

# ***Promises, Promises: AI and Certification***

*Bill Scherlis*  
*[scherlis@cmu.edu](mailto:scherlis@cmu.edu)*

*SCC*  
*Annapolis*  
*15 May 2025*

## Modern AI

- Statistical predictions
- Not logic based

## Dimensions of untrustworthiness

- Reactive
- Statistical errors
- Opacity
- Skew
- Cannot actually plan or reason
- Lacks persistent memory
- Learning



## Symbolic AI and hybridization

- Auditability
- Representation and reasoning
- Heuristic power
- Partner to humans

## Beyond

- Grounding
- Expert behavior

## NN weaknesses and vulnerabilities

### Weaknesses – attributes seen as *intrinsic* to the NN concept

- **Fragility and overfitting** – testing success vs. operational success
  - More training is not necessarily better
- **Inexactness and hallucinations** – LLM recall of approximate facts
  - Retrieval augmented generation (RAG) results are processed within the LLM
  - Half of generated software code is broken
- **Memorization** – unexpected exactness
- **Planning and reasoning failures** – LLM cannot plan or reason
  - Chain of thought (CoT) cannot achieve more than a few dozen steps
- **Auditability** – opacity in models impairs rationale/explanation
  - Uncertainty Quantification (UQ) to model ML predictive uncertainties (aleatoric, epistemic)

### Vulnerabilities – attack types *enabled by* NN weaknesses

- **Data poisoning** – injected training data (possibly public), influencing outputs
- **Misdirection** – input adaptation in operations, influencing outputs
  - Black box and white box
- **Model inversion / membership inference** – elicit confidential training inputs
  - ML is same sensitivity as its data, even as a black box. (NB: PPA and differential privacy)
- **Jailbreak** – overcome guardrails on inappropriate LLM behaviors
  - Many of the attacks are portable across vendors
- **Cyber vulnerabilities** – enabled by model implementations

### Hard problems

- **Interpretability** – what does a given parameter “mean”? (NB: operational)
- **Unlearning** – can a trained net be adapted to “forget” without retraining? (NB: spill scenario)
- **Human anthropomorphization** – magical thinking and over-generous acceptance of LLM outputs
- **Deep fake detection** – how to detect (watermark) multimodal creations (NB: detect vs. label)
- **Privacy-preserving analytics** – protect data thru federation, k-anonymity, ... (NB: PETS role)
- **Bias and alignment** – ensure alignment of training and operations (NB: it's all fake alignment)
- **Smaller models** – Down-scaling models for inference support

## Modern AI

- Statistical predictions
- Not logic based

## Dimensions of untrustworthiness

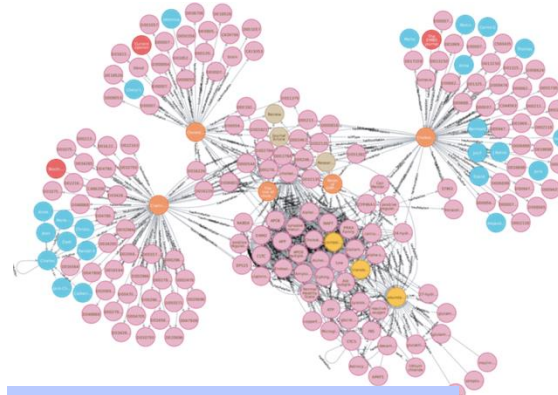
- Reactive
- Statistical errors
- Opacity
- Skew
- Cannot actually plan or reason
- Lacks persistent memory
- Learning

## Symbolic AI and hybridization

- Auditability
- Representation and reasoning
- Heuristic power
- Partner to humans

## Beyond

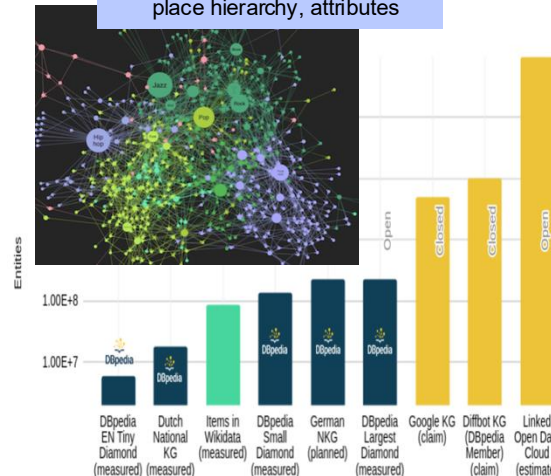
- Grounding
- Expert behavior



Biomedical knowledge graphs

KG: {< uid, node, edge, node >}

Geoname vocabulary  
(25m placenames)  
place hierarchy, attributes



Google knowledge graph  
approaching 10B nodes (2020)

