# *Leave-One-Out Stable Conformal Prediction*

## Kiljae Lee, Yuan Zhang

The Ohio State University

# Why Uncertainty Quantification?

In many machine learning applications, **point predictions** dominate.
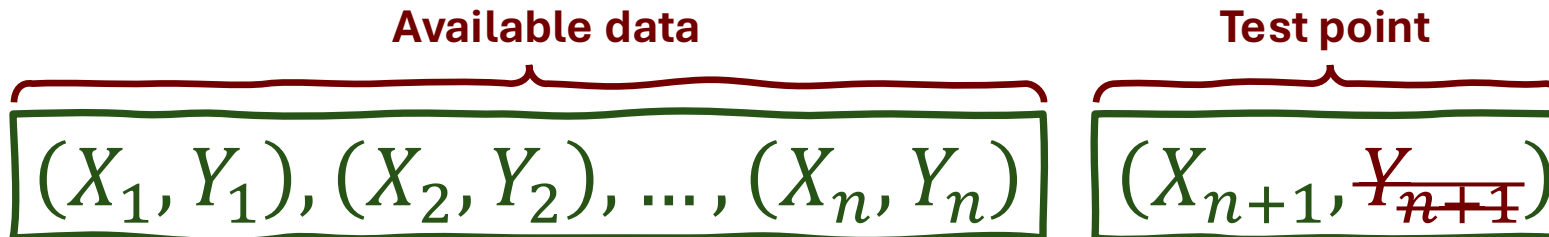
Medical AI: *"This tumor is malignant"*    Autonomous Driving: *"A pedestrian is detected"*

But how it sure?

***Uncertainty quantification*** *is not optional. It is critical.*

# Valid Prediction Set

Available data            Test point

$$(X_1, Y_1), (X_2, Y_2), \ldots, (X_n, Y_n) \quad (X_{n+1}, \cancel{Y_{n+1}})$$

**Our goal**: Construct **a valid set** $C(X_{n+1})$ that contains $Y_{n+1}$ with high probability *without any distribution assumption*.

### Definition (Validity)

*A set $C_\alpha(X_{n+1})$ is **valid** at a given level of $\alpha$ if $P(Y_{n+1} \in C_\alpha(X_{n+1})) \geq 1 - \alpha$ holds.*

- **Conformal Prediction (CP)** provides valid prediction sets, relying solely on the *i.i.d. assumption of data*.

- A core concept, **nonconformity scores**, quantify *how well a given label $Y_{n+1}$ conforms* to a trained model.

**(Examples)**

Regression: $S(Y_i, \hat{f}(X_i)) = \left| Y_i - \hat{f}(X_i) \right|$

Classification: $S(Y_i, \hat{p}(X_i)) = 1 - \hat{p}_{Y_i}(X_i)$

**(WARNING!)** *If some data points are favored* in training, the validity breaks.

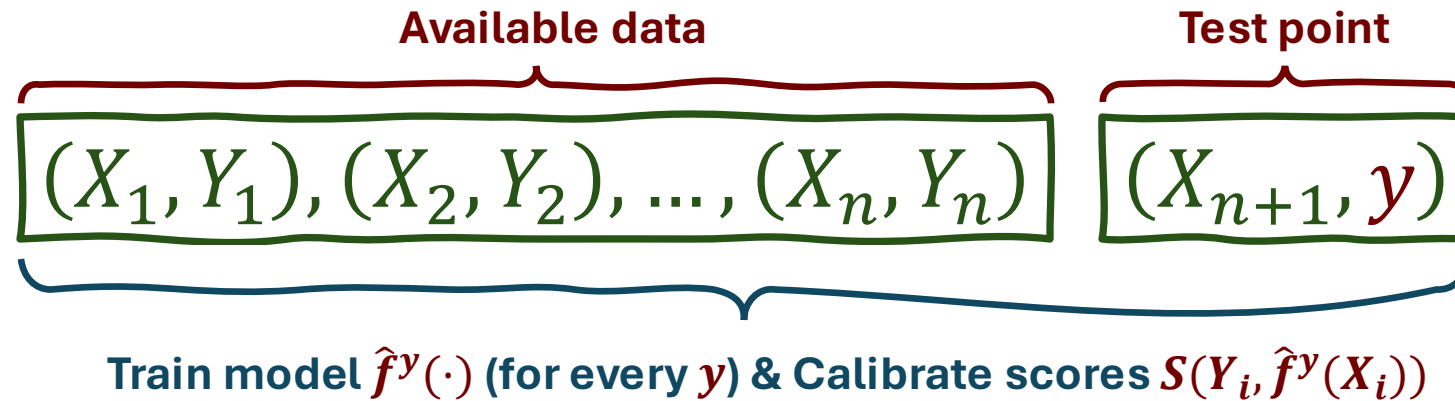There are **two** standard approaches *to guarantee the equivalence (a.k.a. Exchangeability)* across the $S_i$'s.



