# Letting Go of the Numbers: Measuring AI Trustworthiness
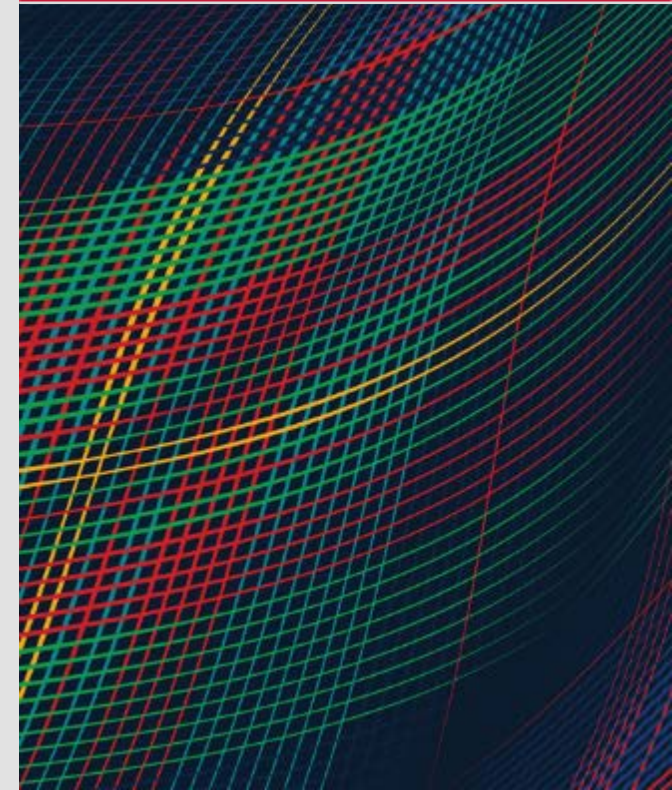
**ICPRAM 2024**

Carol J. Smith

Trust Lab Lead & Principal Research Scientist, AI Division

Carnegie Mellon University
Software Engineering Institute

# Copyright Statement

# About ACM

ACM, the Association for Computing Machinery (www.acm.org), is the premier global community of computing professionals and students with **nearly 100,000 members in more than 170 countries** interacting with more than 2 million computing professionals worldwide.

OUR MISSION: We help computing professionals to be their best and most creative. We connect them to their peers, to what the latest developments, and **inspire them to advance the profession and make a positive impact on society**.

OUR VISION: We see a world where **computing helps solve tomorrow's problems** – where we use our knowledge and skills to advance the computing profession and make a positive social impact throughout the world.

I am proud to be an ACM Member.

# The Distinguished Speakers Program
## is made possible by

**Association for
Computing Machinery**

*Advancing Computing as a Science & Profession*

For additional information, please visit http://dsp.acm.org/
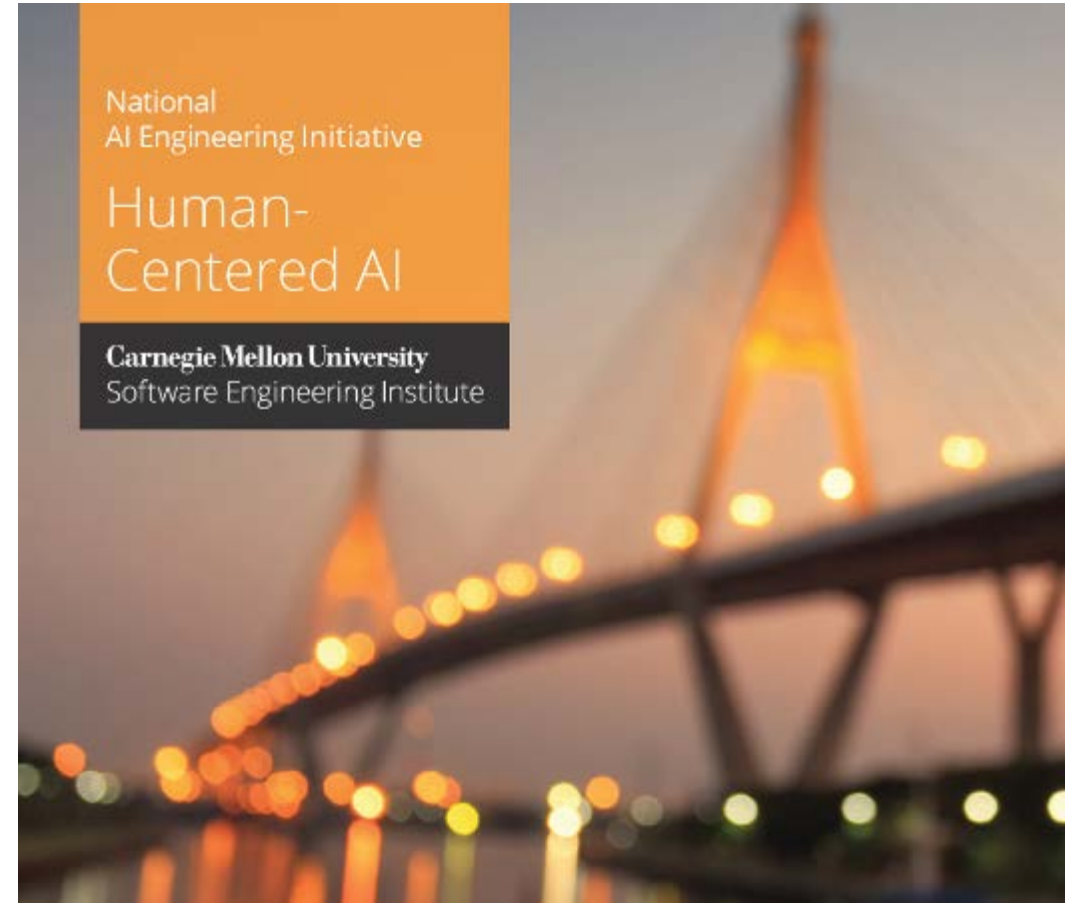
# Human-Machine Teaming



Video: Tesla Autopilot in Heavy LA Traffic by Scott Kubo  https://youtu.be/m3-QzTFxoUg?t=14

# Engineering for Trustworthy AI

Trustworthy AI systems are **designed to work with, and for, people**.

- built for a specific context of use (fit with user needs and tasks)
- with appropriate data, and are
- reliable (robust and secure).

Capabilities are understood, and continuous monitoring and oversight are prioritized.



