

# Implementing Responsible, Ethical, and Human-Centered AI

Alex Van Deusen and Carol J. Smith  
AI Division, CMU SEI  
arvandeusen@sei.cmu.edu and cjsmith@sei.cmu.edu

Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213

# How can teams harness the **power** of AI systems and design them to be **valuable** to humans?

# Advancing human-centered AI

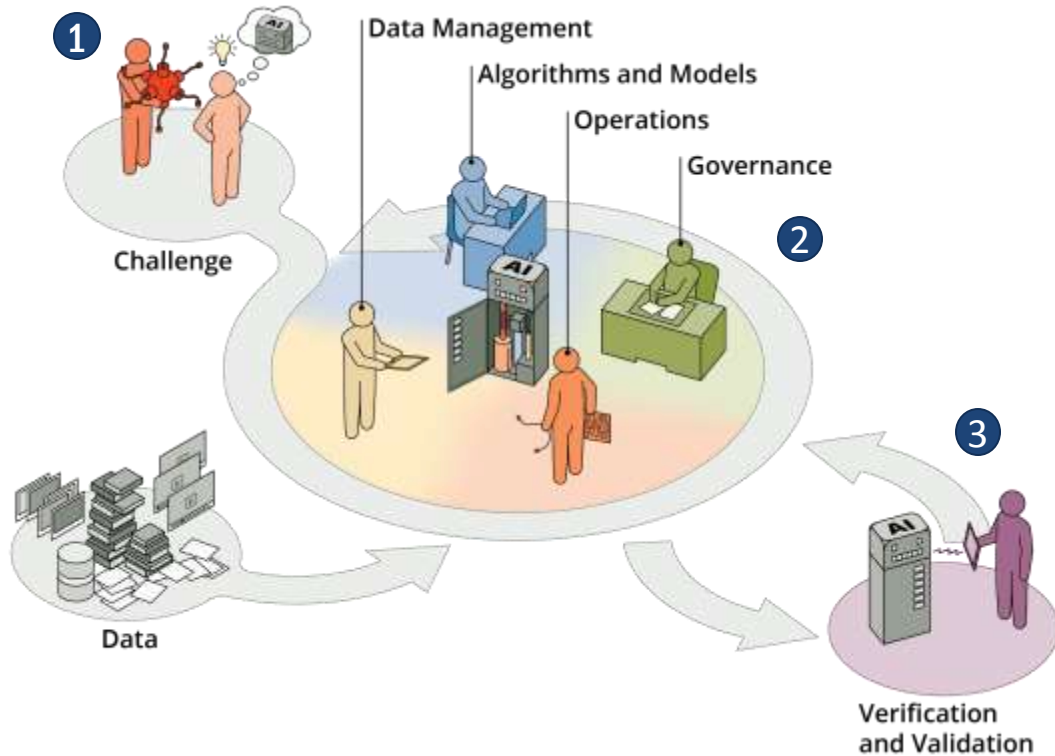
Effective implementations minimize unintended consequences

- Understand complexity of context
- Design for human-machine teaming
- Engage in critical oversight



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

# AI Systems and Human-Machine Teaming



- 1. Understand complexity of context:** focus on user needs
- 2. Design for human-machine teaming:** align on technical ethics
- 3. Engage in critical oversight:** Make AI interpretable, understandable, verifiable