Available data       Test point

$$(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) \quad (X_{n+1}, y)$$

**Full CP (Vovk, 2005)** computes nonconformity scores by *fitting a model over all possible guesses* of the test label.

**Available data**                    **Test point**

$$(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) \quad (X_{n+1}, y)$$

Train model $\hat{f}^y(\cdot)$ (for every $y$) & Calibrate scores $S(Y_i, \hat{f}^y(X_i))$
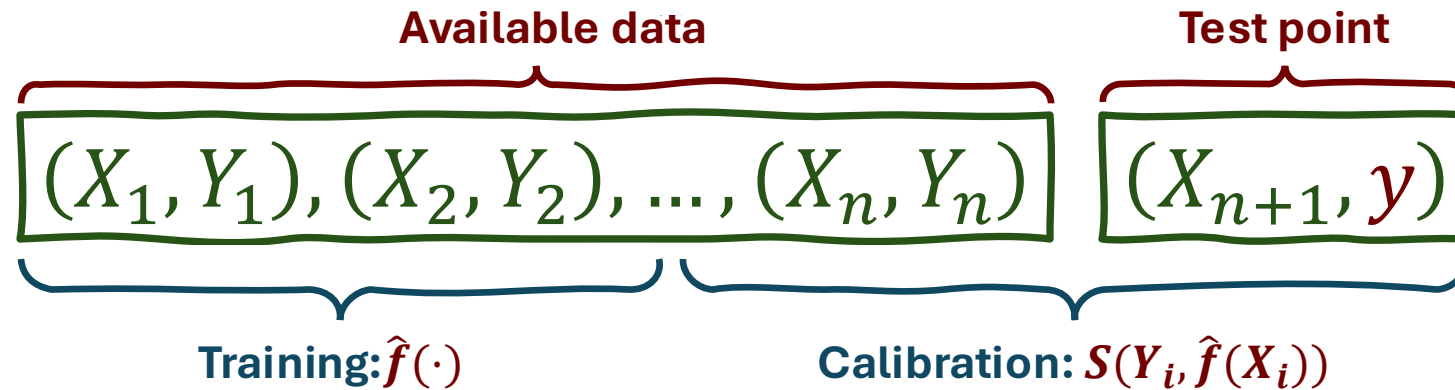
**Pros:** Allows the *use of entire dataset* for both training & calibration

**Cons:** Computational burden due to the *model retraining*

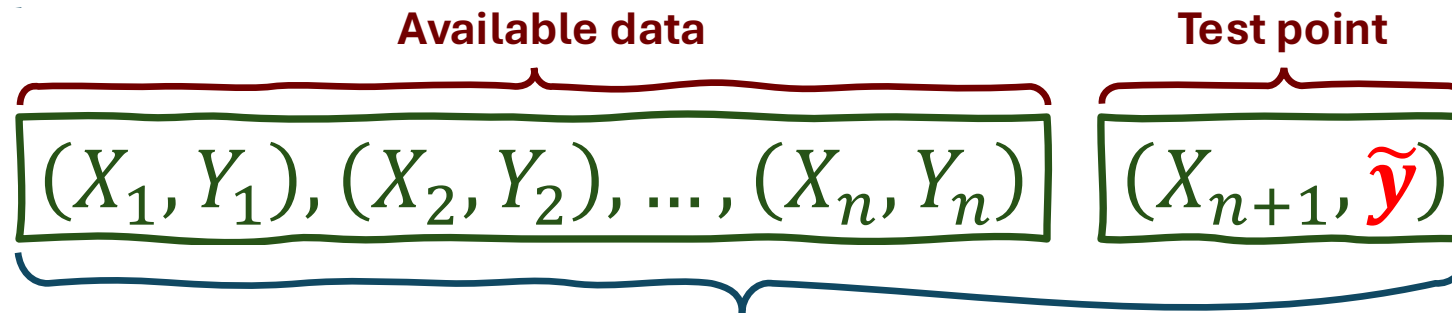**Split CP (Papadopoulos et al., 2002)** *divides the data* into training and calibration sets.

