

Des patrons pour la certification d'IA embarquable

Journée IE + INFORSID

29/03/2024

Bernard BOTELLA ⁽¹⁾, Florent CHENEVIER ⁽²⁾, Stephen CREFF ⁽³⁾, Jean-Loup FARGES ⁽⁴⁾, Anthony FERNANDES PIRES ⁽⁴⁾
Ramon CONEJO LAGUNA ⁽⁵⁾, Eric JENN ⁽⁵⁾, Florent LATOMBE ⁽⁶⁾, Yassir ID MESSAOUD ⁽³⁾, Vincent MUSSOT ⁽⁵⁾

⁽¹⁾ CEA, ⁽²⁾ Thales AVS, ⁽³⁾ IRT SystemX, ⁽⁴⁾ ONERA, ⁽⁵⁾ IRT Saint Exupéry, ⁽⁶⁾ Obeo

Agenda



- Context and Objectives
- The process: Assurance cases and the ML development workflow
- Uncertainty assessment
- Conclusion



Context and Objectives

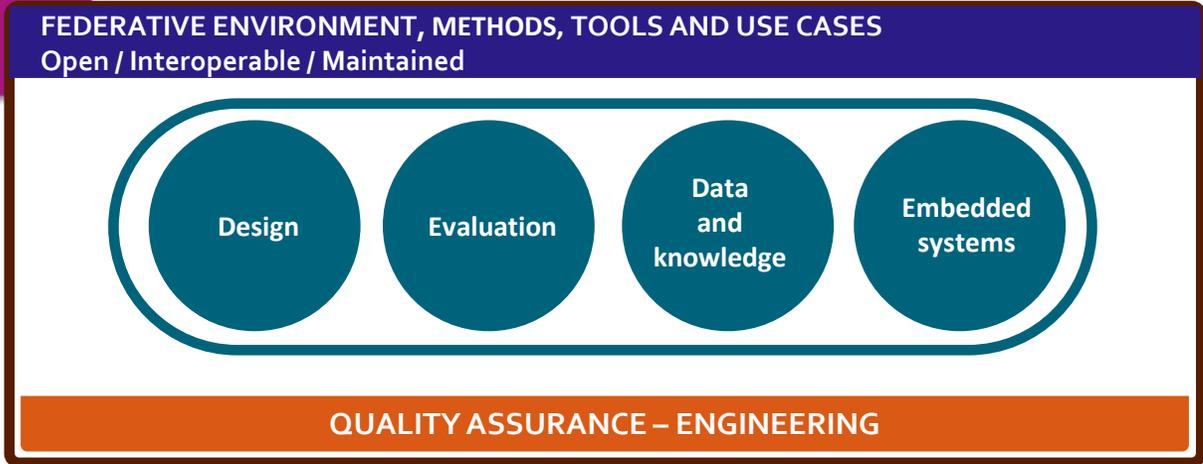
Context

The Confiance.ai Programme (www.confiance.ai)



45M€ budget, 5years (cur. Year 3)

To provide industrial companies with solutions enabling the development of new products and services based on trustworthy AI



Context

The Confiance.ai Programme (www.confiance.ai)



- Program structure : 7 Engineering Challenges, 2021 => 2024

EC	Adressed Topic
EC#1	Integration & Use-Cases, (+ Trusted AI Devops environment)
EC#2	Process, methodology & Guidelines
EC#3	Characterization & Qualification of Trustworthy AI
EC#4	Design for Trustworthy AI @ Algo, Components & System levels
EC#5	Data, Information & knowledge engineering for trusted AI
EC#6	IVV&Q strategy toward homologation/certification
EC#7	Target Embedded AI

We are here

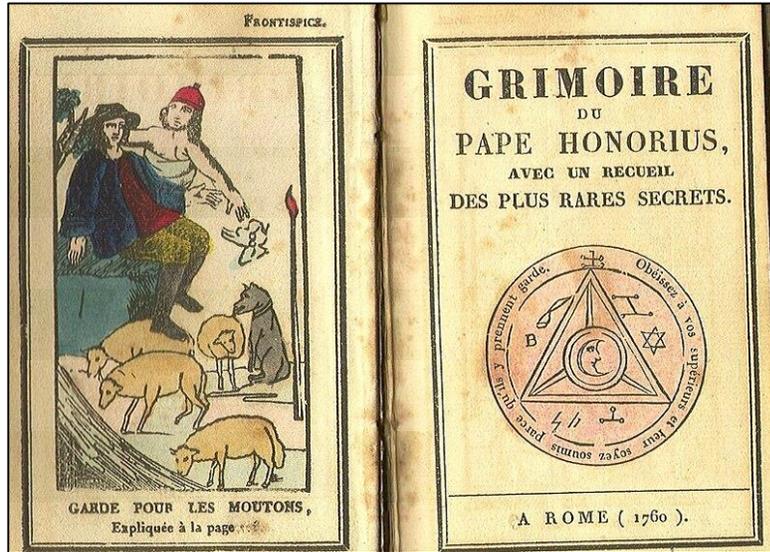
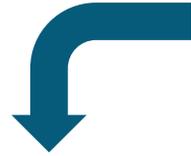
Context

Assurance Cases

Easy path #1

CONFIDENCE

Easy path #2



“Look at the book Chap. 3, Sec. 14, Vs 16”

*“Trussssssssst me...
Trussssssssst me...”*

Context Assurance Cases



1
Dependability

2
Confidence

3
Emerging technology.
No track record.

4
Assurance cases

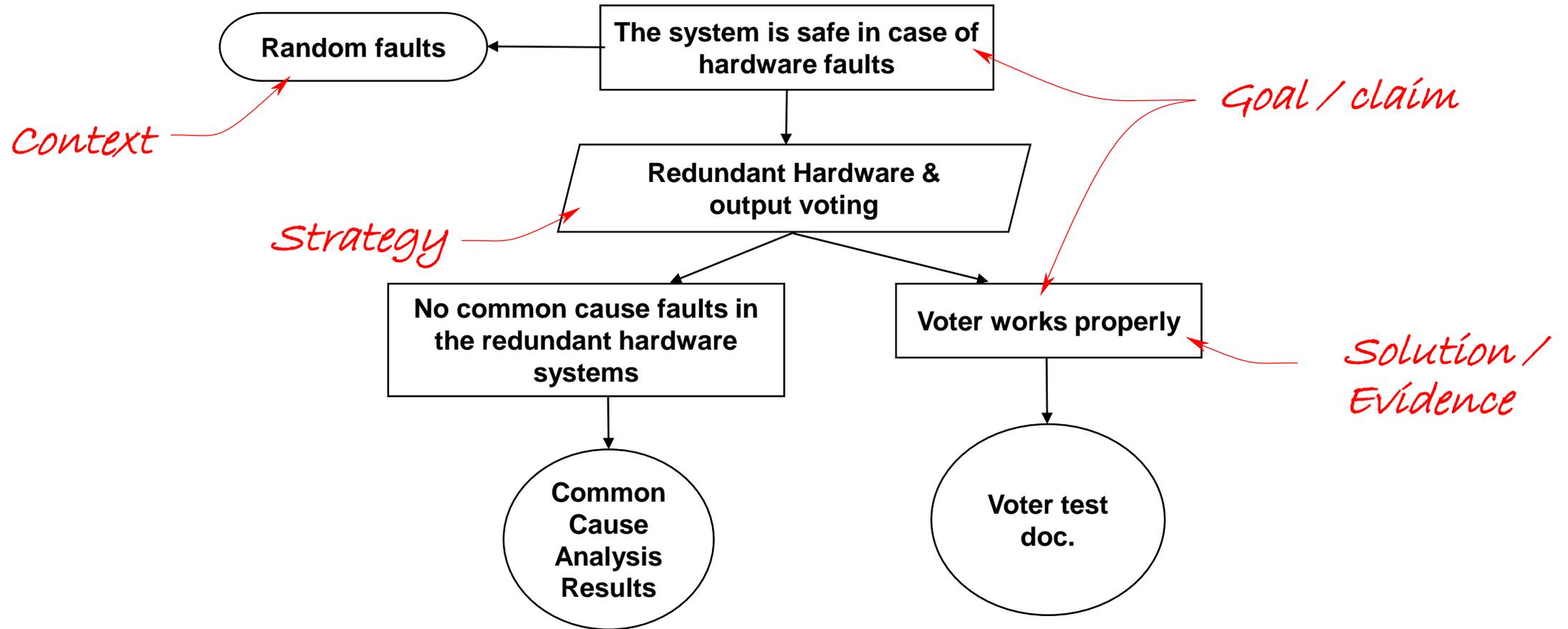
“the trustworthiness of a computer system such that reliance can **justifiably** be placed on the service it delivers” (W.C. Carter, in Laprie *et al.* “*Dependability: Basic Concepts and Terminology*”)

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

“This framework of **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)

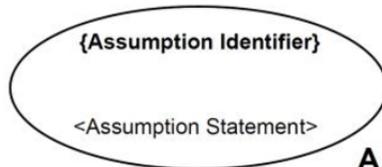
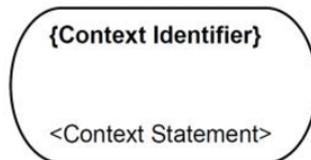
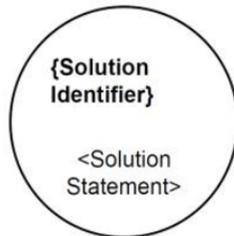
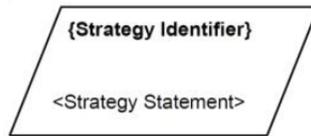
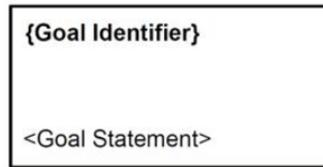
Assurance Cases

Main concepts



Assurance Cases

Main concepts



- **Goal** (& subgoals): affirmation that shall be assessed during the reasoning. Any goal may be refined in several subgoals.
- **Strategy**: justifies the decomposition of claims into sub-claims. It is an additional cue that helps the reader understand the form that an argument is going to take.
- **Solution**: A solution refers to some evidence that is deemed sufficient to establish the truth of the parent claim

- **Context**: define or constrain the scope over which the claim is made.
- **Justification**: describes why a given strategy is proposed as an approach to supporting a particular goal, or provide reasons why the strategy being adopted is adequate.
- **Assumption**: statement about a property considered true. Assumptions must be valid for the related claim/strategy to be valid.

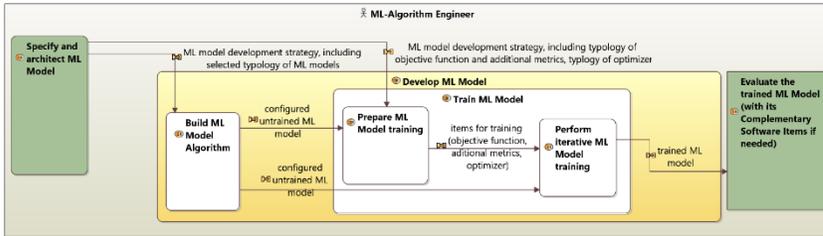
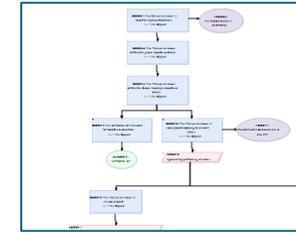


The process

From Engineering Items to
Assurance Cases

Process Overview

[G] The <Trained ML Model> is <robust>



Select an engineering item from the ML workflow



Select a property of interest



Get generic argumentation (AC)



Adapt the argumentation wrt context, cost, confidence, etc.



