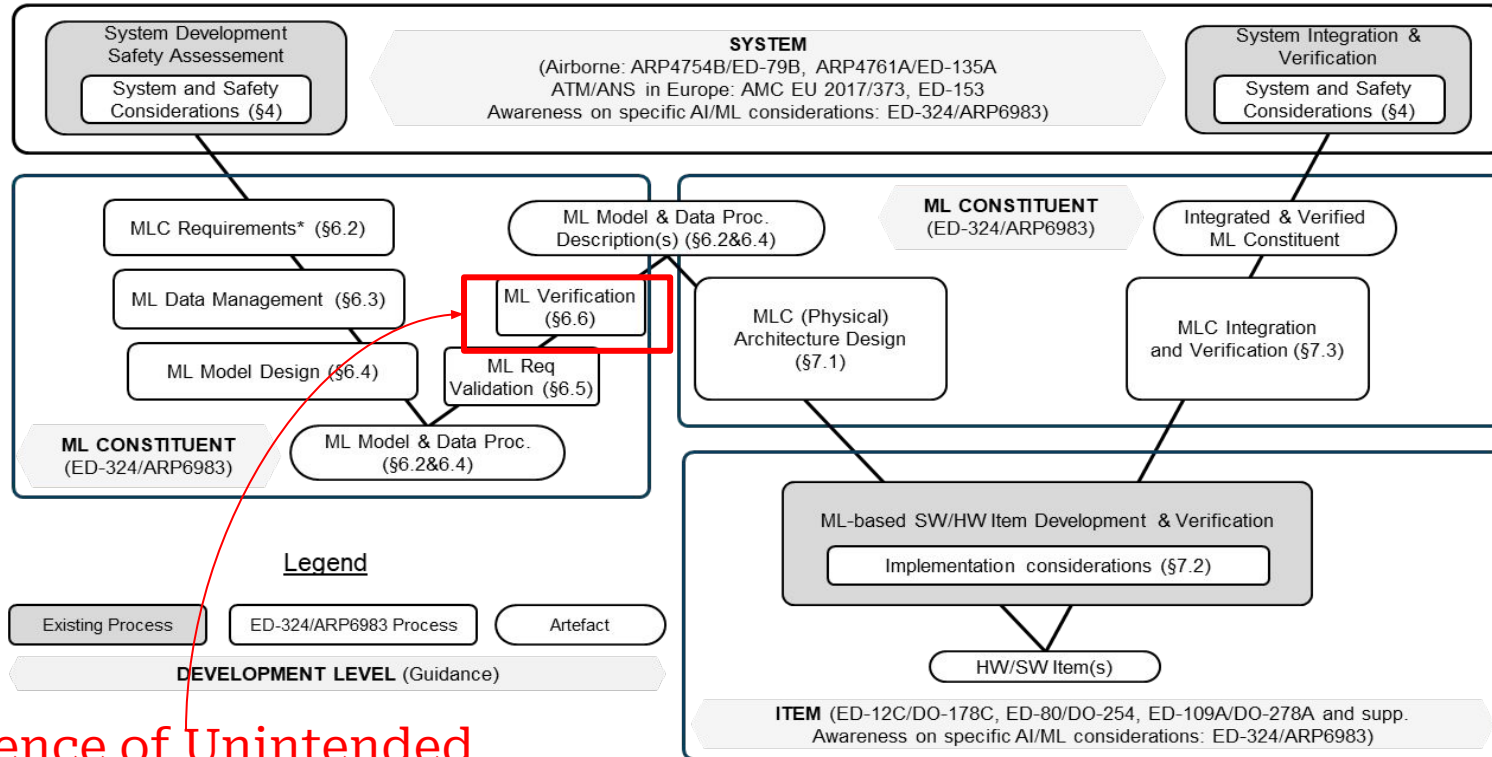# Robustness Verification: Neural Network's Surrogate

Melanie Ducoffe
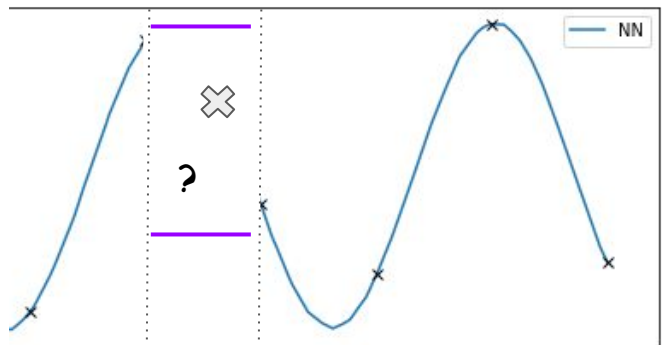
ANITI DAYS
November 2024
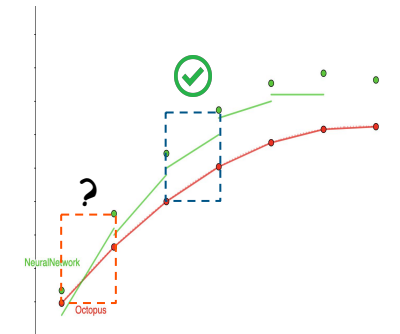
**AIRBUS**



Absence of Unintended Functionality

*ForMuLA: Formal Methods Use for Learning Assurance – EASA & Collins Aerospace partnership*