**Available data**

**Test point**

$$(X_1, Y_1), (X_2, Y_2), \ldots, (X_n, Y_n)$$ $$(X_{n+1}, y)$$

Training: $\hat{f}(\cdot)$

Calibration: $S(Y_i, \hat{f}(X_i))$

**Pros:** Only *a single model fit* is required

**Cons:** *Less data* is available for both training &calibration

**Trade-off!** *Statistical Efficiency vs Computational Efficiency*

Returning to the Full CP, what if we **replace $y$ with an arbitrary guess $\tilde{y}$**?

**Available data**

**Test point**

$$(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) \qquad (X_{n+1}, \tilde{y})$$
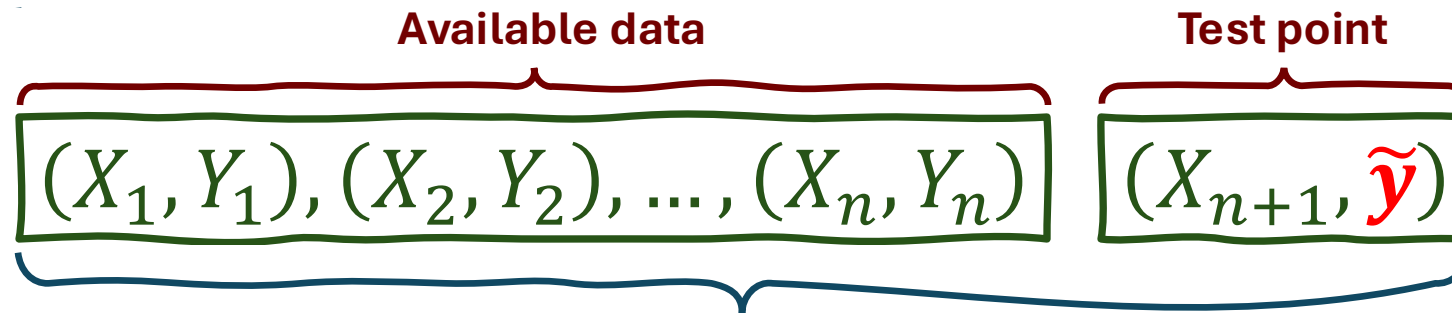
Train model $\hat{f}^{\tilde{y}}(\cdot)$ & Calibrate scores $S(Y_i, \hat{f}^{\tilde{y}}(X_i))$ : **No need to retrain**

Ndiaye (2022) proposed **Replace-One Stable CP (RO-StabCP):**
"A single label replacement yields **controllable changes.**"

Returning to the Full CP, what if we ***replace y with an arbitrary guess*** $\widetilde{y}$?

