

---

# Full Conformal Prediction under Stochastic Non-Conformity Measure

---

Thanawat Sornwanee  
Stanford University

## Abstract

The theory of full conformal prediction uses deterministic non-conformity measure, but modern usage of full conformal prediction often relies on machine learning training, making stochasticity inevitable. A simple sufficient condition of almost sure permutation invariance of the non-conformity measure can be too restrictive, so many have suggested the relaxation to permutation in distribution as a condition for full conformal prediction validity. We, however, show that this commonly known condition is actually insufficient. We then provide a correct sufficient condition: *Conditional Independence & Permutation Invariance in Distribution*, which encompasses several stochastic settings that may be used in machine learning.

## 1 Full Conformal Prediction

Given the dataset of

$$\left( [z_i = (x_i, y_i)]_{i=1}^N, x_{N+1} \right),$$

where  $x_i$  is the input and  $y_i$  is the output, we are interested in constructing a confidence bound on the unobserved output  $y_{N+1}$ .

A standard conformal prediction relies on the routine of retrieving conformal p-value from the rank of the non-conformity score. Afterwards, the conformal prediction returns a confidence set

$$\left\{ \hat{y} \in \mathcal{Y} : \text{conformal p-value} \left( [z_i]_{i=1}^N, (x_{N+1}, \hat{y}) \right) \geq \alpha \right\}.$$

In full conformal prediction setting, we are given a non-conformity score function  $t : \mathcal{Z}^{N+1} \rightarrow \mathcal{S}$ , where

$\mathcal{S}$  is a totally preordered set endowed with a binary relation, which is reflexive, transitive, and complete,  $\succsim$ , and is normally chosen to be  $(\mathbb{R}, \geq)$ . (Vovk et al., 2005, 2016) We then compute the score in leave-one-out manner, and then compute the conformal p-value from the rank of the score as outlined in the algorithm 1.

---

**Algorithm 1** Conformal p-value with deterministic non-conformity measure  $t : \mathcal{Z}^{N+1} \rightarrow \mathcal{S}$

---

**Require:** Full data  $z \in \mathcal{Z}^{N+1}$

- 1: **for**  $i = 1$  to  $N + 1$  **do**
  - 2:   Score  $s_i = t(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_{N+1}, z_i)$
  - 3: **end for**
  - 4: Conformal p-value  $p = \frac{1}{N+1} \left[ 1 + \sum_{i=1}^N \mathbf{1}_{s_i \succsim s_{N+1}} \right]$
  - 5: **return**  $p$
- 

The result is that, if  $Z = [Z_i = (X_i, Y_i)]_{i=1}^{N+1}$  is exchangeable, and the non-conformity score function  $t$  is permutation invariant in the first  $N$  arguments, then the conformal p-value of the actual data  $Z$  is a valid p-value.<sup>1</sup> This implies that the actual output  $Y_{N+1}$  is in the confidence set with respect to  $\left( [Z_i]_{i=1}^N, X_{N+1} \right)$  with probability  $\geq 1 - \alpha$ , where the randomness is from the stochasticity of  $Z$  itself.

In most cases, the function  $t$  uses the first  $N$  data entries, which is  $\{z_j\}_{j \neq i}$ , symmetrically to train a regression model  $\hat{\mu}_{-i} : \mathcal{X} \rightarrow \mathcal{Y}$ , and use the loss  $\|y_i - \hat{\mu}_{-i}(x)\|$  to be the non-conformity score  $S_i$ . If the training returns more estimators, such as variance estimator or quantile estimator, one can modify the score to incorporate them. Papadopoulos et al. (2002); Romano et al. (2019). Overall, the non-conformity score will often be in a form of

$$l \left( \mathcal{A} \left( [z_j]_{j \neq i} \right), x_i, y_i \right),$$

where  $l$  is a loss function.

Note that the trained model  $\hat{\mu}_{-i}$  or  $\mathcal{A} \left( [z_j]_{j \neq i} \right)$  does not have to be good in its performance or satisfy any

---

<sup>1</sup>See appendix A for a formal definition of p-value validity.

conditions apart from that its training process must be deterministic and permutation invariant in the data input  $\{z_j\}_{j \neq i}$ . For example, a ridge regression or machine learning with fixed initialization and stochastic descent will satisfy this.

Permutation invariance seems to be natural to most modern machine learning algorithms, but deterministic training requirements are often violated. Modern applications of full conformal prediction often rely on a training of a machine learning model, which may involve stochasticity via stochastic gradient descent as well as random initialization. Lee and Zhang (2025); Taylor et al. (2025) Some may also employ different stochastic training models such as random forests. (Lillicrap et al., 2014; Gauraha and Spjuth, 2018)

Apart from the stochasticity in the training, we will see that the stochasticity in the inference time is also not uncommon. The emergence of modern LLM (large language model) as well as other generative models can have stochasticity in the inference such as from the API call itself.<sup>2</sup> Angelopoulos and Bates (2022) also suggest adding a small white noise to each score can also be used for tie-breaking.

This suggests that a better understanding of full conformal prediction under a stochastic non-conformity measure is crucial.

## 2 Stochastic Non-Conformity Measure

Before discussing the prior works, we will provide some formulation of full conformal prediction under stochastic measure/score. This is when we replace the routine of computing conformal p-value from using a deterministic function  $t : \mathcal{Z}^{N+1} \rightarrow \mathcal{S}$  to stochastic function(s)  $T$ . Note that we will assume that every function mentioned has to be measurable throughout.

