



# Product Assurance for Al-based Telemetry Health Monitoring of Safety-Critical Space Robotics

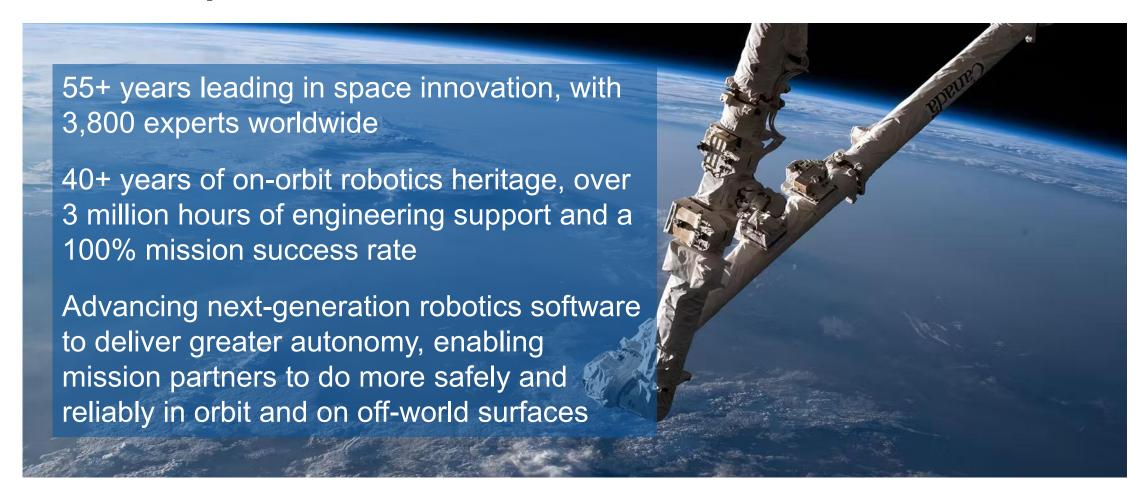
European Space Agency – Software Product Assurance Conference 2025

Kelly Gasperski, Simon Diemert, Emmanuel Lesser, Adam Casey, Nader Abu El Samid, Malav Naik, Justin Kernot, Laure Millet, **Jeff Joyce**, Paul Grouchy





## MDA Space at a Glance

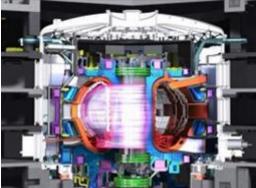


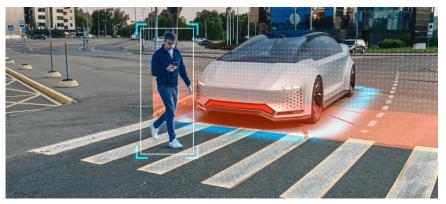




## About Critical Systems Labs



























## Presentation Abstract

This presentation reports on assurance aspects of a study that explores the application of Artificial Intelligence (AI) and Machine Learning (ML) technologies in safety-critical space robotics, with a focus on the use of AI in the ground control segment to enhance operator decision-making. A key challenge is to assure safety and achieve safety certification in a context where established standards and published guidance are not entirely compatible with the use of AI/ML. This study is partially motivated by fundamental uncertainty about how AI-based functionality can be validated to provide high confidence in its safety, i.e., high recall, low false positives. This study also seeks to identify what strategies can be used to mitigate the human factors risk that operators become complacent by trusting the Alfunctionality to monitor operational safety. While focused on a particular application for space-robotics, the results of this study will be broadly applicable beyond the space community including other technical domains such as automotive and medical devices that are rapidly integrating AI/ML into safety-critical technology. Preliminary findings and future research directions will be presented





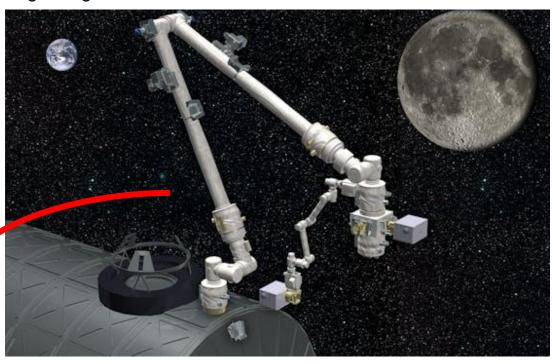
## Problem Formulation

Operators review telemetry data from the Flight Segment to check for **anomalous behaviour** suggestive of system failures.



