# SoK
## Certified Robustness for Deep Neural Networks
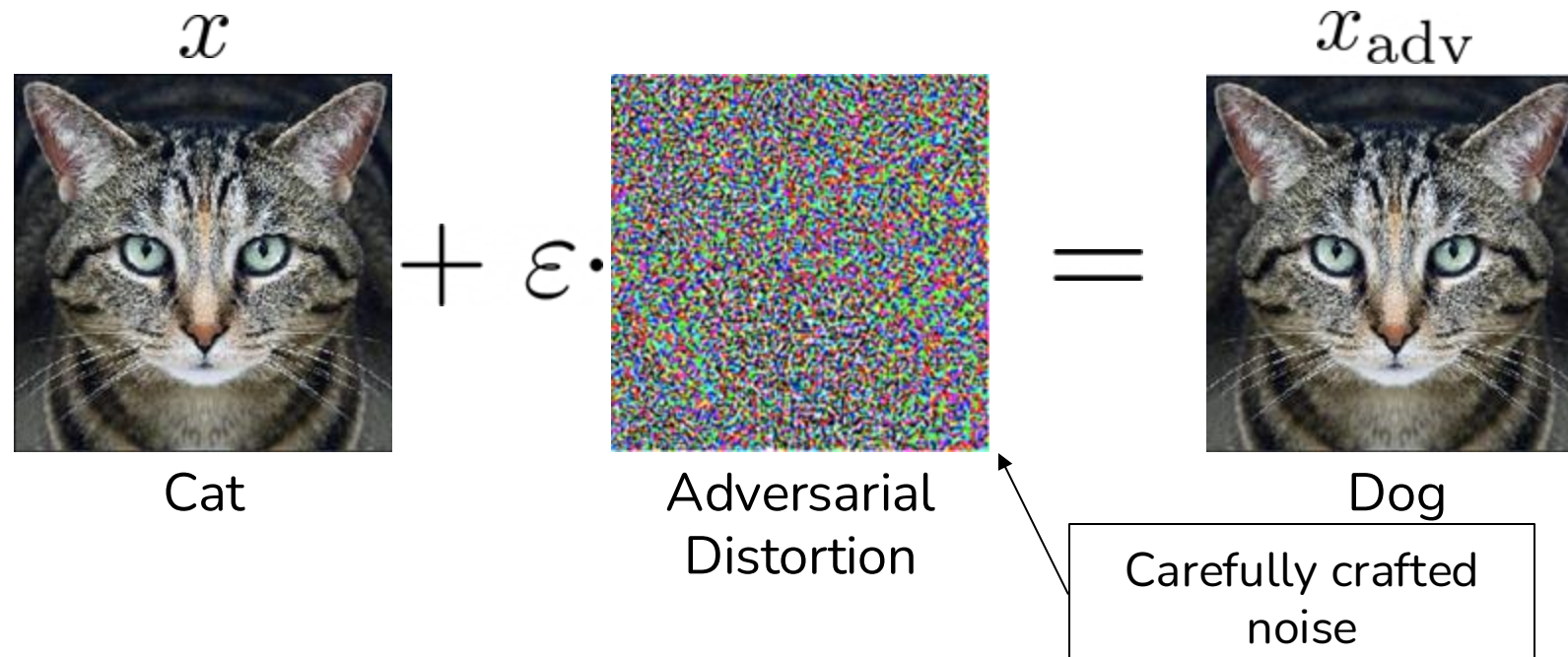
**Linyi Li** (UIUC), Tao Xie (Peking University), Bo Li (UIUC)

UIUC Secure Learning Lab

# Adversarial Robustness – A Lasting Threat

- Deep neural networks (DNNs) can be easily fooled by adversarial examples
    - **Tiny** crafted perturbations can make DNNs give **wrong** predictions

$$x$$



$$+ \; \varepsilon \cdot$$

$$=$$

$$x_{\mathrm{adv}}$$

Cat

Adversarial Distortion

Dog

Carefully crafted noise

# Severe Safety Threats

Example: autonomous driving



Possibly Fatal Accident

# Arm Race

Defenses bypassed by follow-up attacks

discovery
[Szegedy et al. '14]

