# IDA

# Formal and Empirical Verification of Autonomy for Assurance (Conference Briefing)

David M. Tate

**IDA**

The Institute for Defense Analyses is a nonprofit corporation that operates three Federally Funded Research and Development Centers. Its mission is to answer the most challenging U.S. security and science policy questions with objective analysis, leveraging extraordinary scientific, technical, and analytic expertise.

Rigorous Analysis │ Trusted Expertise │ Service to the Nation

# Executive Summary

Assurance cases are a standard technique for establishing the trustworthiness of systems. An assurance case is a structured argument establishing that the system is sufficiently dependable to justify fielding it in a specific, defined operational context. Assurance cases require both evidence and arguments. For verification tools to be useful in helping to field autonomous systems, the tools must produce evidence that supports relevant assurance cases.

Many autonomous systems will be too complex to permit comprehensive, formal verification of all their dependability properties. This is especially true of systems that rely on machine learning (ML) to provide important perception and/or decision-making capabilities, or whose operations involve nontrivial interactions with human beings. For such systems, assurance cases will need to incorporate empirical evidence in addition to formal results to support valid arguments that provide confidence in a system's safety, dependability, and performance.

For a given autonomous system, the structure of the appropriate assurance case arguments is driven by the concept of operations and architecture of the system. This structure can be represented formally. Currently, the most common formal model of argument structure is a graph representation called Goal Structuring Notation (GSN). The evidential nodes of the GSN representation of the assurance argument provide a useful interface between the assurance case and test planning, and between the assurance case and formal verification efforts. The outputs of empirical testing and formal verification activities become the inputs to specific evidential nodes of the assurance case model. How the evidence is used in the assurance argument informs what the outputs of the formal and empirical verification activities need to be. This, in turn, informs the selection of specific formal methods and experimental designs. The specific evidence to be generated also has implications for test instrumentation, physical infrastructure, and personnel needed to execute the test design.

# Formal and Empirical Verification of Autonomy for Assurance

David Tate

Institute for Defense Analyses

February 2021

# The goal is <u>assured</u> effectiveness and dependability

Autonomous capabilities don't help if we're not sufficiently confident to field and employ the systems.

There will always be **some** kind of certification, licensure, or acceptance testing process.

There may be multiple certification authorities (e.g., Safety, Cybersecurity, Effectiveness, Reliability).

# State of the Art:  Assurance Cases

An **assurance case** is a structured argument that the system is sufficiently dependable to permit fielding in a defined operational context.

Existing standards and regulatory bodies **already require** explicit assurance cases for complex systems:
- Safety cases (oldest, most mature literature)
- Software assurance cases (cybersecurity, reliability)
- Robustness cases

Currently, these cases are generally stovepiped.

# Example:  ISO/IEC 15026-2 (2011)

**Systems and Software Engineering—**

**Systems and Software Assurance—**

**Part 2: Assurance Case**

**1 Scope**

This part of ISO/IEC 15026 specifies minimum requirements for the structure and contents of an assurance case. An assurance case includes a top-level claim (or set of claims) for a property of a system or product, systematic argumentation regarding this claim, and the evidence and explicit assumptions that underlie this argumentation. Arguing through multiple levels of subordinate claims, this structured argumentation connects the top-level claim to the evidence and assumptions.

# Assurance cases require both *evidence* and *arguments*

A pile of evidence is not an argument.

An argument without evidence is unconvincing.

The wrong evidence doesn't help.

**The outputs of verification should be evidence that supports needed assurance cases.**

# Where does the evidence come from?

Traditional assurance cases are based on a combination of ***empirical*** and ***formal*** evidence:

~~Exhaustive testing~~

Formal verification

Design of experiments

Run-time monitors

Human in the loop + training

# Work backward to identify infrastructure needs

**Assurance case**

↓

**Supporting evidence**

↓

**Measurements + formal results**

↓

**Instrumentation & models**

↓

**Infrastructure**

1. Given the system assurance case, what evidence will be required?

2. What time series of measurements would produce (some of) that evidence?

3. What instrumentation is required to collect those measurements?

4. What infrastructure is needed to support that instrumentation?

- Simulation models
- Formal models
- Software tools
- Human capital

# Iterate forward to develop assurance cases

Mission requirements

↓

HMT CONOPS

↓

Autonomy requirements

↓

Test oracles

↓

Instrumentation needs

↓

Evidence (time series)

↓

Assurance cases

↓

Certification

1. Make an **initial guess** at how the autonomous system will team with humans.

2. Codify **test oracles** for acceptable behaviors, including internal behaviors.

3. Construct **assurance case outlines** – what arguments will convince? What evidence will they require?

4. Formally **verify** compliance where possible.

5. Derive **evidential requirements** – what measurements will be needed to assess performance against the oracles and provide the empirical evidence?

