# FIIA 2021 Trustworthy Artificial Intelligence in AIRBUS

Romaric Redon AIRBUS
Head Advisor Artificial Intelligence

**AIRBUS**

skywise.

With Skywise, unleash
the full potential of every aircraft

AIRBUS

OneAtlas

Connecting Images from Space
to Decisions on Earth.

sobloo

Beyond the Data
Creative Grounds

up

Build, run,
and scale
geospatial
products

AIRBUS

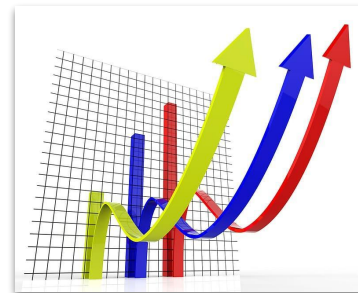# What do we do with Artificial Intelligence @AIRBUS

| Perceive / Observe | Predict / forecast / Orient | Decide/Act |
|---|---|---|

**Computer Vision**

**Pattern recognition and Times Series Analysis**

**Natural Language understanding and Processing**

**Hybrid Modelling**

**Decision making**

Data based AI

Symbolic AI

**AIRBUS**

# What do we do with Artificial Intelligence @AIRBUS

| Perceive / Observe | Predict / forecast / Orient | Decide/Act |

**Computer Vision** TRUST

**Pattern recognition and Times Series Analysis** TRUST

**Hybrid Modelling** TRUST

**Decision making** TRUST

**Natural Language understanding and Processing** TRUST

TRUST
for all application cases
for all AI technologies

TRUST is ensured at
System level

Data based AI TRUST

TRUST Symbolic AI

**AIRBUS**

# Trustworthy AI - Several dimensions

Responsible use of AI



http://www.fcas-forum.eu/en

Safe use of AI

Fairness

Robustness

Explainability

https://arxiv.org/abs/2103.10529

AIRBUS

# Trustworthy AI requirements vs criticality levels

# AI to provide assistance to the pilot?

Hybrid Modelling

Concept of Operations

Certification Basis

Architecture

Crew Awareness Errors Prevention

Engaging Experience

Virtual Assistant

Decision making

Automated Systems

Fatigue Management

Human Health Monitoring

Permanent Auto Flight

Image-based navigation & obstacle detection

Pattern recognition and Times Series Analysis

Computer Vision

Trajectory Predictions & Conflict Resolution

Ground Assistance

Connectivity

Speech to Text

Massive data collection

Massive testing

Trustworthy AI

Natural Language understanding and Processing

AIRBUS

# Trustworthy AI Engineering

Requirement and Data Eng for Trustworthy AI

