sum_{t \in \mathcal{D}_{\text{cal}}} \hat{\epsilon}_t \hat{\epsilon}_t^\top / (n - 1)$  given  $\mathcal{E}$ ;
6:  $\hat{\Sigma}_{\mathcal{G}} \leftarrow$  Compute (26) for  $(\ell_i, \ell_{i'}), \forall i, i' \in \mathcal{I}$  given  $\mathcal{G}$ ;
7:  $A \leftarrow \lambda \hat{\Sigma}_{\mathcal{G}}^{-1} + (1 - \lambda) \hat{\Sigma}_n^{-1}$ ;
8:  $\mathcal{S} \leftarrow \{\hat{\epsilon}_t^\top A \hat{\epsilon}_t\}_{t \in \mathcal{D}_{\text{cal}}}$  given  $\mathcal{E}$ ;
9:  $\hat{Q}_{1-\alpha} \leftarrow$  Compute  $\frac{[(1-\alpha)(n+1)]}{n}$ -th quantile given  $\mathcal{S}$ ;
10: // Testing
11:  $\mathcal{C}(X_{T+1}; \alpha) \leftarrow \{y : s(X_{T+1}, y) \leq \hat{Q}_{1-\alpha}\}$ ;
```

---

166 **Proposition 1.** Consider  $n'$  test data points as the  $n'$  realizations of  $(X_{T+1}, Y_{T+1})$ , denoted by  
167  $\mathcal{D}_{\text{test}}$ . We have the asymptotic coverage guarantee:

$$\lim_{n' \rightarrow \infty} \sum_{t \in \mathcal{D}_{\text{test}}} \mathbb{1}\{Y_t \notin \mathcal{C}(X_t; \alpha_t)\} / n' = \alpha.$$

168 While Proposition 1 provides a weaker coverage guarantee compared to Theorem 1, it offers broader  
169 applicability, remaining valid even in adversarial online settings. Empirical results suggest that when  
170 the error process exhibits minimal distribution shift and the assumptions of Theorems 2 and 1 are only  
171 slightly violated, STACI maintains the predefined coverage level ( $\gamma = 0.01$ ) while achieving efficient  
172 prediction sets. However, when  $\gamma > 0$ , Proposition 1 does not ensure a finite-sample coverage gap.  
173 Understanding this limitation and developing methods to control the finite-sample coverage gap  
174 presents an interesting direction for future research. While Proposition 1 provides a weaker coverage  
175 guarantee compared to Theorem 1, it offers broader applicability and maintains validity even under  
176 adversarial online setting. Our empirical observations indicate that when the error process exhibits  
177 minimal distribution shift and the assumptions of Theorems 2 and 1 are only slightly violated, STACI  
178 achieves efficient prediction sets while maintaining the predefined coverage level with  $\gamma = 0.01$ .  
179 When  $\gamma > 0$ , Proposition 1 does not guarantee a finite coverage gap. Investigating this limitation  
180 further presents a promising direction for future research.

## 181 D Additional Experiment Details and Results

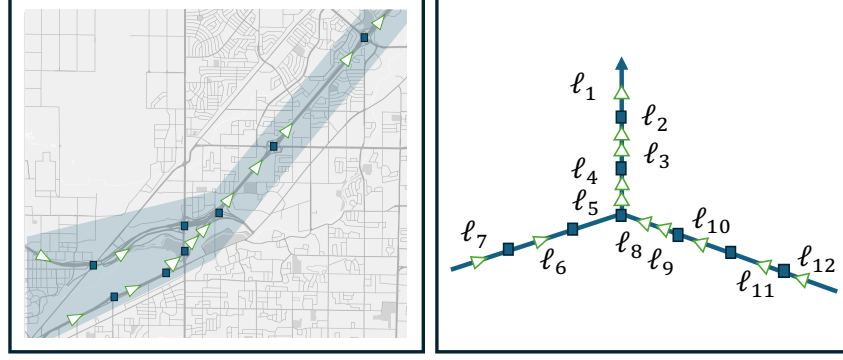
### 182 D.1 Computational Resource

183 Experiments were conducted on a single NVIDIA GeForce RTX 4080 Super GPU, an AMD Ryzen 9  
184 7950X 16-Core Processor CPU, 64GB Memory and 2TB SSD.