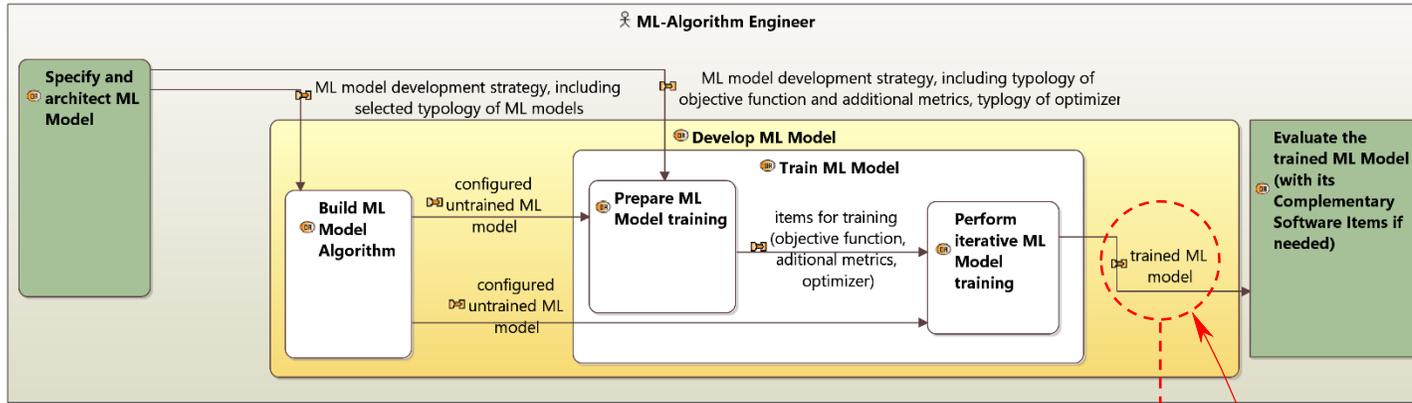
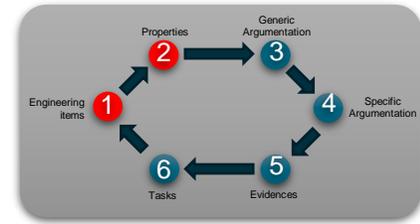
Identify evidences to produce



Update workflow

From Engineering Items to Assurance Cases

Robustness argumentation template



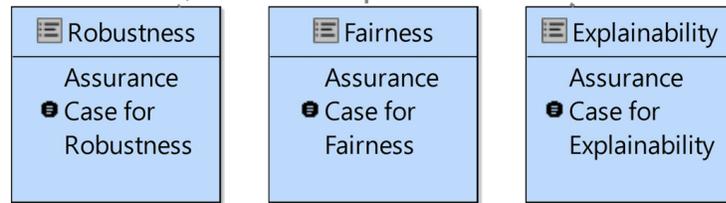
ML workflow

Trained ML model engineering / exchange item



Property

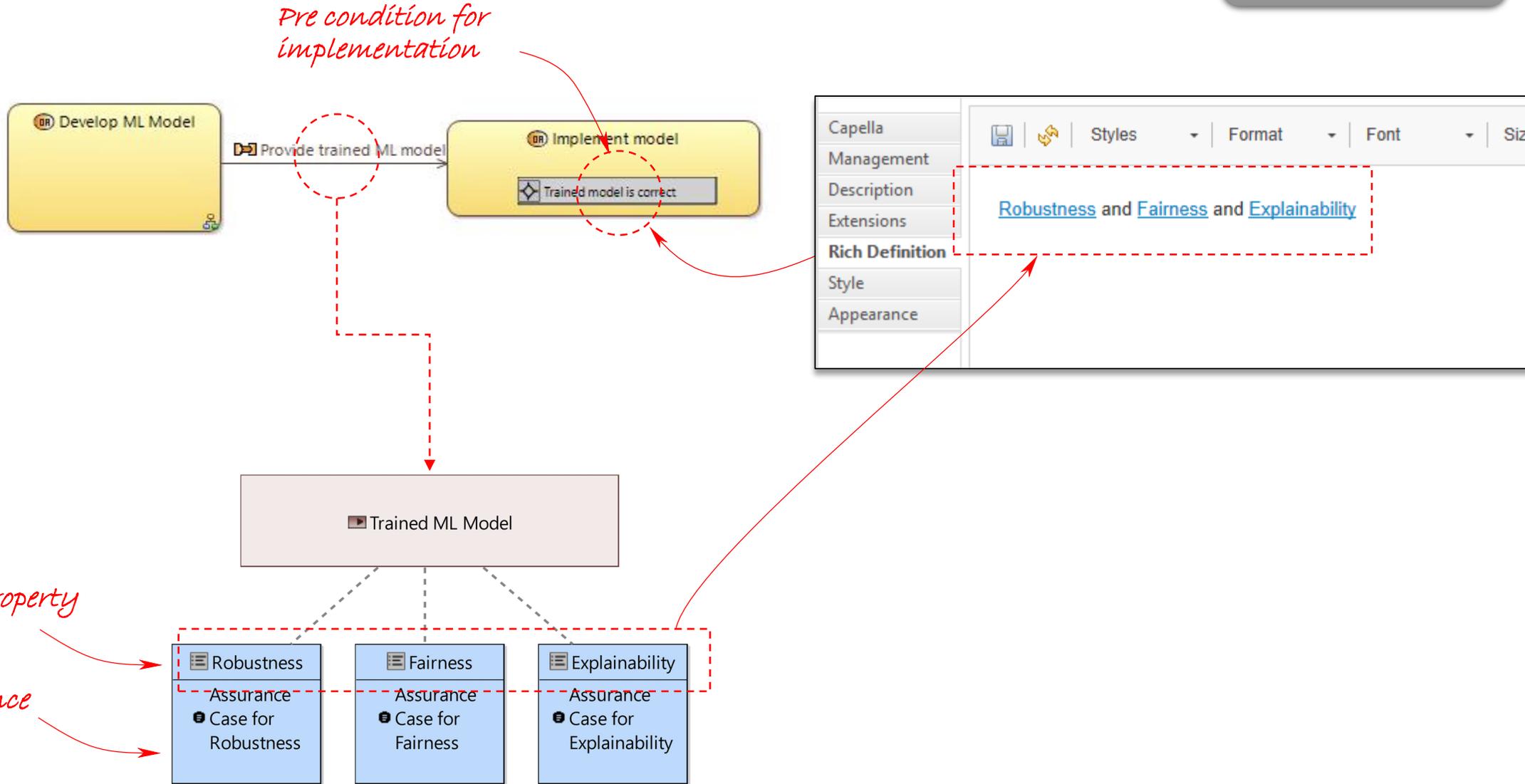
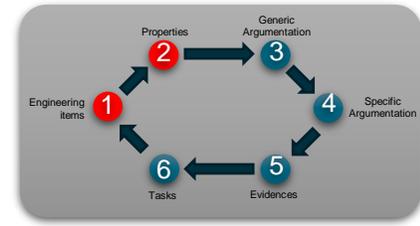
Assurance case



1. Associate Engineering Conditions (e.g. 'Robustness') to Exchange items
2. Define Assurance Cases for each condition

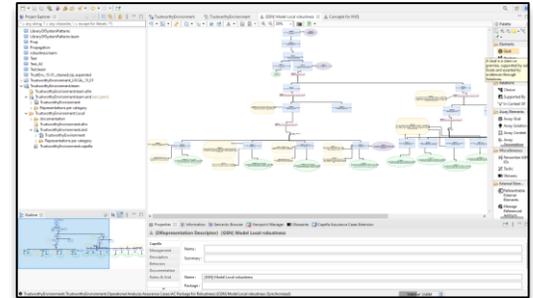
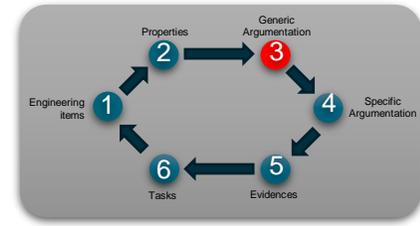
From Engineering Items to Assurance Cases

Robustness argumentation template

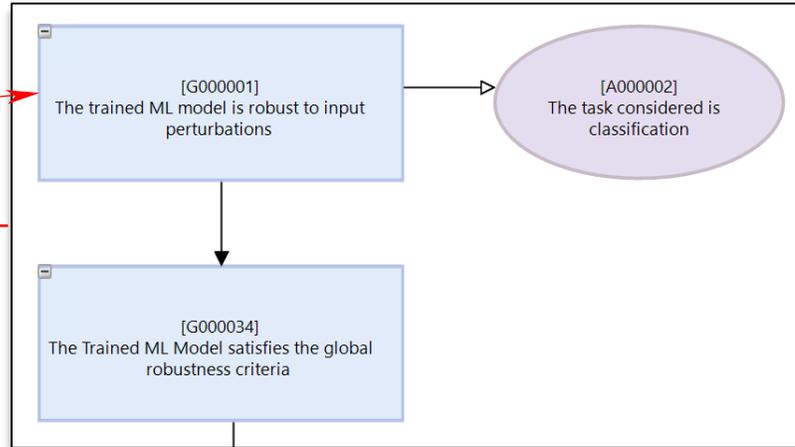


From Engineering Items to Assurance Cases

Robustness argumentation template



GSN extension in the Capella Editor



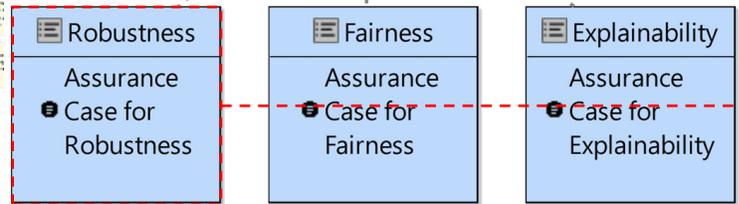
Property of interest

Robustness by design

Robustness by verification

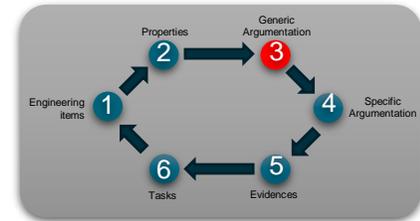
29/03/2024

Assurance case

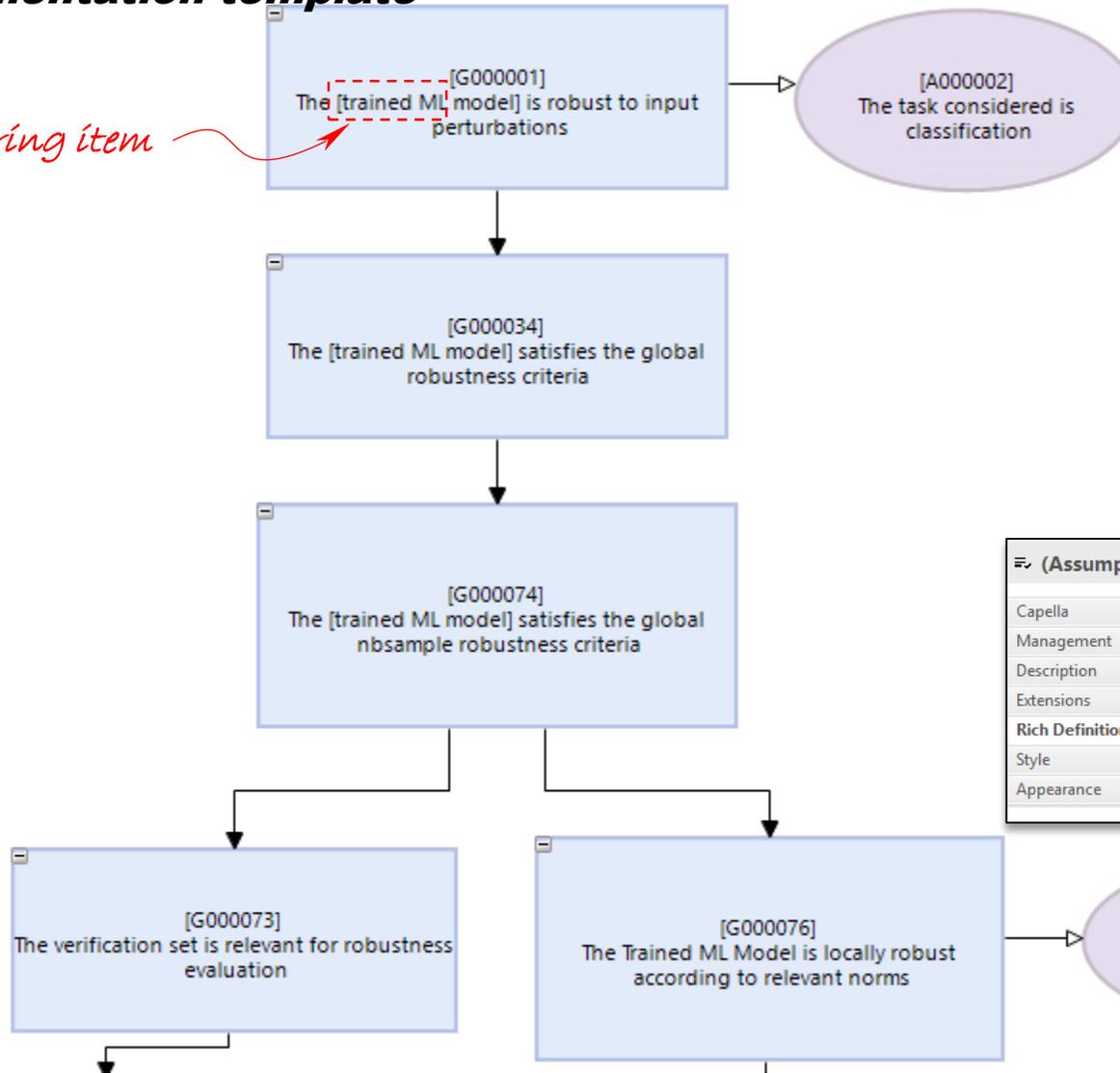


From Engineering Items to Assurance Cases

Robustness argumentation template



Engineering item



Glossaries

- Robustness Glossary
 - locally robust
 - local robustness
 - certified robust training
 - jacobian regularization
 - randomised smoothing
 - applicable
 - Trained ML model
 - L2 norm**
 - L-inf norm
 - Norm

