

UXPA Seattle
January 5, 2022

Designing for Human-Centered AI

Carol J. Smith

Sr. Research Scientist, Human-Machine Interaction, CMU SEI
Adjunct Instructor, CMU Human-Computer Interaction Institute

Twitter: @carologic @SEI_CMU_AI

Software Engineering Institute
Carnegie Mellon University
Pittsburgh, PA 15213

Copyright Statement

Copyright 2022 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.

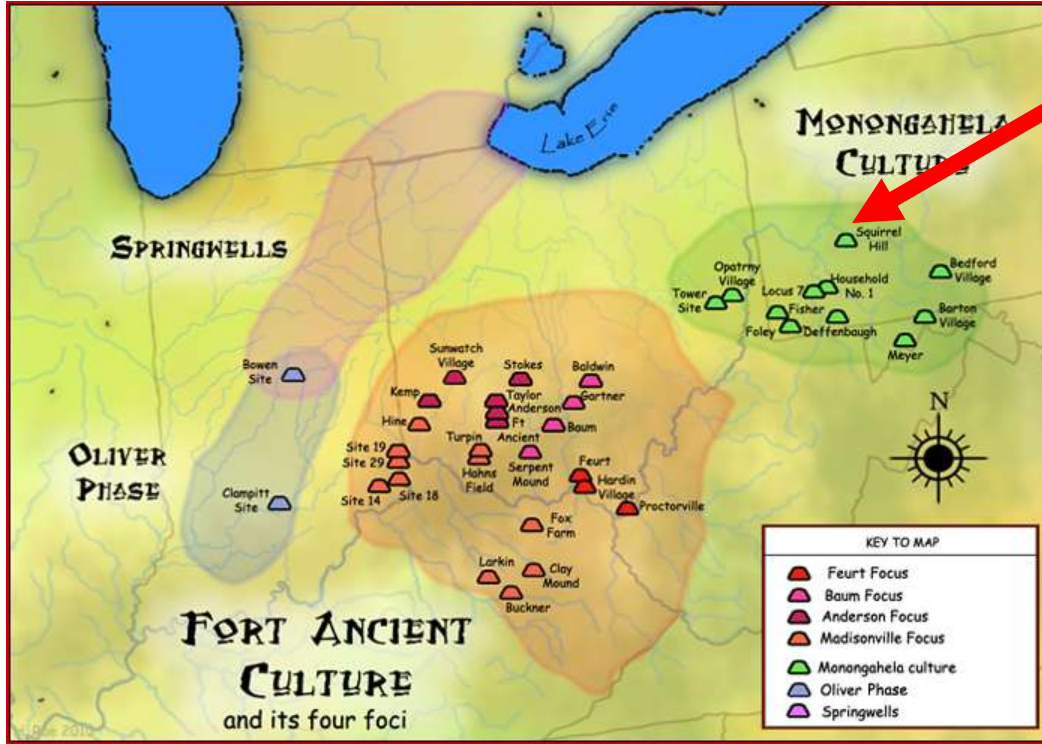
[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

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.

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

DM22-0014

Acknowledgement: The Land I Speak On



Land of Monongahela, Adena and Hopewell Nations;

Seneca, Lenape and Shawnee lands;

Osage, Delaware and Iroquois lands.

Now known as Pittsburgh, PA, USA.

Map by Herb Roe via Wikipedia https://en.wikipedia.org/wiki/Monongahela_culture

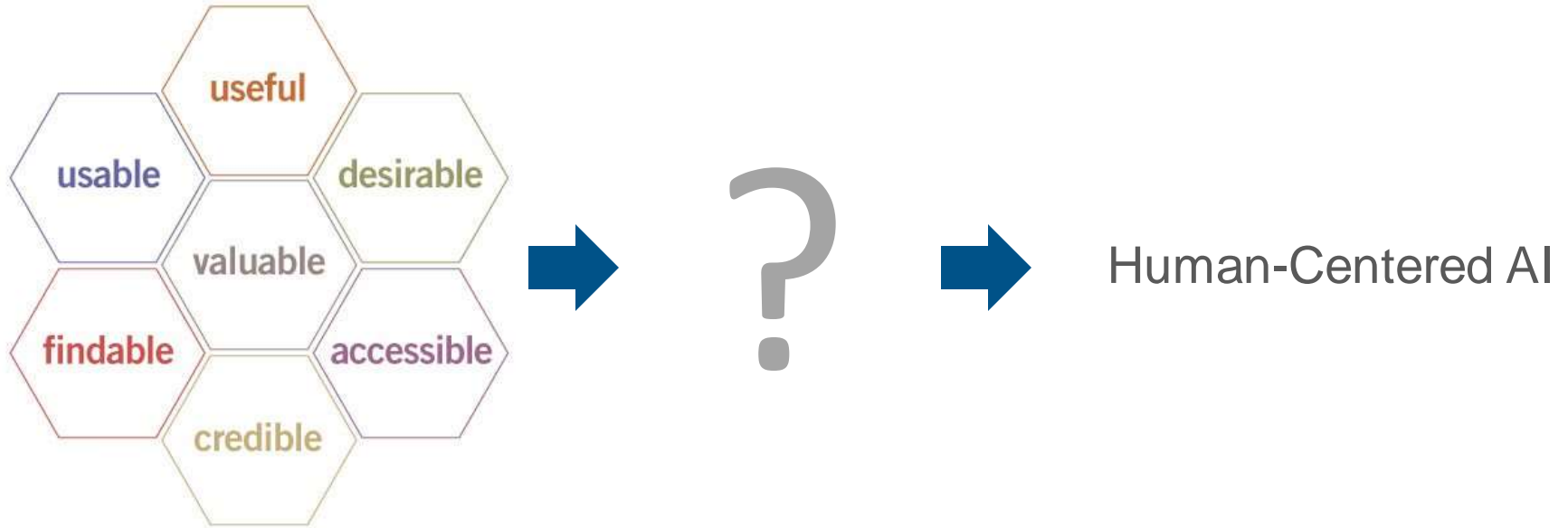
First Machines



Al-Jazari described a water-powered automaton orchestra on a boat in 1206



How do we Make AI, Human-Centered?



