

Yan Scholten<sup>1</sup> Jan Schuchardt<sup>1</sup> Aleksandar Bojchevski<sup>2</sup> Stephan Günnemann<sup>1</sup>  
<sup>1</sup>Technical University of Munich <sup>2</sup>University of Cologne

**tl;dr: Novel robustness certificates for decomposable data**

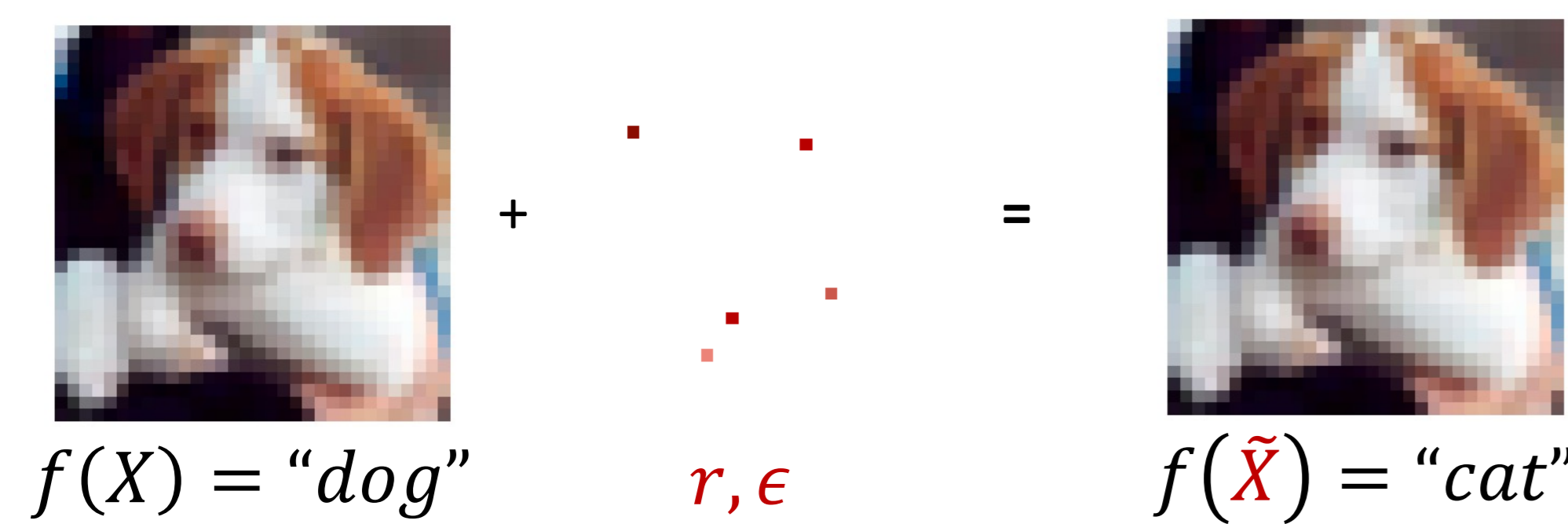
- Adversaries can perturb a subset of all entities of an object (e.g. pixels of an image, nodes of a graph)
- We propose a highly flexible certification framework for continuous and discrete domains
- Superior robustness-accuracy trade-offs under our threat model

## Context

- Machine learning models are susceptible to adversarial perturbations
- Robustness certificates provide provable robustness guarantees

## Problem

Certifying robustness on decomposable data (e.g. images, graphs, ...) is challenging when adversarial perturbations are bounded by both: (1) the number of perturbed entities  $r$ , and (2) perturbation strength  $\epsilon$

