

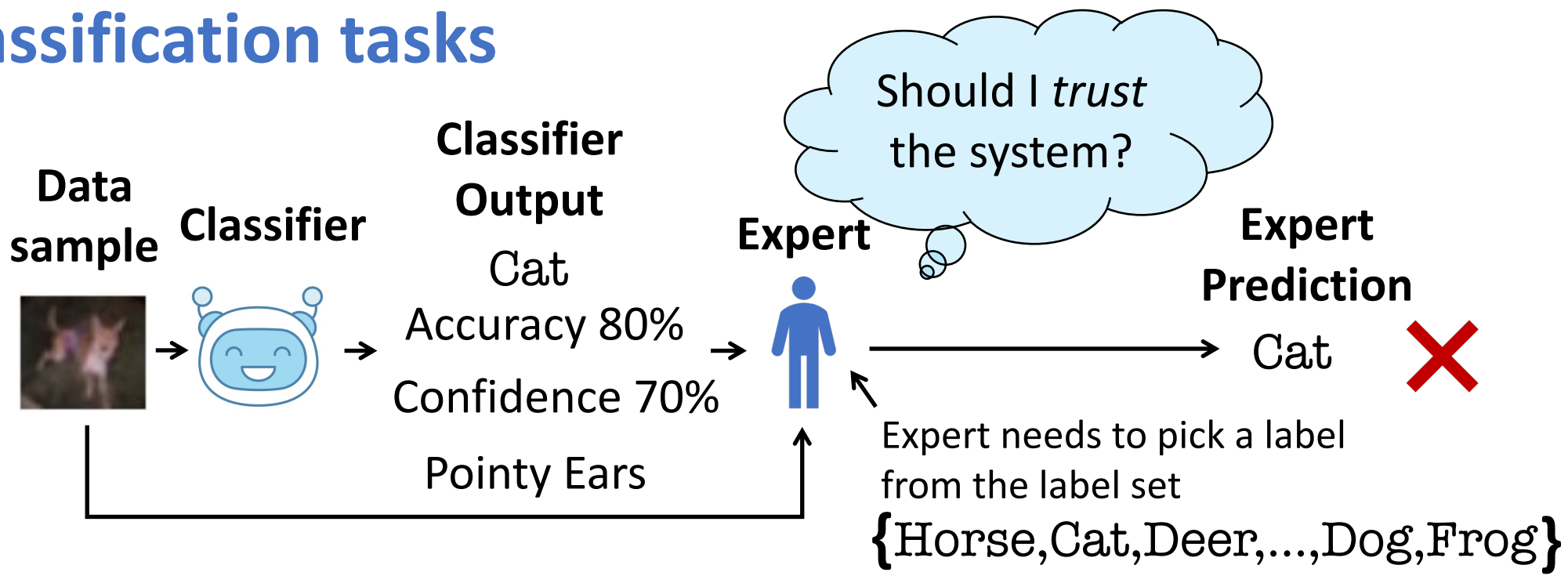


Designing Decision Support Systems with Counterfactual Prediction Sets

Max Planck Institute
for Software Systems

Eleni Straitouri and Manuel Gomez-Rodriguez

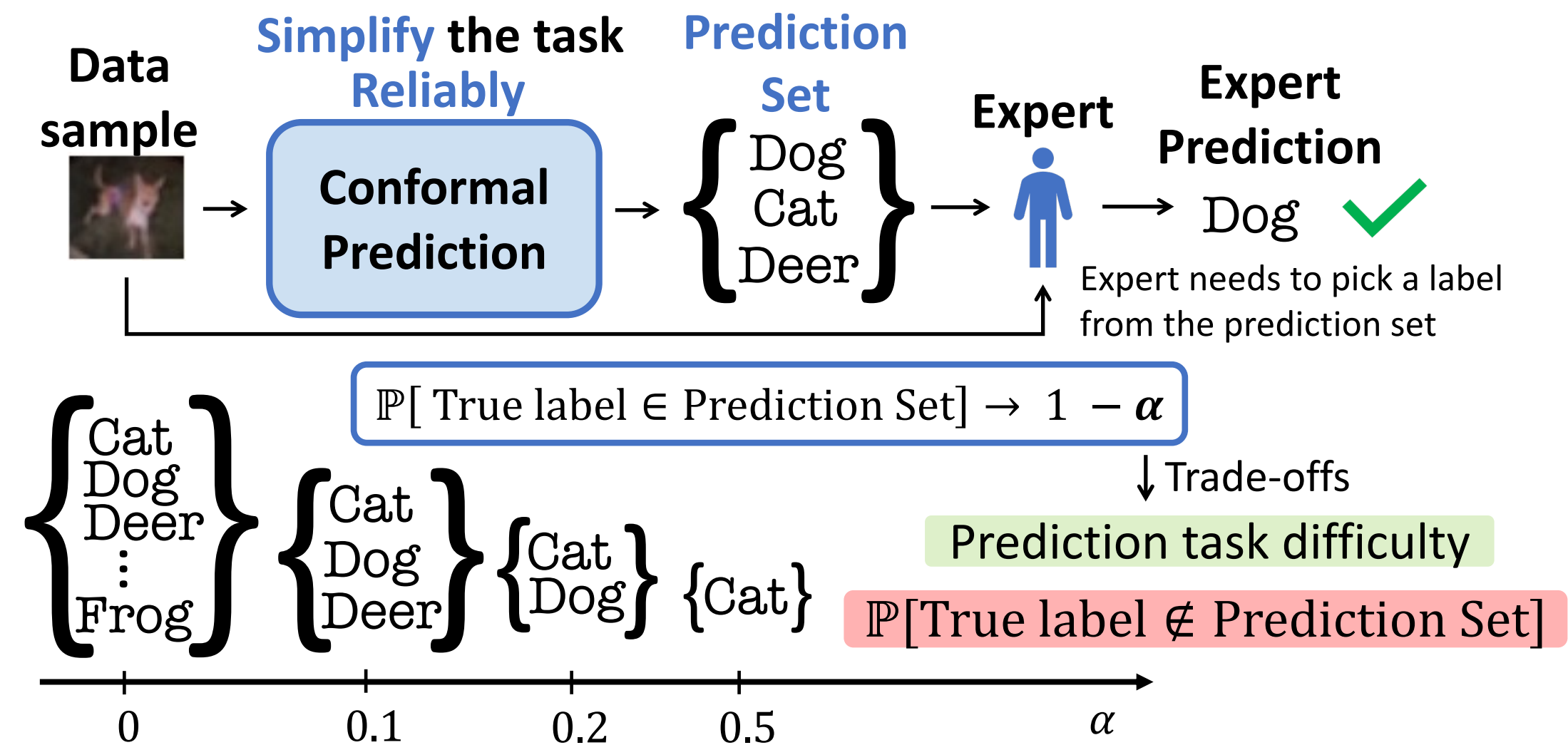
Existing decision support systems for multi-class classification tasks



- ⇒ Experts **have to decide when** to trust the system.