Glossary

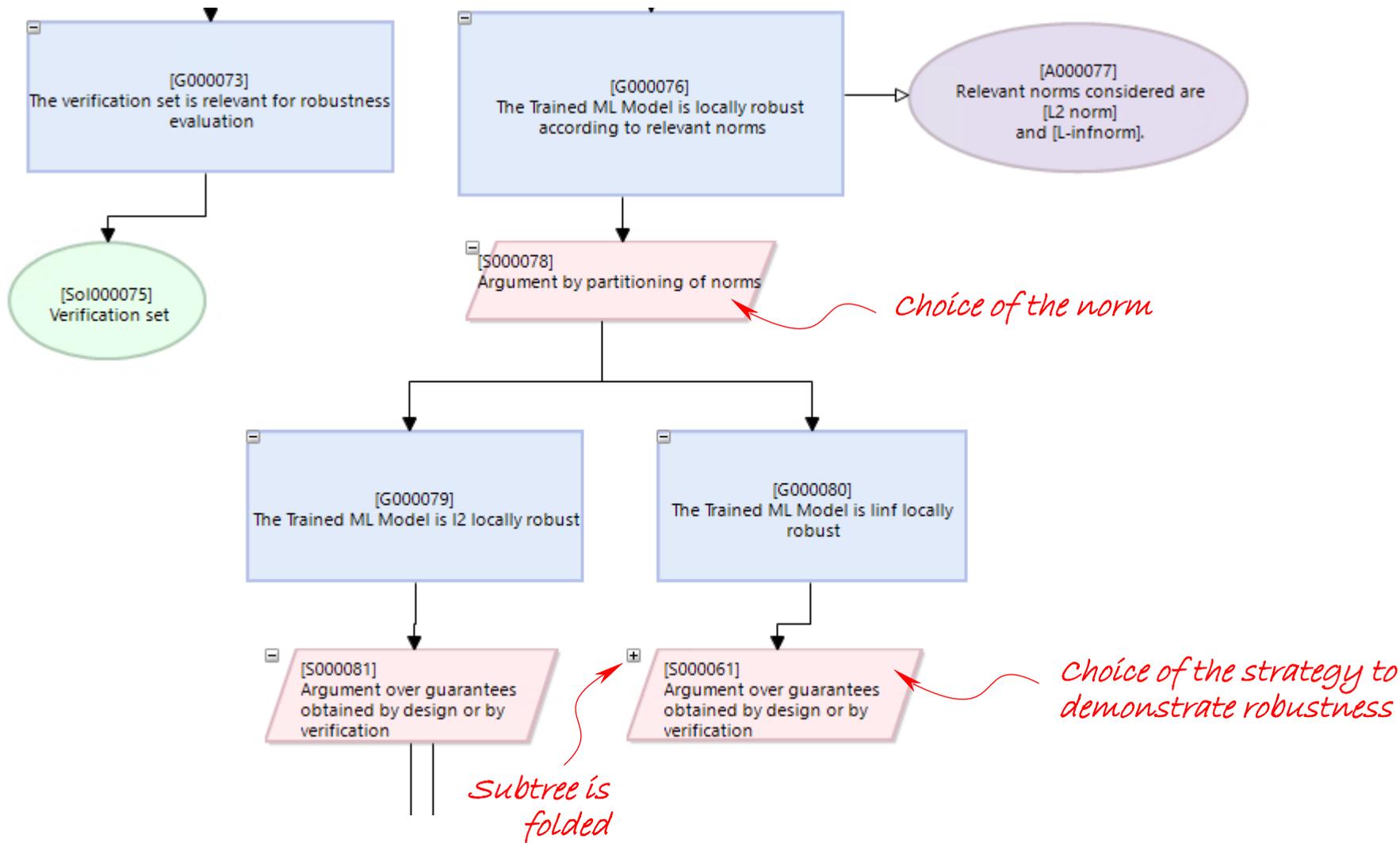
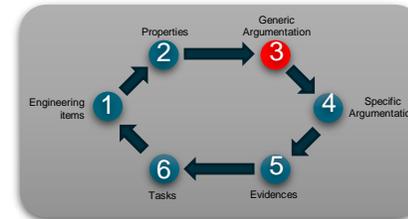
(Assumption)

Relevant norms considered are [L2 norm](#) and [L-inf norm](#).

Glossary term

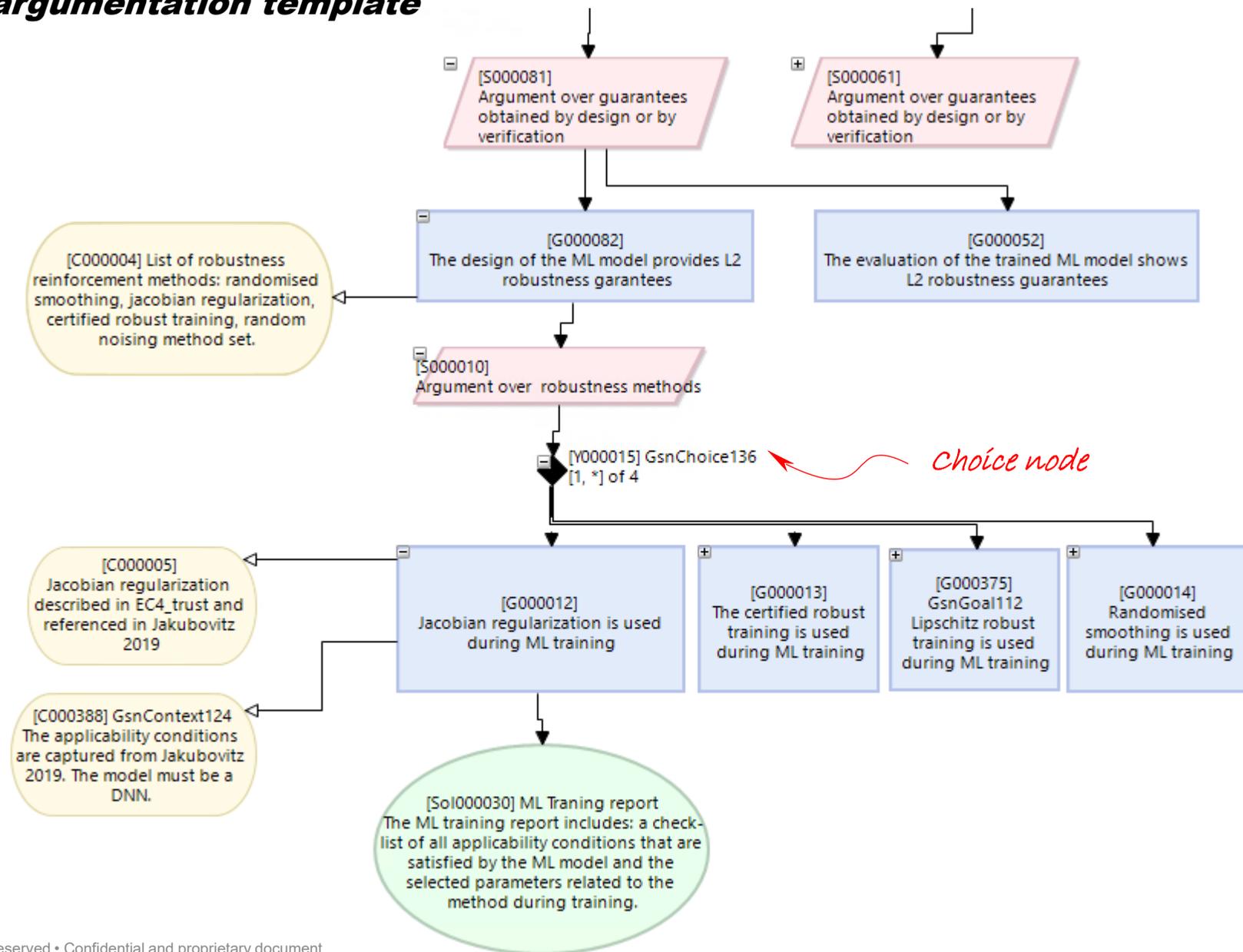
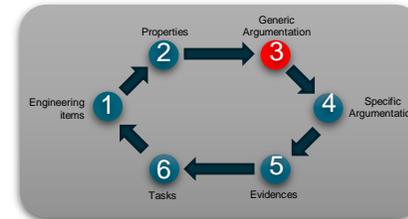
From Engineering Items to Assurance Cases

Robustness argumentation template



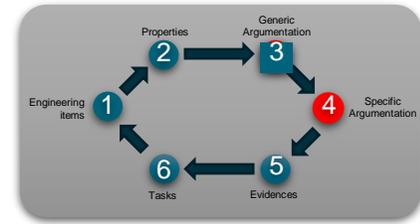
From Engineering Items to Assurance Cases

Robustness argumentation template



From Engineering Items to Assurance Cases

Refinement of requirements



1. Partitioning by robustness criteria

- Percentage of samples that are robust
- Maximal **lambda** for which all samples are robust
- Mean of maximal **lambda** for which each sample is robust

All choices

Robustness criteria

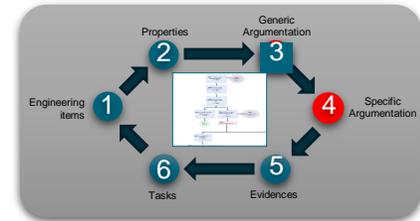
User choice

Configuration menu

Capella Pure::variant configuration wizard

From Engineering Items to Assurance Cases

Refinement of requirements



❑ Partitioning by norms (only l_2 and l_∞ considered)

