





# Cert-RNN: Towards Certifying the Robustness of Recurrent Neural Networks

Tianyu DuShouling JiLujia ShenYao ZhangJinfeng LiJie ShiChengfang FangJianwei YinRaheem BeyahTing Wang

### Background

#### **Speech Recognition**

#### **Sentiment Analysis**





#### **Face Recognition**



#### **Object Detection**



# Deep Learning for Multiple Tasks

#### **Political Content Detection**



## **Breaking Thing Is Easy**

• Adversarial examples exist **EVERYWHERE**!



Task: Toxic Content Detection. Classifier: LSTM. Original label: 96.7% Toxic. Adversarial label: 83.5% Non-toxic.
Text: hello how are you? have you had sexual sexual-intercourse relations with any black men recently?
Task: Toxic Content Detection. Classifier: Perspective. Original label: 92% Toxic. Adversarial label: 78% Non-toxic.
Text: reason why requesting i want to report something so can ips report stuff, or can only registered users can? if only registered users can, then i 'll request an account and it 's just not fair that i cannot edit because of this anon block shit shti c'mon, fucking fucking hell helled.

## **Defense against Adversarial Attacks**

## Empirical Defense



- Example: Adversarial training [Madry et al. 2017]
- Work empirically but no theoretical guarantee
- Attack specific leading to an arms race that attackers are winning

#### ICML 2018

Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples

Anish Athalye<sup>\*1</sup> Nicholas Carlini<sup>\*2</sup> David Wagner<sup>2</sup>

## **Key Questions**

- **Guaranteed accuracy**: what is the minimum accuracy under any attack?
- **Prediction robustness**: given a prediction, can any attack change it?

#### Certified Defense

- Theoretical guarantees against all attacks within a certain threat model
- Robustness certificate  $RC(x, F, \epsilon)$ : for all  $x' \in B(x, \epsilon)$  we have that F(x) = F(x')

## **Robustness Certification**

- The robustness property is determined by the **exact decision boundary**, which can be approximated by upper bound and lower bound.
- Adversarial attacks provide the asymptotic **upper bound**.
- The challenge is to compute the **lower bound** of the minimum adversarial distortions, i.e., certifying the robustness space around the input such that the model's prediction result is consistent within the space.



## **Robustness Certification**

- The basic idea to verify the robustness for a given  $\ell_p$ -norm perturbation space:
  - Compute the lower and upper bounds of the output units for the given perturbation space.
  - If the lower bound of the true label output is larger than the upper bounds of all other labels, the robustness for the given perturbation space is verified.



### **Robustness Certification Methods**

### Exact Certification

- Satisfiability Modulo Theories (SMT) [Ehlers et al. ATVA'17, Huang et al. CAV'17, Katz et al. CAV'17]
- Mixed-Integer Linear Programming (MILP) [Tjeng et al. ICLR '19]
- Accurate but usually computationally expensive, therefore cannot be scaled to large networks

## Relaxed Certification

- Convex Polytope [Wong & Kolter ICML'18]
- Reachability Analysis [Weng et al. ICML'18, Zhang et al. NeurIPS'18]
- Abstract Interpretation [Mirman *et al.* ICML'18, Singh *et al.* POPL'19]
- Efficient but cannot provide precise robustness bounds

### However, they are almost designed for FCNs and CNNs, seldom for RNNs!

## **Challenges for Certifying RNNs**



Figure 2: The architecture of an LSTM.

## **Robustness Certification for RNNs**

#### Current Works (categorized by threat model)

- Symbol/Word Substitutions (limited attackers' ability)
  - Wang et al. NAACL'21
  - Dong et al. ICLR'21
  - Ye et al. ACL'20
  - Huang et al. EMNLP'19
  - Jia et al. EMNLP'19
- Embedding Perturbation (strong attackers' ability)
  - POPQORN [Ko et al. ICML'19]
    - imprecise its linear relaxations do not retain high inter-variable correlations

**Possible? Yes!** 

Efficient

Practical

✓ Precise

- inefficient use gradient-based optimization to compute bounding planes
- impractical only evaluate one single word (one input frame) perturbation