Human-Centered AI, Software Engineering Institute: https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=735362

# AI must be designed to work with, and for, people. Trustworthy, human-centered, and responsible.

# How do we Measure Trustworthiness?

# Can we accurately predict the future?



Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

## Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

# Can we use data to reduce bias in systematically prejudiced organizations?

# Can Anyone?

Amazon scraps secret AI recruiting tool that showed bias against women
By Jeffrey Dastin. October 9, 2018. Reuters.



Google 'fixed' its racist algorithm by removing gorillas from its image-labeling tech. By James Vincent  Jan 12, 2018. The Verge.

# How about generative AI or LLMs?

Jason Allen's A.I.-generated work, "Théâtre D'opéra Spatial," took first place in the digital category at the Colorado State Fair.Credit…
via Jason Allen. https://www.nytimes.com/2022/09/02/technology/ai-artificial-intelligence-artists.html

HAPPINESS    4.185
NEUTRAL      0.901
SURPRISE     89.864
SADNESS      0.01
DISGUST      0.01
ANGER        5.021
FEAR         0.01

Image by Comuzi / © BBC / Better Images of AI / Mirror D / CC-BY 4.0

# Humans create and use imperfect machines.

# Quant Performance Evaluations are Necessary

- Evaluate accuracy, precision, recall
- Ensure it is robust, secure, reliable
- Speed of system
- Scalability

Relative simplicity of these methods is appealing,
but **these are not adequate**.

# Overnight Flight from US to Rome
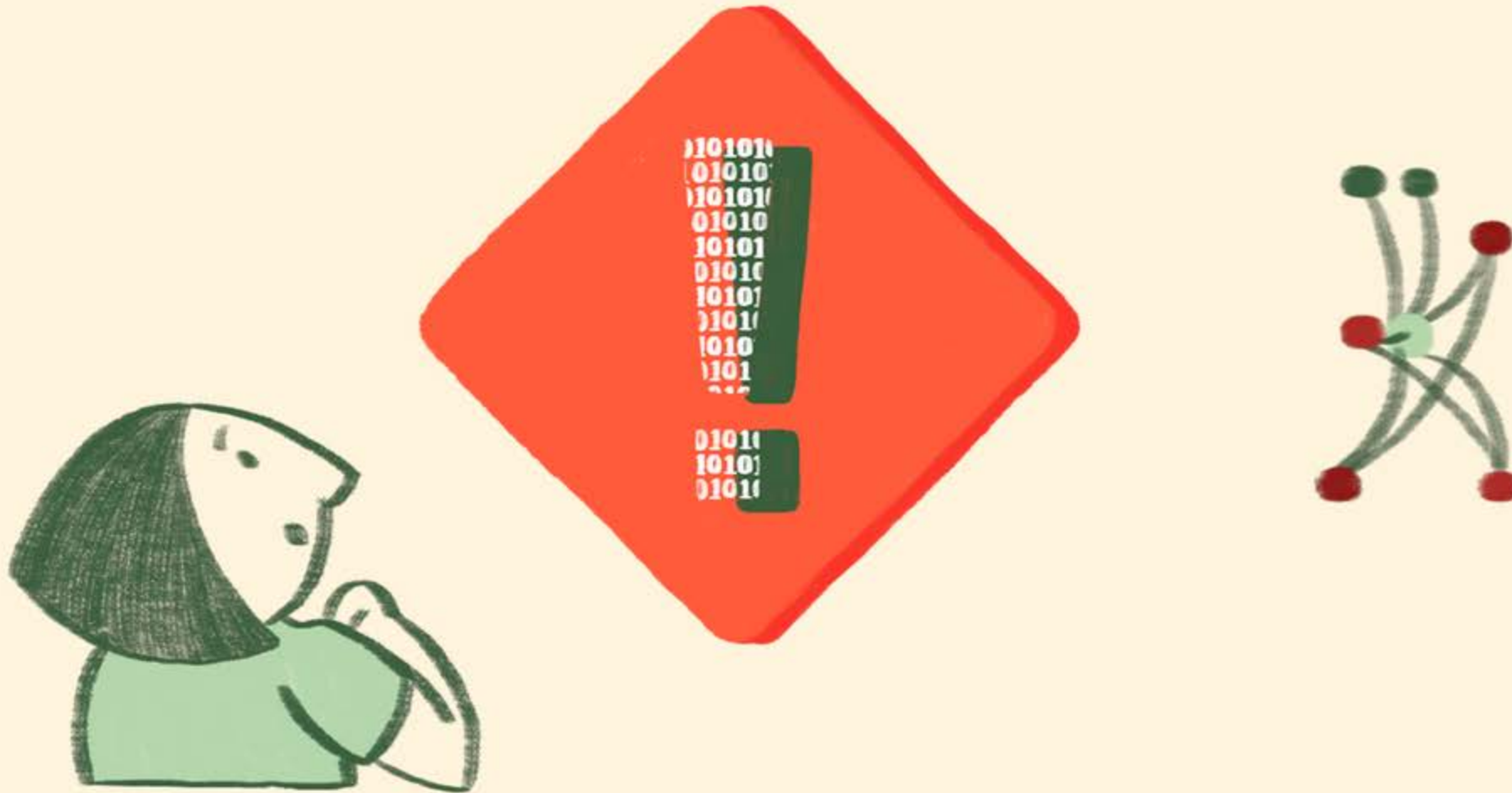
## Quantitative

- Plane arrived 20+ min. early.
- Reduced fuel use.
- Reduced emissions.

## Qualitative

- Delayed meal delivery.
- Reduced sleep time.
- An uncomfortable night.

# Overnight Flight from US to Rome

## Quantitative

- Plane arrived 20+ min. early.
- Reduced fuel use.
- Reduced emissions.

## Qualitative

- Delayed meal delivery.
- Reduced sleep time.
- An uncomfortable night.

# Trustworthiness Requires Qualitative Measures

# All systems will have some form of bias

Complete objectivity is misleading.
Bias can have purpose and can be helpful.
Bias contributes to and is emphasized by decisions.

We must ensure we
- identify and understand bias
- reduce unintended and/or harmful bias.

Risks due to bias are lower when no information about people is present.

# Bias Due to Data, Algorithm, and Training

Photo by sunlightfoundation
https://www.flickr.com/photos/sunlightfoundation/2385174105

> "Data is a function of our history...
> The past dwells within...
> Showing us the inequalities
> that have always been there."

