# *Leveraging Verification to Enhance Formal Explainable AI for Neural Networks*

**Presented by:**

**Ryma Boumazouza**

**Work in collaboration with:**

**Mélanie Ducoffe**

**Raya Elsaleh**

**Shahaf Bassan**

**Guy Katz**

AIR — Airbus AI Research | ANITI | DEEL — Dependable, Explainable & Embeddable Learning — France-Quebec | THE HEBREW UNIVERSITY OF JERUSALEM

1

# Delivering trustworthy AI through XAI



### Need for Trustworthy AI in High-Risk Settings

Ensure safety, compliance, and ethics in critical sectors (e.g., healthcare, aeronautics, finance, autonomous vehicles)



### Guidelines and Regulations Driving Trustworthiness:

EU, OECD, and UNESCO guidelines emphasize the need for AI to be trustworthy, transparent, accountable, and ethically and legally sound



### Challenges with Current eXplainable AI (XAI) Approaches

Scalability limits, high complexity, and difficulty integrating into existing AI systems
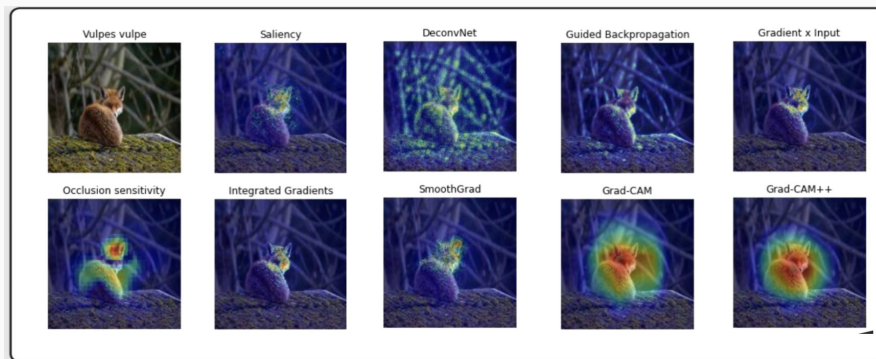
# Can we Truly trust XAI Tools ?

XAI tools promise transparency but…

- are often heuristic
- do not provide guarantees



But,
If we can't trust the explainer, can we trust the model ?

➡ **What about a Formal Explanation?**

# *What Does a Formal Explanation Look Like?*



$$N\left(\ \ \ \ \right) \overset{?}{=} \text{cat}$$

We want an explanation to answer "why" the classifier predicted "cat"

A Sufficient reason/
Abductive explanation would be :

# *Formalizing the Concept of an Explanation*

Formally, an abductive explanation is defined as:

$$\forall (x \in \mathbb{F}). \quad \left[ \bigwedge_{i \in E} (x_i = v_i) \to (N(x) = c) \right]$$

Properties of an abductive explanation:

Minimality:

$$N(\phantom{xxxx}) \neq \text{cat}$$

Sufficiency:

$$N(\phantom{xxxx}) = \text{cat}$$

# *How to compute such abductive explanation ?*

**Algorithm 1** Deletion Algorithm to Find One Abductive Explanation

1: **Input:** Predictive model $M$ and input $x = \langle \chi^1, \ldots, \chi^n \rangle$
2: **Output:** Explanation for class $C$ of $x$
3: $Explanation \leftarrow \emptyset$                                          ▷ Set of relevant features
4: $I \leftarrow \emptyset$                                                          ▷ Set of irrelevant features
5: $\pi \leftarrow TRAVERSALORDER(F)$              ▷ Sort F's features by ascending relevance
6: **for** each feature $x_i \in \pi$ **do**
7:     $\pi' \leftarrow \pi \setminus \{x_i\}$                                          ▷ Removing feature $x_i$
8:     **if** CHECK$(M, \pi')$ **then**                         ▷ Check if prediction changes
9:         $Explanation \leftarrow Explanation \cup \{x_i\}$   ▷ Add feature $x_i$ to relevant features
10:     **else**
11:         $I \leftarrow I \cup \{x_i\}$                            ▷ Add feature $x_i$ to irrelevant features
12:     **end if**
13: **end for**
14: **Return** $Explanation$