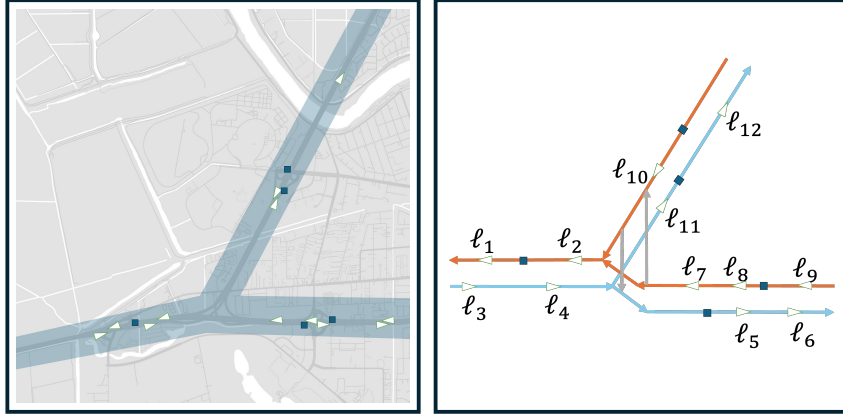
### 185 D.2 Datasets and Codes

186 The PeMS03 dataset used in this paper is collected by California Transportation Agencies (CalTrans)  
187 Performance Measurement System (PeMS). It contains three months of statistics on traffic flow every  
188 5 minutes ranging from Sept. 1st 2018 to Nov. 30th 2018, including 358 sensors. The datasets are  
189 from <https://github.com/guoshnBJTU/ASTGNN/tree/main> without available license.

190 The first backbone ST-Graph model is AGCRN [2], which we adapted the official imple-  
191 mentation from <https://github.com/LeiBAI/AGCRN> under MIT license. ASTGCN [6] was  
192 implemented with PyTorch Geometric Temporal package [12] from [https://github.com/benedekrozemberczki/pytorch\\_geometric\\_temporal](https://github.com/benedekrozemberczki/pytorch_geometric_temporal) under MIT license. STGODE [4] was  
193 adapted from the official implementation <https://github.com/square-coder/STGODE> under  
194 Apache-2.0 license.  
195



(a) PEMS-G1



(b) PEMS-G2

Figure 2: Real-world road network structures and their abstraction for PEMS-G1 and PEMS-G2. In each sub-figure, the left map displays the road network, where freeways are bold gray lines in blue shade, and ramps off the freeway are represented by blue squares. Based on these ramps and road junctions, the network is divided into different segments. Traffic flow monitoring sensors from  $\ell_1$  to  $\ell_{12}$  are placed exclusively on those northbound freeways, marked with green transparent triangles. The right map provides an abstract representation of the road network and sensor locations, using the same symbols for consistency.

### 196 D.3 Simulation Data Generation Details

197 Segment  $r_1$  and  $r_2$  starts with  $(0, 1)$  and  $(0.5, 0.8)$ , respectively, and both end with  $(0.3, 0.5)$ . The  
 198 next segment  $r_3$  also starts with  $(0.3, 0.5)$ , and end with  $(0.2, 0.1)$ . Segment  $r_4$  start from  $(0.6, 0.6)$ ,  
 199 and ends at the same location as  $r_3$ . Starting from this location,  $r_5$  ends at  $(0.4, 0)$ . The weights for  
 200 segment 1 – 5 are set as 0.35, 0.5, 0.85, 0.15 and 1, respectively. Each segment has two observation  
 201 locations – one at the start point, another at the middle point.

202 To approximate the integral, each segment is uniformly divided into 300 smaller sub-intervals. For  
 203 segments without parent nodes ( $r_1$ ,  $r_2$  and  $r_4$  in our example), the source nodes are treated as  
 204 infinitely distant. In implementation, the source node of each segment is extended 10 times in the  
 205 same direction to simulate infinity.

