

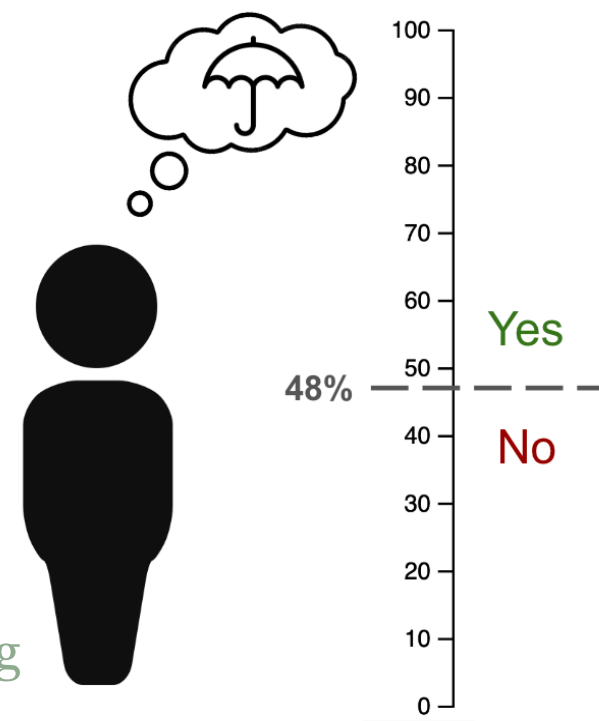
# Calibration and Trustworthy Decision Making



Princewill Okoroafor  
Cornell University  
(Incoming Postdoc at Harvard SEAS)















Based on joint work with Michael P. Kim & Robert Kleinberg

<https://arxiv.org/abs/2501.17205>



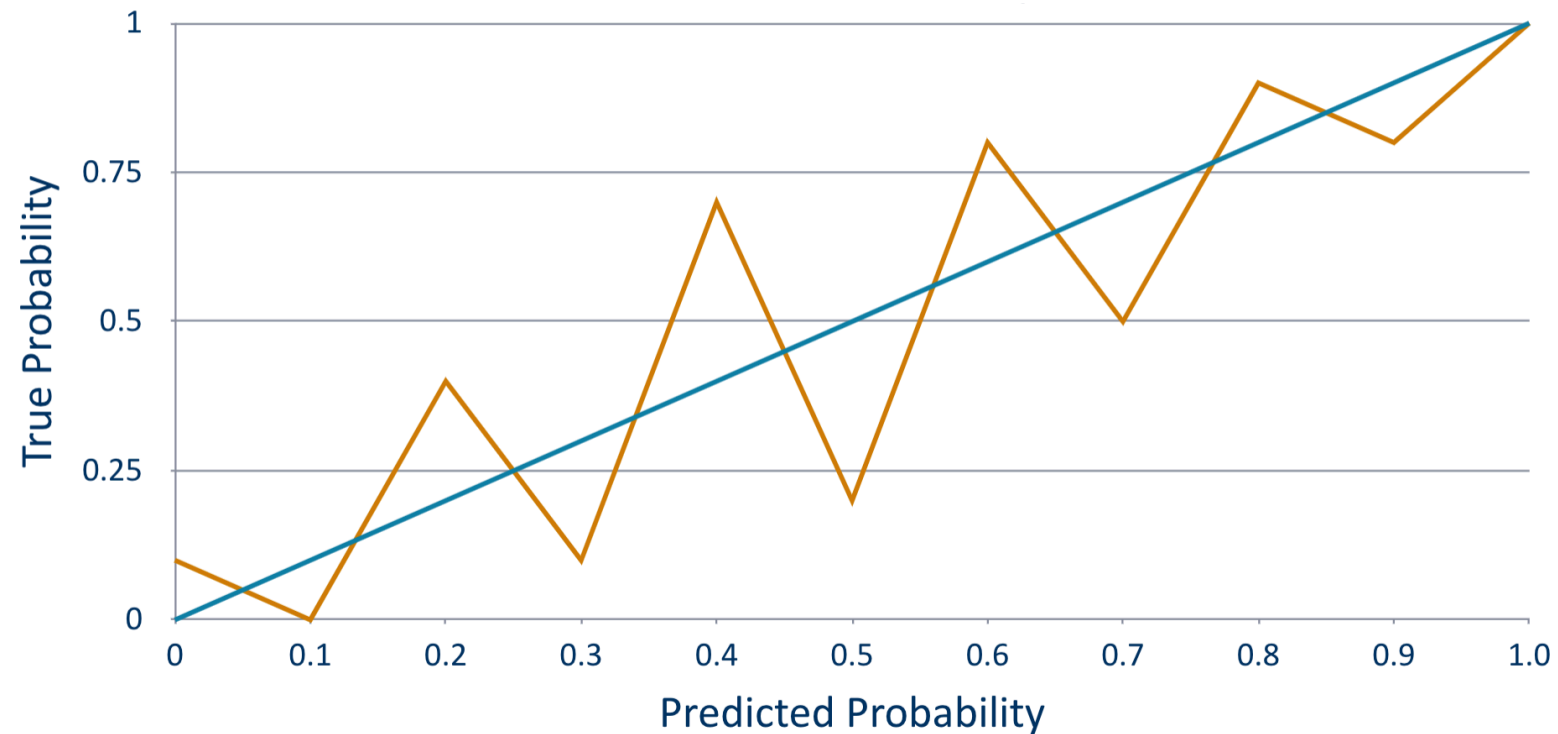
# What is calibration?