# Property Requirement for Surrogate Models

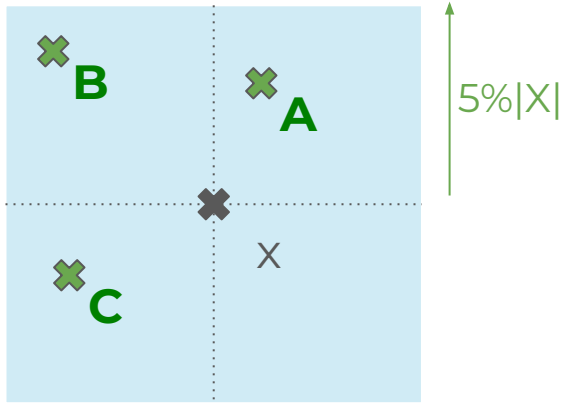$$\max_{x \in \Omega} |\, f(x) - f(x_0)|$$



SAFE SURROGATE    UNSAFE SURROGATE

$$\min_{x \in \Omega} f(x) - f(x_0) \geq 0$$

## Partial Input Monotony



$$x_1 \qquad\qquad x_2 \qquad\qquad f(x_1) \qquad f(x_2)$$
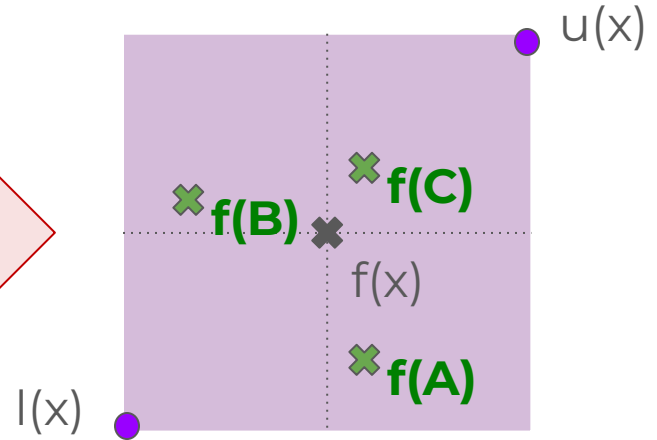
$$\begin{pmatrix} \text{speed} \\ \text{weight} \\ \text{dry runway} \end{pmatrix} \begin{matrix} = \\ = \\ < \end{matrix} \begin{pmatrix} \text{speed} \\ \text{weight} \\ \text{wet runway} \end{pmatrix} \implies \text{BDE}_1 < \text{BDE}_2$$

**AIRBUS**

# INPUT DOMAIN Ω

# OUTPUT DOMAIN

B

A

$5\%|X|$
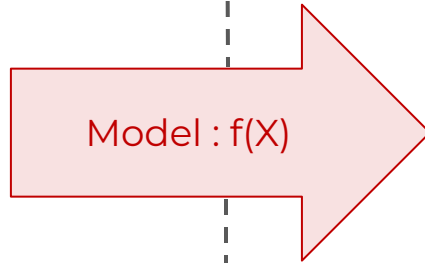
x

C

Model : f(X)

u(x)

f(B)

f(C)

f(x)

f(A)

l(x)

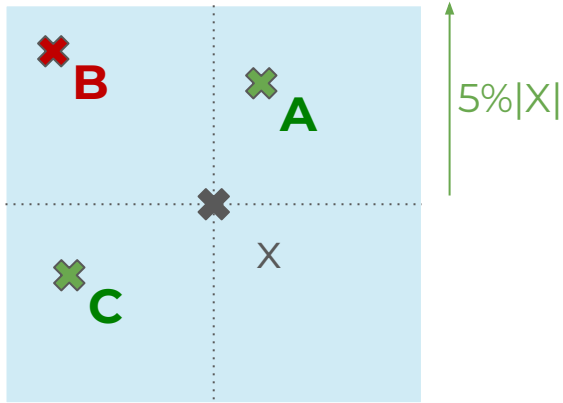The property is **verified** !

$$\max_{z \in \Omega} g(Z;X) \leq 0$$

$g(Z;X) = \max_i \max( f_i(Z) - u_i(X),$
$l_i(X) - f_i(Z) )$

4

**AIRBUS**

## INPUT DOMAIN Ω
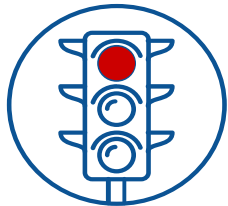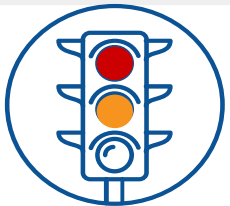
## OUTPUT DOMAIN



$5\%|X|$

Model : f(X)

$y_{max} = u(x)$

$f(C)$

$f(B)$

$f(x)$

$f(A)$

$y_{min} = l(x)$

$\max_{Z \in \Omega} g(Z; X) > 0$

The property is **violated**

$g(Z; X) = \max_i \max( f_i(Z) - u_i(X),$
$l_i(X) - f_i(Z) )$

5

**AIRBUS**

$$\max_{Z \in \Omega} g(Z;X) > 0 \Rightarrow \exists\ Z \in \Omega \text{ s.t } g(Z; X) > 0$$



Direction of gradients with respect to loss (training in action)

Direction of gradients which maximizes the loss (attack)

Model : f(X)

u(x)

l(x)

f(x)

$f(x_1)$

$f(x_2)$

$f(x_n)$

$f(x_0)$

$X_n$  $X_2$  $X_1$  $X_0$

x

Loss

Shows model vulnerabilities but not their absence
>> Do not provide property verification guarantee.

