# Human-Centered AI:
# Reliable, Safe & Trustworthy
# Part 1

Ben Shneiderman   **@benbendc**

Founding Director (1983-2000), Human-Computer Interaction Lab
Professor, Department of Computer Science

Member, National Academy of Engineering

UNIVERSITY OF MARYLAND

PETER WALL
INSTITUTE FOR ADVANCED STUDIES
THE UNIVERSITY OF BRITISH COLUMBIA VANCOUVER

Photo: BK Adams

Interdisciplinary research community
- Computer Science & Info Studies
- Psych, Socio, Educ, Jour & MITH

hcil.umd.edu
vimeo.com/72440805

# *Designing the User Interface*

## Design Theories

Direct manipulation

Menus, speech, search

Social Media

Information Visualization
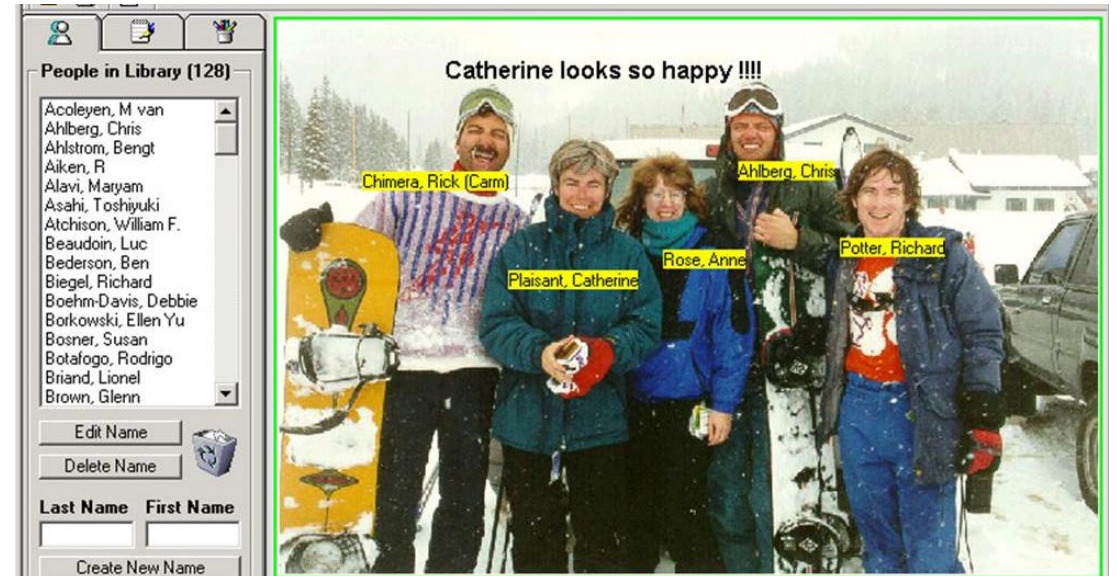
`www.cs.umd.edu/hcil/DTUI6`
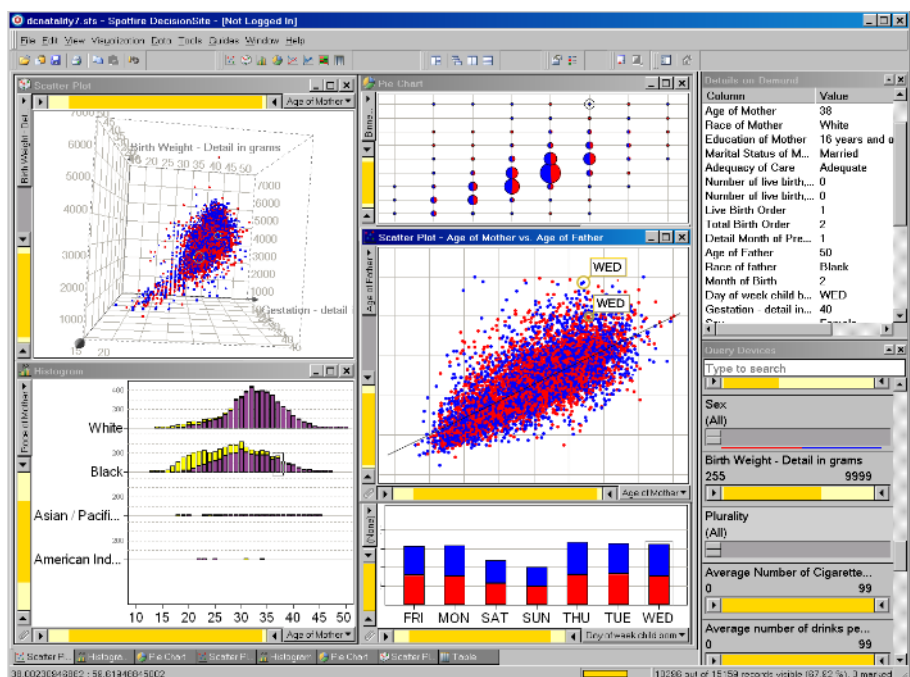
**Sixth Edition: 2016**

