

World Usability Day 2020
Cleveland, Ohio

Implementing Ethics in Emerging Technologies

Carol J. Smith
Sr. Research Scientist - Human-Machine Interaction

Twitter: @carologic @sei_etc

Software Engineering Institute
Carnegie Mellon University
Pittsburgh, PA 15213

Copyright Statement

Copyright 2020 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.

DM20-1051

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.

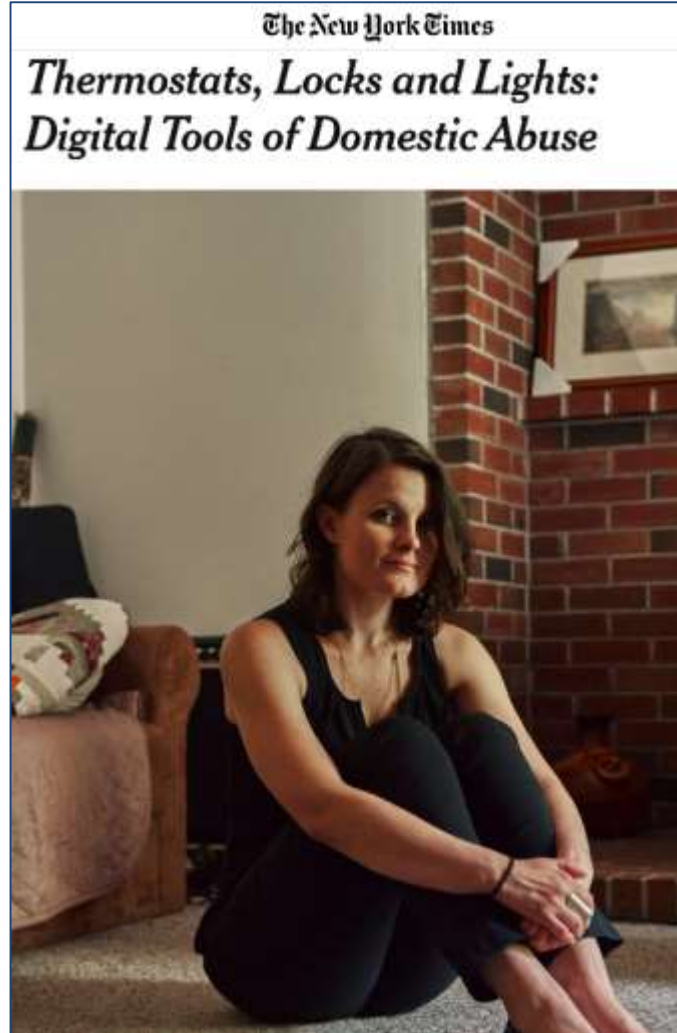
By Mark Hanrahan

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



Ring camera systems being hacked

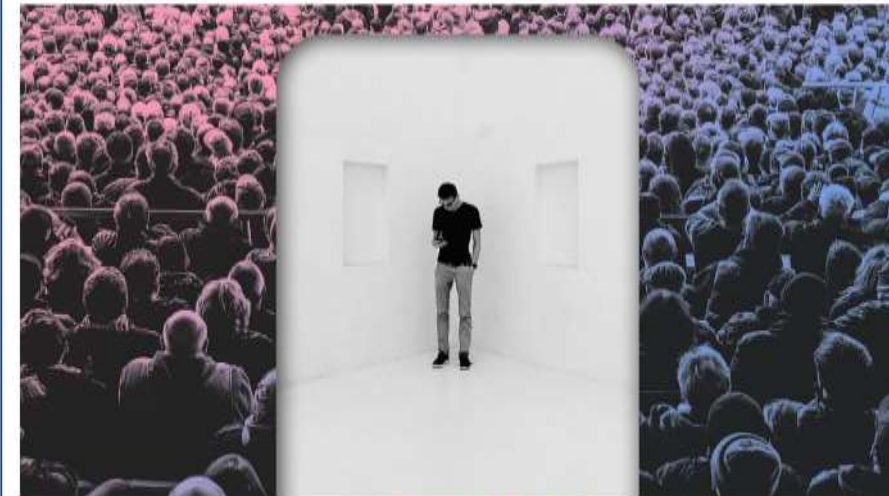
Multiple U.S. families have reported incidents of Ring camera systems being hacked in recent days.



