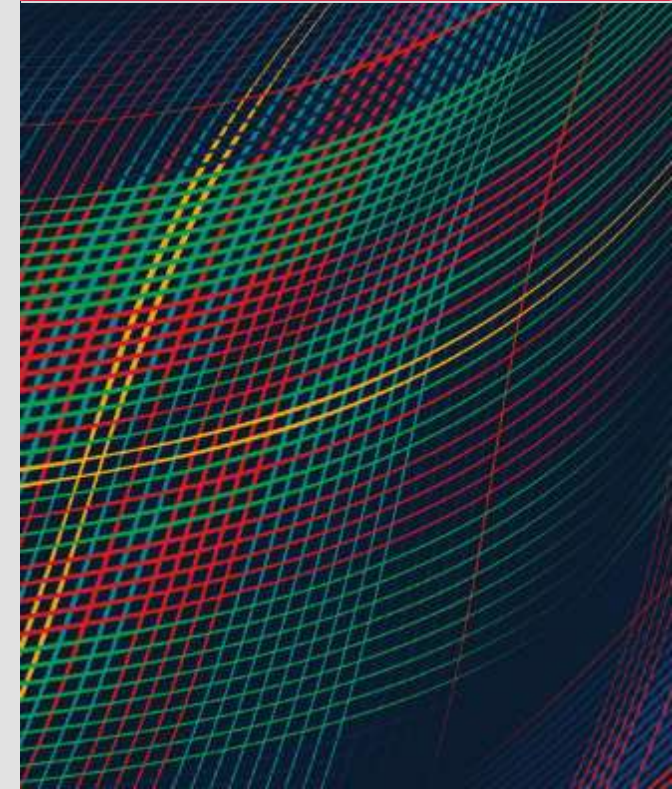


Academic Research Trustworthy AI Systems

TALARIA SUMMER INSTITUTE, JULY 2023

Carol J. Smith
Carnegie Mellon University, Software Engineering Institute
Sr. Research Scientist, Human-Machine Interaction



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.

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

DM23-0720

Carol J. Smith

Software Engineering Institute



AI Division Staff

- Sr. Research Scientist, human-machine interaction
- AI/ML, autonomy, emerging technologies
- Government agencies

Adjunct Instructor

- Interaction Design Overview
- Human-centered design
- Prototyping
- Design and iteration

Research to reduce risk and avoid waste

Software and AI systems

- Decision support
- Recommender systems
- Large Language Models (LLMs)
- Voice recognition (smart speakers)
- Speech to text/text to speech (live captions)
- And more...

**“The biggest waste
of all is building something
no one wants”**

Eric Ries @ericries via @MelBugai on Twitter at LeanStartupMI in 2011

Everything is Designed

Reduce Risk



https://www.reddit.com/r/dangerousdesign/comments/3s75gz/cooking_spray_insect_killer/

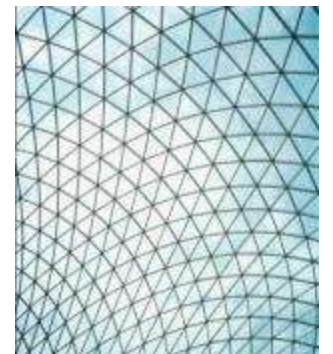
Consider changes over time



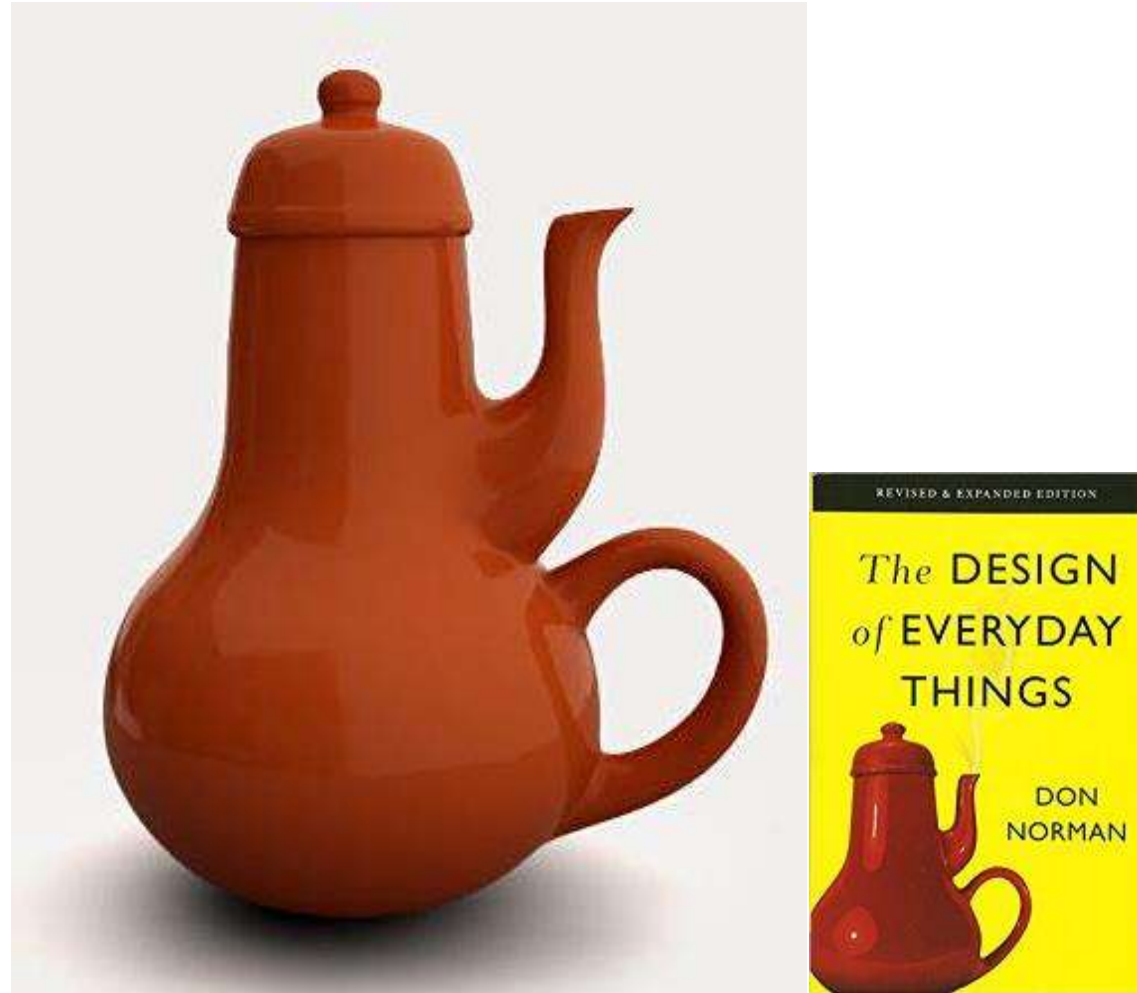
<https://boingboing.net/2014/02/25/the-story-behind-the-too-coo.html>

Long tradition – applying helpful patterns

- Human-Computer Interaction
- User Experience
- Usability
- Interaction Design

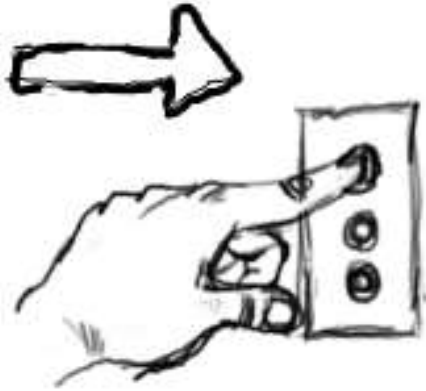


Design Better Experiences



The "masochist's teapot", as borrowed by Don Norman for his book *The Design of Everyday Things*.
Image: <http://playerside.blogspot.com/2012/12/the-design-of-everyday-gaming-feedback.html>