Since we may evaluate different non-conformity scores with different randomness, we then need to provide flexibility for the stochastic function of each coordinate to be different. For example,  $\omega$  may dictate 3 independently drawn random seeds  $\epsilon_1$ ,  $\epsilon_2$ , and  $\epsilon_3$ . The training process to get score 1 may only utilize the seed  $\epsilon_1$ , while that of the score 2 only uses  $\epsilon_2$ . Therefore, mathematically, we need to put the indicator for  $T$ . Formally, we have the algorithmic random space  $(\Omega, \mathcal{F}, \mathbb{P}_{\text{alg}})$ , which will throughout be as-

<sup>2</sup>Currently, the usage of LLM as well as other pretrained models in conformal prediction is more limited to a split conformal setting. (Ravfogel et al., 2023; Quach et al., 2024; Su et al., 2024; Epstein et al., 2026) However, one can imagine the use of full conformal prediction by doing a fine-tuning step, or other post-training alignment methods, with the input data. Though this has not been extensively studied

sumed to be independent from that generating data, and a vector of stochastic functions  $T = [T_i]_{i=1}^{N+1}$  where  $T_i : \mathcal{Z}^{N+1} \times \Omega \rightarrow \mathcal{S}$  for each  $i \in \{1, 2, \dots, N+1\}$ .

The routine of computing non-conformity score becomes the algorithm 2. Note that, for a given  $z$ , the p-value  $P$  still inherits randomness from the algorithm.

---

**Algorithm 2** Conformal p-value with stochastic non-conformity measure  $T = [T_i : \mathcal{Z}^{N+1} \times \Omega \rightarrow \mathcal{S}]_{i=1}^{N+1}$

---

**Require:** Full data  $z \in \mathcal{Z}^{N+1}$ , Algorithmic random space  $(\Omega, \mathcal{F}, \mathbb{P}_{\text{alg}})$

- 1: Realize  $\omega \sim \mathbb{P}_{\text{alg}}$
  - 2: **for**  $i = 1$  to  $N + 1$  **do**
  - 3:      $S_i = T_i(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_{N+1}, z_i; \omega)$
  - 4: **end for**
  - 5:  $P = \frac{1}{N+1} \left[ 1 + \sum_{i=1}^N \mathbf{1}_{S_{N+1} \geq S_i} \right]$
  - 6: **return**  $P$
- 

The question of interest is then

*What is a sufficient condition on the random functions  $T = [T_i]_{i=1}^{N+1}$  for the validity of the conformal prediction?*

This is when the actual observation  $Y_{N+1}$  is included in the confidence set with probability  $\geq 1 - \alpha$ , where now the randomness comes from both the data  $Z$  itself and the algorithm. This is equivalent to asking for a sufficient condition for the conformal p-value in the algorithm 2 to be a valid p-value whenever the random input  $Z$  is an exchangeable random vector. Formally, for any exchangeable  $Z$ , we require, for all  $\alpha \geq 0$ ,

$$\mathbb{P}(\text{conformal p-value}(Z; \omega) \leq \alpha) \leq \alpha.$$

### 2.1 Prior Works

As discussed, the stochasticity in full conformal prediction is not uncommon, but the theoretical study is lacking. To the author's knowledge, only 5 papers explicitly mention about condition for stochastic conformal score. However, the mentions are brief: 3 of them are in footnotes. Moreover, 4 claims are also incorrect. For the excerpt, one can look at the appendix B.

**Correct Claim** Lee et al. (2025) seem to claim that a sufficient condition is when

$$T_i(z; \omega) = T_j \left( [z_{\sigma(i)}]_{i=1}^N, z_{N+1}; \omega \right).$$

almost surely.

This is sufficient but quite strong. For example, it cannot handle stochastic gradient descent in a usual

sense. Still, we provide some examples where this is satisfied in the subsection 4.2.

**Incorrect Claim** Other 4 papers that have mentioned the sufficient condition for the validity of full conformal suggests that the almost sure condition can be relaxed into an in distribution condition.

The paper ‘‘Conformal Prediction Beyond Exchangeability’’ (Barber et al., 2023) has claimed that in a case where the stochastic  $T$  is such that

$$\begin{aligned} T_i(z_1, z_2, \dots, z_{i-1}, z_{i+1}, z_{i+2}, \dots, z_{N+1}, z_i; \omega) \\ = (\mathcal{A}(z_1, z_2, \dots, z_{i-1}, z_{i+1}, z_{i+2}, \dots, z_{N+1}; \omega))(z_i) \end{aligned}$$

for some random algorithm  $\mathcal{A}$  mapping  $z_{-i} \in \mathcal{Z}^N$  to a function  $\hat{\mu} : \mathcal{X} \rightarrow \mathcal{Y}$ , which is assumed to be  $\mathbb{R}$ , then we only require

$$\mathcal{A}(z; \cdot) \stackrel{d}{=} \mathcal{A}([z_{\sigma(i)}]_{i=1}^N; \cdot)$$

for all  $z \in \mathcal{Z}^{N+1}$  and for all permutation  $\sigma \in S_N$ . The same permutation symmetry in distribution sense is also claimed to be a sufficient condition in the paper ‘‘Training-conditional coverage for distribution-free predictive inference’’ (Bian and Barber, 2023) and by Lee et al. (2023). This condition implies that, for all  $i, j \in \{1, 2, \dots, T+1\}$ ,  $z \in \mathcal{Z}^{N+1}$ ,  $\sigma \in S_N$ ,

$$T_i(z; \cdot) \stackrel{d}{=} T_j([z_{\sigma(i)}]_{i=1}^N, z_{N+1}; \cdot),$$

and  $T_i = T_j$ . This condition is claimed to be a sufficient condition by Bai and Jin (2024).

If these claims are true, then we will see that most application of machine learning will work since equality in distribution is much easier to attain. However, all these 4 papers are not correct in these sufficient stochasticity claims, at least not without further structure like independence. We will give out a simple counterexample in the section 3.

