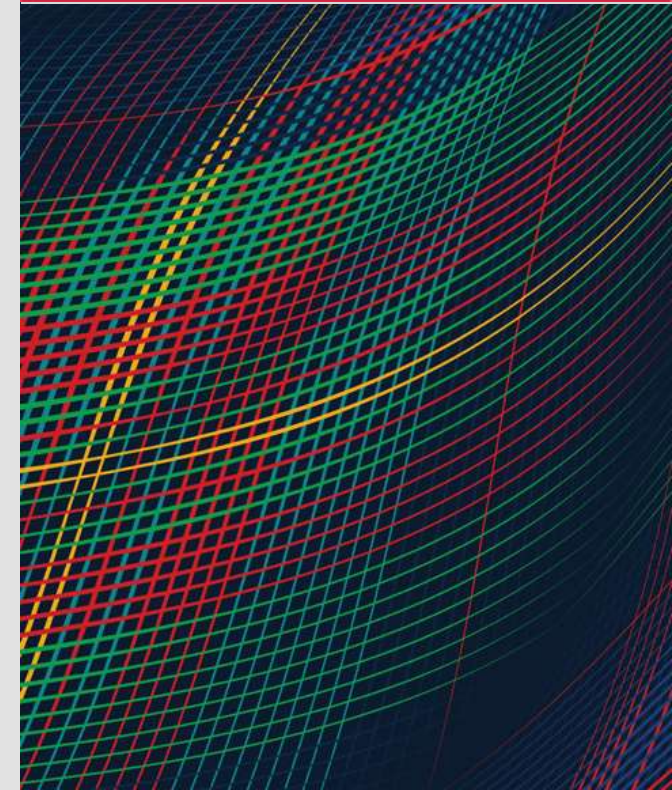


# Ethics and Trust for Emerging Technologies

NTIR 2023 | UALBANY CEHC

**Carol J. Smith**

Sr. Research Scientist, Human-Machine Interaction, AI Division  
Adjunct Instructor, Human-Computer Interaction Institute



# Copyright Statement

Copyright 2023 Carnegie Mellon University.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

References herein to any specific commercial product, process, or service by trade name, trade mark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by Carnegie Mellon University or its Software Engineering Institute.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at [permission@sei.cmu.edu](mailto:permission@sei.cmu.edu).

Carnegie Mellon® is registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.

DM23-0428

# About ACM



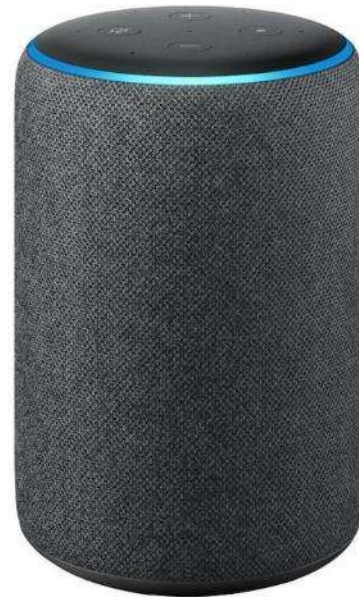
- ACM, the Association for Computing Machinery ([www.acm.org](http://www.acm.org)), is the premier global community of computing professionals and students with nearly 100,000 members in more than 170 countries interacting with more than 2 million computing professionals worldwide.
- **OUR MISSION:**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



For additional information, please visit <http://speakers.acm.org>

# Emerging technology



# Great potential - develop with caution

## Ring security camera hacks see homeowners subjected to racial abuse, ransom demands

*A spate of incidents has seen homeowners in four states fall victim to hackers.*

By Mark Hanrahan

December 12, 2019, 9:56 PM • 7 min read



### Ring camera systems being hacked

*Multiple U.S. families have reported incidents of Ring camera systems being hacked in recent days.*