# Ending Arm Race? Certified Robustness!

- <u>Prove</u> adversarial example doesn't exist

  **Guarantee of Safety**

# Ending Arm Race? Certified Robustness!

- Prove adversarial example doesn't exist

  Guarantee of Safety

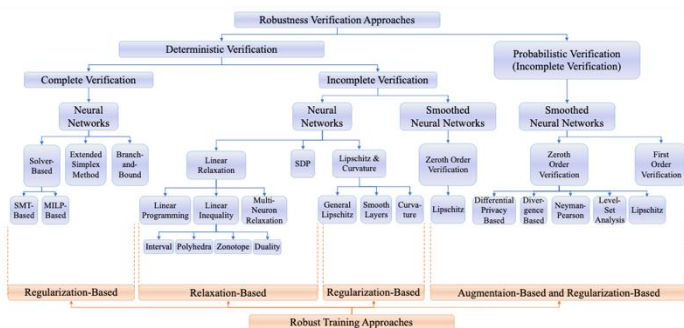**First Systematization of Knowledge on Certified Robustness for Deep Neural Networks!**

# Content

## TAXONOMY



- Characteristics
- Strengths
- Limitations
- Connections
- Generalization
- …

## SUMMARY

## DISCUSSION

- Current Research
- Theoretical Barriers
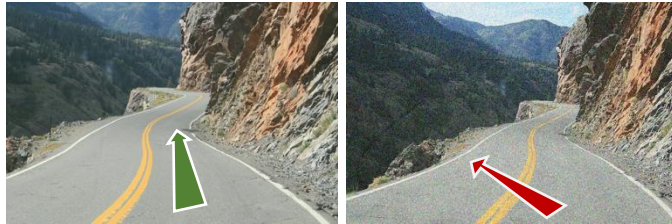- Main Challenges
- Future Directions
- …