### 3 Why Permutation Equivariance in Distribution is Insufficient?

Below we will show a simple example where the equality in distribution under permutation is satisfied for the training algorithm  $\mathcal{A}$ , so the equality in distribution under permutation of the first  $N$  element is satisfied for the stochastic scoring functions  $T$ , but the the full conformal prediction is invalid.

Consider the case where  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} \text{Unif}([0, 1]^2)$ , and  $N = 2$ . Thus, the random vector  $Z = (Z_1, Z_2, Z_3) \in ([0, 1]^2)^3$  is exchangeable.

Let the random space be such that we can define two independent standard Brownian motions

$\{B_t(\omega)\}_{t \in [0, 1]}$  and  $\{B'_t(\omega)\}_{t \in [0, 1]}$ , and we define a random scoring function for all  $i \in \{1, 2, 3\}$  to be

$$T_i(z; \omega) := (\mathcal{A}(z_1, z_2; \omega))(z_3) := B_{y_2}(\omega) + B'_{1-y_2}(\omega),$$

for any  $z \in ([0, 1]^2)^3$ , thereby making

$$T(z; \cdot) \sim N(0, 1),$$

and similarly  $\mathcal{A}$  is permutation invariant in distribution, since it always return a constant function whose value is distributed according to  $N(0, 1)$ .

The non-conformity scores from the algorithm 2 conditioned on the data  $Z$  are

$$\begin{cases} (S_1|Z)(\cdot) &= T(Z_2, Z_3, Z_1; \cdot) = B_{Y_3} + B'_{1-Y_3} \\ (S_2|Z)(\cdot) &= T(Z_1, Z_3, Z_2; \cdot) = B_{Y_3} + B'_{1-Y_3}, \\ (S_3|Z)(\cdot) &= T(Z_1, Z_2, Z_3; \cdot) = B_{Y_2} + B'_{1-Y_2} \end{cases}$$

making the conformal p-value

$$\begin{aligned} (P|Z)(\cdot) &= \left( \frac{1 + \mathbf{1}_{S_1 \geq S_3} + \mathbf{1}_{S_2 \geq S_3}}{3} \middle| Z \right) \\ &= \frac{1}{3} + \frac{2}{3} \left( \mathbf{1}_{B_{Y_2} + B'_{1-Y_2} \geq B_{Y_3} + B'_{1-Y_3}} \middle| Z \right). \end{aligned}$$

Under the case when  $Y_2 \neq Y_3$ , which happens almost surely, we have that  $B_{Y_2} + B'_{1-Y_2} \geq B_{Y_3} + B'_{1-Y_3}$  with probability  $\frac{1}{2}$  conditioned on  $Y_2$  and  $Y_3$  (so the randomness is from the algorithmic randomness). Therefore,  $(P|Z)(\cdot) \sim \frac{1}{2}\delta_{\frac{1}{3}} + \frac{1}{2}\delta_1$  almost surely in  $Z$  (not in the algorithmic random space). Thus, almost surely in  $Z$ ,

$$\mathbb{P}\left(P \leq \frac{1}{3}\right) = \mathbb{E}\left[\mathbb{P}\left(P \leq \frac{1}{3} \middle| Z\right)\right] = \frac{1}{2} > \frac{1}{3},$$

contradicting with that  $\mathbb{P}\left(\hat{P} \geq \frac{1}{3}\right) \geq \frac{1}{3}$ , so  $(\hat{P}|Z)$  cannot be a valid p-value.

Thus, the sufficient condition of permutation invariance in distribution known to the literature is actually *incorrect*.<sup>3</sup>

### 4 Example Sufficient Conditions

Recall that the validity of conformal prediction comes from the validity of permutation test. Thus, we only need that, as long as the data points  $Z$  are exchangeable, the non-conformity scores  $S := [S_i]_{i=1}^{N+1}$  are also exchangeable (unconditionally).

<sup>3</sup>It is easier to construct a counter example when we allow the random function for each score computation to use different part of the randomness. For example, we can consider the case when  $\omega \sim \text{Unif}([0, 1]^2)$ , and  $S_1 = S_2 = \omega_1$ , while  $S_3 = \omega_2$ .

#### 4.1 Deterministic with Permutation Invariance

For deterministic case, a simple sufficient condition is permutation invariant among first  $N$  arguments, meaning that

$$T_i(z; \omega) = T_j([z_{\sigma(i)}]_{i=1}^N, z_{N+1}; \omega') \quad (1)$$

for all  $i, j \in \{1, 2, \dots, N+1\}$ ,  $z \in \mathcal{Z}^{N+1}$ ,  $\sigma \in S_N$ , and for almost surely  $\omega, \omega'$ . This is a standard condition, and can be achieved by symmetrization over a deterministic function  $f: \mathcal{Z}^{N+1} \rightarrow \mathbb{R}$  such as choosing  $T_i(z; \omega) = \frac{1}{N!} \sum_{\sigma \in S_N} f(z_{\sigma(1)}, z_{\sigma(2)}, \dots, z_{\sigma(N)}, z_{N+1})$ , but this requires  $N!$  computation.

#### 4.2 Stochastic with Almost Sure Permutation Invariance

As mentioned, this is when we weaken the degeneracy requirement (by dropping  $\omega'$ ) from the condition 1 into

$$T_i(z; \omega) = T_j([z_{\sigma(i)}]_{i=1}^N, z_{N+1}; \omega) \quad (2)$$

almost surely. This is still hard to achieve, as it is translated into that the permutation invariance holds for each random seed  $\omega$ , and we have to reuse the same seed for each score evaluation process. For example, this can be satisfied when the randomness only comes in through random initialization of neural network parameters and full gradient descent is implemented. Another example is when symmetrization is used with the same fixed seed.

