IRT

SAINT
EXUPÉRY

# Du modèle d'activité au modèle d'argumentation : apport des techniques de MBSE à l'ingénierie des systèmes à base d'IA

AFIS CROcc CNES COMET

# Le programme de recherche Confiance.ai

12/12/2025

# THE EUROPEAN TRUSTWORTHY AI ASSOCIATION

The European Trustworthy AI Association is a non-profit organization established by industrial leaders, building on the legacy of the Confiance.ai programme. It is on a mission to empower the industry with state of the art, open-source methodology and tools, enabling the engineering of AI-based systems that can be trusted and comply with regulations.

The association aims to be a driving force behind an ambitious European strategy for industrial and responsible AI. Its ambition is to propel Europe to the forefront of innovation in trustworthy AI, by making its methodologies and tools an international benchmark and thus, supporting the broader adoption of responsible AI in industy.

12/12/2025

fit | FRENCH INSTITUTES OF TECHNOLOGY

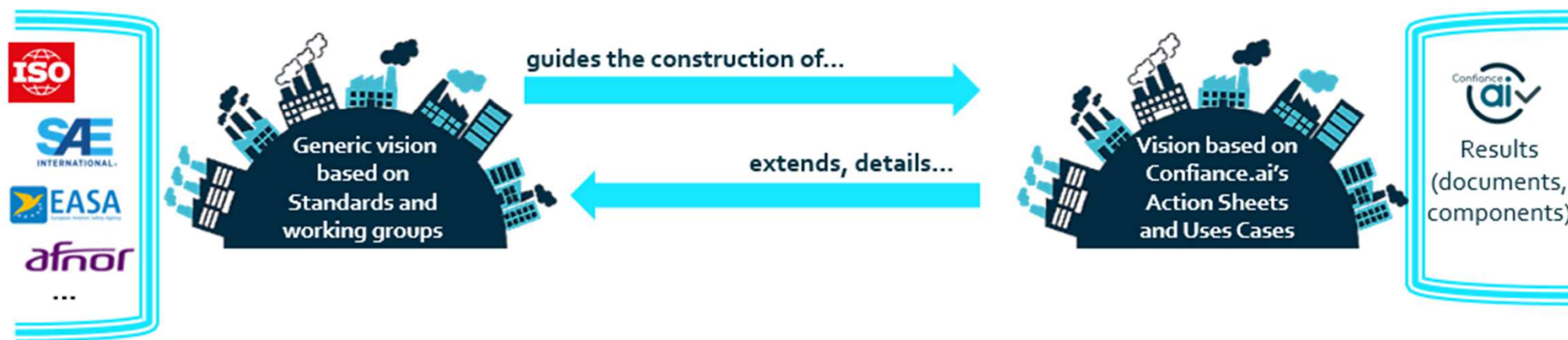# Buts de la méthode « end-to-end » de Confiance.ai

- Complete the « classical » engineering disciplines (Systems Engineering, Software Engineering) to take into account the specificities of ML, with modifications only where necessary
- Structure the results of Confiance.ai (local methods, software components) to facilitate their use
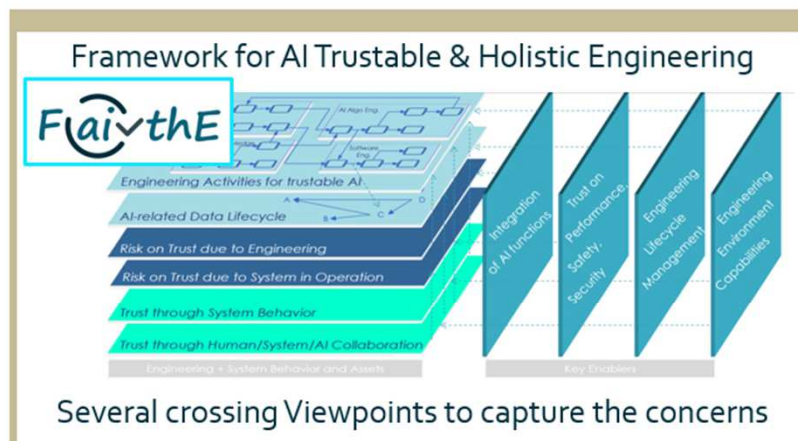
FRENCH
INSTITUTES OF
TECHNOLOGY

# Double approche pour construire la méthode

Top-down approach:
capture of a high-level, holistic vision of an engineering process for trustable AI-based systems