## Challenges to address

➤ Challenge 1: The CHECK function is computationally expensive

➤ Challenge 2: Sequential traversal of the feature set (loop)

➤ Challenge 3: Impact of order on explanation interpretability (size)
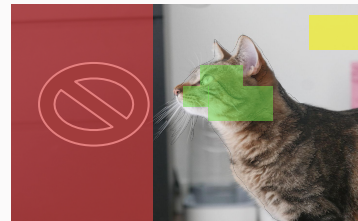
# *About the CHECK method (Verification as a new paradigm for Abductive Explanation)*

**CHECK STEP**

**Property: Local Robustness**

Under test
XAI set
Free set

$$\pi' \leftarrow \pi \setminus \{x_i\}$$
$$\text{if } \text{CHECK}(M, \pi') \text{ then}$$
$$\text{Explanation} \leftarrow \text{Explanation} \cup \{x_i\}$$

$$N \left( \; \right) \overset{?}{=} \text{cat}$$

### *How it works?*
1.  *Fix XAI variables at their nominal values,*
2.  *Allow the removed feature and all other inputs to vary within their valid domains,*
3.  *Verify no property violations occur*

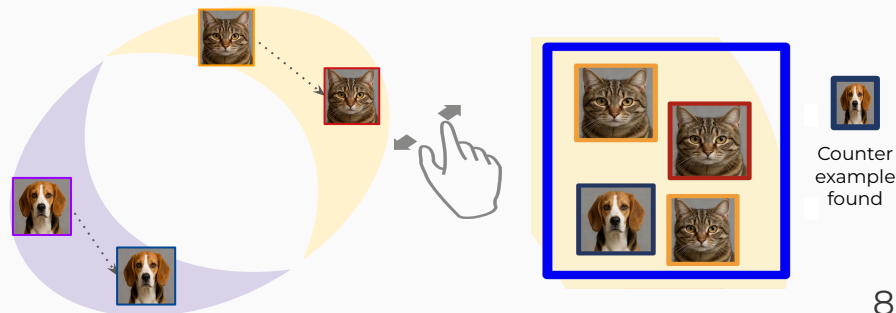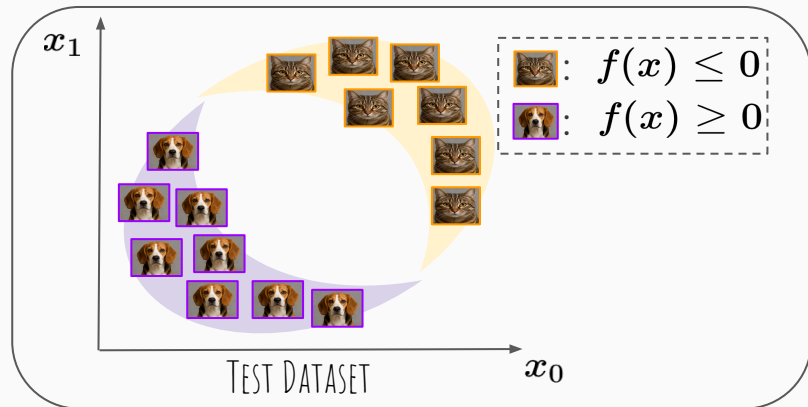# Verification property: Local Robustness

$$N\left( \boxed{\text{[cat image]}} \right) \stackrel{?}{=} \text{cat}$$

An example of a perturbation with epsilon = +/- 1 pixel

[cat image] + [noise image] ∈ [cat image box]