[Term]  
id: DOID:0060013  
name: X-linked severe combined immunodeficiency  
alt\_id: DOID:5811  
def: "A severe combined immunodeficiency that is a X-linked SCID that has material basis in mutations in genes encoding common gamma chain proteins shared by the interleukin (IL-2, 4, 7, 9, 16 and 21) receptors resulting in a non-functional gamma chain, defective interleukin signaling, minimal or absent T- and NK cells and non-functional B-cells."  
[url: [http://en.wikipedia.org/wiki/X-linked\\_severe\\_combined\\_immunodeficiency](http://en.wikipedia.org/wiki/X-linked_severe_combined_immunodeficiency),  
url: <https://ghr.nlm.nih.gov/condition/x-linked-severe-combined-immunodeficiency#synonym>]  
comment: OMM mapping confirmed by DO. [LS]  
subset: DO\_rare\_d  
subset: NCIt\_hesaur  
synonym: "gamma  
synonym: "SCID-X  
synonym: "thymic e  
synonym: "XSCID"  
xref: GARD:5618  
xref: MESH:D053632  
xref: MIM:300400  
xref: NCI:C4682  
xref: SNO MEDCT\_US\_2023\_03\_01203592006  
xref: UMLS\_CUI:C1279481  
is\_a: DOID:6271 severe combined immunodeficiency

Node from a 170kloc  
human disease ontology

way 2: transport a box with a truck  
precondition: truck-at(t,l,f), adjacent(l,f,l,n), box-at(b,l,f),  
in-city(l,f,c,f), in-city(l,n,c,n), same-city(c,f,c,n)  
task network: (load-truck(b,t,l,f), drive(t,l,f,l,n),  
unload-truck(b,t,l,n), deliver(b,l,n,l,t))

Robot planning rules

## Modern AI

- Statistical predictions
- Not logic based

## Dimensions of untrustworthiness

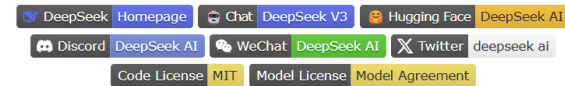
- Reactive
- Statistical errors
- Opacity
- Skew
- Cannot actually plan or reason
- Lacks persistent memory
- Learning

## Symbolic AI and hybridization

- Auditability
- Representation and reasoning
- Heuristic power
- Partner to humans

## Beyond

- Grounding
- Expert behavior

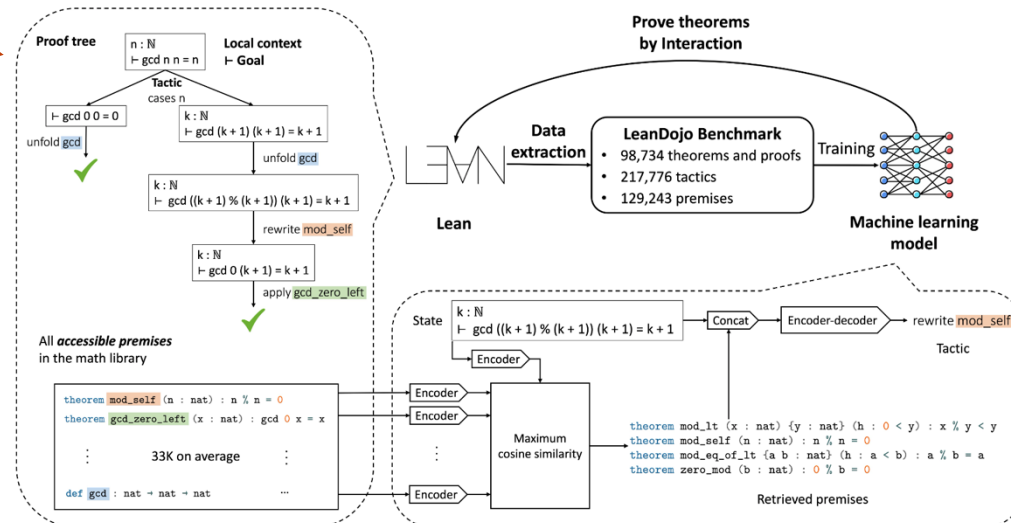


[Paper Link](#)

[Model Summary](#) | [ProverBench](#) | [Model&Dataset Download](#) | [Quick Start](#) | [License](#) | [Contact](#)