**Ground Segment** 

#### Flight Segment



### **Key Operational Challenges**:

- 1. Rule-based monitoring ("if-then-else") is difficult to scale to complex anomalies.
- 2. Manual review by operators can miss subtle anomalies, especially in a demanding operational environment.







Will using AI impact operator capability?

Al does not get "tired"

Can we use AI to detect anomalies in telemetry data sent by the flight segment? How do we certify an AI-enabled system?

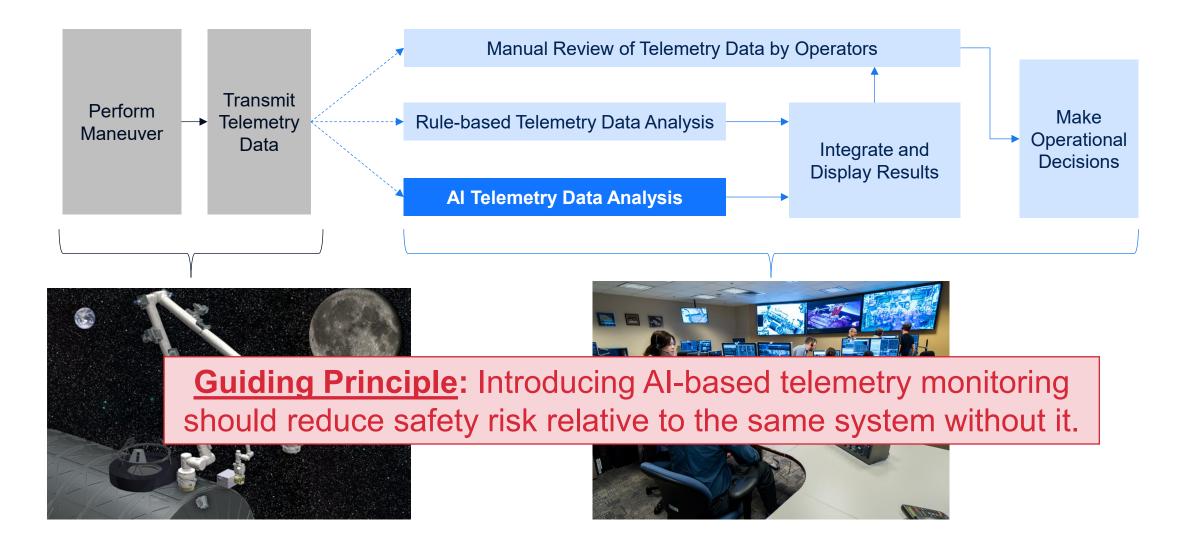
Do we trust an Al to do this task?

Al can analyze data very quickly

Al can detect complex anomalies, based on previous experience (training data)





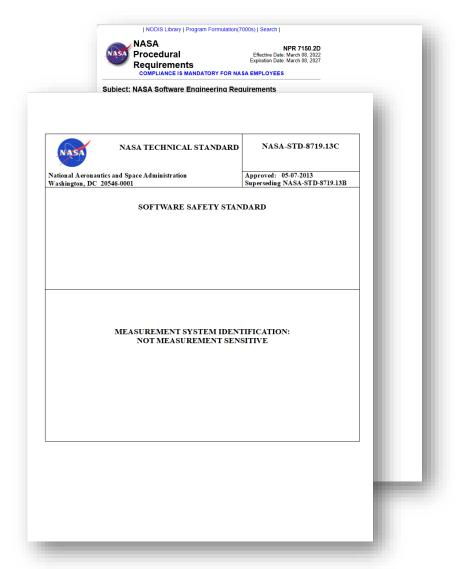




- Safety-related software should be certified:
  - NASA-STD-8719.13
  - NPR-7150.2
  - ECSS-E-ST-40C and ECSS-Q-ST-80C
- But these standards do not account for AI or machine learning technology.
- What do other industries do?



