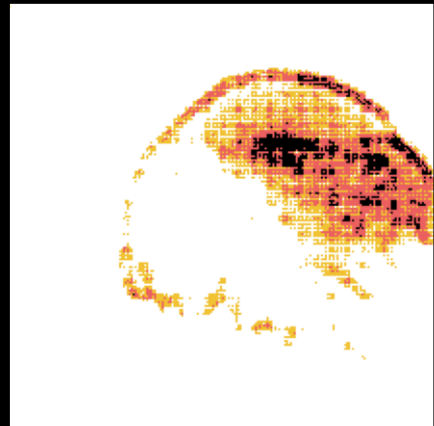


# The First Provably Robust and Interpretable Pixel-Level Explanations



Colored pixels are provably important for predicting "eagle" under bounded input noise



## Pixel-Level Certified Explanations via Randomized Smoothing

Alaa Anani<sup>1,2</sup> Tobias Lorenz<sup>2</sup> Mario Fritz<sup>1</sup> Bernt Schiele<sup>2</sup>

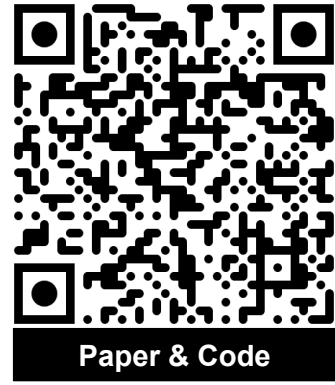


<sup>1</sup> MAX PLANCK INSTITUTE FOR INFORMATICS



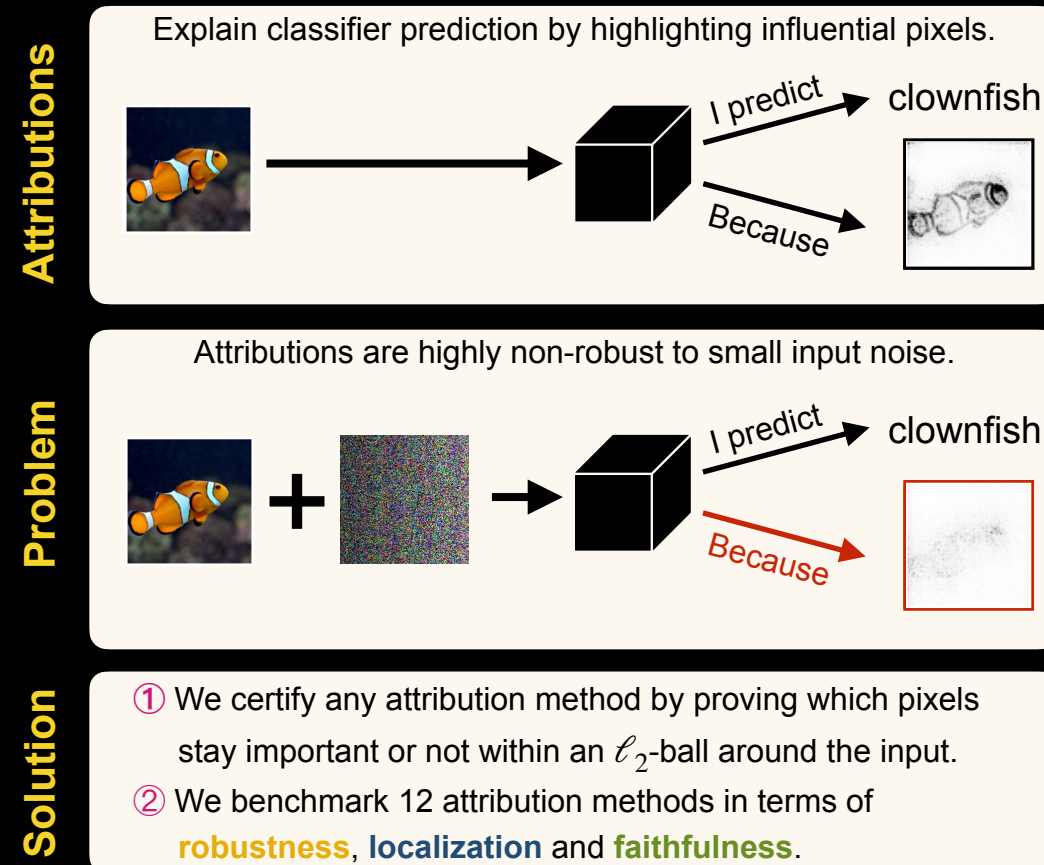
<sup>2</sup> CISPA HELMHOLTZ CENTER FOR INFORMATION SECURITY