- ⇒ Empirical findings [1,2] have reached **no conclusion** on how experts can avoid developing a **misplaced** trust on the system.

Goal: Design a decision support system that does not require humans to decide **when** to trust its predictions

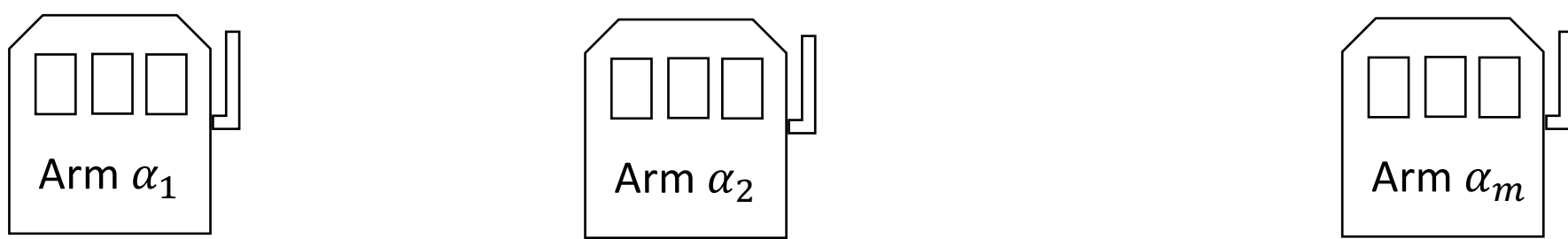
Our decision support system for multi-class classification tasks