Bottom-up approach:
capture of Methods & Processes elaborated by Confiance.ai Projects for specific topics
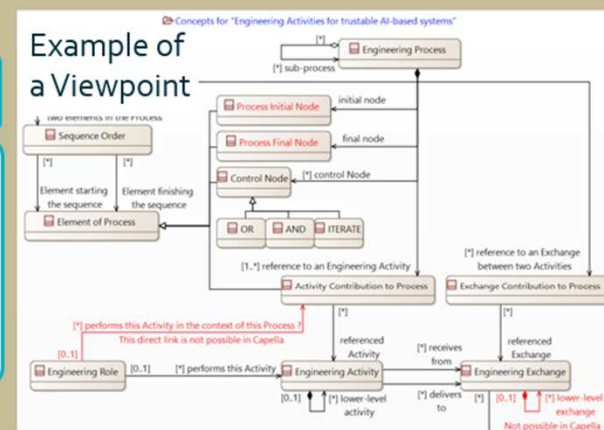
guides the construction of…

extends, details…

Generic vision based on Standards and working groups

Vision based on Confiance.ai's Action Sheets and Uses Cases

Results (documents, components)

FRENCH INSTITUTES OF TECHNOLOGY

# Modélisation des processus et activités d'ingénierie

# Complexité d'une telle méthode

Design breakdown · IVVQ build-up

Simultaneous multi-disciplines / specialties design & assessment

Simultaneous multi-stage / multi-scope feedback loop and continuity (concept, design, IVV, manufacturing, operations)

Simultaneous multi-level integration

12/12/2025

**fit** | FRENCH INSTITUTES OF TECHNOLOGY

# Méthode accessible via la Body Of Knowledge

https://bok.confiance.ai/

12/12/2025

FRENCH
INSTITUTES OF
TECHNOLOGY

# Vue de plus haut niveau

# Zoom sur "Evaluation of ML Model"

# Zoom sur "Evaluation of ML Model robustness"

12/12/2025

FRENCH INSTITUTES OF TECHNOLOGY

# Zoom sur "Test of ML Model robustness"

Related to :

- §B.2. of document "Methodological Guideline for Robustness Functional Set"

- Component 331: Adversarial Attack Characterization Component

- Component 332: AI Metamorphis Observer Component (AIMOS)

- Component 333: Amplification Method for Robustness Evaluation Component

- Component 334: Non-overlapping Corruption Benchmarker Component

- Component 335: Time-series Robustness Characterizer Component

- Component 3141: Chiru

# Zoom sur "Formal verification of ML Model robustness"



Related to :

- §B.1. of document "Methodological Guideline for Robustness Functional Set"

- Component 321: Saimple

- Component 322: nnenum

- Component 323: α-β-crown

- Component 3171: PyRAT

- Component 391: MIP Solver

# Spécificités de l'IS of d'un système basé ML: ODD

# Spécificités de l'IS of d'un système basé ML: écart avec un système classique
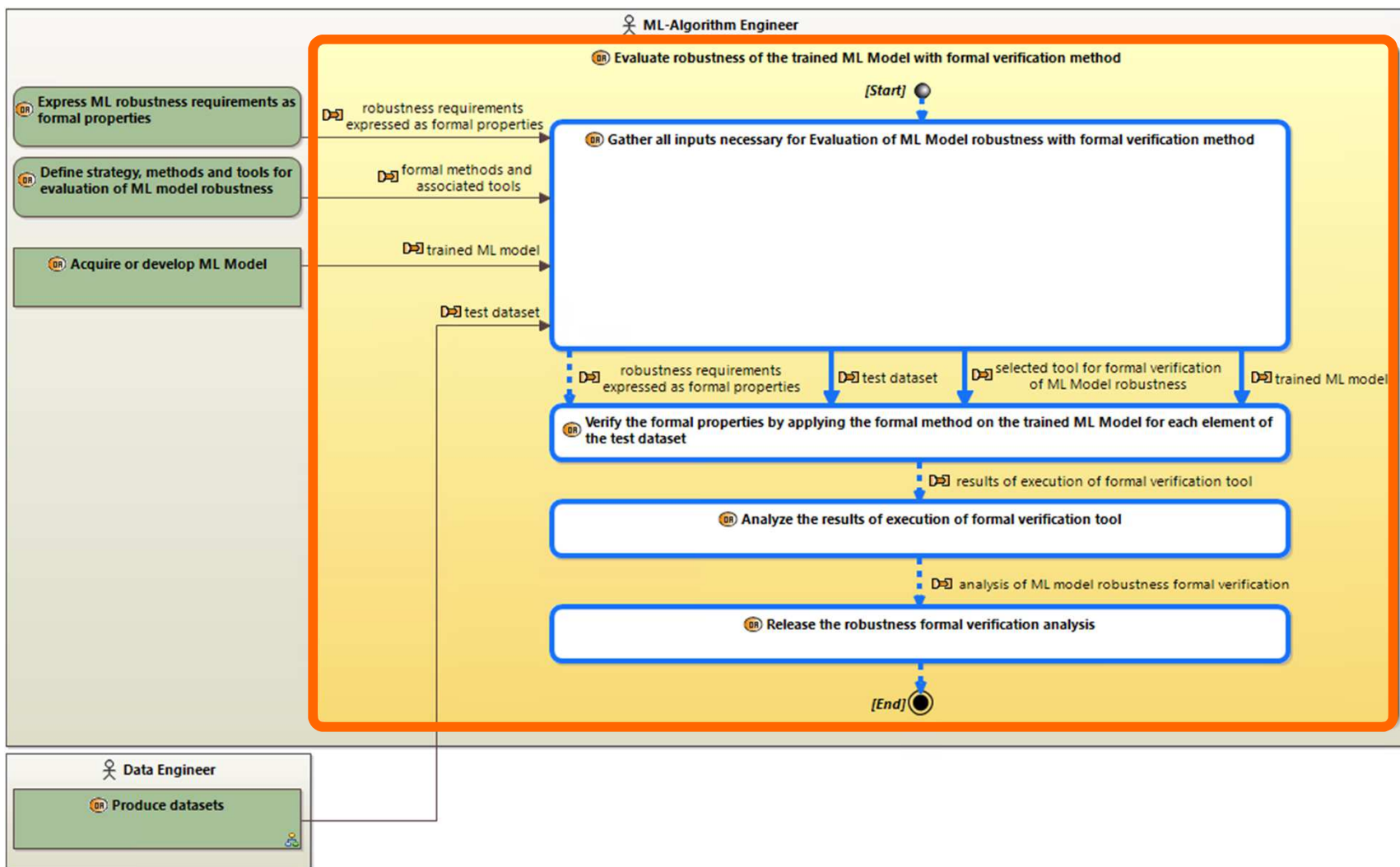
# Spécificités de l'IS of d'un système basé ML: écart avec un système classique