The New York Times

## Thermostats, Locks and Lights: Digital Tools of Domestic Abuse



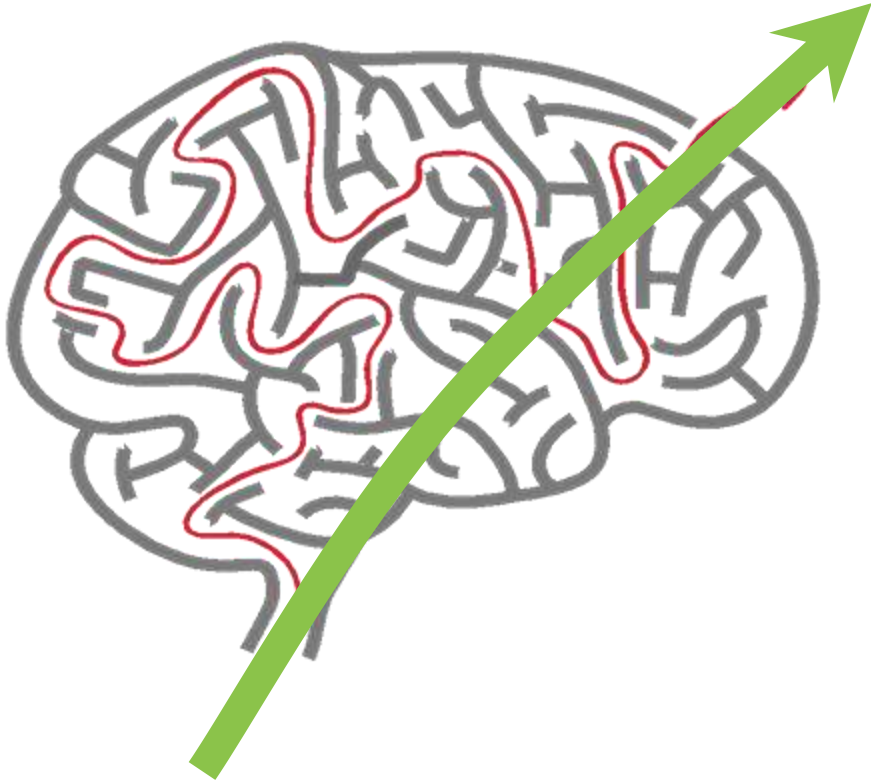
# Responsible, Intentional Design

# Ethics

- Based on well-founded standards of right and wrong
- Standard of expected behavior that guides the correct course of action
- What impact does my work have?

What is Ethics? By Manuel Velasquez, Claire Andre, Thomas Shanks, S.J., and Michael J. Meyer. Markkula Center for Applied Ethics  
<https://www.scu.edu/ethics/ethics-resources/ethical-decision-making/what-is-ethics/>

# To be biased, is to be human



- Bias are shortcuts, to avoid risk and simplify problems.
- Not inherently bad, may be misapplied
- Implicit = invisible
- Not necessarily in sync with our conscious beliefs
- **Can be managed and changed**

# All systems have some form of bias

- Data is collected/curated, by humans, for a purpose.
- Complete objectivity is misleading.
- Bias can have purpose and can be helpful.
- Our Goal: **Reduce unintended and/or harmful bias.**

# Image Recognition

## Train set



## 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

## Train set



Unrepresentative or incomplete training data

## Data encountered



Unlikely to recognize

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

Coded Gaze - 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=hwHnXdoSSFY>

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

## High value in diverse teams

### Diverse teams

- focus more on facts
- process facts more carefully
- are more innovative

“...become more aware  
of their own potential biases”



Photo by Christina @ wocintechchat.com on Unsplash  
[https://unsplash.com/@wocintechchat?utm\\_source=unsplash&utm\\_medium=referral&utm\\_content=creditCopyText](https://unsplash.com/@wocintechchat?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText)

David Rock, Heidi Grant. 2019. Why Diverse Teams Are Smarter. *Harvard Business Review*. November 4, 2019. <https://hbr.org/2016/11/why-diverse-teams-are-smarter>

## Support collaborative work

- Identify biases within ourselves and our teams
- Introduce tools and processes to control for bias
- Be aware of biases emerging in our work

## Early, purposeful work

In addition to understanding people and context

- Are we adding unnecessary additional risk to the situation?
- How will the system partner with people?
- How might these systems be misused/abused?

# Does trust really matter?







# Trust

Trust is complex, transient, and personal.

Trust is a psychological state

- Use evidence to determine risk.
- Gained with enough confidence in positive outcomes, to give control of something significant, to the system.

# Trust is contextual

Calibrated based on personal experiences, current context, and the available evidence of the 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](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>

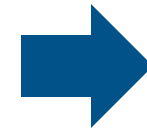
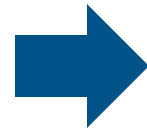
# Trust depends on psychological safety



UX Framework

# Implementing Ethics

# How do we make tech human-centered?



Human-Centered  
Technology

User Experience Honeycomb  
Peter Morville, et al.

# Design to work with, and for, people

User Experience/Human-Computer Interaction research to understand:

- Complexity of current context  
(environmental, human, and via information)
- System requirements
- How system will need to respond to change
- Overall changes, over time and experience



## Speculation keeps people safe - Activate Curiosity

# Make Safe Experiences

Actions to get into or maintain a **safe state** should be **easy** to do.

Actions that can lead to an **unsafe state** (hazard) should be **hard** to do.

Be speculative

- Don't assume only average cases
- Can't verify probabilities



N. G. Leveson. 2017. The Therac-25: 30 Years Later. In Computer, vol. 50, no. 11, (November 2017), 8-11. DOI: 10.1109/MC.2017.4041349  
N. Leveson. 1995. Safeware: System Safety and Computers, Addison Wesley (1995).