Optimizing conformal predictors with bandits

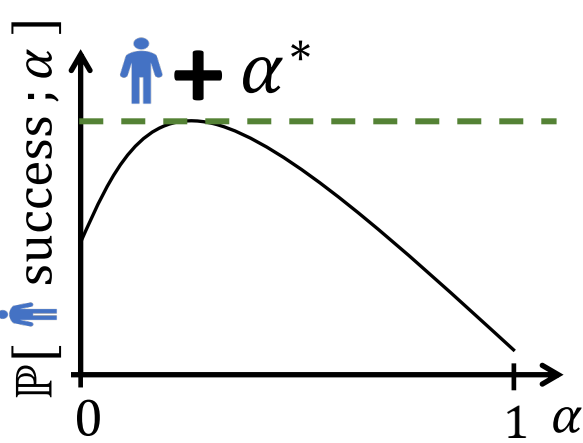
Arms $\Leftrightarrow \alpha$ values

$$\mathbb{P}[\text{success} ; \alpha_1] \quad \mathbb{P}[\text{success} ; \alpha_2] \quad \dots \quad \mathbb{P}[\text{success} ; \alpha_m]$$



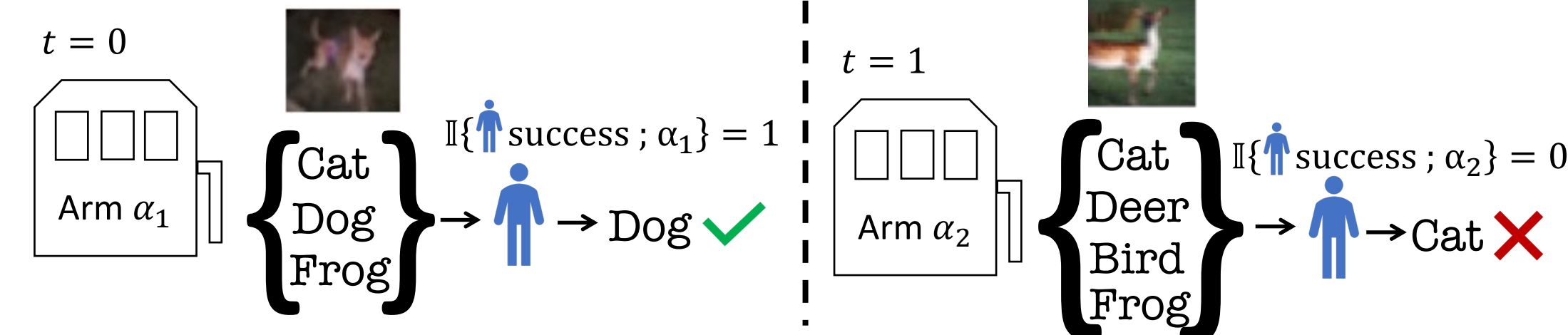
Discover optimal arm

$$\alpha^* = \arg\max_{\alpha} \mathbb{P}[\text{success} ; \alpha]$$