Naïve attacks schemes can be used for regression (FGSM, PGD)

# Casting Local Verification as a Classification Property

$s_0$ $\quad = \min(y - y_{min}, \ y_{max} - y)$

$y$

$s_1$ $\quad = y_{min} - y$

$s_2$ $\quad = y - y_{max}$

$$\max_{Z \in \Omega} g(Z;X) \leq 0$$

$$\text{argmax}_{Z \in \Omega} s(Z;X) = 0$$

**FGSM / PGD**
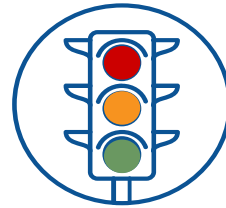
$N_{iter}$ **gradient ascent** on the loss function in x.

**(SUPER) DEEPFOOL (S-DF)**

iterative **projection** on the closest linearized hyperplane boundary.

Low          Runtime

**AUTO ATTACK (AA)**          High          **CARLINI & WAGNER (C&W)**

fixed **combination** of 3 attacks (AutoPGD, FAB, Square) with different losses

**Targeted** attack based on the logits.

**FOOLING RATE (success rate)**

FGSM          PGD          DF/SDF          AA

0%          Carlini

**AIRBUS**
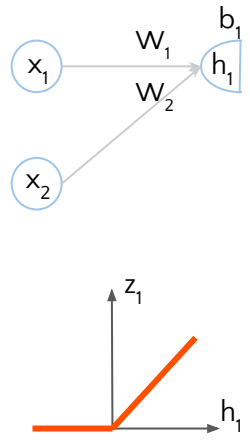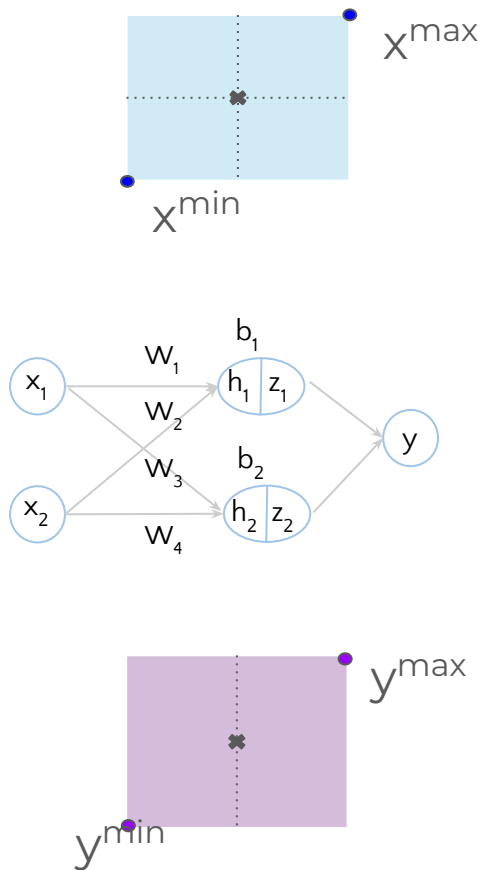


$$\max_{z \in \Omega} g(Z;X) > 0$$

+ convergence to the true optimum that implies robustness or non robustness.
+ Not scalable to larger network
  SMT–solver [Marabou]
  Lipschitz optimization **(Paul Novello)**
  **Mixed Integer Programming (VENUS)**

No magical trick:
white box setting

8

$x^{max}$

$x^{min}$

$x_1$ $W_1$ $W_2$ $b_1$ $h_1$ $z_1$

$b_2$ $h_2$ $z_2$ $W_3$ $W_4$ $x_2$ $y$

$x_1$ $W_1$ $W_2$ $b_1$ $h_1$ $x_2$

$z_1$ $h_1$

$y^{max}$

$y^{min}$

## Bounding the input perturbation

$$x_i^{min} \leq x_i \leq x_i^{max}$$

## Encoding Neural Network

$$h_1 = w_1.x_1 + w_2.x_2 + b_1$$
$$h_2 = w_3.x_1 + w_4.x_2 + b_2$$

### Big M encoding for ReLU

$$z_1 = max(0,h_1)$$
$$z_2 = max(0,h_2)$$

$\equiv$

$z_i \geq 0 \ldots$
$z_i \geq h_i$
$\delta_i \in \{0,1\}^{|xi|}$

$z_i \leq u_i.\delta_i$
$z_i \leq h_i - (1-\delta_i)$

## Encoding property violation (SAT)

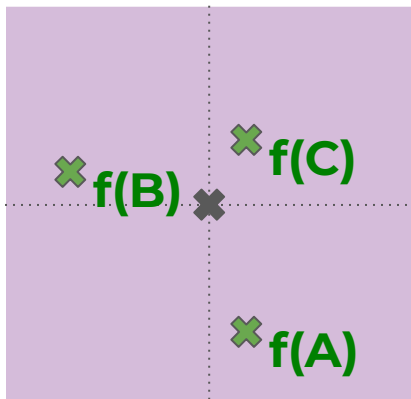$$y_1^{max} < y_1 \text{ or } y_1 < y_1^{min} \text{ etc.}$$

9

Build Under/Over approximation of f

$$\forall z \in \underline{\Omega}\ \underline{f}(z) \leq f(z) \leq \overline{f}(z)$$