## Web links

The **University of Maryland, College Park** (often referred to as the **University of Maryland**, **Maryland**, **UM**, **UMD**, **UMCP**, or **College Park**) is a public research university[10] located in the city of College Park in Prince George's County, Maryland, approximately 4 miles (6.4 km) from the northeast border of Washington, D.C. Founded in 1856, the university is the flagship institution of the University System of Maryland. With a fall 2010 enrollment of more than 37,000 students, over 100 undergraduate majors, and 120 graduate programs,
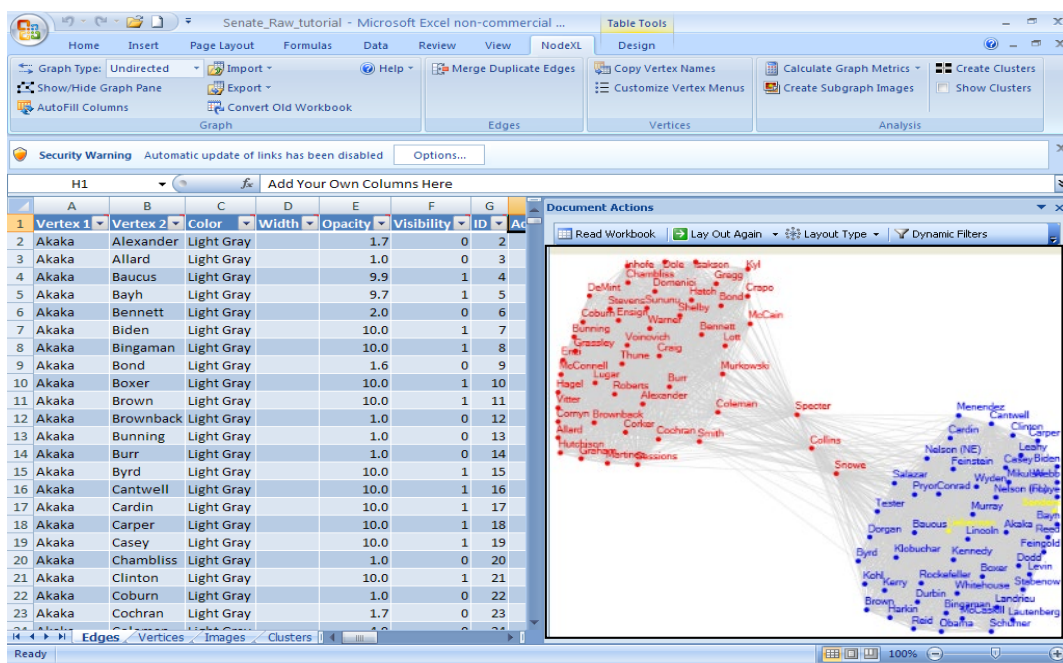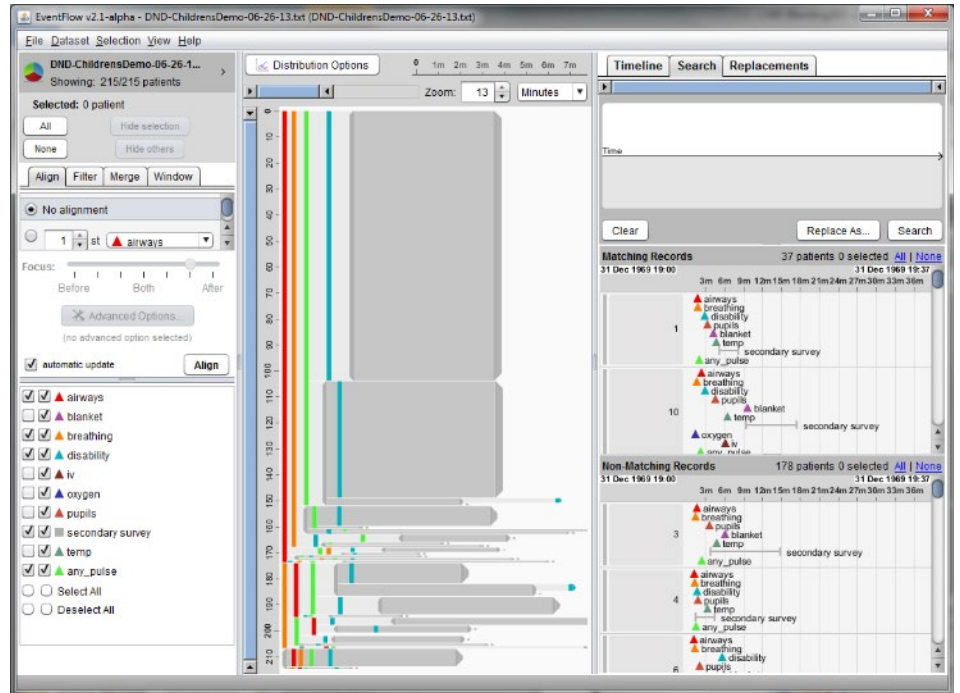
## Tiny touchscreen keyboards