# 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

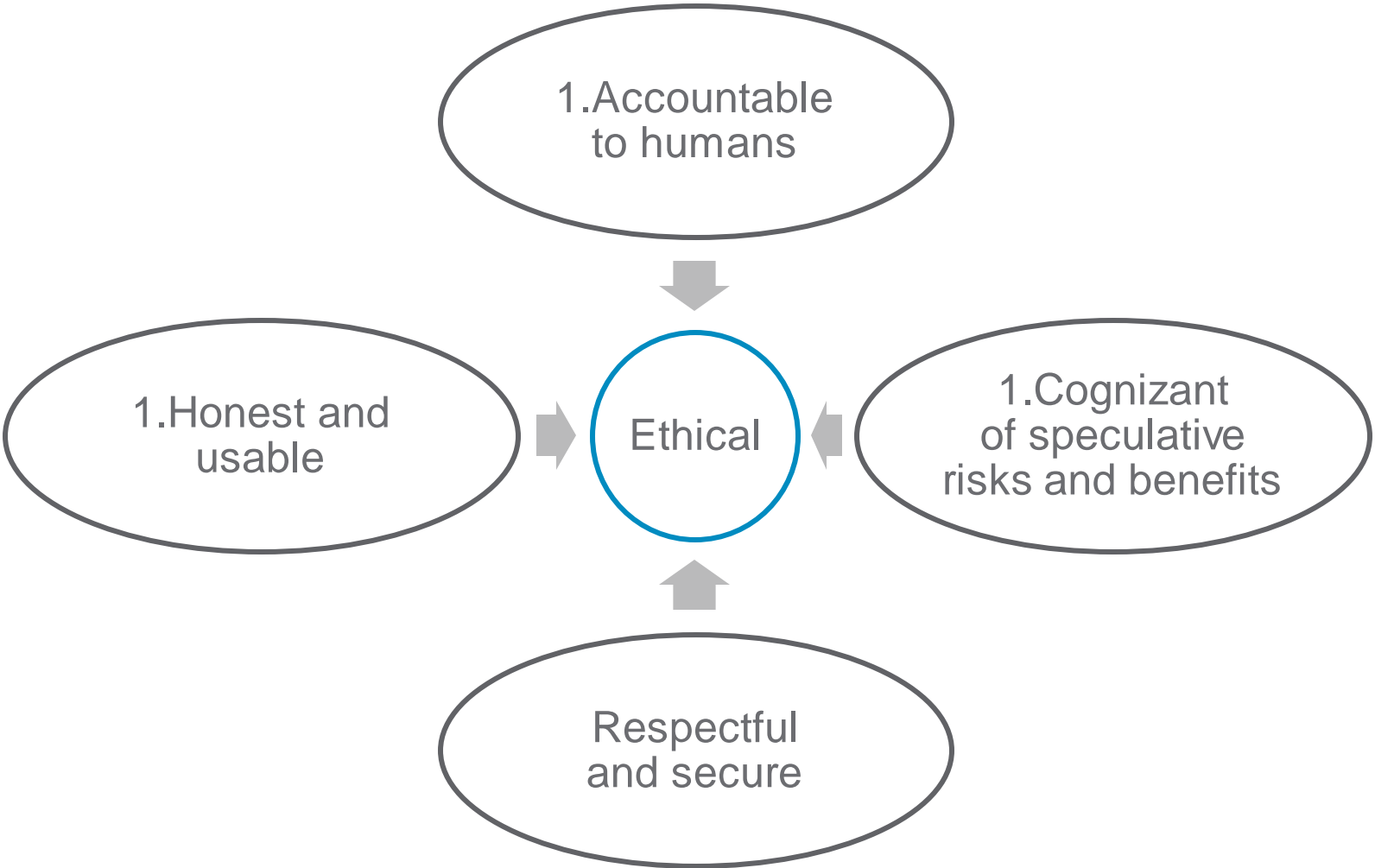
Amanda Muller and Carol Smith. 2022. Perceptions of Function Allocation between Humans and AI-Enabled Systems. UXPA 2022 (pre-print).  
<https://uxpa2022.org/sessions/perceptions-of-function-allocation-between-humans-and-ai-enabled-systems/>

# Adopt Technology Ethics

- Harmonize cultural variations
- Balance to pace of change, industry pressure
- Explicit permission to consider and question breadth of implications



# Prompt conversations – UX Framework



Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020. [https://insights.sei.cmu.edu/sei\\_blog/2020/03/designing-trustworthy-ai-for-human-machine-teaming.html](https://insights.sei.cmu.edu/sei_blog/2020/03/designing-trustworthy-ai-for-human-machine-teaming.html)

# Prompts for conversations

## Pair checklists with technical ethics

- Bridge gaps between “do no harm” and reality
- Reduce risk and unwanted bias
- Support inspection and mitigation planning



Carnegie Mellon University  
Software Engineering Institute

### Designing Ethical AI Experiences: Checklist and Agreement

**USE THIS DOCUMENT TO GUIDE THE DEVELOPMENT** of accountable, de-risked, respectful, secure, honest, and usable 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-Machine Teaming Framework to Guide Development* by Carol Smith, available at <https://andu.org/abs/1910.03515>.