Make informed design choices



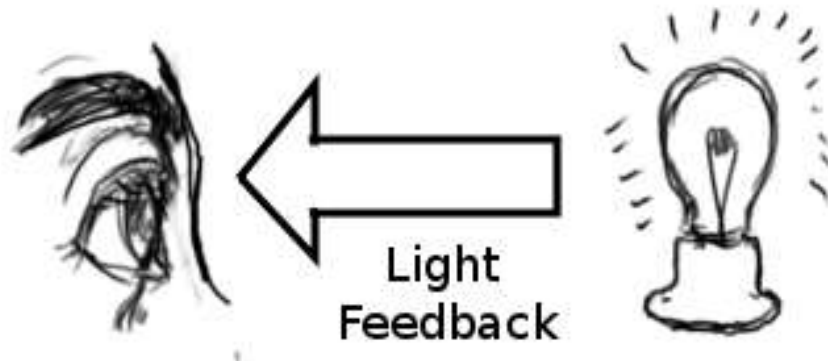
Button - Push



Switch - Flip



Knob - Rotate



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

My Work

Focus on interesting and influential topics

- Someone has an idea
- Explore options
- Plan a research project

Projects are long (1+ years)

Work with software engineers, machine learning researchers, scientists, and more...

Explore the situation

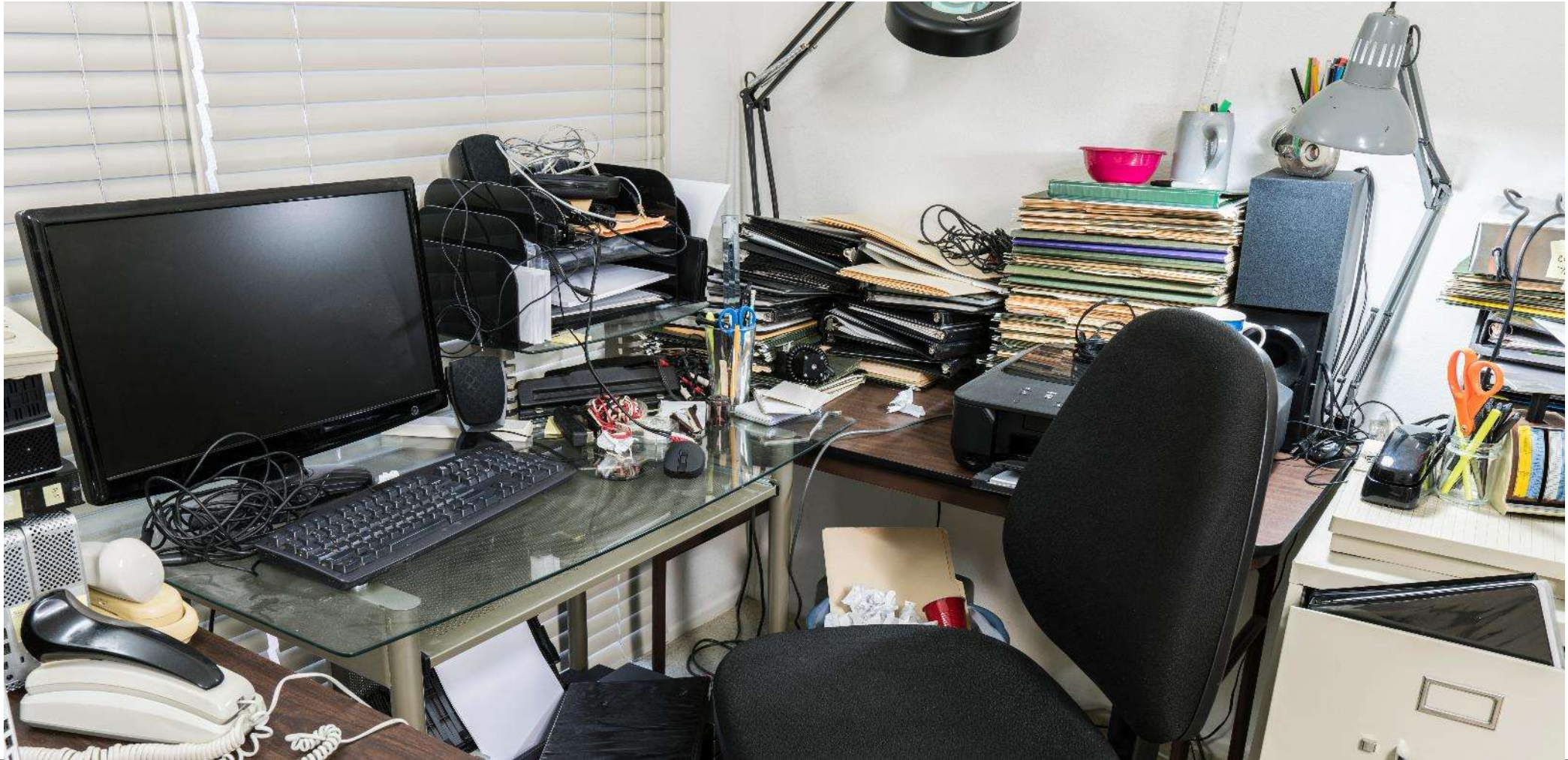
- Who would use the system?
 - What would they do?
 - What are their needs?
-
- Conduct Research
Observations,
Interviews



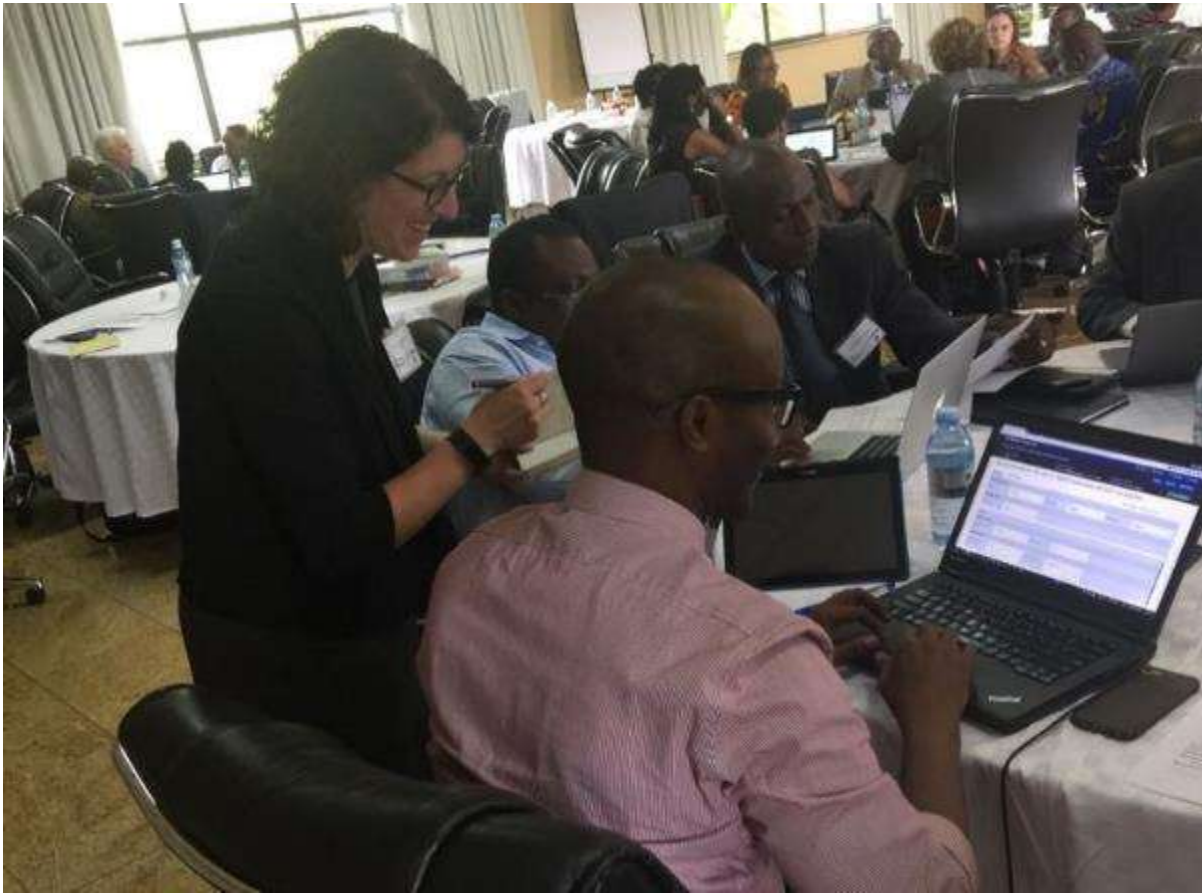
Observations...