Joy Buolamwini, Algorithmic Justice League
Coded Gaze
Movie: Coded Bias on Netflix

Photo: Joy Buolamwini on The Open Mind: Algorithmic Justice.
Jan 12, 2019. https://www.youtube.com/watch?v=hwHnXdoSSFY

# Start with Data

Carnegie
Mellon
University
Software
Engineering
Institute

Need to confirm:

- Appropriateness
- Provenance and understanding of data composition and variance
- Information about people and potential risks



Sample data card, source: Pushkarna et al., 2022

# Spotted Lanternfly

Spotted lanternfly Life Cycle. Published by Oxford University Press on behalf of Entomological Society of America 2021., Public domain, via Wikimedia Commons.

# Data Provenance

- Researcher's motivation
- Collection process
- Data included and excluded
  - Which stage of life cycle?
  - Locations?
- Recommended uses, etc.
- Historical patterns of negative bias
- Sensitivity of data



Spotted lanternfly displaying underwing.
WanderingMogwai, CC BY-SA, via Wikimedia Commons

# Identification of Inherent Bias

Understand inherent bias and amount of variance in dataset due to data provenance.

Bias can be both purposeful and unintended influences
- **Purposeful**: provenance of data, collection process, etc.
- **Unintended**: existing systematic bias that may or may not be known or is only revealed as the system is developed.

**All systems are biased.**

# Each decision creates and affects bias.

# Bias can have purpose and can be helpful.

# Unwanted bias can lead to inequitable outcomes

At the surface, AI systems can seem objective and impartial.

Digging deeper reveals that AI systems can reinforce discrimination against historically marginalized groups

- Alignment problems
- Scale problems
- Multiplicity problems

Fairness research led by Anusha Sinha, AI Division at the SEI

# Bias can result in the right decision for the wrong reasons

## Images correctly classified as "balance beam"



Original images sourced from ImageNet

Fairness research led by Anusha Sinha, AI Division at the SEI

# Wrong reasons can lead to poor real-world performance

A low-stakes example:



Source: ImageNet

Ground truth: horizontal bar

Predicted: balance beam

A high-stakes example:



17:48

Ground truth: Carol in town during protest

Predicted: Carol organized protest

Fairness research led by Anusha Sinha, AI Division at the SEI

# Mitigation of Bias is Complex

**Removing all bias is impossible**
- Removing obvious indicators (gender, zip code, etc.) reduces the ability to track bias.
- Invisible indicators are concealed in the data.
- Share awareness of bias for all audiences (developers, purchasers, users).

# Getting to Trustworthiness

**An AI system's potential is bound to stakeholders' perceptions of its trustworthiness**

32

# Capitalize on Human Strengths

Humans are (still) better
at many activities:

Exposing Bias

Identifying downstream impacts

Judgment
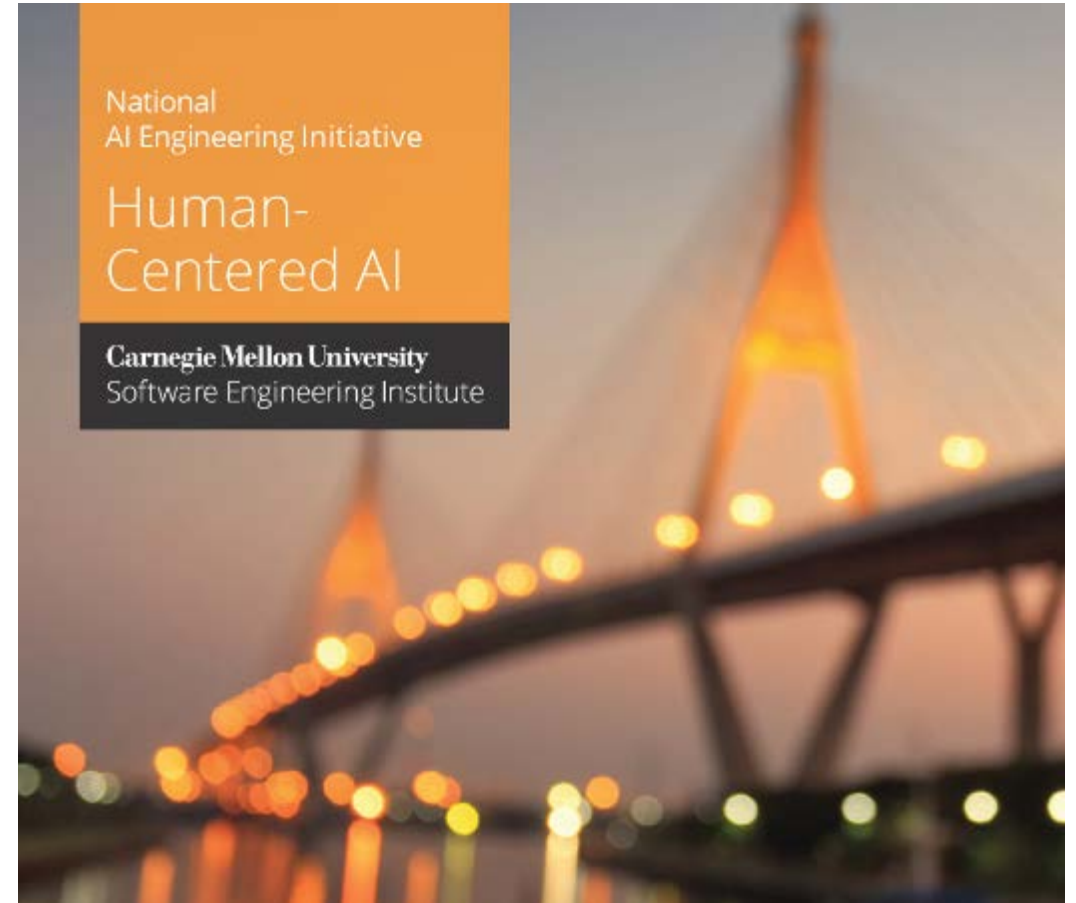
Recognizing Bias

Responding to change

Socio-political nuance

Taking context into consideration

Amanda Muller and Carol Smith. 2022. Perceptions of Function Allocation between Humans and AI-Enabled Systems. UXPA 2022 (pre-print).
https://uxpa2022.org/sessions/perceptions-of-function-allocation-between-humans-and-ai-enabled-systems/

# Trust Should Not be the Default



- ## Dynamic systems
  Data drift, poisoning, system failures

- ## Dynamic contexts
  Weather, adversaries

- ## Human judgement
  Intuition, situational awareness, fatigue

Human-Centered AI, Software Engineering Institute:
https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=735362

# Trust is Contextual

Trust is personal - a dynamic psychological state.
We calibrate trust based on personal experiences, current context, and available evidence of system's capability and integrity.