Partitioning by robustness criteria

Local Robustness Norm Selection

Strategy pattern Process-based Vs. Product-based
Design Method

Local Robustness Norm Selection

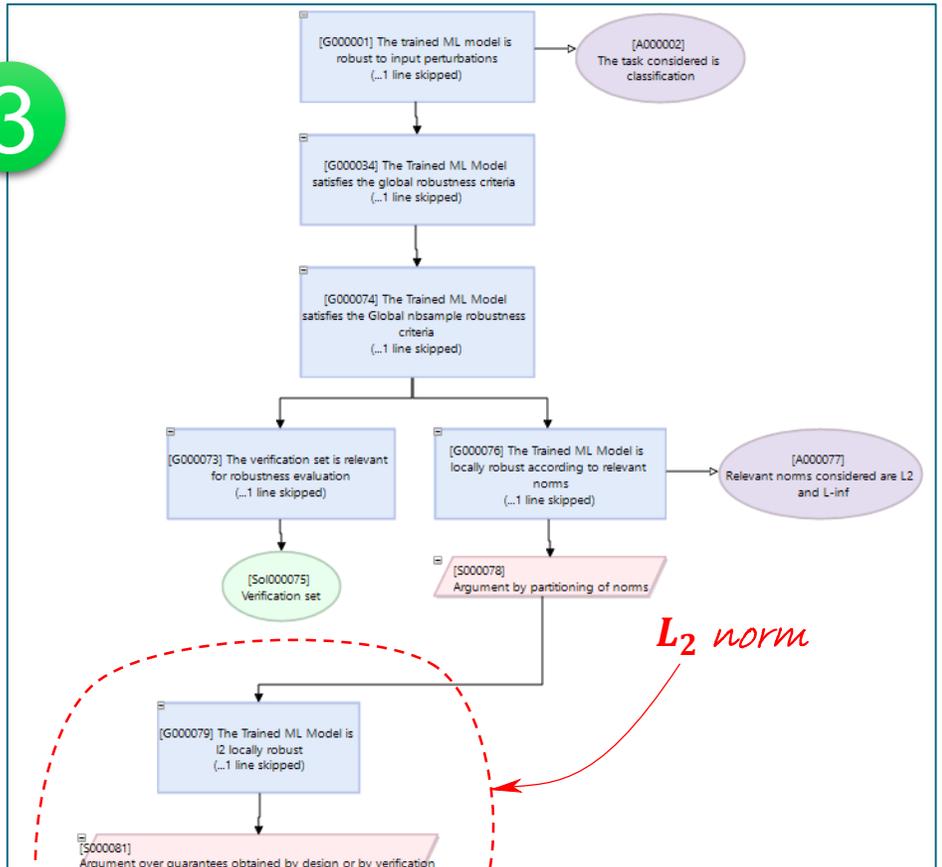
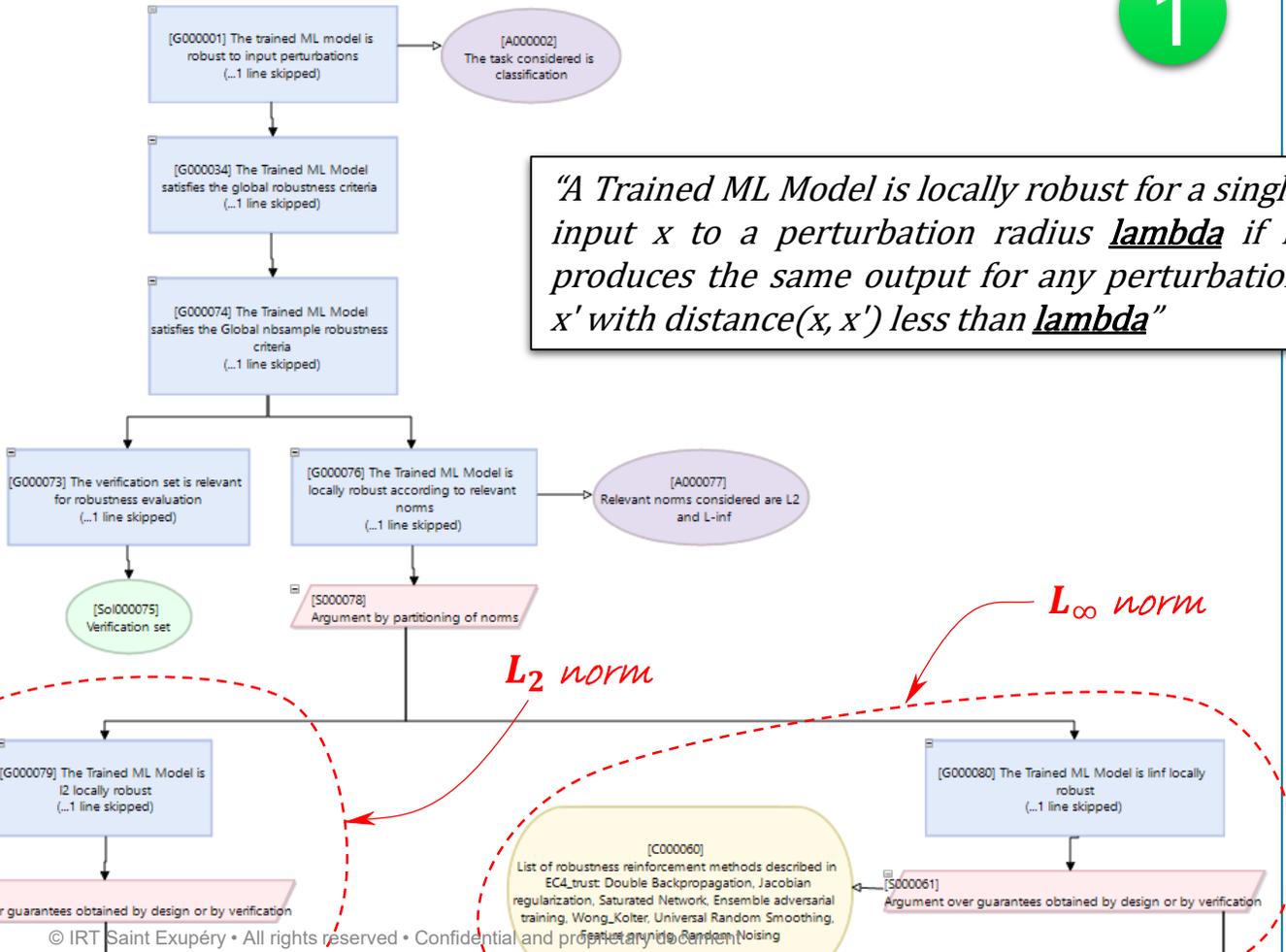
- User choice*
- l_2 locally robust
 - l_∞ locally robust

1

2

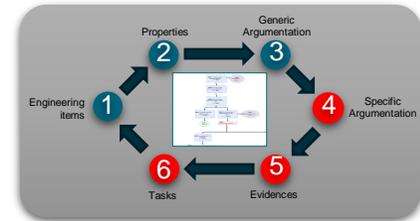
3

"A Trained ML Model is locally robust for a single input x to a perturbation radius λ if it produces the same output for any perturbation x' with distance (x, x') less than λ "



20/03/2024

From Engineering Items to Assurance Cases



□ 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

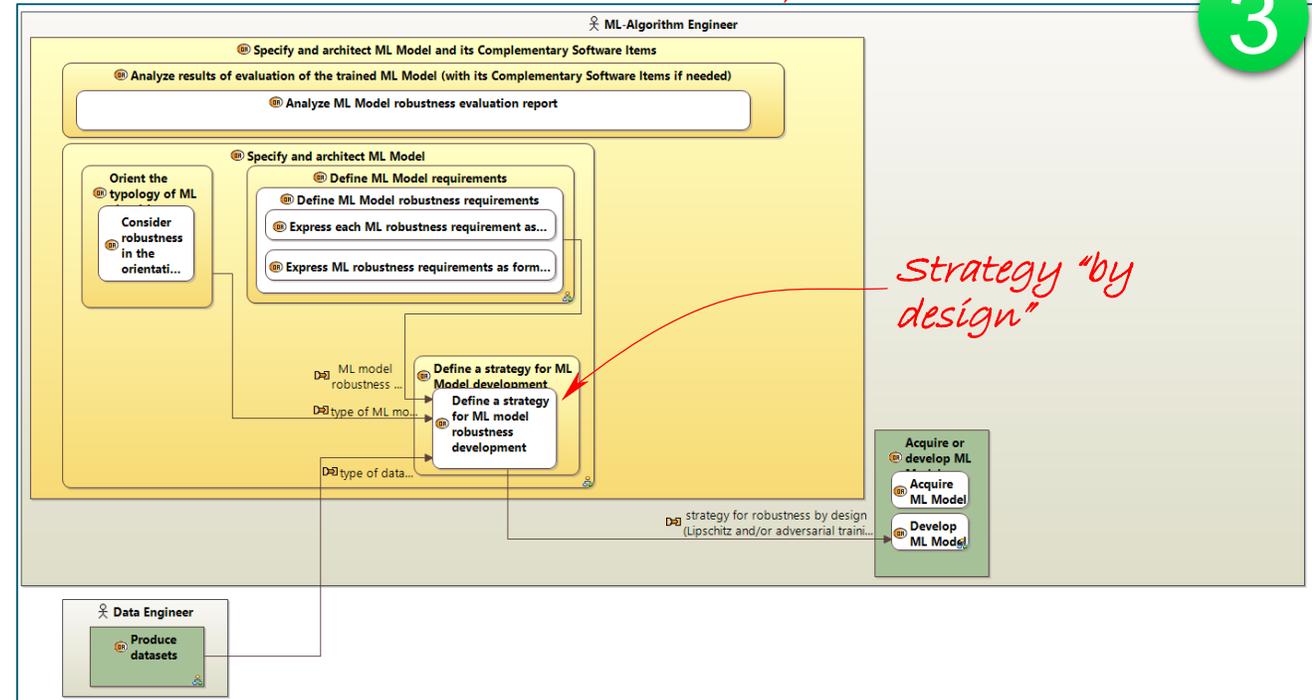
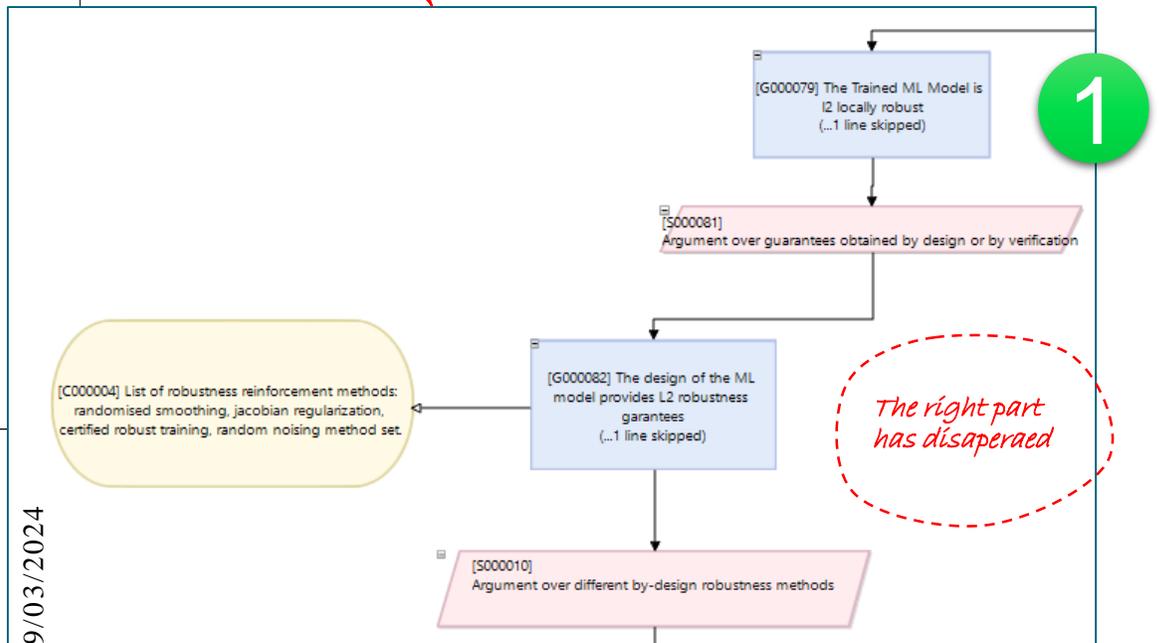
Argumentation

User choice

Workflow

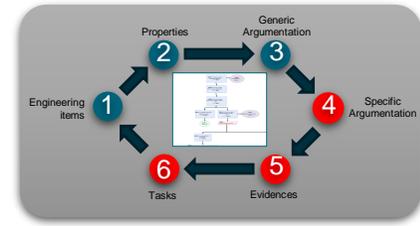
2

3



Strategy "by design"

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

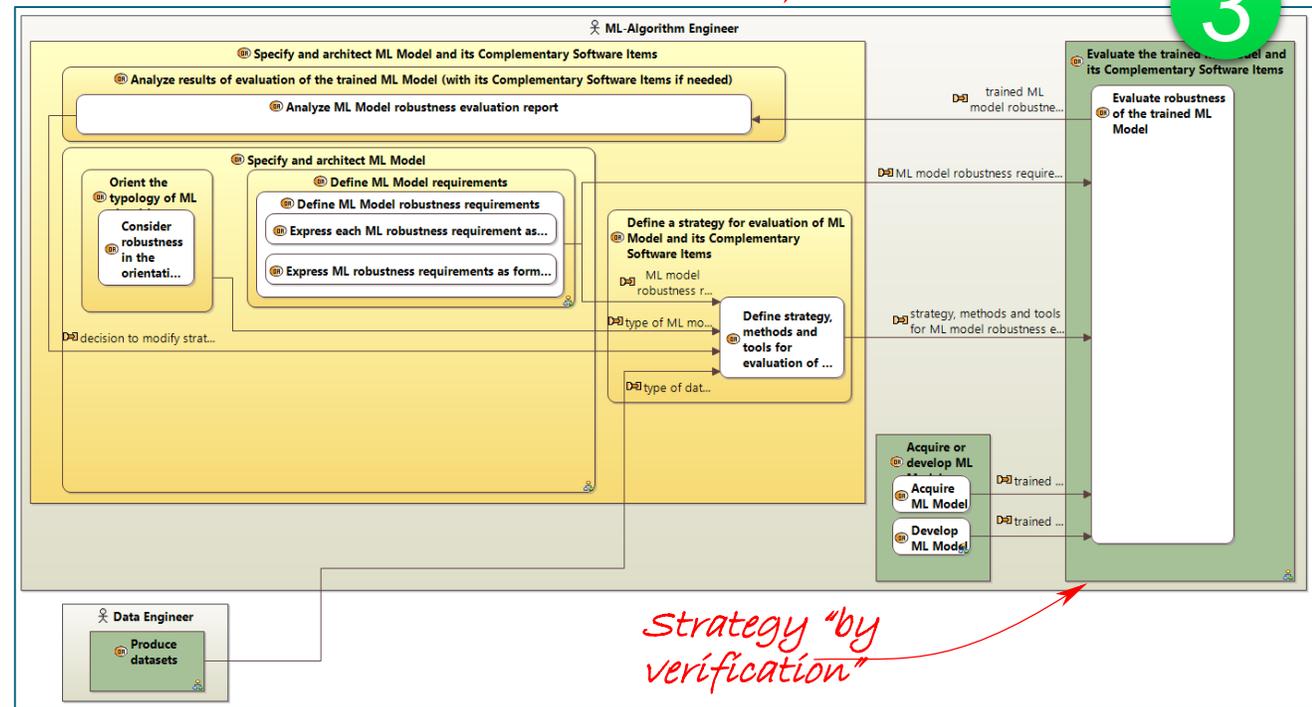
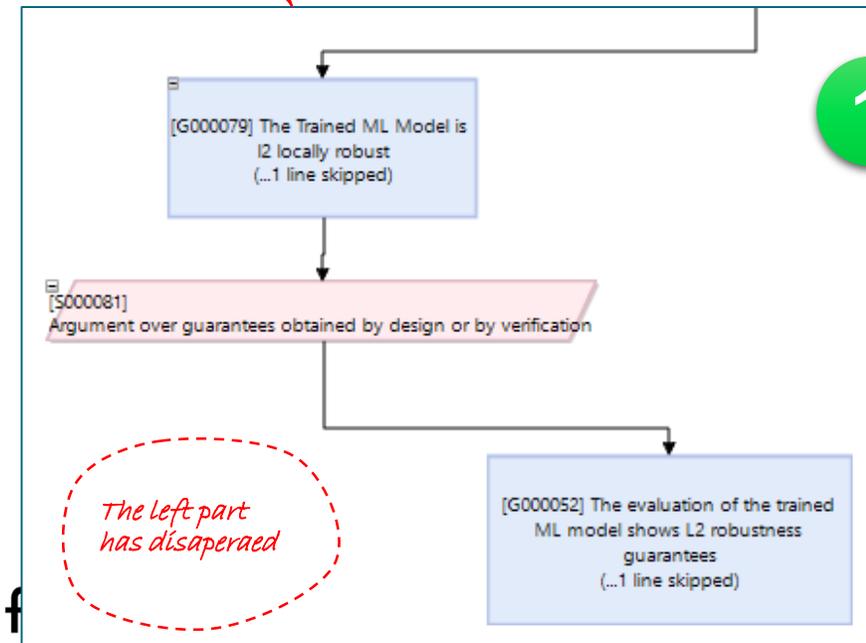
Argumentation