Perturbation domain



$x_1$

$\boxed{\text{[cat]}}: f(x) \leq 0$
$\boxed{\text{[dog]}}: f(x) \geq 0$

TEST DATASET

$x_0$

Counter example found

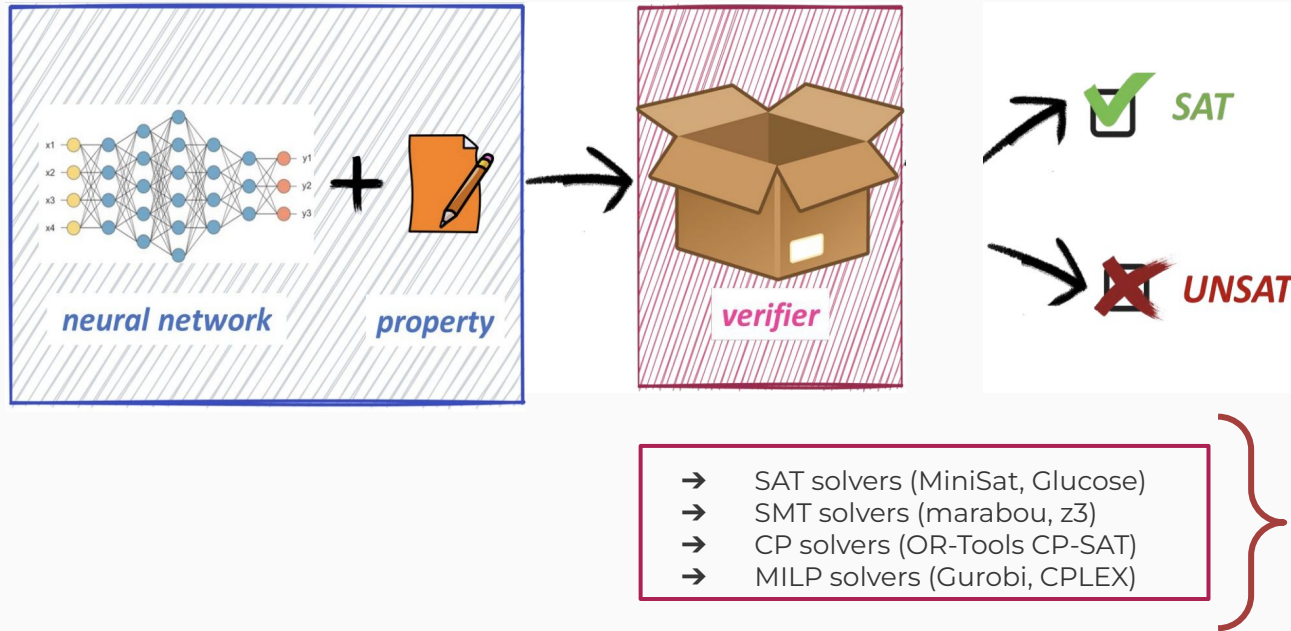# *Scaling up the performance of formal explainers*

---

**Algorithm 1** Deletion Algorithm to Find One Abductive Explanation

1: **Input:** Predictive model $M$ and input $x = \langle \chi^1, \ldots, \chi^n \rangle$
2: **Output:** Explanation for class $C$ of $x$
3: $Explanation \leftarrow \emptyset$        ▷ Set of relevant features
4: $I \leftarrow \emptyset$        ▷ Set of irrelevant features
5: $\pi \leftarrow TRAVERSALORDER(F)$     ▷ Sort F's features by ascending relevance
6: **for** each feature $x_i \in \pi$ **do**
7:     $\pi' \leftarrow \pi \setminus \{x_i\}$     ▷ Removing feature $x_i$
8:     **if** CHECK$(M, \pi')$ **then**     ▷ Check if prediction changes
9:         $Explanation \leftarrow Explanation \cup \{x_i\}$   ▷ Add feature $x_i$ to relevant features
10:     **else**
11:         $I \leftarrow I \cup \{x_i\}$     ▷ Add feature $x_i$ to irrelevant features
12:     **end if**
13: **end for**
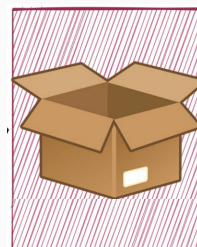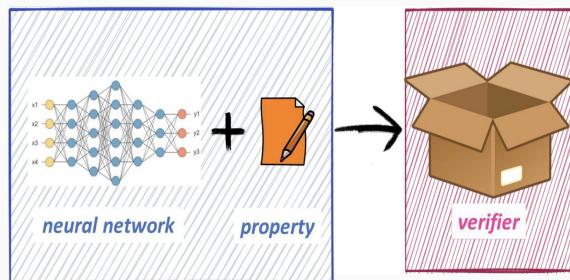14: **Return** $Explanation$