## BENCHMARK

**VeriGauge**
*Open-source platform for 20+ approaches*

# Threat Model

- Various types of adversarial examples exist



Tiny Perturbations
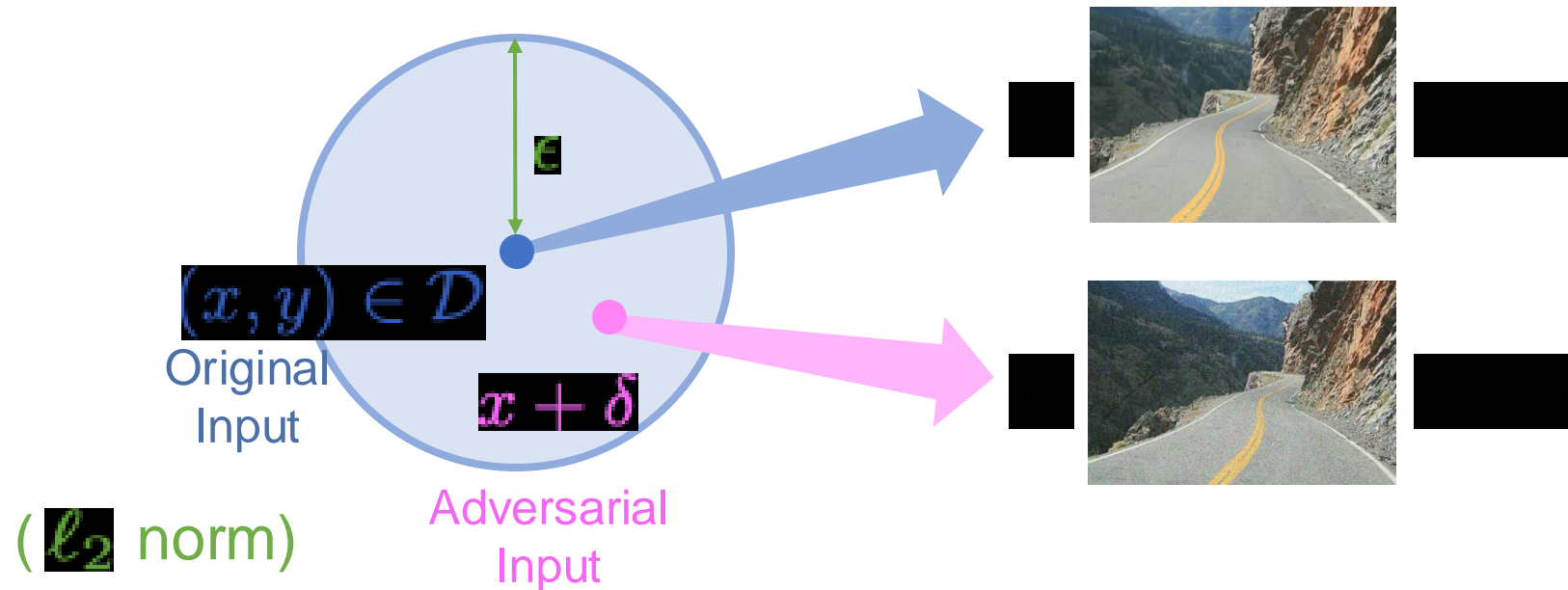


Semantic
Transformations

- Focus on (p-norm constrained) perturbations
  - Widely studied
  - Techniques generalizable to other types of adversarial examples

# Robustness against (p-Norm Constrained) Perturbations

Given a DL model ■, finite test dataset ████████████████

Check: $\forall \delta, \|\delta\|_p \leq \epsilon, f(x+\delta) = y$



$\epsilon$

$(x, y) \in \mathcal{D}$

Original Input

$x + \delta$

$(\ell_2$ norm)

Adversarial Input

# Formal Definition of Certification

For given system █ and data instance ■ with true label █, compute larger $r$, such that

$$\blacksquare\ r\ \blacksquare$$



Larger certified radius

= Tighter certification

= Better certified robustness

# Taxonomy of Verification Method

**Taxonomy Criteria:**

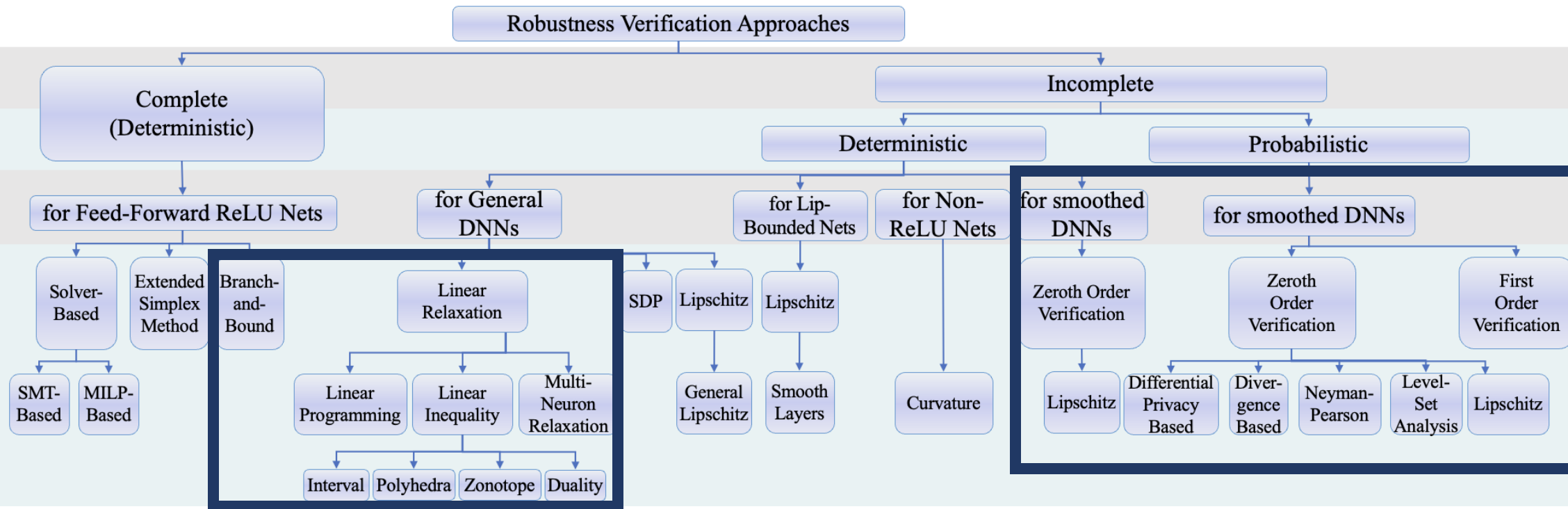Robustness Verification Approaches

**Complete/Incomplete**

Complete (Deterministic)