## Photo tagging

Spotfire

Treemaps
FinViz

NodeXL

EventFlow

# The Goal of Visualization is Insight, Not Pictures

# The Goal of Visualization is Insight, Not Pictures

## Information Visualization Mantra

- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand

(IEEE, *VLC*,1996)

MARTIN FORD

# RISE OF THE ROBOTS

TECHNOLOGY AND THE THREAT OF A JOBLESS FUTURE

---

NEW AFTERWORD

## NICK BOSTROM

# SUPERINTELLIGENCE

### Paths, Dangers, Strategies

'I highly recommend this book'
BILL GATES

TIMES BESTSELLER

---

ARTIFICIAL INTELLIGENCE AND THE END OF THE HUMAN ERA

# OUR FINAL INVENTION

## JAMES BARRAT

WEAPONS OF MATH DESTRUCTION

HOW BIG DATA INCREASES INEQUALITY AND THREATENS DEMOCRACY

CATHY O'NEIL

HARRY COLLINS

ARTIFICTIONAL INTELLIGENCE

Against Humanity's Surrender to Computers

THE AI DELUSION

GARY SMITH

**REBOOTING AI**
Building Artificial
Intelligence We Can Trust

**GARY MARCUS and ERNEST DAVIS**

"so much of what we read about AI strikes us as pure fantasy, predicated on a confidence in AI's imagined strengths that bears no relation to current technological capabilities"

# What is Human-Centered AI?

# What is Human-Centered AI?



## Amplify, Augment, Empower & Enhance People

# Human-Centered AI

**Human Values**
Rights, Justice & Dignity

# Human-Centered AI

**Human Values**
Rights, Justice & Dignity

**Individual Goals**
Self-efficacy, Creativity, Responsibility & Social Connections

# Human-Centered AI

**Human Values**
Rights, Justice & Dignity

**Individual Goals**
Self-efficacy, Creativity, Responsibility & Social Connections

**Design Aspirations**
Reliable, Safe & Trustworthy
Team, Organization, Industry & Government

# Human-Centered AI

**Stakeholders**

- Researchers
- Developers
- Business Leaders
- Policy Makers
- Users

**Human Values**
Rights, Justice & Dignity

**Individual Goals**
Self-efficacy, Creativity, Responsibility & Social Connections

**Design Aspirations**
Reliable, Safe & Trustworthy
Team, Organization, Industry & Government

# Human-Centered AI

**Stakeholders**
- Researchers
- Developers
- Business Leaders
- Policy Makers
- Users

**Human Values**
Rights, Justice & Dignity

**Individual Goals**
Self-efficacy, Creativity, Responsibility & Social Connections

**Design Aspirations**
Reliable, Safe & Trustworthy
Team, Organization, Industry & Government

**Threats**
- Malicious Actors
- Bias
- Flawed Software

# Human-Centered AI



**Stakeholders**
- Researchers
- Developers
- Business Leaders
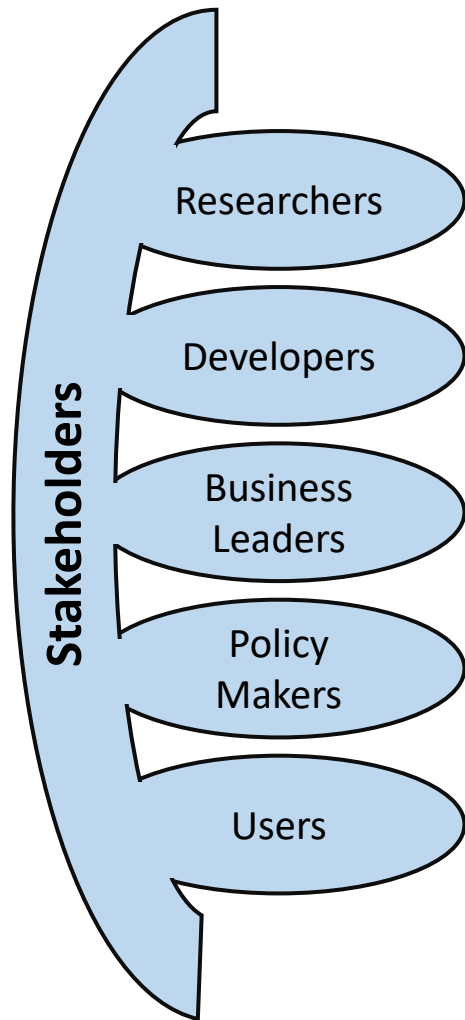- Policy Makers
- Users