**Distrust**
Trust falling short of system capabilities - may lead to disuse.

**Calibrated Trust**
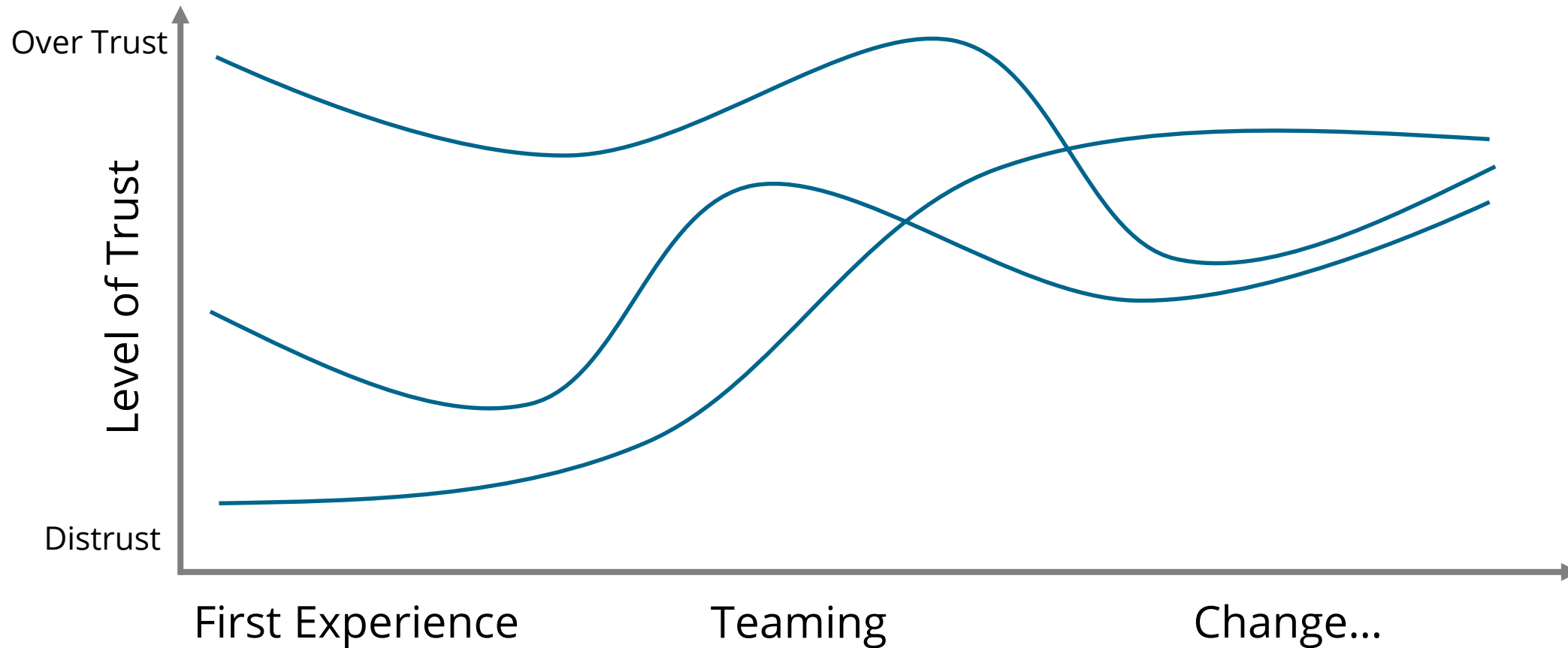Trust matches system capabilities - leading to appropriate use.

**Over Trust**
Trust exceeding system capabilities - may lead to misuse.

**Rejection.**

**Automation bias.**

John D. Lee and Katrina A. See. 2004. Trust in Automation: Designing for Appropriate Reliance. Hum Factors 46, 1 (March 2004), 50–80. DOI:https://doi.org/10.1518/hfes.46.1.50_30392
Bobbie Seppelt and John Lee. 2012. Human Factors and Ergonomics in Automation Design. In Handbook of Human Factors and Ergonomics (Fourth Edition) Chapter 59. Wiley.
DOI:  https://doi.org/10.1002/9781118131350.ch59

# Trust is Complex and Transient

# Design for Use, Context, Trustworthiness

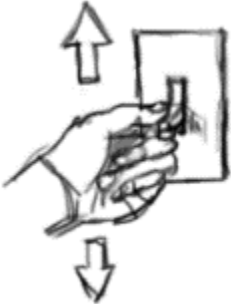Who will use it?

Scientist?

Grocer?

# Understand Stakeholders

- Who will use the system?
- How well are current systems accepted?
- What are the existing issues?
- When and in what context?

# Trustworthy Systems

- Uphold Responsible AI principles
- Utilize data appropriate for task
- Designed for the human-machine team
  to complete their mission
- Augment human teammates and meet their needs
  (human-centered)
- Consistently provide adequate evidence of current
  capabilities and integrity in the current context.
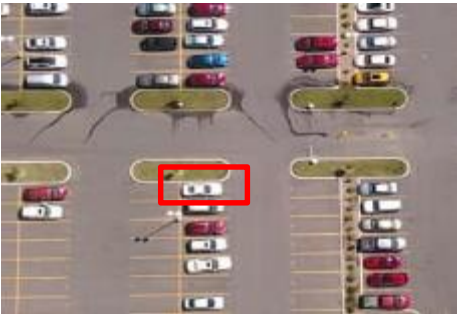
# Measurements of Trustworthiness
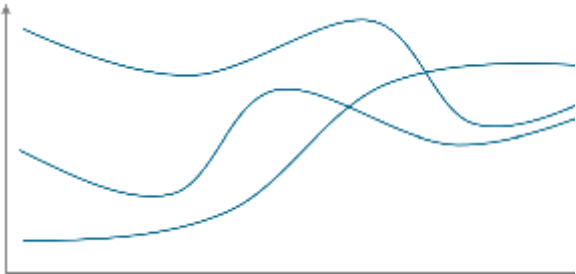
**Usability**