07-07-17

Nest Founder: "I Wake Up In Cold Sweats Thinking, What Did We Bring To The World?"

Tony Fadell, one of the minds behind the iPod and the iPhone, mulls design's unintended consequences.



[Photos: Constantin Renner/EyeEm/Getty Images, davide ragusa/Unsplash]

Just ADAS - #Autonowashing

The image shows a screenshot of a CNN Business webpage. At the top, the navigation bar includes 'CNN BUSINESS', 'Markets', 'Tech', 'Media', 'Success', 'Perspectives', and 'Videos'. On the right, there are links for 'LIVE TV' and 'Edition'. The main headline reads: "'I'm not drunk, it's my car:' Tesla's 'full self-driving' gets mixed reviews". Below the headline, it says 'By Matt McFarland, CNN Business' and 'Updated 2:30 PM ET, Fri October 30, 2020'. A video player is embedded below the text, showing a man in a blue blazer standing next to a red Tesla car. The video player has a 'NOW PLAYING' indicator, a play button, and a progress bar showing 00:16 / 03:14. The video title is 'In 2010, Elon Musk had big plans for Tesla. Listen to his predictions'.

Responsible, Intentional Design

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



Early, purposeful work

In addition to the usual UX work

- What kind of improvements are expected with the emerging technology?
- How will the system partner with people? Compliment?
- What are the obvious risks?
- How might these systems be misused/abused?

Ethics

Based on well-founded standards of right and wrong

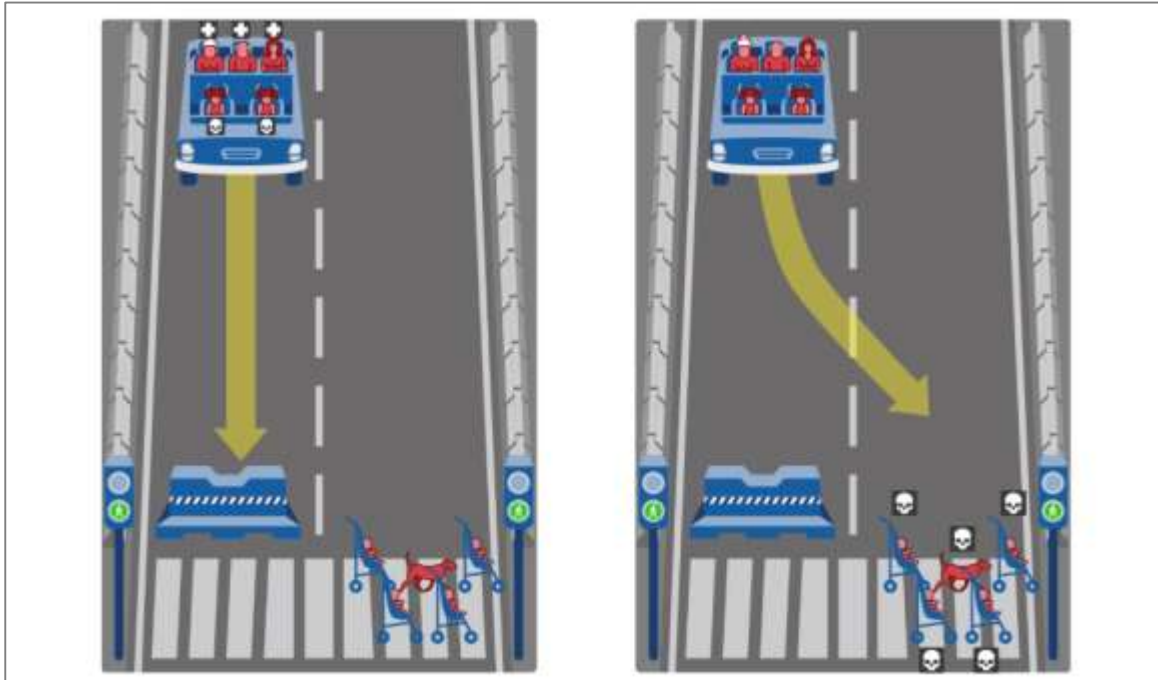
Standard of expected behavior that guides the correct course of action

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/>

Not Trolley Problems (hypotheticals in-the-moment)

MIT Moral Machine Project



Images: 1) MIT Moral Machine Project: <http://moralmachine.mit.edu/>

2) Does the Trolley Problem Have a Problem? <https://slate.com/technology/2018/06/psychologys-trolley-problem-might-have-a-problem.html>



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

Talk about biases in non-threatening, productive ways

Biased due to...