## DeepSeek-Prover-V2: Advancing Formal Mathematical Reasoning via Reinforcement Learning for Subgoal Decomposition

### Overview of LeanDojo



## Modern AI

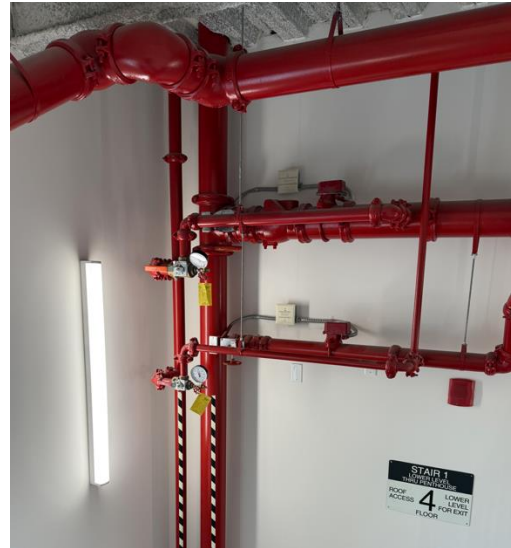
- Statistical predictions
- Not logic based

## Dimensions of untrustworthiness

- Reactive
- Statistical errors
- Opacity
- Skew
- Cannot actually plan or reason
- Lacks persistent memory
- Learning

## Symbolic AI and hybridization

- Auditability
- Representation and reasoning
- Heuristic power
- Partner to humans



## Beyond

- Grounding
- Expert behavior

## Modern AI

- **Statistical predictions**
- **Not logic based**

## Dimensions of untrustworthiness

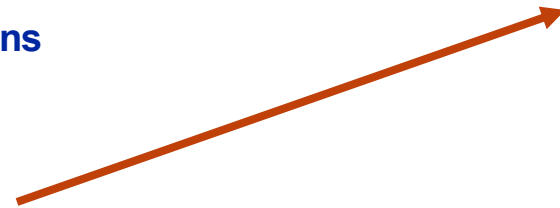
- **Reactive**
- **Statistical errors**
- **Opacity**
- **Skew**
- **Cannot actually plan or reason**
- **Lacks persistent memory**
- **Learning**

## Symbolic AI and hybridization

- **Auditability**
- **Representation and reasoning**
- **Heuristic power**
- **Partner to humans**

## Beyond

- **Grounding**
- **Expert behavior**



## Lessons

- **Dimensions of scale**
  - Scope.
  - Complexity and size.
  - Practice.
  - Evolution.
  - Adoptability.
- **Use case categories**
  - Major systems that matter.
  - Big results.
  - Ecosystems.
  - Broad use.
- **Industry uses**
- **Architecture matters**
- **Productivity can be a selling point**
- **Standards can drive**

## Accepting proofs

- **Social process – but on what?**
- **Physical stuff – the Fallacy of Identification**

## Accepting contingencies

- **What's not modeled**
- **Incorrect behaviors**
- **Validation**
- **Risks**

Table 3-1. Participant Use Cases

Point(s) of Contact	Use Cases	Dimensions of Scale
<ul style="list-style-type: none"> <li>• Clark Barrett</li> </ul>	Solvers for Boolean satisfiability and Satisfiability Modulo Theories (SMT) as industrial workhorses	<ul style="list-style-type: none"> <li>• Complexity and the size of systems</li> <li>• Range of properties and qualities</li> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Darren Cofer</li> <li>• Matt Wilding</li> </ul>	Large aerospace and defense systems	<ul style="list-style-type: none"> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Range of properties and qualities</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Michael Collins</li> <li>• Kristin Giammarco</li> </ul>	Modeling behavior of complex systems concepts supporting reasoning methods critical to national cyber and cryptologic missions	<ul style="list-style-type: none"> <li>• Complexity and size of systems</li> <li>• Range of properties and qualities</li> <li>• Ability to rapidly co-evolve systems and associated evidence within continuous integration/deployment (CI/CD)</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Takeaways from presentation by Byron Cook</li> </ul>	Cloud Services foundational assurance (e.g. cryptography, virtualization, storage)	<ul style="list-style-type: none"> <li>• Range of properties and qualities</li> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Ability to rapidly co-evolve systems and associated evidence within CI/CD</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Mike Dodds</li> <li>• John Launchbury</li> <li>• Stephen Magill</li> </ul>	Cryptography & Formal Methods as a Service	<ul style="list-style-type: none"> <li>• Range of properties and qualities</li> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Ability to rapidly co-evolve systems and associated evidence within CI/CD</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Kathleen Fisher</li> </ul>	Cyber-retrofit of DoD platforms	<ul style="list-style-type: none"> <li>• Range of properties and qualities</li> <li>• Ability to co-evolve systems</li> </ul>
<ul style="list-style-type: none"> <li>• Warren A. Hunt, Jr.</li> <li>• J. Strother Moore</li> </ul>	Creation, analysis, and maintenance of models of computational systems, assist with creating, analyzing, and maintaining models of computational systems developed by in use by companies including AMD, ARM, Centaur Technology, IBM, and Intel	<ul style="list-style-type: none"> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Range of properties and qualities</li> <li>• Complexity and the size of systems</li> <li>• Ability to rapidly co-evolve systems and associated evidence within CI/CD</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Peter O'Hearn</li> </ul>	Scaling the impact of analysis of apps for Android and iOS, Facebook Messenger, Instagram, and other apps	<ul style="list-style-type: none"> <li>• Complexity and size of systems</li> <li>• Range of properties and qualities</li> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Ability to rapidly co-evolve systems and associated evidence within CI/CD</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Ray Richards</li> </ul>	DoD military systems	<ul style="list-style-type: none"> <li>• Complexity and size of systems</li> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Ability to rapidly co-evolve systems and associated evidence within CI/CD</li> <li>• Ease of use</li> </ul>