Design for Trustworthy AI

Monitoring and fail safe Architectures

Trust On-board HW for AI & code generation

V&V methods for critical systems with AI

AI regulation and Standards - certification



**AIRBUS**

# Requirement and Data Eng for Trustworthy AI

With Data Driven AI  Spec ~ Data

Difficulty to define the Operational Design
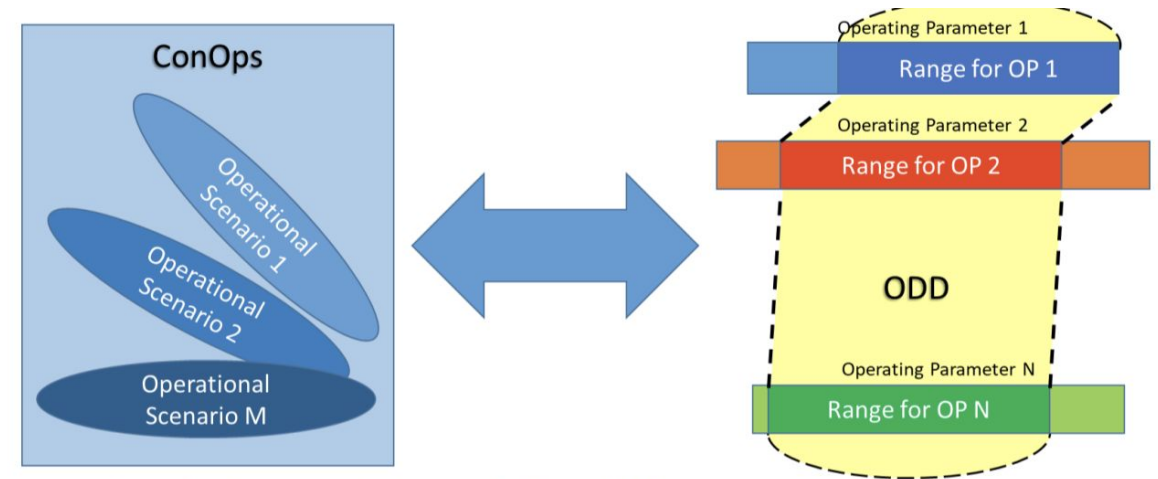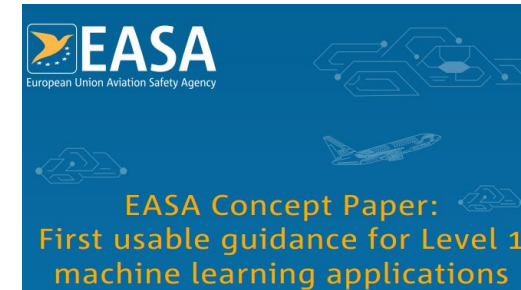Domain for high dimensional input space



EASA Concept Paper:
First usable guidance for Level 1
machine learning applications

How to specify the data needed to cover the ODD?
How to detect undesired biais?
How to choose the right distribution?
How to use a good mix of real/synthetic data to
have the right distribution?



Figure 5 — Interrelationship between ConOps and ODD

Objective CO-03: The applicant should define and document the
ConOpsfor all AI-based (sub)systems. A focus should be put on the
definition of the **operational design domain (ODD)** and on the
capture of specific operational limitations and assumptions.

**AIRBUS**

# Design for Trustworthy AI - Robustness
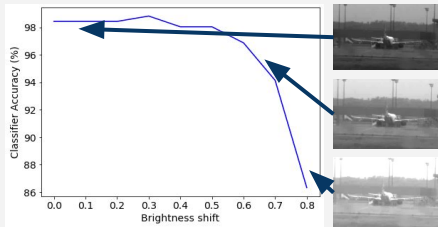
How to improve robustness ?

Generalisation guarantees : how accurately a Machine Learning algorithm is able to predict outcome values for previously unseen data?
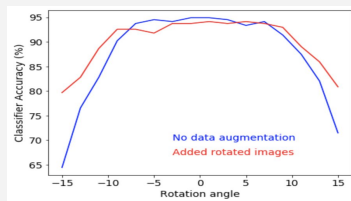
## Natural perturbations

### 1) Robustness assessment

Performance evaluation under artificial corruptions, e.g. brightness change, rotations, occlusions, rain, etc.



### 2) Data augmentation

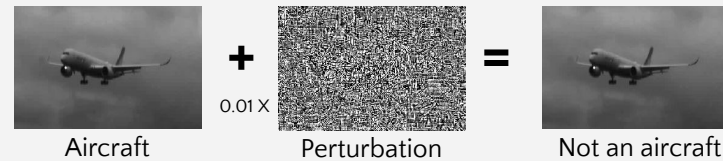Use of perturbed images during training improves the model robustness



## Adversarial attacks

adver torch

IBM ART

### 1) Attack benchmarks

Up to 96% of all inputs can be successfully attacked (i.e. visually indistinguishable but wrong model prediction)



Aircraft          0.01 X   Perturbation          Not an aircraft

### 2) Adversarial defenses

- **Adversarial training**: attacks in training reduces success rate from 96% to 34%
- **Smoothing defense**: trade-off between classification accuracy and robustness