## **Our Contribution**

- Leveraging abstract interpretation, we propose a novel certification framework for RNNs – Cert-RNN, which significantly outperforms prior work in terms of both precision and efficiency.
- We conduct extensive evaluation on **four security-sensitive applications** across **various network architectures** to empirically validate Cert-RNN's superiority.
- The robustness bound certified by Cert-RNN can be practically used as a meaningful quantitative metric for evaluating both the interpretability of RNNs and the provable effectiveness of various defense methods. We also demonstrate Cert-RNN's superiority in improving the robustness of RNNs.



#### **Abstract Interpretation**



#### **Abstract Interpretation**

- Three popular numerical abstract domains
- We choose zonotope abstract domain for the following reasons:
  - Trades off precision and performance
  - Each variable (abstract neuron) captured in an affine form -> exact for linear operations
  - Allows relating variables through parameters



## **Framework of Cert-RNN**



#### **Main Steps**

- 1. A zonotope abstract domain is first defined to capture all potential adversarial inputs
- 2. An abstract transformer is created for each non-linear operation of the RNN
- 3. Propagating the zonotope through all the layers of the target RNN
- 4. The output zonotope of the RNN's last layer is used to certify the robustness

## **Problem Definition**

**Definition 4.2** Given a continuously differentiable non-linear function  $f(x_1, x_2, ...)$  defined in a zonotope, the zonotope approximation for f consists of two parallel planes: the lower bounding plane  $Z^L$  and the upper bounding plane  $Z^U$ . We define  $Z^L$  and  $Z^U$  for any  $(x_1, x_2, ...) \in \mathbf{z}$  as follows:

$$Z^{L} = C_{1} + a_{1} \cdot x_{1} + a_{2} \cdot x_{2} + \cdots,$$
  
$$Z^{U} = C_{2} + a_{1} \cdot x_{1} + a_{2} \cdot x_{2} + \cdots,$$

•  $C_1, C_2, a_i \in \mathbb{R}$ 



- when  $a_i = 0$  (i = 1, 2, ...), the zonotope approximation returns the interval range of f, i.e., [ $C_1, C_2$ ]
- **Problem Definition** Given a non-linear function f and its bounding planes  $Z^L$ ,  $Z^U$ , its output region can be

bounded by a zonotope  $z_o = a_1 \cdot z_1 + a_2 \cdot z_2 + \dots + \frac{C_2 - C_1}{2} \varepsilon_{new}$ , where  $\varepsilon_{new}$  is a new error term which is introduced from the zonotope approximation for f. Thus, the problem to find the tightest bound of  $z_o$  can be formalized as bellow:

$$\min \quad \frac{C_2-C_1}{2}.$$

## **Step 1: Input Region Abstraction**



- Given an input sequence  $X = [x^{(0)}, x^{(1)}, \dots, x^{(t)}, \dots, x^{(T)}]$ , where  $x^{(t)} = [x_1^{(t)}, x_2^{(t)}, \dots, x_K^{(t)}]$  represents the *t*-th input frame.
- Based on Definition 4.1, the input frame  $x^{(t)}$  is mapped to the center coefficient  $\alpha_0$  of a zonotope z.
- For  $\ell_{\infty}$ -norm bounded attack, the adversarial perturbation of the j-th dimension of  $\mathbf{x}^{(t)}$  is mapped to the coefficient  $\alpha_{ij}$ .

#### **Step 2: Intermediate Operation Abstraction**



- Affine Transformation Abstraction
  - Can be exactly captured in our approximation
- Tanh Function Abstract Transformer
  - We propose a new abstract transformer for tanh
  - Tighter than DeepZ[Singh et al. NeurIPS'18]



### **Cert-RNN**

#### Intermediate Operation Abstraction

• Sigmoid • Tanh Abstract Transformer

THEOREM 4.2. Let  $z = \sigma(x) \cdot \tanh(y)$ , where  $(x, y) \in \mathbb{Z} \subseteq [l_x, u_x] \times [l_y, u_y]$ . Then, the fine-grained zonotope approximation planes in  $\mathbb{Z}$  are:

$$Z^{L} = C_1 + Ax + By$$
$$Z^{U} = C_2 + Ax + By$$

where A, B,  $C_1$ ,  $C_2$  have nine different cases as shown in Tab. 8 (deferred to Appendix B) according to the value of  $l_x$ ,  $u_x$ ,  $l_y$  and  $u_y$ .