Use it for dominating g

$$\forall z \in \Omega\ g(z) \leq \overline{g}(z)$$

$$\max_{z \in \Omega} \overline{g}(Z;X) \leq 0 \Rightarrow \max_{z \in \Omega} g(Z;X) \leq 0$$
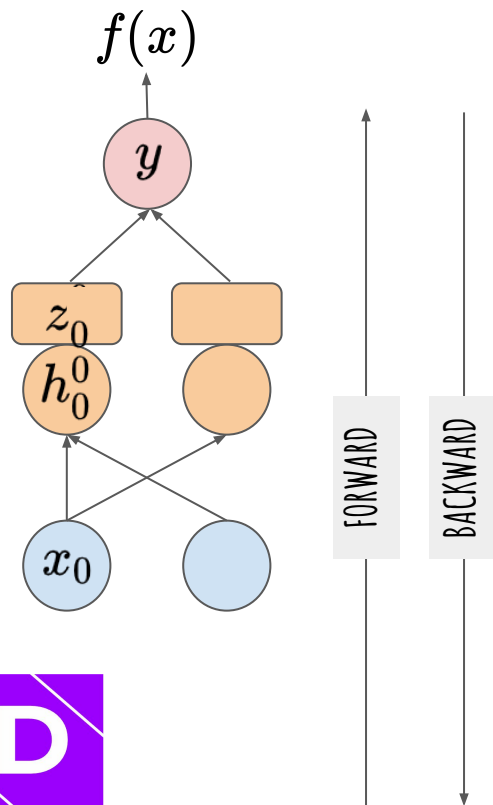
⚠ outer-approximations that only implies robustness:

**Linear Relaxation** [CROWN]
Convex Relaxation [SDP]

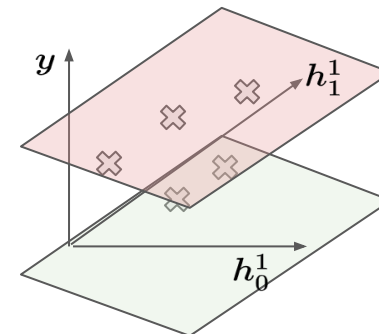$\overline{f}(B)$  $\overline{f}(C)$
$f(C)$
$f(B)$
$\underline{f}(B)$  $\overline{f}(A)$
$f(A)$

$\overline{f}(C)$
$f(C)$
$f(B)$
$f(A)$

# Verification as a Relaxed Optimization problem: LIRPA

$f(x)$

$y$

$z_0$

$h_0^0$

$x_0$

FORWARD

BACKWARD

| VERIFICATION 'FORWARD–FEED' | VERIFICATION 'BACKWARD–FEED' |
|---|---|
| $h_0^1$    $x_1$    $x_0$ | $y$    $h_1^1$    $h_0^1$ |
| +   Self-sufficient | +   Not self-sufficient (pre-processing) |
| +   Complexity = cost of inference | +   Complexity = cost of backpropagation at best (gradient) |
| -   Less accurate | -   More accurate |
| -   Not scalable on large images | -   Not scalable on large outputs (GAN) |

github.com/airbus/decomon

*A Convex Relaxation Barrier to Tigh Robustness Verification of NNs, Salman et al.*

**AIRBUS**

## Aircraft Loads–to–Stress Prediction



**Loads**

216 inputs

**Predicted Stress**

81 outputs



'Knot'

'Wing'

$\Delta$ prediction $f_i(x) - f_i(x)$

(Normalised) model prediction $f_i(x)$

- Model– **Research prototypes**:
  - Two hidden layers (165 neurons)
  - ReLu activation functions
  - Dense output layer (81)
- Test data: 1000 loads/stress points

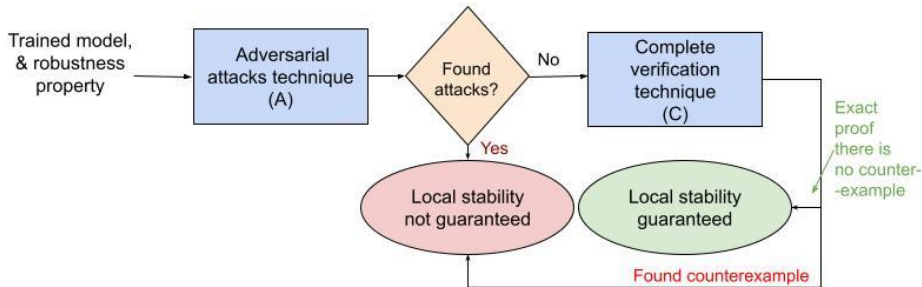# Verification approach – Combinaison



A+B+C



A+C

A: PGD (Cleverhans)
B: CROWN (Decomon)
C: MILP (Gurobi)

# Results

| | (1) A | (2) B | (3) C | (4) A+C | (5) B+C | (6) Pipeline A+B+C |
|---|---|---|---|---|---|---|
| #Tested | 1000 | 1000 | 1000 | 1000/558 | 1000/446 | 1000/558/4 |
| #True | - | 554 | 558 | 558 | 558 | -/554/4 = 558 |
| #False | 442 | - | 442 | 442 | 442 | 442/-/0 = 442 |
| **Runtime** | **10.7** | **3.3** | **267** | **19.8** | **267** | **10.7/1.96/3.91 = 16.6** |