Function of a "conventional" component



Function of an ML-based component



- Expected environmental conditions
- Harmful environmental conditions
- Required intended behavior
- Unwanted or Disturbing Behavior

12/12/2025

**fit** | FRENCH INSTITUTES OF TECHNOLOGY

# Spécificités de l'IS of d'un système basé ML: écart avec un système classique



From Systems Engineering to ML engineering / data engineering: choice of ML algo, choice of data…

# V&V d'un système basé ML

# De la Confiansssssse
# à la Confiance

E. JENN – IRT Saint-Exupery

# Trusting ML-based systems?

```
Compliance to
standards
```

```
Track-record
```

**Assurance cases**

**Confidence**

**Dependability**

"[…] a psychological state which, if rational, must be **based on the reasons**—that is, the justification—for believing the claims." (J. Rushby)

"the trustworthiness of a computer system such that reliance can **justifiably** be placed on the service it delivers" (W.C. Carter, in Laprie *et al. "Dependability: Basic Concepts and Terminology*)

"[…] **claims, argument, and evidence is surely the (perhaps tacit)** intellectual foundation of any rational means for assuring and certifying the safety or other critical property of any kind of system. However, assurance cases differ from other means of assurance, such as those based on standards or guidelines, by making **all three components explicit.**" (J.Rushby)

Ref : S05-TEMP-009
Mise à jour : 05/05/2025

12/12/2025

FRENCH
INSTITUTES OF
TECHNOLOGY

# Robustness AC Template
## ACs and Goal Structuring Notation

- AC formalism uses a set of concepts and notations (cf. GSN 3):
    - **Goal (& subgoals)**: affirmation that shall be assessed during the reasoning.
    - **Solution**: A solution refers to some evidence that is deemed sufficient to establish the truth of the parent claim

    - **Strategy**: justifies the decomposition of goals into sub-goals.

    - And a few other elements (context, assumptions, etc.)

# Global Approach

**Select an engineering item** from the workflow — 1

Identify a **property** of interest on the engineering item — 2

Retrieve **Generic Arguments** (ACs) — 3

7 — Ensure **confidence** in ACs

4 — Adapt the argument w.r.t. context, cost, **confidence**…

5 — Integrate V&V **activities** that generate evidences

6 — Produce **V&V Plan**

**Workflow**

Activity — *Produces* → item

Activity — *Consumed by* — *Verified by* → V&V activity

*Couple*

*Produced by*

*Composes*

V&V plan

**Assurance Case**

Claim

OR

Evidence    Evidence

Propert

23

# Robustness AC Template
## From Engineering Items to Assurance Cases



*ML Workflow*

*Engineering item*

*Property*

*Assurance case*

# Robustness AC Template
## Generic AC



Robustness
Assurance
Case for
Robustness

**[G000001]**
The trained ML model is robust to input perturbations

**[A000002]**
The task considered is classification

**[G000034]**
The Trained ML Model satisfies the global robustness criteria

*Property of interest*

*L2 robustness*

*L∞ robustness*

FRENCH
INSTITUTES OF
TECHNOLOGY

# Robustness AC Template

## 2. Partitioning by norms (only l2 and l∞ considered)



Partitioning by robustness criteria
**Local Robustness Norm Selection**
Strategy pattern Process-based Vs. Product-based
Design Method