## Challenges to address

→ Challenge 1: The CHECK function is computationally expensive

9

# *About the CHECK method (Verification step)*

| | |
|---|---|
| ➔ SAT solvers (MiniSat, Glucose) | |
| ➔ SMT solvers (marabou, z3) | Complete |
| ➔ CP solvers (OR-Tools CP-SAT) | |
| ➔ MILP solvers (Gurobi, CPLEX) | |

Link to figure source

# Different techniques for NN verification



Adversarial | Incomplete | Complete

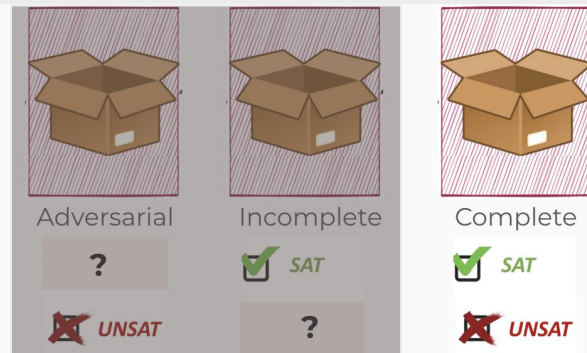**Increasing Runtime**

# Combining methods in one verification 'pipeline'

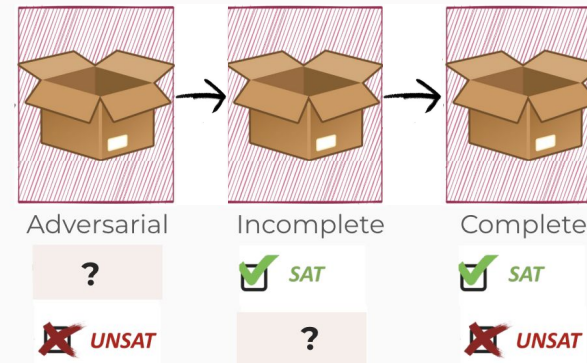**Algorithm 1** Deletion Algorithm to Find One Abductive Explanation

1: **Input:** Predictive model $M$ and input $x = \langle \chi^1, \ldots, \chi^n \rangle$
2: **Output:** Explanation for class $C$ of $x$
3: $Explanation \leftarrow \emptyset$                                   ▷ Set of relevant features
4: $I \leftarrow \emptyset$                                             ▷ Set of irrelevant features
5: $\pi \leftarrow TRAVERSALORDER(F)$           ▷ Sort F's features by ascending relevance
6: **for** each feature $x_i \in \pi$ **do**
7:     $\pi' \leftarrow \pi \setminus \{x_i\}$                          ▷ Removing feature $x_i$
8:     **if** CHECK$(M, \pi')$ **then**                                ▷ Check if prediction changes
9:         $Explanation \leftarrow Explanation \cup \{x_i\}$    ▷ Add feature $x_i$ to relevant features
10:    **else**
11:        $I \leftarrow I \cup \{x_i\}$                               ▷ Add feature $x_i$ to irrelevant features
12:    **end if**
13: **end for**
14: **Return** $Explanation$

➤    Challenge 1: The CHECK function is computationally expensive



Increasing Runtime

12

# *Parallelizing the computation of formal explanation*

**Algorithm 1** Deletion Algorithm to Find One Abductive Explanation

1: **Input:** Predictive model $M$ and input $x = \langle \chi^1, \ldots, \chi^n \rangle$
2: **Output:** Explanation for class $C$ of $x$
3: $Explanation \leftarrow \emptyset$ ▷ Set of relevant features
4: $I \leftarrow \emptyset$ ▷ Set of irrelevant features
5: $\pi \leftarrow TRAVERSALORDER(F)$ ▷ Sort F's features by ascending relevance
6: **for** each feature $x_i \in \pi$ **do**
7:     $\pi' \leftarrow \pi \setminus \{x_i\}$ ▷ Removing feature $x_i$
8:     **if** CHECK$(M, \pi')$ **then** ▷ Check if prediction changes
9:         $Explanation \leftarrow Explanation \cup \{x_i\}$ ▷ Add feature $x_i$ to relevant features
10:    **else**
11:        $I \leftarrow I \cup \{x_i\}$ ▷ Add feature $x_i$ to irrelevant features
12:    **end if**
13: **end for**
14: **Return** $Explanation$

Challenge 2: Sequential traversal of the feature set (loop)

# *Could we Add to explanation a batch of input Features ?*

Adversarial

?