~45% of test data are shown to be non locally stable

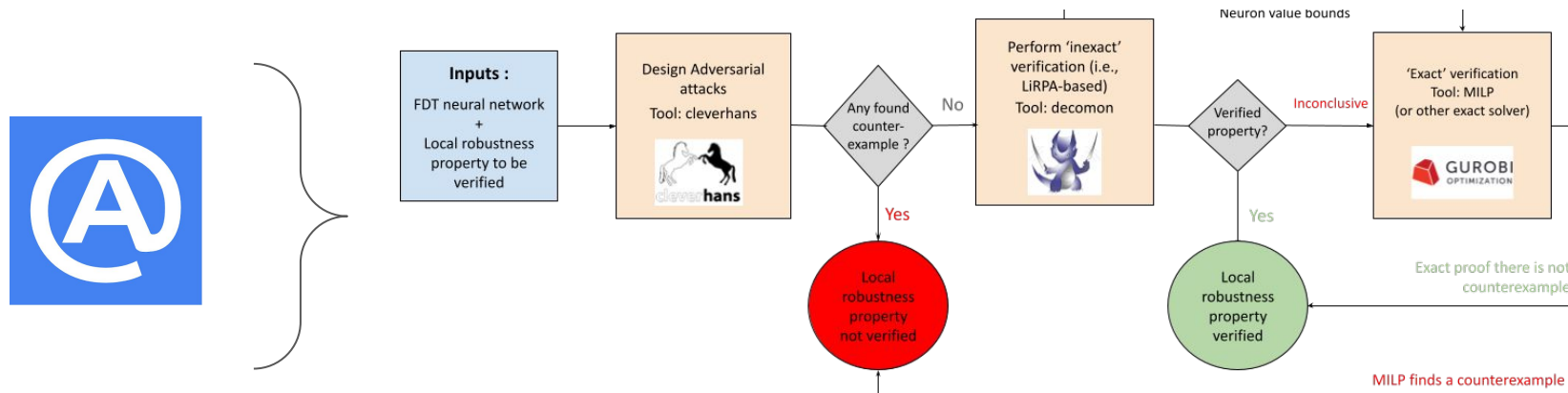The "Adversarial attacks" step was able to find all non-stabilities

Low number of remaining test data to be evaluated by "C" after (A or A+B)

Significant decrease in computational time

Open source library: Airobas

Stability accuracy can be efficiently measured with a verification pipeline



Current models have **deceiving stability accuracy: ~40%**. What tools are at our disposal ?

1) XAI actionability: Reducing the problem complexity (input and output dimensions)
2) **Regularizing the training to balance between good <u>regression performance</u> and good <u>stability accuracy</u>**

**AIRBUS**

Robust ═ Accurate ✚ Stable

### Data Augmentation

Artificially increase the size of the dataset by applying domain-specific transformations on the input and output data. It introduces stability invariance.

### Weights Constraints

Weights constraints limit the Lipschitz constant in a neural network.It is known to increase the model's resilience against adversarial attacks or input perturbations by limiting the model's capacity to fit noise
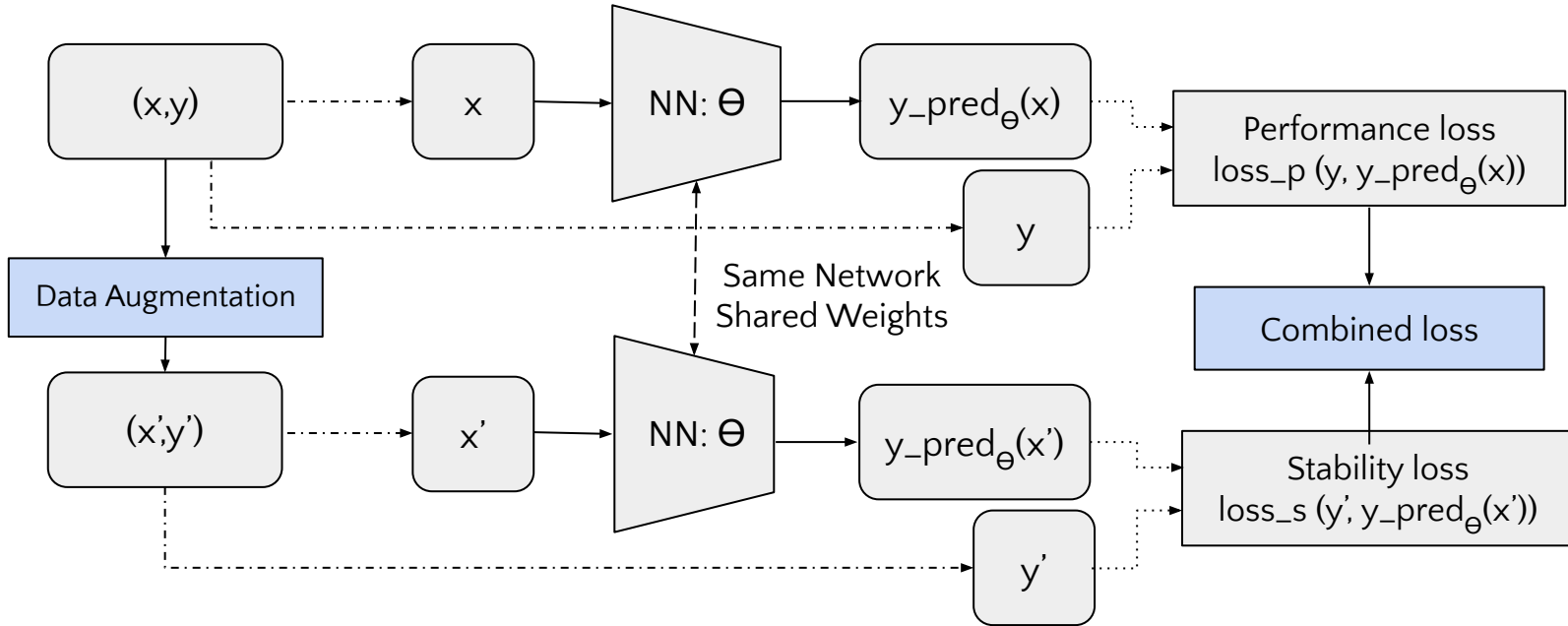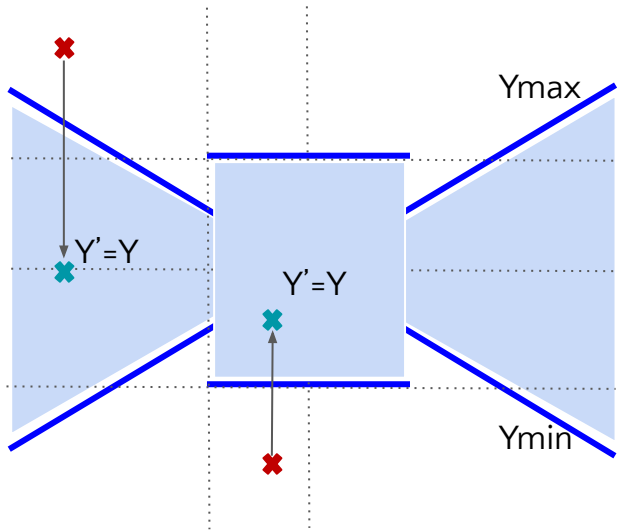
### Certified Training (Meta Networks)