User Experience Honeycomb
Peter Morville, et al.

Broaden our Work

Is this an AI-friendly challenge?

What kind of improvements are expected?

What are the benefits and risks?

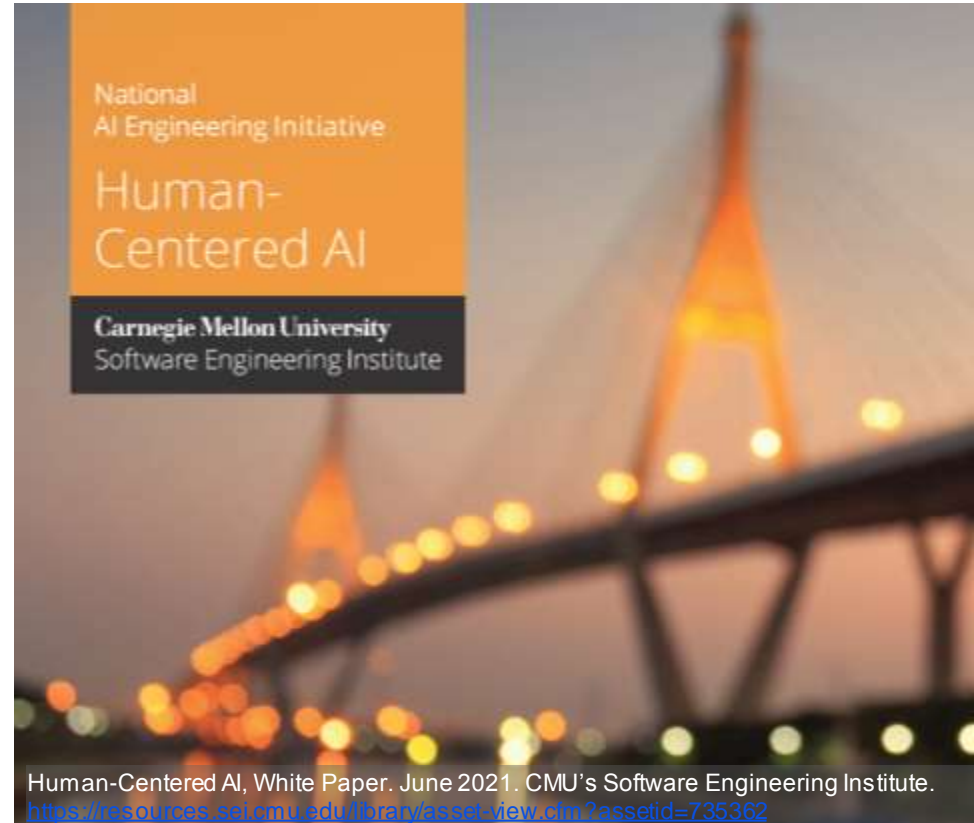
How will we know we've made improvements?

Design to work with, and for, people

Effective implementations

Minimize unintended consequences

1. Understand complexity of context
2. Design for human-machine teaming
3. Engage in critical oversight





Sensing changes over time

Understanding Complexity of Context

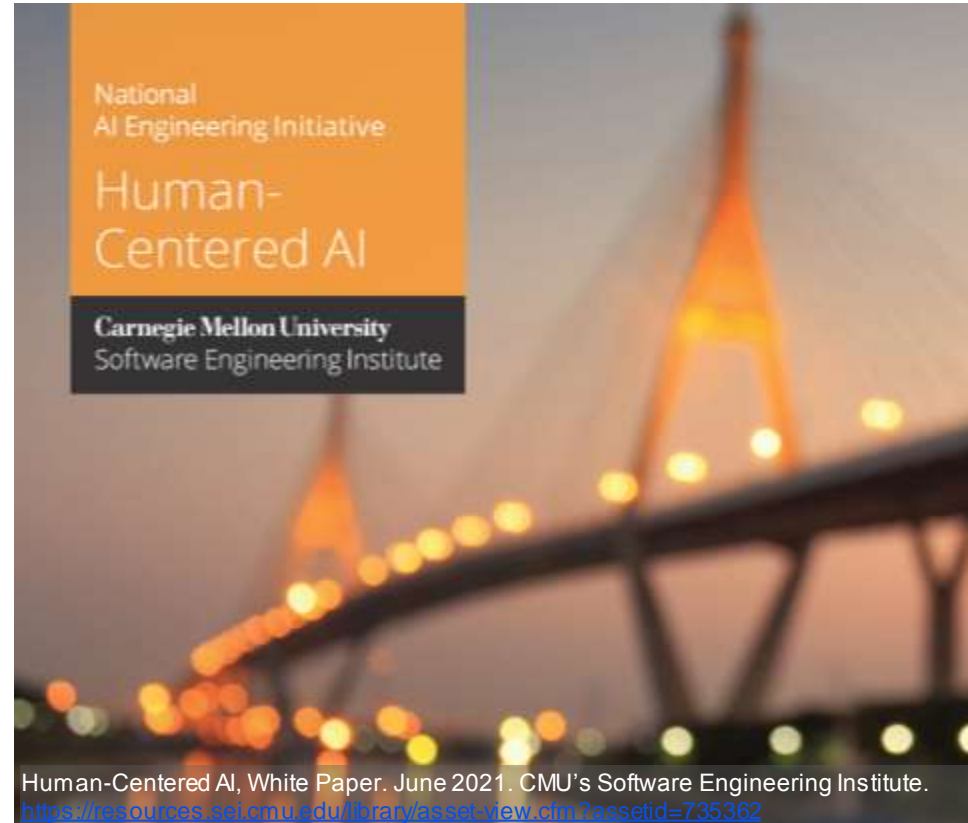
Understanding Complexity of Context

Desired outcome, human's needs

Human and contextual factors
affect outcome

Do human and AI:

- learn when shifts in context have occurred?
- maintain clarity around operational intent?
- adapt and evolve based on dynamic contexts?