- Learn real processes - existing product's actual use
- Workarounds
- Number and type of interruptions
- Remote considerations
- Understand how processes fit together
- What information is lost between systems (print and retype)
- Identify key opportunities

Artifacts extend understanding



Observe use of system and environment



Analyze Information

Map out information

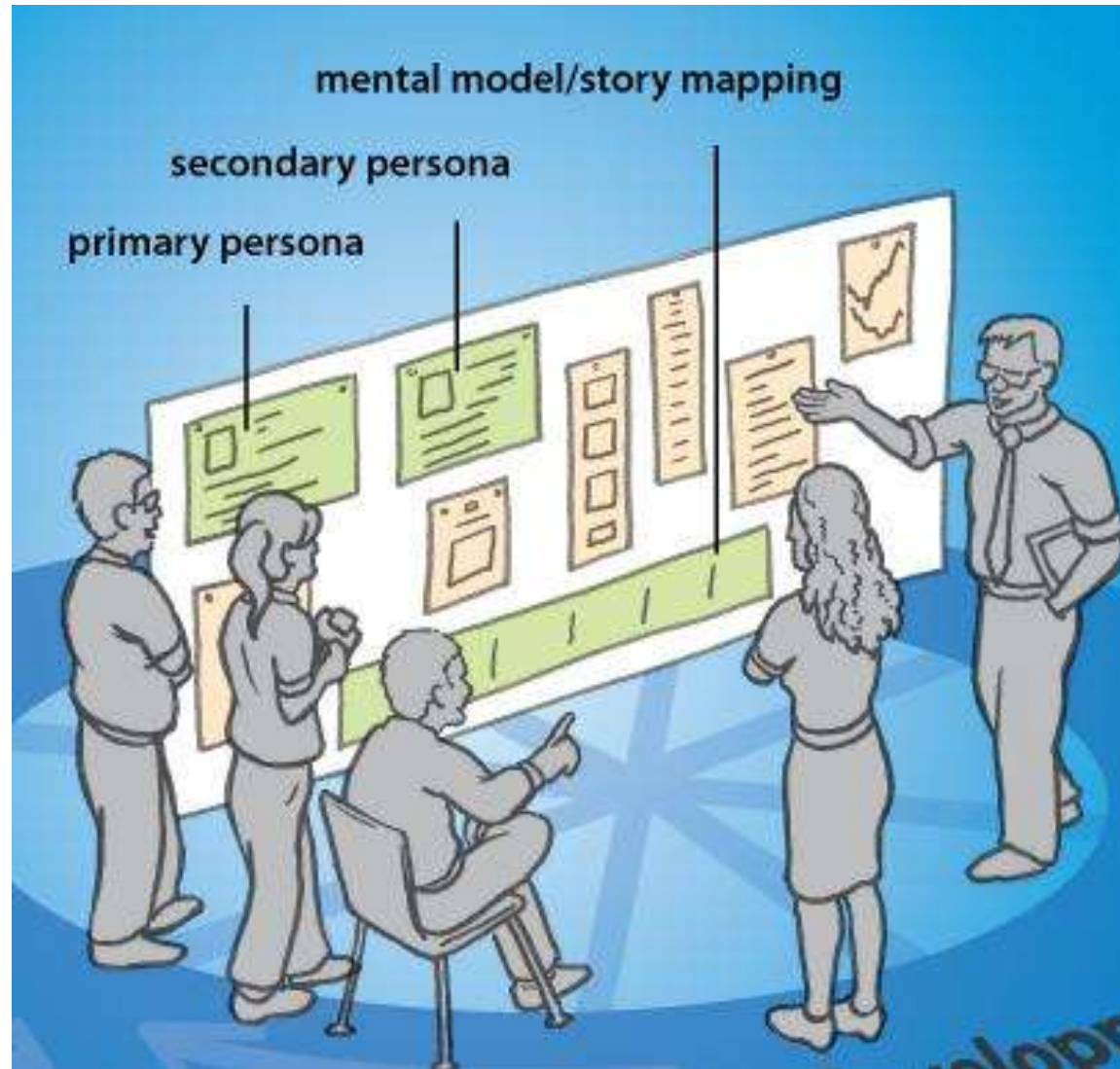
- Look for patterns
- What people do
- Why they do it?
- Create helpful references
- Not “books” for shelving

Sub Tasks	Scenario	Considerations/Influencers	Pain-Points	Functionality
Check out a book from a library	Scenario: A user wants to check out a book from a library. The user is standing at the library counter and is talking to the librarian. The user is holding a book and is looking at the librarian. The user is wearing a white shirt and glasses.	Considerations/Influencers: The user is standing at the library counter and is talking to the librarian. The user is holding a book and is looking at the librarian. The user is wearing a white shirt and glasses.	Pain-Points: The user is standing at the library counter and is talking to the librarian. The user is holding a book and is looking at the librarian. The user is wearing a white shirt and glasses.	Functionality: The user is standing at the library counter and is talking to the librarian. The user is holding a book and is looking at the librarian. The user is wearing a white shirt and glasses.

Example of a Task Analysis by Todd Zaki Warfel from his Agile2010 presentation "Opening the Kimono a look behind the design process."

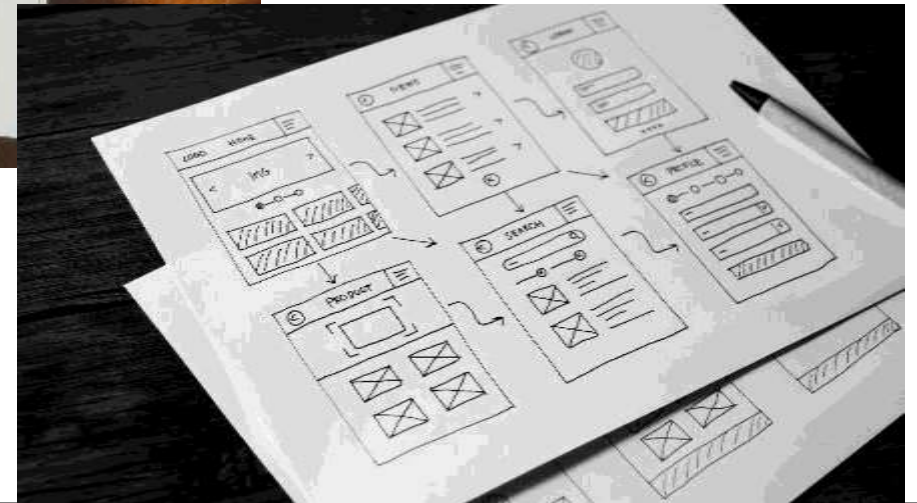
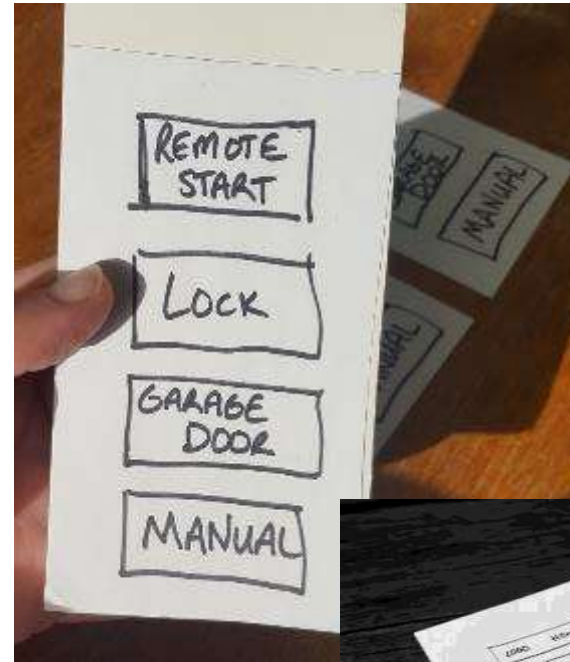


Shared understanding and thoughtful discussion



Prototype – Experiments to Learn

- How can I best support the user?
- How do I solve this problem?
- What interaction will meet the end users' needs?