Social class

Resource availability

Education

Race, gender, sexuality

Culture, theology, tradition

More...

All systems have some form of bias

Complete objectivity is misleading.

Bias can have purpose and can be helpful.

The goal is to reduce unintended and/or harmful bias.

Teams - Diverse, talented and multi-disciplinary

Bringing their varied skills sets, problem framing approaches, and knowledge together.

- Gender, race, culture
- Education (school, program, etc.)
- Experiences
- Thinking process, skill set
- Age, disability and health status, and more...

Representatively diverse leadership for retention



Photo by Christina @ wocintechchat.com on Unsplash
https://unsplash.com/@wocintechchat?utm_source=unsplash&utm_medium=referral&utm_content=creditCopyText

**Not lowering bar
———— extending it**

Great Minds Think Different

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 — entrenched ways of thinking that can otherwise blind them to key information and even lead them to make errors in decision-making processes.”

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>

Coalesce on Shared Set of Technology Ethics



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

Adopt Technology Ethics

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



Association for
Computing Machinery



Microsoft



<>
Montréal Declaration
Responsible AI_
</>

An initiative of Université de Montréal



**Diverse,
inclusive
leaders**

**Diverse,
Multi-
Disciplinary
Teams**

**Shared
Tech Ethics**



UX Framework

Implementing Ethics

Activate curiosity

UX research methods to activate curiosity:

- Abusability Testing
- “Black Mirror” Episodes (inspired by British dystopian sci-fi tv series of same name)
- Flip it to test it
- Implicit Association Test from Harvard

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>

How do we get there?



Trustable,
Ethical
Systems

Conversations for Understanding

UX Framework guides teams

Difficult Topics

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

*"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.

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 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://arxiv.org/abs/1910.03515>.