Card	Cadline	Status	Front
1	$I_1 \supseteq \operatorname{Rend} I_1 \supseteq H$	$\begin{array}{l} A \sim \operatorname{Hom}(\operatorname{Hilling}(a_1,a_2),\operatorname{Hom}(a_1,a_2), \\ C_1 \in f_{\operatorname{Hom}}(C,T') = A^{-1} - A^{-1} - A^{-1} - A^{-1} + A^{-1} - A^$	Appendix T
2	n. 58 mil 1, 54	$\begin{split} & A = \lim_{n \to \infty} \lim_{n \to \infty} \lim_{n \to \infty} \  (x_n) - C_1 & \leq \lim_{n \to \infty} \  (x_n) - A_1 + B_{n_1} - C_2 - S_{n-m}) \  (x') - A_1' + B_{2'} \\ & = \left[ (n, 1_1) - A_2 - \frac{\min\{1_1, 1_1\}}{4} - \frac{\min\{1_1, 1_2\}}{4} - \frac{\min\{1_1, 1_2\}}{2m} (n, 1_1) - A_1 + B_{2'} + B_{$	Appendix F.
		$\begin{split} \mathcal{A} &= \frac{2i_{\text{transf}}}{2} (a_1,b_1) = \mathcal{B} = \begin{cases} \frac{\beta_1(a_1)(f_1,b_1)}{p_1} - \frac{\beta_1(a_1)(f_1,b_1)}{p_2} & \mathcal{A} > \frac{M_{\text{transf}}}{d_2} (\beta_1,a_1) \\ \frac{\beta_1(a_1)(f_1,a_2)}{p_1} - \frac{\beta_1(a_1)(f_1,a_2)}{p_2} & \mathcal{A} > \frac{M_{\text{transf}}}{d_2} (\beta_1,a_2) \\ \mathcal{A} &= \frac{M_{\text{transf}}}{p_1} (a_1,b_2) - 4a_1 - B_1, \\ \frac{\beta_1(a_2)(f_1,a_2)}{p_1} - \frac{\beta_1(a_2)(f_1,a_2)}{p_2} - M_1 - B_2) & \frac{M_{\text{transf}}}{p_2} (a_1,b_2) = \frac{M_{\text{transf}}}{p_2} (a_1,b_2) + \mathcal{B} \\ \max \left(\beta_{1(a_2)}(a_1,a_2) - 4a_2 - B_1,\beta_{1(a_2)}(f_1,a_2) - M_2 - B_2\right) & \frac{M_{\text{transf}}}{p_2} (a_1,b_2) \geq \frac{M_{\text{transf}}}{p_2} (a_1,b_2) + \mathcal{B} \\ \max \left(\beta_{1(a_2)}(a_1,a_2) - 4a_2 - B_1,\beta_{1(a_2)}(f_1,a_2) - M_2 - B_2\right) & \frac{M_{\text{transf}}}{p_2} (a_1,b_2) \geq \frac{M_{\text{transf}}}{p_2} (a_2,b_2) + \mathcal{B} \\ \max \left(\beta_{1(a_2)}(a_1,a_2) - 4a_2 - B_2,\beta_{1(a_2)}(f_1,a_2) - M_2 - B_2\right) & \frac{M_{\text{transf}}}{p_2} (a_2,b_2) + \frac{M_{\text{transf}}}{p_2} (a_2,b_2) + \mathcal{B} \\ \max \left(\beta_1(a_2) - A_2 - B_2,\beta_{1(a_2)}(f_1,a_2) - M_2 - B_2\right) & \frac{M_{\text{transf}}}{p_2} (a_2,b_2) + \mathcal{B} \\ \max \left(\beta_1(a_2) - \beta_2(a_2) - \beta_2$	Appendia N
		$\begin{split} A &= \begin{cases} \frac{\beta (l_1-r^2) + \log_{20}(l_2-r^2) - \log_{20}(l_1,L_1)}{\eta_1} & \beta \geq \frac{\delta (\log_{20}(l_1,L_1))}{\eta_2} \\ &= \frac{\log_{20}(l_1+r^2) + \log_{20}(l_1,r^2) - \log_{20}(l_1,L_1)}{\eta_1} \\ &= \frac{\log_{20}(l_1+r^2) + \log_{20}(l_1,r^2) - \log_{20}(l_1,L_1)}{\eta_2} \\ C_1 &= \int_{20}(\log_{20}(l_1-r^2) + \log_{20}(l_1-r^2) +$	Agendin T