**UNSAT**

Adversarial Attack: Add batch of input feature

**Idea:** Propose a new strategy that breaks the sequential query bottleneck in deletion-based formal XAI

**How?** Enable parallel removal of feature constraints and launch several adversarial attacks at once

**Advantages?** Batch processing & GPU implementation supported by existing adversarial methods, no extra development required!

# *Could we Free a batch of input Features ?*

Abstract Interpretation: free batch of input feature

**Idea:** Go beyond SAT/UNSAT verifier's decision! Leverage solver proofs to pinpoint and free multiple feature indices in a single iteration

**How?** Determine the largest subset of features that can be freed without compromising the property's & soundness

Incomplete

SAT

?

$$\mathbb{I} \text{ s. t. } \forall \mathbb{I}', \ P(\mathbb{I}') \implies |\mathbb{I}'| \leq |\mathbb{I}|,$$

$$\text{where } P(\mathbb{I}) \ : \ \exists \, E \in F \setminus \mathbb{I}, \ \left( \bigwedge_{i \in E} x_i = v_i \right) \implies N(x) = c$$

**Advantages?** One call of abstract interpretation is enough !
Linear complexity + Soundness

# *Statistically-Guided Explanations*

**Algorithm 1** Deletion Algorithm to Find One Abductive Explanation

1: **Input:** Predictive model $M$ and input $x = \langle \chi^1, \ldots, \chi^n \rangle$
2: **Output:** Explanation for class $C$ of $x$
3: $Explanation \leftarrow \emptyset$ ⊳ Set of relevant features
4: $I \leftarrow \emptyset$ ⊳ Set of irrelevant features
5: $\pi \leftarrow TRAVERSALORDER(F)$ ⊳ Sort F's features by ascending relevance
6: **for** each feature $x_i \in \pi$ **do**
7:     $\pi' \leftarrow \pi \setminus \{x_i\}$ ⊳ Removing feature $x_i$
8:     **if** CHECK$(M, \pi')$ **then** ⊳ Check if prediction changes
9:         $Explanation \leftarrow Explanation \cup \{x_i\}$ ⊳ Add feature $x_i$ to relevant features
10:     **else**
11:         $I \leftarrow I \cup \{x_i\}$ ⊳ Add feature $x_i$ to irrelevant features
12:     **end if**
13: **end for**
14: **Return** $Explanation$

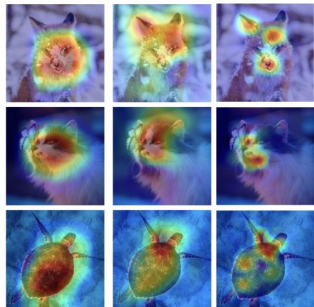Challenge 3: Impact of order on explanation interpretability (size)

