

# Formal Reasoning About Confidence and Automated Verification of NNs

Mohammad Afzal, S. Akshay, Blaise Genest, Ashutosh Gupta  
IIT Bombay India, TCS research Pune India, CNRS@CREATE & IPAL, Singapore

## 1. Motivation

### Neural Networks in safety-critical applications:

- Autonomous vehicles [Bojarski M, et al, 2016]
- Medical diagnosis [Amato F, et al, 2013]
- Speech recognition [Hinton G, et al, 2012]

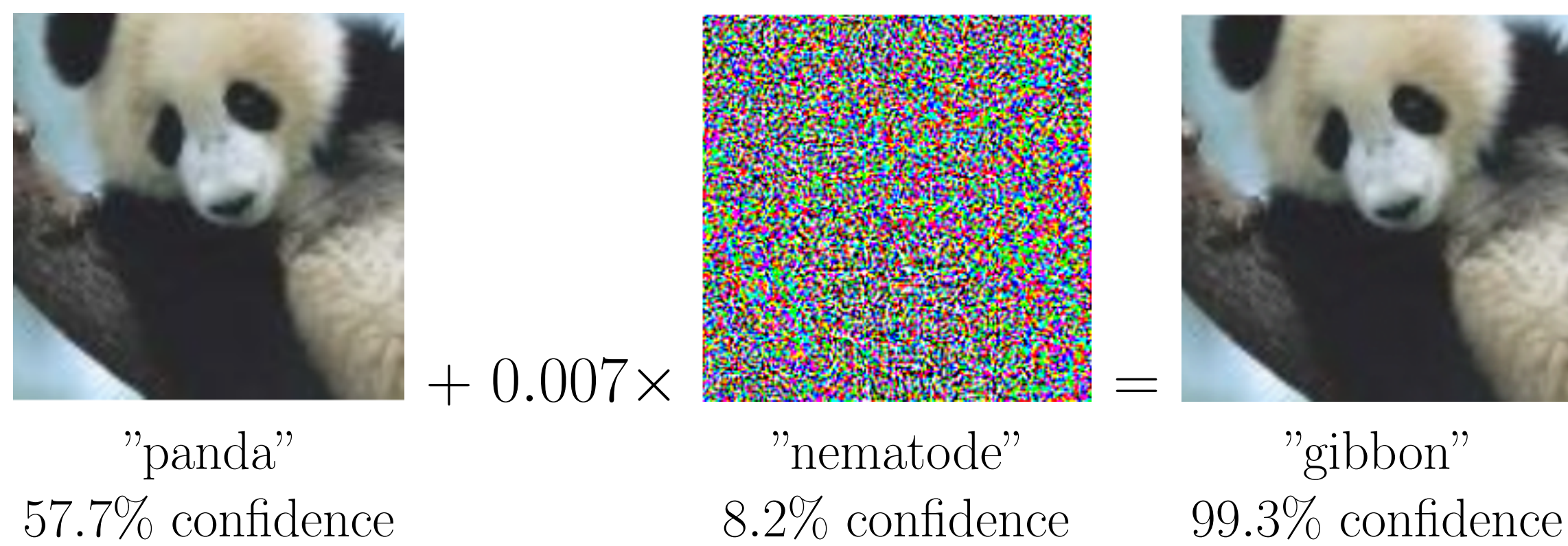


Figure 1: Image source: [Goodfellow I, et al, 2015]

## 2. Local robustness verification

- Neural network  $N: \mathbb{R}^n \rightarrow \mathbb{C}$ ,  $\mathbb{C}$  - classes
- Local robustness property,  $\mathbf{x}$  given image, [Tjeng V, et. al, 2017]:  
 $\neg \forall \mathbf{x}'. dist(\mathbf{x}, \mathbf{x}') \leq \epsilon \implies N(\mathbf{x}) = N(\mathbf{x}')$