For most machine learning training, this rarely holds. Stochastic gradient descent with mini-batch may sample a random index  $I(\omega)$  of the data to be used in the gradient computation. By applying the permutation before running the algorithm, the  $I(\omega)^{\text{th}}$  datapoint in the permuted dataset will correspond to the  $\sigma(I(\omega))^{\text{th}}$  datapoint of the original dataset. Thus, unless we enforce the algorithm to re-sort the datapoints first before using any stochasticity, the data order will change even when the random seed remains fixed.

Moreover, this condition is difficult to satisfy if we have randomness in the inference, since we may not be able to apply a fixed random seed for some API.

#### 4.3 Stochastic with Independence & Permutation Invariance in Distribution

We weaken the almost sure equality 2 to equality in distribution

$$T_i(z; \cdot) \stackrel{d}{=} T_j\left([z_{\sigma(k)}]_{k=1}^N, z_{N+1}; \cdot\right), \quad (3)$$

but with additional independence condition that

$$\{T_i(z_1, z_2, \dots, z_{i-1}, z_{i+1}, z_{i+2}, \dots, z_{N+1}, z_i; \cdot)\}_{i=1}^{N+1}$$

is a collection of independent variables.

This is perhaps the most commonly used method in full conformal prediction. Each time we train a machine learning model, the stochasticity is assumed to be independent from those in the past. If there exists randomness in the inference, as long as such randomness are independent from each other, then the validity of full conformal prediction will still hold. However, we can see that this can be violated in some case, such as when each training share the same random seed, which may be used for reproducibility.

### 5 Conditional Independence & Permutation Invariance in Distribution

Note that the deterministic case 4.1 is a special case for both stochastic cases 4.2 and 4.3, but the two stochastic cases are not a special case of each other. We generalize the combination of the two stochastic cases to show that a more general condition of conditional independence and permutation invariance in distribution is sufficient for full conformal prediction validity.

**Definition 1** (Conditional Independence & Permutation Invariance in Distribution). *The stochastic non-conformity score functions  $T = [T_i]_{i=1}^{N+1}$ , where  $T_i: \mathcal{Z}^{N+1} \times \Omega \rightarrow \mathcal{S}$  for each  $i \in \{1, 2, \dots, N+1\}$  are conditionally independent and permutation invariant in distribution if there exists some random variable  $C$  defined on the algorithmic random space such that, for all  $z \in \mathcal{Z}^{N+1}$  and for almost surely  $C$ ,*

$$(T_i(z; \cdot) | C) \stackrel{d}{=} \left( T_j\left([z_{\sigma(k)}]_{k=1}^N, z_{N+1}; \cdot\right) | C \right) \quad (4)$$

for all  $\sigma \in S_N$ , for all  $i, j \in \{1, 2, \dots, N+1\}$ , and the collection of random variables

$$\{T_i(z_1, z_2, \dots, z_{i-1}, z_{i+1}, z_{i+2}, \dots, z_{N+1}, z_i; \cdot)\}_{i=1}^{N+1}$$

is independent conditioned on  $C$ .

By taking  $C$  to be non-informative, this reverts back to the independence and permutation invariance in distribution setting 4.3, and, by taking  $C$  to be  $\omega$ , this reverts back to the almost sure permutation invariance 4.2.

**Theorem 1.** *The conformal  $p$ -value whose stochastic non-conformity score functions  $T$  are conditionally independent and permutation invariant in distribution is valid.*

The proof, which is in the appendix D, first, shows that the joint score computation

$$A(z; \omega) := [T_i(z_1, z_2, \dots, z_{i-1}, z_{i+1}, z_{i+2}, \dots, z_{N+1}, z_i; \omega)]_{i=1}^{N+1}$$

is permutation equivariance in joint distribution. This means that

$$\left[ (A(z; \cdot))_{\pi(i)} \right]_{i=1}^{N+1} \stackrel{d}{=} A([z_{\pi(i)}]_{i=1}^{N+1}; \cdot)$$

for all  $z \in \mathcal{Z}^{N+1}$  and  $\pi \in S_{N+1}$ .

Afterwards, we use exchangeability of  $Z$  to conclude that the scores

$$S = A(Z; \omega)$$

are exchangeable.

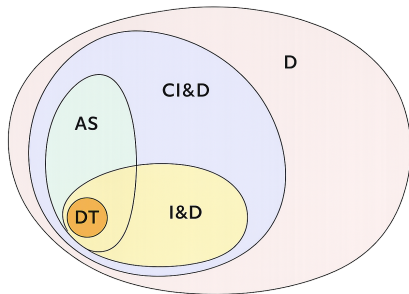


Figure 1: *DT* stands for deterministic permutation invariance, which is a standard assumption in full conformal prediction as outlined in the subsection 4.1. (Vovk et al., 2005) *AS* stands for almost sure permutation invariance as outlined in the subsection 4.2 *I&D* stands for independence and permutation invariance in distribution as outlined in the subsection 4.3. *CI&D* stands for conditional independence and permutation invariance in distribution as defined in the definition 1. We have established through the main theorem 1 that *CI&D* is sufficient to ensure the validity of full conformal prediction when the non-conformity measure is stochastic. *D* stands for permutation invariance in distribution, which is the incorrect sufficient condition in the literature (Barber et al., 2023; Bian and Barber, 2023; Lee et al., 2025; Bai and Jin, 2024). In this figure, we show that  $DT \subseteq AS \cap I\&D$ , while  $AS \cup I\&D \subseteq CI\&D \subseteq D$ . Note that these conditions are defined with respect to an algorithmic random space  $(\Omega, \mathcal{F}, \mathbb{P}_{\text{alg}})$  as well as  $\mathcal{S}$  and  $\mathcal{Z}$ . It is easy to show that, for each  $\subseteq$  relation, we can find a setting such that every  $\subseteq$  relation is indeed  $\subsetneq$ . See appendix E.