SIC Saarland Informatics Campus



Paper & Code

### 1. Take-Home Message



### 2. Randomized Smoothing

Proving model output stability by evaluating it on an input distribution [1,2].

$$f(x) \rightarrow f\left(\bigotimes_{\mathcal{N}(x, \sigma^2 I)} \tilde{x}\right)$$

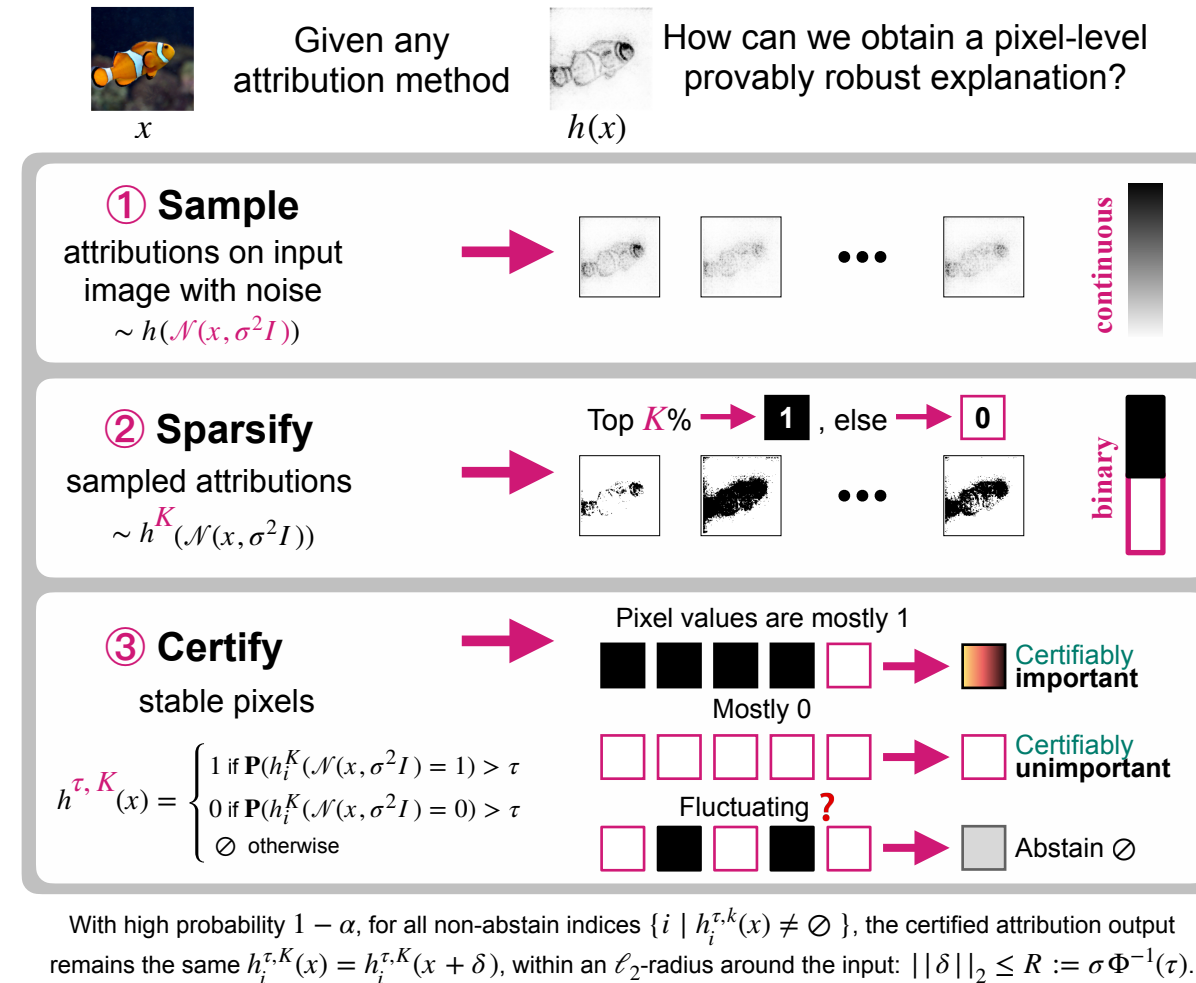
**Certified segmentation** [2]