Incomplete

**Deterministic/Probabilistic**

Deterministic

Probabilistic

**System Model**

for Feed-Forward ReLU Nets

for General DNNs

for Lip-Bounded Nets

for Non-ReLU Nets

for smoothed DNNs

for smoothed DNNs

**Core Methodology**

Solver-Based

Extended Simplex Method

Branch-and-Bound

Linear Relaxation

SDP

Lipschitz

Lipschitz

Zeroth Order Verification

Zeroth Order Verification

First Order Verification

SMT-Based

MILP-Based

Linear Programming

Linear Inequality

Multi-Neuron Relaxation

General Lipschitz

Smooth Layers

Curvature

Lipschitz

Differential Privacy Based

Divergence Based

Neyman-Pearson

Level-Set Analysis

Lipschitz

Interval

Polyhedra

Zonotope

Duality
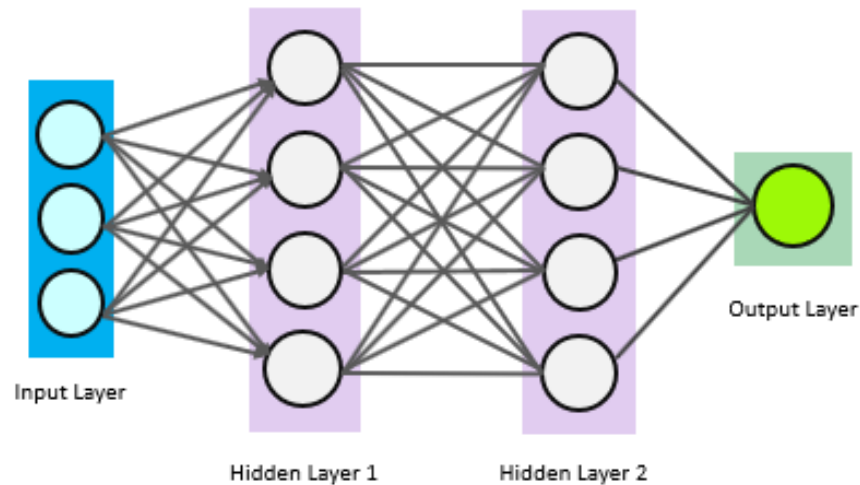
SOTA determinisitc certified robustness for general DNNs

SOTA probabilistic certified robustness
**Only** method supporting large models

11

# DNN Architecture



Input Layer

Hidden Layer 1    Hidden Layer 2

Output Layer

- **Input layer**: vector $x_0$
- **Weights**: $(W_0, b_0), (W_1, b_1), \dots (W_{L-1}, b_{L-1})$.
- **Activation function**:
  - $\text{ReLU}(x) = \max\{x, 0\}$
- **Computation**:
  - $x_1 = \text{ReLU}(W_0 x_0 + b_0)$,
  - $x_2 = \text{ReLU}(W_1 x_1 + b_1)$,
  - …
  - $x_L = W_{L-1} x_{L-1} + b_{L-1}$
- **Output**: $x_L$ - confidence score for each class