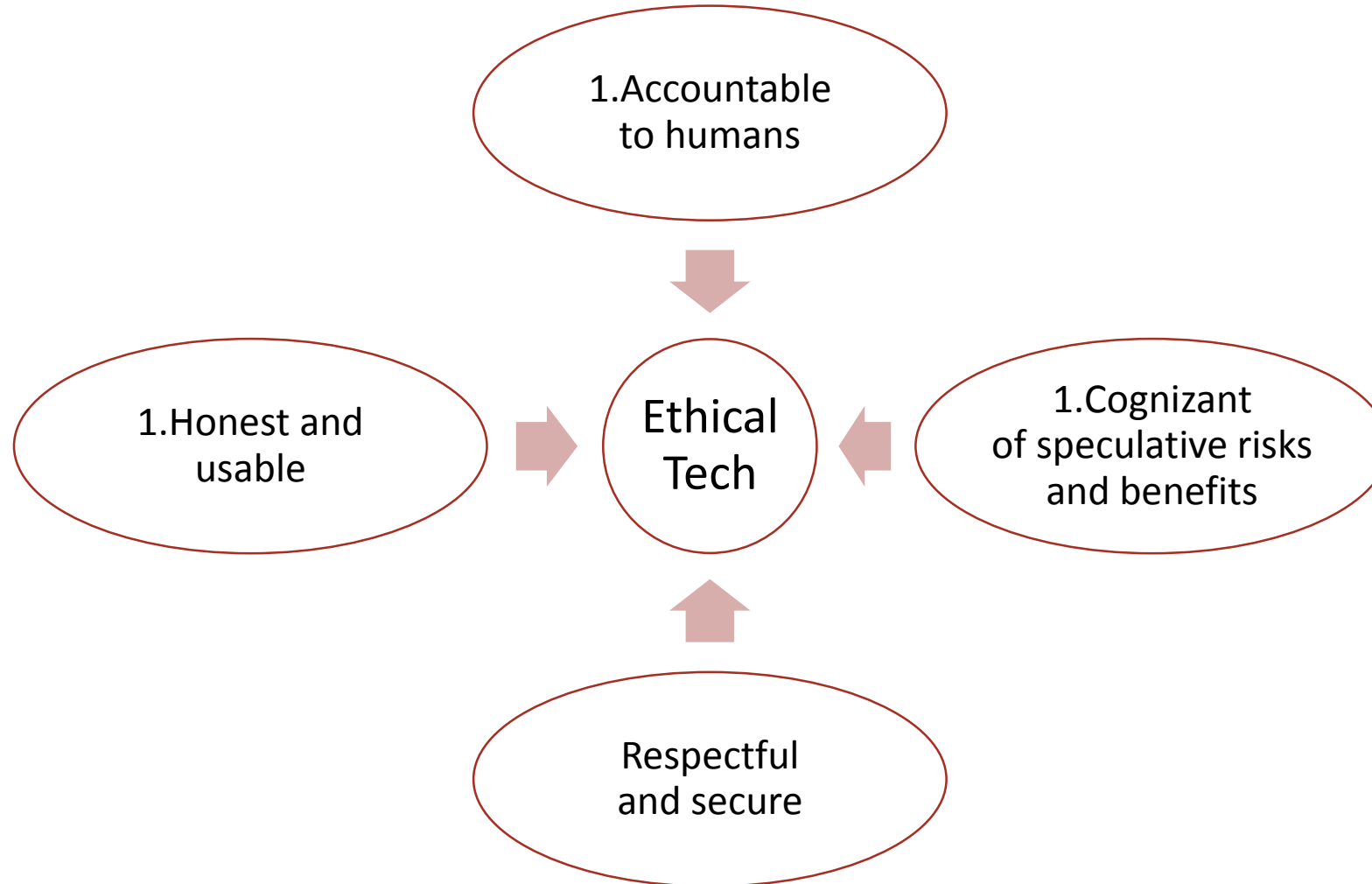
<p>We will design our AI system with the following in mind:</p> <ul style="list-style-type: none"><input type="checkbox"/> Designated humans have the ultimate responsibility for all decisions and outcomes:<ul style="list-style-type: none">• Responsibilities are explicitly defined between the AI system and human(s), 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.<input type="checkbox"/> Significant decisions made by the AI system will be<ul style="list-style-type: none">• explained• able to be overridden• appealable and reversible	<p>We work to speculatively identify the full range of risks and benefits:</p> <ul style="list-style-type: none"><input type="checkbox"/> Harmful, malicious use and consequences, as well as good, beneficial use and consequences.<input type="checkbox"/> We will be cognizant and exhaustively research unintended consequences. <p>We will create plans for the misuse/abuse of the AI system, including the following:</p> <ul style="list-style-type: none"><input type="checkbox"/> communication plans to share pertinent information with all affected people<input type="checkbox"/> mitigation plans for managing the identified speculative risks. <p>We value respect and security:</p> <ul style="list-style-type: none"><input type="checkbox"/> Incorporating our values of humanity, ethics, equity, fairness, accessibility, diversity, and inclusion<input type="checkbox"/> respecting privacy and data rights (Only necessary data will be collected)<input type="checkbox"/> providing understandable security methods<input type="checkbox"/> 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"><input type="checkbox"/> The purpose, limitations, and biases of the AI system are explained in plain language.<input type="checkbox"/> Data sources have unambiguous respected sources, and biases are known and explicitly stated.<input type="checkbox"/> Algorithms and models are appropriate and verifiable.<input type="checkbox"/> Confidence and context are presented for humans to base decisions on.<input type="checkbox"/> Transparent justification for recommendations and outcomes is provided.<input type="checkbox"/> Straightforward and interpretable monitoring systems are provided. <p>We value honesty and usability:</p> <ul style="list-style-type: none"><input type="checkbox"/> Humans can easily discern when they are interacting with the AI system vs. a human.<input type="checkbox"/> Humans can easily discern when and why the AI system is taking action and/or making decisions.<input type="checkbox"/> 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 non-profit research organization dedicated to advancing the state of the art in software engineering and systems research to benefit the public interest. For more information, contact the SEI at seinfo@sei.cmu.edu, www.sei.cmu.edu, or www.sei.cmu.edu.