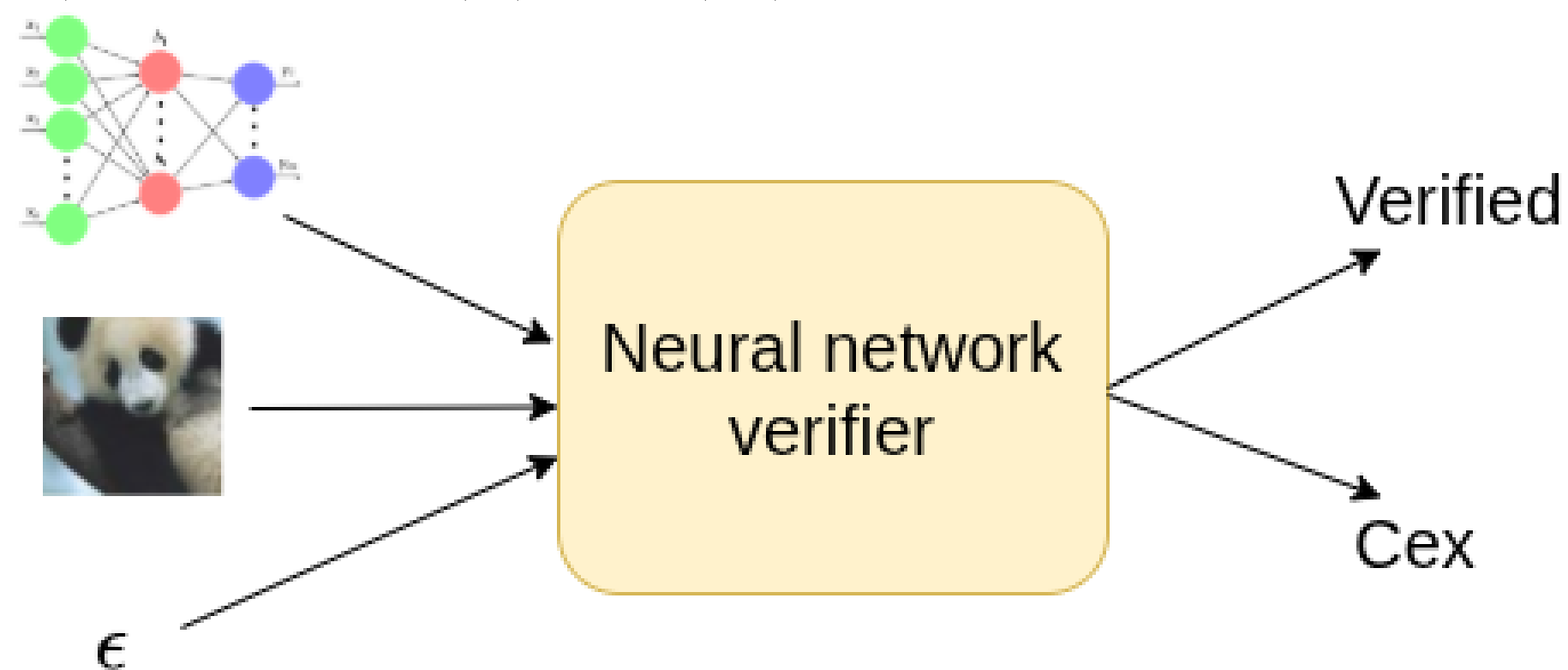
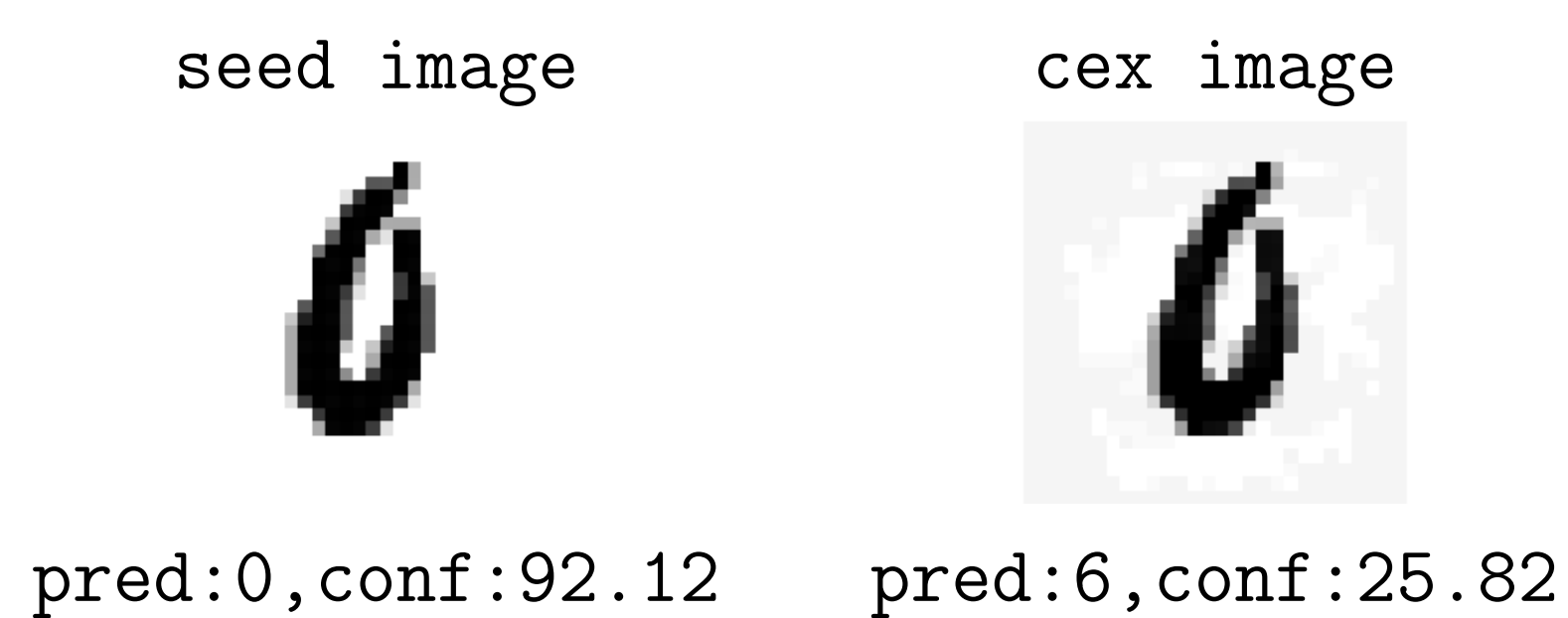