User choice

Workflow

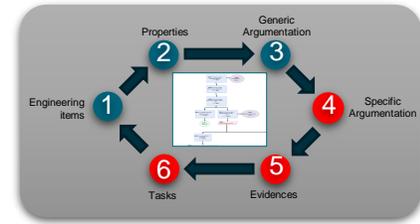
2

3



Strategy "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

Argumentation

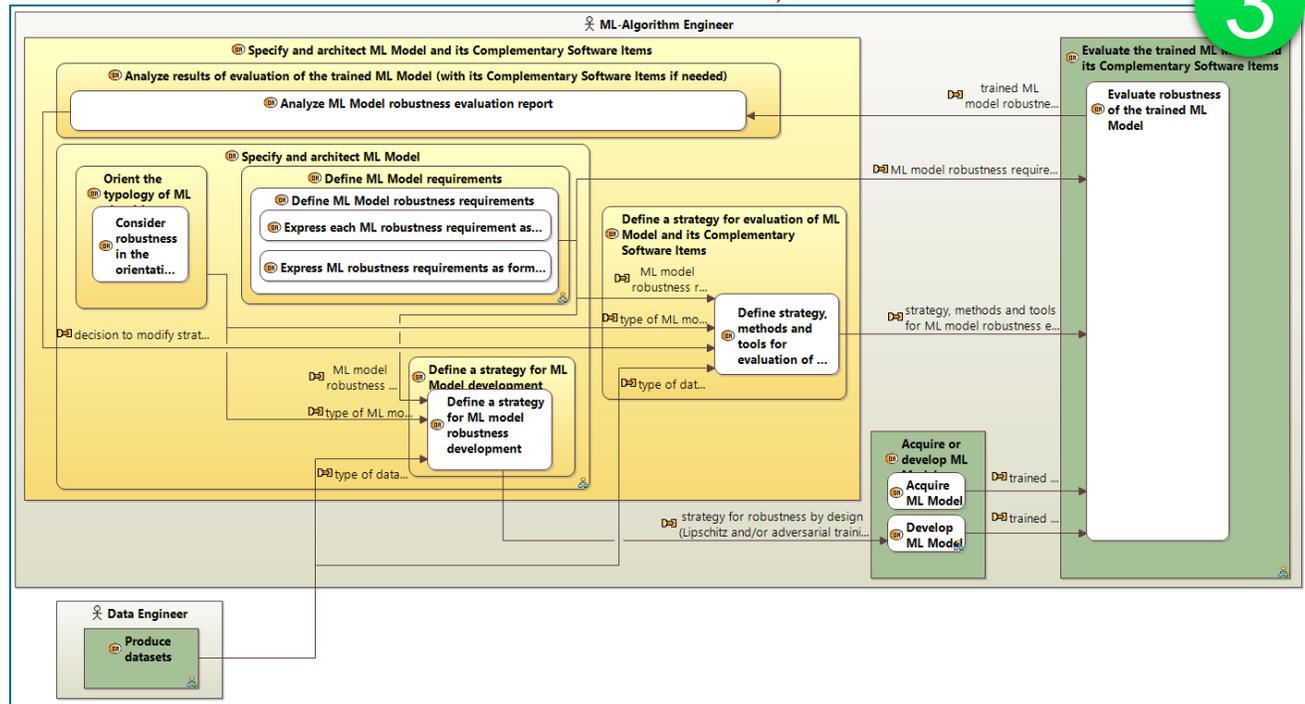
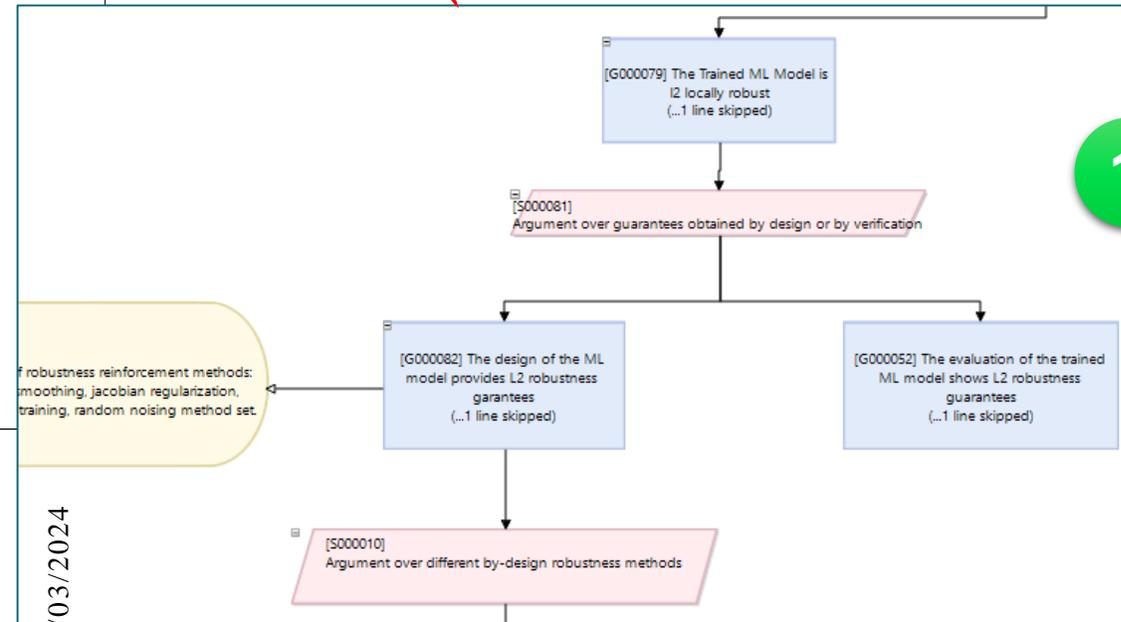
User choice

Workflow

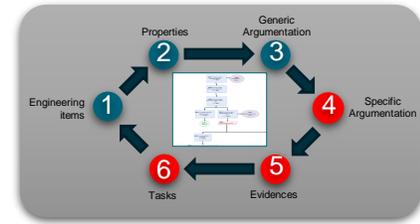
2

3

1



Robustness AC Template



- Families of method (from Confiance SotA: "EC4-Trustworthiness by design"):
 - <Jacobian regularization>, <Lipschitz training>, <Certified robust training>, <Randomised smoothing>, <Random noising>

Partitioning by robustness criteria

Local Robustness Norm Selection

Strategy pattern Process-based Vs. Product-based

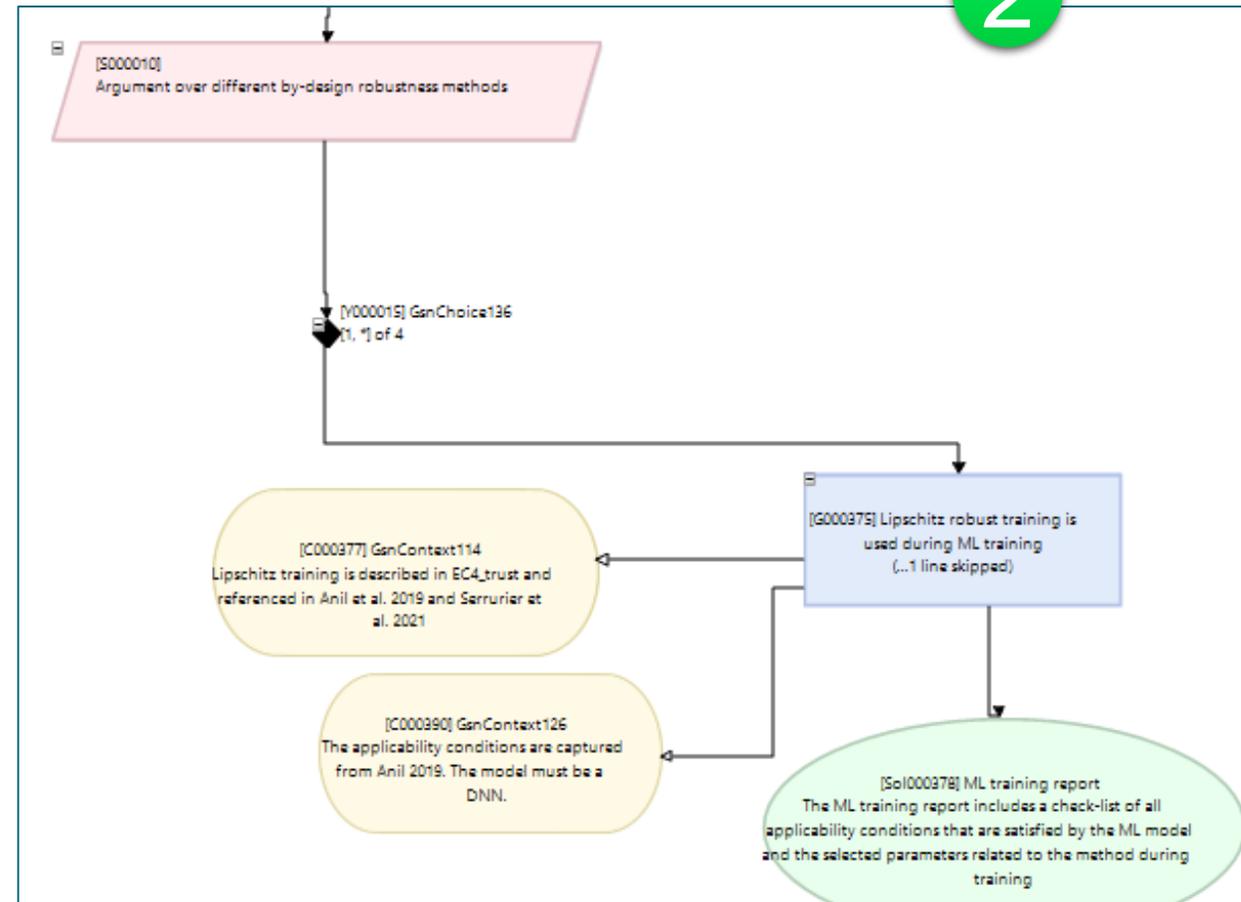
Design Method

1

	Design Method
<input type="checkbox"/>	<input checked="" type="checkbox"/> Jacobian regularization
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Lipschitz training
<input type="checkbox"/>	<input checked="" type="checkbox"/> Certified robust training
<input type="checkbox"/>	<input checked="" type="checkbox"/> Randomised smoothing
<input type="checkbox"/>	<input checked="" type="checkbox"/> Random noising

user choices

2





Uncertainty Assessment & Choice of Strategies

Using Dempster-Shafer theory...

What we would like...

- Choose the most convincing strategy
- Focus the validation effort on the most sensitive parts of the argumentation
 - Assessment performed at each goal provides
 - Goal weakness
 - Contradiction between proof elements
 - For conjunctions
 - Procedure to improve the AC
- Identify the weaknesses of AC structure
 - Not sufficiently convincing strategies associated to a goal whose

Uncertainty in the context of AC

How to establish confidence ?