**Prediction values should mean what they say**

|        |                |   |                  |   |
|--------|----------------|---|------------------|---|
| Sun 07 | <b>51°/30°</b> |    | AM Clouds/PM Sun |  14%   |
| Mon 08 | <b>61°/41°</b> |    | Partly Cloudy    |  7%    |
| Tue 09 | <b>69°/47°</b> |    | Partly Cloudy    |  21%   |
| Wed 10 | <b>59°/47°</b> |    | Showers          |  56%   |
| Thu 11 | <b>62°/52°</b> |    | Showers          |  70%   |
| Fri 12 | <b>56°/41°</b> |  | Showers          |  58% |
| Sat 13 | <b>52°/39°</b> |  | AM Showers       |  32% |

# What is calibration?

Prediction values should mean what they say



|        |         |  |     |
|--------|---------|--|-----|
| Sun 07 | 51°/30° |  AM Clouds/PM Sun | 14% |
| Mon 08 | 61°/41° |  Partly Cloudy    | 7%  |
| Tue 09 | 69°/47° |  Partly Cloudy    | 21% |
| Wed 10 | 59°/47° |  Showers          | 56% |
| Thu 11 | 62°/52° |  Showers          | 70% |
| Fri 12 | 56°/41° |  Showers          | 58% |
| Sat 13 | 52°/39° |  AM Showers       | 32% |

# Binary Prediction Setting

Applications in weather forecasting, prediction markets, etc

Each day  $t = 1, 2, \dots, T$ ,

- Nature chooses an outcome  $y_t \in \{0, 1\}$ ,
- Before observing outcome, Forecaster makes a prediction  $p_t \in [0, 1]$

**Goal:** low calibration error

# Binary Prediction Setting

Applications in weather forecasting, prediction markets, etc

Each day  $t = 1, 2, \dots, T$ ,

- Nature chooses an outcome  $y_t \in \{0, 1\}$ ,
- Before observing outcome, Forecaster makes a prediction  $p_t \in [0, 1]$

**Goal:** low calibration error

$$\text{calerr}(t) = \sum_{p \in [0, 1]} \underbrace{n_t(p)}_{\substack{\text{number of times } p \text{ was predicted}}} \left| p - \underbrace{\frac{m_t(p)}{n_t(p)}}_{\substack{\text{average of outcomes when } p \text{ was predicted}}} \right|$$

sum of outcomes when p was predicted

# Calibration is very non-smooth

| <b>T =</b> | <b>1</b>                 | <b>2</b>                 | <b>...</b>               | <b>T/2</b>               | <b>...</b>               | <b>T</b>                 |
|------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| Nature     | 0                        | 0                        | 0                        | 1                        | 1                        | 1                        |
| Forecaster | $\frac{1}{2} - \epsilon$ | $\frac{1}{2} - \epsilon$ | $\frac{1}{2} - \epsilon$ | $\frac{1}{2} + \epsilon$ | $\frac{1}{2} + \epsilon$ | $\frac{1}{2} + \epsilon$ |

# Calibration is very non-smooth

| <b>T =</b> | <b>1</b>                 | <b>2</b>                 | <b>...</b>               | <b>T/2</b>               | <b>...</b>               | <b>T</b>                 |
|------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| Nature     | 0                        | 0                        | 0                        | 1                        | 1                        | 1                        |
| Forecaster | $\frac{1}{2} - \epsilon$ | $\frac{1}{2} - \epsilon$ | $\frac{1}{2} - \epsilon$ | $\frac{1}{2} + \epsilon$ | $\frac{1}{2} + \epsilon$ | $\frac{1}{2} + \epsilon$ |

Smooth Calibration/Distance to Calibration [KF'o8, FH'18, BGHN'23]

# Calibration is not incentive-compatible

| <b>T =</b> | <b>1</b>                             | <b>...</b>                           | <b>T/m</b>                           | <b>...</b>                           | <b>2T/m</b>                          | <b>...</b> | <b>T</b> |
|------------|--------------------------------------|--------------------------------------|--------------------------------------|--------------------------------------|--------------------------------------|------------|----------|
| Nature     | $\text{Ber}\left(\frac{1}{m}\right)$ | $\text{Ber}\left(\frac{1}{m}\right)$ | $\text{Ber}\left(\frac{2}{m}\right)$ | $\text{Ber}\left(\frac{2}{m}\right)$ | $\text{Ber}\left(\frac{3}{m}\right)$ | ...        | ...      |
| Forecaster | $\frac{1}{m}$                        | $\frac{1}{m}$                        | $\frac{2}{m}$                        | $\frac{2}{m}$                        | $\frac{3}{m}$                        | ...        | ...      |



# Calibration is not incentive-compatible

| <b>T =</b> | <b>1</b>                             | <b>...</b>                           | <b>T/m</b>                           | <b>...</b>                           | <b>2T/m</b>                          | <b>...</b> | <b>T</b> |
|------------|--------------------------------------|--------------------------------------|--------------------------------------|--------------------------------------|--------------------------------------|------------|----------|
| Nature     | $\text{Ber}\left(\frac{1}{m}\right)$ | $\text{Ber}\left(\frac{1}{m}\right)$ | $\text{Ber}\left(\frac{2}{m}\right)$ | $\text{Ber}\left(\frac{2}{m}\right)$ | $\text{Ber}\left(\frac{3}{m}\right)$ | ...        | ...      |
| Forecaster | $\frac{1}{m}$                        | $\frac{1}{m}$                        | $\frac{2}{m}$                        | $\frac{2}{m}$                        | $\frac{3}{m}$                        | ...        | ...      |

Proper Scoring Losses e.g square loss are incentive-compatible

# Why Calibrate?

## **Key Question:**

*If full calibration has all these flaws that can be addressed by weaker notions, why bother with full calibration?*

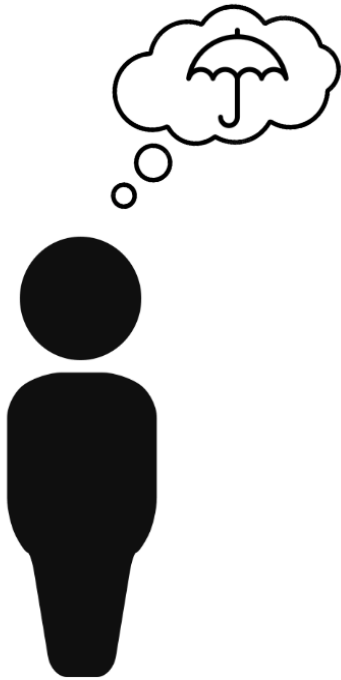
# Why Calibrate?