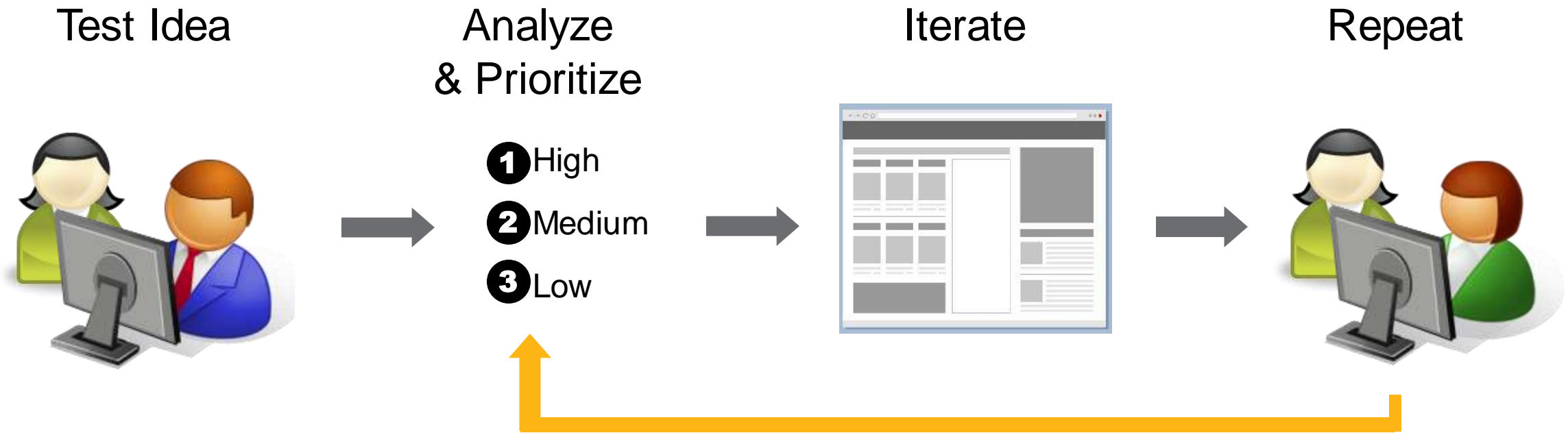


Early Learning: Observe use of prototype

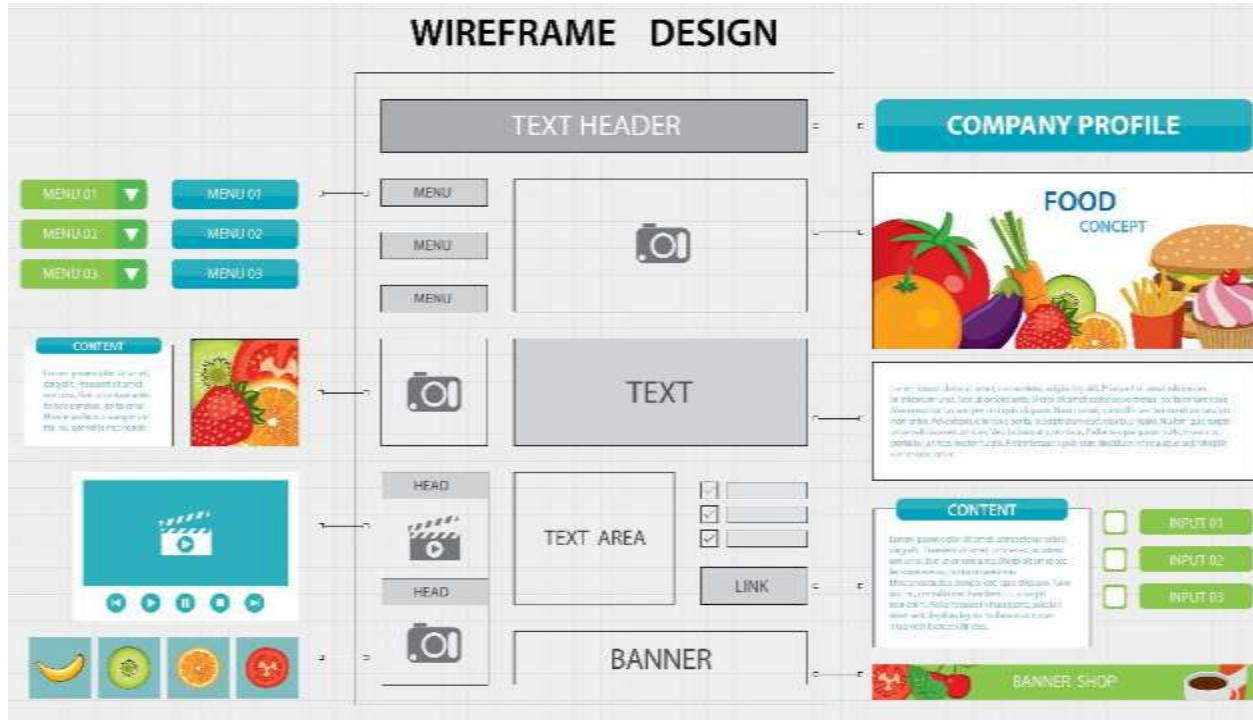


Image from Steve Krug's ["Rocket Surgery Made Easy"](#)

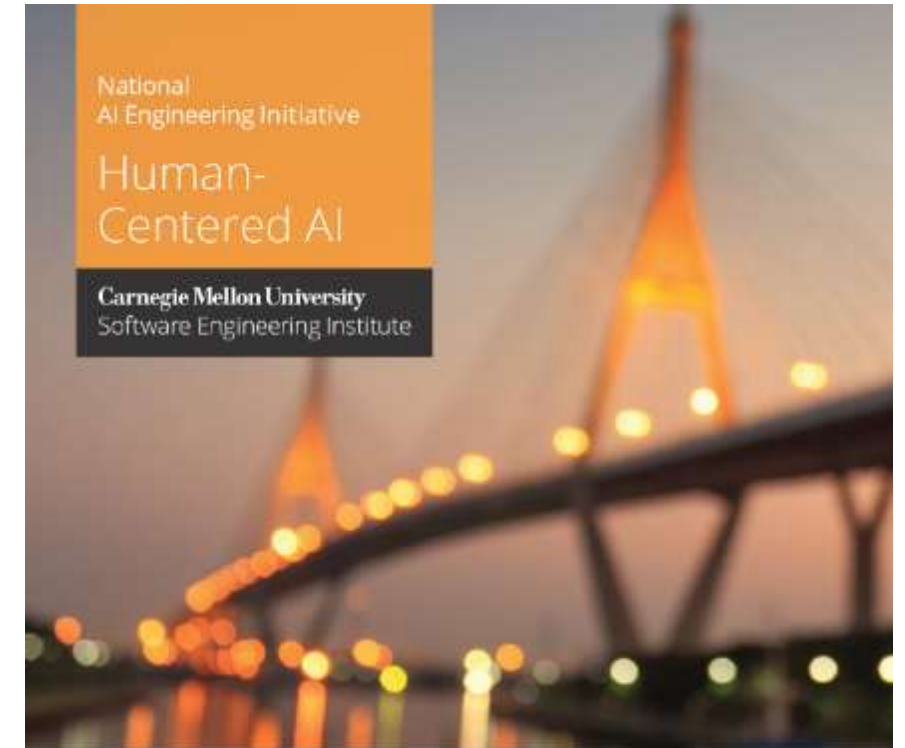
Iterative Cycles: Feedback and Improvement