Figure 2: Neural networks verifier

- State-of-the-arts (SOTA) neural networks verifiers:  $\alpha\beta$ -CROWN, PyRAT, Marabou.
- SOTA ignore the classification confidence.

## 3. Confidence-based Robustness

$CONF(N(\mathbf{x}), c)$ : Define the confidence on class  $c$  with input  $\mathbf{x}$ , using SOFTMAX function.



**Relaxed Robustness:** Avoid low confidence counterexamples.

- $\forall \mathbf{x}'. dist(\mathbf{x}, \mathbf{x}') \leq \epsilon \implies N(\mathbf{x}) = N(\mathbf{x}') \vee CONF(\mathbf{x}', c) \leq th$

**Other Variations:**



**Strong Robustness**

**Smoothness**

**Strong Robustness:** High confidence decision should remain high.

**Smoothness:** Confidence should not change drastically.

Need of a unifying framework to capture all imagined confidence-based properties?

## 4. Unifying Framework for Property Variations

**Grammar:** Let  $\mathbf{y} = y_1, y_2, \dots, y_m = N(\mathbf{x})$

LE ::=  $c_1 y_1 + c_2 y_2 + \dots + c_m y_m + b, \forall i \in [m], c_i, b \in \mathbb{R}$

LC ::= LE > 0 | LE ≥ 0

CC ::= CONF( $\mathbf{y}, t$ )  $\bowtie$   $b, t \in [m], b \in \mathbb{R}^+, \bowtie \in \{\leq, <, \geq, >\}$

PC ::= PC  $\wedge$  PC | PC  $\vee$  PC | CC | LC |  $\neg$ LC

- PC represents the set of post-conditions.
- Need to approximate the CONF.