Given a segmentation model  $f: \mathbb{R}^{c \cdot N} \rightarrow \mathcal{Y}^N$ , threshold  $\tau \in [\frac{1}{2}, 1)$ , and sampling error rate  $\alpha$ , the smoothed (certified) version is defined as

$$g_i^\tau(x) = \begin{cases} c_{A,i} & \text{if } \mathbf{P}(f_i(\mathcal{N}(x, \sigma^2 I)) = c_{A,i}) > \tau \\ \emptyset, & \text{otherwise} \end{cases}$$

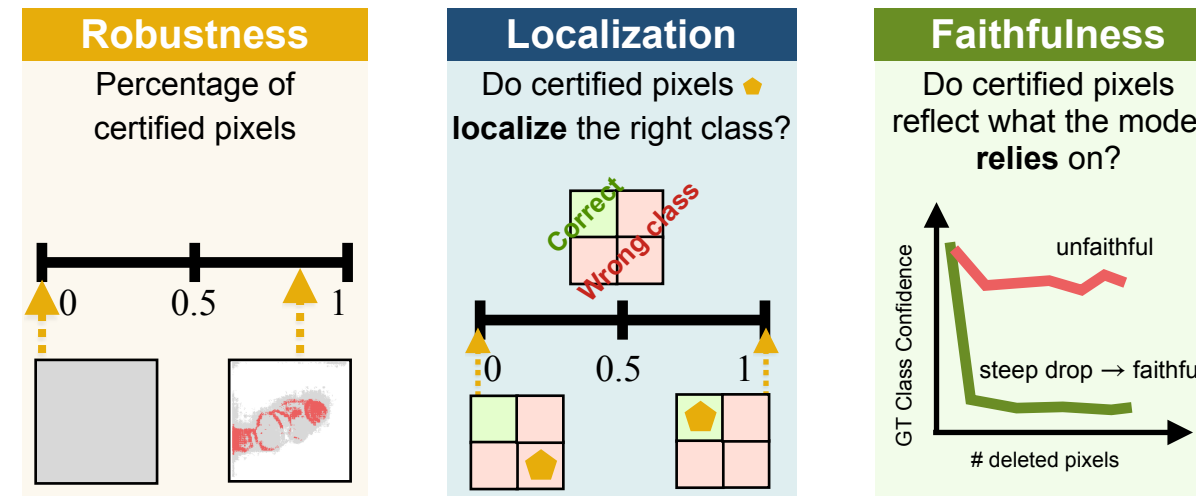
With probability  $1 - \alpha$ , all non-abstain indices  $\{i \mid g_i^\tau(x) \neq \emptyset\}$  remain the same  $g_i^\tau(x) = g_i^\tau(x + \delta)$ , for  $\delta \in \mathbb{R}^{m \cdot N}$  with  $\|\delta\|_2 \leq R := \sigma \Phi^{-1}(\tau)$ .