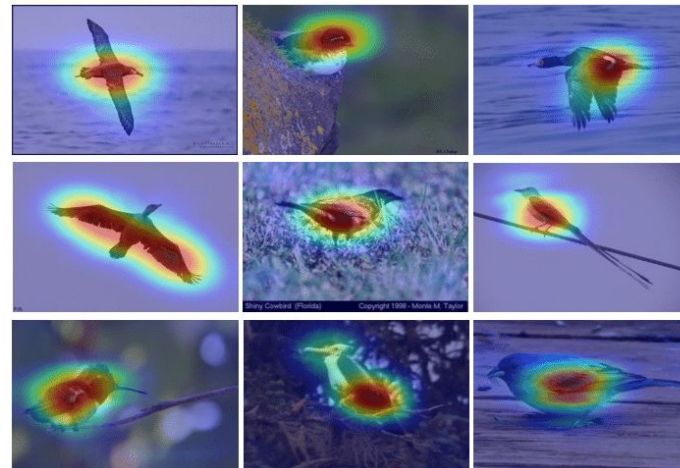Original image          Attacked image used in training

**AIRBUS**

# Design for Trustworthy AI - Explicability

"The AI explainability deals with the capability to provide the human with understandable and relevant information on how an AI/ML application is coming to its results." EASA guidelines level 1
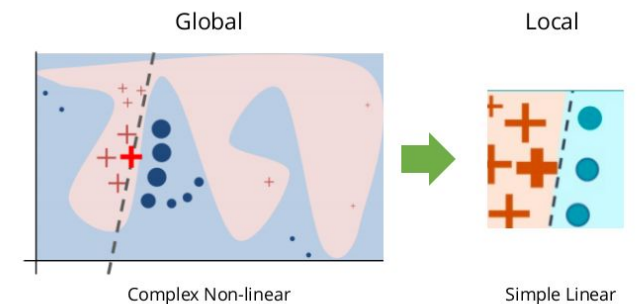
Explainability for

- ML model designer/authorities
- Operator
- forensic investigations

Saliency Maps



He, Xiangteng & Peng, Yuxin & Zhao, Junjie. (2017). Fine-grained Discriminative Localization via Saliency-guided Faster R-CNN.
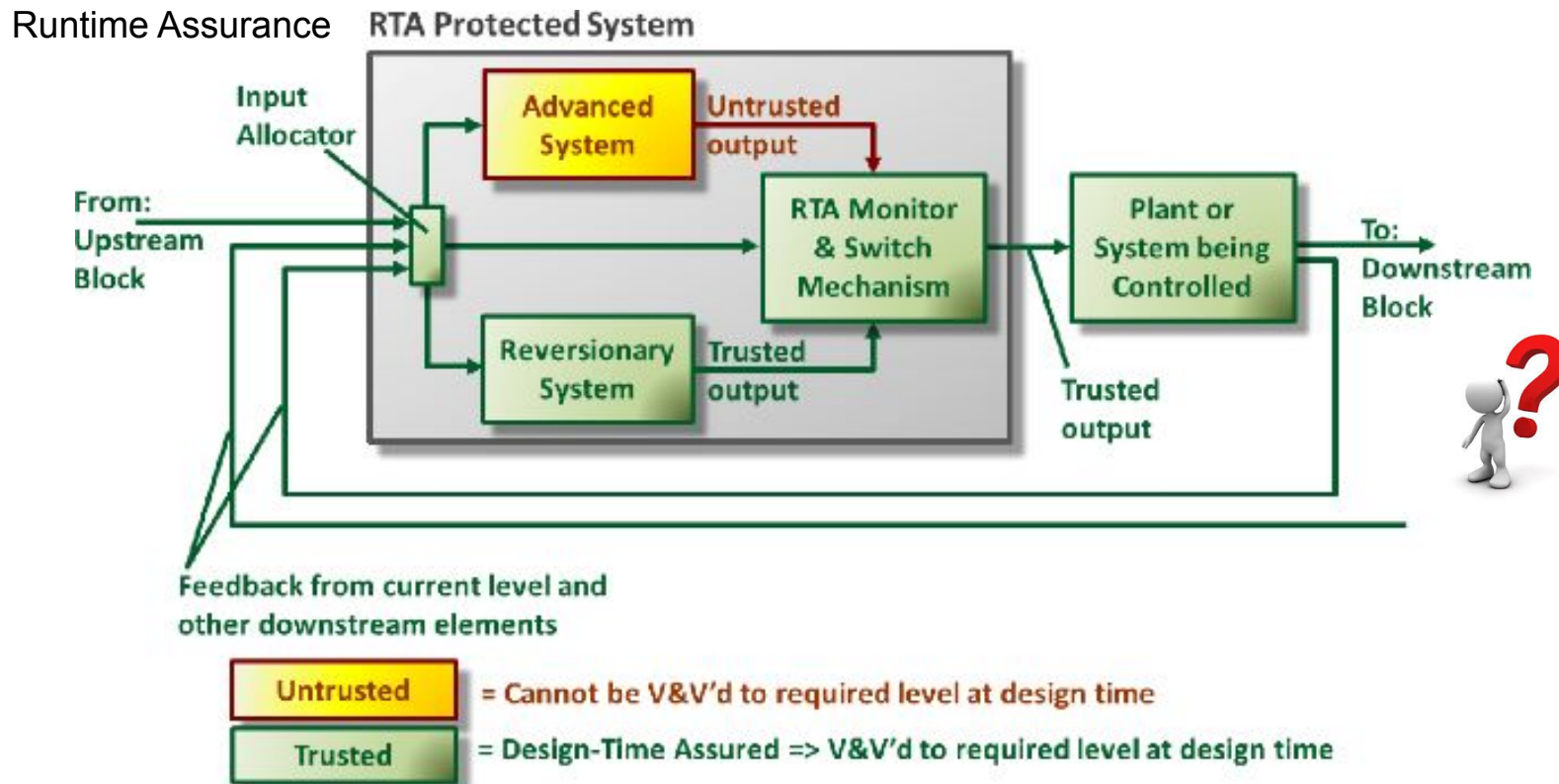
ML interpretation models
SHAP, LIME,...