### 206 D.4 Real-world Data Details

207 As shown in Figure 2, we construct two subgraphs from PeMS03 dataset. Note that road network  
 208 structure are required in stream networks. Among all PeMS datasets, PeMS03 is the only one that  
 209 contains a mapping to real sensor identification number in PeMS. Therefore, other datasets are not  
 210 available for our setting.

## D.5 Running Time

On a single Nvidia Geforce 4080S, the training time of AGCRN with PEMS-G1 is 1419s. In table 1, we show the computation time of our and baseline CP methods on AGCRN model and PEMS-G1 dataset. Our method STACI has comparable running time with other methods.

Table 1: Computation Time (seconds) for Different CP Methods with Different Calibration Set Size  $n$

Calibration Set Size $n$	Sphere	Sphere-ACI ( $\gamma = 0.01$ )	Square	<i>MultiDimSPCI</i>	STACI ( $\gamma = 0$ )	STACI ( $\gamma = 0.01$ )
100	142	144	24	24	25	25
200	152	153	24	24	26	25
300	135	129	24	24	25	26
400	132	135	25	24	24	25
500	135	136	23	24	23	24

## D.6 Method Details

We provide pseudo-codes for our proposed method in 1. This applies all experiments in this paper, excluding offline setting detailed in Section D.8.

## D.7 Hyperparameter Study

The studies for hyperparameter are presented in Figure 3.

## D.8 Additional Ablation Study: Offline Experiment

In both synthetic and real-world data, MultiDimSPCI achieve the closest to our proposed method, STACI, in efficiency. Therefore, we focus our comparison on four specific variants: vanilla MultiDimSPCI( $\gamma = 0$ ), MultiDimSPCI( $\gamma = 0.01$ ), STACI( $\gamma = 0$ ), and STACI( $\gamma = 0.01$ ). In the offline setting, STACI does not update the covariance matrix estimation. To ensure a fair comparison, we similarly fix the covariance matrix for MultiDimSPCI methods at the beginning of the test phase.

The results are illustrated in Figure 4. As seen in the left figure, fixing the covariance matrix significantly improves the coverage rates of all methods, bringing them close to the desired 95% level. However, despite having the same  $\gamma$ , STACI consistently outperforms MultiDimSPCI in efficiency. Notably, when ACI is not applied ( $\gamma = 0$ ), both methods tend to be overly conservative, resulting in coverage rates well above the desired 95%. Therefore, since STACI ( $\gamma = 0$ ) achieves a higher coverage rate, MultiDimSPCI ( $\gamma = 0.01$ ) and STACI( $\gamma = 0$ ) exhibit similar efficiency.

In conclusion, regardless of whether the covariance matrix is fixed or not, STACI consistently surpasses MultiDimSPCI in both coverage and efficiency. Furthermore, to achieve an exact coverage rate, incorporating ACI ( $\gamma = 0.01$ ) is recommended.