## **Key Question:**

*If full calibration has all these flaws that can be addressed by weaker notions, why bother with full calibration?*

Calibrated forecasts implies low regret decision making

# Sequential Decision Making



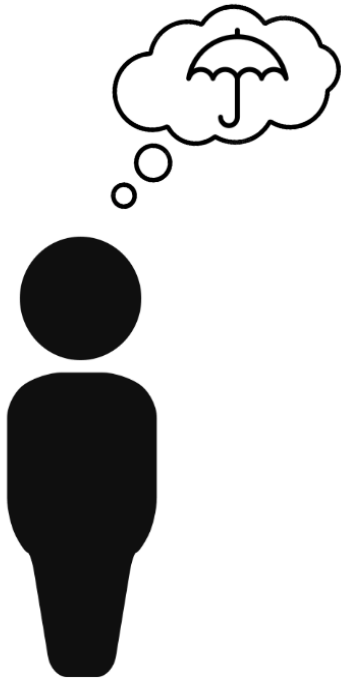
Each day  $t = 1, 2, \dots, T$ ,

- Nature chooses an outcome  $y_t \in \{0, 1\}$ ,
- Decision Maker chooses action  $a_t \in A$
- Decision Maker incurs loss  $\ell(a_t, y_t) \in [0, 1]$

**Goal:** Low Regret

$$\sum_{t=1}^T \ell(a_t, y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y_t) \leq o(T)$$

# Sequential Decision Making



Each day  $t = 1, 2, \dots, T$ ,

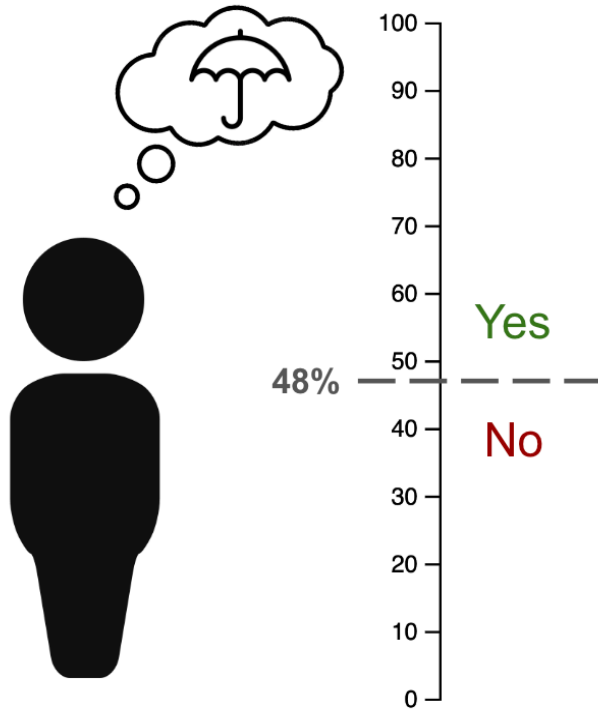
- Nature chooses an outcome  $y_t \in \{0, 1\}$ ,
- Decision Maker chooses action  $a_t \in A$
- Decision Maker incurs loss  $\ell(a_t, y_t) \in [0, 1]$

**Goal:** Low Regret

$$\sum_{t=1}^T \ell(a_t, y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y_t) \leq o(T)$$

Multiplicative Weights guarantees regret of  $O(\sqrt{T \log |A|})$

# Sequential Decision Making



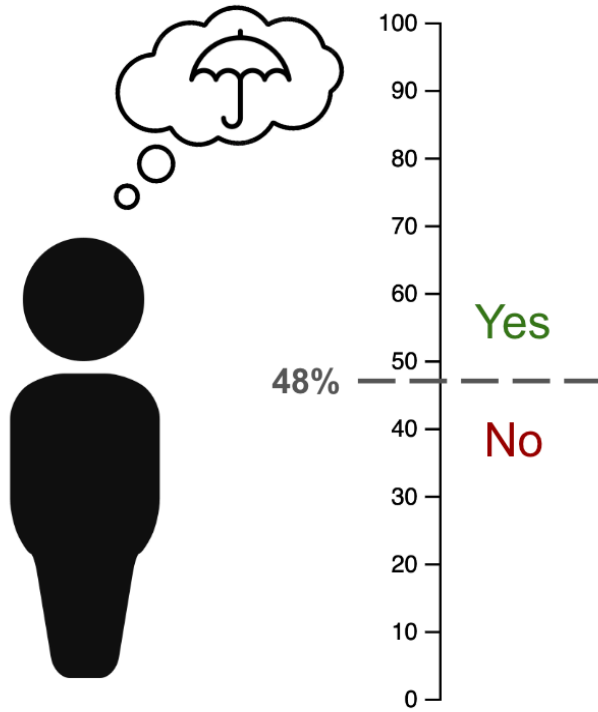
Each day  $t = 1, 2, \dots, T$ ,

- Nature chooses an outcome  $y_t \in \{0, 1\}$ ,
- Forecaster makes a prediction  $p_t \in [0, 1]$
- Decision Maker best-responds to prediction i.e chooses

$$a_t \in \operatorname{argmin}_{a \in A} \mathbb{E}_{y \sim \operatorname{Ber}(p_t)} [\ell(a, y)]$$

$$\text{Low Regret: } \sum_{t=1}^T \ell(a_t, y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y) \leq o(T)$$

# Sequential Decision Making



Each day  $t = 1, 2, \dots, T$ ,

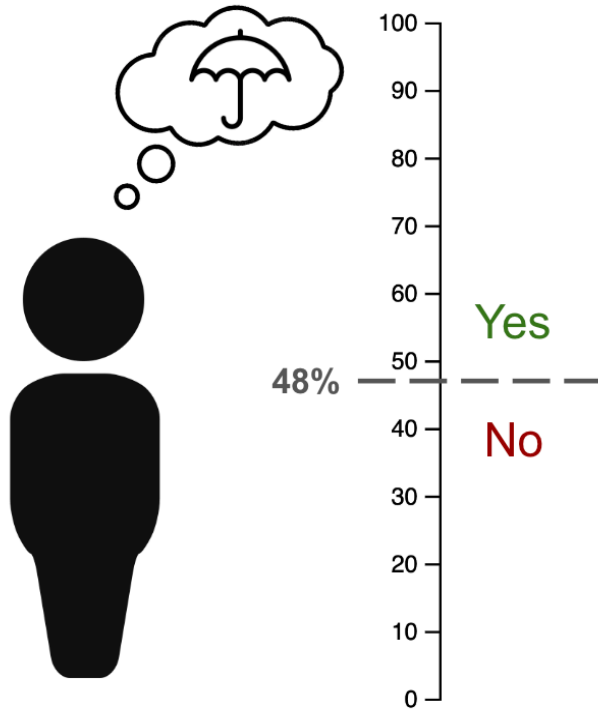
- Nature chooses an outcome  $y_t \in \{0, 1\}$ ,
- Forecaster makes a prediction  $p_t \in [0, 1]$
- Decision Maker best-responds to prediction i.e chooses

$$a_t \in \operatorname{argmin}_{a \in A} \mathbb{E}_{y \sim \operatorname{Ber}(p_t)} [\ell(a, y)]$$

$$\text{Low Regret: } \sum_{t=1}^T \ell(a_t, y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y) \leq o(T)$$

