**Carnegie Mellon University**

Software Engineering Institute

# Explainable Verification for Rapid Certification

**NOVEMBER 13, 2024**

Bjorn Andersson

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

2

# Bottom Line Up Front

Background

- Certification is expensive; testing is expensive.

- Formal methods (FM) potential make it less expensive.

- Practitioners do not trust FM.

**Theorem 4.** *If* $\forall \tau_i \in \tau$ $\forall q \in \text{QSET}_i$ $\exists t \in [0, D_i]$ $\text{reqlpARINC653}(i, t) \leq \text{sbf}(i, q, t)$, *then the system is schedulable.*

Our work

- Create Explanation Methods (XMs) for FM.

- 88% believe XMs increase confidence in FMs.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

3

# DoD Context



Photo Credit: U.S. Army photo by Jay Miller



U.S. Airforce Illustration:
Staff Sgt. Shelby Thurman

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

4

# DoD Context



Photo Credit: U.S. Army photo by Jay Miller



U.S. Airforce Illustration:
Staff Sgt. Shelby Thurman

System includes software that reads sensors and controls physical objects to conduct its mission.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

5

# DoD Context



Photo Credit: U.S. Army photo by Jay Miller



U.S. Airforce Illustration:
Staff Sgt. Shelby Thurman

Critical Systems: Need to be certified/approved before being deployed.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

6

# DoD Context



Photo Credit: U.S. Army photo by Jay Miller



U.S. Airforce Illustration:
Staff Sgt. Shelby Thurman

Correct timing of software in execution is essential.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

7

# DoD Context

Photo Credit: U.S. Army photo by Jay Miller

From page 24 in
*Army Military Airworthiness Certification Criteria (AMACC), Revision A Change 2 (C2)*, 2021:

**3.1.36 Performance.** A quantitative or qualitative measure characterizing a physical or functional attribute relating to the execution of an operation or function. Performance attributes include quantity (how many or how much), quality (how well), coverage (how much area, how far), timeliness (how responsive, how frequent), and readiness (availability, mission/operational readiness). Performance is an attribute for all systems, people, products, and processes including those for development, production, verification, deployment, operations, support, training, and disposal. Supportability parameters, manufacturing process variability, reliability, and so forth are all performance measures.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

8

# DoD Context



Photo Credit: U.S. Army photo by Jay Miller

From page 33 in
*Army Military Airworthiness Certification Criteria (AMACC), Revision A Change 2 (C2), 2021:*

| WCET | Worst-Case Execution Time |
|------|---------------------------|

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

9

# DoD Context



U.S. Airforce Illustration:
Staff Sgt. Shelby Thurman

From page 46 in
*DAFMAN91-119_DAFGM2023-01*,
June 8 2023:

**8.10.5** The software shall detect and appropriately respond to Critical Signal-impacting hardware and software faults during all stages of execution; including startup, operation, and shutdown; within the lesser of (1) the maximum operationally acceptable time and prior to the time limit; and (2) the time to any associated irreversible adverse system event. **(T-1)** Verification activities should prove that the software detects hardware and software faults during all stages of execution within an acceptable time and responds to those within an acceptable time.

**8.10.6** The software shall report the events defined in the following subparagraphs to the operator within the lesser of (1) the maximum operationally acceptable time and prior to the time limit, and (2) the time to any associated irreversible adverse system event. **(T-1)** Verification activities should prove that the software reports the listed events and automated actions to the operator within an acceptable time.

# DoD Context



U.S. Airforce Illustration:
Staff Sgt. Shelby Thurman

From page 52 in
*DAFMAN91-119_DAFGM2023-01*,
June 8, 2023:

**9.4 Real-Time Embedded Systems.** An embedded system is a processor with software that usually serves a dedicated function within a larger system. Embedded systems rarely have traditional user Interfaces typical of computing systems and are not as multipurpose as other computing systems. Embedded systems have specific software safety concerns. Real-time software has strict timing constraints that require the system to react to inputs within a certain amount of time. Real-time software uses operating systems and software paradigms to maintain the timing constraints.