<p><b>We will design our AI system with the following in mind:</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> 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 human(s), and how they are shared.</li> <li>• Human responsibility will be preserved for final 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><input type="checkbox"/> 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 work to speculatively identify the full range of risks and benefits:</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Harmful, malicious use and consequences, as well as good, beneficial use and consequences</li> <li><input type="checkbox"/> We will be cognizant and exhaustively research unintended consequences.</li> </ul> <p><b>We will create plans for the misuse/abuse of the AI system, including the following:</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> communication plans to share pertinent information with all affected people</li> <li><input type="checkbox"/> mitigation plans for managing the identified speculative risks.</li> </ul> <p><b>We value respect and security:</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> incorporating our values of humanity, ethics, equity, fairness, accessibility, diversity, and inclusion</li> <li><input type="checkbox"/> respecting privacy and data rights (Only necessary data will be collected.)</li> <li><input type="checkbox"/> providing understandable security methods</li> <li><input type="checkbox"/> making the AI system robust, valid, and reliable</li> </ul>	<p><b>We value transparency with the goal of engendering trust:</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> The purpose, limitations, and biases of the AI system are explained in plain language.</li> <li><input type="checkbox"/> Data sources have unambiguous, respected sources, and biases are known and explicitly stated.</li> <li><input type="checkbox"/> Algorithms and models are appropriate and verifiable.</li> <li><input type="checkbox"/> Confidence and consent are presented for humans to base decisions on.</li> <li><input type="checkbox"/> Transparent justification for recommendations and outcomes is provided.</li> <li><input type="checkbox"/> Straightforward and interpretable monitoring systems are provided.</li> </ul> <p><b>We value honesty and usability:</b></p> <ul style="list-style-type: none"> <li><input type="checkbox"/> Humans can easily discern when they are interacting with the AI system vs. a human.</li> <li><input type="checkbox"/> Humans can easily discern when and why the AI system is taking action and/or making decisions.</li> <li><input type="checkbox"/> 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 federally chartered non-profit organization under 501(c)(3) that works to advance and govern the use of software, hardware, and systems in support of the nation's defense and industry. The SEI is a national leader in providing state-of-the-art research, development, and education in software engineering.

**Contact Us**  
CARRIE-LEE WILSON  
SOFTWARE ENGINEERING INSTITUTE  
4800 WALKER DRIVE, PITTSBURGH, PA 15262  
412.263.1479  
seil@se.cmu.edu

©2023 Carnegie Mellon University | 0270 - 0107-2003 | 1.0 (03/2023)