Contact Us
Carnegie Mellon University
Software Engineering Institute
4800 Forbes Avenue, Pittsburgh, PA 15213-1502
412.263.1000
412.263.1000 | 800.351.3333
© 2019 SEI

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

SmartPackage - Scenario

Track online orders, shipping progress, and receipt.

Users: Consumers at home

Goals of SmartPackage:

- No worries - expected delivery is easy to track
- Alert when arrive and location via images
- Manages returns

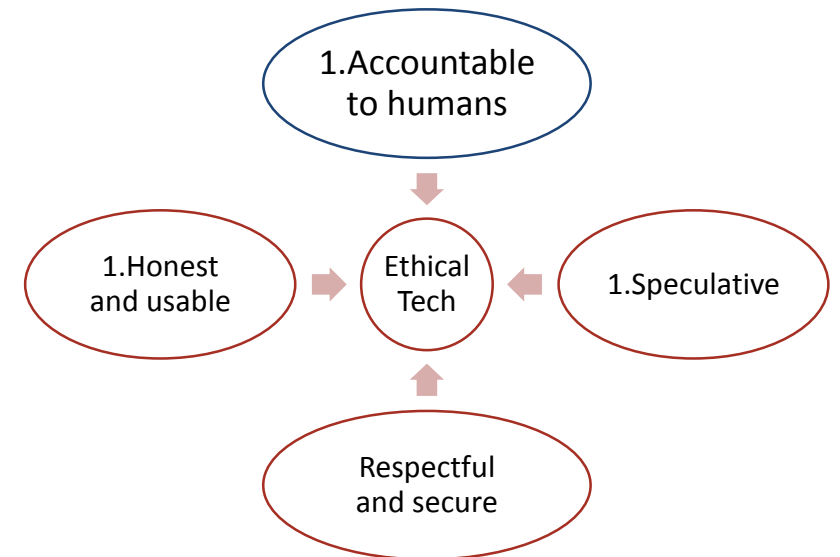
Accountable to Humans

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 system will be

- explained
- able to be overridden
- appealable and reversible

SmartPackage

- Ability to turn on and off notifications
- Ability to control camera content/capture

Responsibilities explicitly defined

Between system and human(s)

SmartPackage (System or Consumer?)

- Integrates new purchases?
- Integrates new vendors?
- Determines when to send alert?
- How many to keep available?

Abusability Testing

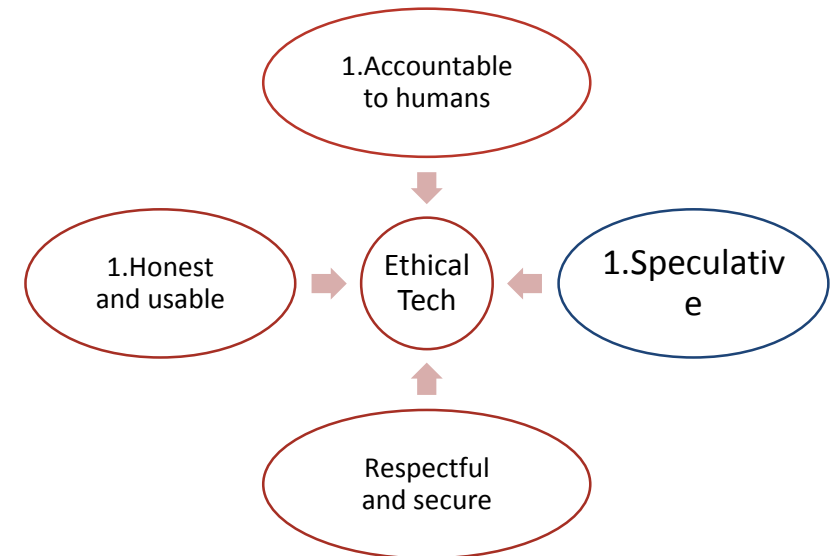
Feature added to enable SmartPackage to report delivery issues