6. Collect evidence and **iterate**.

# What new (and existing) tools might be useful?

Formal methods

Instrumenting cognition/explainable AI

Intelligent adversarial testing

Assurance case development tools

# Examples:  Formal Methods

*VERIFAI: A Toolkit for the Design and Analysis of Artificial Intelligence-Based Systems*

Tommaso Dreossi, Daniel J. Fremont, Shromona Ghosh, Edward Kim, Hadi Ravanbakhsh, Marcell Vazquez-Chanlatte, Sanjit A. Seshia

*We present VERIFAI, a software toolkit for the formal design and analysis of systems that include artificial intelligence (AI) and machine learning (ML) components. VERIFAI particularly seeks to address challenges with applying formal methods to perception and ML components, including those based on neural networks, and to model and analyze system behavior in the presence of environment uncertainty.*

*Using Formal Verification to Evaluate Human-Automation Interaction: A Review*

Bolton, Bass, and Siminiceanu

*IEEE Transactions on Systems, Man, and Cybernetics: Systems* 43 #3, May 2013

# Example: Empirical Evidence Using XAI



Salient pixel analysis of the NVIDIA PilotNet self-steering system shows that the system all but ignores the road surface itself, focusing instead on features that indicate not-road. This system does not maintain an internal representation of the terrain; the neural net generates steering commands based on the real-time camera inputs.

Image from Bojarski et al., *Explaining How a Deep Neural Network Trained with End-to-End Learning Steers a Car.* arXiv:1704.07911v1 [cs.CV] April 25, 2017

# Examples:  Assurance Case Development Tools

Assurance Case Editor  (**ACedit)**

Assurance Case Automation Toolset (**AdvoCATE**)

Evidence Confidence Assessor (**EviCA**)

**Astah GSN** (commercial product, see [astah.net](astah.net))
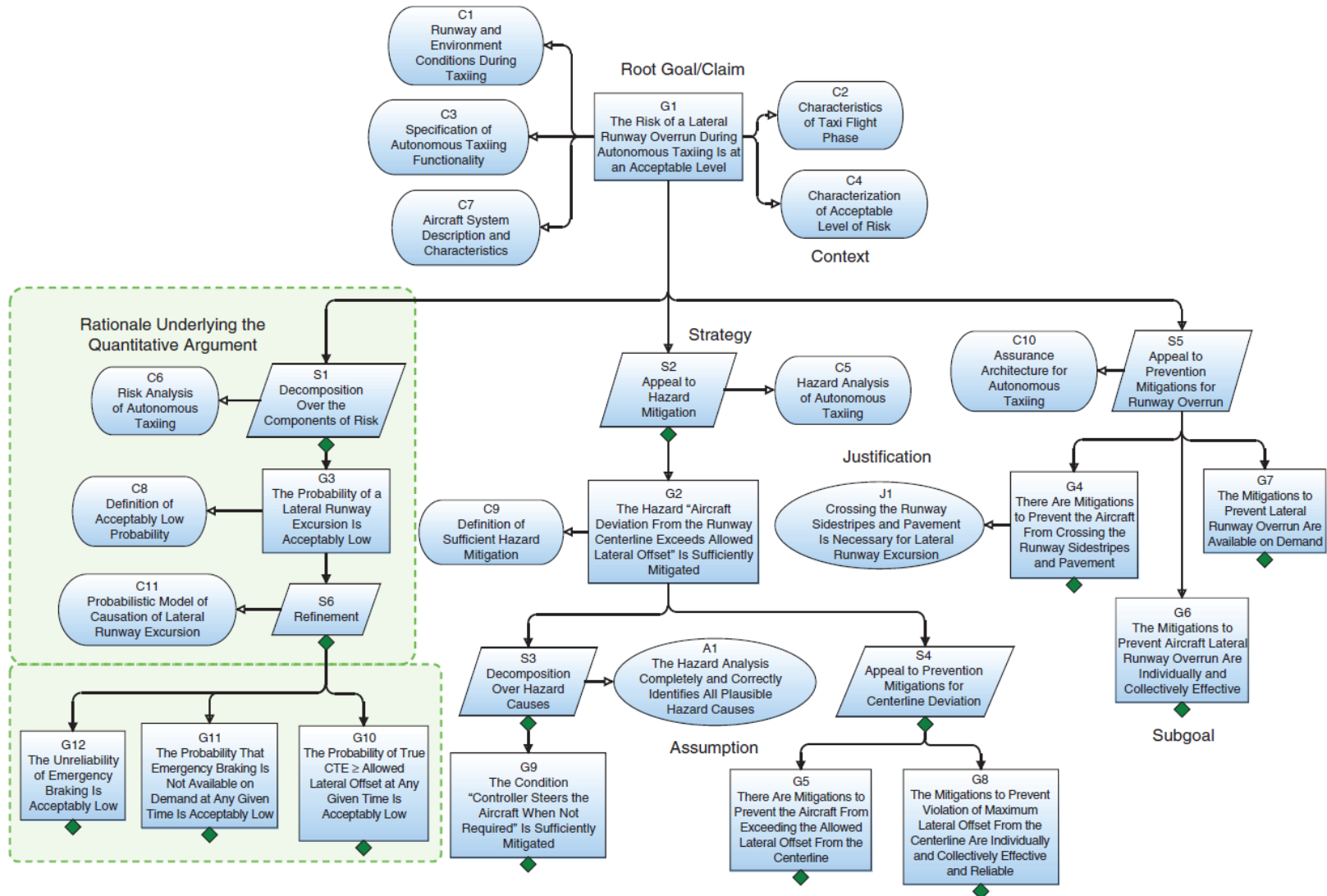
Each tool uses **Goal Structuring Notation** (GSN) as the graphical language for describing and manipulating arguments.