**Explainability**



**Fairness**



**Likelihood of Failure**
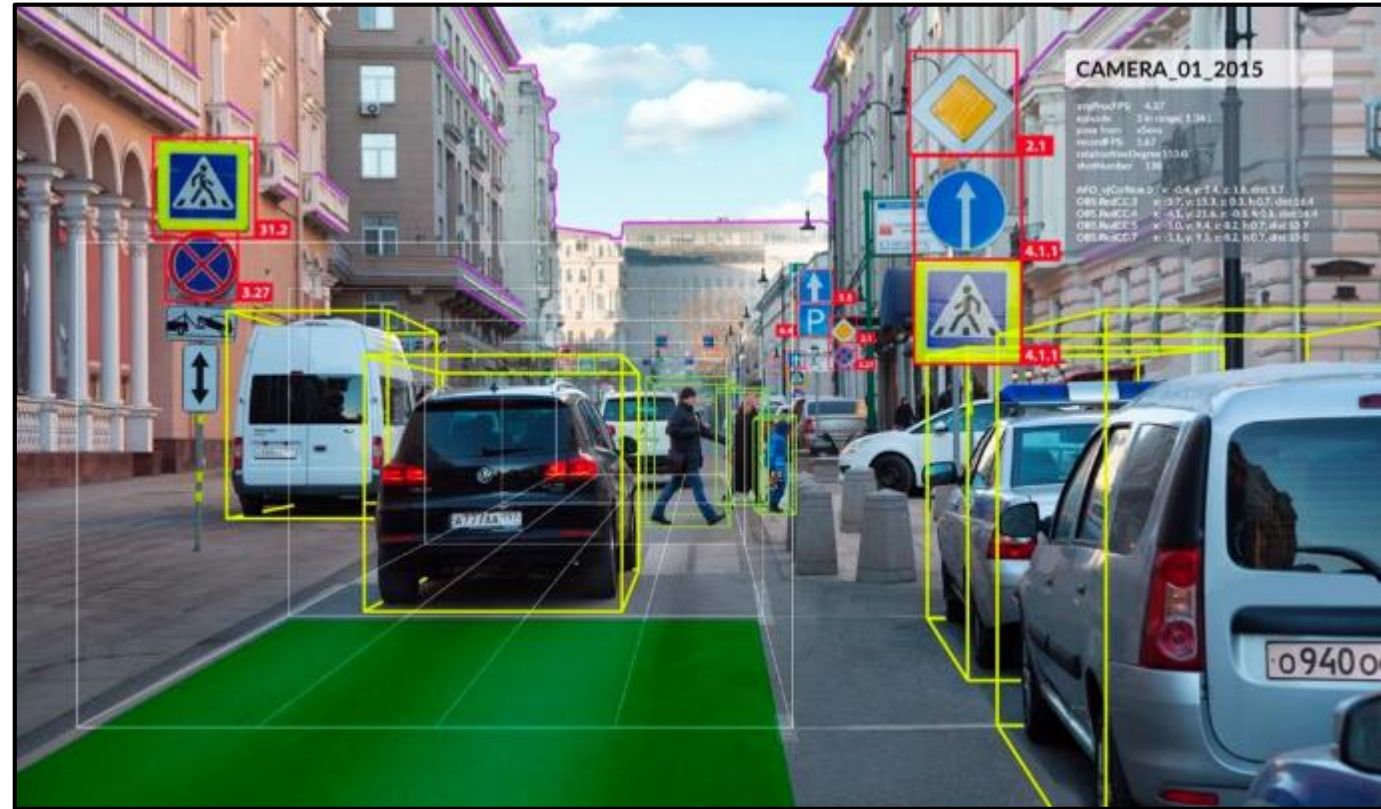


**Usage**

**How is the AI system making decisions?**
# Explainability

Photo GOODFELLOW AIR FORCE BASE, TX, UNITED STATES, 11.05.2020by Airman 1st Class Ethan Sherwood, 17th Training Wing Public Affairs https://www.dvidshub.net/image/6443325/drones-goodfellow

# Explainability reveals decision-making processes

Interpretability facilitates optimization and evaluation
([Doshi-Velez & Kim, 2017](#))

- Safety
- Ethics
- Mismatched objectives
- Multi-objective tradeoffs



Example of a computer vision system, Source: [Welker Media](#)

Explainability research led by Violet Turri, AI Division Trust Lab at the SEI

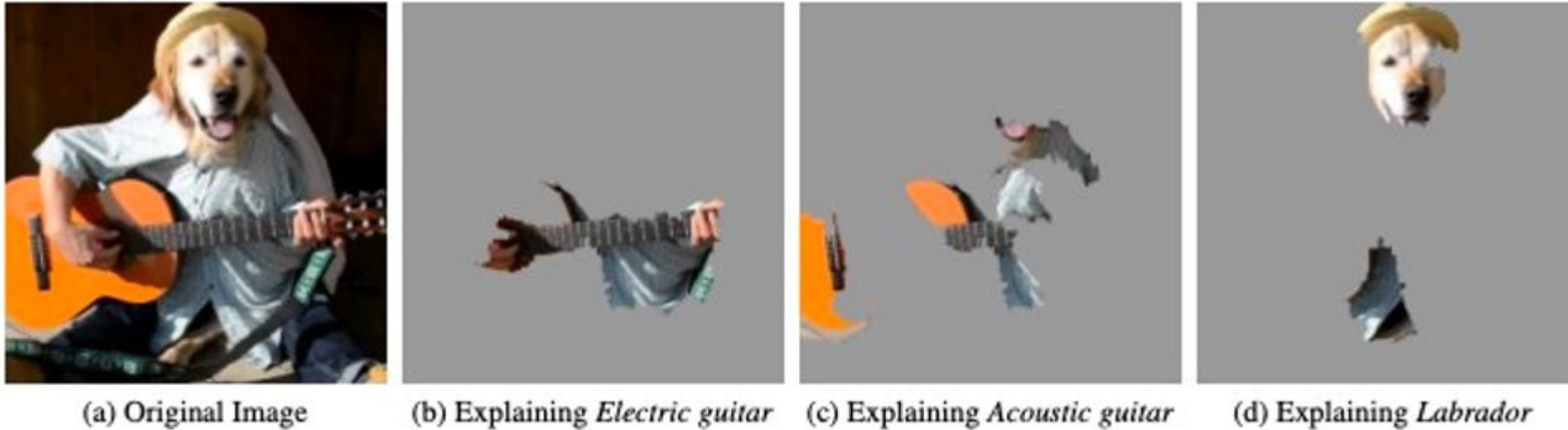# Explanations can Illuminate Unintended System Behavior