Complexity

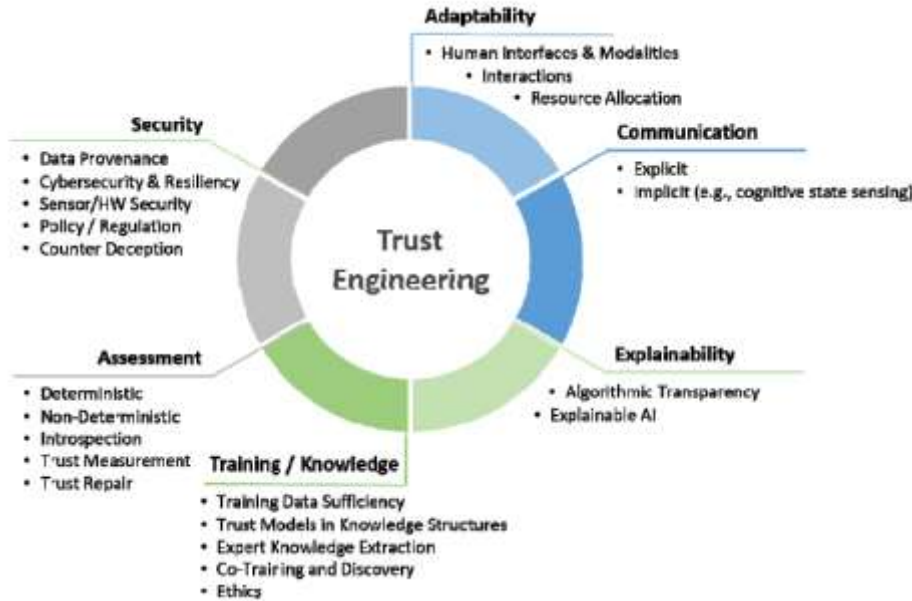
Environmental context

Human factors

Information



Collaborative Activities and Interactions - Design Components



- Security
- Adaptability
- Communication
- Explainability
- Training/Knowledge
- Assessment

Neta Ezer, Sylvain Bruni, Yang Cai, Sam J. Heppenstall, Christopher A. Miller, and Dylan D. Schmorow. 2019. Trust Engineering for Human-AI Teams. Proceedings of the Human Factors and Ergonomics Society Annual Meeting 63, no. 1 (November 2019): 322–26. <https://doi.org/10.1177/1071181319631264>.

What Changes Across Time Cycles?

Length of interactions

- Short and hectic
- Longer, cyclical

Iterative

Require clear communication,
negotiation,
and coordination.



How IAs Can Shape the Future of Human-AI Collaboration
Presented on April 28-30, 2021 at the Information Architecture Conference (IAC21)

Example: Semi-autonomous vehicles

Potential factors

- Driver behavior
- Weather changes
- Street painting changes
- Change in desired route
- Highway vs. city driving
- Emergent situations



Example: Decision making for medical treatment

Potential factors

- How much information is already known?
- Stage of disease?
- Specifics of the patient's health, stage of disease, family situation, insurance status, etc.?
- How is new information integrated and how does that change interactions?

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.

Don't rely on operators to detect errors and recover before an accident – it isn't realistic.



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).

Make Systems Effective Team Players

Easy to direct

- How observable is its behavior?
- How easily and efficiently allows itself to be directed?
- Even (or especially) during busy, novel episodes?

S. W. A. Dekker and D. D. Woods. 2002. MABA-MABA or Abracadabra? Progress on Human–Automation Co-ordination. *Cognition Tech Work* 4, (2002) 240–244. DOI: <https://doi.org/10.1007/s101110200022> Note: MABA-MABA (Men-Are-Better-At/Machines-Are-Better-At lists)

Capitalize on Human Strengths

Humans are better at:

- Perceiving patterns
- Improvising and using flexible procedures
- Recalling relevant facts at the appropriate time
- Reasoning inductively
- Exercising judgment

Mary Cummings. 2004. Automation Bias in Intelligent Time Critical Decision Support Systems. AIAA 2004-6313. AIAA 1st Intelligent Systems Technical Conference. (September 2004). DOI: <https://doi.org/10.2514/6.2004-6313>



Context

Human-centered research to

- Understand complexity (environmental, human and information)
- Changes over time

Inform and support designs that provide clear communication, negotiation, and coordination