Local Robustness Norm Selection

*User choice*

⊙ ⬌ **l2 locally robust**
○ ⬌ **l∞ locally robust**

[G000001] The trained ML model is robust to input perturbations (...1 line skipped)

[A000002] The task considered is classification

[G000034] The Trained ML Model satisfies the global robustness criteria (...1 line skipped)

[G000074] The Trained ML Model satisfies the Global nbsample robustness criteria (...1 line skipped)

[G000073] The verification set is relevant for robustness evaluation (...1 line skipped)

[G000076] The Trained ML Model is locally robust according to relevant norms (...1 line skipped)

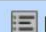[A000077] Relevant norms considered are L2 and L-inf

[Sol000075] Verification set

[S000078] Argument by partitioning of norms

*L∞ norm*

*L2 norm*

[G000079] The Trained ML Model is l2 locally robust (...1 line skipped)

[S000081] Argument over guarantees obtained by design or by verification

**FRENCH INSTITUTES OF TECHNOLOGY**

# Robustness AC Template

## Strategy pattern Process-based (By Design) Vs. Product-based (By verification)

# Robustness AC Template

## Strategy pattern Process-based (By Design) Vs. Product-based (By verification)

# Robustness AC Template

## Strategy pattern Process-based (By Design)  Vs. Product-based (By verification)

Partitioning by robustness criteria
Local Robustness Norm Selection
**Strategy pattern Process-based Vs. Product-based**
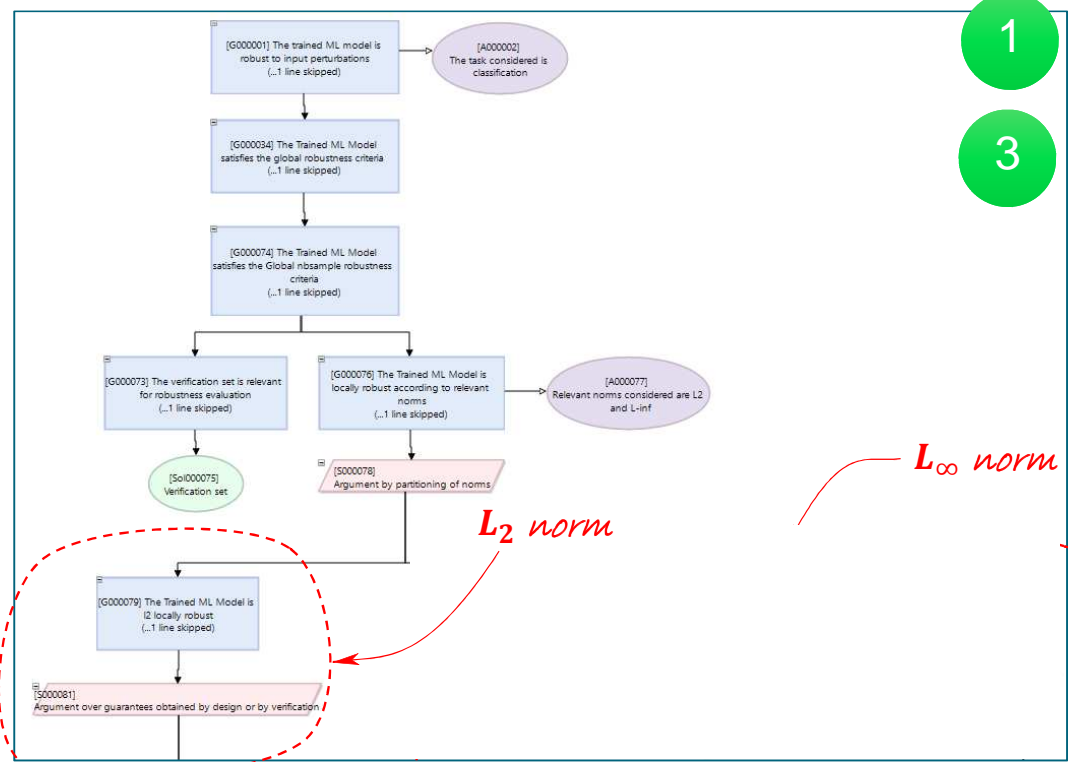Design Method