16

# *Statistically-Guided Explanations*

**DEEL**
**XPLIQUE**
Explainability Toolbox for Neural Networks



Feature attributions

**Idea:** Address the cardinality bottleneck in formal XAI by leveraging statistical explanation orderings

**How?** Synergies with statistical XAI to guide formal search, test different feature ordering and chose the best
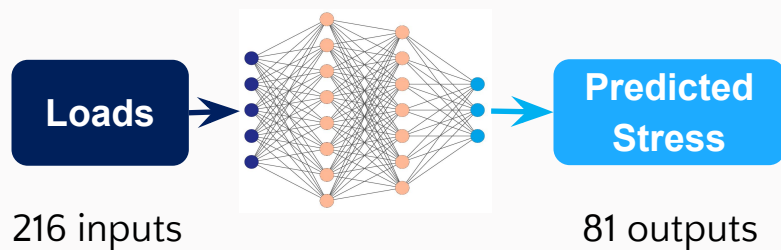
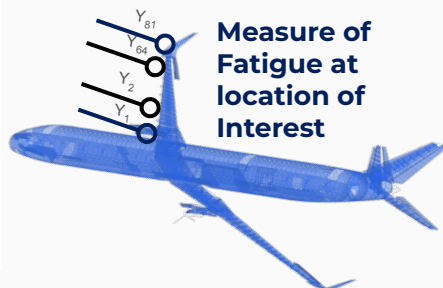**Advantages?** Take advantage of the several XAI statistical-based techniques available and create synergies between the two communities!
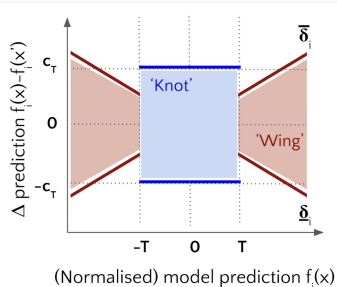
*Preliminary results*

# Uses cases

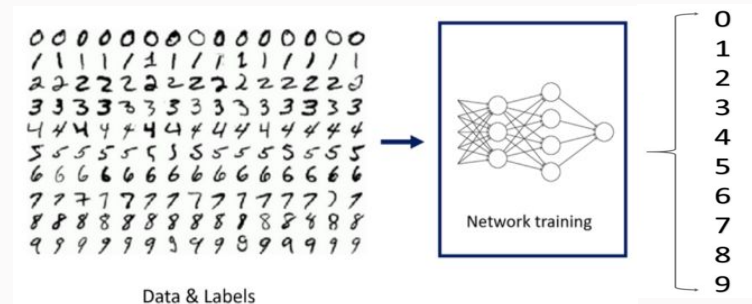## Local Stability
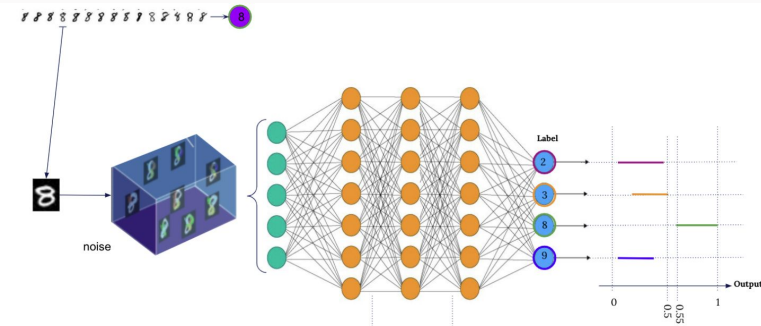
### Industrial use case : Fatigue Digital Twin



216 inputs          81 outputs

### Property



Measure of Fatigue at location of Interest

## Local Robustness

### Academic use case



Data & Labels

Network training

### Property



noise

# *Libraries & Tools*

**Explanation computation**



### *XAirobas*

*AIRBUS/ANITI*



**Verification Pipeline**



### *Airobas*

*AIRBUS*
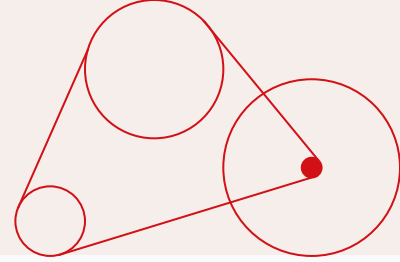


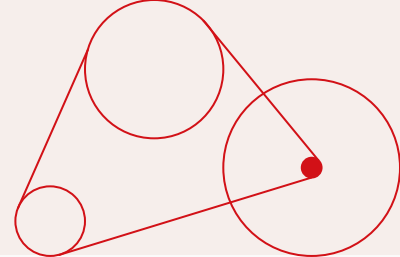**Abstract Interpretation**



### *Decomon*

*AIRBUS/ANITI*

# *Results on industrial use case (Fatigue Digital Twin)*

| CHECK config / Metrics | Only complete **Baseline** | Verification pipeline **Challenge 1** | Pipeline & batch processing **Challenge 1+2** |
|---|---|---|---|
| Runtime per explanation | 18mn48sec | **10mn79sec** | **2.14sec!** |
| Average explanation size | 107 | **72** | **36** |
| #Calls to adv attacks | 0 | 104 | 34 |
| #Calls to incomplete solver | 0 | 99 | 2 + **1** |
| #Calls to complete solver | 216 | **13** | **4** |