Challenges

Demystifying AI to colleagues

Getting time and budget for UX research

Using existing use patterns wisely

Understanding what changes with regard to time cycles



Development of tools, processes, and practices

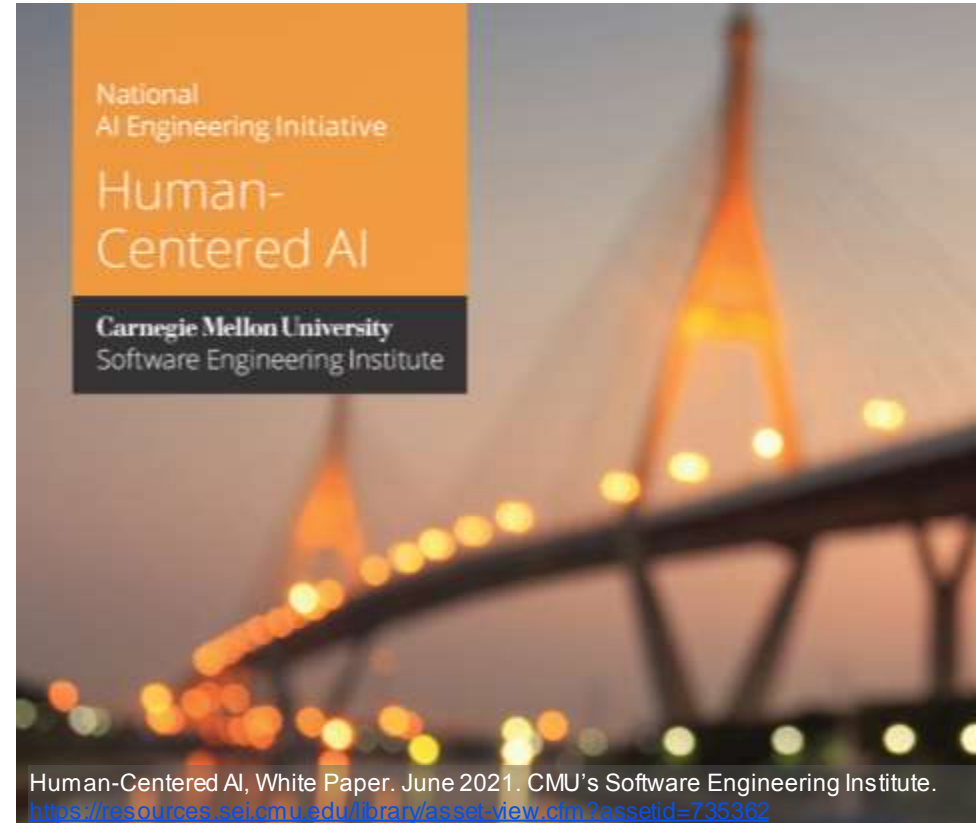
Design for Human-Machine Teaming

Design for HMT - Interdependence

People interacting with
and understanding systems

Gaining *appropriate* levels
of trust

Design AI system to provide
transparency regarding AI limitations



AI is NOT sentient or unknowable

No one should use (or buy, or maintain),
important systems that are described as indecipherable
People must have ability to control and monitor systems

Provide transparency regarding AI limitations
– boundaries and unfamiliar scenarios.

Provenance of data
– what is the system basing recommendations on?

Automation Bias

Propensity for humans to **favor suggestions** from automated decision-making systems and to **ignore contradictory information** made without automation, even if it is correct.

Mary Cummings. 2004. Automation Bias in Intelligent Time Critical Decision Support Systems. AIAA 2004-6313. AIAA 1st Intelligent Systems Technical Conference. (September 2004). DOI: <https://doi.org/10.2514/6.2004-6313>

Trust is a Continuum

Distrust

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

Calibrated Trust

Trust matches system capabilities leading to appropriate use.