Work with Visual Designers



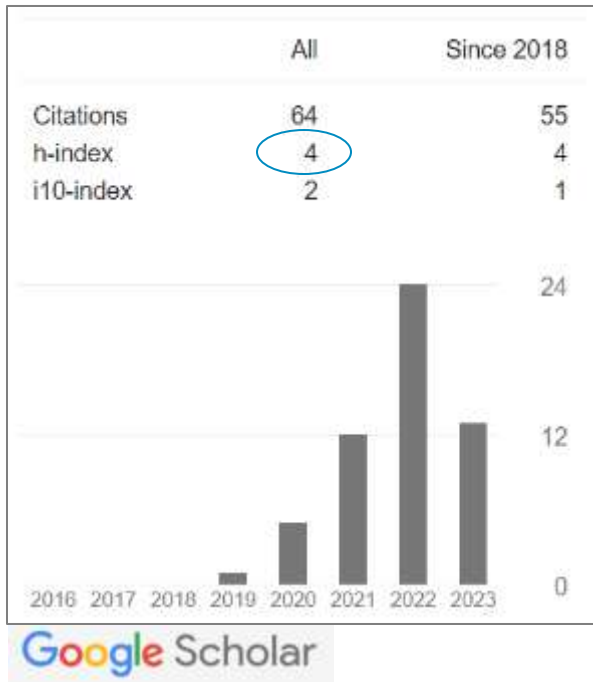
Presentations and Publications



Human-Centered AI, White Paper. June 2021. CMU's Software Engineering Institute. <https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=735362>

Carol J. Smith. Designing Trustworthy AI: A Human-Machine Teaming Framework to Guide Development. <https://arxiv.org/abs/1910.03515>
 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>

Publishing for an Academic Footprint



Publications	5
h-index	3
Citations	105
Highly Influential Citations	6



Research Interest Score	13.3
Citations	7
h-index	2

ResearchGate

Articles	Citation numbers
1	33
2	30
3	20
4	15
5	7
6	6 = h-index
7	5
8	4

Univ of Waterloo: <https://subjectguides.uwaterloo.ca/calculate-academic-footprint/YourHIndex>

Making Responsible Artificial Intelligence



AI systems can

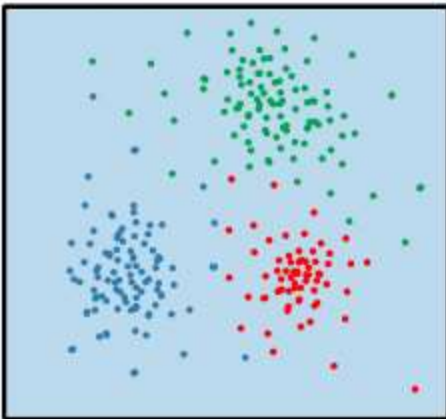
- recognize patterns
- create predictions
- make decisions, and/or
- generate new content

without being explicitly programmed

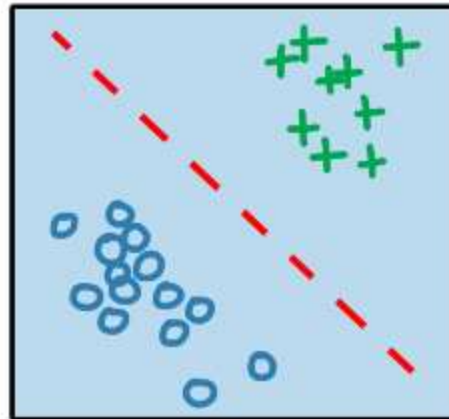
Artificial Intelligence

machine learning

unsupervised learning



supervised learning



reinforcement learning

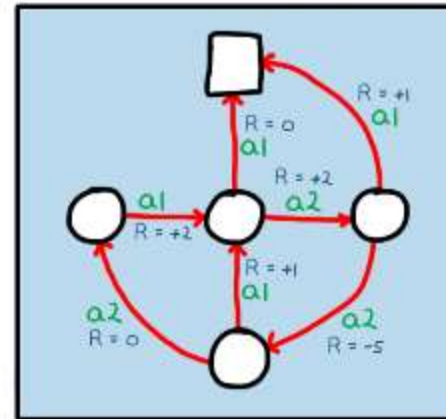
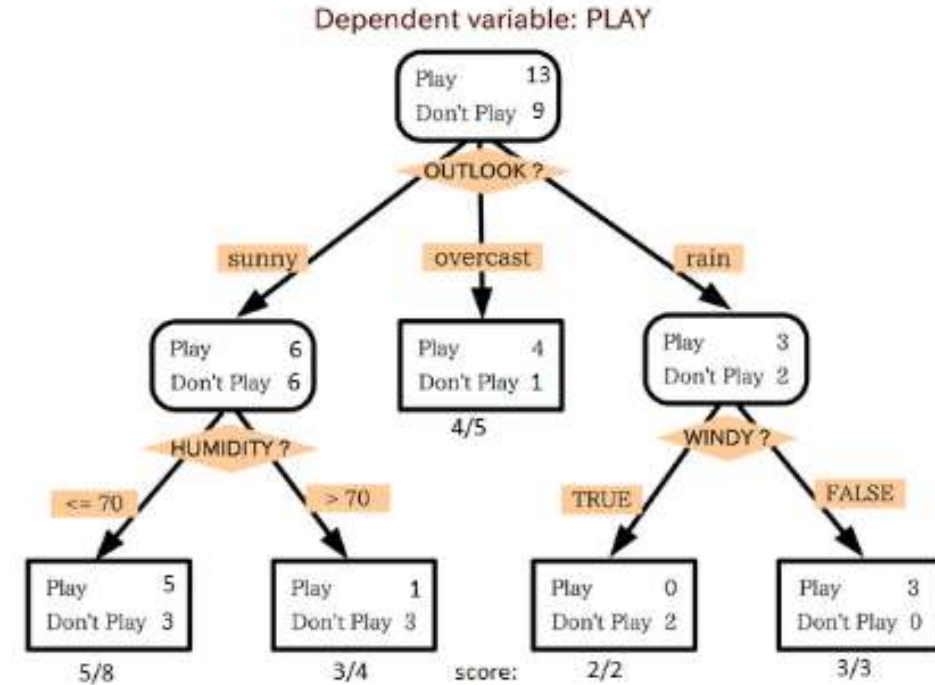


Image © 1994-2022 The MathWorks, Inc.

Deep Learning, Neural Networks, Large Language Models (LLMs)

AI / Machine Learning

- Algorithms
 - math + programming
- Model (AI)
 - Algorithms + data
- Know ONLY what taught
- Control ONLY given control of
- Aware of nuances
 - can continue to learn



source: [statsexchange](https://statsexchange.com)

<https://www.analyticsvidhya.com/blog/2015/08/common-machine-learning-algorithms/>

Taxonomies and Ontologies coming to life (NOT like humans learn)



AI is NOT sentient

Not unknowable

Never Enough Time

- Physician: ~90 hours reading a week*
- AI could bring that information to the physician
- Enabling more evidence-based decisions

Alper, Brian S. et al. "How Much Effort Is Needed to Keep up with the Literature Relevant for Primary Care?"
Journal of the Medical Library Association 92.4 (2004): 429–437. Print.
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC521514/>



Transfer human concepts and relationships



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



What is a tomato?