Sensing changes over time

# Understand Complexity of Context

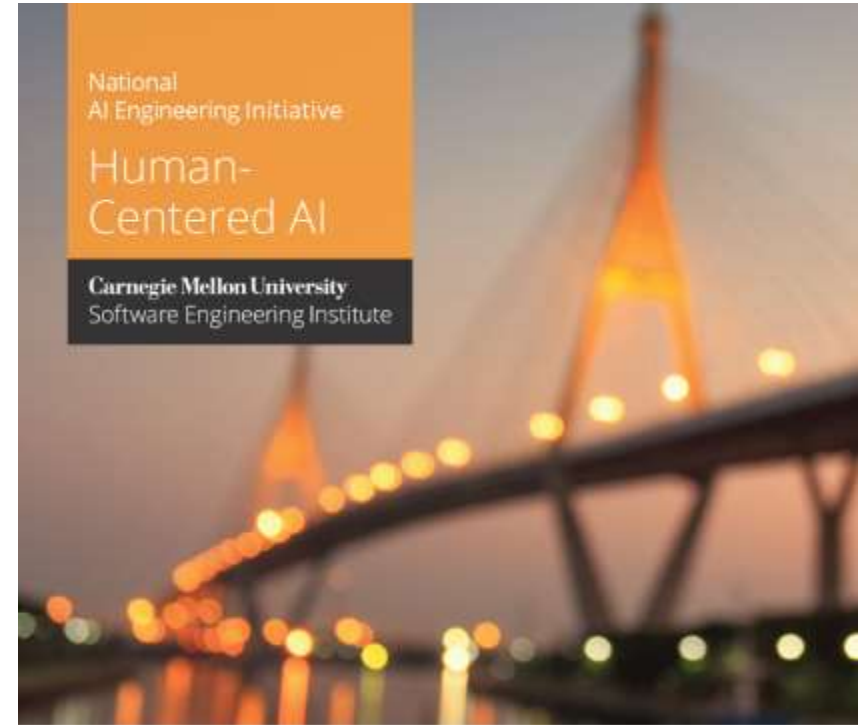


# Understanding context and sensing changes over time

- What is the desired system outcome?
- Human and contextual factors affect outcome

Does the human-machine team:

- Learn when shifts in context have occurred?
- Maintain clarity around operational intent?
- Adapt and evolve based on dynamic contexts and user needs?



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

# Complexity of context

- Environmental
- Human
- AI system capabilities
- Information



# Define the desired outcome

- What problem are you solving?
- For whom?
- What will help them?
- What kind of improvements are expected?
- What might a machine do better or faster?
- What happens if the AI system is *not* used?
- What is *not* going to be improved (out of scope)?

# Example: Semi-autonomous vehicles

## Adaptability to change

- Road conditions
- Weather
- Desired route
- Street painting
- Emergent situations



# Understanding context of use

- **Conduct human-centered research with primary users prior to development efforts**
  - Identify humans' actual objectives
  - Sense how objectives evolve, and continuously improve approximations of desired outcomes
- **Integrate insights from human-centered research and translate into**
  - Design and operations of system
  - AI System Reward functions
- **Research needed on opportunities and mechanisms to collect information and sense or infer user intention from other parts of larger system**