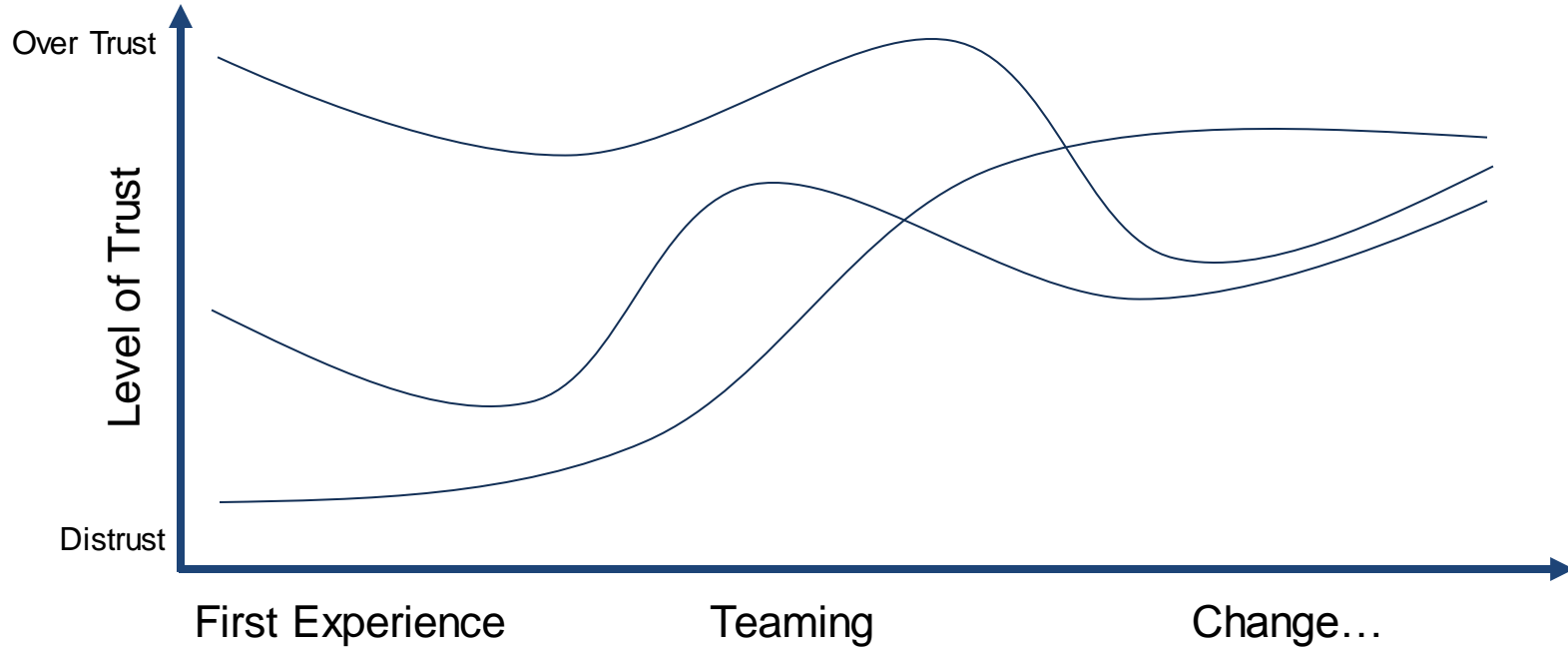
Over Trust

Trust exceeding system capabilities - may lead to misuse



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>

Skepticism and Speculation Keep People Safe

Lt Col Stanislav Yevgrafovich Petrov
Soviet Nuclear false alarm incident,
Sept 26, 1983



Activate Curiosity

UX research methods and activities to activate curiosity:

- Abusability Testing ([Dan Brown](#))
- “Black Mirror” Episodes ([Casey Fiesler](#))
(inspired by British dystopian sci-fi tv series of same name)

Speculate about system misuse and abuse

- What are potential unintended/unwanted consequences?

More methods to “Outsmart Your Own Biases.”: <https://hbr.org/2015/05/outsmart-your-own-biases>

Implicit Association Test (IAT): <https://implicit.harvard.edu/implicit/takeatest.html>

Reward team members for finding ethics bugs

Dr. Ayanna Howard

- on the Artificial Intelligence Podcast with Lex Fridman



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?*
- 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 On Unsplash - https://unsplash.com/s/photos/business-woman-smiling?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText



New uncomfortable work

“*Be uncomfortable*”

- Laura Kalbag

Ethical design is not superficial.

Adopt Technology Ethics

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



Prompt conversations

Pair Checklist with Technical Ethics

- Bridges gap 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 AN ETHICAL, DE-RISKED, RESPECTFUL, HARM-FREE, AND USEFUL ARTIFICIAL INTELLIGENCE (AI) SYSTEMS WITH A DIVERSE TEAM ALIGNED ON SHARED VALUES. An initial version of this document was presented with the paper *Designing Thoughtful AI: A Human-Machine Learning Framework to Guide Development* by Carol Smith, available at <https://arxiv.org/abs/1910.03016>.

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

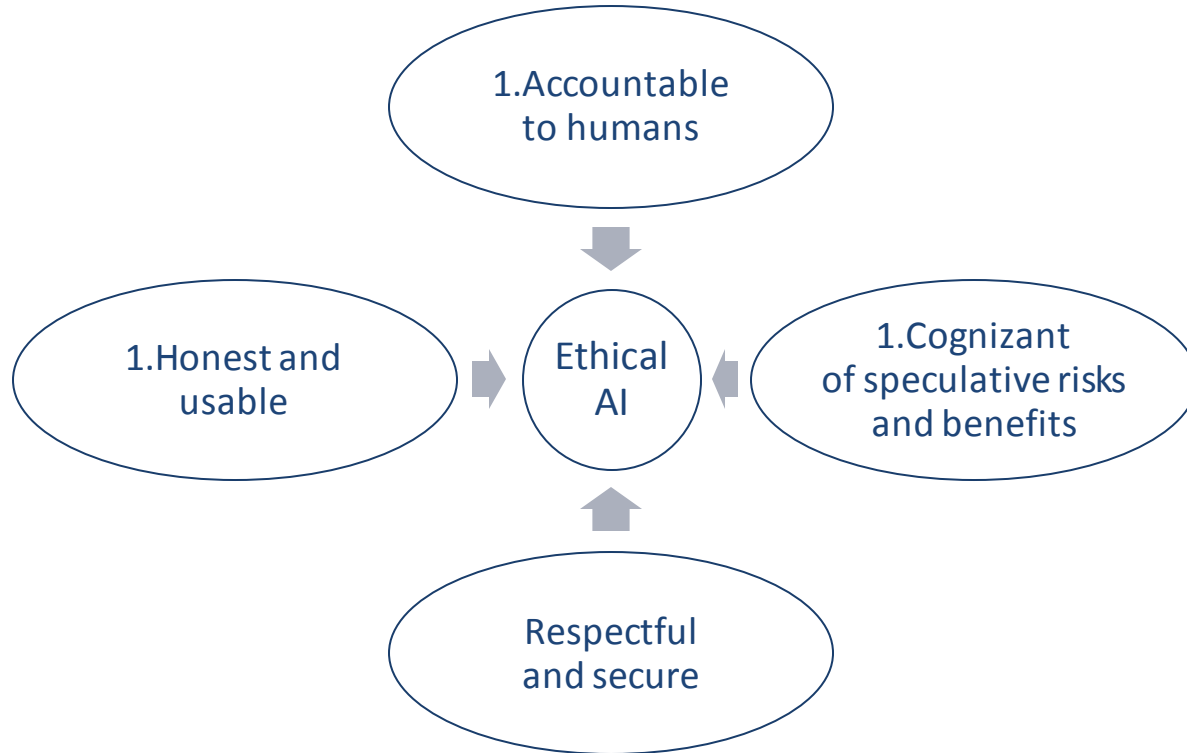
Team Signatures and Date

About the SEI
The Software Engineering Institute is a federally funded research and development organization (FRDO) whose mission is to advance the state-of-the-art in software engineering and to disseminate the results of that research to the software engineering community. For more information, visit <http://www.sei.cmu.edu>.