		$\begin{split} A &= \begin{cases} \frac{\delta(l_1-r^2) + \int \log_2(h_1, r^2) - \int \log_2(h, r^2)}{h_1} & \# > \frac{\delta(l_1, r_1)}{h_1} \\ \# = \frac{\int \log_2(r + r^2) + \int \log_2(h, r^2) - \int \log_2(h, r^2)}{h_1} & \# = \frac{\int \log_2(h, r_1)}{h_1} \\ \frac{\delta(l_1-r^2) + \int \log_2(h, r^2) - \int \log_2(h, r^2)}{h_1} \\ F_1 &= \int \log_2(h, r_1) - \int \log_2(h, r^2) + \int \log_2(h, r^2)$	4,quendis 1
		$A_{ij} = \frac{\pi (z_i) + ((1/(m_{ij}, -m_{ij})))}{2\pi},  B_{ij} = \frac{(m_{ij}, -(z_i)) + (m_{ij}, -(m_{ij}))}{2\pi},  C_1 = \begin{cases} \min\{P(i, 1), P(i^*, n_i)\},  G_2(z_i, n_i) < A_i G(z_i, l_i) > A \\ P(i^*, l_i),  G_2(z_i, n_i) < A_i G(z_i, l_i) > A \\ P(i^*, l_i),  G_2(z_i, n_i) < A_i G(z_i, l_i) > A \\ P(i^*, l_i),  G_2(z_i, n_i) < A_i G(z_i, l_i) > A \\ G_1(z_i, n_i) < G_2(z_i, n_i) < G_2(z_i, n_i) < A_i G(z_i, l_i) > A \\ G_1(z_i, n_i) < G_2(z_i, n_i) < G_2(z_i, n_i) < A_i G(z_i, l_i) > A \\ F_1(z_i, l_i),  G_2(z_i, n_i) < G_2(z_i, n_i) < G_2(z_i, n_i) > A \\ G_2(z_i, n_i) < G_2(z_i, n_i) > A \\ F_1(z_i, l_i) < G_2(z_i, n_i) < G_2(z_i, n_i) > A \\ G_2(z_i, n_i) < G_2(z_i, n_i) > A \\ G_2(z_i, n_i) < G_2(z_i, n_i) < G_2(z_i, n_i) > A \\ G_2(z_i, n_i$	tgeals)
Л.	as ST and as ST	The same or systematics to the same I, and the proof can be deduced to the same way	
	$l_{x} \geq 0$ and $a_{x} < 0$ $d_{x} \geq 0$ and $l_{y} < 0$	This case or spectration is the case $L$ and the proof cas by deducted in the same way $A = 0,  B = \min\{\frac{f_{max}(u_n, u_n) - f_{max}(u_n, u_n)}{u_n}, \frac{f_{max}(u_n, u_n)}{u_n}, f_$	Agreads 11
	with a school of the	This case we presented to the case A and the proof can be defaulted to the same way	_
0			
9	$\delta_{\rm p}\!\geq\!6, a_{\rm q}\!>\!6{\rm mat}\delta_{\rm q}\!<\!6$	$\mathbf{g} = \frac{g_{\mathrm{sum}}(\mathbf{r}_{1}^{+}) \cdot g_{\mathrm{sum}}(\mathbf{r}_{1}^{+})}{g_{1}},  \mathbf{g} = \frac{g_{\mathrm{sum}}(\mathbf{r}_{1}^{+}) \cdot g_{\mathrm{sum}}(\mathbf{r}_{1}^{+})}{g_{1}} \cdot \mathbf{A}(\mathbf{r}_{1}^{+}) \cdot \mathbf{A}(\mathbf{r}_{$	Appendia 1