9.4.1. Embedded systems in the nuclear mission usually contain real-time software. Per **paragraph 8.16**, Real-Time Processing is only nuclear critical if a failure of the software can lead to the failure to cancel a Critical Signal after an operator has withdrawn human intent or if a failure of the software can lead to a failure to issue a Safing Command. If the design precludes these scenarios, then Real-Time Processing requirements are not applicable.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.
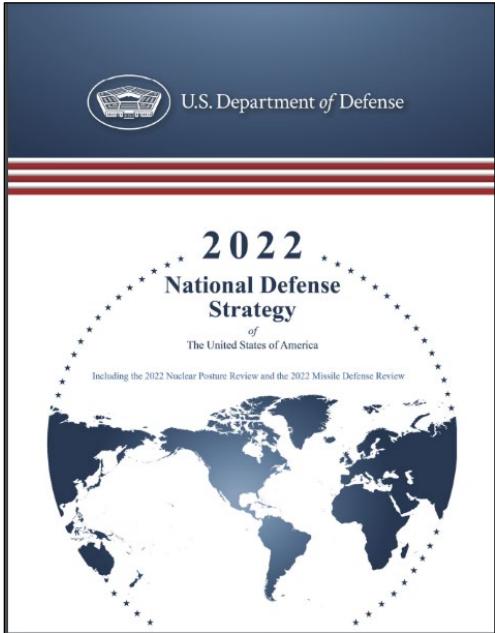
11

# DoD Context



Photo Credit: U.S. Army photo by Jay Miller



U.S. Airforce Illustration:
Staff Sgt. Shelby Thurman

Critical Systems: Need to be certified/approved before being deployed. Time consuming and expensive.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

12

# The Need for Speed (of Software Development of Weapon Systems)



"The Department will instead reward rapid experimentation, acquisition, and fielding…."

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

13

How do we combine the need for safety with the need for speed?

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

14

How do we combine the need for safety with the need for speed?



Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

15

MORE SAFETY

FASTER FIELDING

Developers

How do we combine the need for safety with the need for speed?



Even if we decrease the development time to zero, then certification time is still there and becomes a bottleneck for our ability to field new systems rapidly.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

16

**MORE SAFETY** ← → **FASTER FIELDING**

Developers

How do we combine the need for safety with the need for speed?

Test data as cert evidence
- Exhaustive takes too long
- Non-exhaustive is unsafe

→

Formal methods have the potential to avoid these.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

17

# Formal Methods



**Question:** Given model M, is correctness property $\varphi$ true for all executions?

FM tool

Yes/No/ Undecided

Input: (i) a model of a system and (ii) a correctness condition.

Output: True/False/Undecided

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

18

# Formal Methods for Real-Time Systems



**Question:** Given taskset $\tau$, the following run-time scheduler S, and assumptions A, is it schedulable?

FM tool

Yes/No/ Undecided

## Real-Time Systems

- Software systems where computing the right result is not enough. It needs to be computed at the right time.

- It is often implemented as a set of concurrent tasks.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

19

# The Role of Explanation in Certification



**Question:** Given taskset $\tau$, the following run-time scheduler S, and assumptions A, is it schedulable?

FM tool → Yes/No/Undecided

Developer → I have run a Formal Methods (FM) tool and it outputs SAFE → Certifier

Why should I trust this tool?

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

20

# The Role of Explanation in Certification



Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

21

**Our New Tool**

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

22

Carnegie
Mellon
University
Software
Engineering
Institute

**Satisfying Real-Time Requirements of Multicore Software on ARINC 653**

This program implements a schedulability test for tasks with co-runner dependent execution times. The model is the one in B. Andersson et al., "Satisfying Real-..." but the schedulability test is different.

| | | | |
|---|---|---|---|
| Number of tasks | 5 | Number of processors | 2 |
| Number of partitions | 3 | Number of partition windows | 4 |
| Major time frame | 0.002 | Interpretation | 0 |