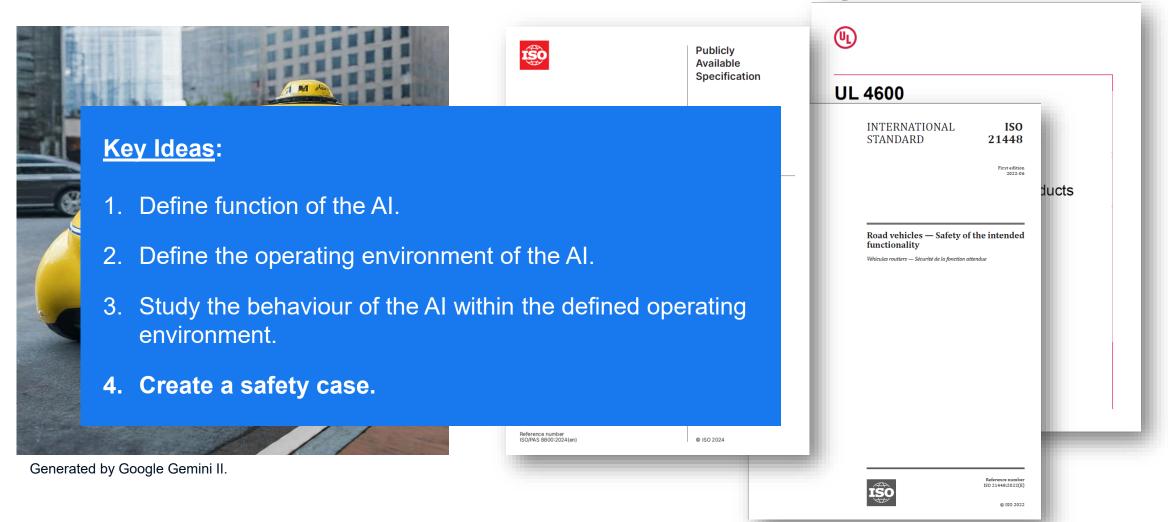








# Autonomous vehicles use AI, right?







## Three Objectives

- 1. Develop a **safety case argument** for Al-based telemetry monitoring.
- 2. Review literature on **human factors** related to Al-based monitoring systems.
- 3. Identify and apply validation methods for Al models.