#### **Cert-RNN**

#### Intermediate Operation Abstraction

• Sigmoid • Identity Abstract Transformer

THEOREM 4.3. Let  $z = x \cdot \sigma(y)$ , where  $(x, y) \in \mathbb{Z} \subseteq [l_x, u_x] \times [l_y, u_y]$ . Then, the zonotope approximation planes in  $\mathbb{Z}$  are:

$$Z^{L} = C_1 + Ax + By$$
$$Z^{U} = C_2 + Ax + By$$

where A, B,  $C_1$ ,  $C_2$  have three different cases as shown in Tab. 9 (deferred to Appendix C) according to the value of  $l_x$ ,  $u_x$ .

Case	Conditions	Solutions	Proof
1	$l_X \ge 0$	$A = \frac{(\sigma(u_x) - \sigma(l_x))(\tanh(u_y) + \tanh(l_y))}{2w_x},  B = \frac{(\sigma(u_x) + \sigma(l_x))(\tanh(u_y) - \tanh(l_y))}{2w_y},  C_1 = f_{x \cdot \sigma}(x^{\star \star}, y^{\star \star}) - Ax^{\star \star} - By^{\star \star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_1 = f_{x \cdot \sigma}(x^{\star \star}, y^{\star \star}) - Ax^{\star \star} - By^{\star \star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_1 = f_{x \cdot \sigma}(x^{\star \star}, y^{\star \star}) - Ax^{\star \star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - Ax^{\star} - By^{\star},  C_2 = f_{x \cdot \sigma}(x^{\star}, y^{\star}) - A$	Appendix C.1
2	$u_x \leq 0$	In this case, we use the same method used in Case 1.	
3	$l_x < 0$ and $u_x > 0$	$A = \min\{\frac{f_{x \cdot \sigma}(u_x, u_y) - f_{x \cdot \sigma}(l_x, l_y)}{w_x}, \frac{f_{x \cdot \sigma}(u_x, l_y) - f_{x \cdot \sigma}(l_x, u_y)}{w_x}\},  B = 0,  C_1 = f_{x \cdot \sigma}(l_x, u_y) - Al_x,  C_2 = f_{x \cdot \sigma}(u_x, u_y) - Au_x$	Appendix C.2

### **Cert-RNN**

#### Certifying the Robustness Bound

 Specifically, finding the largest robustness bound ε<sub>c</sub> for the input sequence with true label c can be formalized as the following optimization problem:

max  $\varepsilon_c$ 

s.t. 
$$\alpha_{0c} - \sum_{j=1}^{p} |\alpha_{jc} \cdot \varepsilon_j| \ge \alpha_{0i} + \sum_{j=1}^{p} |\alpha_{ji} \cdot \varepsilon_j|, \quad \forall i \neq c$$

Algorithm 1: Computing the robustness bound. **Result:** Certified robustness bound  $\varepsilon_c$ **Data:** model  $\mathcal{F}$ , input sequence  $\mathbf{X}_0$ , true label c 1 for t in T do  $e^{(t)} = 0.5$ 2 for l = 2 to 13 do 3  $\mathbf{z}_o = \text{CERT-RNN}(t, \mathcal{F}, \mathbf{X}_0, \boldsymbol{\varepsilon}^{(t)});$ 4 if  $\alpha_{c0} - \sum_{j=1}^{p} |\alpha_{c_j} \cdot \varepsilon_j| \ge \alpha_{i0} + \sum_{j=1}^{p} |\alpha_{ij} \cdot \varepsilon_j|$ 5 then  $\mathbf{\varepsilon}^{(t)} = \mathbf{\varepsilon}^{(t)} + 0.5^l;$ 6 else 7  $\mathbf{\varepsilon}^{(t)} = \mathbf{\varepsilon}^{(t)} - 0.5^l;$ 8 9  $\varepsilon_c = min(\varepsilon^{(1)}, \varepsilon^{(2)}, \cdots, \varepsilon^{(T)})$ 