If predictions are calibrated, then Agent achieves low regret

# Sequential Decision Making



Each day  $t = 1, 2, \dots, T$ ,

- Nature chooses an outcome  $y_t \in \{0, 1\}$ ,
- Forecaster makes a prediction  $p_t \in [0, 1]$
- Decision Maker best-responds to prediction i.e chooses

$$a_t \in \operatorname{argmin}_{a \in A} \mathbb{E}_{y \sim \operatorname{Ber}(p_t)} [\ell(a, y)]$$

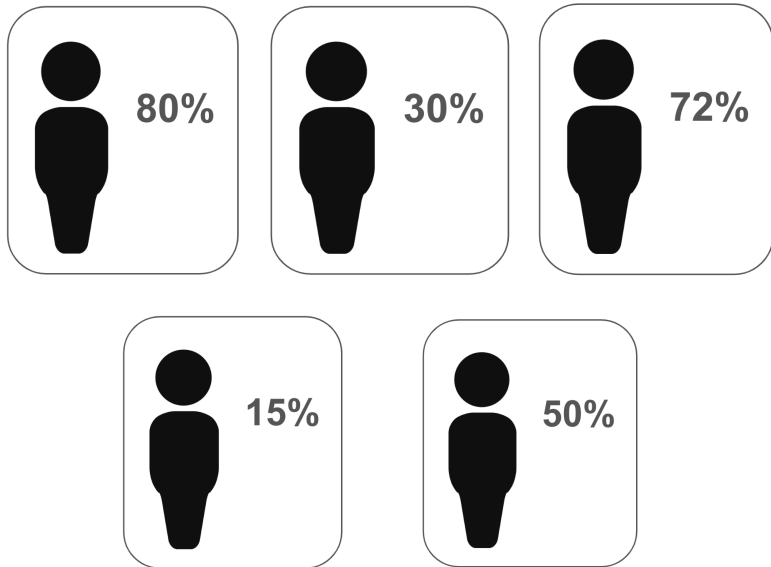
$$\text{Low Regret: } \sum_{t=1}^T \ell(a_t, y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y) \leq o(T)$$

Agent's regret is bounded by Calibration error rate  $f(T)$



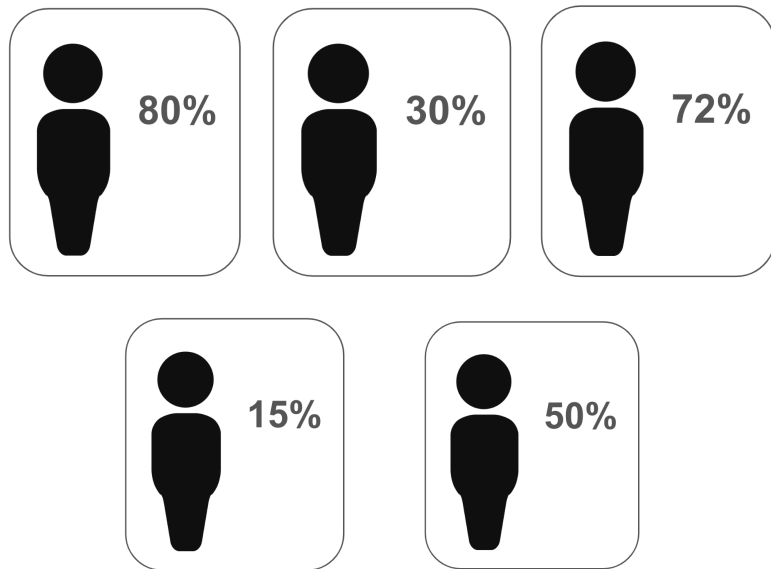
# Trustworthy Decision Making

Participants at COLT Workshop



# Trustworthy Decision Making

Participants at COLT Workshop



Every agent wants to obtain low regret with respect to their loss function  $\ell \in L$

$$\sup_{\ell \in L} \sum_{t=1}^T \ell(a_t, y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y) \leq o(T)$$

**Theorem:** (KLST'23)

Calibrated forecasts guarantee low regret (at a rate of  $f(T)$ ) simultaneously for every decision maker

# The rate of online calibration

**Lower Bound:** (QV'21, DDFGKO'25)

No algorithm can guarantee calibration at a rate better than  $\Omega(T^{0.543})$

**Upper Bound:** (FV'98, DDFGKO'25)

There exists an algorithm that guarantees calibration at a rate of  $O(T^{2/3-\epsilon})$

QV'21 - Stronger Lower Bounds for Calibration via Sidestepping

DDFGKO'25 - Breaking the  $T^{2/3}$  Barrier for Sequential Calibration

FV'98 - Forecast Hedging and Calibration

# The rate of online calibration

**Lower Bound:** (QV'21, DDFGKO'25)

No algorithm can guarantee calibration at a rate better than  $\Omega(T^{0.543})$

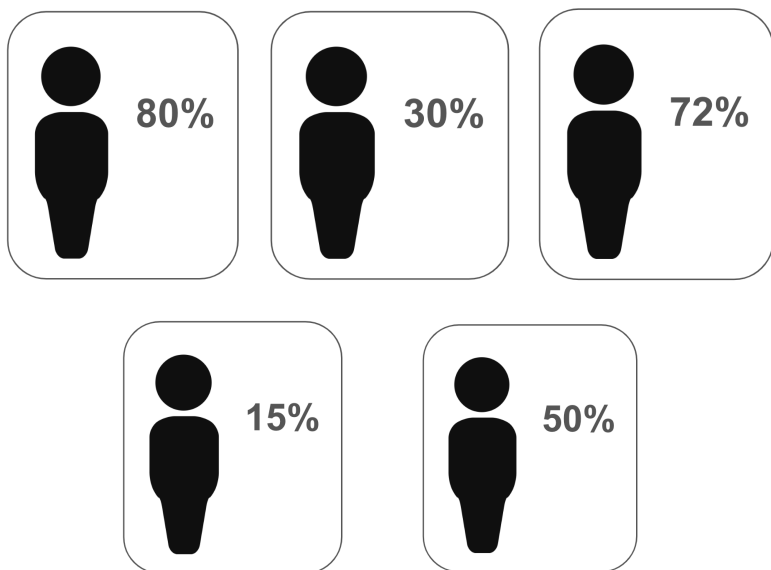
**Upper Bound:** (FV'98, DDFGKO'25)