Certified training use Incomplete Formal Methods as a Meta Model to provide formal guarantees about a model's robustness against domain-specific perturbations.

17

**AIRBUS**

## Data Augmentation

Artificially increase the size of the dataset by applying domain–specific transformations on the input and output data. It introduces stability invariance.
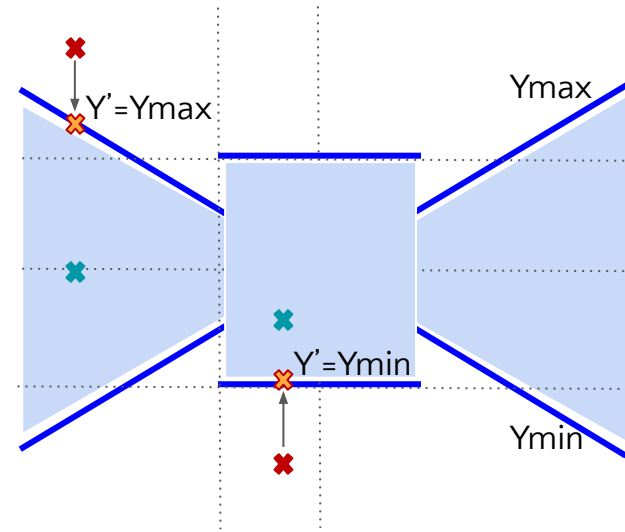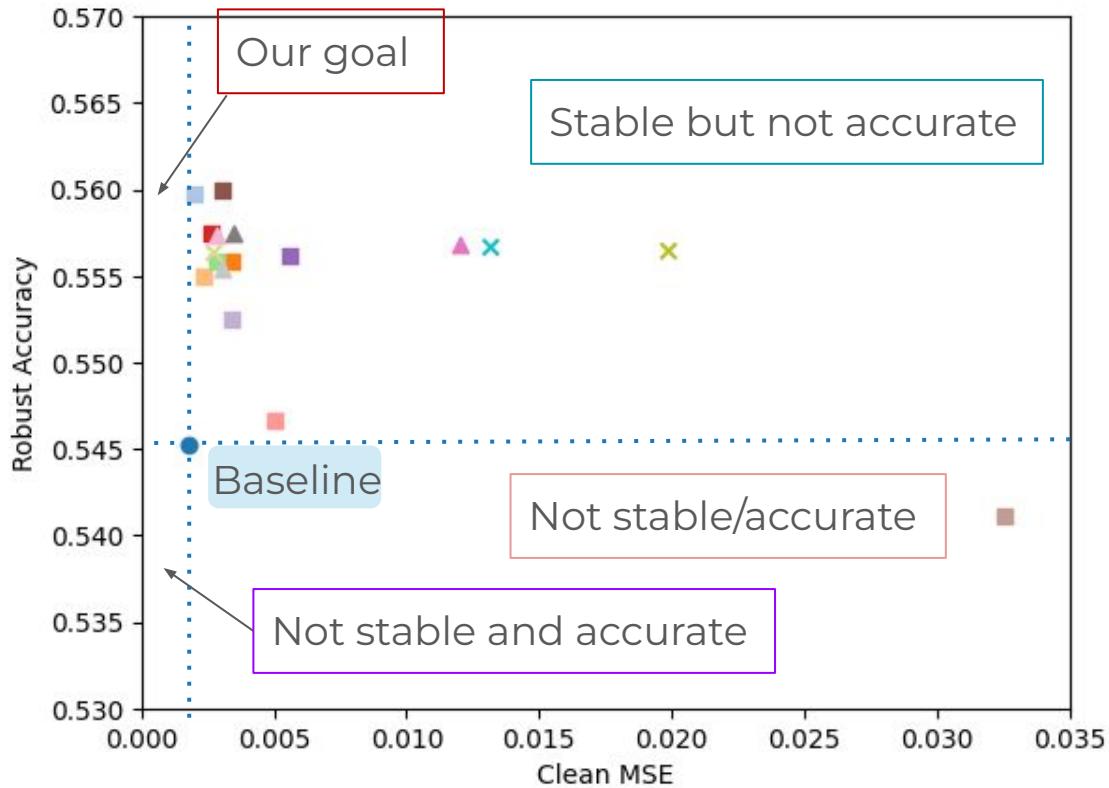
**AIRBUS**