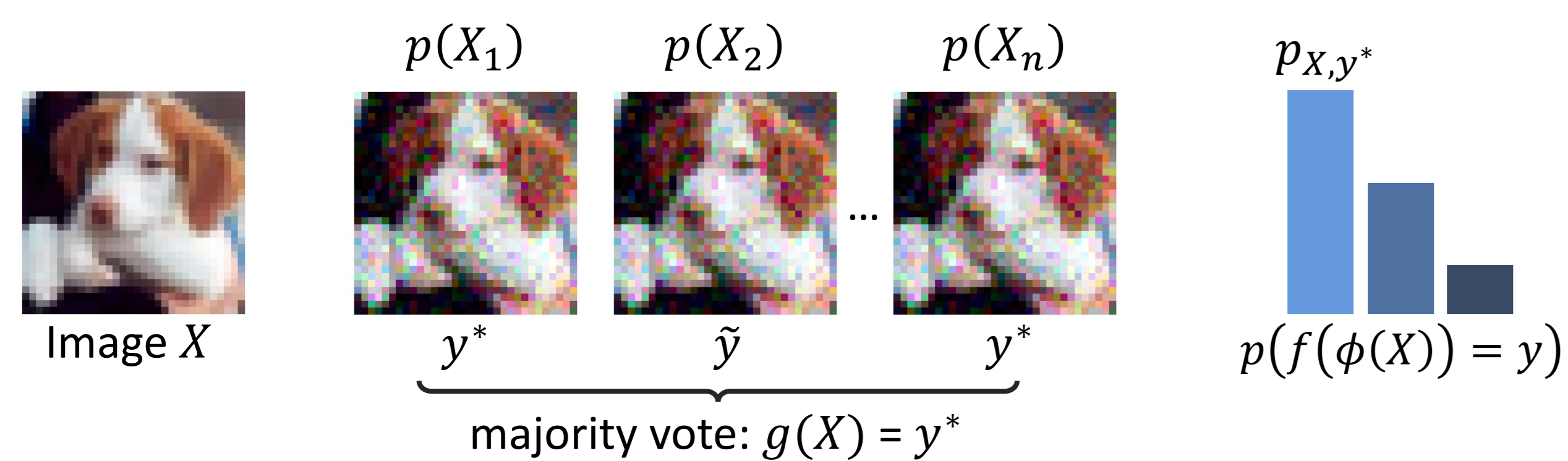


How can we guarantee robustness under such adversarial perturbations?

Existing approaches sacrifice robustness over accuracy or vice versa

## Background: Randomized smoothing

- Sample smoothed images  $X_i \sim \phi(X)$  from smoothing distribution  $\phi$
- Classify them with base classifier  $f$  and certify the majority vote



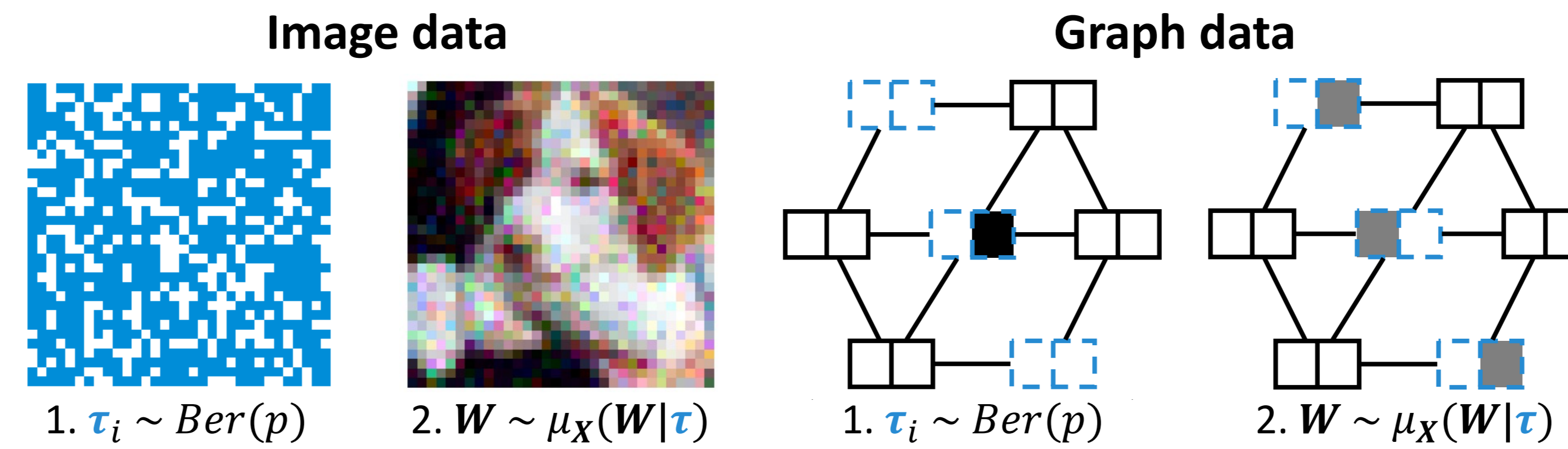
### How to certify robustness under randomized smoothing?

