# Building specifications for perception systems : formal proofs of deep networks trained with simulators

Julien Girard-Satabin (CEA LIST), Guillaume Charpiat (INRIA TAU),
Zakaria Chihani (CEA LIST), Marc Schoenauer (INRIA TAU)

25 février 2020

# Outline

Necessity to certify deep neural networks and challenges

    Glory and faults of deep learning software

    How to certify classical software ?

    Formal methods

    Challenges of deep neural networks verification

Deep learning verification : a review

Verification of perception models trained with simulators

Proof of concept and future works

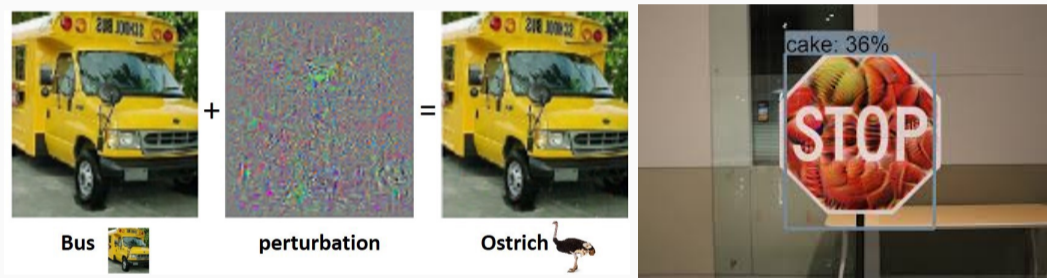# Necessity to certify deep neural networks and challenges

# Deep neural networks are awesome. . .

Active research community, profusion of tools, lot of industrial applications. . .
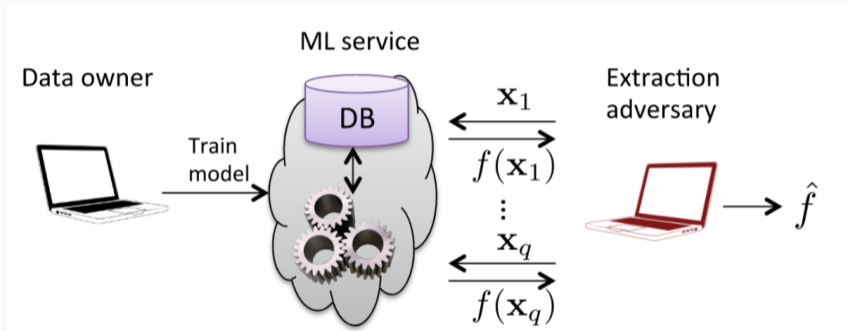
# Deep neural networks are awesome...

Active research community, profusion of tools, lot of industrial applications... ... yet they are not perfect

Bus  perturbation  Ostrich

Innocuous to humans, transferable between datasets, not systematic detection method

*A critical system is a system whose failure may cause physical harm, economical losses or damage the environment*

*A critical system is a system whose failure may cause physical harm, economical losses or damage the environment*

**Goal** : guarantee that the system respects a *safety specification*

$\phi$ : *an autonomous car will not run over pedestrians*

| Test on real environment | real conditions | cumbersome, potentially hazardous, non exhaustive |
|---|---|---|
| Test on virtual environment | can be automated, easy to integrate in existing workflow | non exhaustive, biased towards success |

And more (fuzzing...)

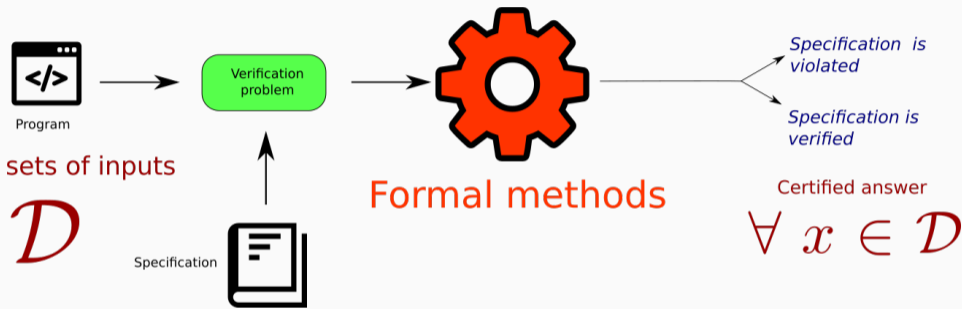Useful technique, widely used, enough for most of use cases