Reference: *Tool Support for Assurance Case Development*, Ewen Denney and Ganesh Pai, NASA Ames Research Laboratory

# Example: A Partial Safety Case in GSN

From "Dynamic Assurance Cases: A Pathway to Trusted Autonomy," Asaadi, Denny, Menzies, Pai, Petroff

# Bottom Line at the Bottom

Assurance cases for autonomous systems will strain human cognitive abilities.

Evidence to support the arguments will require a mix of formal and empirical verification techniques.

Tools exist to support the automated development and management of assurance cases and the incorporation of both formal and empirical evidence.

# Backup

# Case Study:  Assurance Case Development

*A Case Study in Assurance Case Development*

*for Scientific Software*

Mojdeh Sayari Nejad

MS (Computer Science) Thesis

McMaster University, 2017


Assurance case for 3dfim+ software for analyzing

functional MRI images of the brain


https://macsphere.mcmaster.ca/handle/11375/23075

# "Levels of autonomy" is a red herring

"It takes more sophisticated technology to keep the humans in the loop than it does to automate them out… On a commonly used scale of levels of autonomy, level 1 is fully manual control and level 10 is full autonomy… history and experience show that the most difficult, challenging, and worthwhile problem is not full autonomy but the perfect five—a mix of human and machine and the optimal amount of automation to offer trusted, transparent collaboration, situated within human environments."

— David Mindell, MIT

# Autonomy TEV&V R&D Priorities



**Mission requirements**

**HMT CONOPS**

**Autonomy requirements**

**Test oracles**

**Instrumentation needs**

**Evidence (time series)**

**Assurance cases**

**Certification**

**Specifying testable cognitive requirements**

- Perception, reasoning, …
- Teaming and self-organizing
- Negative requirements

**Formal verification methods**

**Instrumenting cognitive functions**

- Aligned with oracles
- Support both assurance and trust

**Virtual test environments**

**V&V of ML training and training data**

**Intelligent adversarial testing**

**Test oracle specification**

- Technology (in)dependent
- Human-machine teaming
- Automated generation of oracles

**Logic of assurance cases**

- Multi-legged arguments
- Combining formal and empirical
- Composability criteria

**Regression testing criteria**

# Examples:  Intelligent Adversarial Testing

The **Range Adversarial Planning Tool** (RAPT) developed at Johns Hopkins Applied Physics Laboratory automates adversarial testing of autonomous systems using simulations of the autonomy software and environment.  RAPT builds a machine-learning model of the autonomy performance and then identifies regions of the configuration space with steep response gradients, indicating possible edge cases.  RAPT then generates test designs that oversample the identified regions.

Similarly, the IBM **Adversarial Robustness Toolbox** (ART) supports the verification of robustness and hardening for machine learning models.

# REPORT DOCUMENTATION PAGE

The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing the burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.
**PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.**

| 1. REPORT DATE *(DD-MM-YYYY)* | 2. REPORT TYPE | 3. DATES COVERED *(From - To)* |
|---|---|---|

**4. TITLE AND SUBTITLE**

**5a. CONTRACT NUMBER**

**5b. GRANT NUMBER**

**5c. PROGRAM ELEMENT NUMBER**

**6. AUTHOR(S)**

**5d. PROJECT NUMBER**

**5e. TASK NUMBER**

**5f. WORK UNIT NUMBER**

**7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)**

**8. PERFORMING ORGANIZATION REPORT NUMBER**

**9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)**

**10. SPONSOR/MONITOR'S ACRONYM(S)**

**11. SPONSOR/MONITOR'S REPORT NUMBER(S)**

**12. DISTRIBUTION/AVAILABILITY STATEMENT**

**13. SUPPLEMENTARY NOTES**

**14. ABSTRACT**

**15. SUBJECT TERMS**

| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT | 18. NUMBER OF PAGES | 19a. NAME OF RESPONSIBLE PERSON |
|---|---|---|---|---|---|
| a. REPORT | b. ABSTRACT | c. THIS PAGE | | | |
| | | | | | 19b. TELEPHONE NUMBER *(Include area code)* |