There exists an algorithm that guarantees calibration at a rate of  $O(T^{2/3-\epsilon})$

*Is calibration even necessary for predictions to be trustworthy for decisions?*

# Trustworthy Decision Making

Participants at COLT Workshop



Every agent wants to obtain low regret with respect to their loss function  $\ell \in L$

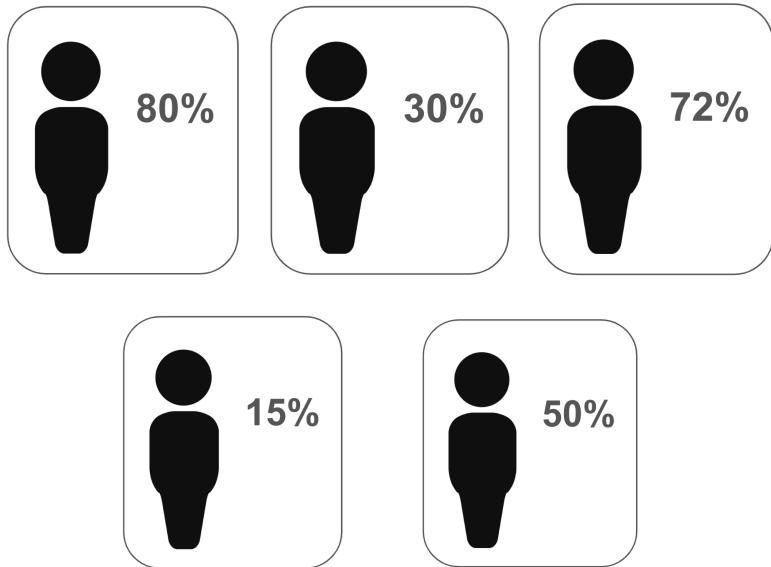
$$\sup_{\ell \in L} \sum_{t=1}^T \ell(k_{\ell}(p_t), y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y_t) \leq o(T)$$

**Theorem:** (U-Calibration - KLST'23)

There exists an algorithm that makes predictions such that regret of every Agent is bounded by  $O(\sqrt{T})$

# Trustworthy Decision Making

Participants at COLT Workshop



Every agent wants to obtain low regret with respect to their loss function  $\ell \in L$

$$\sup_{\ell \in L} \sum_{t=1}^T \ell(k_{\ell}(p_t), y_t) - \min_{a \in A} \sum_{t=1}^T \ell(a, y_t) \leq o(T)$$

**Theorem:** (U-Calibration - KLST'23)

There exists an algorithm that makes predictions such that regret of every Agent is bounded by  $O(\sqrt{T})$

*Calibration is NOT necessary for predictions to be trustworthy for decisions*

# Decision Making with Expert Advice

Class of Experts (or Hypothesis class)

 Take action A



 Take action C



 Take action B



In online learning, we have a class of experts to help inform our decisions.

In supervised learning, we use a hypothesis class or a class of models for decision making.

# Omnipredictors



Single predictor, ***simultaneous loss minimizer*** for many losses

(Gopalan, Kalai, Reingold, Sharan, Wieder; ITCS'22)

**Omniprediction:** For loss class  $L$ , hypothesis class  $H$ ,  $\varepsilon > 0$ ,  
find predictor  $p$  such that ***for every***  $\ell \in L$

$$\mathbb{E} \left[ \ell \left( k_\ell \circ p(x), y \right) \right] \leq \min_{h \in H} \mathbb{E} \left[ \ell \left( h(x), y \right) \right] + \varepsilon$$

$k_\ell$  is the “best response” function



# Omnipredictors



Single predictor, ***simultaneous loss minimizer*** for many losses

(Gopalan, Kalai, Reingold, Sharan, Wieder; ITCS'22)

**Omniprediction:** For loss class  $L$ , hypothesis class  $H$ ,  $\varepsilon > 0$ ,  
find predictor  $p$  such that ***for every***  $\ell \in L$

$$\mathbb{E} \left[ \ell \left( k_\ell \circ p(x), y \right) \right] \leq \min_{h \in H} \mathbb{E} \left[ \ell \left( h(x), y \right) \right] + \varepsilon$$

Recovers Loss Minimization

$k_\ell$  is the “best response” function

# Omnipredictors



Single predictor, ***simultaneous loss minimizer*** for many losses

(Gopalan, Kalai, Reingold, Sharan, Wieder; ITCS'22)

**Omniprediction:** For loss class  $L$ , hypothesis class  $H$ ,  $\varepsilon > 0$ ,  
find predictor  $p$  such that ***for every***  $\ell \in L$

$$\mathbb{E} \left[ \ell \left( k_\ell \circ p(x), y \right) \right] \leq \min_{h \in H} \mathbb{E} \left[ \ell \left( h(x), y \right) \right] + \varepsilon$$

No Retraining Necessary!

$k_\ell$  is the “best response” function

# Towards Optimal Omniprediction

## **Key Question:**

*What is the Complexity of Omniprediction,  
and how does it compare to Loss Minimization?*

# Towards Optimal Omniprediction

## Key Question:

*What is the Sample Complexity of Omniprediction,  
and how does it compare to Loss Minimization?*

Optimal Loss Minimization:

$$\Theta(d_{\ell \circ H}/\varepsilon^2)$$

(GHKRW'23)

Omniprediction, so far:

$$O(d_{L \circ H}/\varepsilon^6 + 1/\varepsilon^{10})$$

Is the gap in complexity inherent?

# Towards Optimal Omniprediction

## Key Question:

*What is the Sample Complexity of Omniprediction,  
and how does it compare to Loss Minimization?*

(O., Kleinberg, Kim'25)

Optimal Loss Minimization:

$$\Theta(d_{\ell \circ H}/\varepsilon^2)$$

Omniprediction:

$$\Theta(d_{L \circ H}/\varepsilon^2)$$

Sample Complexity of Omniprediction  $\approx$  Minimization of Single Loss!

# Towards Optimal Omniprediction

## Key Question:

*What is the Sample Complexity of Omniprediction,  
and how does it compare to Loss Minimization?*

(O., Kleinberg, Kim'25)

Optimal Loss Minimization:

$$\Theta(d_{\ell \circ H}/\varepsilon^2)$$

Omniprediction:

$$\Theta(d_{L \circ H}/\varepsilon^2)$$

Similar result for RKHS by DHIPT'25

Sample Complexity of Omniprediction  $\approx$  Minimization of Single Loss!