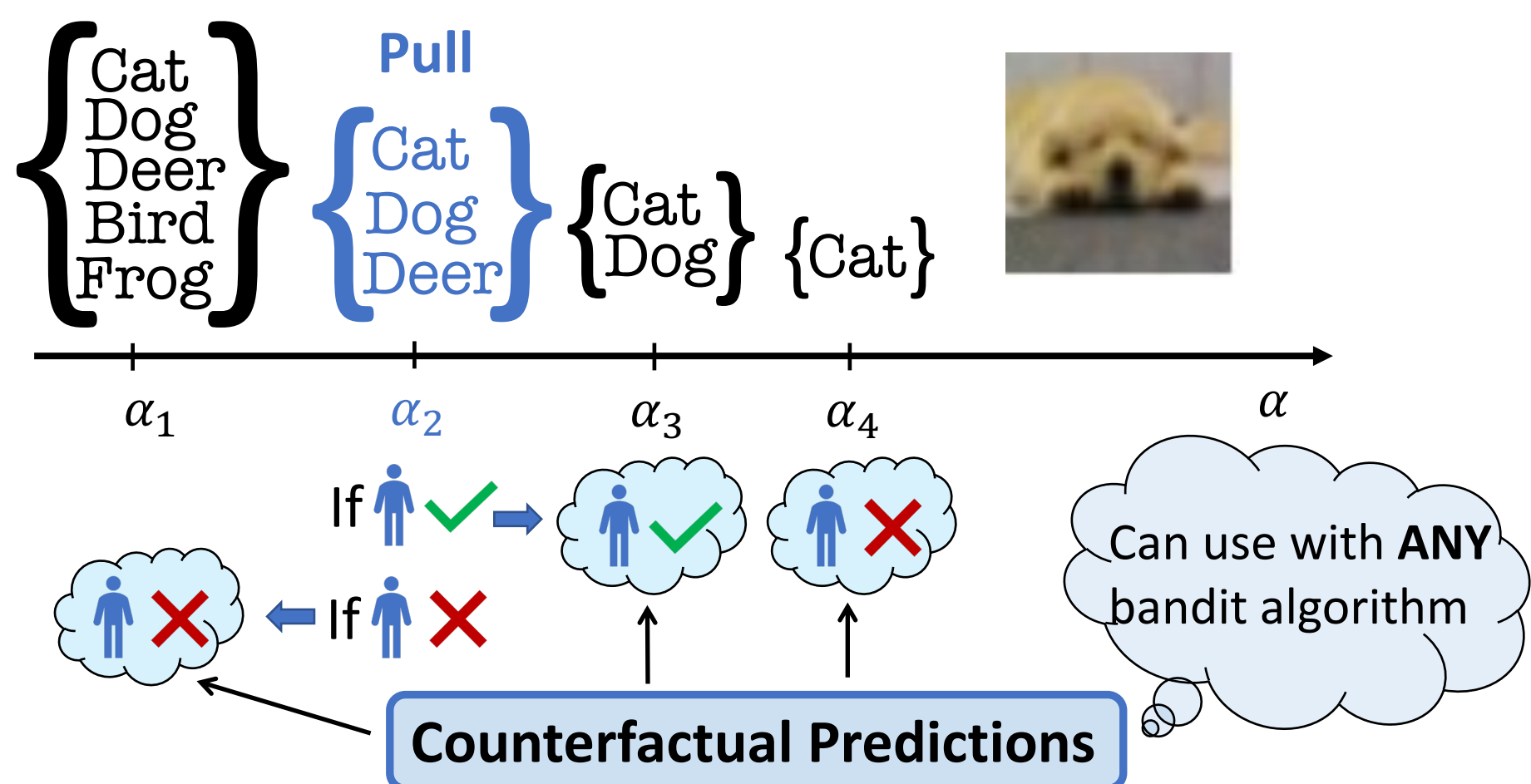


Successive Elimination

Pull arms sequentially



Counterfactual Monotonicity Assumption



Counterfactual Successive Elimination finds α^* exponentially faster

Large-scale human subject study

- ⇒ 194,407 predictions
- ⇒ 2,751 human subjects
- ⇒ 19,200 different pairs of natural images and sets

Example

Which one of the following categories fits better the image below?