### 3. Method

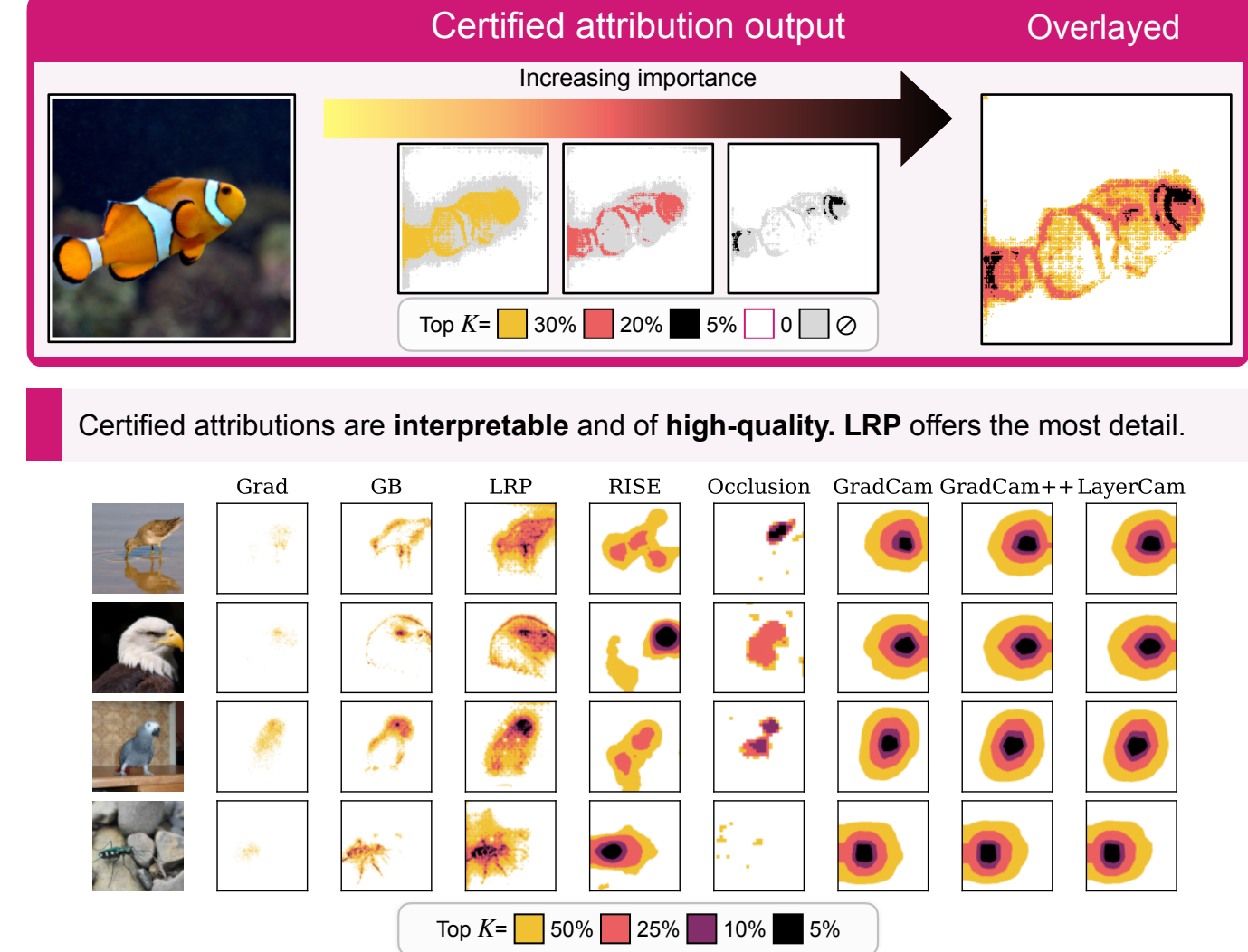


### 4. Evaluation Setup

We evaluate 12 attribution methods from 3 families (gradient, perturbation, activation) on 5 models using 3 novel metrics:



### 5. Results



**Certified Robustness**

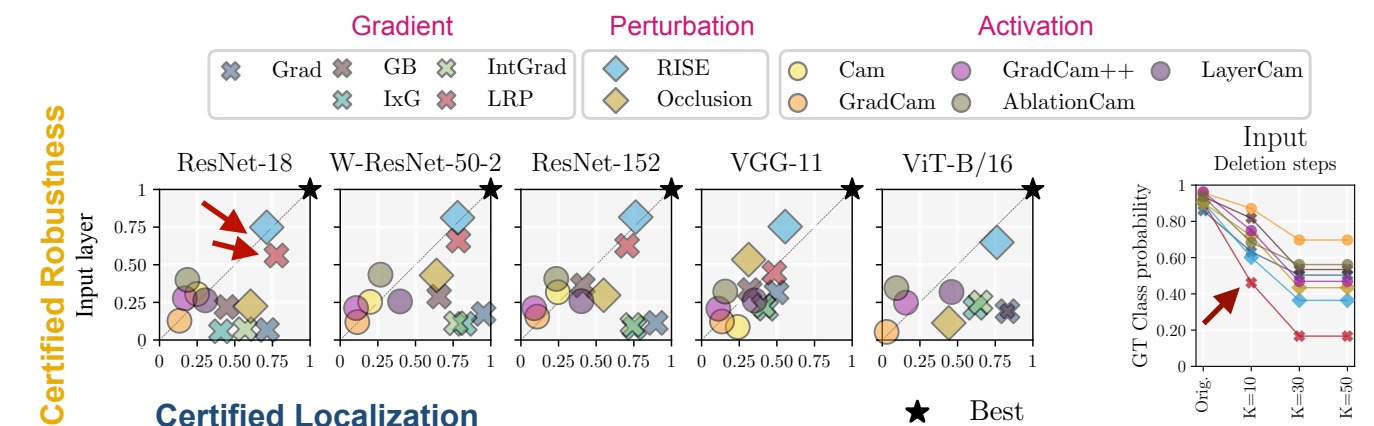
LRP and RISE are the most robust.

**Certified Localization**

LRP, RISE and Occlusion localize best.

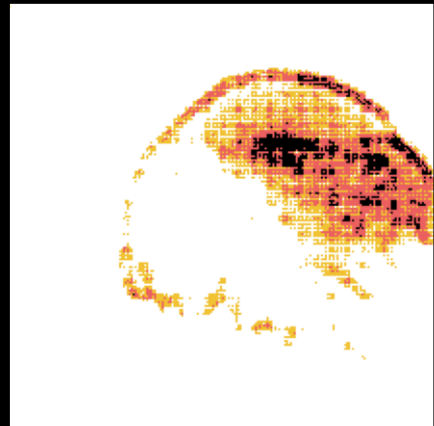
**Faithfulness**

LRP and RISE are the most faithful.





# The First Provably Robust and Interpretable Pixel-Level Explanations



Colored pixels are provably important for predicting "eagle" under bounded input noise



## Pixel-Level Certified Explanations via Randomized Smoothing

Alaa Anani<sup>1,2</sup> Tobias Lorenz<sup>2</sup> Mario Fritz<sup>2</sup> Bernt Schiele<sup>1</sup>