## References

- Angelopoulos, A. N. and Bates, S. (2022). A gentle introduction to conformal prediction and distribution-free uncertainty quantification.
- Bai, T. and Jin, Y. (2024). Optimized conformal selection: Powerful selective inference after conformity score optimization. *arXiv preprint arXiv:2411.17983*.
- Barber, R. F., Candes, E. J., Ramdas, A., and Tibshirani, R. J. (2023). Conformal prediction beyond exchangeability. *The Annals of Statistics*, 51(2):816–845.
- Bian, M. and Barber, R. F. (2023). Training-conditional coverage for distribution-free predictive inference. *Electronic Journal of Statistics*, 17(2):2044–2066.
- Epstein, E. L., Winnicki, J., Sornwanee, T., and Dwaraknath, R. V. (2026). LLMs are overconfident: Evaluating confidence interval calibration with fermieval. In *AAAI 2026 Workshop on Assessing and Improving Reliability of Foundation Models in the Real World*.
- Gauraha, N. and Spjuth, O. (2018). conformalclassification: A conformal prediction r package for classification.
- Lee, J., Popov, I., and Ren, Z. (2025). Full-conformal novelty detection: A powerful and non-random approach. *arXiv preprint arXiv:2501.02703*.
- Lee, K. and Zhang, Y. (2025). Leave-one-out stable conformal prediction. In *The Thirteenth International Conference on Learning Representations*.
- Lee, Y., Barber, R. F., and Willett, R. (2023). Distribution-free inference with hierarchical data. *arXiv preprint arXiv:2306.06342*.
- Lehmann, E. L. and Romano, J. P. (2005). *Testing statistical hypotheses*. Springer.
- Linusson, H., Johansson, U., Boström, H., and Löfström, T. (2014). Efficiency comparison of unstable transductive and inductive conformal classifiers. In *IFIP International conference on artificial intelligence applications and innovations*, pages 261–270. Springer.
- Papadopoulos, H., Proedrou, K., Vovk, V., and Gamerman, A. (2002). Inductive confidence machines for regression. In *European conference on machine learning*, pages 345–356. Springer.
- Quach, V., Fisch, A., Schuster, T., Yala, A., Sohn, J. H., Jaakkola, T. S., and Barzilay, R. (2024). Conformal language modeling. In *The Twelfth International Conference on Learning Representations*.

- Ravfogel, S., Goldberg, Y., and Goldberger, J. (2023). Conformal nucleus sampling. In *Findings of the association for computational linguistics: ACL 2023*, pages 27–34.
- Romano, Y., Patterson, E., and Candes, E. (2019). Conformalized quantile regression. *Advances in neural information processing systems*, 32.
- Su, J., Luo, J., Wang, H., and Cheng, L. (2024). Api is enough: Conformal prediction for large language models without logit-access. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 979–995.
- Taylor, D., Correia, A., Nalisnick, E., and Louizos, C. (2025). Approximating full conformal prediction for neural network regression with gauss-newton influence. In *The Thirteenth International Conference on Learning Representations*.
- Vovk, V., Fedorova, V., Nouretdinov, I., and Gammerman, A. (2016). Criteria of efficiency for conformal prediction. In *Symposium on conformal and probabilistic prediction with applications*, pages 23–39. Springer.
- Vovk, V., Gammerman, A., and Shafer, G. (2005). *Algorithmic learning in a random world*. Springer.

## A p-Value via Order Test

We define a valid p-value, and provide background on permutation test, which is the backbone conformal prediction.

**Definition 2** (Valid p-Value). *A random variable  $P$  is a valid p-value if*

$$\mathbb{P}(P \leq \alpha) \leq \alpha$$

for all  $\alpha \geq 0$ .

**Lemma 1.** *Let  $S := [S_i]_{i=1}^n$  be an exchangeable random variable where  $S_i \in \mathcal{S}$  for each  $i \in \{1, 2, \dots, n\}$  and  $(\mathcal{S}, \succ)$  is a total preorder. The distribution of*

$$P_1 := \frac{\sum_{i=1}^n \mathbf{1}_{S_i \succ S_i}}{n}$$

*first order stochastically dominates  $\text{Unif}(\{\frac{1}{n}, \frac{2}{n}, \dots, \frac{n}{n}\}) \succ_{\text{FOSD}} \text{Unif}([0, 1])$ .*

*Proof.* Define  $P_i := \frac{\sum_{j=1}^n \mathbf{1}_{S_i \succ S_j}}{n}$ . Exchangeability of  $S$  implies exchangeability of  $P$ . Note that, for any integer  $k \in \{1, 2, \dots, n\}$ , there can be at most  $k$  values of  $i \in \{1, 2, \dots, n\}$  where  $P_i \leq \frac{k}{n}$ , so

$$\mathbb{P}(P_1) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{P_i \leq \frac{k}{n}} \leq \frac{k}{n},$$

implying the first order stochastic dominations. □

**Corollary 1.**  *$P_1$  in the lemma 1 is a valid p-value if  $S$  is exchangeable.*

*Proof.* From the lemma 1, we have that  $\mathbb{P}(P_1 \leq \alpha) \leq (\text{Unif}([0, 1]))([0, \alpha]) = \alpha$ , so  $P_1$  is a valid p-value. □

In permutation testing (Lehmann and Romano, 2005), we test a hypothesis of whether random variables  $Z := [Z_i]_{i=1}^n$  are exchangeable, by randomly permuting the variables and compute test statistic each time. Under such hypothesis, the test statistics should be exchangeable, so  $P_1 \leq \alpha$  with probability  $\leq \alpha$ , meaning that it is a valid p-value.

## B Excerpts from the Related Works

We provide an excerpt as well as some explanation of the context for each related work.