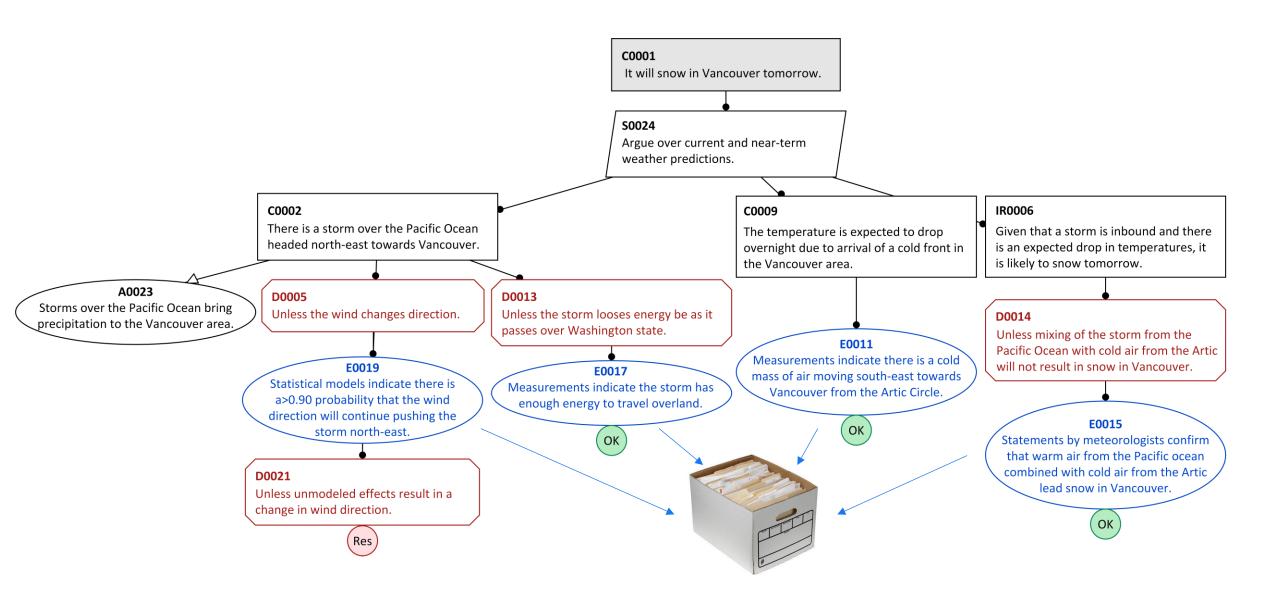




























Introducing an Al-based telemetry monitoring function to the Ground Segment will reduce safety risk relative to the same system without the Al-based telemetry monitoring function.

**Top-level Claim** 

#### S0005

Argue that the actions taken by the Albased monitor will influence the overall system in a way that reduces overall safety risk.

Strategy to argue this claim

#### C0100

2025-09-24

Bounded Actions: The only way the Albased monitor can influence the system are: 1) inform the operator that it has detected an anomaly, and 2) do nothing.

#### C1000

Actions are Safe: The net effect of the actions taken by the Al-based monitor (report anomaly or don't report anomaly) is to reduce the overall safety risk.

#### Two branches:

- **Bounded Actions**
- 2. Net Reduction in Safety Risk

(109)











#### C1000

Actions are Safe: The net effect of the actions taken by the <u>Al-based monitor</u> (report anomaly or don't report anomaly) is to reduce the overall safety risk.

Divide this claim into two cases....

#### C1001

When an anomaly is present, actions taken by the monitor will reduce risk overall.

93

Case: Anomaly Present

#### C1002

When an anomaly is not present, actions taken by the monitor will only minimally increase overall risk, if at all.



Case: Anomaly Not Present

#### IR1008

"An anomaly is present" and "an anomaly is not present" covers all possible cases. If risk is reduced overall when an anomaly is present, and no more than minimally increased when an anomaly is not present, then overall risk is reduced.













#### C1001

When an anomaly is present, actions taken by the monitor will reduce risk overall.

Divide this claim into two cases....

#### C1003

If the <u>Al-based monitor</u> correctly reports an anomaly (true positive), safety risk is maintained or decreased.



Case: True Positive

#### C1007

If the <u>Al-based monitor</u> fails to report an anomaly (false negative), safety risk is minimally increased



(10)

#### IR1016

If correctly reporting an anomaly maintains or decreases safety risk, and failing to report an anomaly only minimally increases safety risk, then safety risk overall is maintained or decreased when an anomaly is present.