Developing tools, processes, and practices

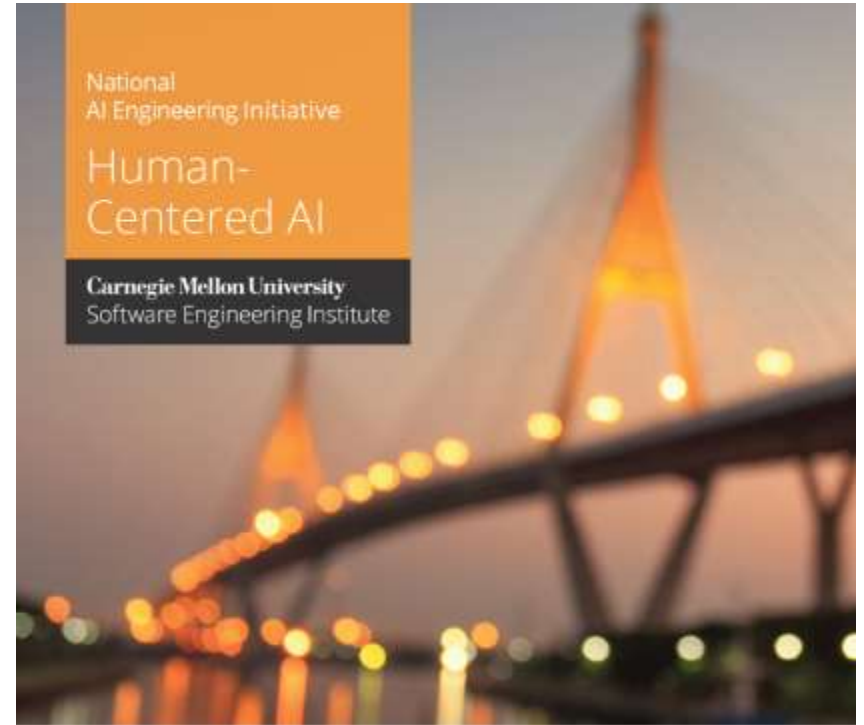
# Designing for Human-Machine Teaming



# Facilitating interdependence

## Human and machine complimenting each other

- Requires primary user to interact with and understand system
- Gaining *calibrated* level of trust
- AI system designed provide transparency regarding its limitations



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

# Trust is personal

Calibrated 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.