The proof of the theorem 2 of (Lee et al., 2025) states that “The scoring function  $V^{(k)}(\cdot)$  was trained invariant to the order of elements inside  $E_j$  (this assumes that  $V^{(k)}$  does not use external randomness during training; similar results can be attained for random training procedures by conditioning on the random seed),”  $V^{(k)}$  is equivalent to  $T_i$  in our notation. This suggests that the almost sure permutation invariance condition is sufficient, which is true.

Next, we look at the claim that is more in line of in distribution permutation invariance.

The footnote 2 of (Barber et al., 2023) states that “If  $\mathcal{A}$  is a randomized algorithm, then this equality is only required to hold in a distributional sense,” where the original equality refers to  $\mathcal{A}(z) = \mathcal{A}\left([z_{\sigma(i)}]_{i=1}^N\right)$ , which is a commonly used sufficient condition in deterministic full conformal prediction.

The subsection 2.1.3 “A note on randomized algorithm” of (Bian and Barber, 2023) states that “The background given above implicitly the algorithm  $\mathcal{A}$  as a deterministic function of the training data—that is, we view  $\mathcal{A}$  as a function  $((X_1, Y_1), \dots, (X_n, Y_n)) \mapsto \hat{\mu}$ . In many settings, however, it is common to use a randomized regression algorithm—for instance, stochastic gradient descent. In this setting, we can formally view  $\mathcal{A}$  as a function  $((X_1, Y_1), \dots, (X_n, Y_n), \xi) \mapsto \hat{\mu}$ , where the term  $\xi$  introduces stochastic noise (effectively, a random seed). All the results described above hold for both the deterministic and randomized settings. (For results that assume  $\mathcal{A}$  is symmetric, the symmetry condition (6) should be understood in the distributional sense—that is, the

training data points are treated symmetrically with respect to the randomized training procedure. For example, for stochastic gradient descent, if data points are drawn uniformly at random during the training epochs, then symmetry is satisfied.)”

The footnote 12 of Lee et al. (2023) mentions that “The framework also allows for a randomized algorithm  $\mathcal{A}$ , in which case the symmetry condition is required to hold in a distribution sense.”

The footnote 1 of (Bai and Jin, 2024) mentions that “In the general case where  $\mathcal{V}$  is a randomized algorithm, we require ... Definition 4 to hold in a distributional sense, i.e.,  $=$  could be replaced by  $\stackrel{d}{=}$ ,” where the equality refers to the one used in (Barber et al., 2023).

## C Permutation Equivariance in Distribution Algorithm

To prove the main theorem 1, we recall that  $Z$  is an exchangeable random variable whose stochasticity is independent from that of the stochastic algorithm and  $(\mathcal{S}, \succ)$  is already assumed to be a total preorder. Thus, we only need to ensure that a stochastic algorithm  $A : \mathcal{Z}^{N+1} \times \Omega \rightarrow \mathcal{S}^{N+1}$  that jointly return the score

$$S := [S_i]_{i=1}^{N+1} := A(Z; \omega)$$

will be able to transfer the exchangeability in  $Z$  into the exchangeability in  $S$ , both of which in a distribution sense. We can then construct a p-value

$$P := \frac{1}{N+1} \left[ 1 + \sum_{i=1}^N \mathbf{1}_{S_{N+1} \succ S_i} \right],$$

and get its validity by the corollary 1.

In this section, we first show that permutation equivariant in joint distribution is a sufficient condition (in the definition 3) for the stochastic algorithm  $A$  to preserve exchangeability. Next, we will show a condition on the non-conformity score functions  $T$  to achieve permutation equivariant in joint distribution.

### C.1 Definition & Property

A convenient sufficient condition for exchangeability to be inherited through  $A$  is that the stochastic algorithm  $A$  is permutation equivariant in joint distribution.

**Definition 3** (Permutation Equivariance in Joint Distribution). *A random function  $A : \mathcal{Z}^{N+1} \times \Omega \rightarrow \mathcal{S}^{N+1}$  is permutation equivariant in joint distribution if, for any permutation  $\pi \in S_{N+1}$ , for any  $z \in \mathcal{Z}^{N+1}$ ,*

$$A \left( [z_{\pi(i)}]_{i=1}^{N+1}; \cdot \right) \stackrel{d}{=} [A(z; \cdot)]_{\pi(i)} \Big|_{i=1}^{N+1}.$$

Recall that exchangeable is the same as permutation invariance in distribution. It is then unsurprising that  $S = A(Z; \omega)$  can inherit the exchangeability in distribution from  $Z$ . Recall that the distribution of  $Z$  is fixed to be an exchangeable one, so we leave it out of the lemma statement.

**Lemma 2** (Exchangeability & Permutation Equivariance imply Exchangeability). *If  $A : \mathcal{Z}^{N+1} \times \Omega \rightarrow \mathcal{S}^{N+1}$  is permutation equivariant in joint distribution, then  $A(Z; \omega)$  is exchangeable.*

*Proof.* Consider a permutation  $\pi \in S_{N+1}$ . From exchangeability of  $Z$ , and the permutation equivariance of  $A$ , we have that

$$A(Z) \stackrel{d}{=} A \left( [Z_{\pi(i)}]_{i=1}^{N+1} \right) \stackrel{d}{=} [A(Z)]_{\pi(i)} \Big|_{i=1}^{N+1},$$

meaning that it is exchangeable. □

Therefore, we get the following corollary.