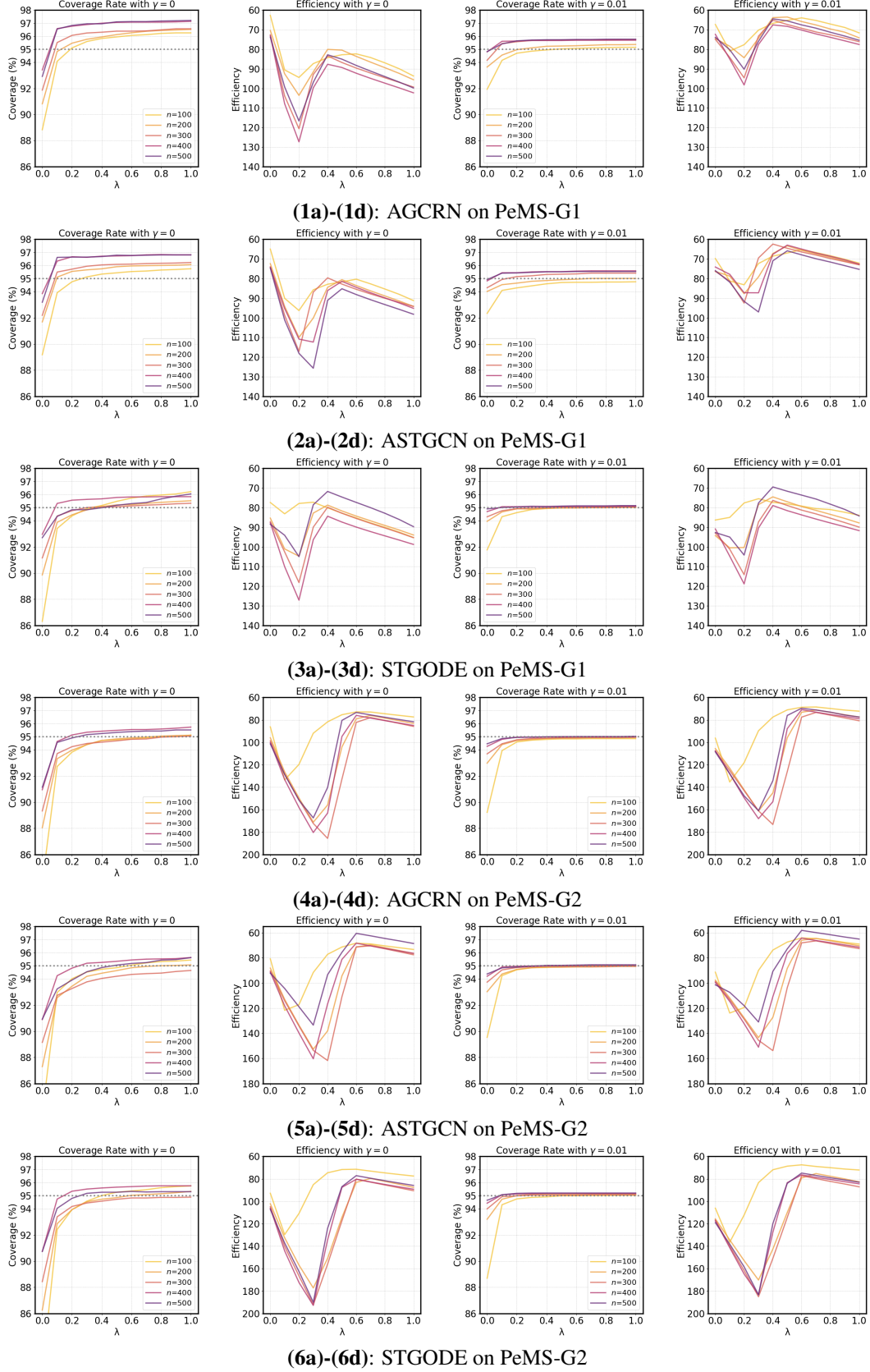


Figure 3: Comparison of Coverage and Efficiency for all PeMS experiments with different backbone GNN models and datasets. Each figure show results for different belief weight  $\lambda$  and calibration set size  $n$ , with adaptive step size  $\gamma = 0$  (a&b) and 0.01 (c&d). Pre-determined coverage threshold of 95% is shown by horizontal gray dotted lines. 12

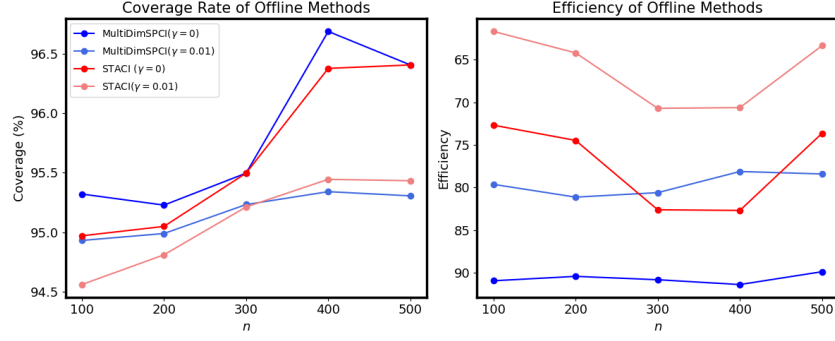


Figure 4: Coverage and efficiency of different methods with different calibration set size  $n$ . Multi-DimSPCI methods are in blue, and our methods are in red. Methods with  $\gamma = 0$  are in darker colors; while those with adaptive coverage,  $\gamma = 0.01$ , are shown in shallow colors.