- Derive lower bound  $\underline{p}_{\tilde{X}, y^*}(p_{X, y^*})$  on probability  $p_{\tilde{X}, y^*}$  to classify  $\tilde{X}$  as  $y^*$
- Smoothed classifier  $g$  is certifiably robust if

$$\underline{p}_{\tilde{X}, y^*}(p_{X, y^*}) > 0.5 \text{ for any perturbed } \tilde{X} \in \mathcal{B}(X)$$

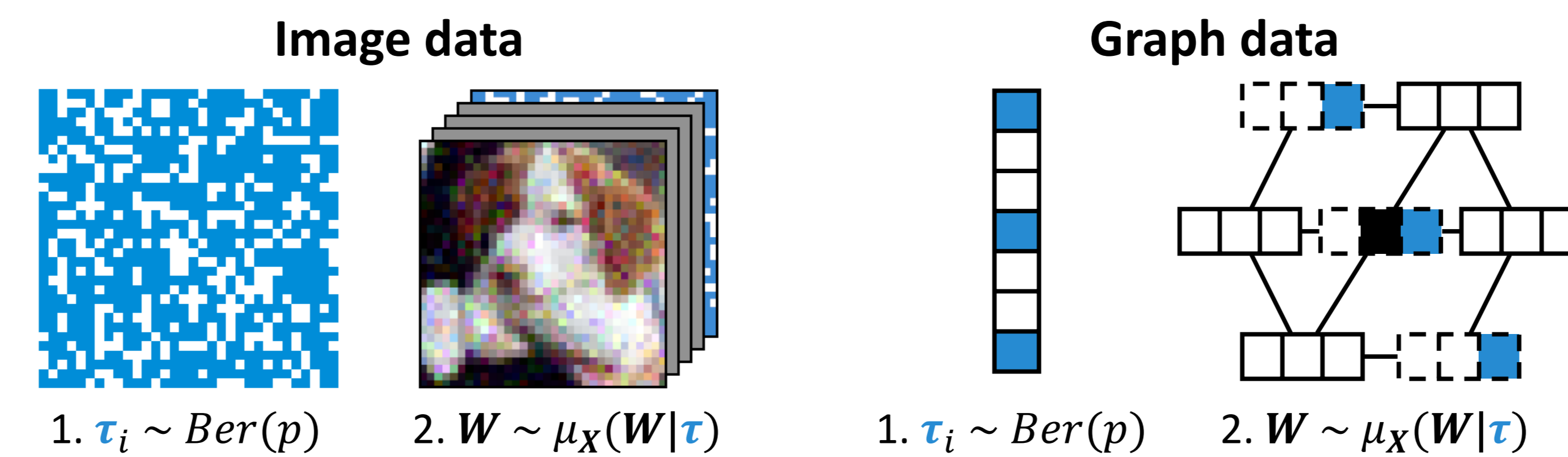
## Hierarchical smoothing distribution

1. **Upper-level smoothing:** Sample indicator  $\tau_i \sim \text{Ber}(p)$  with probability  $p$
2. **Lower-level smoothing  $\mu$ :** Sample additive noise for indicated entities only



## How to certify robustness?

- Append indicator  $\tau$  to the object  $X$
- Construct a new base classifier  $f$  operating on this higher-dimensional space