**Tasks**

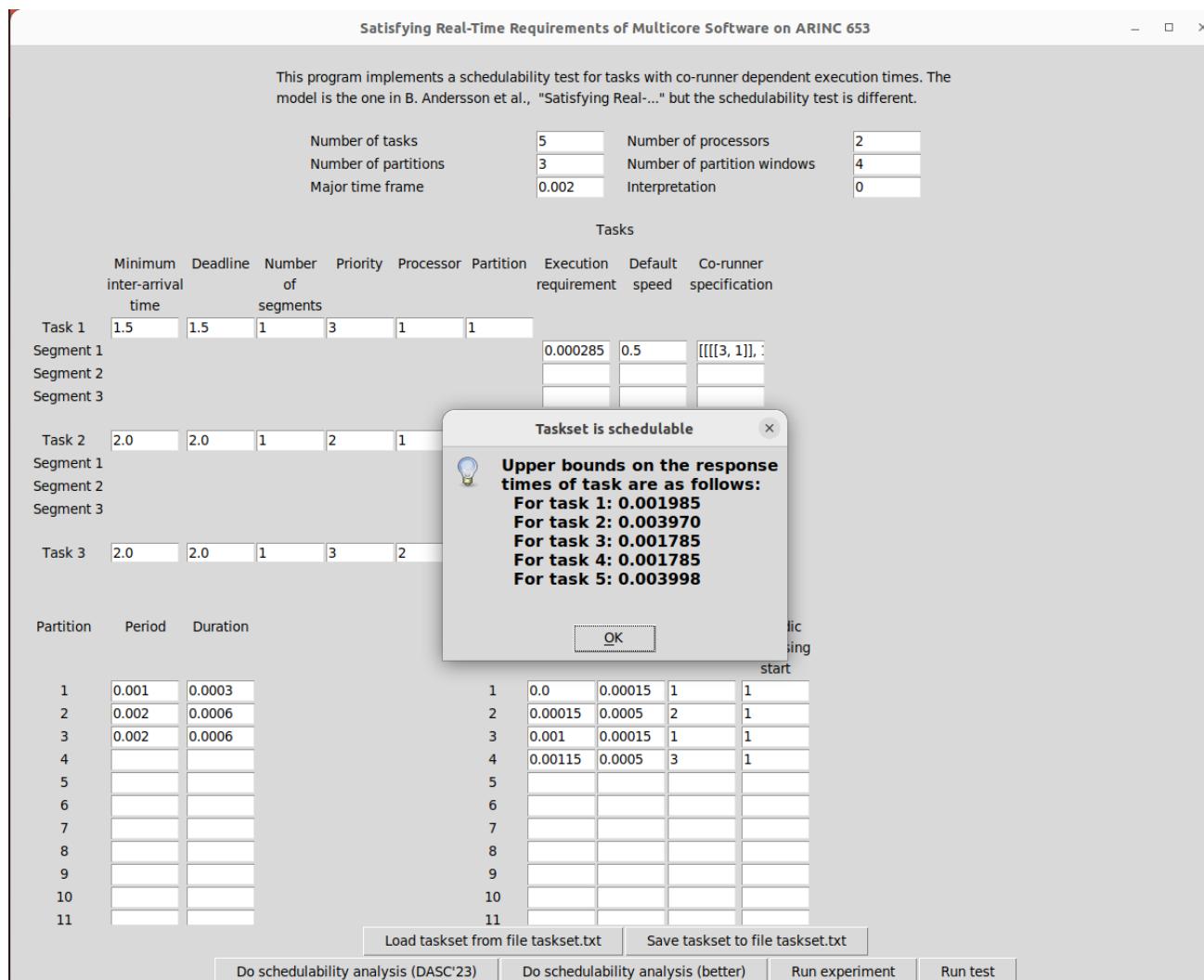| | Minimum inter-arrival time | Deadline | Number of segments | Priority | Processor | Partition | Execution requirement | Default speed | Co-runner specification |
|---|---|---|---|---|---|---|---|---|---|
| Task 1 | 1.5 | 1.5 | 1 | 3 | 1 | 1 | | | |
| Segment 1 | | | | | | | 0.000285 | 0.5 | [[[[3, 1]], |
| Segment 2 | | | | | | | | | |
| Segment 3 | | | | | | | | | |
| Task 2 | 2.0 | 2.0 | 1 | 2 | 1 | 1 | | | |
| Segment 1 | | | | | | | 0.000285 | 0.5 | [[[[3, 1]], ( |
| Segment 2 | | | | | | | | | |
| Segment 3 | | | | | | | | | |
| Task 3 | 2.0 | 2.0 | 1 | 3 | 2 | 2 | | | |