## References

- [1] Anastasios Angelopoulos, Emmanuel Candes, and Ryan J Tibshirani. Conformal pid control for time series prediction. *Advances in neural information processing systems*, 36:23047–23074, 2023.
- [2] Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. Adaptive graph convolutional recurrent network for traffic forecasting. *Advances in neural information processing systems*, 33:17804–17815, 2020.
- [3] Jacopo Diquigiovanni, Matteo Fontana, and Simone Vantini. Distribution-free prediction bands for multivariate functional time series: an application to the italian gas market. *arXiv preprint arXiv:2107.00527*, 2021.
- [4] Zheng Fang, Qingqing Long, Guojie Song, and Kunqing Xie. Spatial-temporal graph ode networks for traffic flow forecasting. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pages 364–373, 2021.
- [5] Isaac Gibbs and Emmanuel Candes. Adaptive conformal inference under distribution shift. *Advances in Neural Information Processing Systems*, 34:1660–1672, 2021.
- [6] Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):922–929, Jul. 2019.
- [7] Michael R Kosorok. *Introduction to empirical processes and semiparametric inference*, volume 61. Springer, 2008.
- [8] Michael Landsberger and Isaac Meilijson. Mean-preserving portfolio dominance. *The Review of Economic Studies*, 60(2):479–485, 1993.
- [9] Albert W Marshall, Ingram Olkin, and Barry C Arnold. *Inequalities: theory of majorization and its applications*. 1979.
- [10] Soundouss Messoudi, Sébastien Destercke, and Sylvain Rousseau. Copula-based conformal prediction for multi-target regression. *Pattern Recognition*, 120:108101, 2021.
- [11] Michael Rothschild and Joseph E Stiglitz. Increasing risk: I. a definition. In *Uncertainty in economics*, pages 99–121. Elsevier, 1978.
- [12] Benedek Rozemberczki, Paul Scherer, Yixuan He, George Panagopoulos, Alexander Riedel, Maria Astefanoaei, Oliver Kiss, Ferenc Beres, Guzman Lopez, Nicolas Collignon, and Rik Sarkar. PyTorch Geometric Temporal: Spatiotemporal Signal Processing with Neural Machine Learning Models. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management*, page 4564–4573, 2021.

- 269 [13] Moshe Shaked and J George Shanthikumar. *Stochastic orders*. Springer, 2007.
- 270 [14] Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random*  
271 *world*, volume 29. Springer, 2005.
- 272 [15] Chen Xu, Hanyang Jiang, and Yao Xie. Conformal prediction for multi-dimensional time series  
273 by ellipsoidal sets. In *Forty-first International Conference on Machine Learning*, 2024.
- 274 [16] Chen Xu and Yao Xie. Conformal prediction for time series. *IEEE Transactions on Pattern*  
275 *Analysis and Machine Intelligence*, 45(10):11575–11587, 2023.
- 276 [17] Zitong Yang, Emmanuel Candès, and Lihua Lei. Bellman conformal inference: Calibrating  
277 prediction intervals for time series. *arXiv preprint arXiv:2402.05203*, 2024.
- 278 [18] Margaux Zaffran, Olivier Féron, Yannig Goude, Julie Josse, and Aymeric Dieuleveut. Adaptive  
279 conformal predictions for time series. In *International Conference on Machine Learning*, pages  
280 25834–25866. PMLR, 2022.