Contact Us
Carnegie Mellon University
Software Engineering Institute
4800 Forbes Avenue, Pittsburgh, PA 15213-1502
Phone: 412.263.1000
Fax: 412.263.1001
Email: sei@cmu.edu

©2020 Carnegie Mellon University. SEI-20-1107 (v1.0)

UX Framework for Designing Trustworthy AI



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

RightStaff Scenario

AI shift scheduling system

Users: Store managers of fast food restaurants

Goals of RightStaff:

- Faster staffing decisions and scheduling
- Reduced bias of shift selection

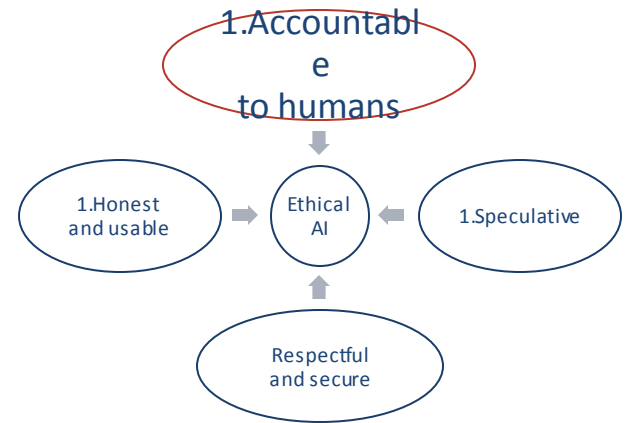
Accountable to Humans

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



Significant decisions

Significant decisions made by the AI system will be

- explained
- able to be overridden
- appealable and reversible

RightStaff

- Manager able to reschedule people as needed

Responsibilities and limitations explicitly defined

For AI system and human(s)

RightStaff (*AI System or Manager?*)

- Picks employees to schedule?
- Defines shifts?
- Method to integrate new information?
 - Sick time
 - Resignations

Abusability Testing

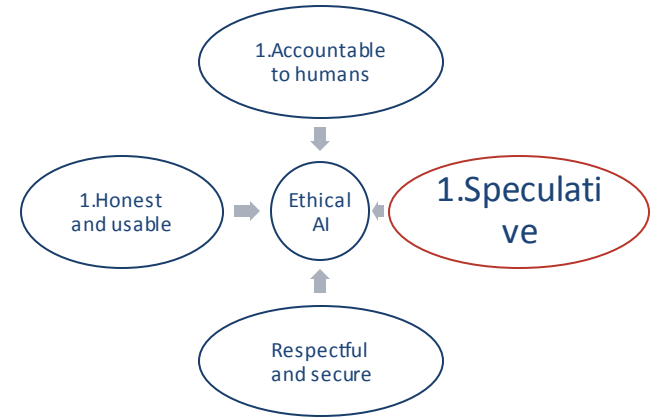
Feature added to enable RightStaff to turn off by itself

- What are limits to functionality?
- How is the situation communicated?
- How could this be abused/misused?
- Implications?
- Risks?

Cognizant of Speculative Risks and Benefits

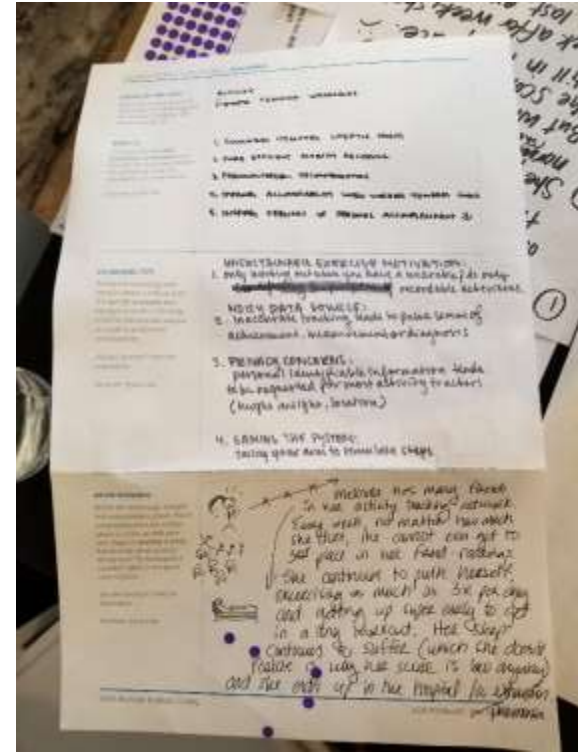
Identify full range of

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



Speculative: Conduct UX research - activate curiosity

- Speculate about misuse and abuse
- Potential severe abuse and consequences
- Perspective of people in frequently marginalized groups
- “Black Mirror” episodes



“Black Mirror” episode

- RightStaff begins prioritizing people with easier schedules
- Managers approve these schedules, reinforcing bias
- People who were previously discriminated against are *still* discriminated against

- What else?

Speculative: Create communication & mitigation plans

Plan for unwanted consequences

Misuse and abuse of AI system

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

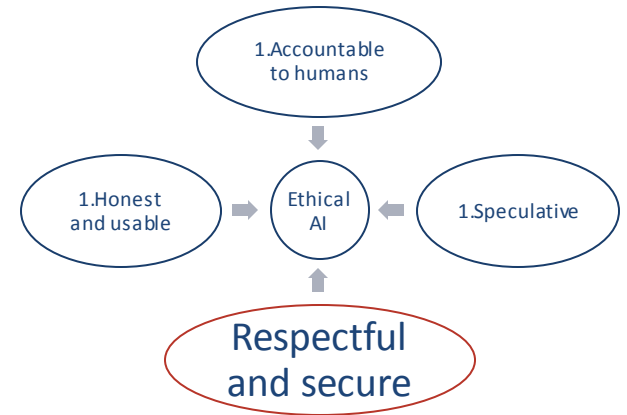
Respectful and Secure

Values of humanity, ethics, equity, fairness, accessibility, diversity and inclusion