**Partitions and partition windows**

| Partition | Period | Duration | | Partition window | Offset | Duration | Belongs to partition | Periodic processing start |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.001 | 0.0003 | | 1 | 0.0 | 0.00015 | 1 | 1 |
| 2 | 0.002 | 0.0006 | | 2 | 0.00015 | 0.0005 | 2 | 1 |
| 3 | 0.002 | 0.0006 | | 3 | 0.001 | 0.00015 | 1 | 1 |
| 4 | | | | 4 | 0.00115 | 0.0005 | 3 | 1 |
| 5 | | | | 5 | | | | |
| 6 | | | | 6 | | | | |
| 7 | | | | 7 | | | | |
| 8 | | | | 8 | | | | |
| 9 | | | | 9 | | | | |
| 10 | | | | 10 | | | | |
| 11 | | | | 11 | | | | |

| Load taskset from file taskset.txt | Save taskset to file taskset.txt |
|---|---|

| Do schedulability analysis (DASC'23) | Do schedulability analysis (better) | Run experiment | Run test |
|---|---|---|---|

Enter data to describe system.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

23

Click button.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

24

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

25

Carnegie
Mellon
University
Software
Engineering
Institute

**Satisfying Real-Time Requirements of Multicore Software on ARINC 653**

This program implements a schedulability test for tasks with co-runner dependent execution times. The model is the one in B. Andersson et al., "Satisfying Real-..." but the schedulability test is different.

| | | |
|---|---|---|
| Number of tasks | 5 | |
| Number of partitions | 3 | |
| Major time frame | 0.002 | |
| Number of processors | 2 | |
| Number of partition windows | 4 | |
| Interpretation | 0 | |

**Tasks**

| | Minimum inter-arrival time | Deadline | Number of segments | Priority | Processor | Partition | Execution requirement | Default speed | Co-runner specification |
|---|---|---|---|---|---|---|---|---|---|
| Task 1 | 1.5 | 1.5 | 1 | 3 | 1 | 1 | | | |
| Segment 1 | | | | | | | 0.000285 | 0.5 | [[[3, 1]], |
| Segment 2 | | | | | | | | | |
| Segment 3 | | | | | | | | | |
| Task 2 | 2.0 | 2.0 | 1 | 2 | 1 | 1 | | | |
| Segment 1 | | | | | | | 0.000285 | 0.5 | [[[3, 1]], ( |
| Segment 2 | | | | | | | | | |
| Segment 3 | | | | | | | | | |
| Task 3 | 2.0 | 2.0 | 1 | 3 | 2 | | | | |

**Do you want to have an explanation?** ×

Yes   No

| Partition | Period | Duration | | Partition window | Offset | Duration | Belongs to partition | Periodic processing start |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.001 | 0.0003 | | 1 | 0.0 | 0.00015 | 1 | 1 |
| 2 | 0.002 | 0.0006 | | 2 | 0.00015 | 0.0005 | 2 | 1 |
| 3 | 0.002 | 0.0006 | | 3 | 0.001 | 0.00015 | 1 | 1 |
| 4 | | | | 4 | 0.00115 | 0.0005 | 3 | 1 |
| 5 | | | | 5 | | | | |
| 6 | | | | 6 | | | | |
| 7 | | | | 7 | | | | |
| 8 | | | | 8 | | | | |
| 9 | | | | 9 | | | | |
| 10 | | | | 10 | | | | |
| 11 | | | | 11 | | | | |

Load taskset from file taskset.txt   Save taskset to file taskset.txt

Do schedulability analysis (DASC'23)   Do schedulability analysis (better)   Run experiment   Run test

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

26

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

27

**Satisfying Real-Time Requirements of Multicore Software on ARINC 653**

This program implements a schedulability test for tasks with co-runner dependent execution times. The model is the one in B. Andersson et al., "Satisfying Real-..." but the schedulability test is different.

| | | | |
|---|---|---|---|
| Number of tasks | 5 | Number of processors | 2 |
| Number of partitions | 3 | Number of partition windows | 4 |
| Major time frame | 0.002 | Interpretation | 0 |

**Tasks**

| | Minimum inter-arrival time | Deadline | Number of segments | Priority | Processor | Partition | Execution requirement | Default speed | Co-runner specification |
|---|---|---|---|---|---|---|---|---|---|
| Task 1 | 1.5 | 1.5 | 1 | 3 | | 1 | | | |
| Segment 1 | | | | | | | | | |
| Segment 2 | | | | | | | | | |
| Segment 3 | | | | | | | | | |
| Task 2 | 2.0 | 2.0 | 1 | 2 | | 1 | | | |
| Segment 1 | | | | | | | | | |
| Segment 2 | | | | | | | | | |
| Segment 3 | | | | | | | | | |
| Task 3 | 2.0 | 2.0 | 1 | 3 | | 2 | | | |

| Partition | Period | Duration |
|---|---|---|
| 1 | 0.001 | 0.0003 |
| 2 | 0.002 | 0.0006 |
| 3 | 0.002 | 0.0006 |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |
| 8 | | |
| 9 | | |
| 10 | | |
| 11 | | |