(a) Original Image     (b) Explaining *Electric guitar*     (c) Explaining *Acoustic guitar*     (d) Explaining *Labrador*

Figure 4: Explaining an image classification prediction made by Google's Inception network, high-lighting positive pixels. The top 3 classes predicted are "Electric Guitar" ($p = 0.32$), "Acoustic guitar" ($p = 0.24$) and "Labrador" ($p = 0.21$)

Sample explanation of an image classifier, Source: Explainable AI: current status and future directions

Explainability research led by Violet Turri, AI Division Trust Lab at the SEI

# Understanding
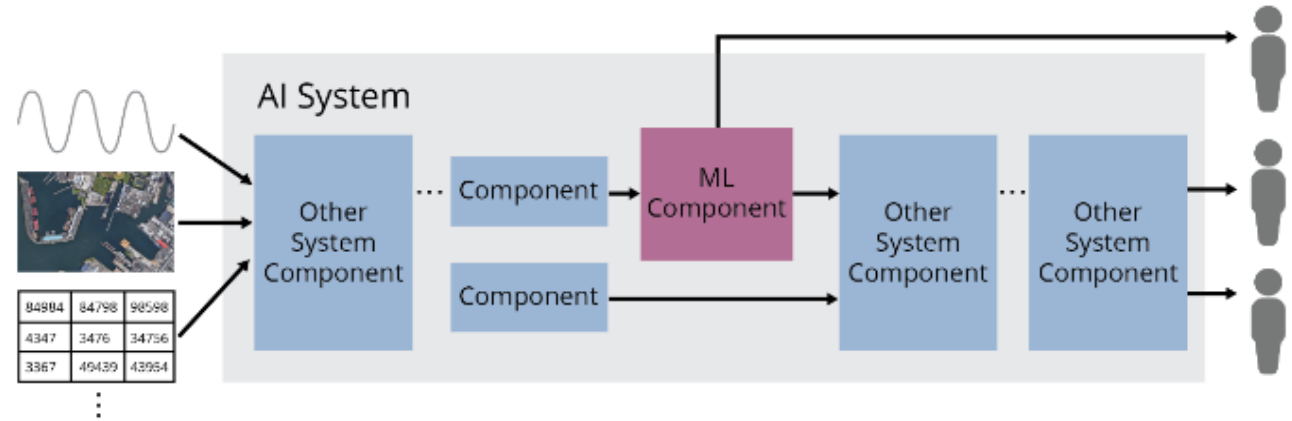# the Likelihood of Failure

# Accurate confidence measures can inform better decision making in complex contexts

## People need
- situational awareness (system and context), and
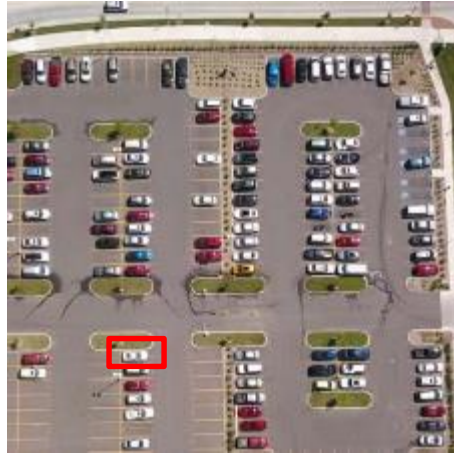- probability of failure.

## Decide what to do next
- inform other parts of the system
- alert an analyst, use another sensor, etc.



Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why is confidence important?

0.2203 Confident        0.9637 Confident

- More informed decision making and prioritization
- Focus on the car on the right
- Use additional resources to confirm

Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why are context and shifting environments important?

Many training data sets do not provide sufficient coverage of cases that can be encountered in the deployment environment.
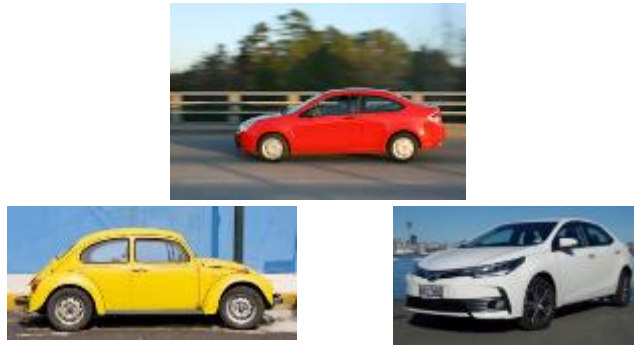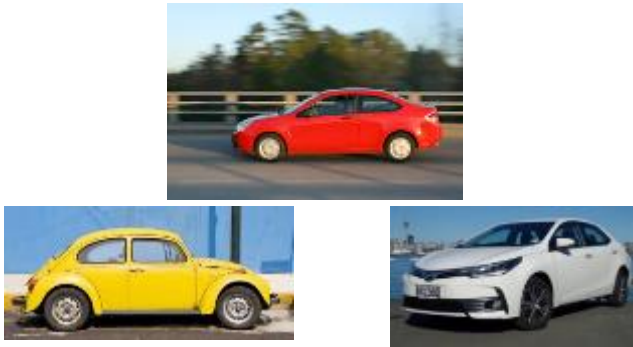


Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why are context and shifting environments important?

Many training data sets do not provide sufficient coverage of cases that can be encountered in the deployment environment.

Train Set

Encountered During Deployment
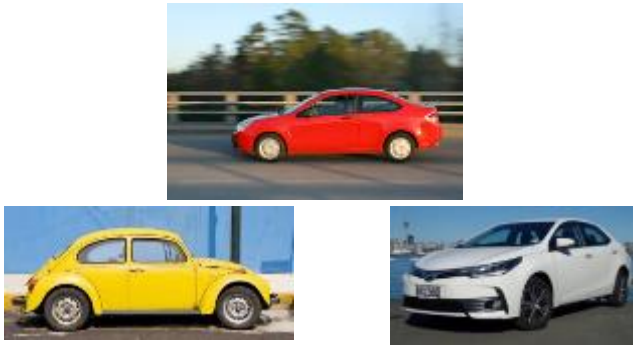


Potential causes of shifts:
- Sensor failure or degradation

Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why are context and shifting environments important?

Many training data sets do not provide sufficient coverage of cases that can be encountered in the deployment environment.