**Approx of  $conf(\bar{y}, t) \leq b$**

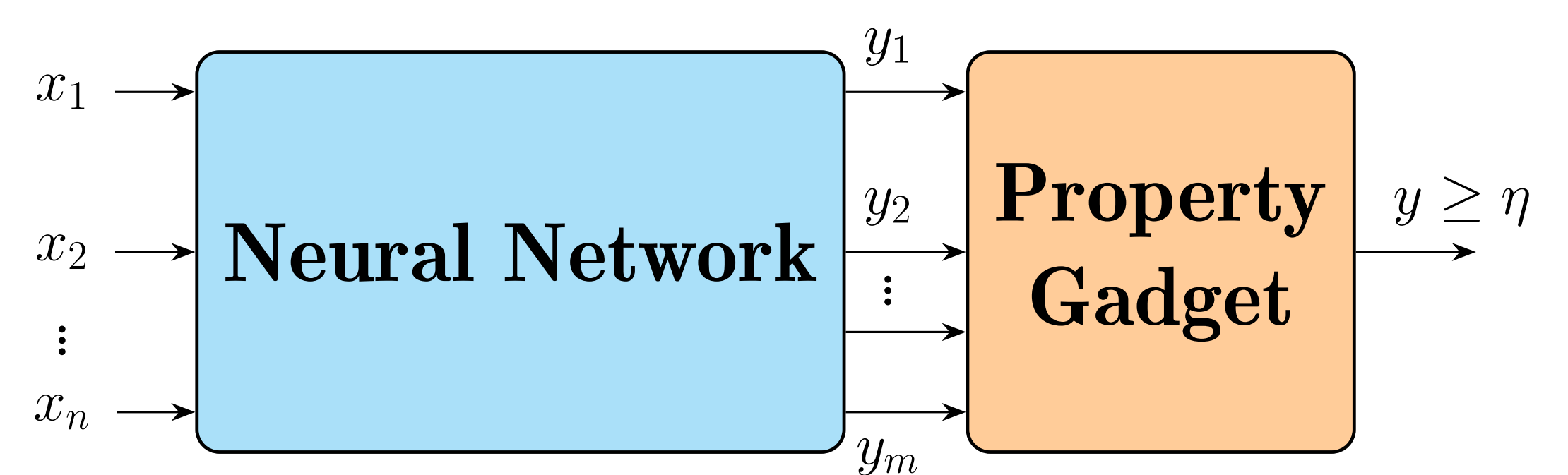
Let  $y_t = \max_{i=1, i \neq t}^m (y_i)$  and  $\delta = -\ln(\frac{100}{b} - 1)$ .

**Lemma.** If  $y_t \leq y_t + \delta$ , then  $CONF(\bar{y}, t) \leq b$ ; otherwise,  $CONF(\bar{y}, t) \geq \frac{100}{1+(m-1)e^{-\delta}}$ .

$$y_t \leq y_t + \delta \iff \bigvee_{j=1, j \neq t}^m y_t \leq y_j + \delta$$

Each post condition can only be encoded manually to constraints-based solver.