Fruit?

Vegetable?

Computer Vision - Image Recognition

Train set



Data encountered



Use case courtesy of Dr. Eric Heim, CMU Software Engineering Institute

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.”

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>

THE
OPEN MIND

Data is biased

- Not inherently neutral.
- Reflect priorities, preferences, and prejudices of people making them

The image shows a video frame of a woman with a red bounding box around her face. Below it is a smaller version of the same image. To the right is a software interface displaying demographic data and a gaze score.

Gender: Female
Age: 22
Ethnicity: Black

Coded Gaze Score: 4/13

	Gender	Age*	Detected
IBM	M	21**	✓
Microsoft	✗	✗	✗
Face++	✗	✗	✗
Kairos	M	29	✓

Gender Shades Project by Joy Buolamwini. MIT Media Lab. <https://www.media.mit.edu/projects/gender-shades/overview/>

Algorithms of Oppression, Dr. Safiya Noble

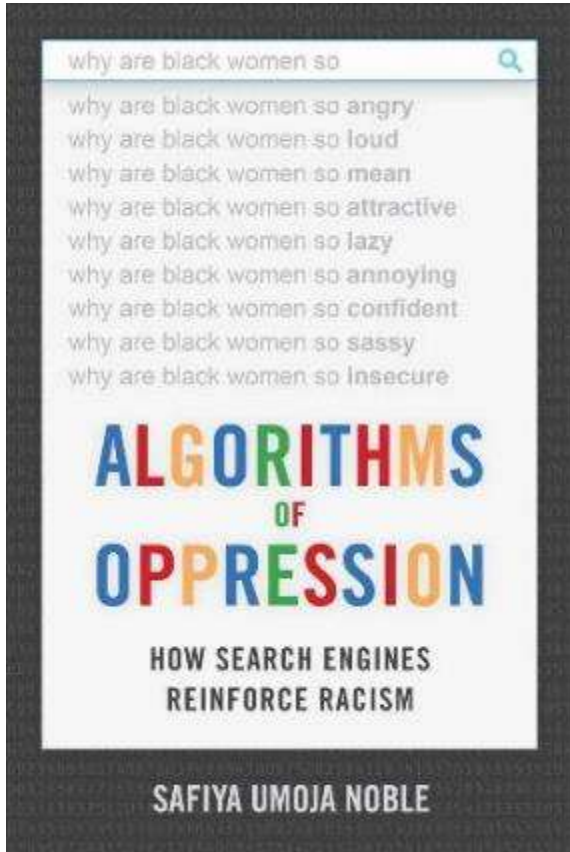


Photo from article: Google This: Algorithmic Oppression. ACLU News & Commentary. February 19, 2021. <https://www.aclu.org/news/privacy-technology/google-this-algorithmic-oppression/>



Responsible,
Intentional
Design

**Just because you can,
doesn't mean you should.**

Conversations for Understanding

Difficult Topics

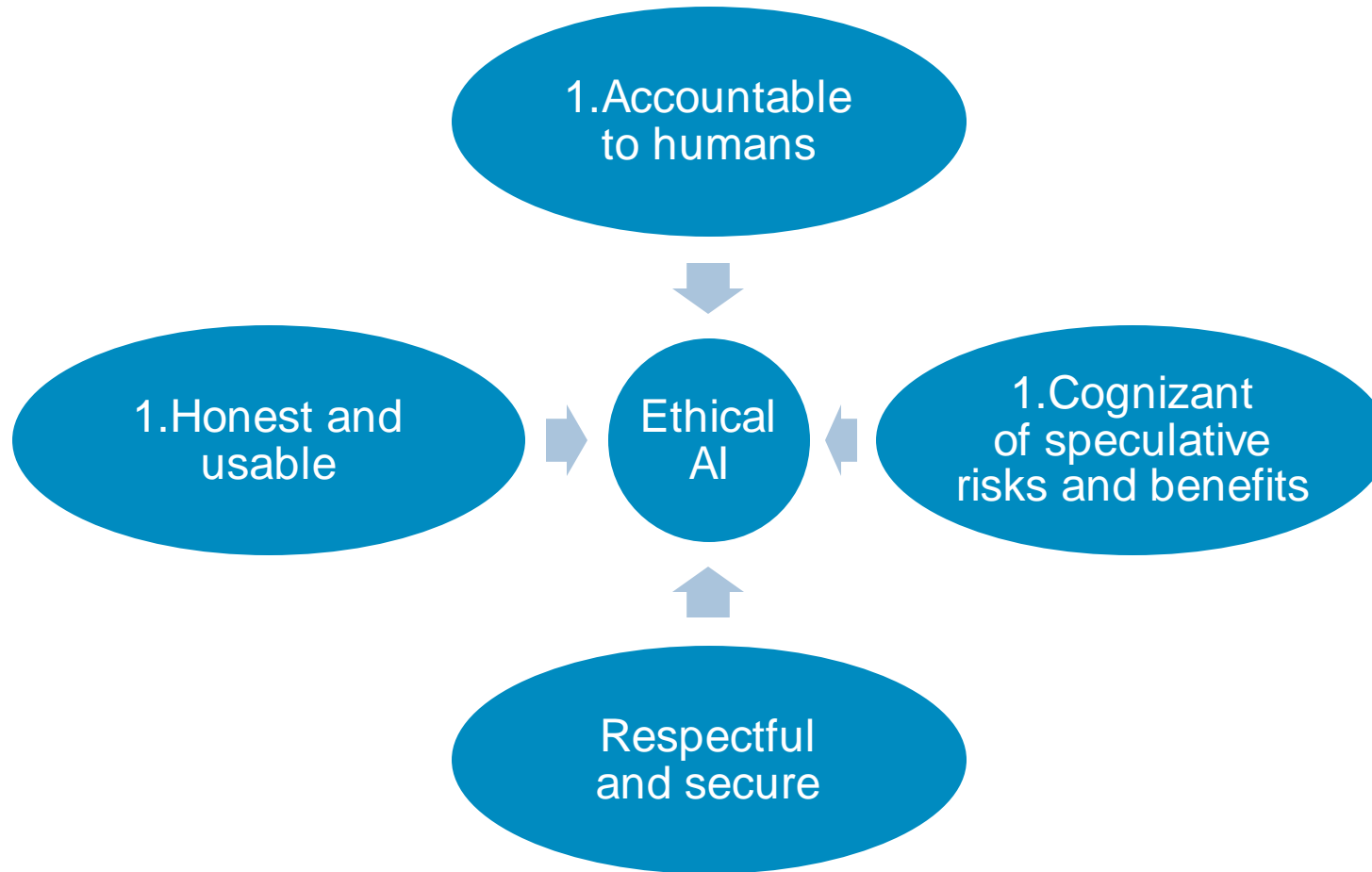
- What do we value?
- Who could be hurt?
- What lines won't our AI cross?
- How are we shifting power?*

*"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> **"How is this ML model shifting power?" @riakall #NeurIPS2019