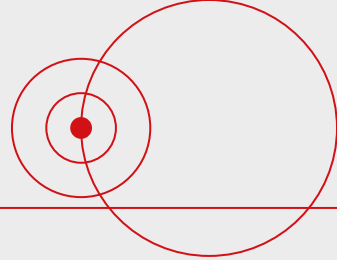1 call to free in average **176 out of 216 at once !!!!**

# *Results on academic use case (MNIST)*

| CHECK config / Metrics | Only complete **Baseline** | Verification pipeline **Challenge 1** | Pipeline & batch processing **Challenge 1+2** |
|---|---|---|---|
| Runtime per explanation | 3876.61sec | **1783.28sec** | **495.5sec** |
| Average explanation size | 111.85 | 111.85 | **110.93** |
| #Calls to adv attacks | 0 | 58 | 61 |
| #Calls to incomplete solver | 0 | 583 | 303 + **1** |
| #Calls to complete solver | 784 | **143** | 134 |

1 call to free in average **286 out of 784 at once !!!!**

# *Takeaways*

**Algorithm 1** Deletion Algorithm to Find One Abductive Explanation

1: **Input:** Predictive model $M$ and input $x = \langle \chi^1, \ldots, \chi^n \rangle$
2: **Output:** Explanation for class $C$ of $x$
3: $Explanation \leftarrow \emptyset$        ▷ Set of relevant features
4: $I \leftarrow \emptyset$        ▷ Set of irrelevant features
5: $\pi \leftarrow TRAVERSALORDER(F)$      ▷ Sort F's features by ascending relevance
6: **for** each feature $x_i \in \pi$ **do**
7:     $\pi' \leftarrow \pi \setminus \{x_i\}$        ▷ Removing feature $x_i$
8:     **if** CHECK$(M, \pi')$ **then**       ▷ Check if prediction changes
9:        $Explanation \leftarrow Explanation \cup \{x_i\}$    ▷ Add feature $x_i$ to relevant features
10:     **else**
11:        $I \leftarrow I \cup \{x_i\}$       ▷ Add feature $x_i$ to irrelevant features
12:     **end if**
13: **end for**
14: **Return** $Explanation$

**XAirobas**

COMING SOON

## Contributions

→ Contribution 1: Introduced a modern formal verification pipeline tailored to the scalability demands

→ Contribution 2: Propose a novel distributed strategy that breaks the sequential query bottleneck

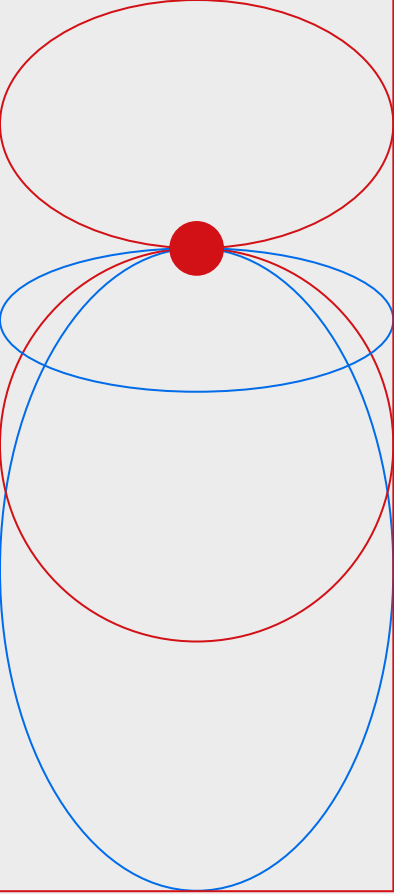→ Contribution 3: Statistically-Guided Explanations

**Contacts:**

Ryma Boumazouza :
ryma.boumazouza@airbus.com

Mélanie Ducoffe :
melanie.ducoffe@airbus.com

*Thank you !*