# Learning Omnipredictors

**Theorem:** (Gopalan, Hu, Kim, Reingold, Wieder; ITCS'23)

Calibration +  $L \circ H$ -Multiaccuracy  $\implies$  Omniprediction

# Learning Omnipredictors

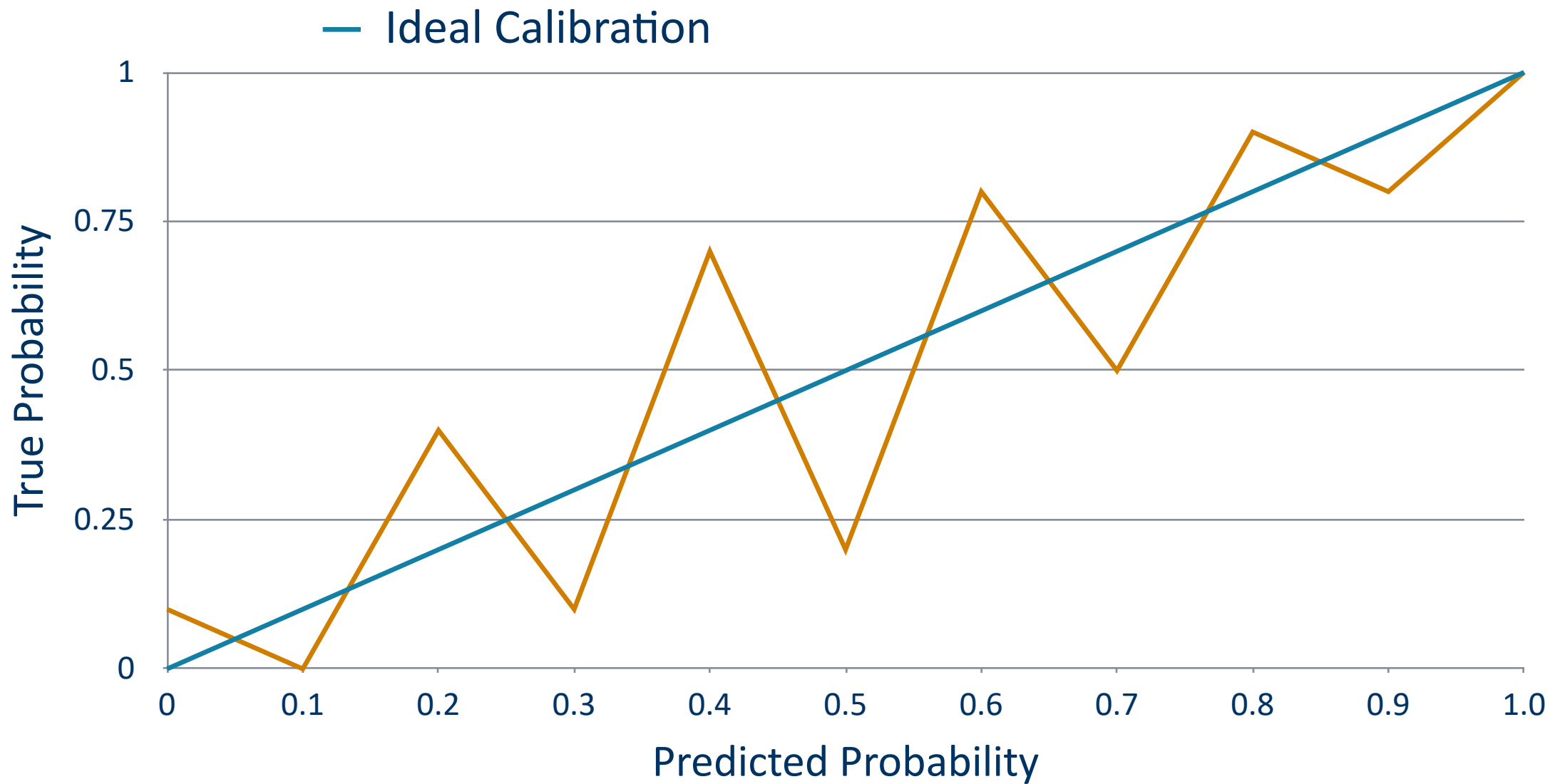
**Theorem:** (Gopalan, Hu, Kim, Reingold, Wieder; ITCS'23)

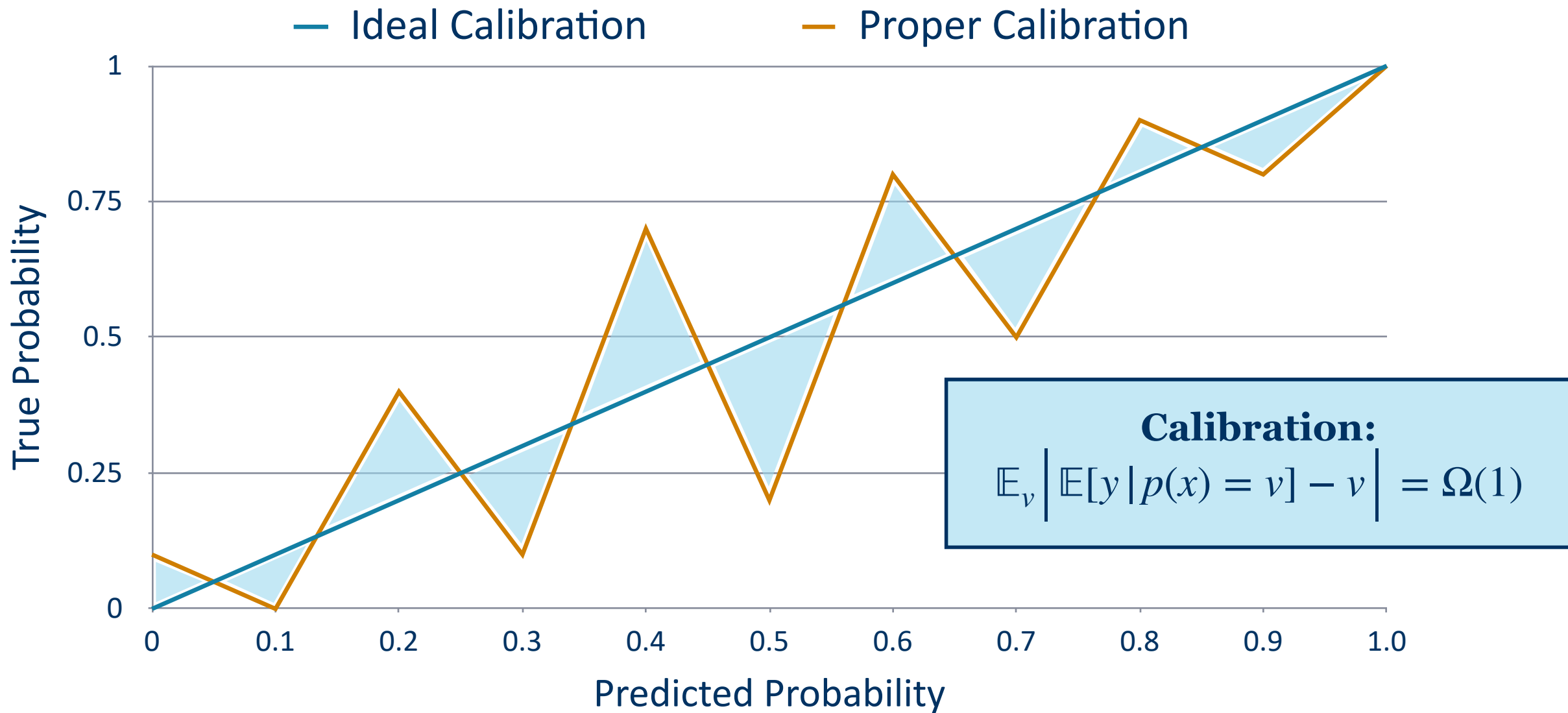
Calibration +  $L \circ H$ -Multiaccuracy  $\implies$  Omniprediction