Photo by Pam On Unsplash



UX Framework for Designing Trustworthy AI



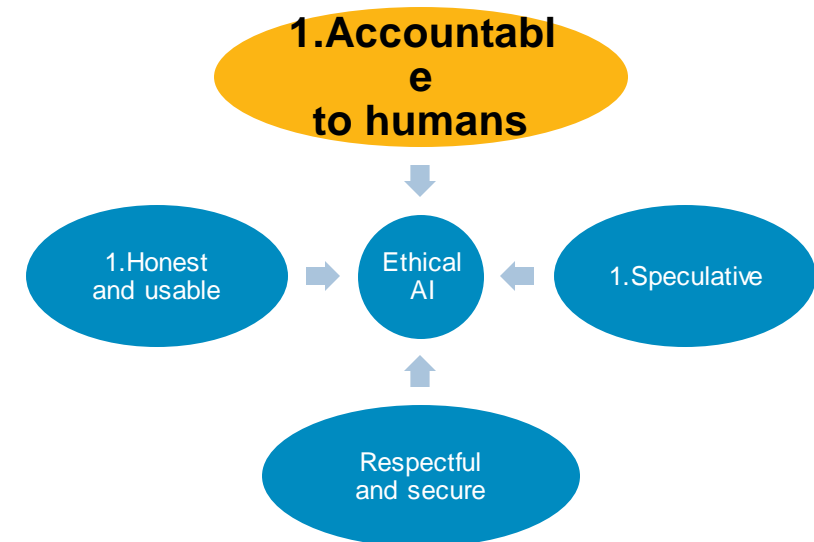
Designing Trustworthy AI for Human-Machine Teaming. By Carol Smith. Software Engineering Institute Blog. March 9, 2020.
https://insights.sei.cmu.edu/sei_blog/2020/03/designing-trustworthy-ai-for-human-machine-teaming.html

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



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

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 explicitly defined

- Between AI system and human(s)
- **RightStaff** (*AI System or Manager?*)
 - Picks employees to schedule?
 - Defines shifts?

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

Checklist and Agreement - Downloadable PDF:
<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 accountable, de-risked, respectful, secure, timely, and usable AI in enterprise AI systems with diverse, well-aligned, trusted ethics. An initial version of this document was presented with the paper *Designing Trustworthy AI: Human-machine Teaming Framework to Guide Development* by Carol Smith, available at <https://arxiv.org/abs/1910.03545>.

<p>We will design our AI system with the following in mind:</p> <ul style="list-style-type: none"> Designers/developers have the ultimate responsibility for a decision and outcome: <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 preserved for final decisions that affect a person's life, quality of life, health, or reputation. Humans are always able to monitor, control, and deactivate systems. Significant decisions made by the AI system will be explainable <ul style="list-style-type: none"> and also understandable appealable and reversible 	<p>We will to specifically identify the full range of risks and benefits:</p> <ul style="list-style-type: none"> Harmful, malicious use and consequences as well as good, life-affirming use and consequences We will be cognizant and exhaustively research the related consequences. <p>We will create plans for the misuse/abuse of the AI system, including the following:</p> <ul style="list-style-type: none"> communication and responsible information with all affected people mitigation plans for managing the identified specific risks <p>We value respect and security:</p> <ul style="list-style-type: none"> Incorporating our values of humanity, ethics, equity, fairness, accessibility, diversity and inclusion reporting privacy and data rights (only necessary data will be collected) providing understandable security methods making the AI system robust, valid, and reliable 	<p>We value 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. Disclosures, laws, branding, and trademarks, and known and explicitly stated algorithms and models are appropriate and available. Centered and do not create possible or harmful consequences. Transparent justification for recommendations and outcomes is provided. Strong password and interpretable monitoring systems are available. <p>We value honesty and usability:</p> <ul style="list-style-type: none"> Humans can easily discern when they are interacting with the AI system vs. a human Humans can easily discern when and why the AI system is taking action and/or making decisions. Improvements will be made regularly to meet human needs and technical standards.
---	--	--

Team Signatures and Date

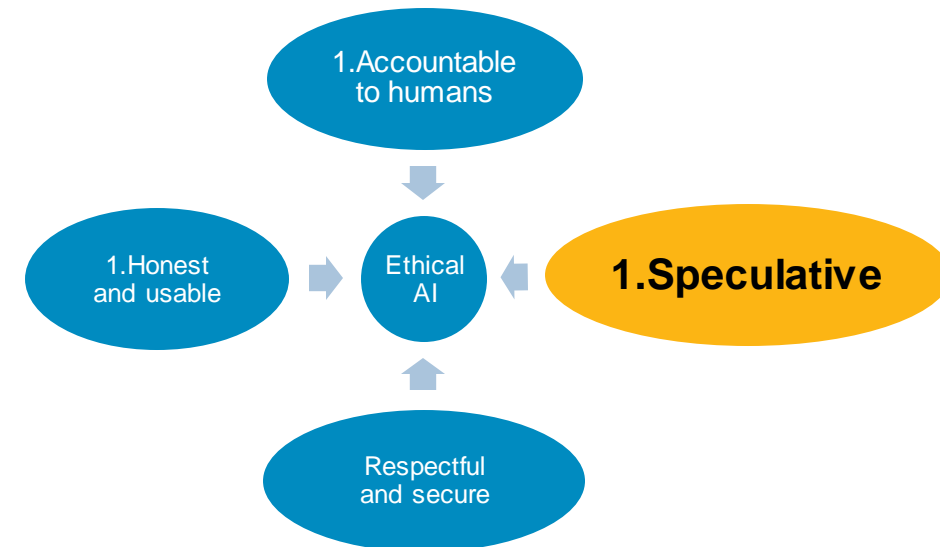
About the SEI
 The Software Engineering Institute (SEI) is a non-profit organization that provides research, education, and training in software engineering and related fields. SEI is a leader in the field of software engineering and has been instrumental in the development of many of the standards and best practices that are used in the industry today.

Contact Us
 4800 Forbes Avenue, Pittsburgh, PA 15260-1501
 412.263.1000
www.sei.cmu.edu
sei@sei.cmu.edu

© 2023 Carnegie Mellon University. All rights reserved.

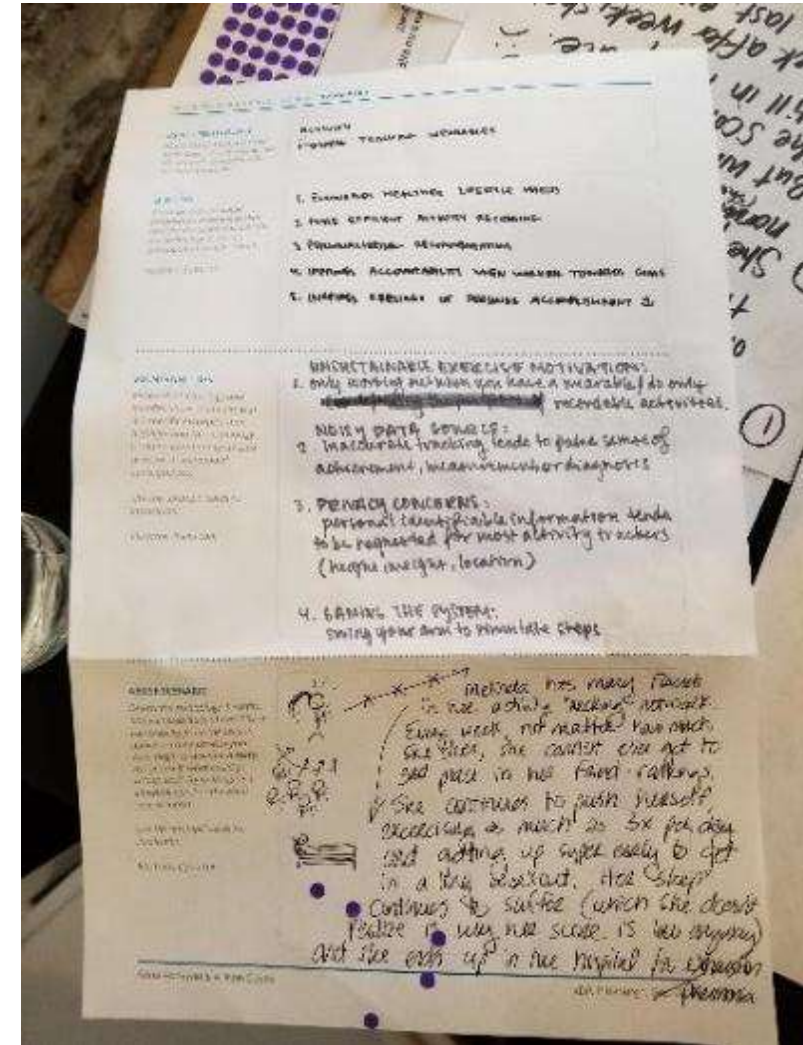
Cognizant of Speculative Risks and Benefits