# Linear Relaxation of ReLU



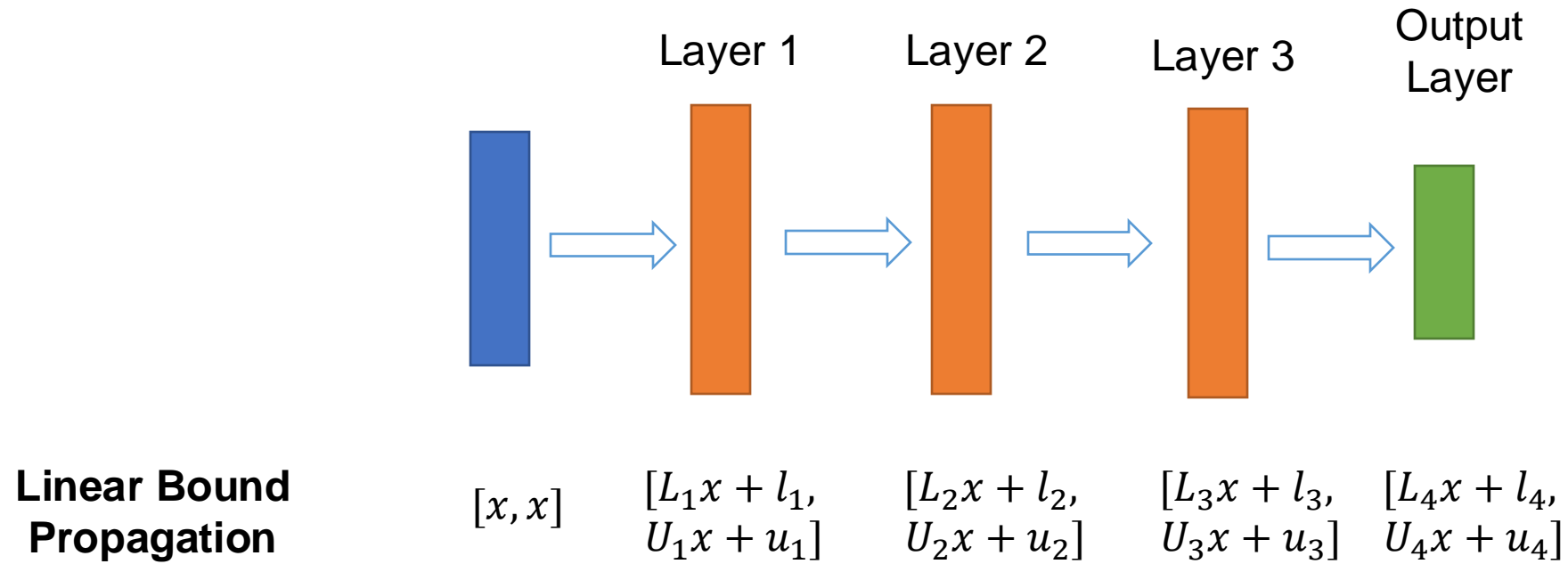Different Linear Relaxations for $\text{ReLU}(x) = \max\{x, 0\}$

*Weng, Lily, et al. "Towards fast computation of certified robustness for relu networks." ICML 2018*
*Wong, Eric, and Zico Kolter. "Provable defenses against adversarial examples via the convex outer adversarial polytope." ICML 2018*
*Singh, Gagandeep, et al. "Fast and Effective Robustness Certification." NIPS 2018*

# Linear Relaxation Induces Linear Inequality

Propagate linear inequalities that bound possible output region



Layer 1     Layer 2     Layer 3     Output Layer

**Linear Bound Propagation**

$[x, x]$     $[L_1 x + l_1, U_1 x + u_1]$     $[L_2 x + l_2, U_2 x + u_2]$     $[L_3 x + l_3, U_3 x + u_3]$     $[L_4 x + l_4, U_4 x + u_4]$

# Combat Over-Relaxation with Branch-and-Bound

Conditioned on two branches: $x \leq 0$ and $x > 0$,

each ReLU neuron is reduced to **linear constraints**: $y = 0$ or $y = x$

Select some neurons to condition on, and solve two subproblems

- If x<=0, y=0 (linearized subproblem)
- If x> 0, y=x (linearized subproblem)
- Relax other neurons by linear relaxation and bound propagation