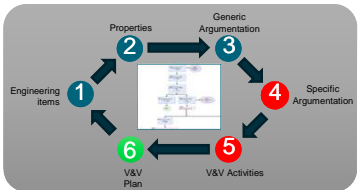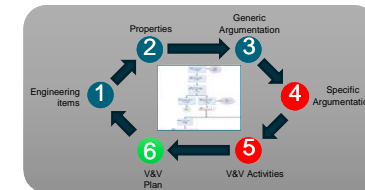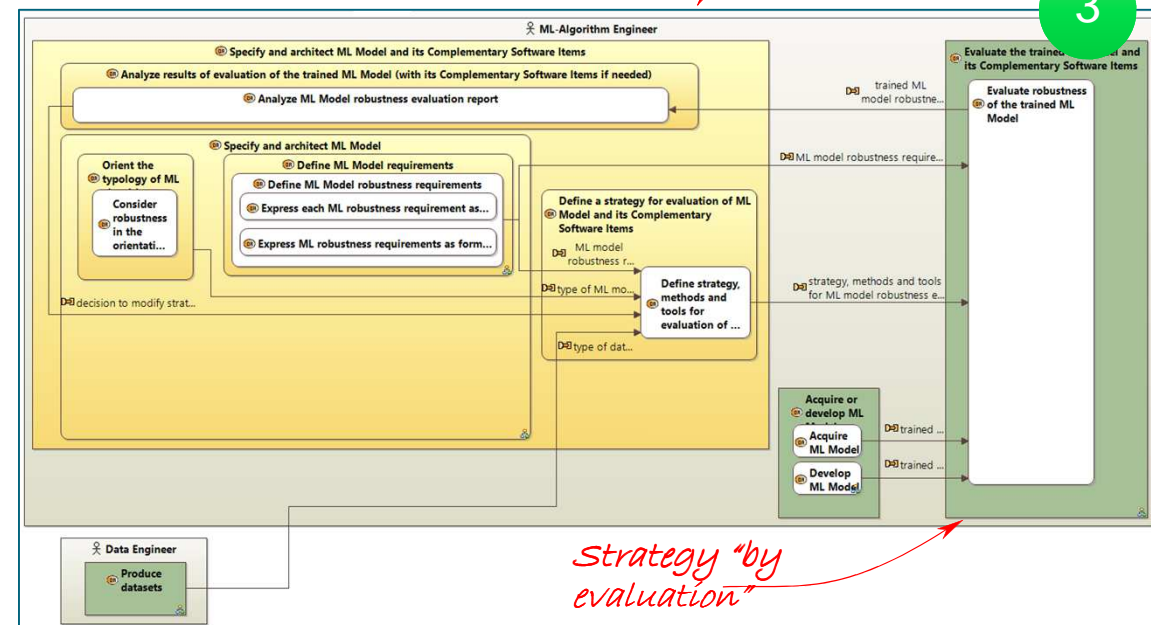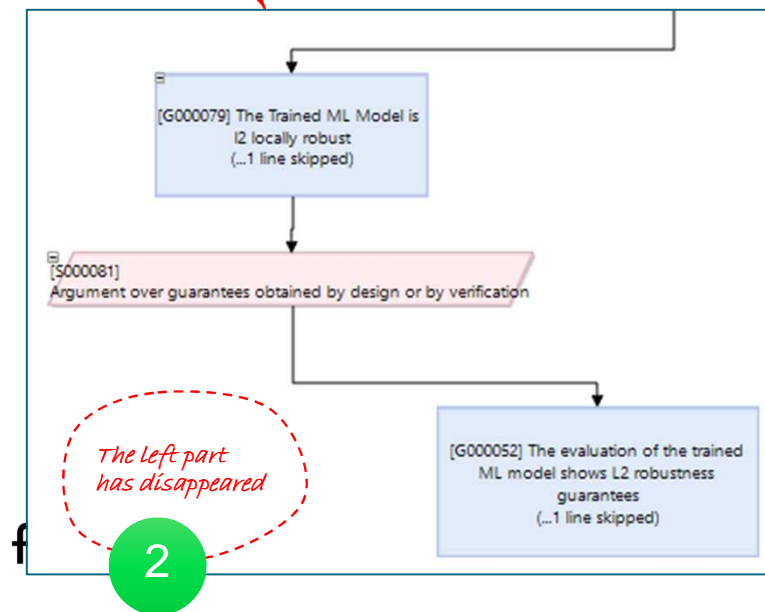Strategy pattern Process-based Vs. Product-based