Train Set

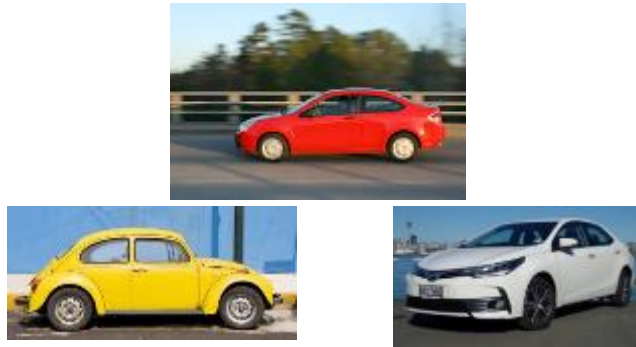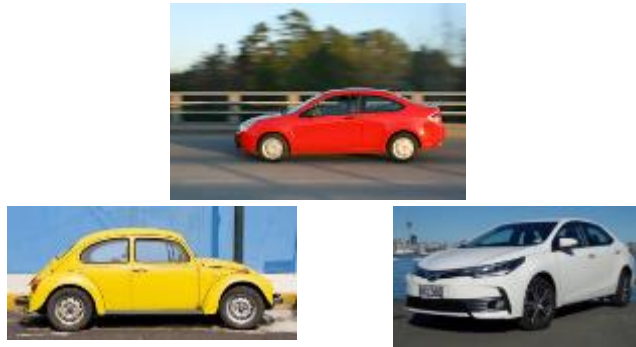Encountered During Deployment



Potential causes of shifts:
- Sensor failure or degradation
- Unidentified biases

Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why are context and shifting environments important?

Many training data sets do not provide sufficient coverage of cases that can be encountered in the deployment environment.

Train Set                                              Encountered During Deployment
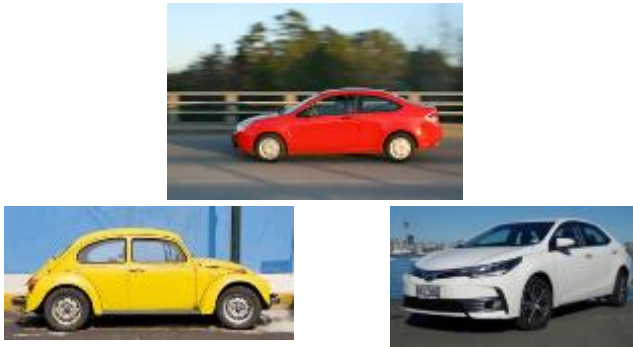


Potential causes of shifts:
- Sensor failure or degradation
- Unidentified biases
- Unaccounted for changes in data pipelines

Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why are context and shifting environments important?

Many training data sets do not provide sufficient coverage of cases that can be encountered in the deployment environment.



Train Set



Encountered During Deployment

Potential causes of shifts:
- Sensor failure or degradation
- Unidentified biases
- Unaccounted for changes in data pipelines
- Change in context

Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why are context and shifting environments important?

Many training data sets do not provide sufficient coverage of cases that can be encountered in the deployment environment.

Train Set

Encountered During Deployment



Potential causes of shifts:
- Sensor failure or degradation
- Unidentified biases
- Unaccounted for changes in data pipelines
- Change in context
- Rare, but possible cases

Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why are context and shifting environments important?

Many training data sets do not provide sufficient coverage of cases that can be encountered in the deployment environment.

Train Set



Encountered During Deployment



Potential causes of shifts:
- Sensor failure or degradation
- Unidentified biases
- Unaccounted for changes in data pipelines
- Change in context
- Rare, but possible cases
- Novel, but relevant classes

Uncertainty Quantification research led by Eric Heim, AI Division at the SEI

# Why is the model uncertain?

Carnegie
Mellon
University
Software
Engineering
Institute

- What is the cause?
- How do we want the system to respond when it encounters new information – new situations?
- What is the appropriate way to communicate the likelihood of failure?

# Testing for Trustworthiness

# Measurements of Trustworthiness

**Usability**



**Explainability**



**Fairness**



**Likelihood of Failure**



**Usage**

# Conversations for Understanding



Difficult Topics
- What do we value?
- Who could be hurt?
- What lines won't our AI cross?
- How are we shifting power?*

*"Don't ask if artificial intelligence is good or fair, ask how it shifts power." Pratyusha Kalluri.
https://www.nature.com/articles/d41586-020-02003-2

# Speculation Keeps People Safe

# Activate Curiosity

Speculate about misuse and abuse
- Unintended and unwanted consequences
- Negative consequences for people
  who are frequently marginalized

Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020.

# Card Game: What Could Go Wrong?
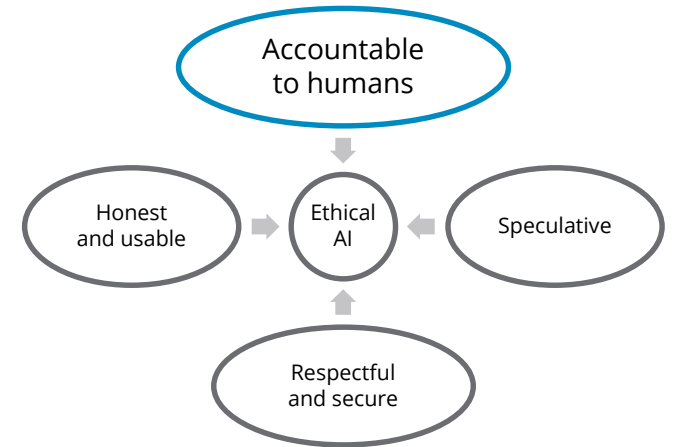
Foster conversations around potential challenges and issues with complex technologies.



Figure 1: What Could Go Wrong? A digital card game to foster conversation around potential challenges and issues of autonomous vehicle technology.

Nikolas Martelaro and Wendy Ju. 2020. What Could Go Wrong? Exploring the Downsides of Autonomous Vehicles. In 12th International Conference on Automotive User Interfaces and Interactive Vehicular Applications (AutomotiveUI '20). Association for Computing Machinery, New York, NY, USA, 99–101. https://doi.org/10.1145/3409251.3411734

# Abusability Testing

1) **Value proposition**
   Benefits tech brings to individuals, society

2) **Vulnerabilities**
   How tech could be misused

3) **Abuse scenario**

Provocation via prompt statements

UX in the Age of Abusability. The role of Composition, Collaboration, and Craft in building ethical products. Dan Brown. Sep 18, 2018. https://greenonions.com/ux-in-the-age-of-abusability-797cd01f6b13
Photo from workshop organized by Anna Abovyan, Theora Kvitka and Allison Cosby of the Pittsburgh IxDA Chapter for World Interaction Design Day 2019.