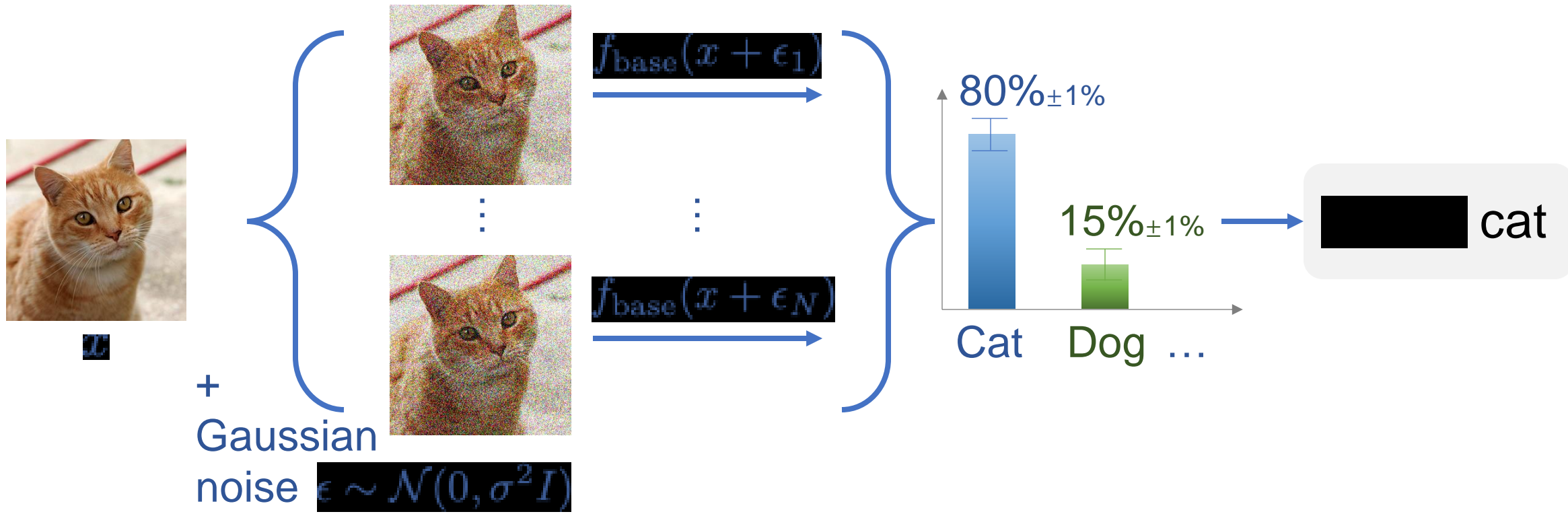- **Most scalable verification method so far**

# Randomized Smoothing

1. Train a model $f_{\text{base}}$ ("base classifier") under some known noise
2. Smooth $f_{\text{base}}$ into a new classifier $f$ ("smoothed classifier"), such that

   $f(x)$ = the most probable prediction by $f_{\text{base}}$ under noised corruptions of ■
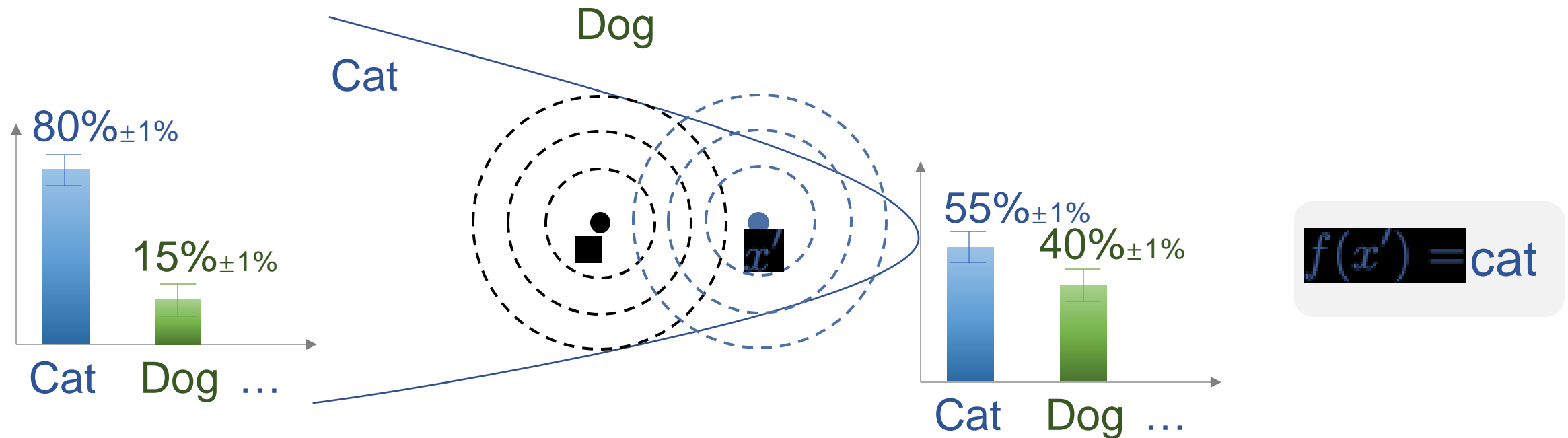
In deployment, use smoothed classifier $f$