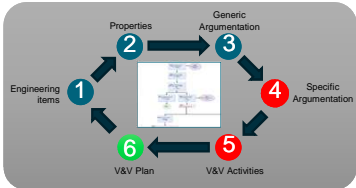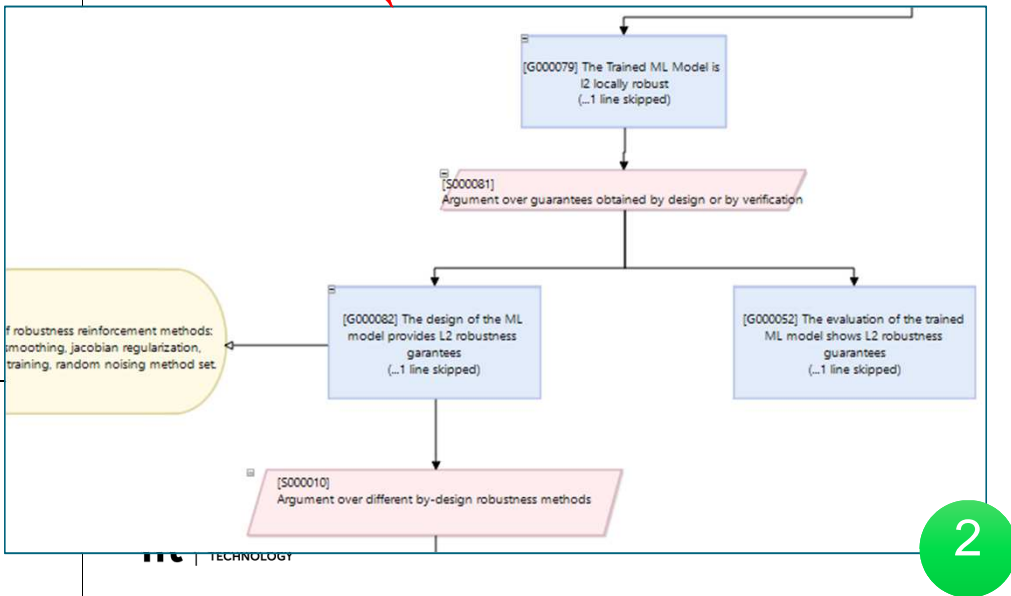☑ ✖ Property satisfied by design
☑ ✖ Property satisfied by verification

*User choice*

*Argumentation*

**1**

*Workflow*

**3**

[G000079] The Trained ML Model is l2 locally robust
(...1 line skipped)

[S000081]
Argument over guarantees obtained by design or by verification

f robustness reinforcement methods: smoothing, jacobian regularization, training, random noising method set.

[G000082] The design of the ML model provides L2 robustness garantees
(...1 line skipped)

[G000052] The evaluation of the trained ML model shows L2 robustness guarantees
(...1 line skipped)

[S000010]
Argument over different by-design robustness methods

**2**

👤 ML-Algorithm Engineer

Specify and architect ML Model and its Complementary Software Items

Analyze results of evaluation of the trained ML Model (with its Complementary Software Items if needed)

Analyze ML Model robustness evaluation report

trained ML model robustne...

Specify and architect ML Model

Orient the typology of ML

Consider robustness in the orientati...

Define ML Model requirements

Define ML Model robustness requirements

Express each ML robustness requirement as...

Express ML robustness requirements as form...

Define a strategy for evaluation of ML Model and its Complementary Software Items

ML model robustness r...

type of ML mo...

Define strategy, methods and tools for evaluation of ...

decision to modify strat...

ML model robustness ...

Define a strategy for ML Model development

Define a strategy for ML model robustness development

type of ML mo...

type of dat...

type of data...

ML model robustness require...

strategy, methods and tools for ML model robustness e...

Evaluate the trained ML Model and its Complementary Software Items

Evaluate robustness of the trained ML Model

Acquire or develop ML

Acquire ML Model

Develop ML Model

strategy for robustness by design (Lipschitz and/or adversarial traini...

trained ...

trained ...

👤 Data Engineer

Produce datasets

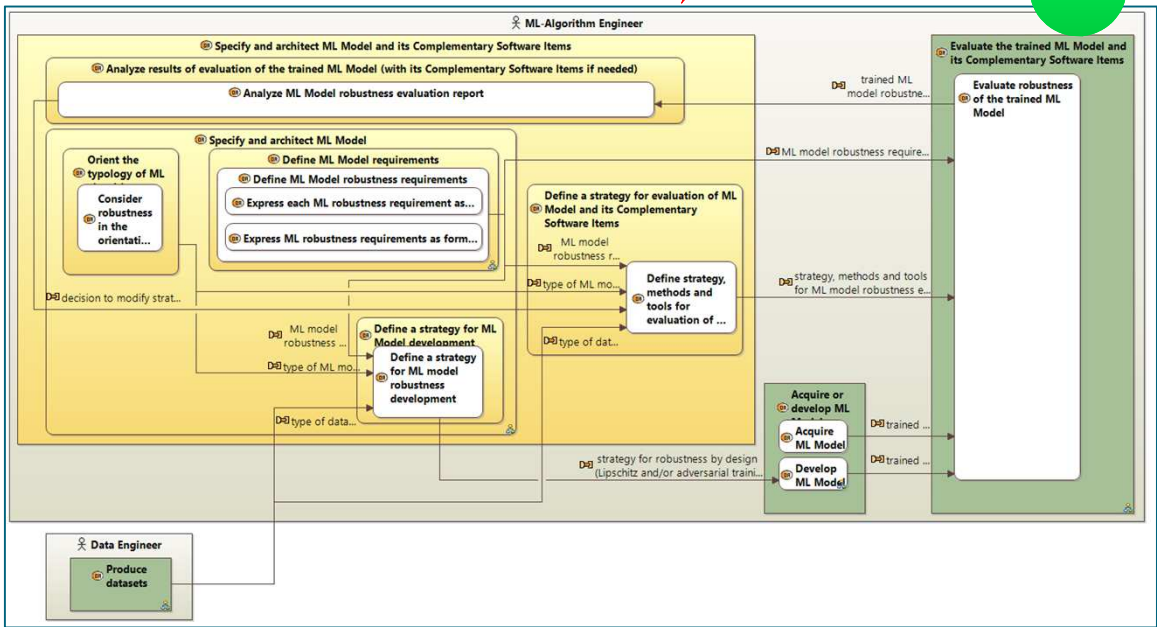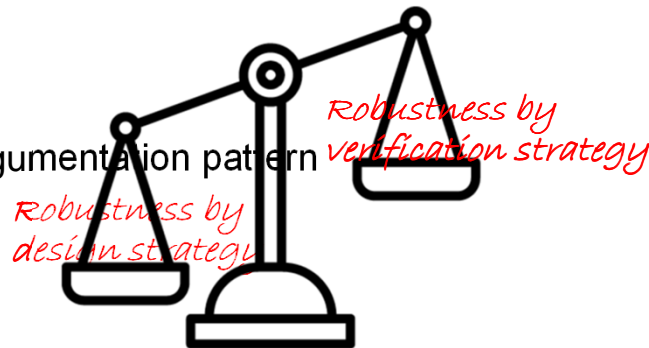TECHNOLOGY