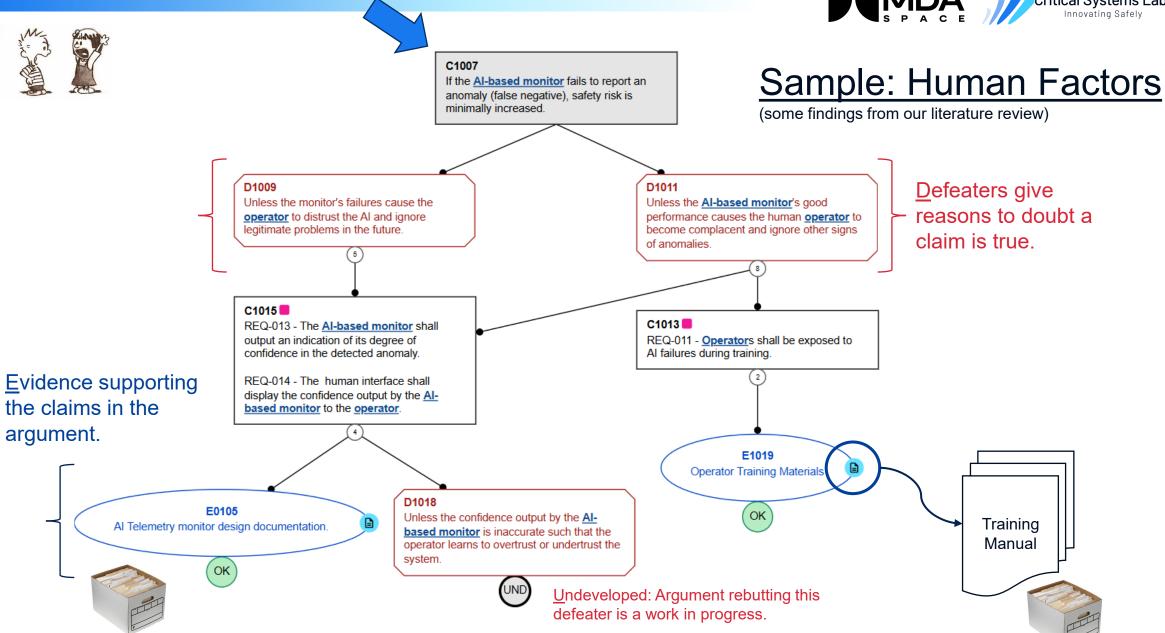




the claims in the

2025-09-24

argument.



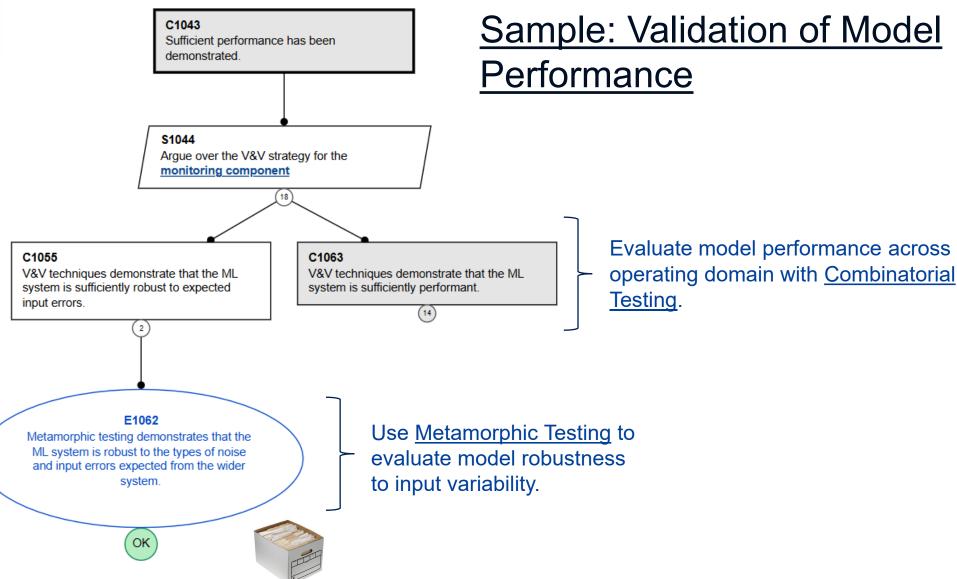
Copyright MDA Space and Critical Systems Labs















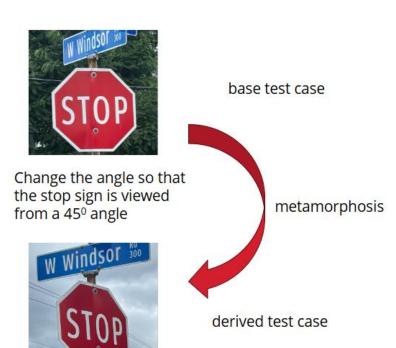
## Metamorphic Testing: Overview



Z. Q. Zhou and L. Sun, "Metamorphic testing of driverless cars," *Commun. ACM*, vol. 62, no. 3, pp. 61–67, Feb. 2019.



L. Millet, S. Diemert, R. Debouk, R. S, M. Delgado, and J. Joyce, "Specifying Functional Safety Requirements for Al/ML in terms of Metamorphic Relations," presented at the International System Safety Conference, Minneapolis, MN, United states, Aug. 2024.