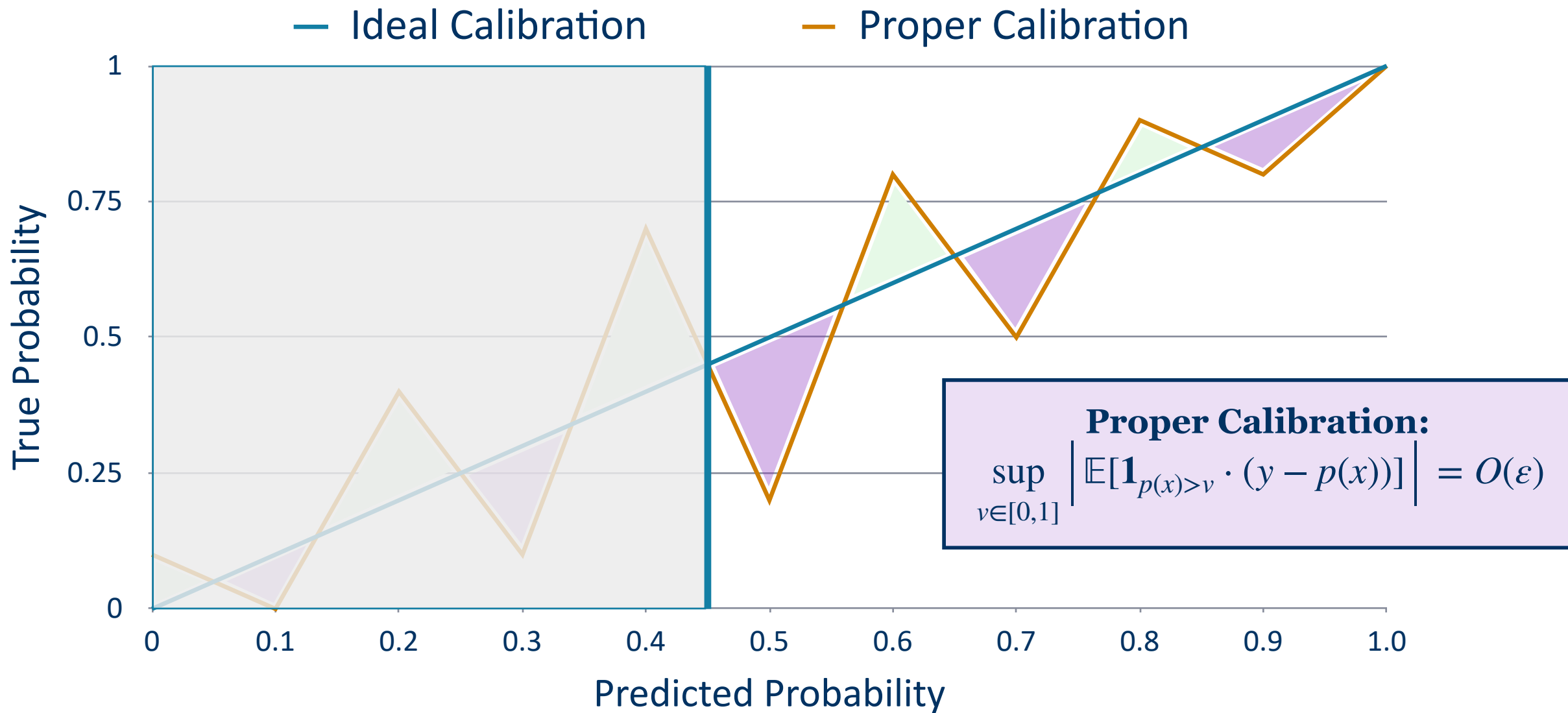
**Theorem:** (O., Kleinberg, Kim'25)