Template by: Anna Abovyan & Allison Cosby, IxDA Pittsburgh, Sep 2019

# Prototype: Make Informed Design Choices



Button - Push    Switch - Flip    Knob - Rotate

Light Feedback



REMOTE START
LOCK
GARAGE DOOR
MANUAL



Drawings of Affordance: http://paaralan.blogspot.com/2010/09/affordance-and-educational-games.html

# Significant Decisions

Made by system
- explained
- able to be overridden
- appealable and reversible

Responsibilities are explicitly defined
between people and systems.

# Humans are Accountable

Ensure humans have ultimate control.
Able to monitor and control risk.

A person is always responsible for final decisions:

- Person's life
- Quality of life
- Health
- Reputation



Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020.

"Ensure humans can unplug the machines"
– Grady Booch

TED Talk, Grady Booch, Scientist, Philosopher, IBM'er
https://www.ted.com/talks/grady_booch_don_t_fear_superintelligence

# Iterative Cycles: Feedback and Improvement

**Test Prototype with Users**

**Analyze & Prioritize**

1 High

2 Medium

3 Low

**Iterate**

**Repeat**

# Reward team members for finding ethics bugs

Ayanna Howard

Let go of Some of the Numbers

# Rely More on Qualitative Information

# What do the people who will use the system expect?

# Provide Evidence

# Design AI to work with, and for, people

# Carol J. Smith

AI Division Trust Lab Lead
Principal Research Scientist

Email: cjsmith@andrew.cmu.edu
LinkedIn: https://www.linkedin.com/in/caroljsmith/

# Appendix

# Adopt Technology Ethics

Harmonize cultural variations.

Balance to pace of change.

Explicit permission to consider
and question breadth
of implications.

# Prompt conversations

Checklists, frameworks, and guidelines – pair with technical ethics.

- Bridge gaps between "do no harm" and reality
- Support inspection and mitigation planning

Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020. Checklist and Agreement - Downloadable PDF: https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=636620
Defense Innovation Unit. Artificial Intelligence Portfolio, Responsible AI Guidelines. https://www.diu.mil/responsible-ai-guidelines

# Tools to Support Conversations for Understanding



Pair DoD Ethical Principles for AI (or another set) with frameworks and tools that provoke discussion on relevant topics.





Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020.
Checklist and Agreement - Downloadable PDF: https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=636620
Photo by Pam Sharpe https://unsplash.com/@msgrace?utm_source=unsplash&amp;utm_medium=referral&amp;utm_content=creditCopyText On Unsplash
https://unsplash.com/s/photos/business-woman-smiling?utm_source=unsplash&amp;utm_medium=referral&amp;utm_content=creditCopyText

# Publications

# Defense Innovation Unit
## RAI Report, Guidelines, Worksheets, and Workshops



Defense Innovation Unit. Artificial Intelligence Portfolio, Responsible AI Guidelines. https://www.diu.mil/responsible-ai-guidelines

# Usable Hazard Analysis for AI Engineering

Support teams making complex systems,
in early risk identification.



**Exploring Opportunities in Usable Hazard Analysis Processes for AI Engineering**

Nikolas Martelaro,[1] Carol J. Smith,[2] Tamara Zilovic[1]

HCI Institute - Carnegie Mellon University,[1] Software Engineering Institute, Carnegie Mellon University[2]
nikmart@cmu.edu, cjsmith@sei.cmu.edu, tzilovic@andrew.cmu.edu

**Abstract**

Embedding artificial intelligence into systems introduces significant challenges to modern engineering practices. Hazard analysis tools and processes have not yet been adequately adapted to the new paradigm. This paper describes initial research and findings regarding current practices in AI-related hazard analysis and on the tools used to conduct this work. Our goal with this initial research is to better understand the needs of practitioners and the emerging challenges of considering hazards and risks for AI-enabled products and services. Our primary research question is: *Can we develop new structured thinking methods and systems engineering tools to support effective and engaging ways for preemptively considering failure modes in AI systems?* The preliminary findings from our review of the literature and interviews with practitioners highlight various challenges around

implications for the organizations that develop these products. While the use of new technologies always comes with the possibility of unintended consequences, we believe that many of these examples could have been prevented through strategic and thoughtful consideration when these systems are being designed and engineered.

Within systems engineering, strategies for hazard analysis can be used by teams to identify risks and potential failures with the goal of developing more robust and safe engineered systems. While many formal hazard analysis techniques exist, these activities largely center around helping teams determine potential risks and/or sources of failure *before* products have begun the development

tinyurl.com/hazards-ai-eng

Nikolas Martelaro, Carol J. Smith, and Tamara Zilovic. 2022. Exploring Opportunities in Usable Hazard Analysis Processes for AI Engineering. Presented at 2022 AAAI Spring Symposium Series Workshop on AI Engineering: Creating Scalable, Human-Centered and Robust AI Systems. arXiv:2203.15628 [cs] (March 2022).

Pair DoD Ethical Principles for AI (or another set) with frameworks and tools that provoke discussion on relevant topics.



Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020.
Checklist and Agreement - Downloadable PDF: https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=636620
Photo by Pam Sharpe https://unsplash.com/@msgrace?utm_source=unsplash&amp;utm_medium=referral&amp;utm_content=creditCopyText On Unsplash - https://unsplash.com/s/photos/business-woman-smiling?utm_source=unsplash&amp;utm_medium=referral&amp;utm_content=creditCopyText

# Additional Publications

J. Dunnmon, B. Goodman, P. Kirechu, C. Smith, A. Van Deusen. "Responsible AI Guidelines in Practice: Lessons Learned from the DIU AI Portfolio."  DIU.



H. Barmer; R. Dzombak; M. Gaston; V. Palat; F. Redner; C. Smith; et al. (2021): "Human-Centered AI." SEI, CMU.

- Blog: <u>Contextualizing End-User Needs: How to Measure the Trustworthiness of an AI System</u>
- Checklist: <u>Designing Ethical AI Experiences: Checklist and Agreement</u>
- Whitepaper: <u>SEI: Human-Centered AI</u>
- Blog: <u>What is explainable AI?</u>
- Video: <u>Collaboration Conversation: Human-Centered AI</u>
- Video: <u>Implementing the DoD's Ethical AI Principles</u>
- Video: <u>Bias in AI: Impact, Challenges, and Opportunities</u>