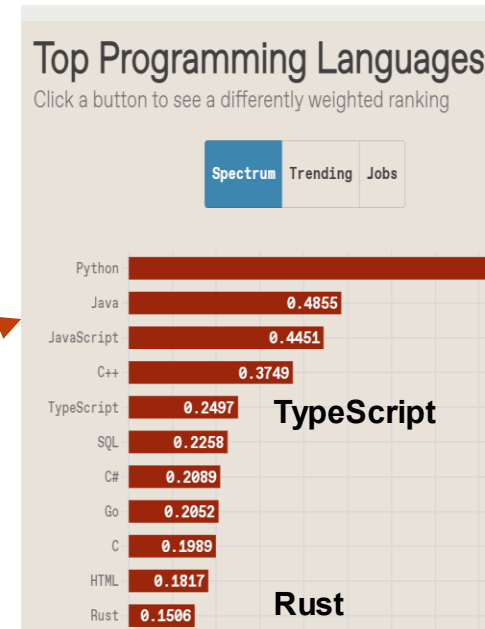
Formal Methods at Scale: 2019 Workshops Report

Point(s) of Contact	Use Cases	Dimensions of Scale
<ul style="list-style-type: none"> <li>• Natarajan Shankar</li> </ul>	Designing, analyzing, and creating computer systems	<ul style="list-style-type: none"> <li>• Complexity and size of systems</li> <li>• Range of properties and qualities</li> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Ability to rapidly co-evolve systems and associated evidence within CI/CD</li> <li>• Ease of use</li> </ul>
<ul style="list-style-type: none"> <li>• Nikhil Swamy</li> </ul>	Components in HTTPS ecosystem, including transport layer security, the main protocol at the heart of HTTPS, as well as the main underlying cryptographic algorithms, such as AES, SHA2 or X25519	<ul style="list-style-type: none"> <li>• Complexity and size of systems</li> <li>• Range of properties and qualities</li> <li>• Effectiveness and efficiency of formal methods-related modeling and tooling</li> <li>• Ability to rapidly co-evolve systems and associated evidence</li> <li>• Ease of use</li> </ul>



## Lessons

- **Dimensions of scale**
  - Scope.
  - Complexity and size.
  - Practice.
  - Evolution.
  - Adoptability.
- **Use case categories**
  - Major systems that matter.
  - Big results.
  - Ecosystems.
  - Broad use.
- **Industry uses**
- **Architecture matters**
- **Productivity can be a selling point**
- **Standards can drive**



## Accepting proofs

- **Social process – but on what?**
- **Physical stuff – the Fallacy of Identification**

## Accepting contingencies

- **What's not modeled**
- **Incorrect behaviors**
- **Validation**
- **Risks**



# Information Loss

## Software Understanding

- **NB: Report and workshop**
- **Broadly inclusive problem**


## Nature of evidence


- **Code and beyond**
- **Rules of the road**
- **Higher level models and patterns**
- **Development artifacts, informal and formal**
- **Argumentation**

## Tightening feedback loops

- **Defer stimulus**
- **Advance response**

## Roles for FM and heuristic AI

TLP: CLEAR



## Closing the Software Understanding Gap

Publication: January 16, 2025

Cybersecurity and Infrastructure Security Agency  
Defense Advanced Research Projects Agency  
Office of the Under Secretary of Defense for Research and Engineering  
National Security Agency

This document is marked TLP: CLEAR. Disclosure is not limited. Sources may use TLP: CLEAR when information carries minimal or no foreseeable risk of misuse, in accordance with applicable rules and procedures for public release. Subject to standard copyright rules, TLP: CLEAR information may be distributed without restriction. For more information on the Traffic Light Protocol, see [cisa.gov/tlp](https://cisa.gov/tlp).