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>

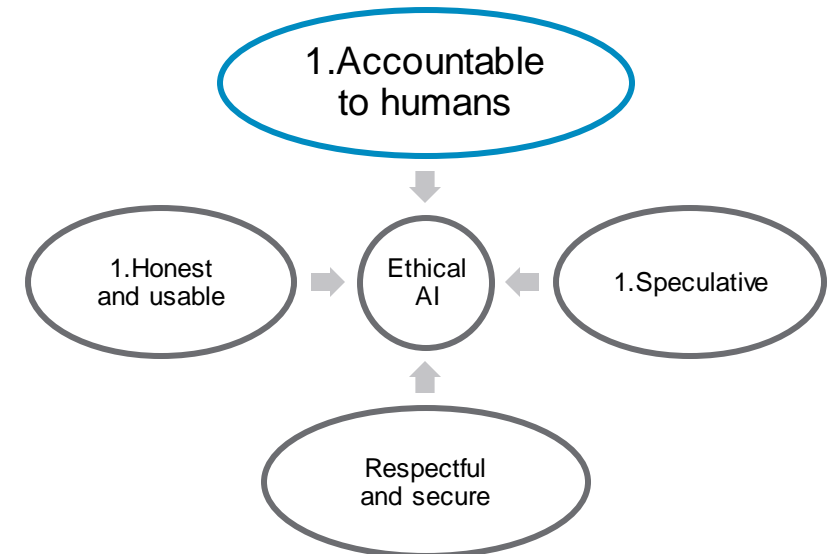
# Accountable to Humans

Ensure humans have ultimate control

- Able to monitor and control risk

Human responsibility for final decisions

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



# “Ensure humans can unplug the machines”

– Grady Booch



TED Talk, Grady Booch, Scientist, Philosopher, IBM'er

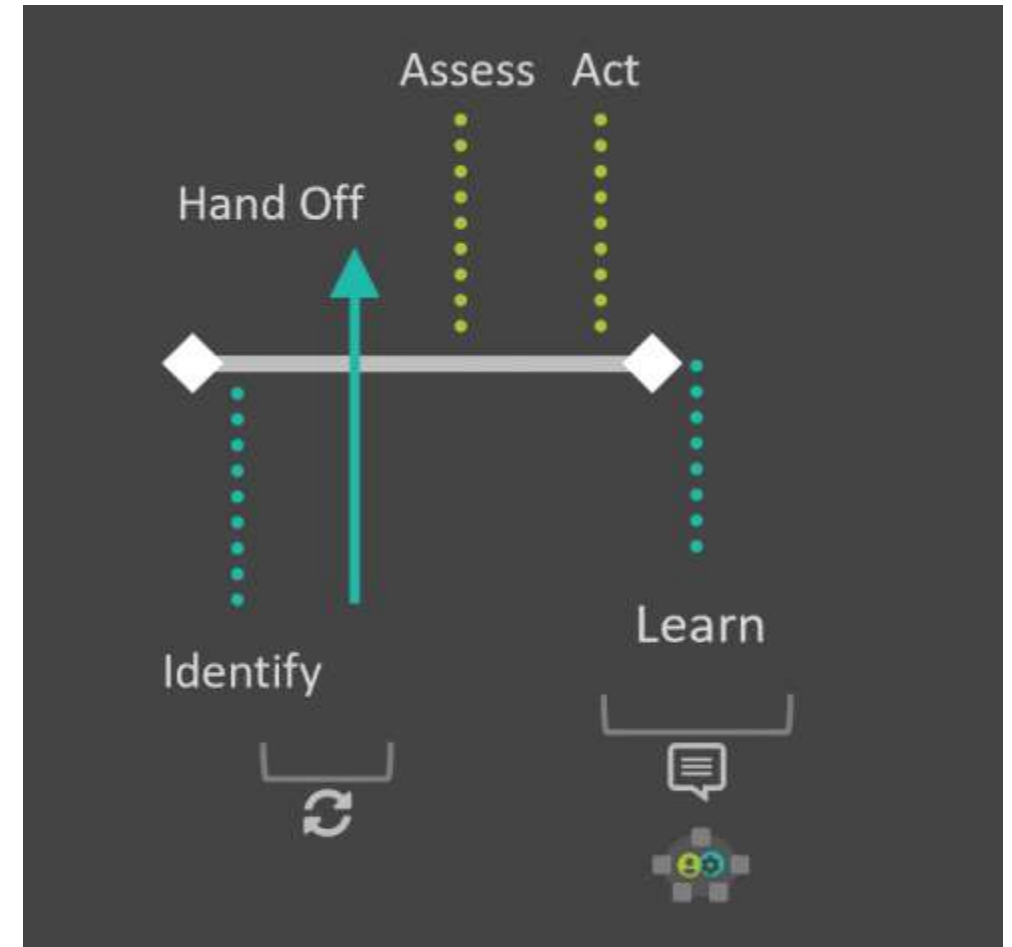
[https://www.ted.com/talks/grady\\_booch\\_don\\_t\\_fear\\_superintelligence](https://www.ted.com/talks/grady_booch_don_t_fear_superintelligence)

# Significant decisions

## Significant decisions made by system

- explained
- able to be overridden
- appealable and reversible

Responsibilities explicitly defined between people and systems.