**Corollary 2.** *If  $A : \mathcal{Z}^{N+1} \times \Omega \rightarrow \mathcal{S}^{N+1}$  is permutation equivariant in joint distribution, and  $S = A(Z; \omega)$ , then, by defining*

$$P := \frac{1}{N+1} \left[ 1 + \sum_{i=1}^N \mathbf{1}_{S_{N+1} \succ S_i} \right],$$

*we will have that  $P$  is a valid p-value.*

*Proof.* From the lemma 2, we have that  $S$  is exchangeable. Note that the construction of  $P$  is the same (but use the last index instead of first index) as  $P_1$  in the lemma 1, so by the corollary 1,  $P$  is a valid p-value.  $\square$

## C.2 Construction & Sufficient Condition

In the last subsection, we have established that permutation equivariance will be sufficient for the validity of the rank p-value. Note that this is more general than a conformal p-value, since  $A$  does not have to create non-conformity scores like in the algorithm 2. In this subsection, we will restrict the attention when the algorithm  $A$  generates non-conformity scores under full conformal prediction framework. This is when

$$A(z; \omega) := [T_i(z_1, z_2, \dots, z_{i-1}, z_{i+1}, z_{i+2}, \dots, z_{N+1}, z_i; \omega)]_{i=1}^{N+1}$$

for any  $z \in \mathcal{Z}^{N+1}$ .

Consider

$$A\left([z_{\pi(i)}]_{i=1}^{N+1}; \omega\right) = \begin{bmatrix} T_1(z_{\pi(2)}, z_{\pi(3)}, \dots, z_{\pi(N+1)}, z_{\pi(1)}; \omega) \\ T_2(z_{\pi(1)}, z_{\pi(3)}, z_{\pi(4)}, \dots, z_{\pi(N+1)}, z_{\pi(2)}; \omega) \\ \vdots \\ T_i(z_{\pi(1)}, z_{\pi(2)}, \dots, z_{\pi(i-1)}, z_{\pi(i+1)}, z_{\pi(i+2)}, \dots, z_{\pi(N+1)}, z_{\pi(i)}; \omega) \\ \vdots \\ T_{N+1}(z_{\pi(1)}, z_{\pi(2)}, \dots, z_{\pi(N+1)}; \omega) \end{bmatrix},$$

while

$$[(A(z; \omega))_{\pi(i)}]_{i=1}^{N+1} = \begin{bmatrix} T_{\pi(1)}(z_1, z_2, \dots, z_{\pi(1)-1}, z_{\pi(1)+1}, z_{\pi(1)+2}, \dots, z_{N+1}, z_{\pi(1)}; \omega) \\ T_{\pi(2)}(z_1, z_2, \dots, z_{\pi(2)-1}, z_{\pi(2)+1}, z_{\pi(2)+2}, \dots, z_{N+1}, z_{\pi(2)}; \omega) \\ \vdots \\ T_{\pi(i)}(z_1, z_2, \dots, z_{\pi(i)-1}, z_{\pi(i)+1}, z_{\pi(i)+2}, \dots, z_{N+1}, z_{\pi(i)}; \omega) \\ \vdots \\ T_{\pi(N+1)}(z_1, z_2, \dots, z_{\pi(N+1)-1}, z_{\pi(N+1)+1}, z_{\pi(N+1)+2}, \dots, z_{N+1}, z_{\pi(N+1)}; \omega) \end{bmatrix}.$$

Thus, for each coordinate, the two random vectors share the same last element, but with different ordering and index of the score function. If we want permutation equivariance in joint distribution, we then need the marginal distribution of each coordinate to be the same. This can be easily satisfied by enforcing permutation invariance in distribution 4. However, as we have shown in the section 3, this is not sufficient, since this will only give permutation equivariance in marginal distributions but not joint distribution. Thus, a stronger condition is required.

For example, if we have almost sure permutation invariance 2, we can rearrange the first  $N$  terms and change the index from  $T_i$  to  $T_{\pi(i)}$  for each coordinate. This gives  $A\left([z_{\pi(i)}]_{i=1}^{N+1}; \omega\right) = [(A(z; \omega))_{\pi(i)}]_{i=1}^{N+1}$  almost surely, which is stronger than what required in permutation equivariance in joint distribution.

Next, we will show that conditional independence and permutation invariance in distribution condition of the stochastic scoring functions  $T$  ensures that the stochastic algorithm  $A$  is permutation equivariant in joint distribution.

**Lemma 3** (Conditional Independence & Permutation Invariance in Distribution implies Permutation Equivariance). *If the stochastic non-conformity score functions  $T$  are conditionally independent and symmetric, then the stochastic algorithm  $A(z; \omega) := [T_i(z_1, z_2, \dots, z_{i-1}, z_{i+1}, z_{i+2}, \dots, z_{N+1}, z_i; \omega)]_{i=1}^{N+1}$  almost surely for all  $z \in \mathcal{Z}^{N+1}$  is permutation equivariant in joint distribution.*

*Proof.* From Dynkin's  $\pi - \lambda$  theorem, it is sufficient to show the equality in probability over rectangular set. Consider an arbitrary  $B := \bigotimes_{i=1}^{N+1} B_i$ , where each  $B_i$  is a measurable subset of  $\mathcal{S}$ . From the definition 1 and the definition of  $A$ , we have that the collection

$$\{(A(z; \cdot))_i\}_{i=1}^{N+1}$$

is independent conditioned on  $C$ , so

$$\begin{aligned} \mathbb{P}\left(A\left(\left[z_{\pi(i)}\right]_{i=1}^{N+1}\right) \in B \mid C\right) &= \mathbb{P}\left(\left(A\left(\left[z_{\pi(j)}\right]_{j=1}^{N+1}\right)\right)_i \in B_i; \forall i \in \{1, 2, \dots, N+1\} \mid C\right) \\ &= \prod_{i=1}^{N+1} \mathbb{P}\left(\left(A\left(\left[z_{\pi(j)}\right]_{j=1}^{N+1}\right)\right)_i \in B_i \mid C\right). \end{aligned}$$