**Available data**  **Test point**

$$(X_1, Y_1), (X_2, Y_2), \ldots, (X_n, Y_n) \qquad (X_{n+1}, \widetilde{y})$$

Train model $\hat{f}^{\widetilde{y}}(\cdot)$ & Calibrate scores $S(Y_i, \hat{f}^{\widetilde{y}}(X_i))$ : No need to retrain
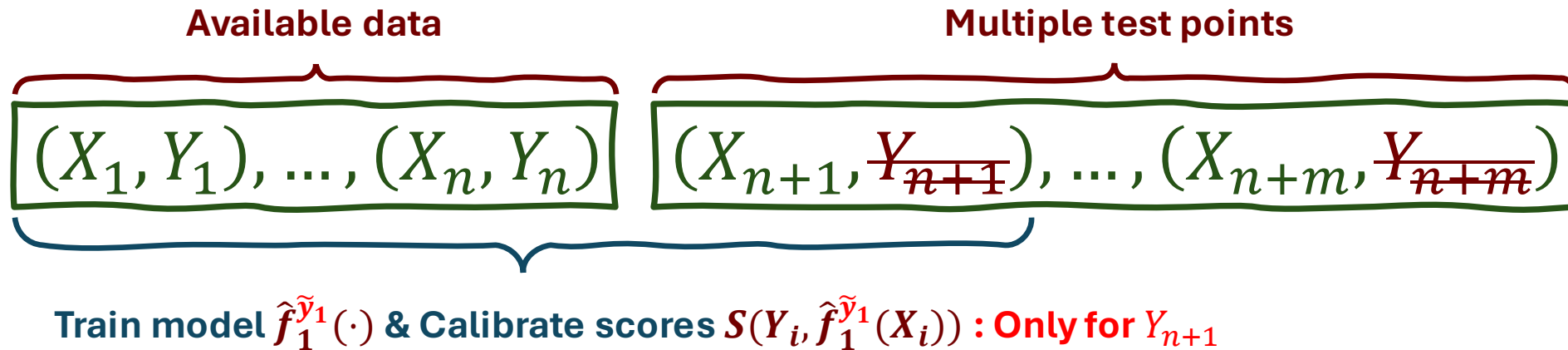
**Definition (Algorithmic Stability: Replace-One)**

*A model $\hat{f}$ is **replace-one stable**, if for all $i$, there exists $\tau^{RO}$ such that*

$$\sup_{z, y, \widetilde{y}} \left| S\left(z, \hat{f}^{y}(X_i)\right) - S\left(z, \hat{f}^{\widetilde{y}}(X_i)\right) \right| < \tau^{RO}.$$

However, in real-world scenarios, we often have ***more than one test point***.

**Available data**

**Multiple test points**

$$(X_1, Y_1), \ldots, (X_n, Y_n) \quad (X_{n+1}, \cancel{Y_{n+1}}), \ldots, (X_{n+m}, \cancel{Y_{n+m}})$$

**Train model** $\hat{f}_1^{\tilde{y}_1}(\cdot)$ **& Calibrate scores** $S(Y_i, \hat{f}_1^{\tilde{y}_1}(X_i))$ **: Only for** $Y_{n+1}$

In the above case, **RO-StabCP** requires $m$ model trainings since ***the model still includes*** $X_{n+j}$.

**THE OHIO STATE UNIVERSITY**

*"What if we train a model <u>without any test point information?</u>"*

**Available data**

**Multiple test points**

$$(X_1, Y_1), \dots, (X_n, Y_n) \quad (X_{n+1}, \cancel{Y_{n+1}}), \dots, (X_{n+m}, \cancel{Y_{n+m}})$$

**Train model** $\hat{f}(\cdot)$
**Calibrate scores** $S(Y_i, \hat{f}(X_i))$

**: For every** $Y_{n+j}, \ j = 1, \dots, m$

*"**Excluding a single point** still yields controllable changes."*

*"What if we train a model **without any test point information?**"*

**Available data**

**Multiple test points**

$$(X_1, Y_1), \ldots, (X_n, Y_n) \quad (X_{n+1}, \cancel{Y_{n+1}}), \ldots, (X_{n+m}, \cancel{Y_{n+m}})$$

**Train model** $\hat{f}(\cdot)$
**Calibrate scores** $S(Y_i, \hat{f}(X_i))$

: **For every** $Y_{n+j}, \, j = 1, \ldots, m$

**Definition (Algorithmic Stability: Leave-One-Out)**

*A model $\hat{f}$ is **leave-one-out stable**, if for all $i$, there exists $\tau^{LOO}$ such that*

$$\sup_{z,y} \left| S\left(z, \hat{f}^y(X_i)\right) - S\left(z, \hat{f}(X_i)\right) \right| < \tau^{LOO}.$$

| Method | # Model Fits | # Prediction Calls | Stable |
|---|---|---|---|
| FullCP | $\lvert \mathcal{Y} \rvert \cdot m$ | $(n+1) \cdot \lvert \mathcal{Y} \rvert \cdot m$ | ✓ |
| SplitCP | 1 | $n+m$ | ✗ |
| RO-StabCP | $m$ | $(n+1) \cdot m$ | ✓ |
| **LOO-StabCP** | **1** | $\boldsymbol{n+m}$ | ✓ |