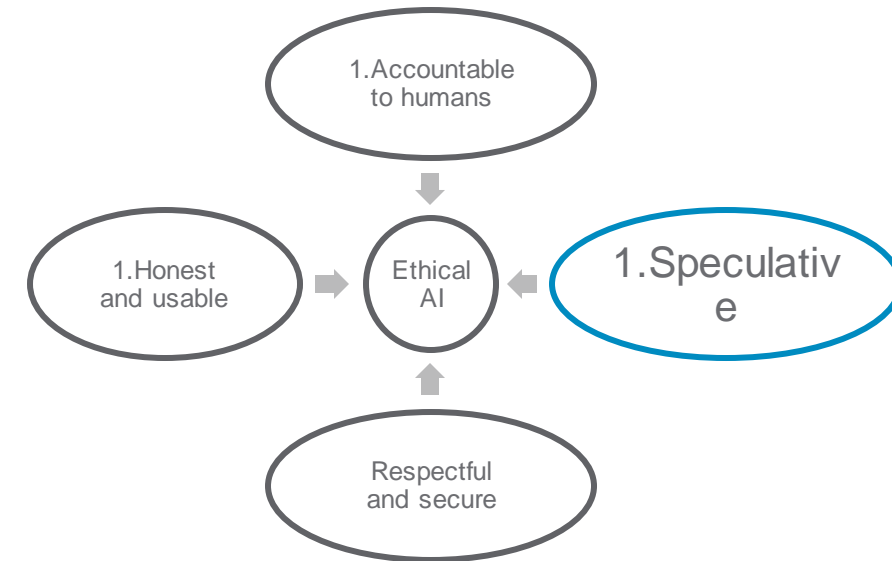


# Cognizant of Speculative Risks and Benefits

Identify full range of

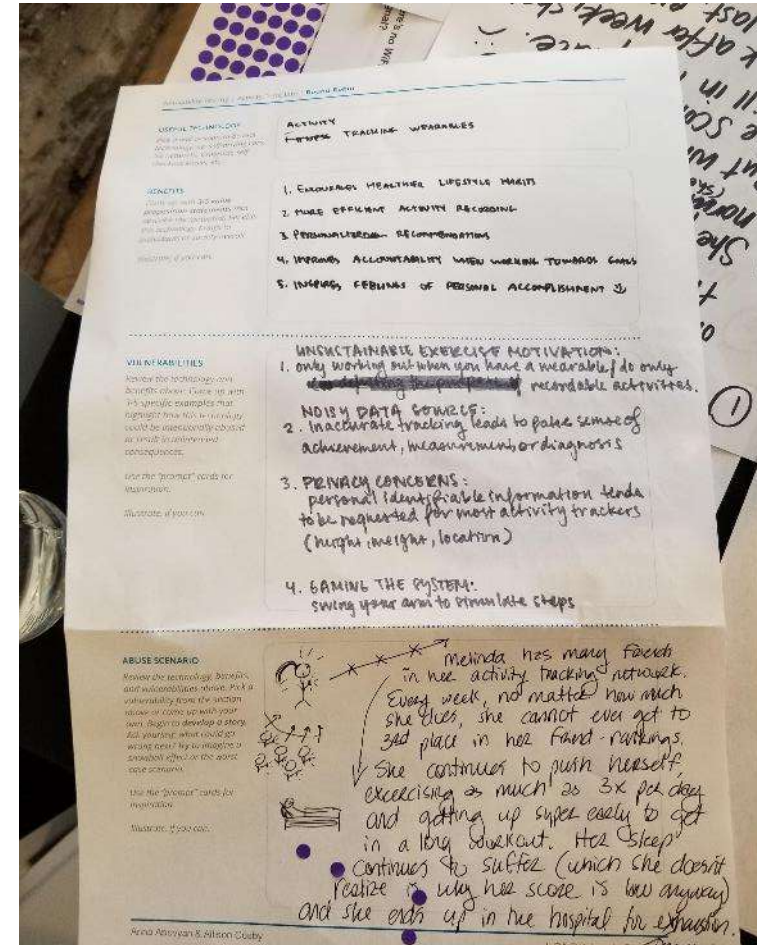
- Harmful, malicious use, as well as good, beneficial use.
- Unwanted/unintended consequences.

Prevent potential harms.



# Conduct UX research - activate curiosity

- Speculate about misuse and abuse – abusability testing
- Potentially severe consequences (even if rare)
- Perspective of people in frequently marginalized groups



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

# Create communication & mitigation plans

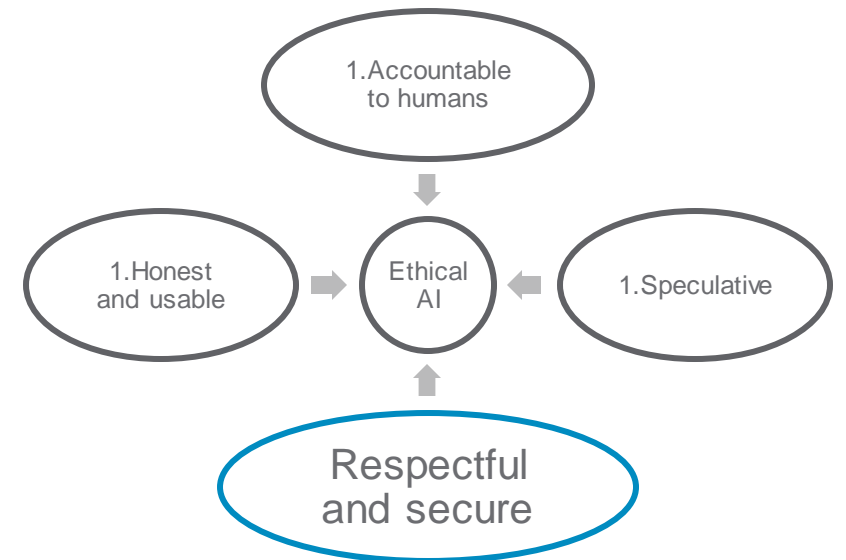
Plan for unwanted consequences.