- Identify full range of
 - Harmful, malicious use, as well as good, beneficial use
 - Unwanted/unintended consequences



Prevent bias - speculate about misuse and abuse

- At earliest moments (before coding)
- Activate curiosity
- Perspective of people in frequently marginalized groups



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

Reward team members for finding ethics bugs

**Ayanna
Howard**



Bias Scenario

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

Abuse Scenario

- Managers want to avoid providing benefits to employees
- RightStaff adjusted to ensure that no one has full shifts
- No regular staff can keep benefits

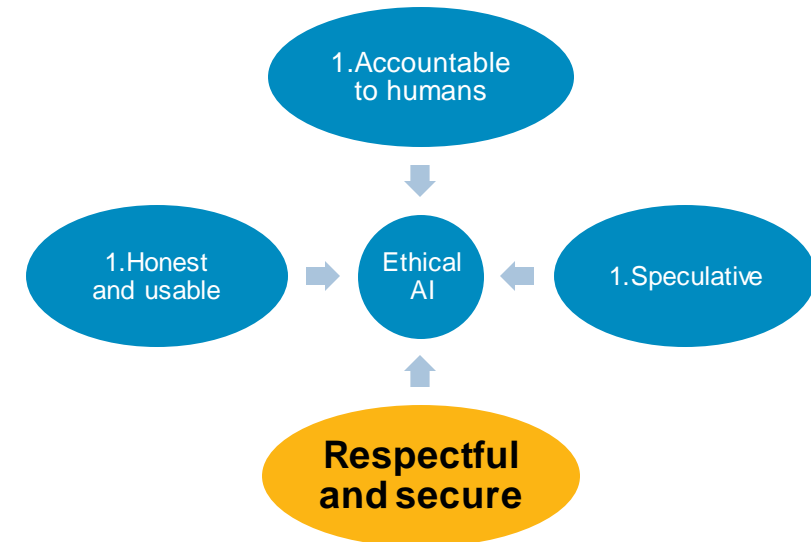
- What else?

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



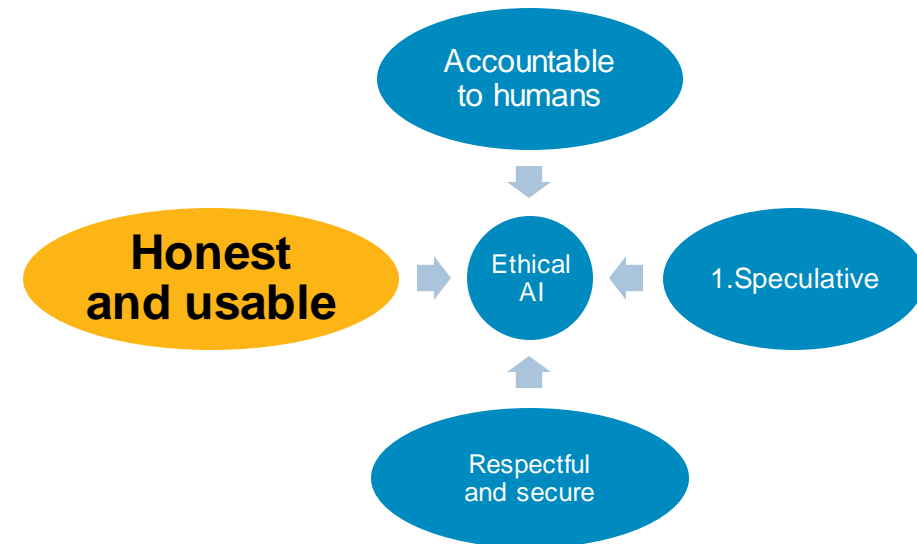
Respectful and Secure

- **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 trust
- Explicitly state identity as an AI system



Fair: Identify bias in data

- Show awareness of known and desirable bias
- Acknowledge and communicate about unwanted issues

RightStaff

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

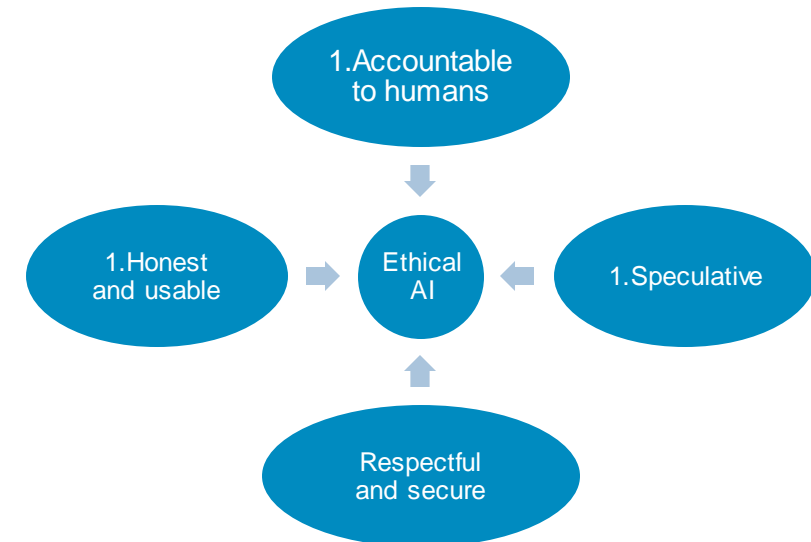
AI has great potential, develop with caution

Future AI's may be trusted to substitute human cognition and abilities.

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

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

- Empower diverse teams, inclusive environments
- Adopt technical ethics
- Encourage deep conversations (Checklist)
- Activate curiosity; be speculative; imaginative



Evangelize for human values

Ethical.
Transparent. Fair.

Advice

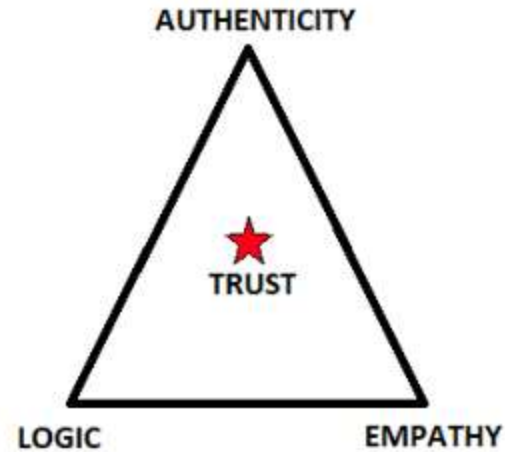
Make Relationships.

Find people who believe
in you - keep
in touch with them.

And...

- Trust your instincts
- Don't plan too far out
- Learn what you can
- Ask for feedback
- Work is long - look for fun, interesting work

Authenticity, Empathy, Logic



Understand where you [and others] wobble
- Frances Frei, Harvard Business School professor



Frances Frei, Harvard Business School, <https://blog.ted.com/how-to-rebuild-trust-frances-frei-speaks-at-ted2018/>

Build a manager Voltron

Develop diverse manager crew/coaching network

Look for people who:

- push you out of your comfort zone
- have different levels of experience than you
- are from a different background
- are good at things you're terrible at

When your manager isn't supporting you, build a Voltron by Lara Hogan Originally posted Jan 4, 2018
<https://larahogan.me/blog/manager-voltron/>



Your Responsibility

- Question status quo
- Stand for ethics – do the right thing
- Talk about bias
- Be inclusive and kind

Find your priorities

- What you can put up with and what you cannot.
- Invest in experiences
- Learn what you can for as long as you can
- Not learning, being challenged?
 - Change is ok - need to be able to tell the story
 - Look/ask for new opportunities
 - Go and try something new

Q & A Time!