# How to choose the argumentation strategy...

- Focus the validation effort on the most sensitive parts of the argumentation

  - Identify insufficiently convincing strategies associated to a goal

  - Identify contradiction between proof elements

  - Improve the argumentation

$\Rightarrow$ Estimate level of confidence in the argumentation pattern

*Robustness by verification strategy*

*Robustness by design strategy*

12/12/2025

# What does confidence mean in our framework ?

- Level of confidence ≈ Amount of information to justify a judgment about a proposition or, reciprocally, level of uncertainty about a judgment
  - Choice of an uncertainty representation
  - Elicitation of uncertainty associated to atomic elements
  - Propagation of uncertainty through the AC

- Complete information consists of what is known, and what is unknown (uncertainty/ignorance) about a proposition A:

$$Conf(A) + Uncer(A) = 1$$

Uncertainty is a general description of a state of knowledge that makes it difficult/impossible to assess the truth or the falsity of a piece of information (or a proposition).
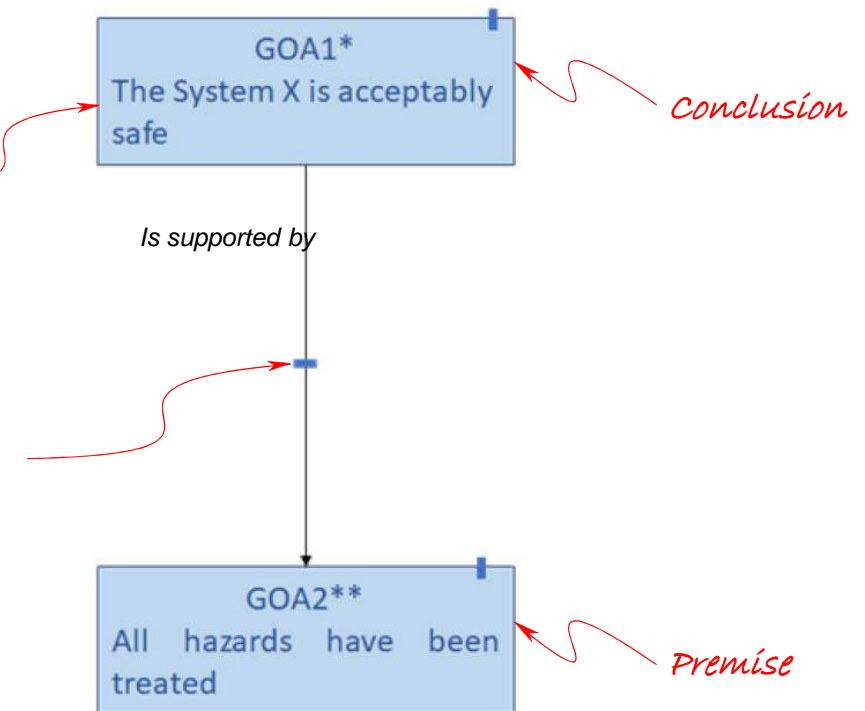
# Sources of uncertainty in AC