How good is your Explanation?
How could you use it to provide guarantees?

What about NLP and Speech, Planning & Scheduling?

AIRBUS

# Monitoring and fail safe architecture

- ODD Monitoring verifies that the ML-based system is operated in its usage domain
- OOD Monitoring ensures that the ML Model operates in the distribution defined during the training process.
- Attacks monitoring allows to detect adversarial attacks.
- Robustness monitoring ensures that the ML Model is used in a stable area.
- Consistency monitoring analyzes the consistency of outputs.
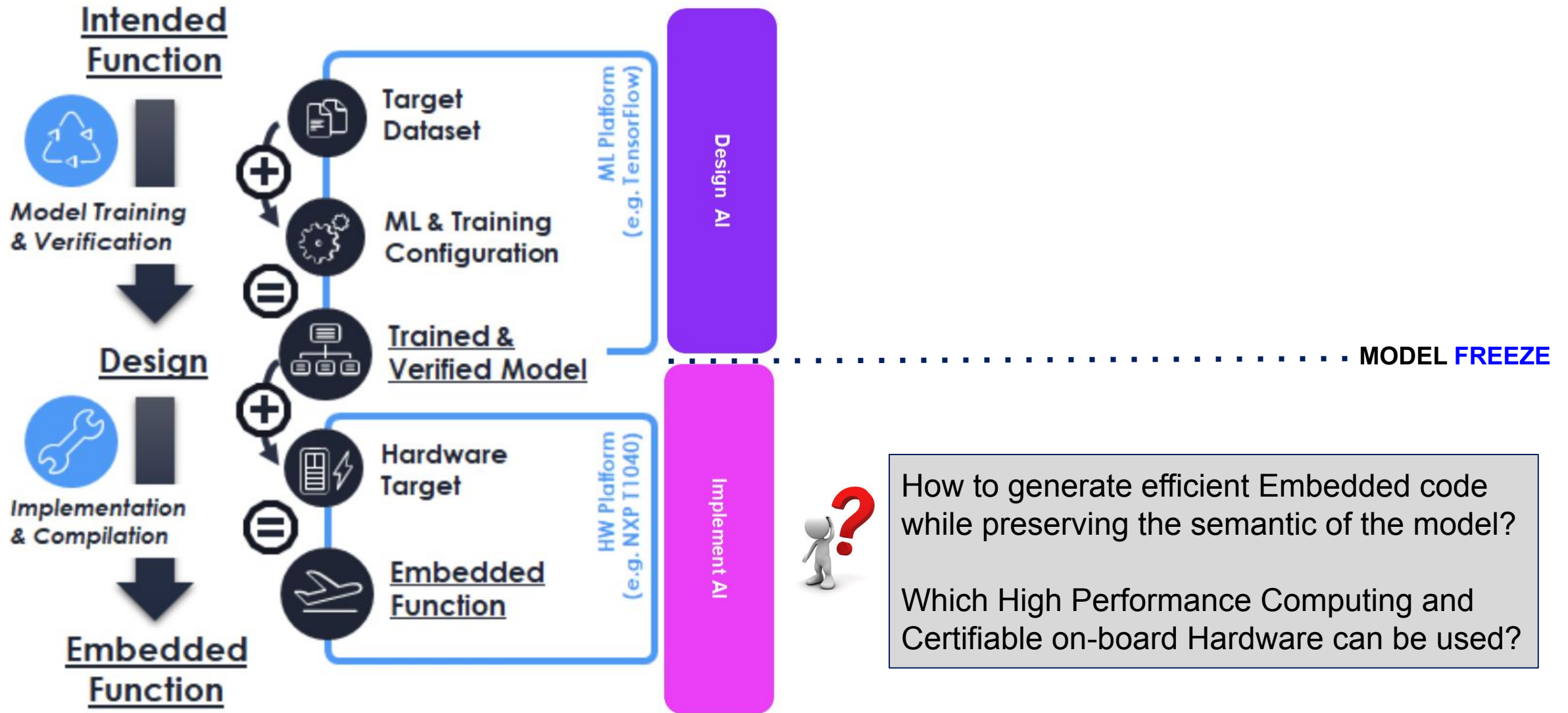