## Misuse and abuse of system

- Who can report?
- To whom?
- Turn off?
- Who notified?
- Consequences?

# Respectful and Secure

- Values of humanity, ethics, equity, fairness, accessibility, diversity and inclusion.
- Respect privacy and data rights (only collect what is necessary).
- Make systems robust, valid, and reliable.
- Provide understandable security.



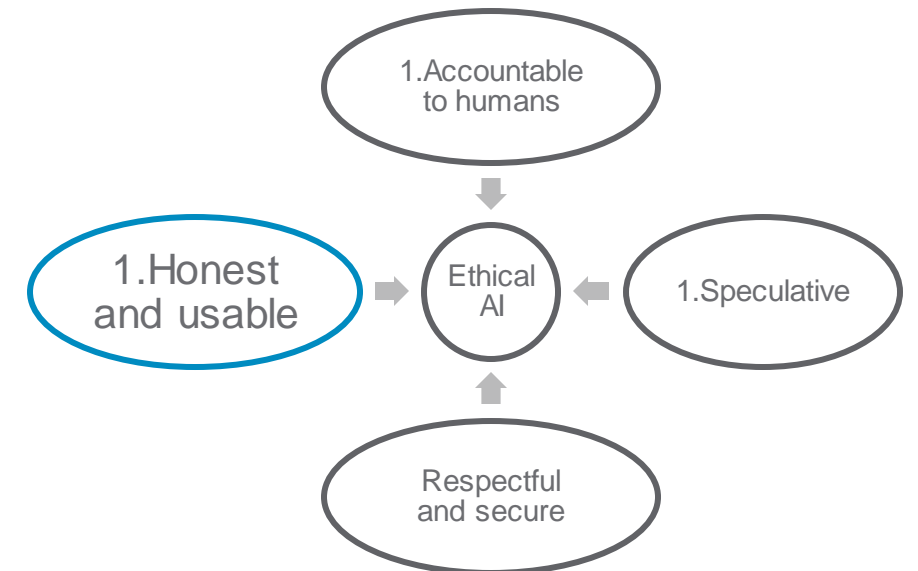
Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020.

## Honest and Usable

- Value transparency with the goal of engendering calibrated trust.
- Provide transparency regarding boundaries and unfamiliar scenarios.
- Explicitly state identity as an AI system.

## Fairness

- Show awareness of purposeful bias.
- Provide AI system limitations.
- Overcommunicate on issues.