TLP: CLEAR

# Information Loss

## Software Understanding

- **NB: Report and workshop**
- **Broadly inclusive problem**

## Nature of evidence

- **Code and beyond**
- **Rules of the road**
- **Higher level models and patterns**
- **Development artifacts, informal and formal**
- **Argumentation**

## Tightening feedback loops

- **Defer stimulus**
- **Advance response**

## Roles for FM and heuristic AI



# Risk Modeling and Mitigation

Integration context for overall risk management

- Evidence role

The inadequacy of traditional models

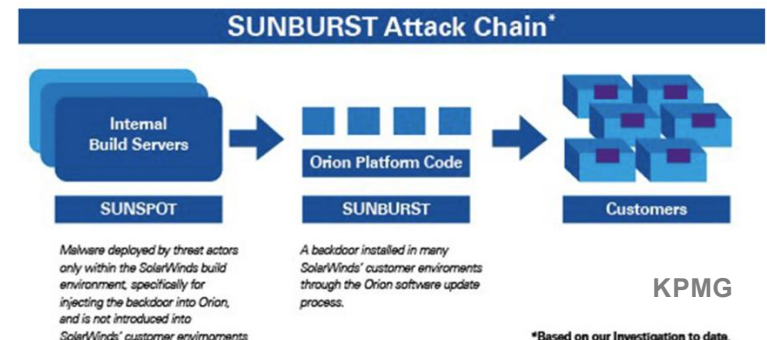
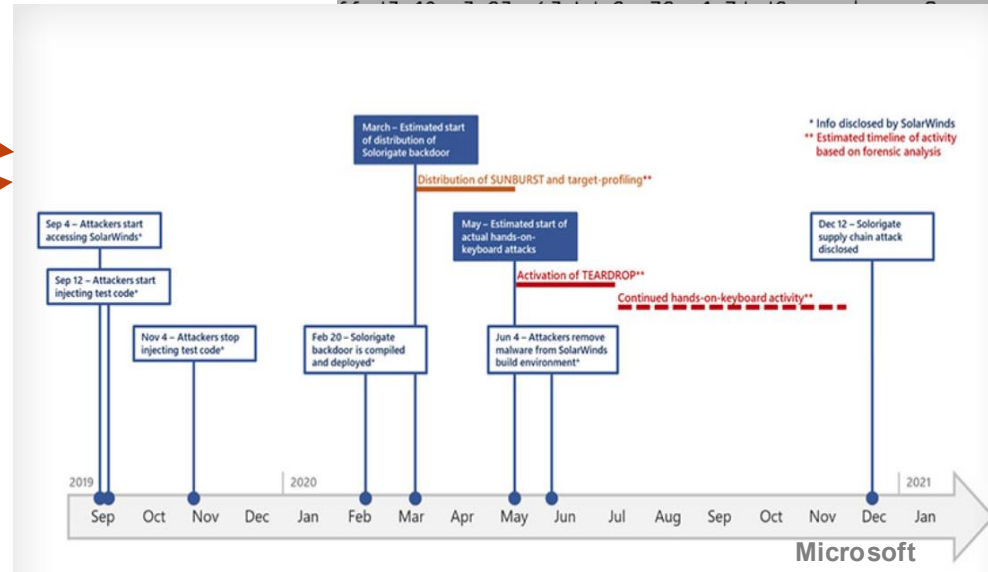
- Actuarial risk

Cyber difficulties

- Adversaries
- Opacity
- Hidden correlations
- Context of system operations
- Vulnerabilities
- Attack surfaces

AI difficulties

- Abundant attack surfaces
- Dimensions of consequence
- Opacity and auditing



# Opportunities for Advancing FM and AI

---

## AI

- AI for FM
- Modeling and proving
- Proof maintenance
- Hybrid AI for high impact
  - Auditability and trustworthiness

## FM

- Usability and integration into process
- Reduce information loss -- the role of evidence
- Find the math; hide the math
- AI assist for model development and provers
- Model scope and expressiveness

## Architecture

- Design for verifiability
- Resilience
- Maximizing architecture benefits
- Modeling and analysis -- move to left

## Attack surface modeling

- Modeling
- Mitigation
- Work factor interventions