# Illustration of Randomized Smoothing
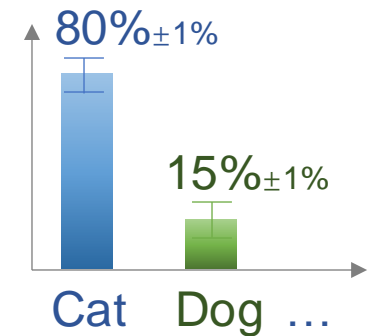
# Randomized Smoothing Enables Certified Robustness

Shift center of the distribution cannot change probability much
- If order doesn't change, then consistent prediction guaranteed
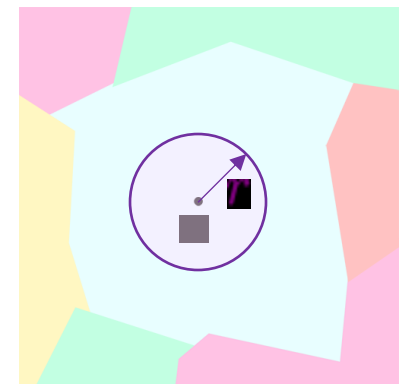


$f(x') = $ cat

# Closed-form Robustness Guarantee

- ████████'s probability of the top class (cat)
- ████████'s probability of the runner-up class (dog)

$f$ certifiably returns top class within an $\ell_2$ ball around $x$ of radius

$$r = \frac{\sigma}{2}\left(\Phi^{-1}(P_A) - \Phi^{-1}(P_B)\right)$$

- $\sigma$ : variance of Gaussian smoothing noise
- $\Phi^{-1}$: the inverse standard Gaussian CDF

80%±1%

15%±1%

Cat    Dog   …

Jeremy M Cohen, Elan Rosenfeld, J. Zico Kolter. Certified Adversarial Robustness via Randomized Smoothing. ICML 2019

# Certification Induces Robust Training

- Training DNN in specific ways can improve certified robustness

**For linear relaxation + branch-and-bound:**
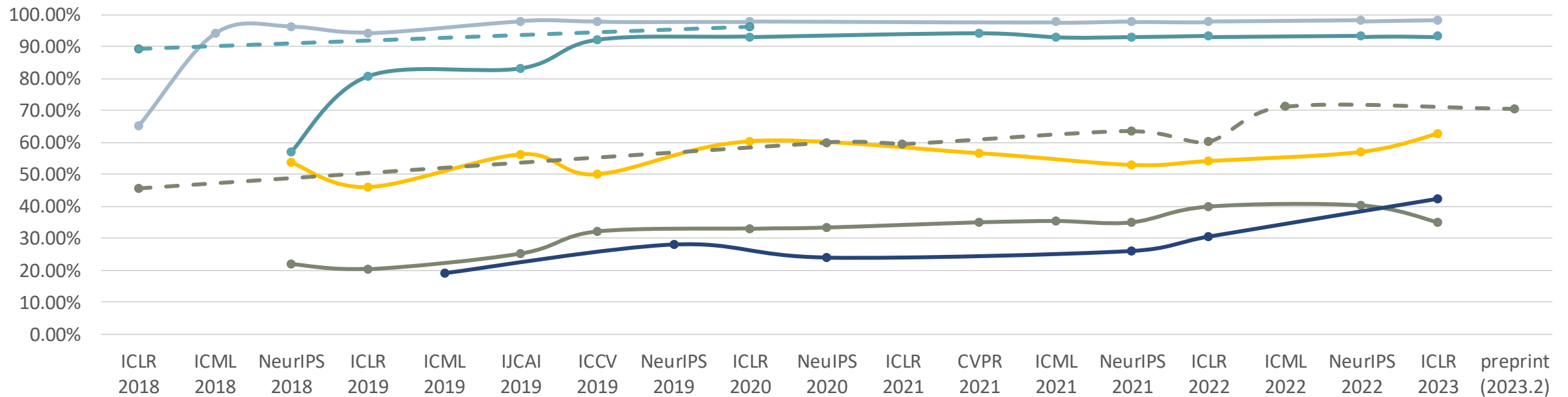training to **reduce upper bound of loss function** computed from over-approximation