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

# Honest?



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>

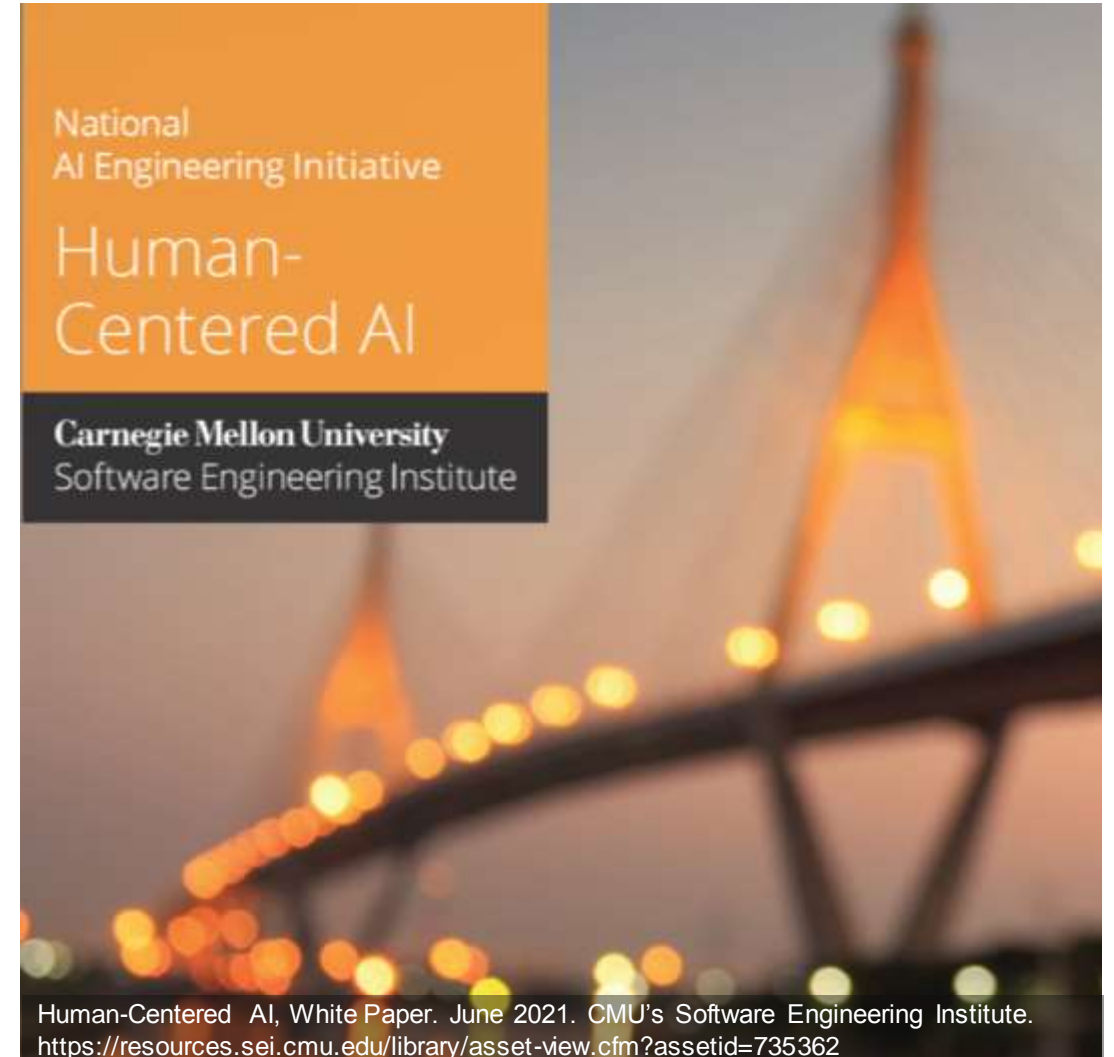
# Design for Calibrated Trust

- Provide transparency regarding system limitations
  - boundaries and unfamiliar scenarios
- Speculate about misuse and abuse
- Prevent or plan to mitigate situation

## Engage in critical oversight

“What are we doing?  
Why are we doing it,  
and for whom?”

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



# Establish psychological safety for diverse teams



# Change is constant



# 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?\*
- Perspective of frequently marginalized groups

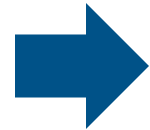
\*"Don't ask if artificial intelligence is good or fair, ask how it shifts power." Pratyusha Kalluri.

<https://www.nature.com/articles/d41586-020-02003-2>

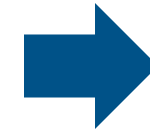
Photo by Pam Sharpe [https://unsplash.com/@msgrace?utm\\_source=unsplash&utm\\_medium=referral&utm\\_content=creditCopyText](https://unsplash.com/@msgrace?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText) On Unsplash - [https://unsplash.com/s/photos/business-woman-smiling?utm\\_source=unsplash&utm\\_medium=referral&utm\\_content=creditCopyText]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



# Design to work with, and for, people



1. Capitalize on human strengths
2. Encourage deep conversations
3. Be intentional, keep people safe



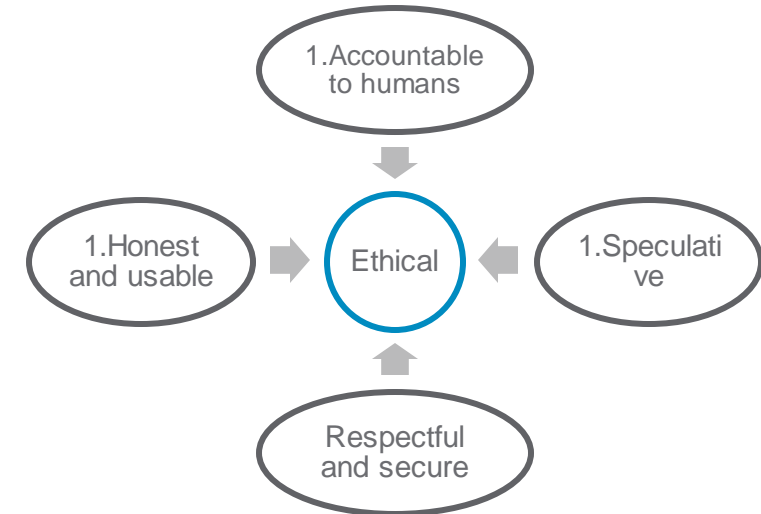
Human-Centered Technology

User Experience Honeycomb  
Peter Morville, et al.

New uncomfortable work

# Be uncomfortable and kind.

We aren't perfect, tech won't be perfect.



**Carol J. Smith**

LinkedIn: <https://www.linkedin.com/in/caroljsmith/>

**AI DIVISION  
SOFTWARE ENGINEERING INSTITUTE**