**Groundtruth**
use the groundtruth label of
the initial input

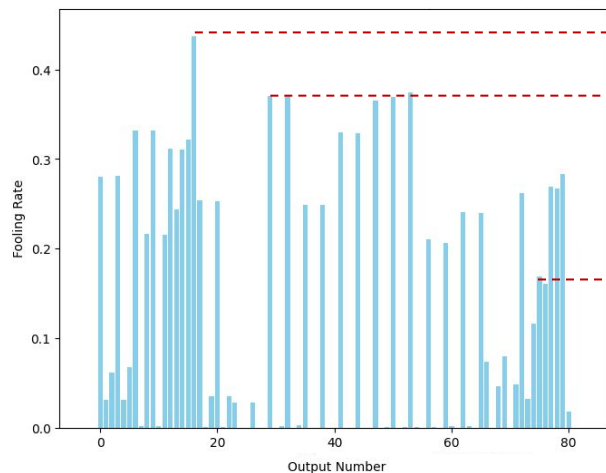**Stability clipping**
clip the prediction of the
adversarial input to lie within
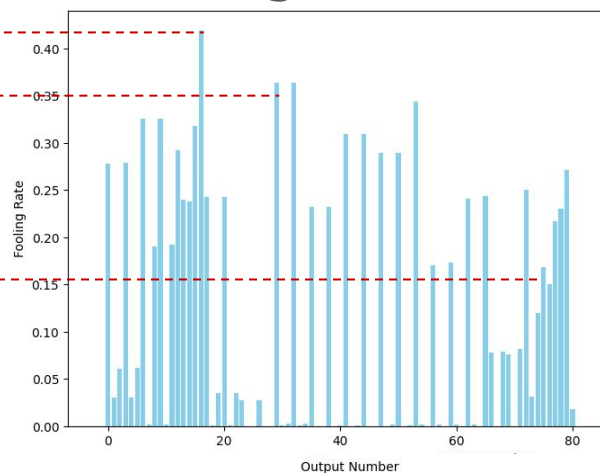the stability bounds
(Ymin, Ymax)