## A Toy Example



The true label's confidence value  $z_{o2}$  always larger than  $z_{o1}$ , thus the robustness is verified for  $\varepsilon = 1$ .



## **Experimental Setting**

#### **Dataset & Models**

Dataset	MNIST Se	equence	Rotten Tomatoes			Toxi	ic Commen	t Detection	Malicious URL Detection			
Duusot	# of Images	Size	Positive	Negative	Avg Length	Toxic	Normal	Avg Length	Malicious	Benign	Avg Length	
Training	60,000	$28 \times 28$	23,498	15,564	23 words	6,720	6,720	32 words	60,450	275,921	48 chars	
Validation	1	1	3,362	1,562	23 words	1,280	1,280	32 words	7,567	34,479	48 chars	
Testing	10,000	$28 \times 28$	3,016	1,867	22 words	1,280	1,280	34 words	7,625	34,420	48 chars	

- 8 vanilla RNNs and 9 LSTMs with different hidden units and layers for MNIST Sequence
- an RNN and an LSTM with 32 hidden units for the other three datasets

#### **Baseline Method**

• POPQORN [Ko et al., ICML'19]

## **Experimental Setting**

#### **Evaluation Metrics**

#### Certified robustness bound

• the certified robustness bound of a particular sample x is the maximum  $\varepsilon$  for which we can certify that the model f(x') will return the correct label, where x' is any adversarially perturbed version of x such that  $||x - x'||_{\infty} \le \varepsilon$ 

#### Verified accuracy

• the verified accuracy at  $\varepsilon$  of a dataset if the fraction of data items in the dataset with certified robustness bound of at least  $\varepsilon$ 

## **Effectiveness & Efficiency**

Dataset	Model	Acc		POPQORN			CERT-RNN		
Dataset	Woder	The	Mean	Std	Time (min)	Mean	Std	Time (min)	
	RNN-2-32	96.8%	0.0084	0.0037	0.13	0.0157	0.0077	0.61	
	RNN-2-64	94.4%	0.0084	0.0033	0.12	0.0152	0.0076	0.63	
	RNN-4-32	95.4%	0.0168	0.0058	0.30	0.0222	0.0074	1.72	
	RNN-4-64	94.8%	0.0034	0.0018	0.40	0.0056	0.0032	1.70	
	RNN-7-32	89.0%	0.0027	0.0016	0.64	0.0037	0.0025	4.01	
	RNN-7-64	92.2%	0.0012	0.0012	0.60	0.0018	0.0012	4.21	
	RNN-14-32	92.2%	0.0190	0.0064	1.44	0.0270	0.0075	13.44	
MNIST	RNN-14-64	95.8%	0.0089	0.0030	2.31	0.0166	0.0044	14.38	
Sequence	LSTM-1-32	98.0%	0.0152	0.0071	46.78	0.0187	0.0087	2.66	
1	LSTM-1-64	99.0%	0.0152	0.0064	53.09	0.0178	0.0075	4.92	
	LSTM-1-128	98.0%	0.0143	0.0065	53.09	0.0184	0.0074	3.98	
	LSTM-2-32	96.0%	0.0147	0.0062	150.00	0.0176	0.0080	8.42	
	LSTM-2-64	98.0%	0.0145	0.0063	246.50	0.0167	0.0067	11.92	
	LSTM-2-128	97.4%	0.0129	0.0052	192.77	0.0143	0.0056	12.77	
	LSTM-4-32	95.0%	0.0093	0.0045	551.70	0.0095	0.0045	29.24	
	LSTM-4-64	97.8%	0.0088	0.0040	593.31	0.0092	0.0039	37.13	
	LSTM-7-32	96.6%	0.0054	0.0017	1522.77	0.0056	0.0015	90.99	
DT	RNN	76.0%	0.0091	0.0049	1342.20	0.0207	0.0098	40.20	
RI	LSTM	82.0%	. e	*		0.0080	0.0026	2464.2	
TO	RNN	90.0%	0.0190	0.0107	2070.60	0.0332	0.0243	98.40	
ic	LSTM	93.0%	5 <b>.</b> -	*	1	0.0117	0.0068	3903.60	
MalURL	RNN	94.0%	0.0282	0.0132	2923.80	0.0361	0.0203	243.60	
	LSTM	98.0%		-		0.0097	0.0044	9851.40	