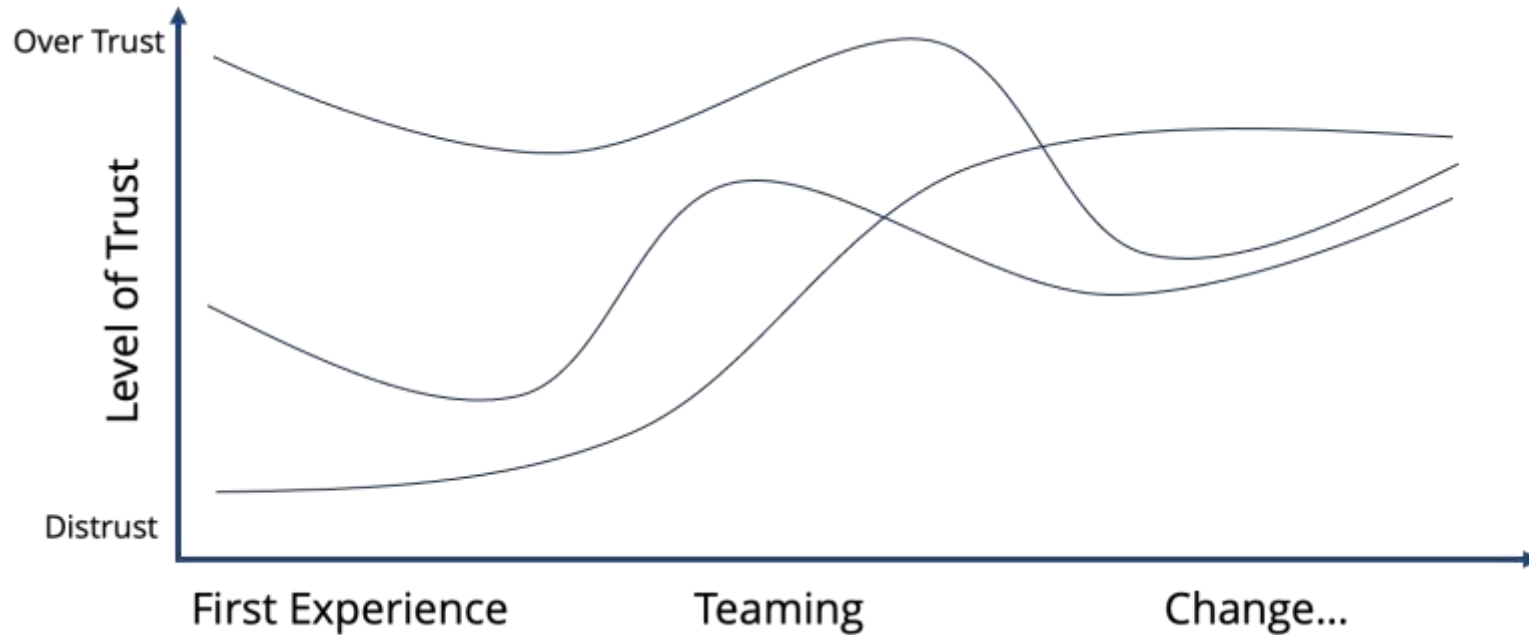


Rejection

Automation Bias

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 changes over time



Kun Yu, Shlomo Berkovsky, Ronnie Taib, Dan Conway, Jianlong Zhou, and Fang Chen. 2017. User Trust Dynamics: An Investigation Driven by Differences in System Performance. IUI 2017 (March 2017), 307-317. DOI: <http://dx.doi.org/10.1145/3025171.3025219>

# Activate curiosity

Incorporate user experience (UX) research and human-computer interaction (HCI) methods to activate curiosity.

- Speculate about misuse and abuse.
- Identify potential unintended/unwanted consequences.
- Potential severe abuse and consequences.
- Perspectives of people in frequently marginalized groups.



# Share understanding

## Converse on difficult topics:

- What do we value?
- Who could be hurt?
- What lines won't our AI cross?
- How are we shifting power?\*
- How will we track our progress?
- Perspective of frequently marginalized groups

\*"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>

Photo by Pam Sharpe [https://unsplash.com/@msgrace?utm\\_source=unsplash&utm\\_medium=referral&utm\\_content=creditCopyText](https://unsplash.com/@msgrace?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText) On Unsplash - [https://unsplash.com/s/photos/business-woman-smiling?utm\\_source=unsplash&utm\\_medium=referral&utm\\_content=creditCopyText](https://unsplash.com/s/photos/business-woman-smiling?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText)



# Integrate technical ethics

- Pick a set of ethics that works well for you and your team, adjust as necessary
- Pair with checklists and other RAI tools
  - Reduce risk and unwanted bias
  - Support inspection and mitigation planning

Checklist and Agreement - Downloadable PDF at SEI:  
<https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=636620>



Carnegie Mellon University  
Software Engineering Institute

## Designing Ethical AI Experiences: Checklist and Agreement

**USE THIS DOCUMENT TO GUIDE THE DEVELOPMENT OF AN AI SYSTEM THAT IS FAIR, RESPECTFUL, ACCURATE, TRUSTWORTHY, AND USEFUL. ARTIFICIAL INTELLIGENCE (AI) SYSTEMS WITH A DIVERSE TEAM ALIGNED ON SHARED ETHICS. AN INITIAL VERSION OF THIS DOCUMENT WAS PRESENTED WITH THE PAPER "DESIGNING TRUSTWORTHY AI: A HUMAN-CENTERED TRAINING FRAMEWORK TO GUIDE DEVELOPMENT" BY CAROL SMITH, AVAILABLE AT <https://arxiv.org/abs/1910.03016>.**