# Metamorphic Testing: Example Relation

We have identified 15 MRs for our study.

MR act as a type of functional requirement for an ML model.

Research Question: Do these MRs discriminate between models of different quality?

## MR #12:

GIVEN that the model correctly detects a nominal event (no anomaly present),

APPLYING gaussian noise with standard deviation 1 is to Input A,

RESULTS IN the model detecting an anomaly.

**Base Case** 

**Transformation** 

**Expected Result** 





## Metamorphic Testing: Preliminary Results

Metamorphic Relation	Number of Failing Tests (by highest-quality models)
Given that the model correctly detects a nominal event, <b>if the time profile is shifted by +/-1s</b> , the model should still detect it as nominal.	30 / 150
Given that the model correctly detects a nominal event, <b>if the Input B profile is shifted by +/- 0.1 units</b> , the model should still detect this change as nominal.	30 / 150
Given that the model correctly detects a nominal event, if the gaussian noise with standard deviation 0.25 is added to Input A, the model should still detect no anomaly.	30 / 150
Given that the model correctly detects a nominal event, <b>if gaussian noise with standard deviation 0.25 is added to Input B</b> , the model should still detect no anomaly.	30 / 150
Given that no anomaly is detected and there are fewer than 2.5% outliers, <b>if 5% random outliers are added to Inputs A and B</b> , the model should detect an anomaly.	30 / 150

**Preliminary Observation:** MRs can reveal inconsistencies between an engineer's expectations of a model and the model's actual behaviour.





## Validation Method: Combinatorial Testing

- Approach developed by researchers U.S. NIST
- Generate a set of test cases that contain a t-way tuples of the input set.
  - For conventional software,  $t \le 6$  is sufficient to find \*most\* (~99%) of software defects [1].
- Research Question: Does CT discriminate between ML models of varying quality.
  - Experiment in progress!

Given  $A, B, C \in \{0,1\}$  for t = 2:

Α	В	С
0	0	1
0	1	0
1	1	1
1	0	0

[1] R. Kuhn, R. Kacker, and Y. Lei, "Practical Combinatorial Testing," National Institute of Standards and Technology, NIST Special Publication (SP) 800-142, Oct. 2010.





## Combinatorial Testing: Preliminary Results

Preliminary results averaged across high-quality models, on a subset of the identified operating domain.

t value	# Test Cases	# Passed Test Cases	Sensitivity	Specificity
t=1	21	21 (100%)	100%	0%
t=2	441	421.5 (95.6%)	95.7%	0%
t=3	1384	1334.9 (96.5%)	96.5%	50.0%
t=4	3970	3834.4 (96.6%)	96.7%	50.0%

**Preliminary Observation:** Higher combinatorial strengths begin to reveal model insufficiencies. On-going experiments are exploring this question more deeply.





## Summary and Takeaways

The results of this project will enable MDA Space to:

- Use safety case arguments as a key tool for certifying Al-enabled safety-critical systems, driving more reliable and consistent operations
- Incorporate human factors as part of the safety argument to better integrate AI with operators, handle greater data volumes, and enhance preventative maintenance to reduce downtime and costs
- 3. Apply the safety argument to **organize and motivate** verification methods (e.g., metamorphic and combinatorial testing), creating a **scalable framework** to extend these benefits across future missions and platforms







Emmanuel Lesser
Software S&MA Lead
emmanuel.lesser@mda.space

Jeff Joyce CEO / President jeff.joyce@cslabs.com