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

Cognizant of Speculative Risks and Benefits

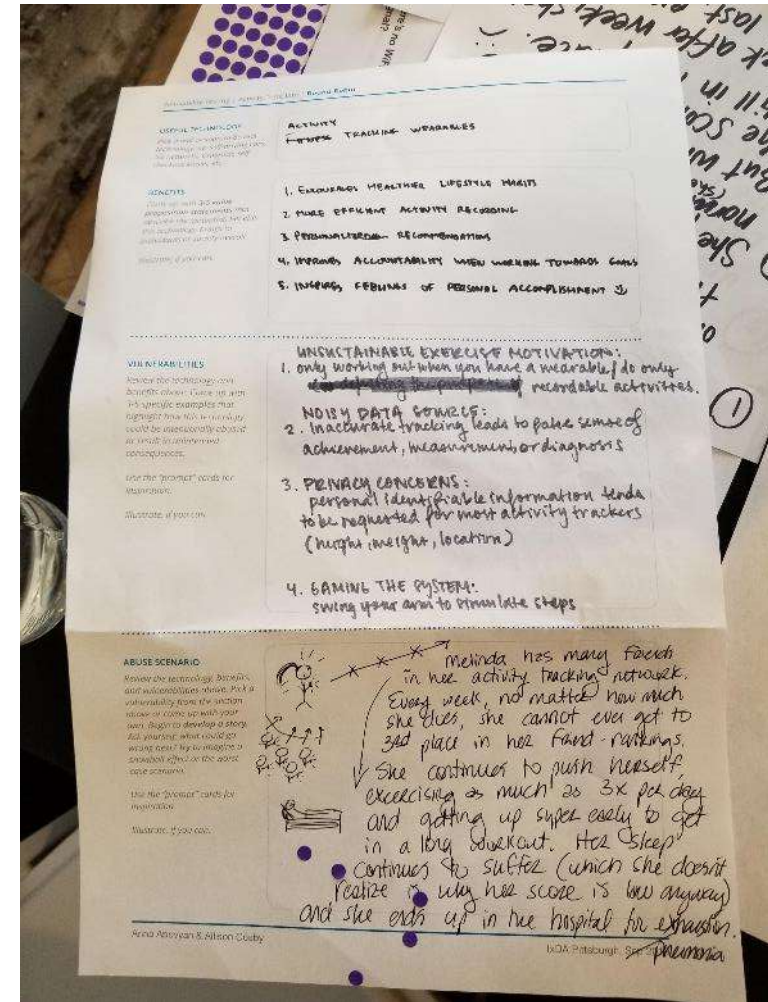
Identify full range of

- Harmful, malicious use, as well as good, beneficial use
- Blind spots and 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

- SmartPackage was designed for individual use.
- Households with multiple people receiving packages, find that SmartPackage starts to refuse deliveries meant for other household members
- SmartPackage reports them as errors to the shipping carriers, creating further frustration and grief

Speculative: Create communication & mitigation plans

Plan for unwanted consequences

Misuse and abuse of system

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

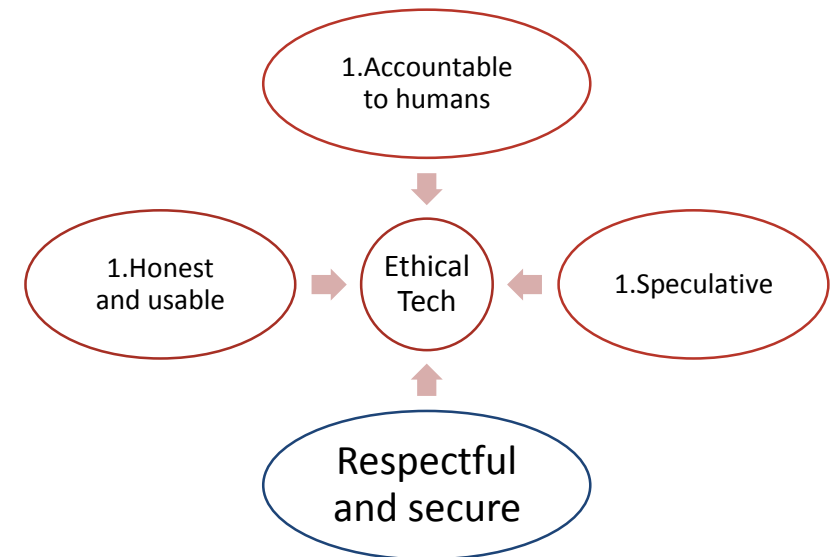
Respectful and Secure

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

Respect privacy and data rights

Make system robust, valid and reliable

Provide understandable security



Respectful and Secure

SmartPackage

- Who has access to shipments and contents?
- Who has access to images of shipments and address?
- How is that information used?
- How is PII* of consumers protected?