<sup>1</sup> MAX PLANCK INSTITUTE FOR INFORMATICS



<sup>2</sup> CISPA  
HELMHOLTZ CENTER FOR INFORMATION SECURITY

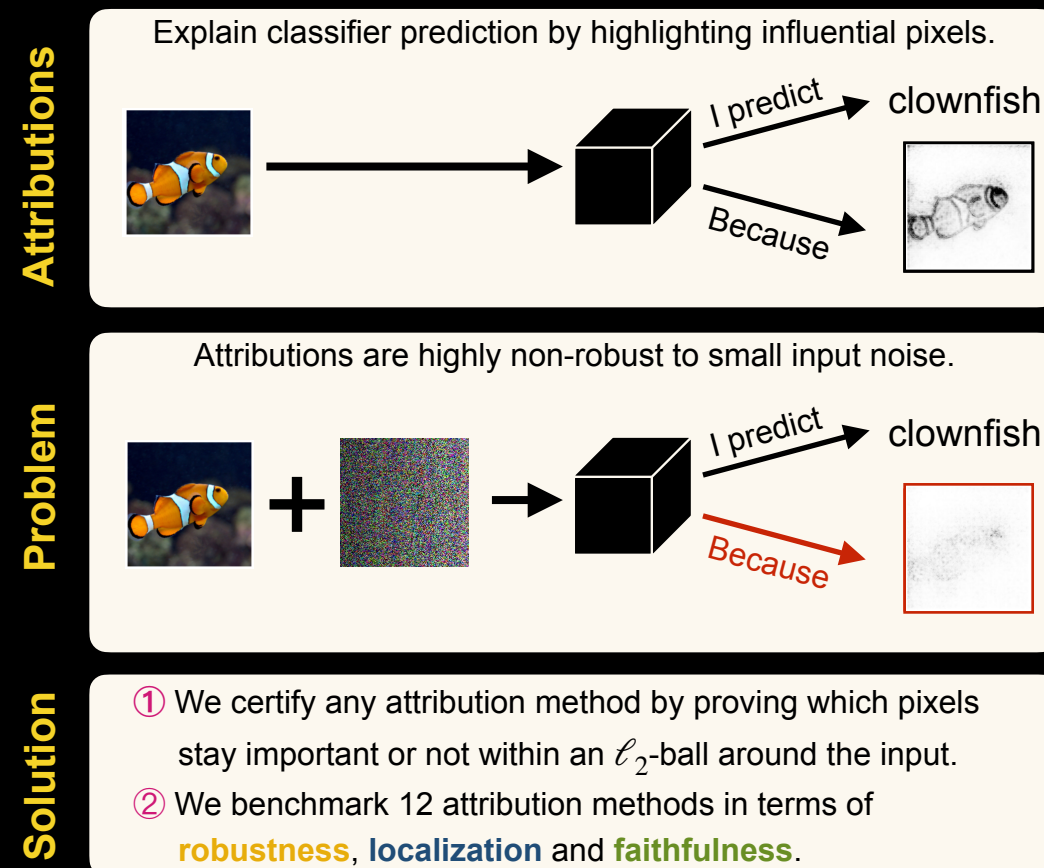


SIC Saarland Informatics Campus



Paper & Code

### 1. Take-Home Message



### 2. Randomized Smoothing

Proving model output stability by evaluating it on an input distribution [1,2].

$$f(x) \rightarrow f\left(\bigcirc_{\mathcal{N}(x, \sigma^2 I)} \dot{x}\right)$$

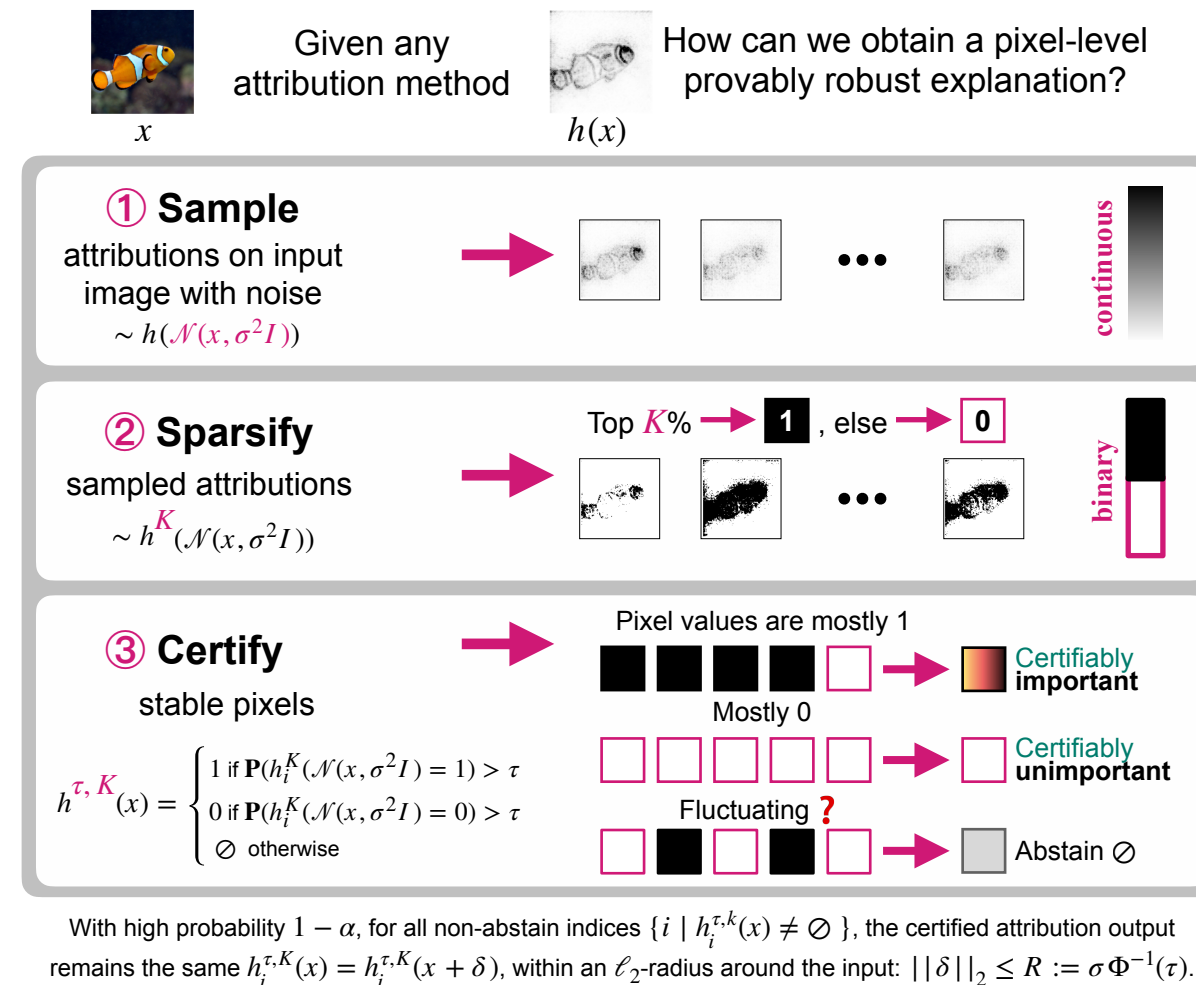
**Certified segmentation** [2]