Runtime Assurance

What is the more efficient monitoring architecture?

What would be a Dissimilar Architecture with ML?
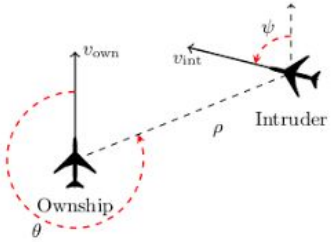Is Same data set - different models enough ?

How to find the best trade-of
Impact of monitoring on availability?

Schierman, John D. et al. "Runtime Assurance Framework Development for Highly Adaptive Flight Control Systems." (2015).

**AIRBUS**

# Trusted on-board hardware for AI + code generation



How to generate efficient Embedded code while preserving the semantic of the model?

Which High Performance Computing and Certifiable on-board Hardware can be used?
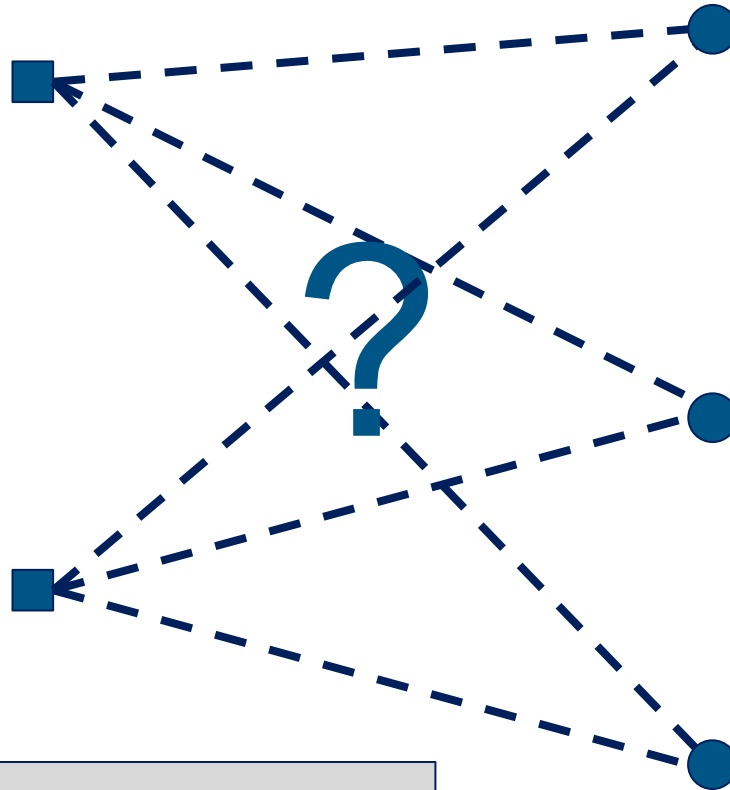
**AIRBUS**

# V&V methods for critical systems with AI



NN to approximate complex functions (table, data, physical models)

Deep NN for Vision Based Navigation

Formal Verification

Adversarial ML

Massive Testing with real and synthetic data

How to improve scalability of formal methods?
How to define the good properties and perturbations?
How to perform more efficient/frugal massive Testing?
How to qualify synthetic data generator?
How to combine/Hybridize methods?