20

**AIRBUS**

Fooling Rate = 1 - Robust_accuracy
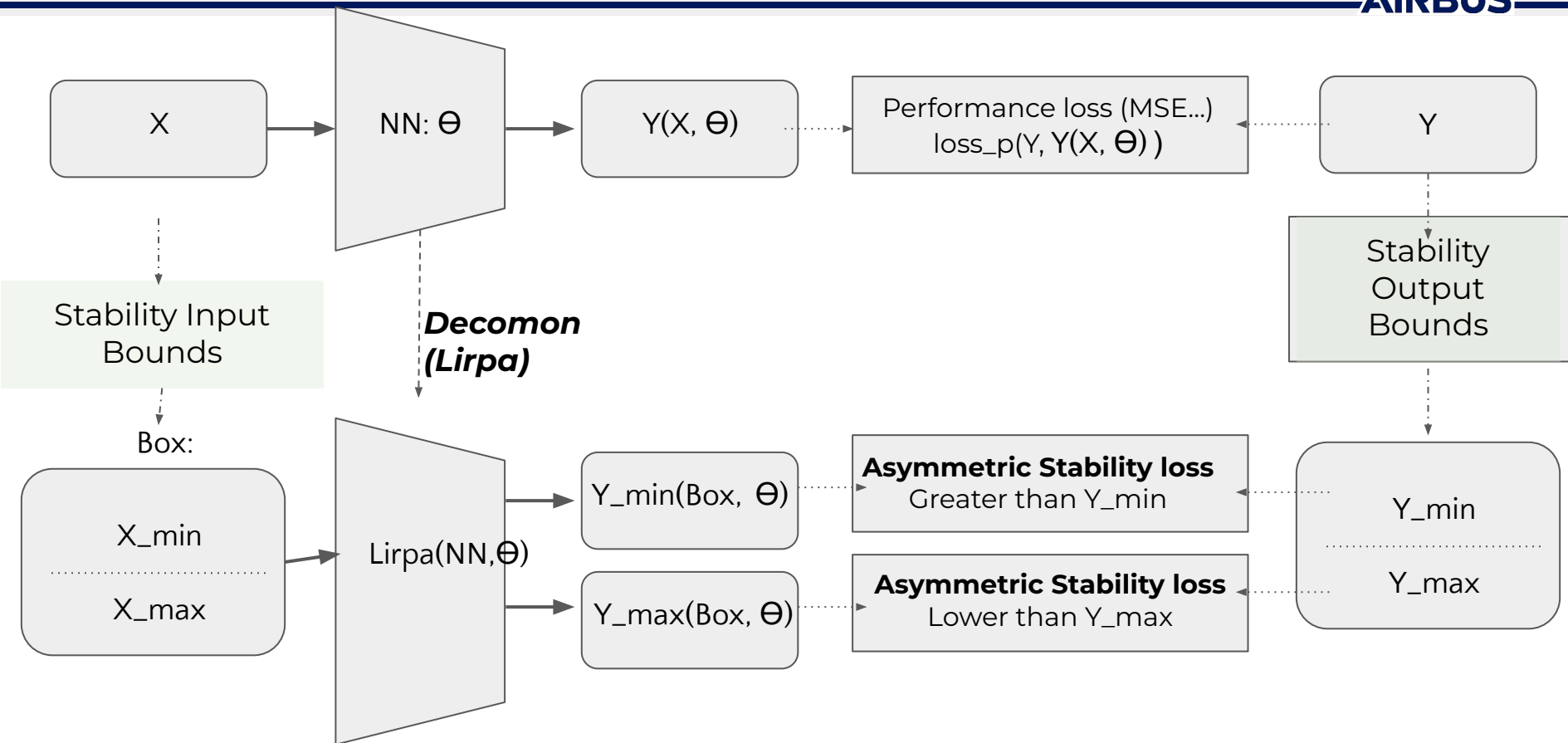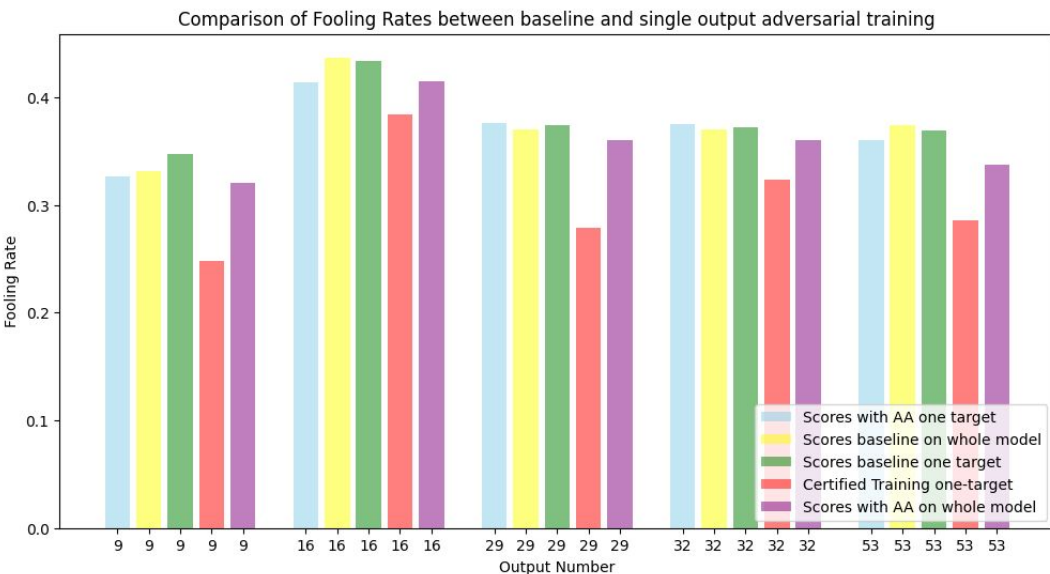
## Baseline                    ## Data Augmentation



- Most of the outputs are **naturally Robust**
- About **40%** of outputs are **problematic**

→ Can we **target** those outputs ?
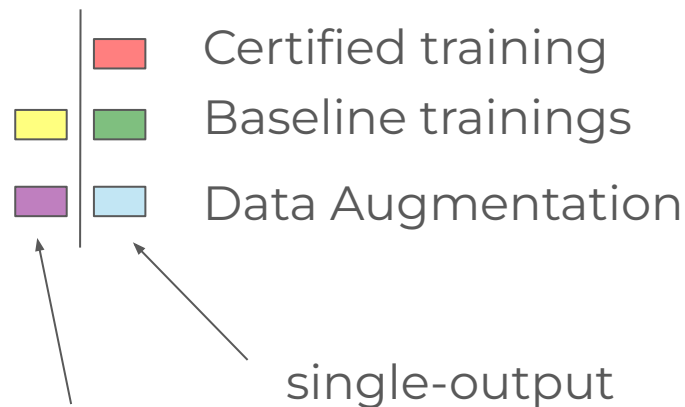
22

**AIRBUS**

**Certified Training
(Meta Networks)**

Certified training use <u>Incomplete Formal Methods</u> as a Meta Model to provide formal guarantees about a model's robustness against <u>domain-specific perturbations</u>.

**AIRBUS**



Comparison of Fooling Rates between baseline and single output adversarial training

- **CT for single-output models** for the 5 problematic outputs

  Certified training

  Baseline trainings

  Data Augmentation

  single-output

multi-output

Promising results ! 5-10% drop in the fooling rates compared to the previous models

*Surrogate Neural Networks Local Stability for Aircraft Predictive Maintenance, FMICS 2024*



Thomas Deltort    Ryma Boumazouza    Guillaume Poveda    Marion Cécile Martin    Audrey Galametz

ANITI EVENT: Hands on Verification
6th March 2025



https://github.com/airbus/Airobas

https://github.com/airbus/decomon