Given a segmentation model  $f: \mathbb{R}^{c \cdot N} \rightarrow \mathcal{Y}^N$ , threshold  $\tau \in [\frac{1}{2}, 1)$ , and sampling error rate  $\alpha$ , the smoothed (certified) version is defined as

$$g_i^\tau(x) = \begin{cases} c_{A,i} & \text{if } \mathbf{P}(f_i(\mathcal{N}(x, \sigma^2 I)) = c_{A,i}) > \tau \\ \emptyset, & \text{otherwise} \end{cases}$$

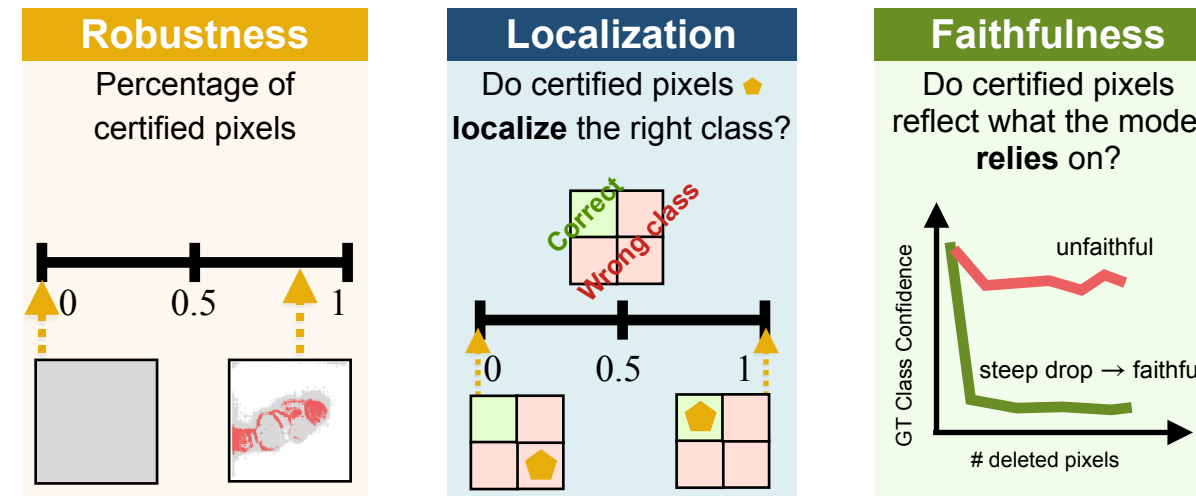
With probability  $1 - \alpha$ , all non-abstain indices  $\{i \mid g_i^\tau(x) \neq \emptyset\}$  remain the same  $g_i^\tau(x) = g_i^\tau(x + \delta)$ , for  $\delta \in \mathbb{R}^{m \cdot N}$  with  $\|\delta\|_2 \leq R := \sigma \Phi^{-1}(\tau)$ .