**AIRBUS**
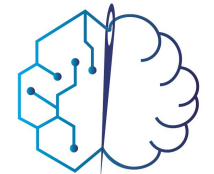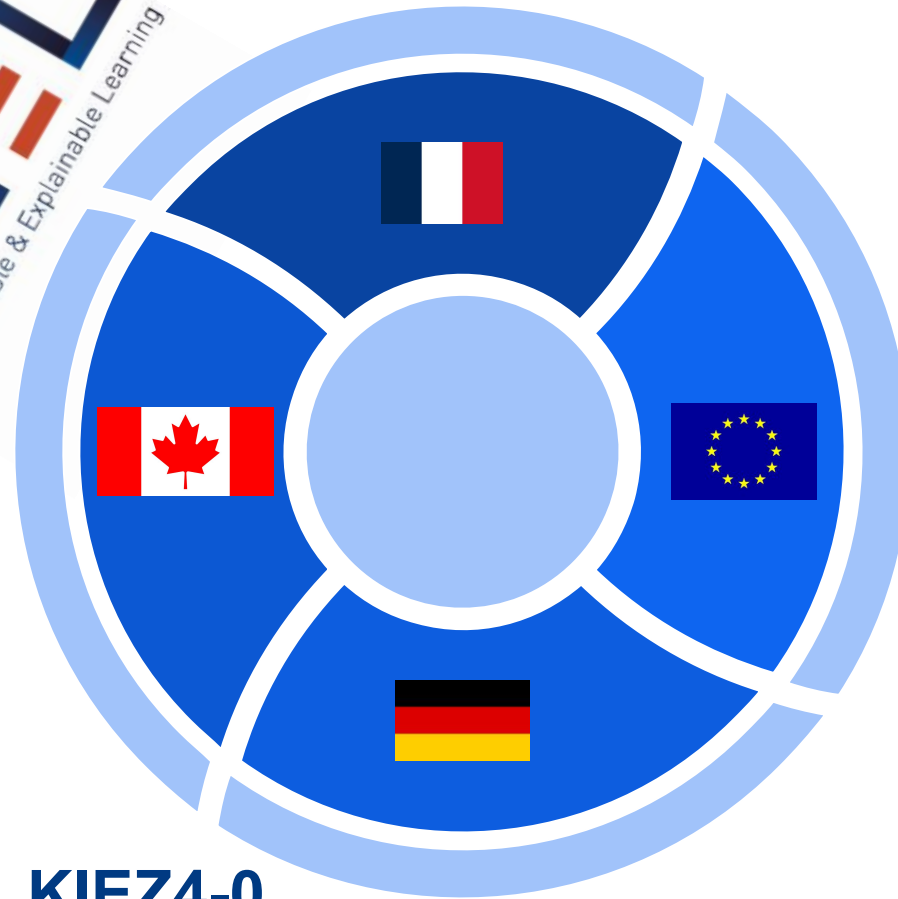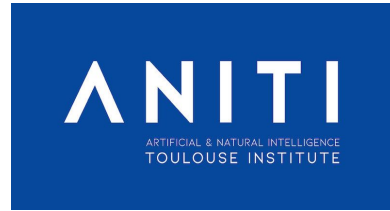
# AI regulation and Standards



SECTORIAL

How General and Sectorial regulations will be articulated ?

GENERAL

AIRBUS

# Conclusion (1/2)

**AIRBUS**

Trusted AI methods

Integrate AI in critical system engineering

What we can demonstrate

What we should demonstrate

Ensure Safe and Efficient operations

**Academics + Tech companies**

**Industrials**

**Regulation authorities**