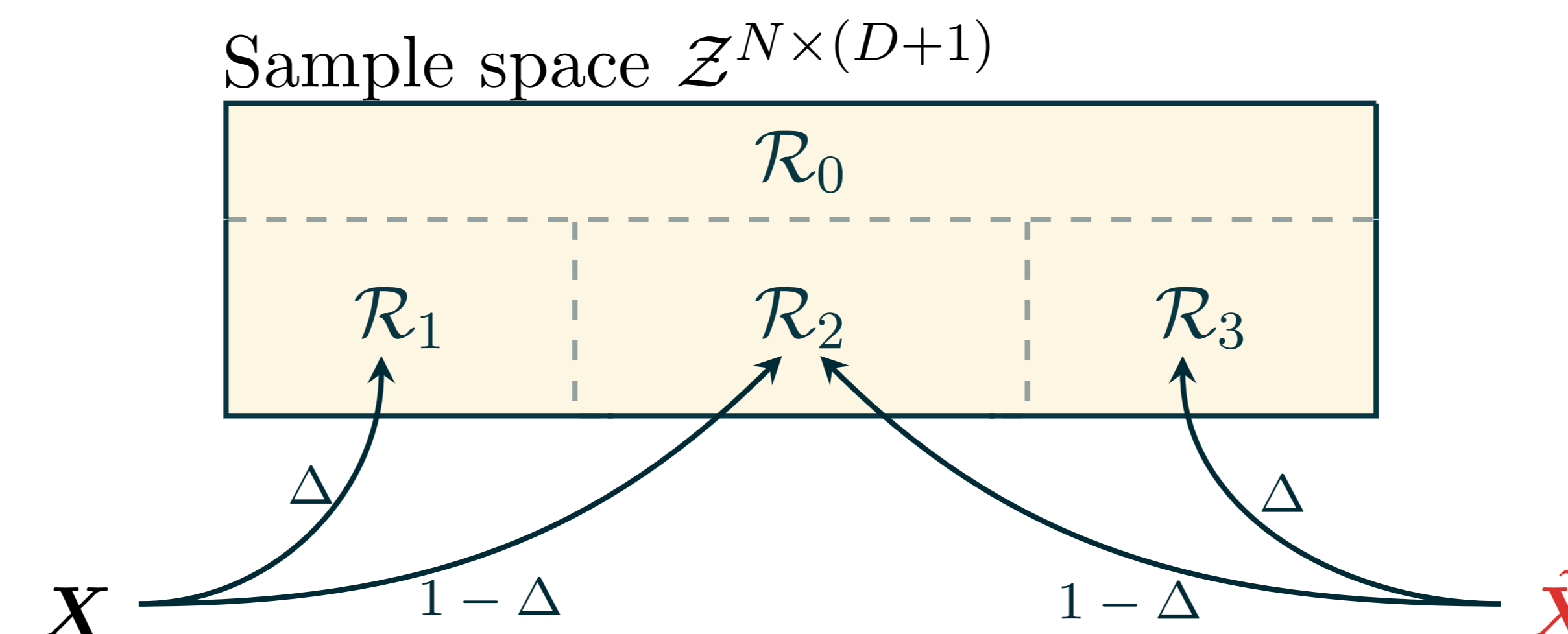


### Certify robustness by reusing existing bounds

1. Compute constant  $\Delta = 1 - p^r$  for selection probability  $p$  and radius  $r$
2. Plug  $\Delta$  into existing lower bound  $\underline{p}_{\tilde{X}, y^*}(p_{X, y^*})$  for the lower-level distribution