<p><b>We will design our AI system with the following in mind:</b></p> <ul style="list-style-type: none"><li>□ Designated humans have the ultimate responsibility for all decisions and outcomes.<ul style="list-style-type: none"><li>• Responsibilities are explicitly defined between the AI system and humans, and how they are shared.</li><li>• Human responsibility will be provided for those decisions that affect a person's life, quality of life, health, or reputation.</li><li>• Humans are always able to monitor, control, and deactivate systems.</li></ul></li><li>□ Significant decisions made by the AI system will be:<ul style="list-style-type: none"><li>• explained.</li><li>• able to be overridden.</li><li>• appealable and reversible.</li></ul></li></ul>	<p><b>We will create plans for the misuse of the AI system, including the following:</b></p> <ul style="list-style-type: none"><li>□ communication plans to share partners information with affected areas.</li><li>□ mitigation plans for managing the identified speculative risks.</li></ul> <p><b>We value respect and security of humanity, ethics, equity, privacy, accessibility, diversity, and inclusion.</b></p> <ul style="list-style-type: none"><li>□ respecting privacy and data rights (Only necessary data will be collected).</li><li>□ providing understandable security methods.</li><li>□ making the AI system robust, valid, and reliable.</li></ul>	<p><b>We make transparency with the goal of engendering trust:</b></p> <ul style="list-style-type: none"><li>□ The purpose, limitations, and biases of the AI system are explained in plain language.</li><li>□ Data sources have unambiguous (mapped) sources, and biases are known and explicitly stated.</li><li>□ Algorithms and models are open source and verifiable.</li><li>□ Certificates and cookies are presented for humans to make decisions on.</li><li>□ Transparency justification for recommendations and outcomes is provided.</li><li>□ Strong feedback and measurable monitoring systems are provided.</li></ul> <p><b>We value honesty and usability:</b></p> <ul style="list-style-type: none"><li>□ Humans can easily discern when they are interacting with an AI system, a human.</li><li>□ Humans can easily discern when and why the AI system is taking action or making decisions.</li><li>□ Improvements will be made regularly to meet human needs and technical standards.</li></ul>
---	---	--

**Team Signatures and Date**

---

**About the SEI**  
The Software Engineering Institute is a leader in the development of software engineering research, education, and professional development. For more information, visit [www.sei.cmu.edu](http://www.sei.cmu.edu).  
© 2022 Carnegie Mellon University. All rights reserved. This document is the property of Carnegie Mellon University. All rights reserved.

**Contact Us**  
Carnegie Mellon University  
Software Engineering Institute  
4800 Forbes Avenue, Pittsburgh, PA 15288-0151  
Phone: (412) 263-1000  
Fax: (412) 263-1001  
Email: [sei@cmu.edu](mailto:sei@cmu.edu)

© 2022 Carnegie Mellon University. SEI-22-1108-01-1000-01

# Designing for Human-Machine Teaming

- Design AI systems to provide transparency regarding AI limitations - boundaries and unfamiliar scenarios.
- Humans will gain *calibrated* levels of trust.
- Speculate about misuse and abuses
- Prevent or plan to mitigate situation.

Methods, Mechanisms, and Mindsets

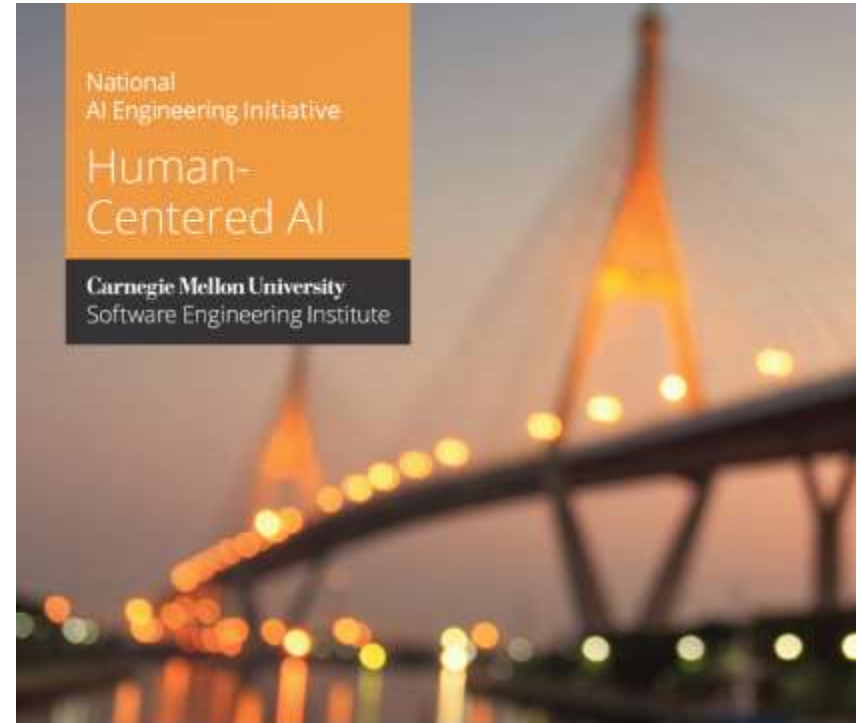
# Engage in Critical Oversight



# Engage in critical oversight

AI systems learn through data and observations

- Continuous human oversight
- Proactively identify risks of bias, misuse, abuse, and unintended consequences



