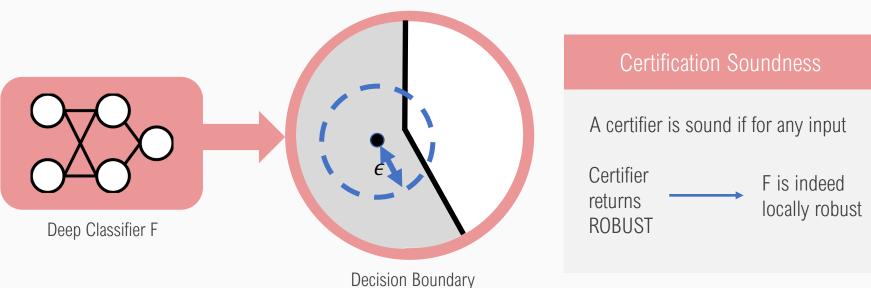# On the Perils of Cascading Robust Classifiers

Ravi Mangal*, Zifan Wang*, Chi Zhang*, Klas Leino, Corina Pasareanu, Matt Fredrikson
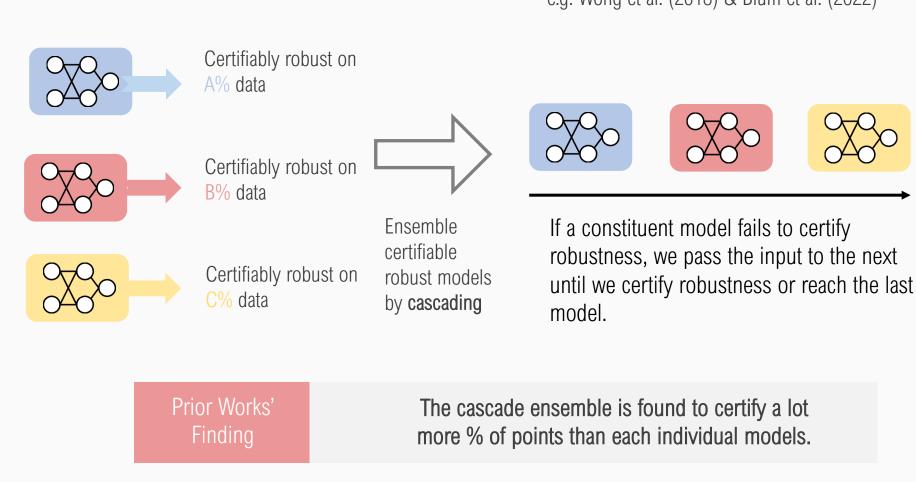
**Carnegie Mellon University**

## Introduction

### Certifying *Local Robustness* for A Deep Classifier

$$\forall x', ||x' - x|| \leq \epsilon \rightarrow F(x) = F(x')$$



Deep Classifier F

Decision Boundary

**Certification Soundness**

A certifier is sound if for any input

Certifier returns ROBUST → F is indeed locally robust

### One Strategy to Promote Certified Robust Is Using Cascading Ensemble

e.g. Wong et al. (2018) & Blum et al. (2022)



Certifiably robust on A% data

Certifiably robust on B% data

Certifiably robust on C% data

Ensemble certifiable robust models by **cascading**

If a constituent model fails to certify robustness, we pass the input to the next until we certify robustness or reach the last model.

**Prior Works' Finding** — The cascade ensemble is found to certify a lot more % of points than each individual models.

## Contribution

We prove that cascading certification is **UNSOUND**, even though each constituent model has a sound robust certifier.

Our Cascade Attack (CasA) successfully finds points that are claimed to be certifiable ROBUST according to the cascading ensemble created by Wong et al. (2018) but are **NOT** locally robust in fact.

Compared to cascading, ensemble by (weighted) voting is proved to be sound; however, gaining robustness improvement from voting requires model diversity to begin with.

## Cascading: An Unsound Certification

### Proof by Counter Example



*Cascading Direction*

*Output of Ensemble*

x

x'

Decision Boundaries

ROBUST as $y_{gray}$

ROBUST as $y_{blue}$

$$\exists x', ||x - x'|| \leq \epsilon \nRightarrow$$
$$Ensemble(x) = Ensemble(x')$$
$$\quad\quad y_{gray} \quad\quad\quad y_{blue}$$

**A violation of local robustness**

(see formal proof in the paper)

Constituent Models are not guaranteed to output the same classes for neighboring points without additionally memorizing other models' outputs.