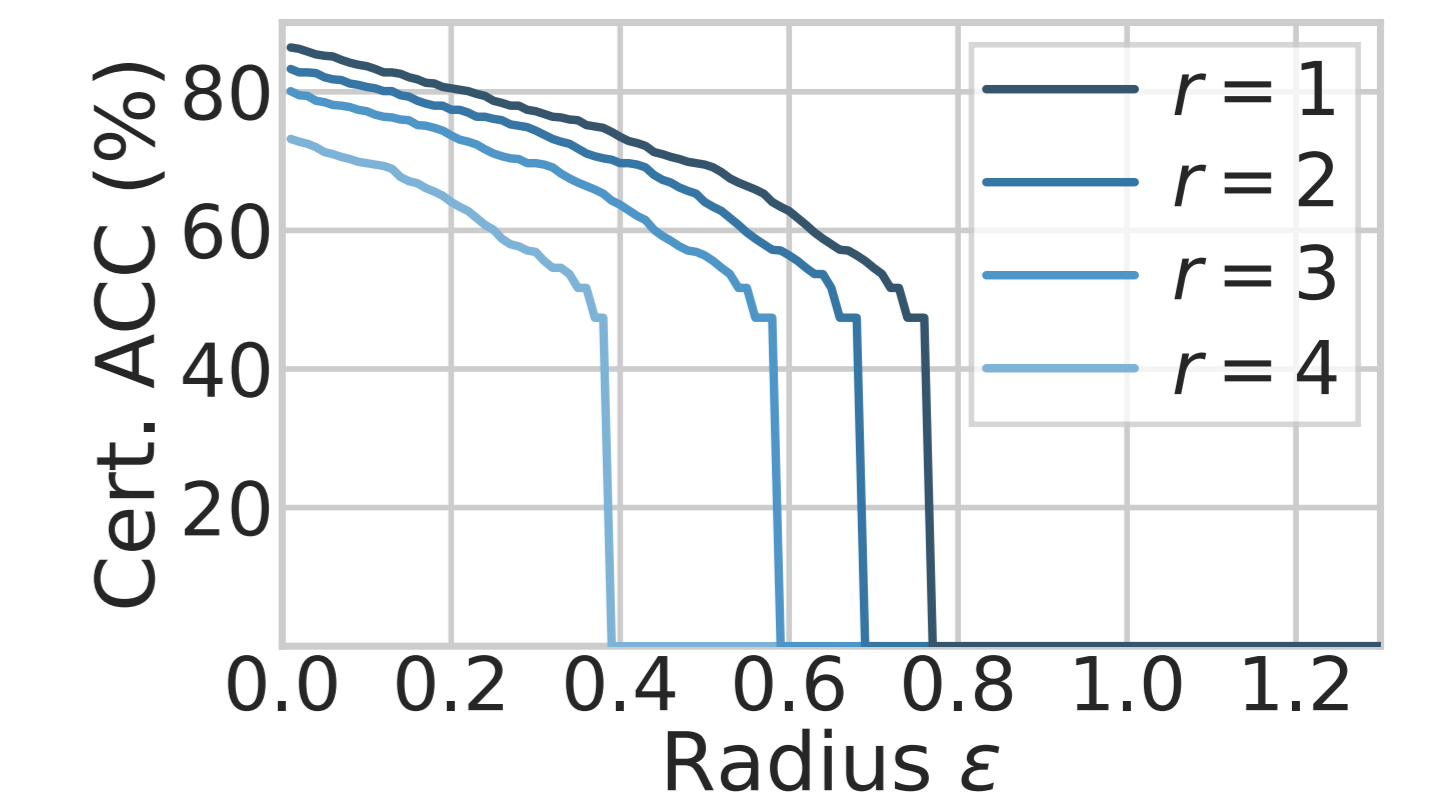
## Why can we integrate robustness guarantees for the lower-level smoothing distribution?

- Partition the sample space into disjoint regions  $R_0, R_1, R_2, R_3$
- Supports  $S_X$  and  $S_{\tilde{X}}$  for hierarchical smoothing around  $X$  and  $\tilde{X}$  intersect only for samples where all perturbed entities are selected by  $\tau$  (Region  $R_2$ )
- This allows the certificate to separate clean from perturbed entities