**Explanation window**

Click on numbers below to get explanation

| Task | Meets deadlines | Upper bound on response time through simulation | Response time obtained | Explain adverse phasing | Explain not dominated ptw | Explain iteration |
|---|---|---|---|---|---|---|
| 1 | Yes | 0.001985 | 0.001985 | Do it! | Do it! | Do it! |
| 2 | Yes | 0.003970 | 0.003970 | Do it! | Do it! | Do it! |
| 3 | Yes | 0.001785 | 0.001785 | Do it! | Do it! | Do it! |
| 4 | Yes | 0.001785 | 0.001785 | Do it! | Do it! | Do it! |
| 5 | Yes | 0.003998 | 0.003998 | Do it! | Do it! | Do it! |

Load taskset from file taskset.txt    Save taskset to file taskset.txt

Do schedulability analysis (DASC'23)    Do schedulability analysis (better)    Run experiment    Run test

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

30

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

# Evaluation of Our Tool

Three different groups of people

Certification authorities and experienced software developers of safety-critical systems

These groups are representative users of our tool, and they include: Galois, AvMC, AFSEC.



Number of persons: 35

Number of persons who found our explanations helpful to gain confidence in output of FM tool: 31

$31/35 \approx 88.5\%$

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

32

# Without Explanation

# With Explanation



Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

33

# The Big Picture and this Project



Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

34

# The Big Picture and this Project

Explanation of properties

Explanation of timing correctness

Explanation of program correctness

Explanation of link between software correctness to properties of physical world

Explanation of WCET analysis

Explanation of schedulability analysis

**This project**

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

35

# The Big Picture and this Project

Explanation of Yes or No

Explanation of Yes

Explanation of No

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

36

# The Big Picture and this Project

Explanation of Yes or No

Explanation of Yes

Explanation of No

This project

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

37

# The Big Picture and this Project

**Question:** Given model M, is correctness property $\varphi$ true for all executions?

FM tool

Yes/No/ Undecided

I am feeding a model M into this tool. What does it (M) actually mean? Does it describe the system that I want?

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

38

# The Big Picture and this Project



**Question:** Given model M, is correctness property $\varphi$ true for all executions?

FM tool

Yes/No/ Undecided

I need input explanation.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

39

# The Big Picture and this Project

**Question:** Given model M, is correctness property φ true for all executions?

→ FM tool → Yes/No/ Undecided →

This tool gives me an answer to the question whether φ holds. What does it actually mean? Does it describe the correctness property that I care about?

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

40

# The Big Picture and this Project



**Question:** Given model M, is correctness property $\varphi$ true for all executions?

FM tool

Yes/No/ Undecided

I need output explanation.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

41

# The Big Picture and this Project

**Question:** Given model M, is correctness property $\varphi$ true for all executions?

FM tool

Yes/No/Undecided

Why did this tool produce this output for this input?

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

42

# The Big Picture and this Project



**Question:** Given model M, is correctness property φ true for all executions?

FM tool

Yes/No/ Undecided

I need input-output relation explanation.

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

43

# The Big Picture and this Project

Explanation

Input explanation

Input-output
explanation

Output explanation

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

44

# The Big Picture and this Project

Explanation

Input explanation

Input-output
explanation

Output explanation

This project

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

45

# Conclusion

- Cert expensive; testing expensive
- Formal methods (FM) potential make it less expensive
- Practitioners do not trust FM

**Theorem 4.** *If* $\forall \tau_i \in \tau \; \forall q \in \mathrm{QSET}_i \; \exists t \in [0, D_i]$ $\mathrm{reqlpARINC653}(i, t) \leq \mathrm{sbf}(i, q, t)$, *then the system is schedulable.*

Our work
- Create Explanation Methods (XMs) for FM
- 88% believe XMs increase confidence in FMs
- Focus on timing
- Broad applicability for critical systems

# Team



**Bjorn Andersson**

**Dionisio de Niz**

**Mark Klein**

Explainable Verification for Rapid Certification
©2024 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

47

Carnegie
Mellon
University
Software
Engineering
Institute

# Thanks!