Respect privacy and data rights

Make system robust, valid and reliable

Provide understandable security



Respectful and Secure

RightStaff

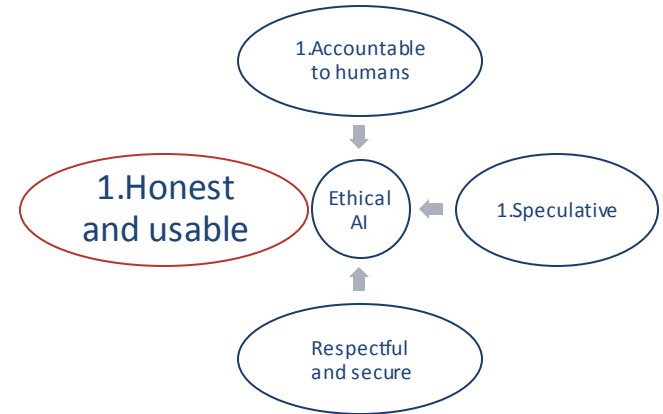
- Who has visibility to reasons for changing schedules?
- How is that information used?
- How is PII* of employees protected?

*PII is Personally Identifiable Information (social security number, address, etc.)

Honest and Usable

Value transparency with the goal of engendering appropriate trust

Explicitly state identity as an AI system



Fair: Unwanted bias in data

Show awareness of known and desirable bias

Acknowledge issues

Overcommunicate on issues

RightStaff

- System built to reduce the known bias in existing data
- Make it easy to report bias (or prevent it)

Design for Human-Machine Teaming

Provide transparency regarding AI limitations

- boundaries and unfamiliar scenarios

Encourage appropriate trust

Speculate about misuse and abuse

Prevent or plan to mitigate situation

Challenges

Need more speculative activities

Engage people in this hard and necessary work

A challenge - building on our experiences with UX/HCI and accessibility...

AI Systems are not fully able to team with humans yet,
but we need to be ready!



Methods, Mechanisms, and Mindsets

Engage in Critical Oversight

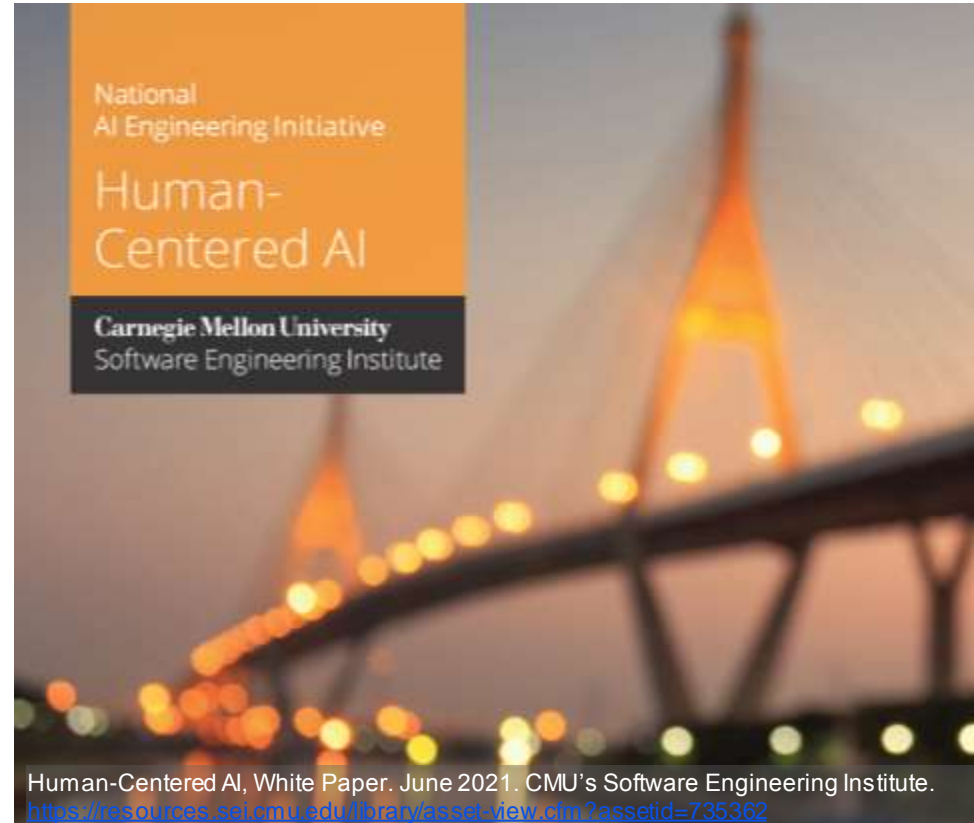
Engage in critical oversight

“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



Data transparency

Must understand data at a deep level

Provenance and creator's motivation

- What data included? Why?
- What not included? Why?

Use

- Datasheets for Datasets¹
- Model Cards for ML systems²

1. T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. Daumé III, and K. Crawford. Datasheets for Datasets. The latest version of this paper can be found online at <https://arxiv.org/abs/1803.09010>

2. M. Mitchell et al., "Model Cards for Model Reporting," Proc. Conf. Fairness Account. Transpar., pp. 220–229, Jan. 2019, doi: 10.1145/3287560.3287596

How can Data be biased?

Lawn Care Company

System: Select the right lawn care treatment. Save time.

Data: Multiple data source choices?



Selecting Data Source

Company A

- Primarily uses chemicals to treat lawns
- Data likely biased towards chemical use

Company B

- Only “all natural” treatments
- Data likely biased against chemical use

Selecting Data Source

Company A

- Primarily uses chemicals to treat lawns
- Data likely biased towards chemical use

Company B

- Only “all natural” treatments
- Data likely biased against chemical use

**Neither are wrong.
Both are biased.**

Bias in data, algorithm selection, and training

Unintended and purposeful bias

Misuse and abuse of the system

Understand inherent bias and amount of variance – data:

- Motivation
- Composition
- Collection process
- Recommended uses, etc.