## Experimental evaluation

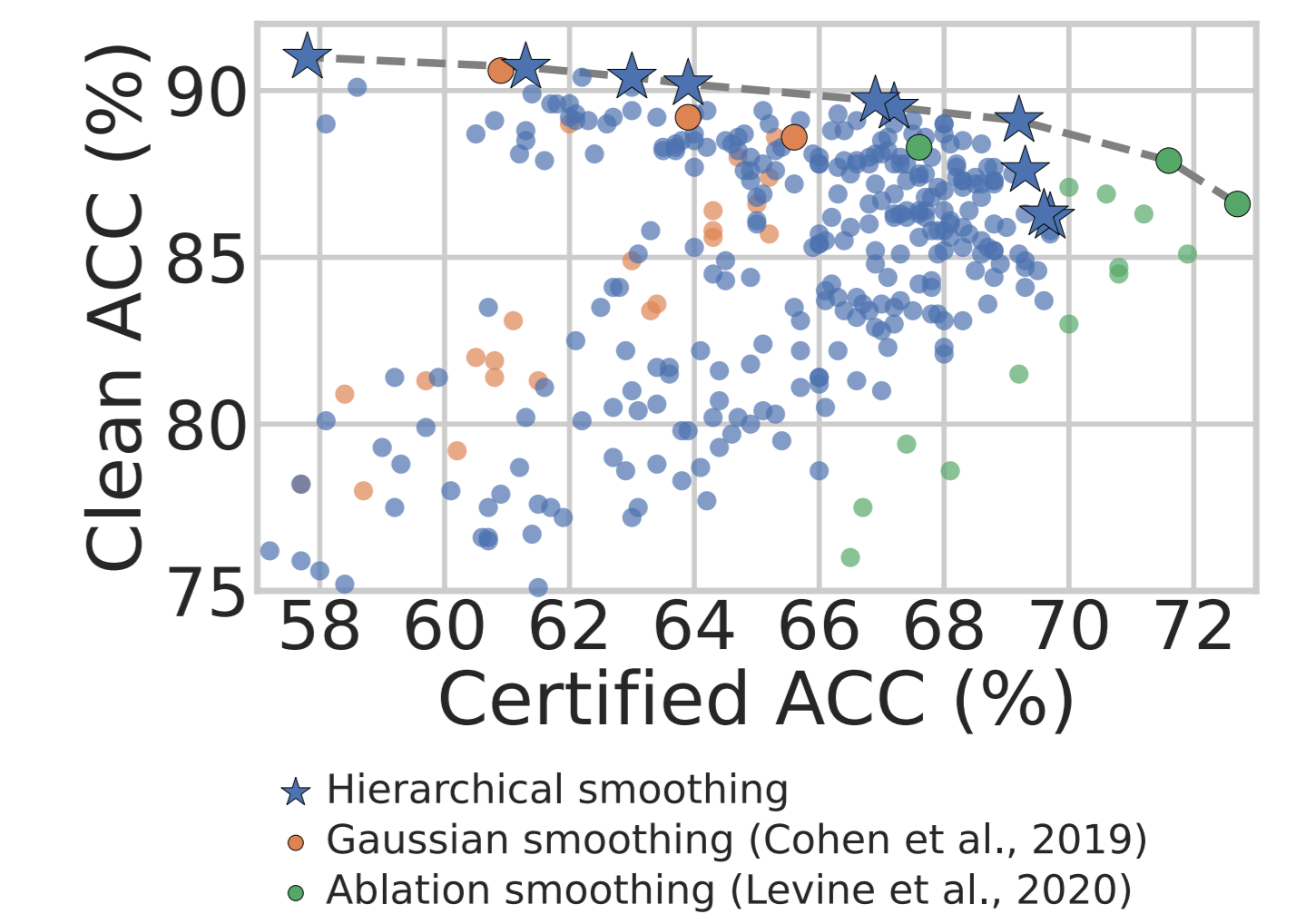
Certificate captures robustness w.r.t. both radii  $r$  and  $\epsilon$