# Sometimes, tests alone are not enough !

| Claim | Discussion |
|-------|------------|
| "A car drove 5,472km, 99% in autonomous mode"[1] | If it translate to a failure rate, $10^{-2}$, insufficient compared to requirements in other critical systems (about $10^{-6}$ in aerospace) |
| "Our test cases are exhaustive" | Testing sets tend to be biased towards "normal" operation (accidents are rare)[2]  |

1. https ://www.wired.com/2015/04/delphi-autonomous-car-cross-country/
2. https ://arstechnica.com/cars/2019/05/feds-autopilot-was-active-during-deadly-march-tesla-crash/

- Studied in the academics since 1930 ($\lambda-$calculus, Church, Turing)
- Different techniques : abstract interpretation (Cousot and Cousot 1977), SAT/SMT (Davis and Putman 1960 ; Tinelli 2009), deductive verification (Coquand 1989), etc.
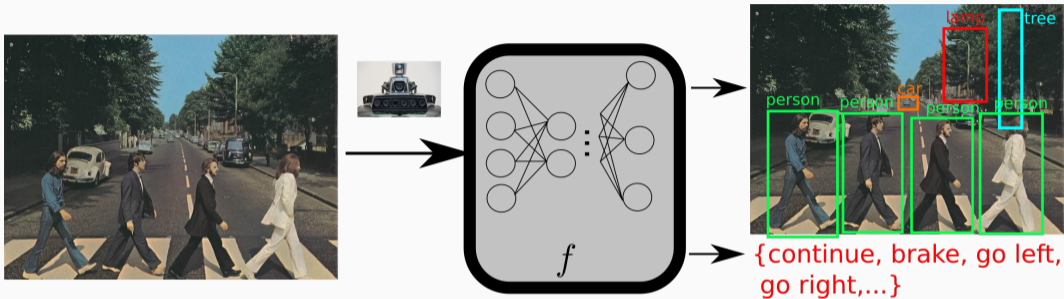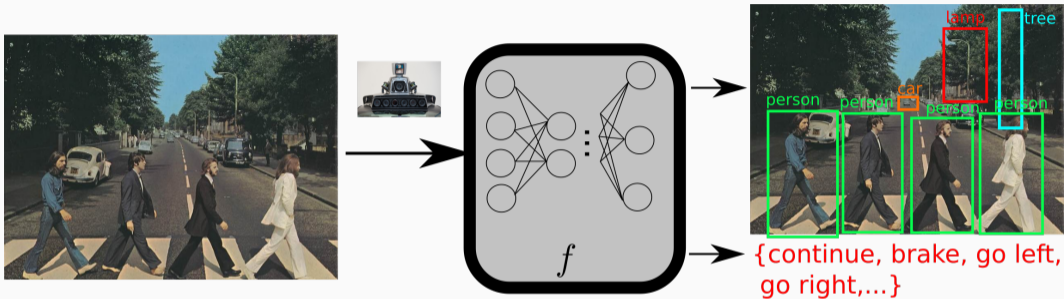- Used in industrial settings such as aerospace, automated transports, energy to *formally* certify

Work on domains $\mathcal{D}$ of inputs (global properties)

Answer is sound, formally guaranteed by mathematical logic

Dream property $\phi$ : *the autonomous car never run over pedestrians*

Dream property $\phi$ : *the autonomous car never run over pedestrians*

no *formal characterization* of what a pedestrian is !

Dream property $\phi$ : *the autonomous car never run over pedestrians*

no *formal characterization* of what a pedestrian is !

*Lack of formal definition on inputs prevents from formulating interesting safety properties*

# It's hard to use formal methods on deep learning

| Classical software | Machine learning |
|---|---|
| Explicit control flow | Generated control flow |
| Explicit specifications | Data-driven specifications (lack of generality) |
| Abstractions and well known concepts | Very few abstractions and reusability |
| Documented and understood vulnerabilites | Flaws without systematic characterization |

Some differences between classical software and machine learning

2 cases per ReLU node for the solvers

Several million ReLU nodes $\rightarrow$ $2^{O(10^6)}$ case splits

Combinatory explosion (if done naively)