- **Two aspects to estimate uncertainty**

  - **Trustworthiness** which quantifies the truth (with belief measures) and the falsity (with disbelief measures) in propositions (i.e., goals).
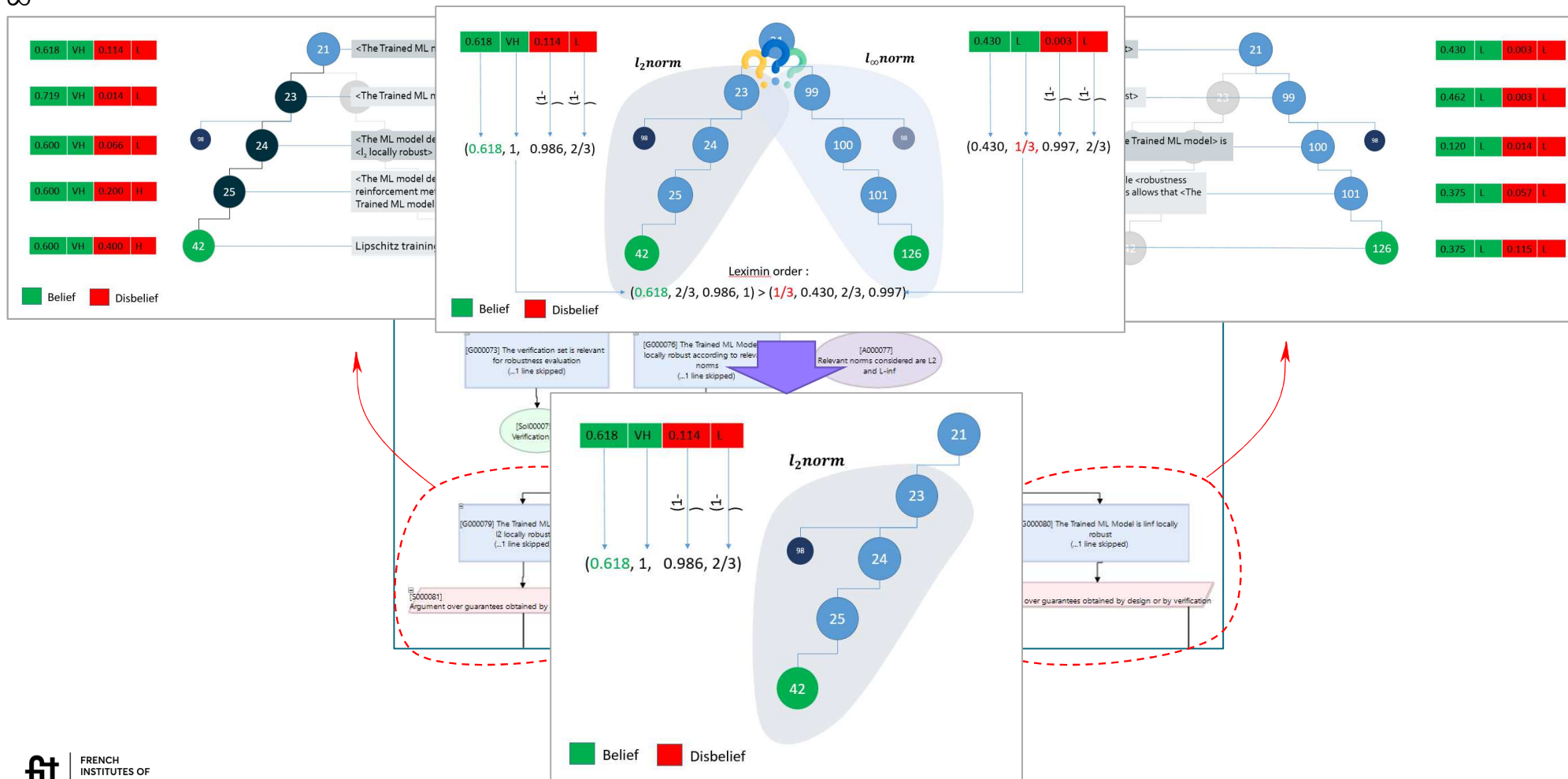
$$Conf(G) = Bel(G) + Disb(G)$$

  - **Appropriateness** which quantifies the truth about the inference (i.e., "supported by" relation) between a parent goal and its child goal(s).

GOA1*
The System X is acceptably safe

*Conclusion*

*Is supported by*

GOA2**
All hazards have been treated

*Premise*

FRENCH
INSTITUTES OF
TECHNOLOGY

ERTS – Uncertainty in Assurance Case Template for Machine Learning

# Status

- A method to link design and argumentation
- A tool (Capella plugin using pure::variant) to implement the method
- A method to evaluate confidence in the argumentation

V. Mussot *et al.*, 'Assurance Cases to face the complexity of ML-based systems verification', in *Embedded Real Time System Congress, ERTS'24*, Toulouse, France, June 2024. Accessed: Sept. 03, 2025. [Online]. Available: https://hal.science/hal-04588599

Y. I. Messaoud, J.-L. Farges, E. Jenn, and V. Mussot, 'Uncertainty in Assurance Case Pattern for Machine Learning', in *Embedded Real Time System Congress, ERTS'24*, Toulouse, France, June 2024. [Online]. Available: https://hal.science/hal-04584490v1/document

12/12/2025

FRENCH INSTITUTES OF TECHNOLOGY

12/12/2025



Merci pour votre *attensssssion!*