We expand the Pareto-front w.r.t. robustness and accuracy

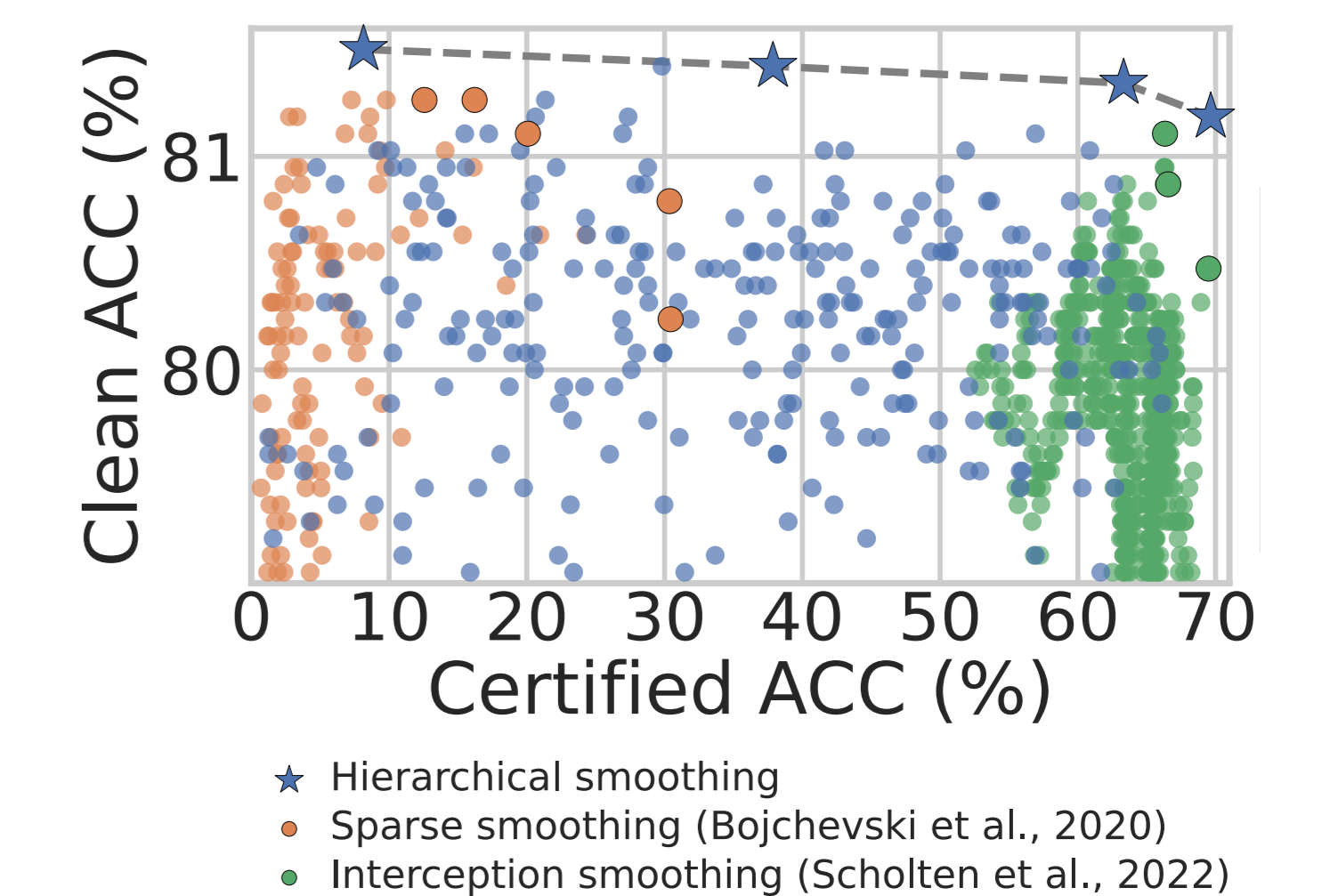
### Image classification

- Perturbation strength bounded under  $\ell_2$ -norm ( $r = 3, \epsilon = 0.35$ )
- Hierarchical smoothing with Gaussian smoothing
- Smoothed ResNet50 on CIFAR10



### Node classification

- Sparse attribute perturbations ( $r = 1, r_a = 0, r_d = 40$ )
- Hierarchical smoothing with sparse smoothing
- Smoothed GAT on CoraML



## Paper, code, and more