Note that, from the conditional symmetry in 4,

$$\begin{aligned} \left(\left(A\left(\left[z_{\pi(j)}\right]_{j=1}^{N+1}\right)\right)_i \mid C\right) &= (T_i(z_{\pi(1)}, z_{\pi(2)}, \dots, z_{\pi(i-1)}, z_{\pi(i+1)}, z_{\pi(i+2)}, \dots, z_{\pi(N+1)}, z_{\pi(i)}) \mid C) \\ &\stackrel{d}{=} (T_{\pi(i)}(z_1, z_2, \dots, z_{\pi(i)-1}, z_{\pi(i)+1}, z_{\pi(i)+2}, \dots, z_{N+1}, z_{\pi(i)}) \mid C) \\ &= ((A(z))_{\pi(i)}) \mid C, \end{aligned}$$

so

$$\begin{aligned} \mathbb{P}\left(A\left(\left[z_{\pi(i)}\right]_{i=1}^{N+1}\right) \in B \mid C\right) &= \prod_{i=1}^{N+1} \mathbb{P}\left((A(z))_{\pi(i)} \in B_i \mid C\right) \\ &= \mathbb{P}\left((A(z))_{\pi(i)} \in B_i; \forall i \in \{1, 2, \dots, N+1\} \mid C\right) \\ &= \mathbb{P}\left(\left[(A(z))_{\pi(i)}\right]_{i=1}^{N+1} \in B \mid C\right). \end{aligned}$$

Note that all equalities here are held almost surely  $C$ . We then take expectation, so

$$\mathbb{P}\left(A\left(\left[z_{\pi(i)}\right]_{i=1}^{N+1}\right) \in B\right) = \mathbb{P}\left(\left[(A(z))_{\pi(i)}\right]_{i=1}^{N+1} \in B\right).$$

Thus,  $A$  is permutation equivariant in joint distribution.  $\square$

## D Proof of the Main Theorem

Now that we have proved other claims, we can put them together to prove the main theorem.

*Proof.* From the lemma 3, we have that the joint score function is permutation invariant in joint distribution. From the corollary 2, we then have that  $P$  is a valid p-value.  $\square$

## E Strict Relationship

As mentioned in the figure 1, the relationships

$$\text{DT} \subseteq \text{AS} \cap \text{I\&D},$$

and

$$\text{AS} \cup \text{I\&D} \subseteq \text{CI\&D} \subseteq \text{D}$$

can be strengthened from  $\subseteq$  to  $\subsetneq$  for some appropriate setting  $((\Omega, \mathcal{F}, \mathbb{P}_{\text{alg}}), \mathcal{Z}, \mathcal{S})$ .

Consider the case when  $N = 3$ ,  $\mathcal{Z}$  is  $\{0, 1\} \times [0, 1]$ , and  $\mathcal{S} = \mathbb{R}$ , while the algorithmic random space sufficiently rich that we can define 2 independent Brownian motions.

### E.1 CI&D $\subsetneq$ D

Since the random space is sufficiently rich to define 2 independent Brownian motions, the same construction in the section 3 will provide the example of  $T$  that is permutation invariance in distribution but does not yield a valid conformal p-value.<sup>4</sup> From the theorem 1, we can then conclude that  $T$  is not conditionally independence and permutation invariance in distribution. Therefore,

$$\text{CI\&D} \subsetneq \text{D}.$$

### E.2 DT $\subsetneq$ AS $\cap$ I&D

Let  $U(\cdot)$  be a random variable distributed according to  $\text{Unif}([0, 1])$ .

We define

$$T_i(z; \omega) := \begin{cases} U(\omega) & \text{if } x_1 = x_2 = 0 \text{ and } x_3 = 1 \\ U(\omega) & \text{if } x_1 = x_2 = 1 \text{ and } x_3 = 0 \\ 0 & \text{otherwise} \end{cases}$$

for all  $i \in \{1, 2, 3\}$ ,  $z \in \mathcal{Z}^3$ , and  $\omega \in \Omega$ .

Thus,  $T$  is not deterministic, but it is almost surely permutation invariant<sup>5</sup>

Next, we will show that this is independent and permutation invariant in distribution. As we have discussed, permutation invariant in distribution is weaker than the almost sure condition, so it is automatically satisfied. For any  $(x_1, x_2, x_3) \in \{0, 1\}^3$ , we have that

$$\{T_1(x_2, x_3, x_1), T_2(x_1, x_3, x_2), T_3(x_1, x_2, x_3)\}$$

will have at least two of them being a degenerate random variable taking a value of 0. Thus, they are independent.

Therefore, we have shown that  $T$  is almost surely permutation invariant and independent and permutation invariant in distribution, but is not deterministically permutation invariant. Thus,

$$\text{DT} \subsetneq \text{AS} \cap \text{I\&D}.$$

### E.3 AS $\cup$ I&D $\subsetneq$ CI&D

Consider the case when  $U_i(\cdot) \stackrel{\text{iid}}{\sim} \text{Unif}([0, 1])$  for all  $i \in \{0, 1, 2, 3\}$ .

Define

$$T_i(z; \omega) = C_0(\omega) + C_i(\omega)$$

for all  $i \in \{1, 2, 3\}$ ,  $z \in \mathcal{Z}^3$ , and for almost surely  $\omega$ .

Thus, by choosing  $C(\cdot)$  to be  $C_0(\cdot)$ , we will then have that the collection of scores are independent conditioned on  $C$ . Thus,  $T$  is conditionally independent and permutation invariant in distribution.

However, it is obvious that  $T$  is not permutation invariant almost surely nor independent and permutation invariant in distribution. Thus,

$$\text{AS} \cup \text{I\&D} \subsetneq \text{CI\&D}.$$

---

<sup>4</sup>The construction in the section 3 assumes  $\mathcal{X} = [0, 1]$ , but  $\mathcal{X}$  is not used in the definition, so the result will be the same.

<sup>5</sup>Recall that we only require the permutation invariance over the first  $N = 2$  arguments.