### Visualizing Decision Boundaries for 2-D Datasets



Data

Decision boundaries of constituent models
(Darker regions denote certified robustness)

Decision boundaries of cascading ensemble

Zoom-in view

Boundaries are actually much "thinner" compared to each constituent model

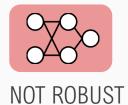## Weighted Voting: A Sound Certification

### Definition of Weighted Voting



ROBUST (w=0.3)   NOT ROBUST (w=0.4)   ROBUST (w=0.3)

Votes for ROBUST
0.3 + 0.3 = 0.6

>

Votes for Not ROBUST
0.4

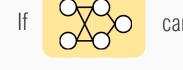**Certification of Ensemble: ROBUST**
(See proof of soundness in the paper)

### Can Weighted Voting Certify More Points?
(a thought experiment)
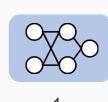
Let examples sorted as $x_1, x_2, ..., x_{100}$

If [model] can certify $x_0, x_2, ..., x_{49}$ and is correct.
Certified Acc. = 50%

If [model] can certify $x_{25}, x_{26}, ..., x_{74}$ and is correct.
Certified Acc. = 50%

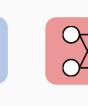If [model] can certify $x_0 ... x_{24}, x_{50} ... x_{74}$ and is correct.
Certified Acc. = 50%

And we let the weight for each model to be 1/3, then for each of the first 75 examples the voting ensemble is both correct and certifiably robust. → Certified Acc. = **75%**
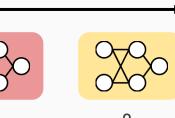
## Attacking Cascading Ensemble

### Basic Idea

To attack a cascading ensemble at $x$, we find a neighbor of $x$, $x'$, which can be certified with a different label on another constituent model $F'$
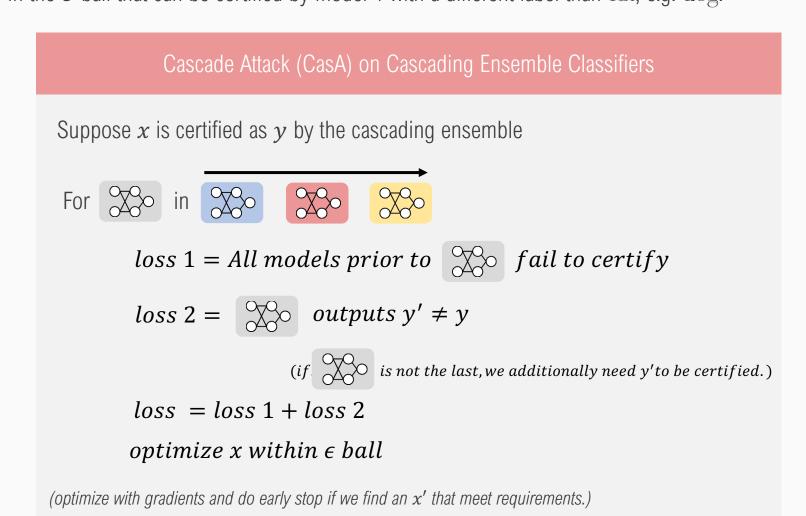- **prior to** the one that certifies $x$;
- or all other models **prior to** F' can not certify $x'$.



1   2   3

### Example

Suppose $x$ can not be certified by model 1 and is certified by model 2 as label cat, we find a neighbor $x'$ in the $\epsilon$-ball that can be certified by model 1 with a different label than cat, e.g. dog.

**Cascade Attack (CasA) on Cascading Ensemble Classifiers**

Suppose $x$ is certified as $y$ by the cascading ensemble

For [model] in [models]

$loss\ 1 = All\ models\ prior\ to$ [model] $fail\ to\ certify$

$loss\ 2 =$ [model] $outputs\ y' \neq y$

(if [model] is not the last, we additionally need y' to be certified.)

$loss\ = loss\ 1 + loss\ 2$

$optimize\ x\ within\ \epsilon\ ball$

(optimize with gradients and do early stop if we find an $x'$ that meet requirements.)

### Experiments

We test CasA on the cascading ensembles used in Wong et al. (2018) from their public repository. We copy and paste the Certified Acc. from their publication and show % of points that are claimed to be robust but later found to be not using CasA (False Positives). We report the **Empirical Robust Accuracy** of these cascading ensembles under CasA.

| Dataset | $\ell_p, \epsilon$ | Single Model Certified Acc. | Cascading Ensemble | | | |
|---|---|---|---|---|---|---|
| | | | Unsound Certified Acc. Reported | % Of False Positive Results. | Clean Acc. | Empirical Robust Accuracy Under CasA |
| MNIST | $\ell_\infty, 0.1$ | 95.54 | 96.63 | 88.71 | 96.62 | 11.17 |
| CIFAR-10 | $\ell_\infty, 2/255$ | 52.65 | 64.87 | 10.47 | 65.13 | 58.15 |
| MNIST | $\ell_2, 1.58$ | 43.52 | 75.58 | 44.72 | 80.43 | 43.46 |
| CIFAR-10 | $\ell_2, 36/255$ | 50.26 | 58.72 | 2.70 | 58.76 | 57.17 |

learn more | check out our talk and the full paper for more!
code available on GitHub