- Use of assurance case to justify the well-founded development of systems integrating machine learning

What is an assurance case?

- A structured argument used to justify a desired claim (safe, reliable, robust ...), based on evidence(s) concerning both the system and the environment in which it operates.

Issue

- **What are the sources of uncertainty in a structured argument?**
- **How to measure and propagate uncertainty in these structures?**

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).

What does confidence mean in our framework ?

❑ The concept of “**Confidence**”, in our context (i.e., argumentation), reflects the amount of information an expert has that can justify his/her judgment about a proposition.

❑ A justification can be **for** or **against** a proposition.
Formally, it’s defined as:

$$Conf(A) = Bel(A) + Disb(A).$$

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

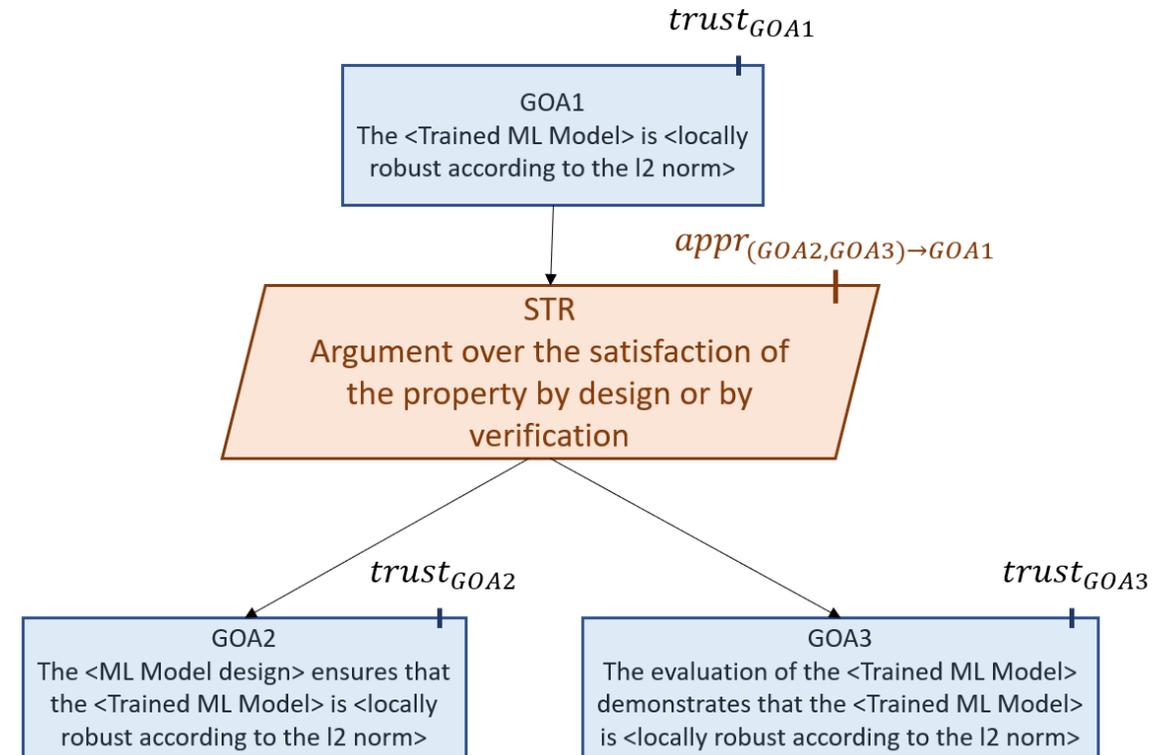
$$Conf(A) + Uncer(A) = 1.$$

Sources of uncertainty in AC



□ Two factor to estimate uncertainty

- **Trustworthiness** which quantifies the **truth** (with belief measures) and the **falsity** (with disbelief measures) in propositions (i.e., goals).
- **Appropriateness** which quantifies the **truth** about the inference (i.e., supported by relation) between a parent goal and its child goal(s). This is related to the strategy deployed by the AC designer to develop his/her reasoning.



- $trust_i \equiv (Bel_i, Disb_i, Uncer_i), i = \{GOA1, GOA2, GOA3\}$
- $appr \equiv (Bel_{(GOA2,GOA3) \rightarrow GOA1}, Uncer_{(GOA2,GOA3) \rightarrow GOA1})$

Source of uncertainty

- ❑ **Aleatoric** uncertainty (or Randomness) due to the variability of natural phenomena. E.g., rolling a dice.
- ❑ **Epistemic** uncertainty (or Incompleteness) due to lack of information. E.g., “The crime suspect fled in a grey car”. This information is not that sufficient to track down the suspect. What kind of car was it? In which direction did he/she flee?
- ❑ **Inconsistency** due to misinformation and contradiction. E.g., pro- and anti-vaccine arguments in a global pandemic situation.
- ❑ **Fuzziness** or **vagueness** due to imprecise information. E.g., Pierre is tall. The borderline between “tall” and “not tall” is not well-defined.

Measuring uncertainty – Probability Theory

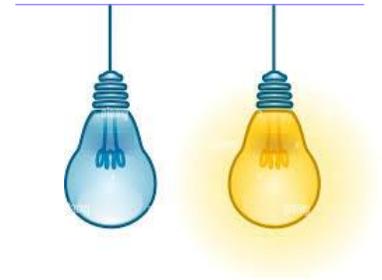
□ **Probability theory** deals well with random events (frequencies), but less well with singular events due to a lack of information.

□ It represents uncertainty by assuming even distribution over the whole frame of discernment Ω , such that: $P(\{\omega_i\}) = \frac{1}{|\Omega|}$.

▪ Example: Case of a light bulb, $\Omega = \{On, Off\}$

▪ I have no idea of the state of the light bulb: $P(\{On\}) = P(\{Off\}) = \frac{1}{2}$

▪ There's an equal chance of the light bulb being on or off: $P(\{On\}) = P(\{Off\}) = \frac{1}{2}$



Both situations are described using the same model

Measuring uncertainty – Dempster-Shafer Theory

□ Dempster-Shafer theory (DST) is a generalization of probability theory that deal well with both **epistemic** and **aleatory** uncertainties.

□ It defines the concepts of: **Mass function** (BPA) $m: 2^\Omega \rightarrow [0,1]$ such that:
 $\sum_{E \subseteq \Omega} m(E) = 1.$

□ Example: Case of a light bulb, $\Omega = \{On, Off\}$

- $m(\{On\})$: Quantifies the probability that the light bulb is “On”.
- $m(\{Off\})$: Quantifies the probability that the light bulb is “Off”.
- $m(\Omega)$: Quantifies **ignorance** on the state of the light bulb “On” or “Off”.
- $m(\emptyset)$: Quantifies **contradiction**. I.e., “On” and “Off” at the same time.



Uncertainty



Measuring uncertainty – Dempster-Shafer Theory

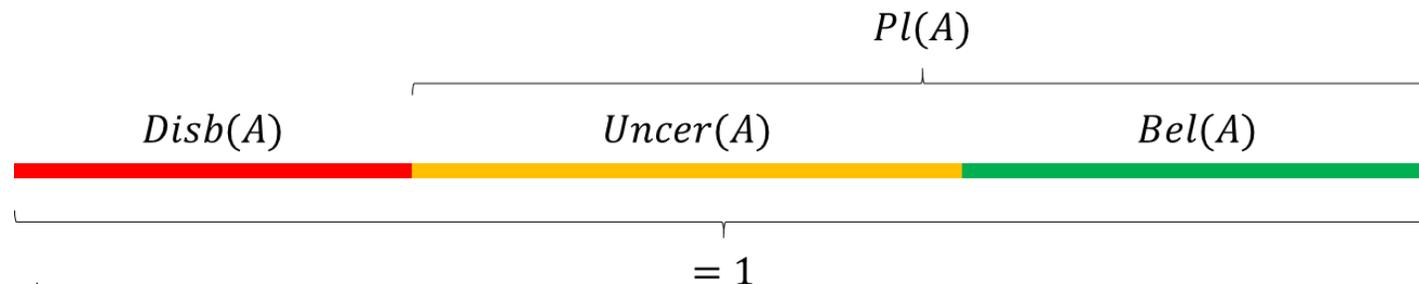
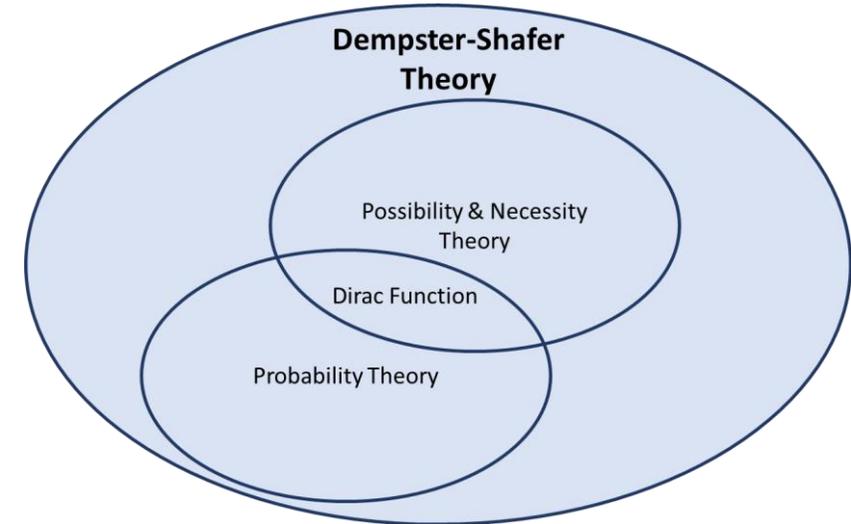
□ From a mass function, we define the concepts of:

- Belief function:

$$Bel(A) = \sum_{E \subseteq A, E \neq \emptyset} m(E) \text{ and } Disb(A) = Bel(\bar{A})$$

- Plausibility function:

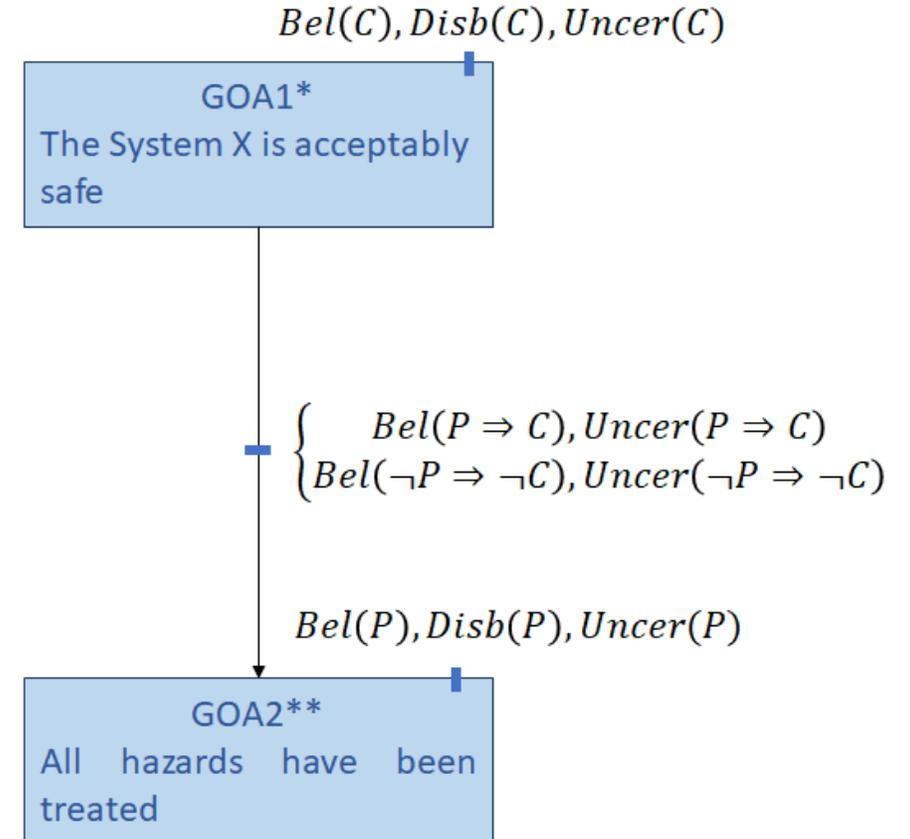
$$Pl(A) = \sum_{E \cap A \neq \emptyset} m(E) = 1 - Disb(A)$$