Table 2: Evaluation results in the four scenarios, including model accuracy (Acc), mean value and standard deviation of the certified robustness bound (where a large mean implies a large robustness space), and running time.

#### Table 3: Mann-Whitney U test results.

Model	RNN-2-32	RNN-4-32	RNN-7-32	RNN-14-32
p-value	6.93×10 <sup>-9</sup>	$1.91 \times 10^{-22}$	$2.10 \times 10^{-29}$	$1.11 \times 10^{-30}$
Model	RNN-2-64	RNN-4-64	RNN-7-64	RNN-14-64
p-value	$1.12 \times 10^{-20}$	$1.83 \times 10^{-12}$	$4.81 \times 10^{-7}$	$2.76 \times 10^{-12}$

#### > Remarks

٠

٠

In all cases, Cert-RNN can obtain larger robustness bounds than that of POPQORN, i.e., the result of Cert-RNN is more accurate.

Cert-RNN is much more efficient than POPQORN in general, especially for large and complex networks.

For Mann-Whitney U test, the p-values of all models are small enough to reject the null hypothesis, which further demonstrates the superiority of Cert-RNN.

## **Effectiveness & Efficiency**



Figure 6: Certified robustness bound in the four scenarios. The violin plot shows the data distribution shape and its probability density, which combines the features of box and density charts. The thick black bar in the middle indicates the quartile range, the thin black line extending from it represents the 95% confidence interval, and the white point is the median.

#### > Remarks

- When the number of hidden units is the same, LSTMs with less layers would be more robust.
- When the number of layers is same, LSTMs with less hidden units would be more robust.
- Too many hidden units may increase the attack surface and decrease the generalizability

(i.e., have a high variance) of the model, which makes it less robust.

### **Verified Accuracy**



Fig. 7. Verified accuracy of the four datasets for each bound  $\epsilon \in \Delta$  (the x axis). The subfigures (a) to (q) are the results of the MNIST sequence dataset.

> Remarks

• The verified accuracy of Cert-RNN is much higher than that of POPQORN in most cases.

### **A More Threatening Scenario**

#### > Perturbing All Frames

- Cert-RNN can handle this threat model while POPQORN cannot.
- Compared with perturbing one single frame, the robustness bounds for perturbing all frames decrease to some extent.

1		Cert-RNN	
	Mean	Std	Time (sec)
RNN-2-32	0.0126	0.0055	6.8420
RNN-2-64	0.0130	0.0056	9.0874
RNN-4-32	0.0044	0.0044	0.0044
RNN-4-64	0.0047	0.0023	14.3441
RNN-7-32	0.0044	0.0044	20.5882
RNN-7-64	0.0017	0.0009	15.4963
RNN-14-32	0.0127	0.0036	31.8162
RNN-14-64	0.0074	0.0020	34.0596

Table 6: Results for perturbing all frames on the MNIST sequence dataset.



## **Certifying Adversarial Defenses**



#### Defense Methods

- FGSM-AT (Fast Gradient Sign Method-based Adversarial Training) (Goodfellow et al. ICLR'15)
- PGD-AT (Projected Gradient Descent-based Adversarial Training) (Madry et al. ICLR'18)
- IBP-VT (Interval Bound Propagation-based Verified Training) (Gowal et al. ICCV'19)

#### > Remarks

• Cert-RNN can provide an accurate qualitative metric to evaluate the provable effectiveness of various defenses, which would be more reliable than previous empirical metrics, e.g., the attack success rate after applying a defense method.