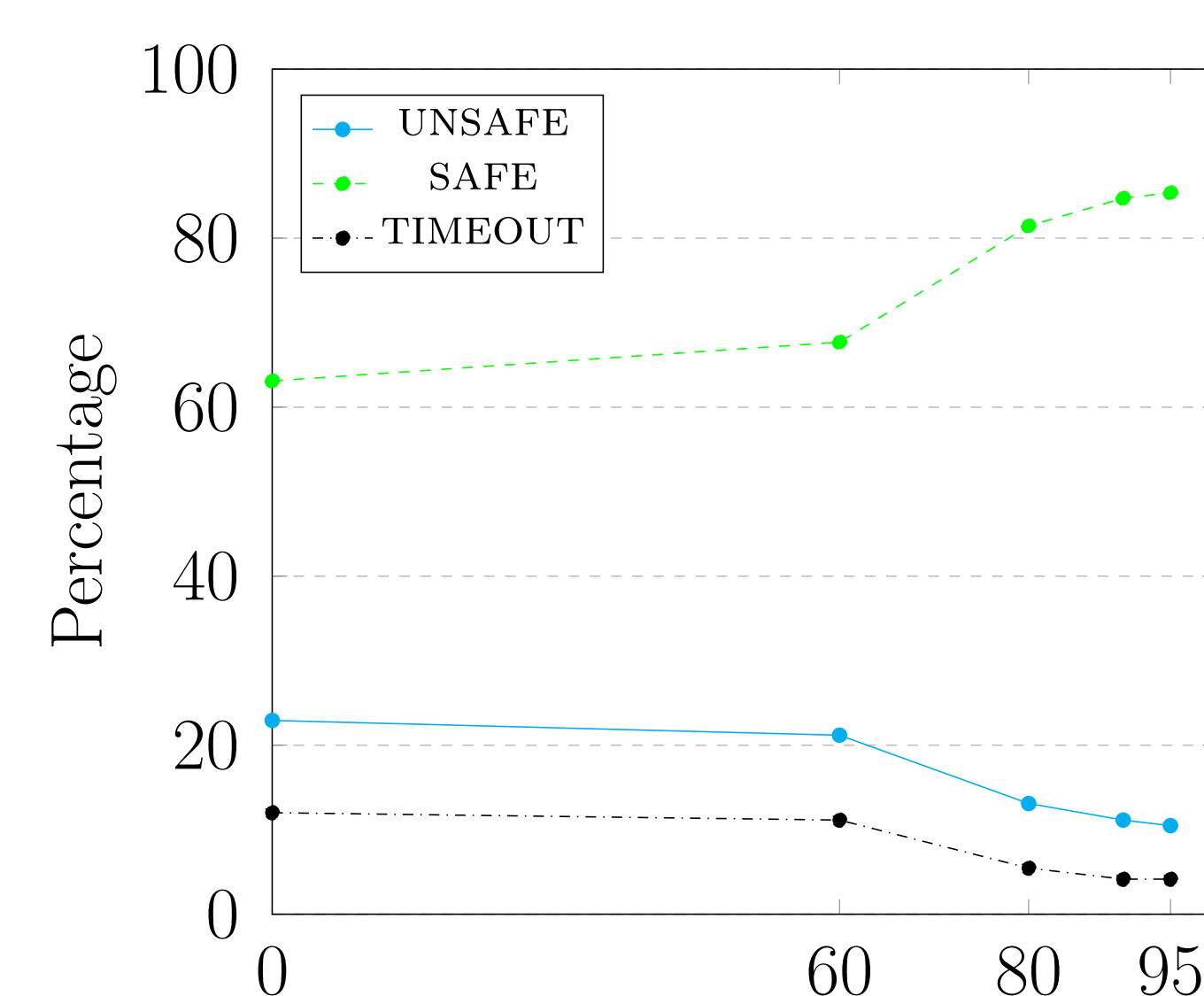
## 5. Layered Encoding



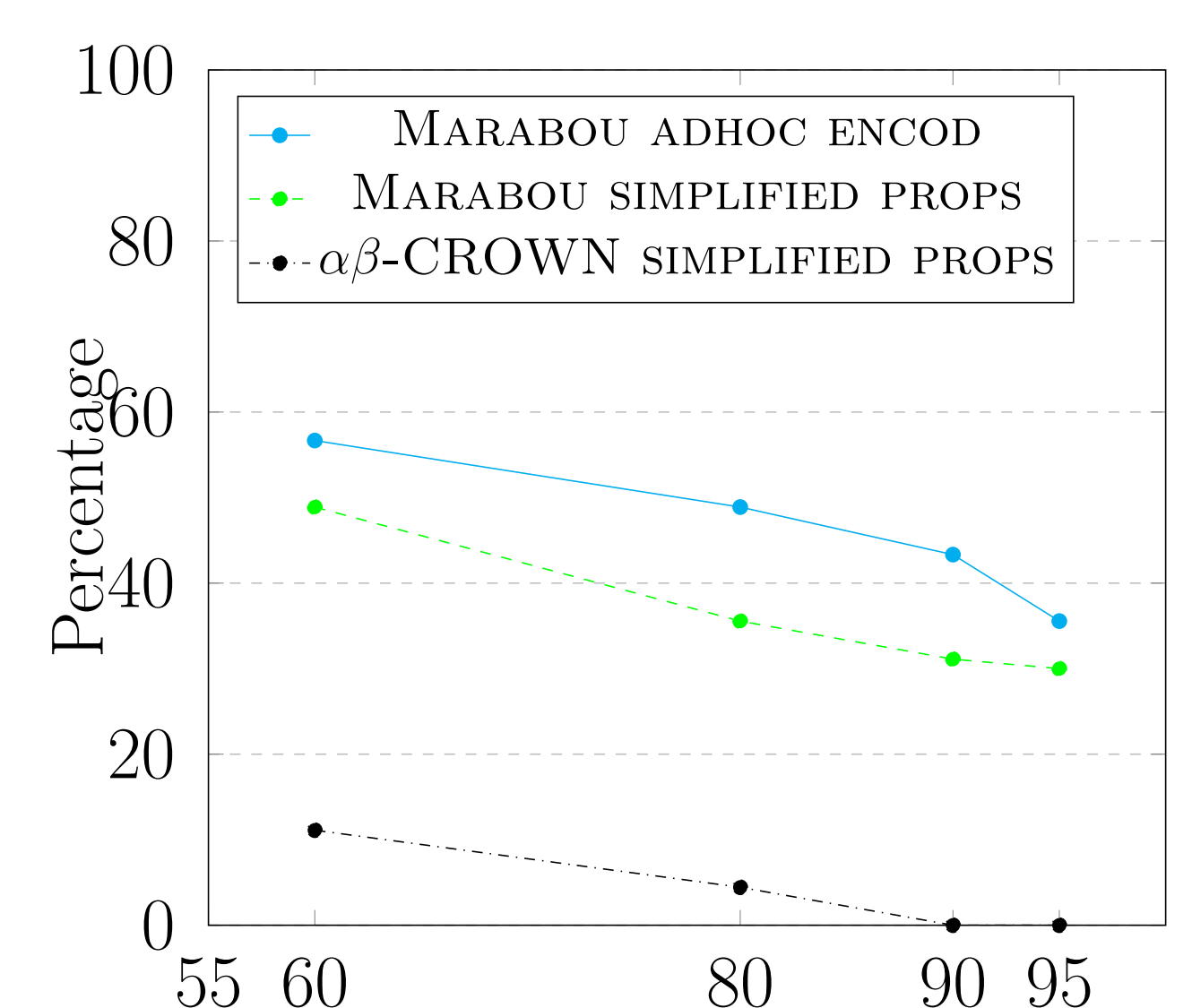
- Converts property to gadget and append to the network.
- Property gadget use only RELU activations.
- Property is simplified to  $y \geq \eta$ .
- Enables the use of SOTA, i.e.  $\alpha\beta$ -CROWN.

Property to gadget is a non-trivial conversion if  $\wedge$  and  $\vee$  both occurs in the property

## 6. Experiments



Scalability analysis



Comparison with adhoc encoding

- Datasets: MNIST, CIFAR-10, GTSRB, and Image-1k. Sourced from VNNCOMP.
- Network size: 512 to 11.16M relus.
- Tool:  $\alpha\beta$ -CROWN.

**Conclusion**

- Existing works ignore the classification confidence.
- Introduced novel definition of robustness – Relaxed robustness.
- We introduced a grammar that captures a wide range of robustness properties.
- Provided a layered encoding that enabled the use of SOTA.
- Our approach scales on large networks.

## References