Uncertainty Evaluation – Mathematical Background

□ Hypothesis

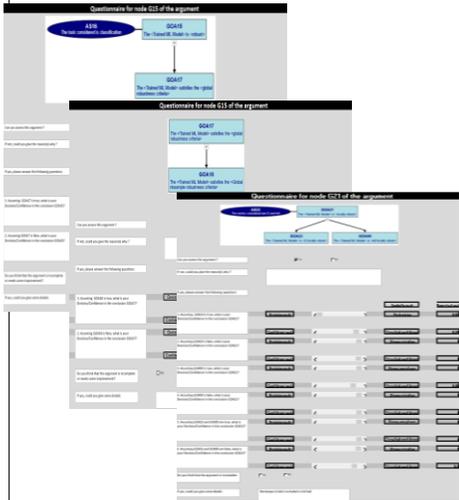
- **Goals** directly supported by a **Solution** can be:
 - Believed, i.e. $m(g)g$ can be different from zero
 - Disbelieved, i.e. $m(g)not\ g$ can be different from zero
 - Epistemically uncertain, i.e. $m(g)g$ or $not\ g$ can be different from zero
- Rules involved in a **Strategy** can be:
 - Believed, i.e. $m(r)r$ can be different from zero
 - Epistemically uncertain, i.e. $m(r)r$ or $not\ r$ can be different from zero
 - But cannot be disbelieved, i.e. $m(g)not\ g = 0$
- Rules are: $pi \Rightarrow C$, $not\ pi \Rightarrow not\ C$
 - Provides a formal and flexible definition of *Is* supported by



*GOA1 = C (Conclusion)

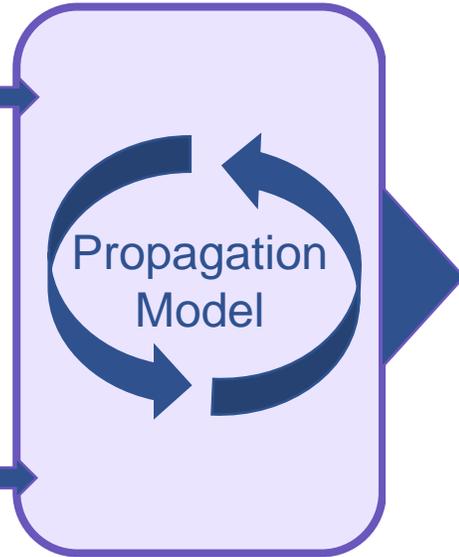
**GOA2 = P (Premise)

Confidence & Uncertainty in DST/AC framework

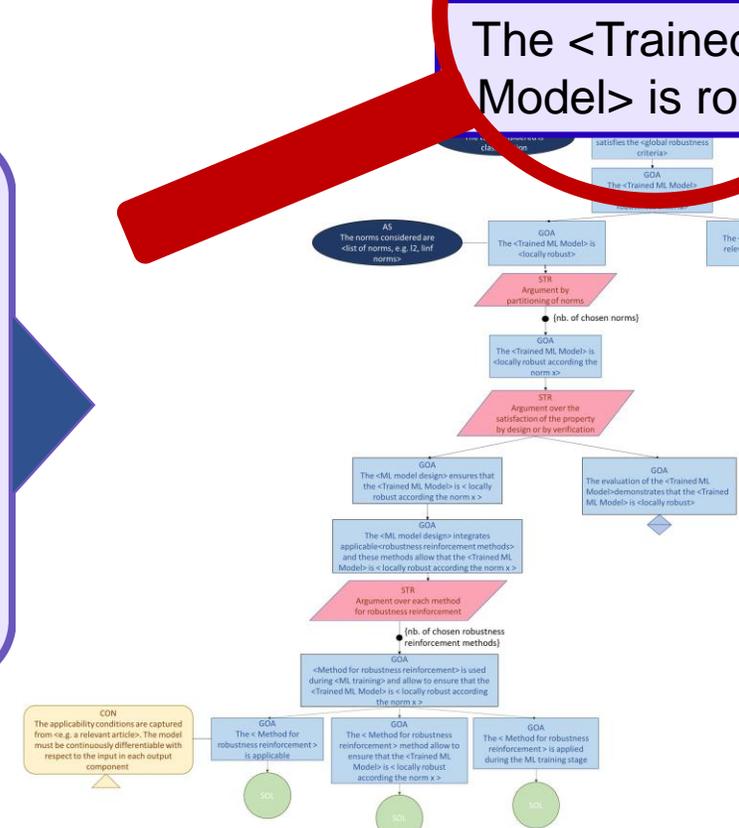
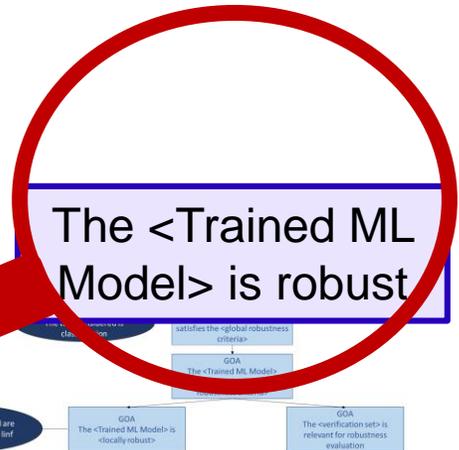


Belief values on rules (i.e., appropriateness)

Belief and disbelief values on premises (i.e., trustworthiness)



Belief and disbelief values on "Top-goal" (i.e., trustworthiness)

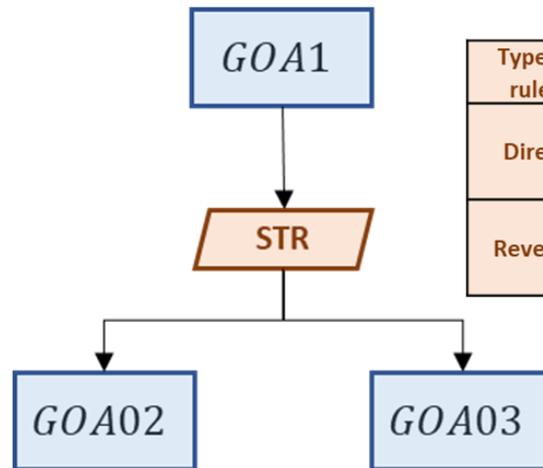


Uncertainty metrics

Visualization format

□ Uncertainty metrics are displayed in terms of belief and disbelief degrees.

Metric (GOA01)	Numerical value	Qualitative value
Belief degree	0.95	Very high
Disbelief degree	0.00	Very low
Conflict degree	0.00	Very low



Type of rules	Rules (STR)	Belief Values	Qualitative belief
Direct	All → GOA01	1.00	Very high
	GOA02 → GOA01	0.95	Very high
	GOA03 → GOA01	0.70	High
Reverse	None → ¬GOA01	1.00	Very high
	¬GOA02 → ¬GOA01	0.50	High
	¬GOA03 → ¬GOA01	0.50	High

Metric (GOA02)	Numerical value	Qualitative value
Belief degree	0.90	Very high
Disbelief degree	0.00	Very low
Conflict degree	0.00	Very low

Metric (GOA03)	Numerical value	Qualitative value
Belief degree	0.75	Very high
Disbelief degree	0.00	Very low
Conflict degree	0.00	Very low

+ *Nota.* Belief (resp. disbelief) degree, noted $Bel(\{A\})$ (resp. $Disb(\{A\}) = Bel(\{\neg A\})$) represents the sum of all evidence in favour of (resp. against) an assertion (A). While uncertainty degree is noted $Uncer(\{A\}) = 1 - Bel(\{A\}) - Disb(\{A\})$. The strength of an evidence for or against A is called a mass and is resp. noted $m(\{A\})$ to quantify the probability that A is True or $m(\{\neg A\})$ when A is False, while $m(\{A, \neg A\})$ quantifies ignorance.

Uncertainty Evaluation – Mathematical Background

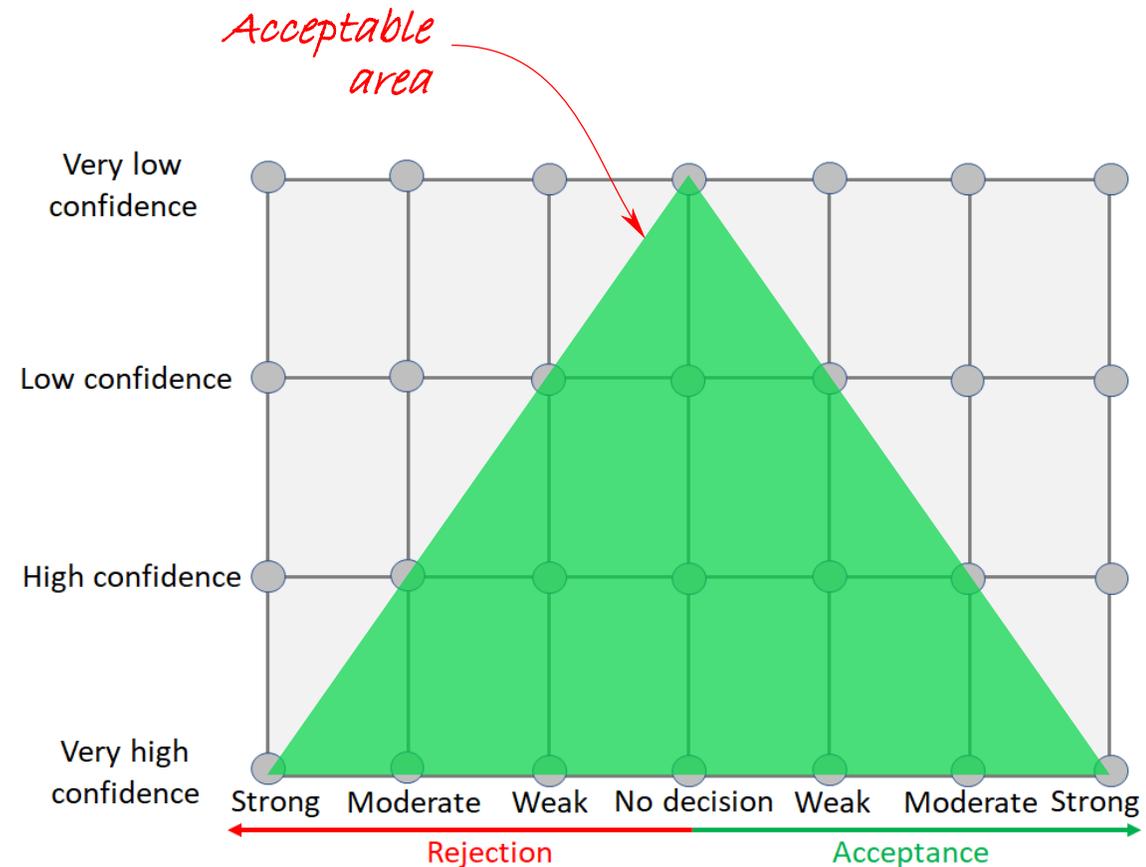
□ Elicitation

- **Decision**, $Dec(A)$: given by an expert to accept or reject a proposition (A)
 - $Dec(A)=[1+Bel(A)-Disb(A)]/2$

- **Confidence**, $Conf(A)$: the amount of information the expert needs to justify his/her decision
 - $Conf(A)=Bel(A)+Disb(A)$

□ A constraint is added to ensure that strong decisions are not taken in cases of significant uncertainty:

- $[1-Conf(A)]/2 \leq Dec(A) \leq [1+Conf(A)]/2$



Uncertainty metrics

Entry format

- ❑ Uncertainty metrics are pre-entered by the developer

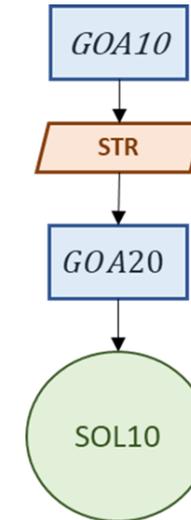
Strategy STR:

Q1. Assuming GOA20 is valid, what is your assessment of the conclusion GOA10?

Confidence	<input type="range"/>	Numerical value	Qualitative value
Decision	<input type="range"/>	0.75	Very high
Confidence	<input type="range"/>	1.00	Strong acceptance

Q2. Assuming GOA20 is invalid, what is your assessment of the conclusion GOA10?