### 3. Method

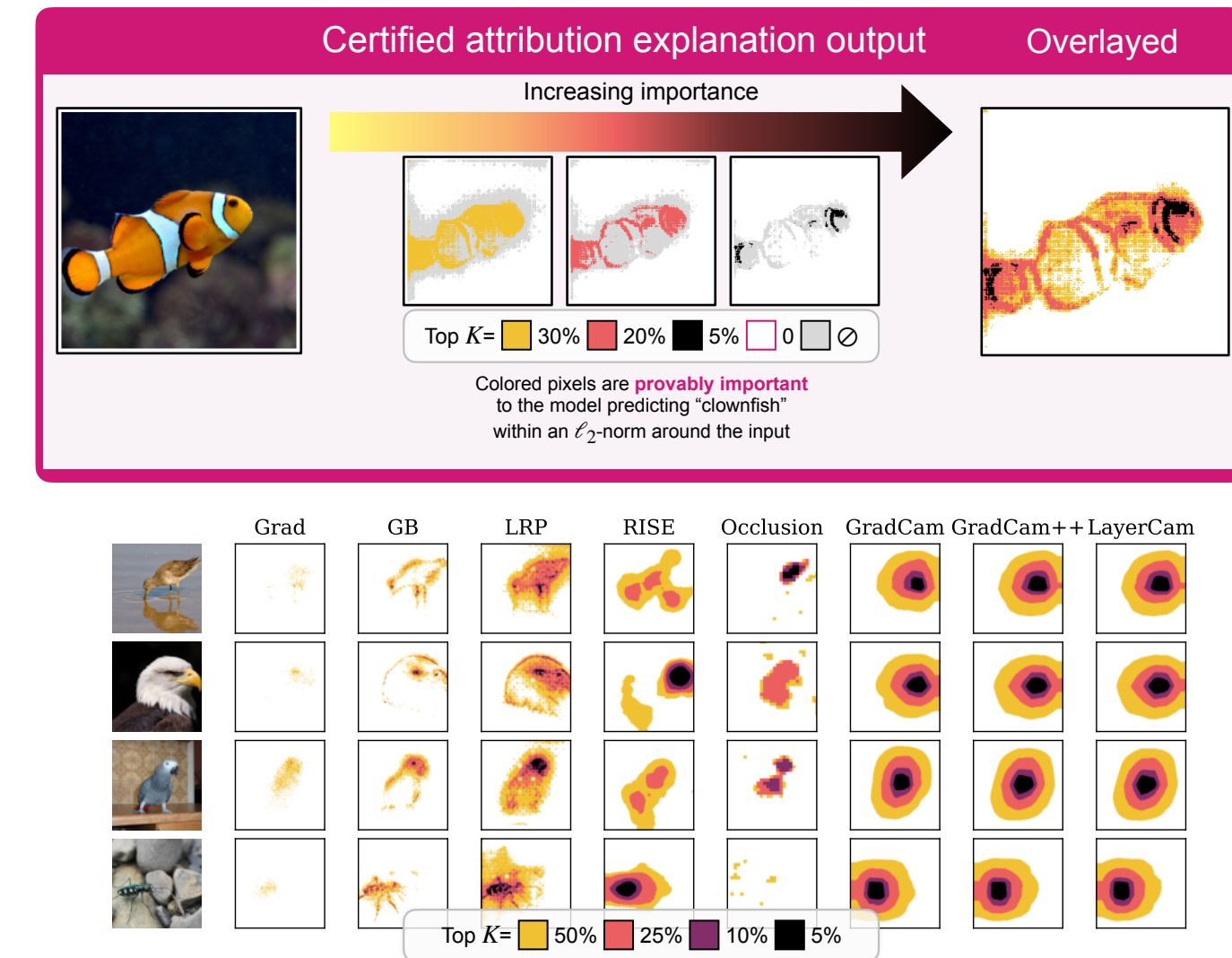


### 4. Evaluation Setup

We evaluate 12 attribution methods from 3 families (gradient, perturbation, activation) on 5 models using 3 novel metrics:



### 5. Results



Certified Robustness	Certified Localization	Faithfulness
LRP and RISE are the most robust.	LRP, RISE and Occlusion localize best.	LRP and RISE are the most faithful.