Goal: Transparency and accountability.



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

Joy Buolamwini, Algorithmic Justice League

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

Coded Gaze

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

THE
OPEN MIND



Leaders must establish psychological safety



Regular Auditing

AI should not be assumed to be “stable” – it is dynamic

AI systems require continuous human oversight to identify risks of bias, misuse, abuse, and unintended consequences

Probe with hypothetical cases

Checks for bias, brittleness or potential distribution shift

Access history of system operation and usage*

*Consider ethical principles when determining what data needs to be collected.

Challenge: Broaden our work

Examining dynamic data and evaluating dynamic outcomes

- Is this the right data? What has changed?
- Is there appropriate 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



ARTIFICIAL INTELLIGENCE PORTFOLIO

Responsible AI Guidelines

Operationalizing DoD's Ethical Principles
for AI

Download DIU's
Responsible AI
Guidelines report and
learn how to
implement ethical AI
principles.

[Responsible AI Guidelines](#)

<https://www.diu.mil/responsible-ai-guidelines>

Phase I: Planning



Phase I: Planning Worksheet for DIU AI Guidelines

Have you clearly defined tasks, quantitative performance metrics, and a baseline against which to evaluate system performance?

responsible, reliable, governable

Have you evaluated ownership of, access to, provenance of, and relevance of candidate data/models?

traceable

Are end users, stakeholders, and the responsible mission owner identified?

responsible, equitable

Have you conducted harms modeling to assess likelihood and magnitude of harm?

equitable, governable

Have you identified the process for system rollback and error identification/correction?

responsible, governable

PROCEED TO DEVELOPMENT

.....1

.....2

.....4

Review
Read the AI Guidelines and process to help guide thinking and then later to avoid unintended consequences in creating AI systems. Worksheets for planning, development and deployment efforts. These worksheets are intended to supplement or replace existing laws and

AI Principles for the development and use of artificial intelligence in Defense in 2020:
exercise appropriate levels of judgment and care, while remaining transparent, and use of AI capabilities.
take deliberate steps to minimize unintended bias in AI capabilities.

Traceable. The Department's AI capabilities will be developed and deployed such that relevant personnel possess an appropriate understanding of the technology, development processes, and operational methods applicable to AI capabilities, including with transparent and available methodologies, data sources, and design procedure and documentation.

Reliable. The Department's AI capabilities will have explicit, well-defined uses, and the safety, security, and effectiveness of such capabilities will be subject to testing and assurance within those defined uses across their entire life-cycles.

Governable. The Department will design and engineer AI capabilities to fulfill their intended functions while possessing the ability to detect and avoid unintended consequences, and the ability to disengage or deactivate deployed systems that demonstrate unintended behavior.

<https://www.diu.mil/responsible-ai-guidelines>

AI has great potential, develop with caution

“AI will ensure appropriate human judgement and not replace it”

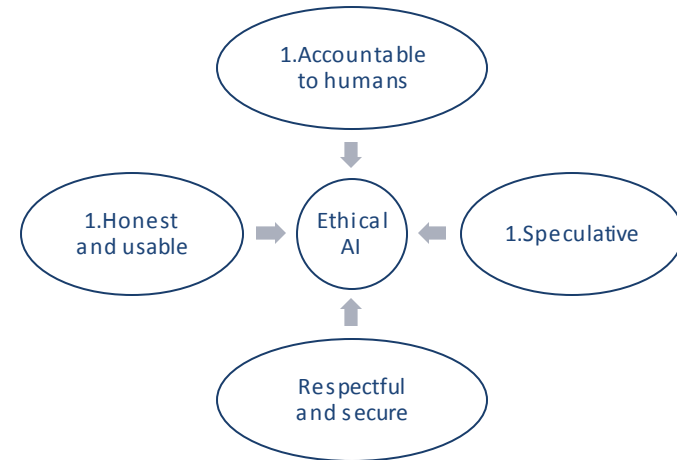
- Defense Innovation Board. 2019

We aren't perfect, AI won't be perfect

Empower diverse teams, inclusive environments

Encourage deep conversations

Activate curiosity; be speculative; imaginative



Design to work with, and for, people



User Experience Honeycomb
Peter Morville, et al.

1. Understand complexity of context
2. Design for human-machine teaming
3. Engage in critical oversight

Human-Centered AI



Carol J. Smith

Twitter: @carologic

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

CMU Software Engineering Institute,
AI Division

Twitter: @SEI_CMU_AI

Montréal Declaration for a Responsible Development of AI



1. Well-being
2. Respect for autonomy
3. Protection of privacy and intimacy
4. Solidarity
5. Democratic participation
6. Equity
7. Diversity inclusion
8. Prudence
9. Responsibility
10. Sustainable development

Categories of Harm (ServiceEase)

Use these categories of harm to evaluate how a product, service, or technology could cause harm.

CATEGORIES OF HARM	DEFINITION
FINANCIAL	Negative impact on finances, property, or other resources
HEALTH	Negative impact on mental, emotional, or physical health
TIME	Inefficient or unproductive activities, processes, or systems
FAIRNESS/EQUITY	Perpetuating or facilitating prejudice, bias and/or unfairness
SAFETY	Physical and/or emotional wellbeing compromised by fear, danger, or uncertainty
PRIVACY	Lack of control over personal information
MISINFORMATION	The creation, spread and/or amplification of false or inaccurate information intended to deceive
CONTROL	Inability to freely direct information, activities, or systems
TRANSPARENCY	Lack of disclosure of information, activities, or systems

ServiceEase