**For randomized smoothing:**
training to **predict correctly for noised inputs**

# How Far Are We on Real-World Datasets?



Progress of Robustness on Typical Datasets and Settings

# On MNIST

$\ell_\infty$ norm, $r = 0.3$

- SOTA Certified Robust Accuracy: **94.02%**
  - *[CVPR 2021] Towards Evaluating and Training Verifiably Robust Neural Networks*

- SOTA Empirical Robust Accuracy (against existing attacks): **96.34%**
  - https://github.com/MadryLab/mnist_challenge
  - *Uncovering the Limits of Adversarial Training against Norm-Bounded Adversarial Examples*
  - *ArXiv: 2010.03593*

➢Not much difference

# On CIFAR-10

$\ell_\infty$ norm, $r = 8/255$:

- SOTA Certified Robust Accuracy: **40.39%**
  - *[NeurIPS 2022] Rethinking Lipschitz Neural Networks and Certified Robustness: A Boolean Function Perspective*

- SOTA Empirical Robust Accuracy (against existing attacks): **71.29%**
  - *[ICML 2022] Diffusion Models for Adversarial Purification*

➢Still a gap

# On ImageNet

$\ell_2$ norm, $r = 2.0$

- SOTA Certified Robust Accuracy: **30.4%**
  - *Our paper at [ICLR 2022] On the Certified Robustness for Ensemble Models and Beyond*

- SOTA empirical robustness accuracy: **43.18%**
  - Against $\ell_\infty$ norm, $r = \frac{4}{255}$
  - *[ICML 2022] Diffusion Models for Adversarial Purification*

- Hard to achieve robustness

# Key Messages

- Since 2017, **many methods proposed** to provide & improve DNN certified robustness
  - Linear relaxation
  - Branch-and-bound
  - Randomized smoothing
  - …
- **Remarkable certified robustness** achieved on small datasets, **but still challenging** on large ones
  - Good on MNIST
  - To be improved on CIFAR-10 and ImageNet

- Certification for p-norm bounded adversary **generalizable** for other threat models
  - Semantic adversary
  - Patch adversary
  - Word substitution adversary
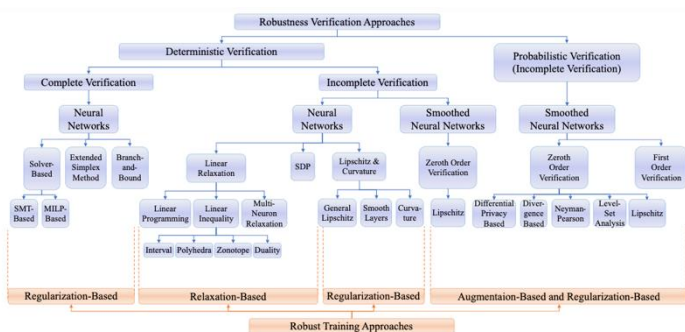  - Control state perturbation
  - Poisoning attack
  - …

[sokcertifiedrobustness.github.io](sokcertifiedrobustness.github.io)

SoK: Certified Robustness for Deep Neural Networks      Benchmark   Leaderboard   Paper   Website Repo   Toolbox Repo

## TAXONOMY

## SUMMARY

- Characteristics
- Strengths
- Limitations
- Connections
- Generalization
- …

## DISCUSSION

- Current Research
- Theoretical Barriers
- Main Challenges
- Future Directions
- …

## BENCHMARK

**VeriGauge**
*Open-source platform
for 20+ approaches*