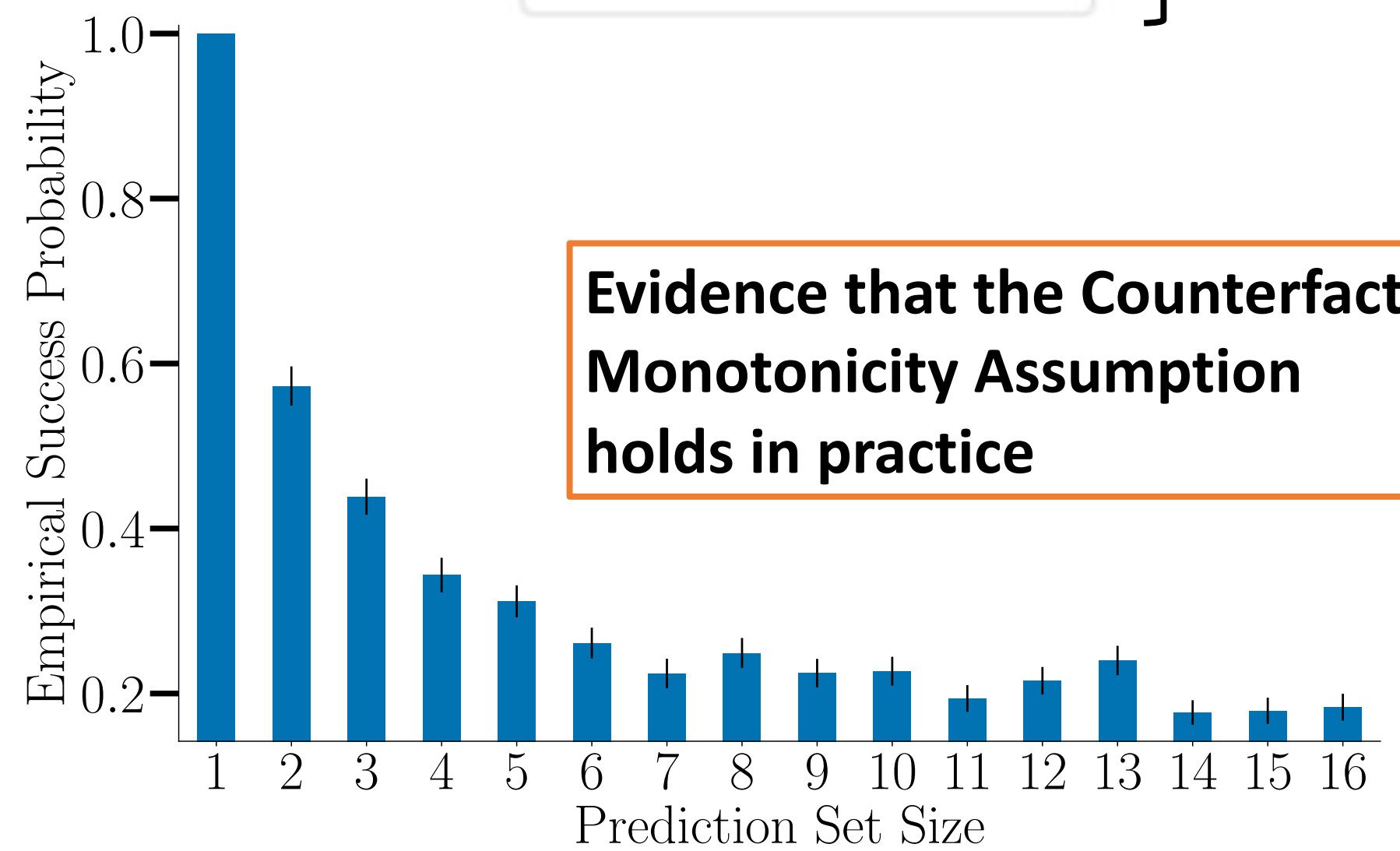
- Car
- Airplane
- Truck

Strict implementation
picks only from the prediction set

- Car
- Airplane
- Truck
- Other

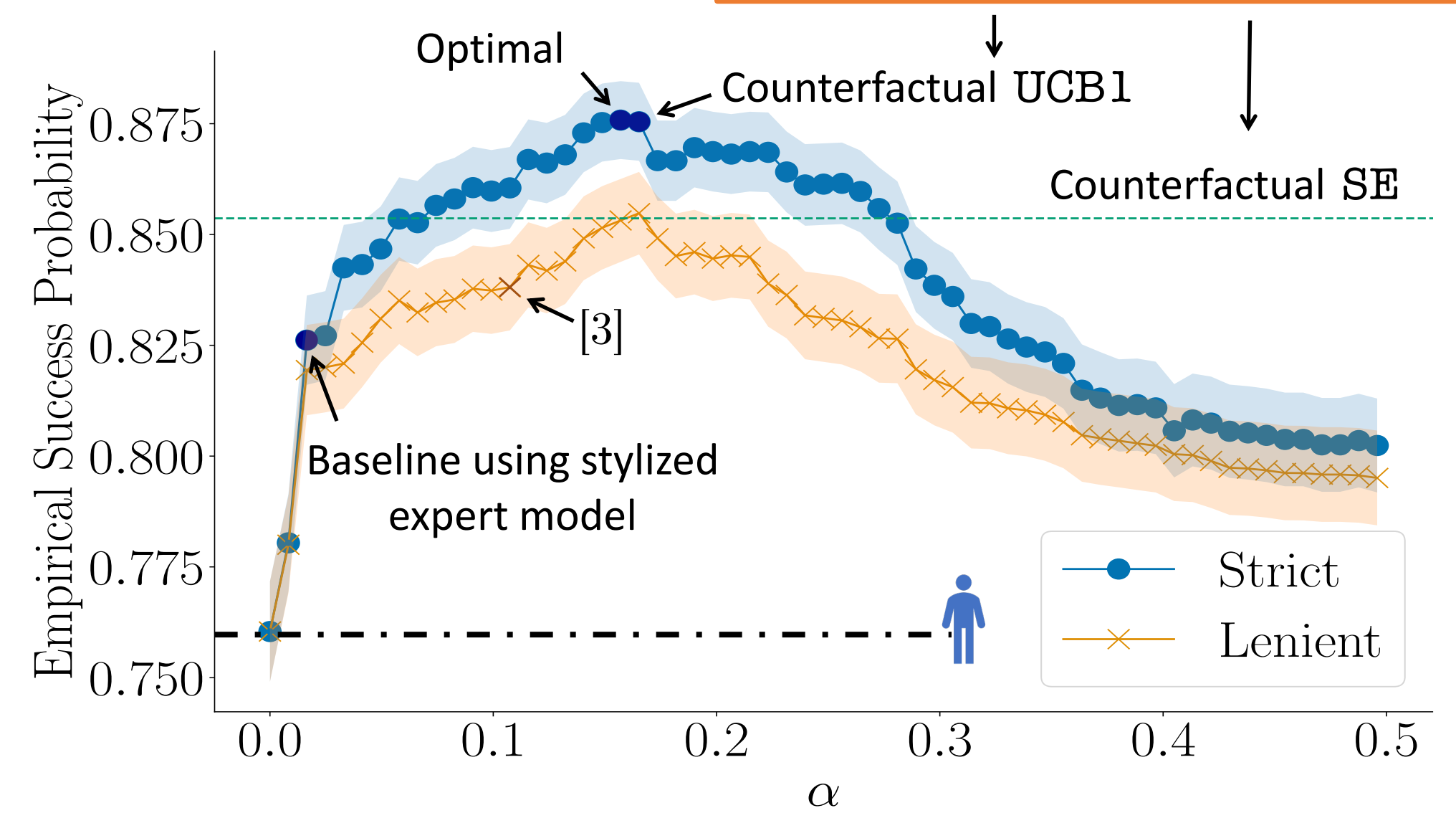
Lenient implementation
can pick any label

If you chose 'Other' above, please choose a category:



Evidence that the Counterfactual Monotonicity Assumption holds in practice

Bandit algorithms using Counterfactual Monotonicity



Lenient implementation

Disadvantage over strict

Number of predictions in which the expert **misplaces** their trust on the system

Advantage over strict

Number of predictions in which the expert predicts correctly from outside the prediction set