# Deep learning verification : a review

For a *given* input $x$, a classification function $f$, an adversarial perturbation $\delta$ :

find delta satisfying

classifier misclassification

such that

perturbation stays below a certain threshold

For a *given* input x, a classification function $f$, an adversarial perturbation $\delta$ :
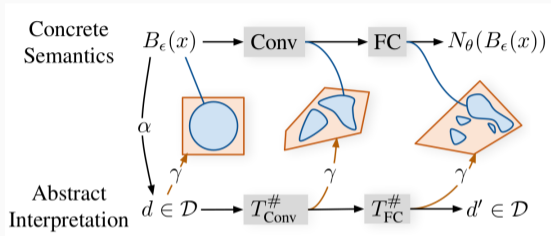
find delta
satisfying
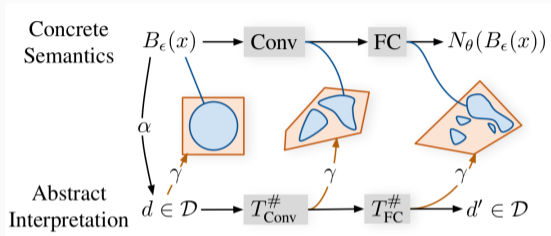
$$f(x) \neq f(x + \delta)$$

such that

$$\|\delta\|_p \leq \varepsilon$$

1. *abstract* the network

2. *propagate* perturbations

3. *assess* robustness properties
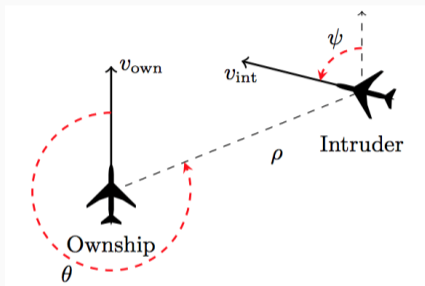
4. *learn to minimize* adversarial loss



Improve adversarial robustness on 100 samples from CIFAR-10 from 0 to 80%, $\varepsilon = 8/255$, 3 hidden layers, convolutional network

1. *abstract* the network

2. *propagate* perturbations

3. *assess* robustness properties
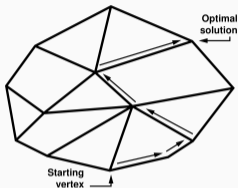
4. *learn to minimize* adversarial loss



Improve adversarial robustness on 100 samples from CIFAR-10 from 0 to 80%, $\varepsilon = 8/255$, 3 hidden layers, convolutional network

## Scalable verification, but local properties

*If the intruder is distant and is significantly slower than the ownship, the score of a COC advisory will always be below a certain fixed threshold. Bounds : $\rho \geq 55947.691$, $v_{own} \geq 1145$, $v_{int} \leq 60$*

Critical system

Core of most SMT solvers working with number values

Modified to lazily evaluate ReLUs

Exact verification of several properties on a ACAS-Xu implementation

## Global properties
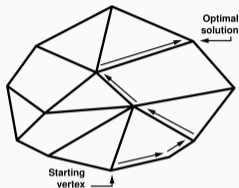
Core of most SMT solvers working with number values

Modified to lazily evaluate ReLUs

Exact verification of several properties on a ACAS-Xu implementation

# Global properties

Assumes perfect plane detection beforehand

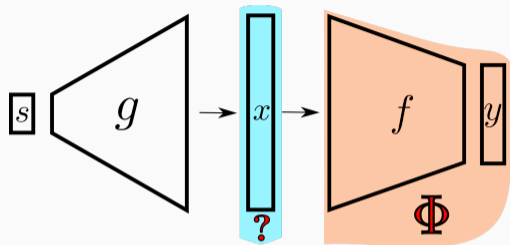How do we verify perception ? What is an intruder ?

# Verification of perception models trained with simulators

# Example of simulator

Industry relies more and more on simulators to generate scenarios to train and evaluate deep learning models



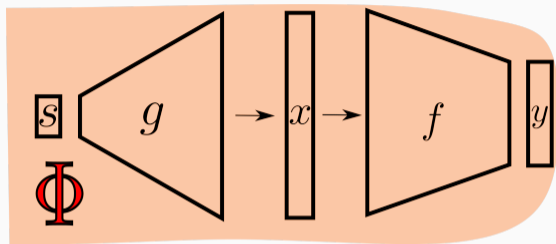Screenshot from the CARLA open source simulator

# Simulators as data providers



- $s$ : parameters (obstacles, weather conditions...)

- $g$ : simulator

- $x$ : perceptual input (images)

- $f$ : model

- $y$ : decision output (brake...)