### Improving RNN Robustness

#### Implementation

- Our training follows [Gowal et al. CVPR'19, Mirman et al. ICML'18] we perturb the input signal and propagate interval bounds obtained by Cert-RNN through the RNN stages.
- To train, we combine standard loss with the worst case loss obtained using interval propagation.

#### > Experimental Results

• The RNNs trained with Cert-RNN-VT achieve larger robustness bounds, outperforming the RNNs trained with IBP-VT on all three datasets. This is because the interval bounds obtained by our approximation of the tanh function is more accurate than that obtained by the IBP method.

# Table 7: Certified robustness bounds for verified robustly trained RNNs.

Dataset	Original	IBP-VT	Cert-RNN-VT
RT	0.0207	0.0219	0.0224
TC	0.0332	0.0428	0.0436
MalURL	0.0361	0.0702	0.0730

## **Identifying Sensitive Words**

Example	<u>this</u>	is	a	<u>stupid</u>	<u>idea</u>	all	it	is	doing	<b>g is</b>	adding	junk	to	an	alread	y good	page
Bound	0.0183	0.0188 (	).0222	0.0178	0.0183	0.0232	0.0315	0.0320	0.031	5 0.0334	0.0320	0.0334	0.0398	0.0427	7 0.0457	0.0496	0.0564
Example Bound	you 0.0178	are 0.017	<b>'</b> 3 0.	an .0188	<u>idiot</u> 0.0149	nothing 0.0188	sugge 0.019	sts th 3 0.0	nat 212	she 0.0247	needs 0.0247	to 0.0305	atte 5 0.03	nd 05 0	a ).0305	hearing 0.0217	
Example	hi	,	<u>i</u>	<u>diot</u>	<b>,</b>	why	are	y	'ou	delate	my	talking	,	93 0	just	come	out
Bound	0.0134	0.011(	0 0.	.0071	0.0085	0.0115	0.013	9 0.0	183	0.0154	0.0208	0.0159	0.019		.0232	0.0256	0.0291
Example	oh	yeah	ר	,	you	ʻre	really	, pro	oof	of	the	<b>hypocr</b>	<u>isy</u> of	f wi	kipedia	right	here
Bound	0.0090	0.009	ס5 0.	.0110	0.0120	0.0129	0.012	9 0.0	090	0.0110	0.0129	0.0085	5 0.01	20 C	).0153	0.0193	0.0242
Example	you	mus	st	be	a	real	<u>lose</u>	<u>r</u> a	ind <u>I</u>	<u>mental</u>	infant	to	try	,	to	block	me
Bound	0.0105	5 0.009	95 0.	.0125	0.0144	0.0105	0.008	10.0	134	0.0081	0.0139	0.0183	3 0.01	98 (	).0212	0.0247	0.0237

#### Remarks

• The words with smaller certified robustness bounds tend to be more important for the final prediction result, i.e., more sensitive.

### **Limitation & Discussion**

#### >Improving Zonotope Approximation

• Explore alternative zonotope approximations which lead to tighter robustness bounds

#### Supporting Other Norm-bounded Attacks

• The perturbations bounded by other norms can be considered as the subsets of  $\ell_{\infty}$  in Cert-RNN

#### Supporting More Network Types

- Directly applicable to Gated Recurrent Unit (GRU) model
- New abstract transformers for attention module in Transformers
- The possibility for certifying sequence-to-sequence models

#### Supporting Other Threat Models

• Word substitution perturbation

## Conclusion

## Cert-RNN has three important advantages:

- a) Effectiveness it provides much tighter robustness bounds.
- **b)** Efficiency it scales to much more complex models.
- c) Practicality it enables a range of practical applications including evaluating the provable effectiveness for various defenses, improving the robustness of RNNs and identifying sensitive words.



zjradty@zju.edu.cn