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



What is a tomato?

Fruit?

Vegetable?

# Bias in Image Recognition

## Training data



## Data encountered



Use case courtesy of Dr. Eric Heim, CMU SEI <https://resources.sei.cmu.edu/library/author.cfm?authorid=542374>

# Only know what taught

## Training data



Unrepresentative  
or incomplete training data

## Data encountered



Unlikely to recognize

# All systems have some form of bias

Complete objectivity is misleading.

Unintended and purposeful bias

- Bias can have purpose
- Bias can be helpful

Reduce unintended/unwanted and/or harmful bias.

**“Data is a function of our history...  
The past dwells within our algorithms...  
Showing us the inequalities that have always been  
there.”**

**Dr. Joy Buolamwini  
Algorithmic Justice League  
Movie: Coded Bias on Netflix**

Photo: Joy Buolamwini on The Open Mind: Algorithmic Justice.  
Jan 12, 2019. <https://www.youtube.com/watch?v=hwHnXd0SSFY>

THE  
**OPEN MIND**



# Transparency and accountability

Understand inherent bias and amount of variance in the data:

- Creator's motivation
- Collection process
- Data included, and excluded
- Recommended uses, etc.

Provide evidence as appropriate – use context and knowledge of people using system.

# Plan for long term implementation and oversight

- Plan for AI system training and management of new data.
- Conduct continuous monitoring, evaluation, and audits for bias, brittleness, or potential distribution shift.



Nacho Kamenov & Humans in the Loop / Better Images of AI /  
A trainer instructing a data annotator on how to label images / CC-BY 4.0

# Engage in Critical Oversight

## Examine dynamic data and evaluate dynamic outcomes

- Is this the right data? What has changed?
- Is there enough evidence for calibrated trust?
- Did the system respond appropriately given the situation?
- Is the AI an effective collaborator?

We must work to define standard methods and processes for evaluating system outcomes