- $\phi$ : "$\forall\ x$ that contains a pedestrian, do not roll over it"
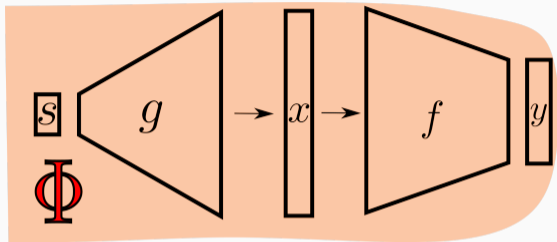
How to formulate $\phi$? What is an image $x$ with a pedestrian?

Modify the verification problem formulation to include $g$ and $s$

$\phi$ now encompasses $s$ and can now be expressed : *For all values of s that are translated by g as the presence of pedestrians into x, do not run over those pedestrians*
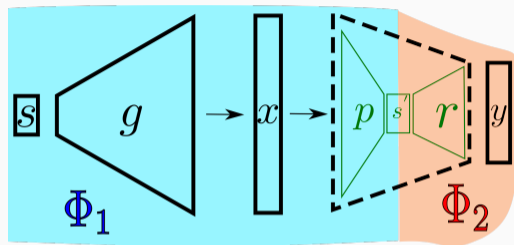
Modify the verification problem formulation to include *g* and *s*

*φ* now encompasses *s* and can now be expressed : *For all values of s that are translated by g as the presence of pedestrians into x, do not run over those pedestrians*

We now have a property to verify a perceptive unit !

$f$ splits in *p*erception and *r*easoning, $p$ learns $s$
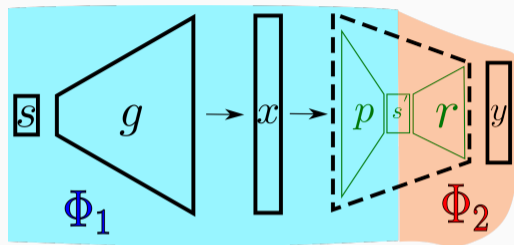
$\phi_1$ on $p$ : guarantee of no information loss : *reconstruct s from x*
$s^{'} = s \;\forall\; s \rightarrow p \circ g = Id$

$\phi_2$ on $r$ : do not kill pedestrians (assuming perfect perception)

# Refinement : splitting perception and reasoning



$f$ splits in *p*erception and *r*easoning, $p$ learns $s$

$\phi_1$ on $p$ : guarantee of controlled information loss : *reconstruct $s$ from $x$*
$s^{'} \simeq s \ \forall \ s \rightarrow p \circ ||g - id < \varepsilon||$

$\phi_2$ on $r$ : do not kill pedestrians (assuming perfect perception)

How to express $\phi$, $g$, $f$, $\mathcal{X}$, $\mathcal{Y}$, $\mathcal{S}$ ?

How to express $\phi$, $g$, $f$, $\mathcal{X}$, $\mathcal{Y}$, $\mathcal{S}$ ?



SMTLIB : Tinelli et al., 2017, https ://onnx.ai/

High-level workflow

From all mainstreams deep learning frameworks to all mainstreams SMT solvers

# Proof of concept and future works

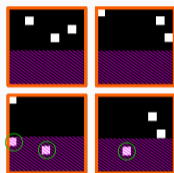Train a simple model to output a single command directive if a simplified input is in a pre-defined danger zone

Train a simple model to output a single command directive if a simplified input is in a pre-defined danger zone



$s =$ (position of obstacles)

$x$

output scalar (obstacle detected if $> 0$)

$y$

Network has 16 neurons, 2 hidden layers

Train a simple model to output a single command directive if a simplified input is in a pre-defined danger zone



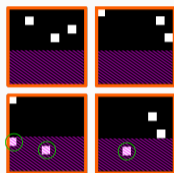$s = $ (position of obstacles)

output scalar (obstacle detected if $> 0$)

$y$

$x$

Network has 16 neurons, 2 hidden layers

We prove the given trained network will *never fail*

```
            .  .  .
( declare −fun  l_1 . weight38  ()  Real )
( assert  (=  l_1 . weight38  (/  (−  13947381)  120892581961462917470G176)))
( declare −fun  l_1 . weight37  ()  Real )
( assert  (=  l_1 . weight37  (/  (−  405697  )  2251799813685248)))
( declare −fun  l_1 . weight36  ()  Real )
            .  .  .
```

Part of the network's translation

1. Include noise and incomplete reconstruction in the framework

2. Add rewriting rules

3. Release and enhance the toolkit

4. Add a systematic representation of the simulator

5. Integration of state-of-the art verification tools

# Any questions ?

contact for the paper : julien.girard2@cea.fr