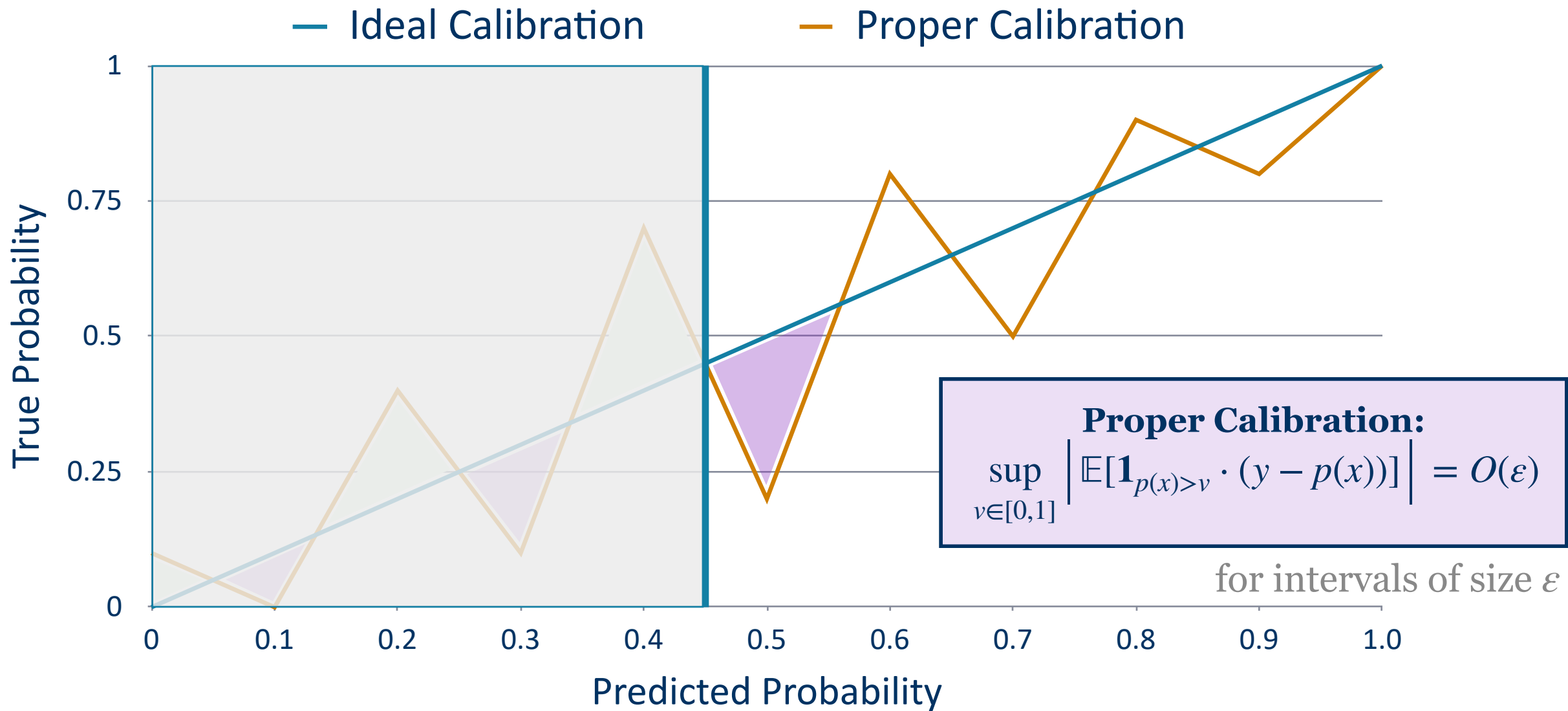
*Proper Calibration* +  $L \circ H$ -Multiaccuracy  $\implies$  Omniprediction











# Achieving Proper Calibration

**Theorem:** (O., Kleinberg, Kim'25)

Proper Calibration is efficiently testable using  $O(1/\epsilon^2)$  samples

# Achieving Proper Calibration

Calibration is not even statistically testable  
for arbitrary predictors

**Theorem:** (O., Kleinberg, Kim'25)

Proper Calibration is efficiently testable using  $O(1/\epsilon^2)$  samples

# Achieving Proper Calibration

Calibration is not even statistically testable for arbitrary predictors

**Theorem:** (O., Kleinberg, Kim'25)

Proper Calibration is efficiently testable using  $O(1/\epsilon^2)$  samples

**Theorem:** (O., Kleinberg, Kim'25)

There exists an algorithm that achieves proper calibration at a rate of  $O(\sqrt{T})$

# Achieving Proper Calibration

Calibration is not even statistically testable for arbitrary predictors

**Theorem:** (O., Kleinberg, Kim'25)

Proper Calibration is efficiently testable using  $O(1/\epsilon^2)$  samples

Recall Calibration requires  $\Omega(T^{0.534})$

**Theorem:** (O., Kleinberg, Kim'25)

There exists an algorithm that achieves proper calibration at a rate of  $O(\sqrt{T})$



# Achieving Proper Calibration

Calibration is not even statistically testable for arbitrary predictors

**Theorem:** (O., Kleinberg, Kim'25)

Proper Calibration is efficiently testable using  $O(1/\epsilon^2)$  samples

Recall Calibration requires  $\Omega(T^{0.534})$

**Theorem:** (O., Kleinberg, Kim'25)

There exists an algorithm that achieves proper calibration at a rate of  $O(\sqrt{T})$

Contemporary Work by (QZ'25,RSBRW'25)

QZ'25 - Truthfulness of Decision-Theoretic Calibration Measures

RSBRW'25 - Can a calibration metric be both testable and actionable?

# Online Omniprediction

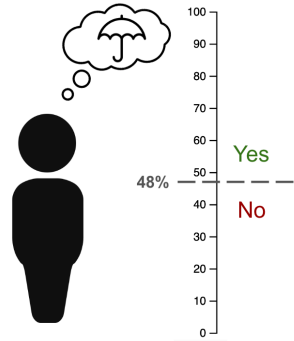
(O., Kleinberg, Kim'25)

**Theorem:** Given a loss class  $L$ , a hypothesis class  $H$  and a sequence of adversarially chosen pairs  $(x_t, y_t)$ , there exists an online algorithm that outputs predictors that achieve expected regret

$$\tilde{O} \left( \sqrt{T \cdot d_{L \circ H}^{\text{seq}}} \right)$$

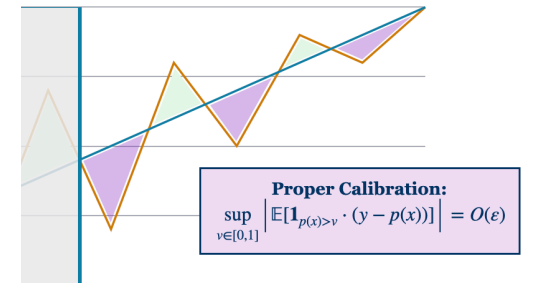
where  $d_{L \circ H}$  = sequential dimension of function class derived from  $L \circ H$

# Key Takeaway



While Calibration provides lots of guarantees, understanding the tasks that we need it for helps us design better calibration measures for those tasks

Proper calibration is a calibration measure powerful enough to guarantee omniprediction while still efficiently achievable



# Future Work

Some progress by LRS'25

## **What about settings with multiple outcomes?**

*Can we generalize these results to multi-class and real valued settings? Can we achieve proper calibration efficiently in higher dimensions?*

## **For what tasks is full calibration necessary?**

*While calibration is sufficient for trustworthy decision making, it's not necessary.*