*PII is Personally Identifiable Information (name, address, etc.)

Honest and Usable

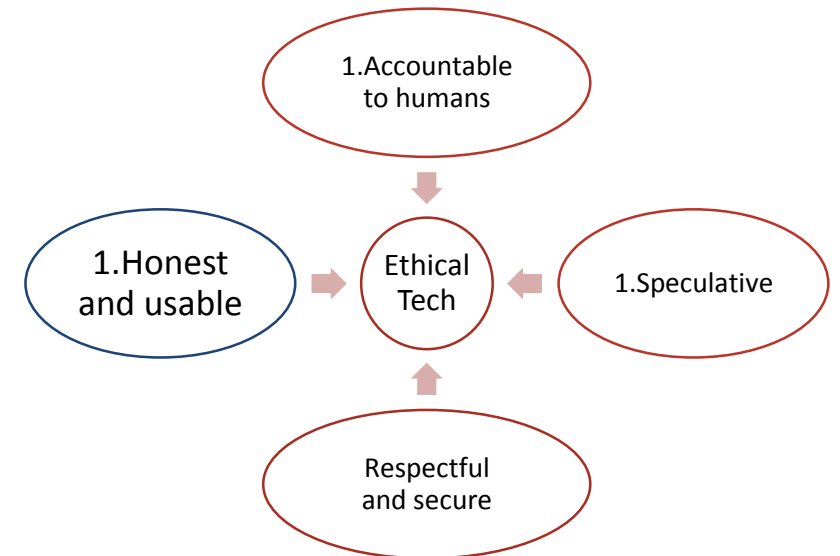
Value transparency with the goal of engendering trust

Smart speakers explicitly state identity as an AI system

Remove unwanted bias in data

Show awareness of known and desirable bias

Acknowledge issues



Reward team members for finding ethics bugs

Dr. Ayanna Howard

- on the Artificial Intelligence Podcast with Lex Fridman



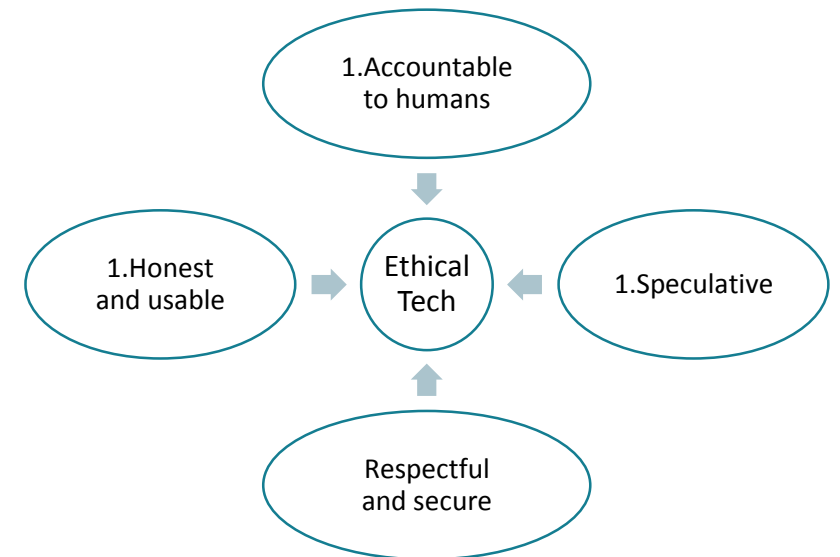
We aren't perfect, tech won't be perfect

Empower diverse teams, inclusive environments

Adopt technical ethics

Encourage deep conversations (Checklist)

Activate curiosity; be speculative; imaginative



**Evangelize
for human values**

**Ethical.
Transparent. Fair.**

Carol J. Smith

Twitter: @carologic

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

CMU's Software Engineering Institute,
Emerging Technology Center

Twitter: @sei_etc