Implementing Responsible, Ethical, and Human-Centered AI

# Moving Forward

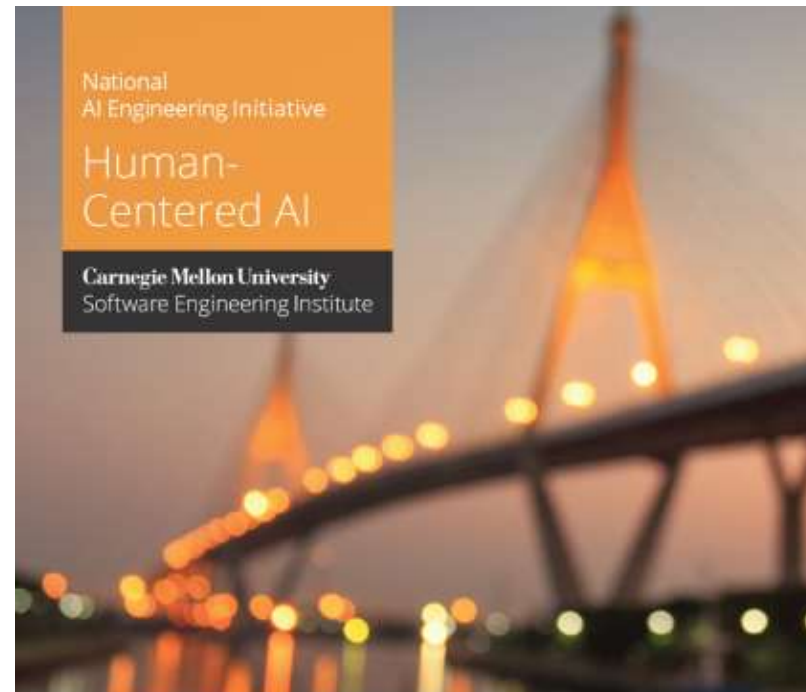


# How can teams harness the **power** of AI systems and design them to be **valuable** to humans?

# Advancing human-centered AI

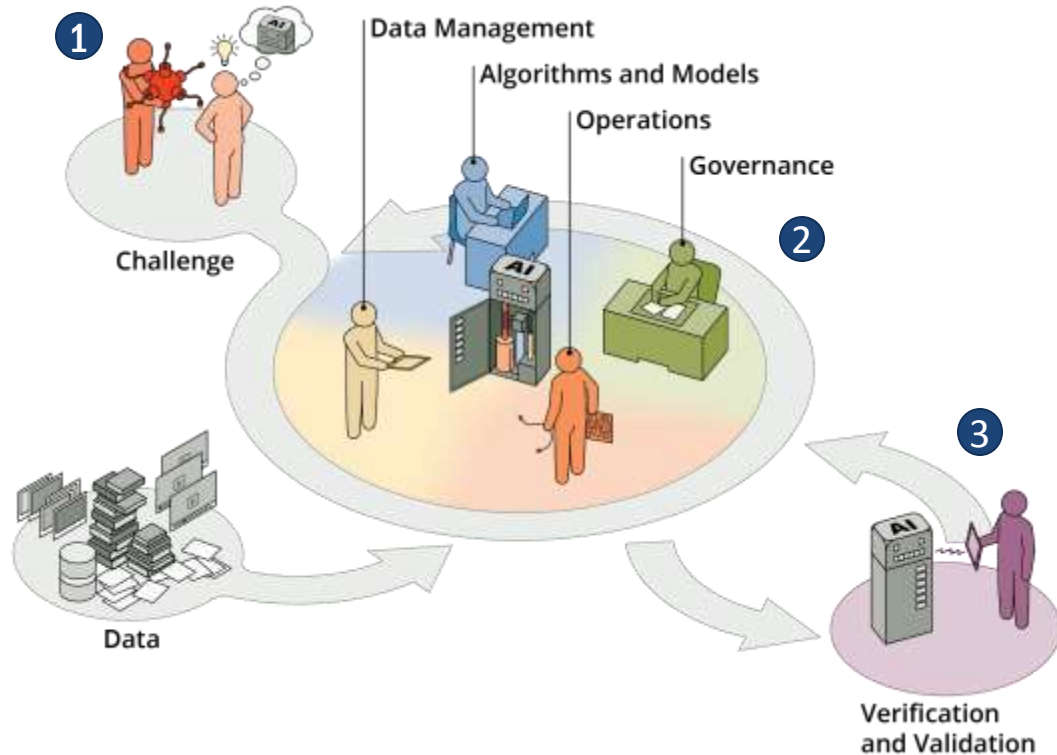
Effective implementations minimize unintended consequences

- Understand complexity of context
- Design for human-machine teaming
- Engage in critical oversight



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

# AI Systems and Human-Machine Teaming



- 1. Understand complexity of context:** focus on user needs
- 2. Design for human-machine teaming:** align on technical ethics
- 3. Engage in critical oversight:** Make AI interpretable, understandable, verifiable

**Encourage deep conversations, speculation, imaginative thinking, and active curiosity.**

# Thank You

**Alex Van Deusen**  
Design Researcher, CMU SEI  
arvandeusen@sei.cmu.edu

**Carol J. Smith**  
Sr. Research Scientist, CMU SEI  
cjsmith@sei.cmu.edu

Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213