Confidence	<input type="range"/>	Numerical value	Qualitative value
Decision	<input type="range"/>	1.00	Very high
Confidence	<input type="range"/>	1.00	Strong rejection



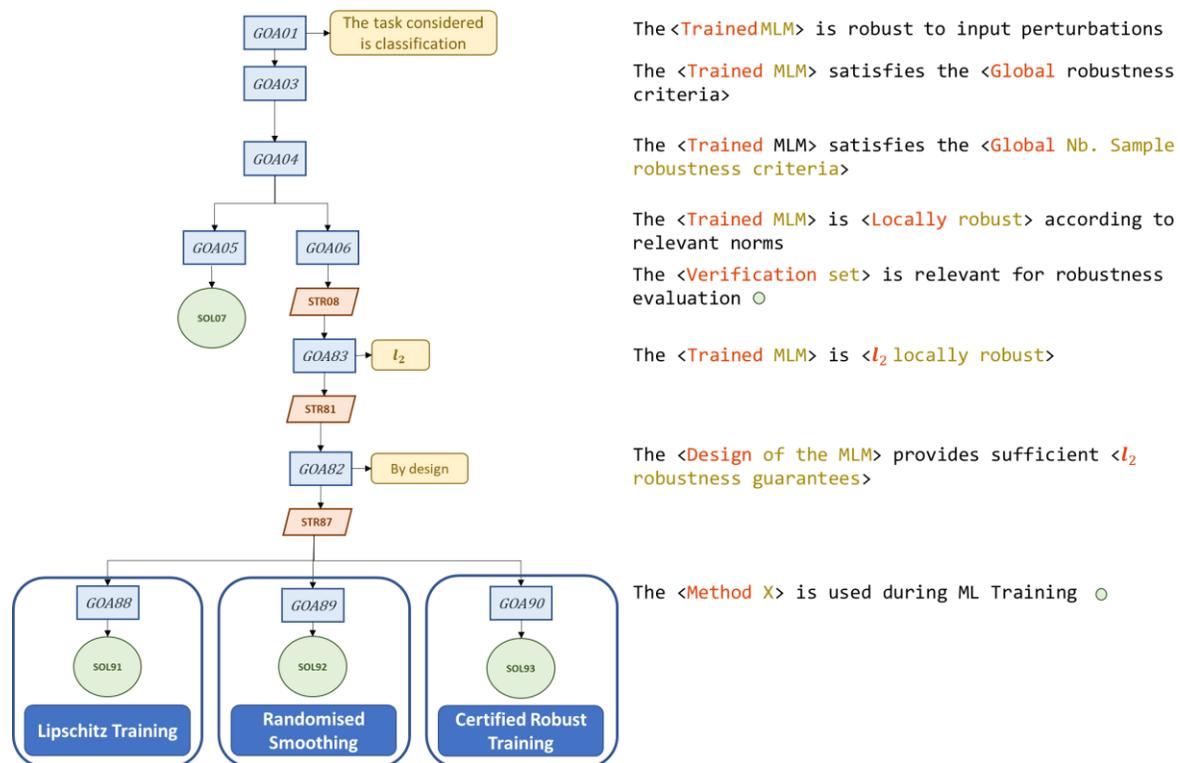
Solution (SOL10)

Considering the provided solution(s), what is your assessment of the conclusion GOA20?

Confidence	<input type="range"/>	Numerical value	Qualitative value
Decision	<input type="range"/>	1.00	Very high
Confidence	<input type="range"/>	0.60	Moderate acceptance

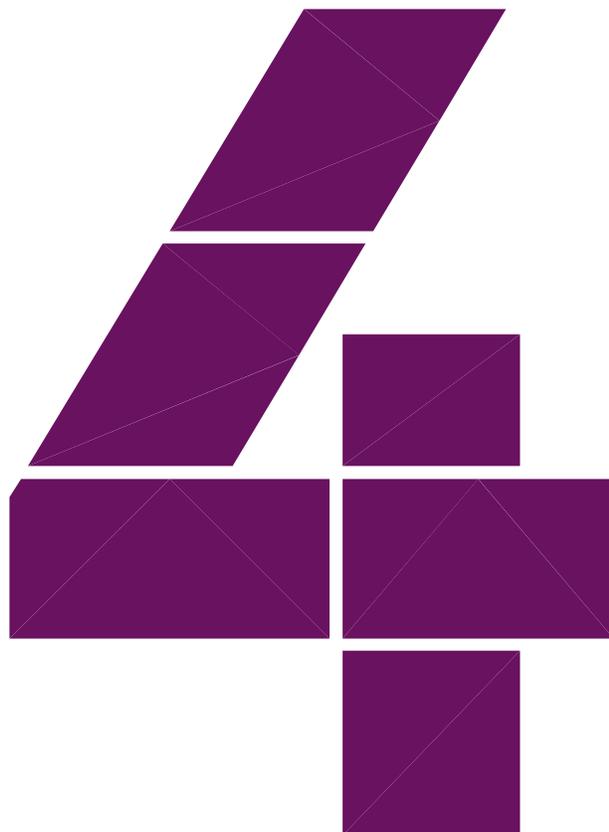
Confidence metrics propagation

Method selection on the basis of propagation results:



AC after requirements specification

Methods	Propagation results on the top-goal	
	Belief degree	Disbelief degree
Lipschitz Training	0,92	0,01
Randomised Smoothing	0,78	0,02
Certified Robust Training	0,89	0,01



Tool support

Assurance case viewpoint in Capella

Environment



Glossary entries

GSN palette

The screenshot displays the Capella software interface. On the left, the Project Explorer shows a tree structure with 'Test' expanded to 'Assurance Cases' and 'AssuranceCase 1'. Under 'Glossaries', 'Main Glossary' is expanded, and 'Global nbsample robustness criteria' is selected. A red arrow points from the handwritten text 'Glossary entries' to this entry. The main workspace shows a diagram with three nodes: a goal '[G000001]The <trained ML Model> is robust', a strategy '[A000007]The task considered is classification', and a goal '[G000002]The <Trained ML Model> satisfies the <global robustness criteria>'. A red dashed box highlights the 'GSN palette' on the right, which contains various GSN elements like 'Strategy', 'Context', 'Relations', 'Choice', 'Supported By', 'Away Elements', 'Away Goal', 'Away Solution', 'Miscellaneous', 'Renummer GSN IDs', 'Tactic', 'Glossary', 'External Elements', 'Referentiable External Elements', and 'Manage Referenced Artifacts'. A red arrow points from the handwritten text 'GSN palette' to this box. At the bottom, the Properties window shows details for the selected goal, including its name and summary.

Assurance case viewpoint in Capella Environment

Content assist

EventExchangeItem 1

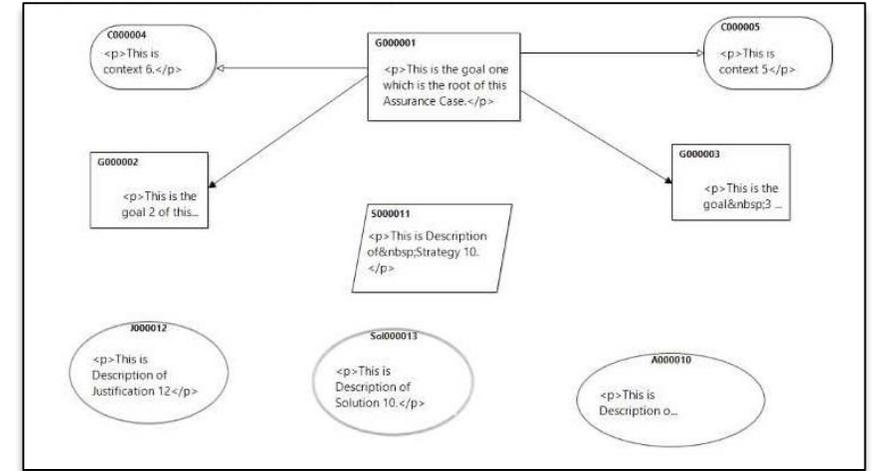
#Item|

- EventExchangeItem 1
- FlowExchangeItem 2
- ExchangeItem 3

Property \Leftrightarrow Exchange item

The diagram shows a 'DataPackageWithExchangeItems' containing three exchange items: 'EventExchangeItem', 'OperationExchangeItem', and 'FlowExchangeItem'. A property 'EngProperty 1' is associated with the 'EventExchangeItem'. A palette is open, showing 'Engineering Properties' selected, with options to 'Insert/Remove Engineering Properties'.

Assurance Cases



Assurance case viewpoint in Capella Tactics

Tabular definition of tactics

Parent Node	Choice	Selected Branches	Rationale
TacticalDecision1	[G000005]GOA18	[Y000011]GsnChoice13 - [1, *] of 2	This rationale explains why we choose to go with this selected branch. May use terms from the glossaries: [abcde (GOA14 Glossary)]
TacticalDecision2	[S000015]GsnStrat	[Y000019]GsnChoice32 - [1, *] of 2	

Selected Branches -- TacticalDecision2

Filter Available Choices
Choice Pattern (* or ?)

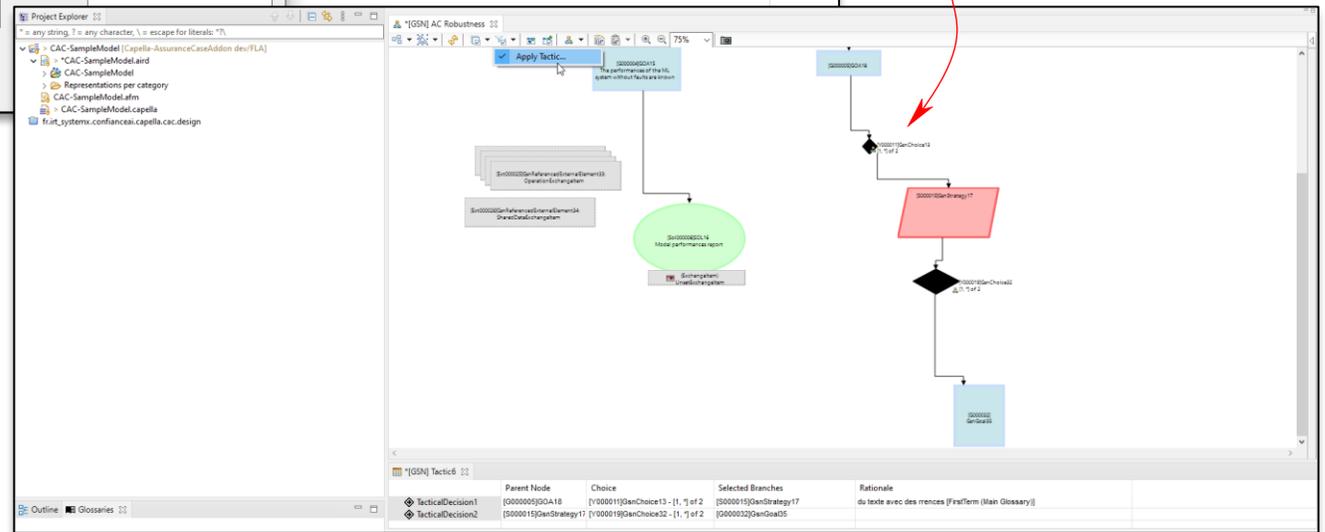
Choices

- [G000030]GsnGoal34
- [G000032]GsnGoal35

Feature

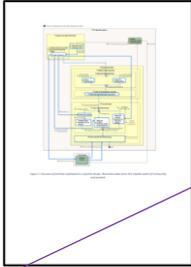
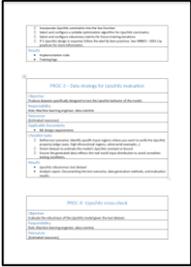
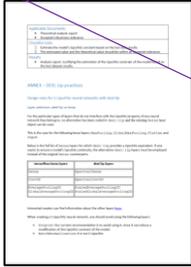
Buttons: Add, Remove, Up, Down

Filtered view of tactics



Assurance case viewpoint in Capella

V&V Plan

		<p>PROC-1B – Transform to Lipschitz design</p>
		<p>Objective Transform an existing model that adheres to Lipschitz continuity constraints.</p>
		<p>Responsibility Role: Machine learning engineer, mathematician.</p>
		<p>Resources {Estimated resources}</p>
		<p>Applicable Documents</p> <ul style="list-style-type: none"> ML model algorithm definition Deel-lip Tool: https://github.com/deel-ai/deel-lip
		<p>Checklist tasks</p> <ul style="list-style-type: none"> <input type="checkbox"/> Define the mathematical requirements for Lipschitz continuity <input type="checkbox"/> Design a model architecture using Lipschitz-compliant layers/operators <input type="checkbox"/> All k-Lipschitz layer types adhere to the design requirements, ensuring desired properties like input and output dimensions <input type="checkbox"/> All k-Lipschitz layers are compatible with the overall model architecture <input type="checkbox"/> If 1-Lipschitz design is required, follow the <i>deel-lip</i> best practices. See ANNEX – DEEL-Lip practices for more information.
		<p>Results</p> <ul style="list-style-type: none"> Architecture design document Theoretical analysis report



Conclusion

Where are we now?

What next?

Status and next steps...

□ Status

- A (small) set of **Assurance Cases** on “important” properties for the development of systems embedding ML components (robustness, explicability, ODD correctness and completeness,...)
- A Model-based approach integrating and linking workflow and assurance case models
- A Capella **GSN viewpoint** with extensions supporting the approach

□ Next steps

- Improve integration within the “Confiance.ai” workbench
 - Links with the **set of solutions** proposed by the project
 - Links with the “**Body of Knowledge**” created by the project
- Extension of Assurance Cases to **other properties**
- Addition of **new features**
 - Impact analysis
 - Dependencies between strategies



THANKS FOR YOUR
ATTENTION