**HCAI Framework**

**Human Values**
Rights, Justice & Dignity

**Individual Goals**
Self-efficacy, Creativity, Responsibility & Social Connections

**Design Aspirations**
Reliable, Safe & Trustworthy
Team, Organization, Industry & Government

**Threats**
- Malicious Actors
- Bias
- Flawed Software

# Human-Centered AI

# Human-Centered AI



**Stakeholders**
- Researchers
- Developers
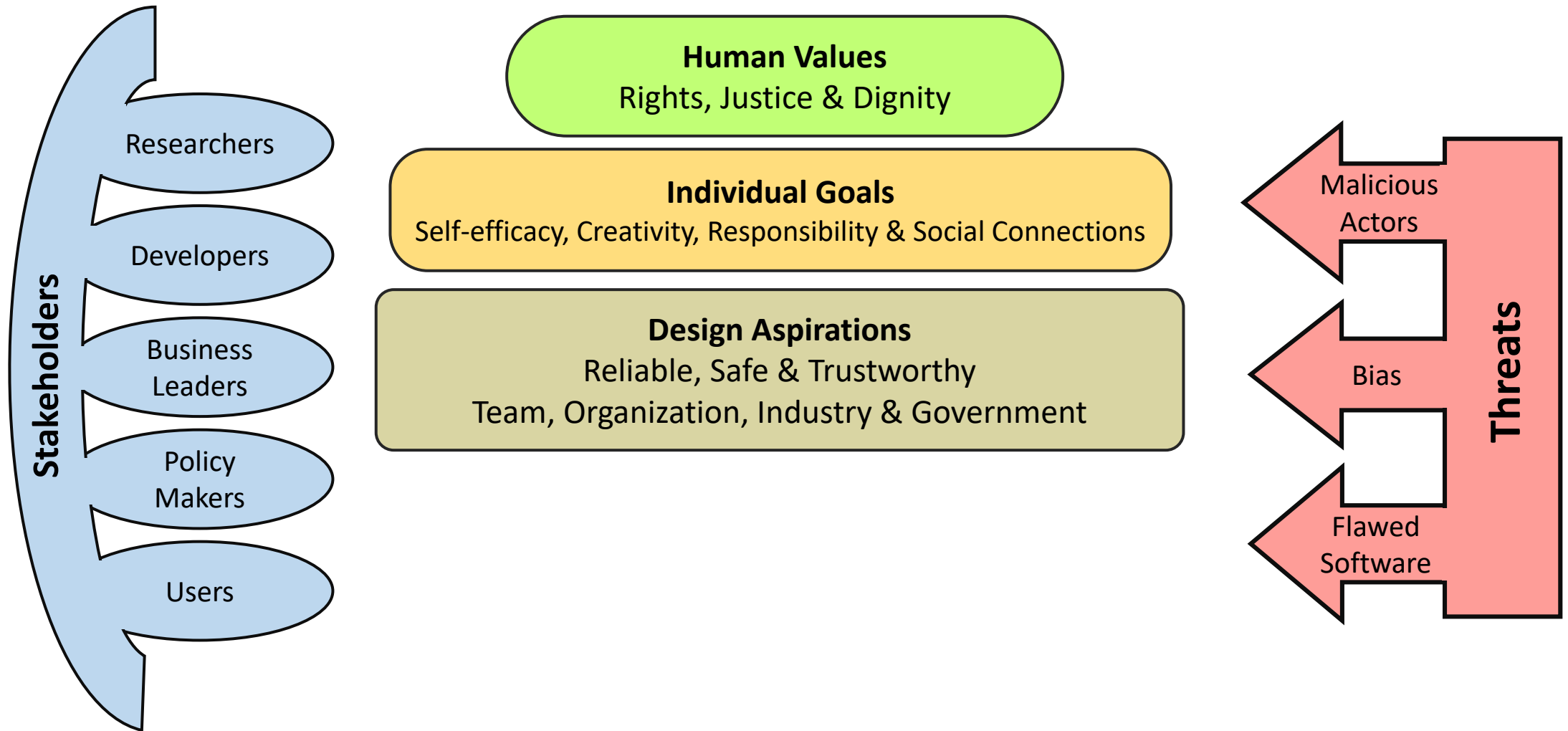- Business Leaders
- Policy Makers
- Users

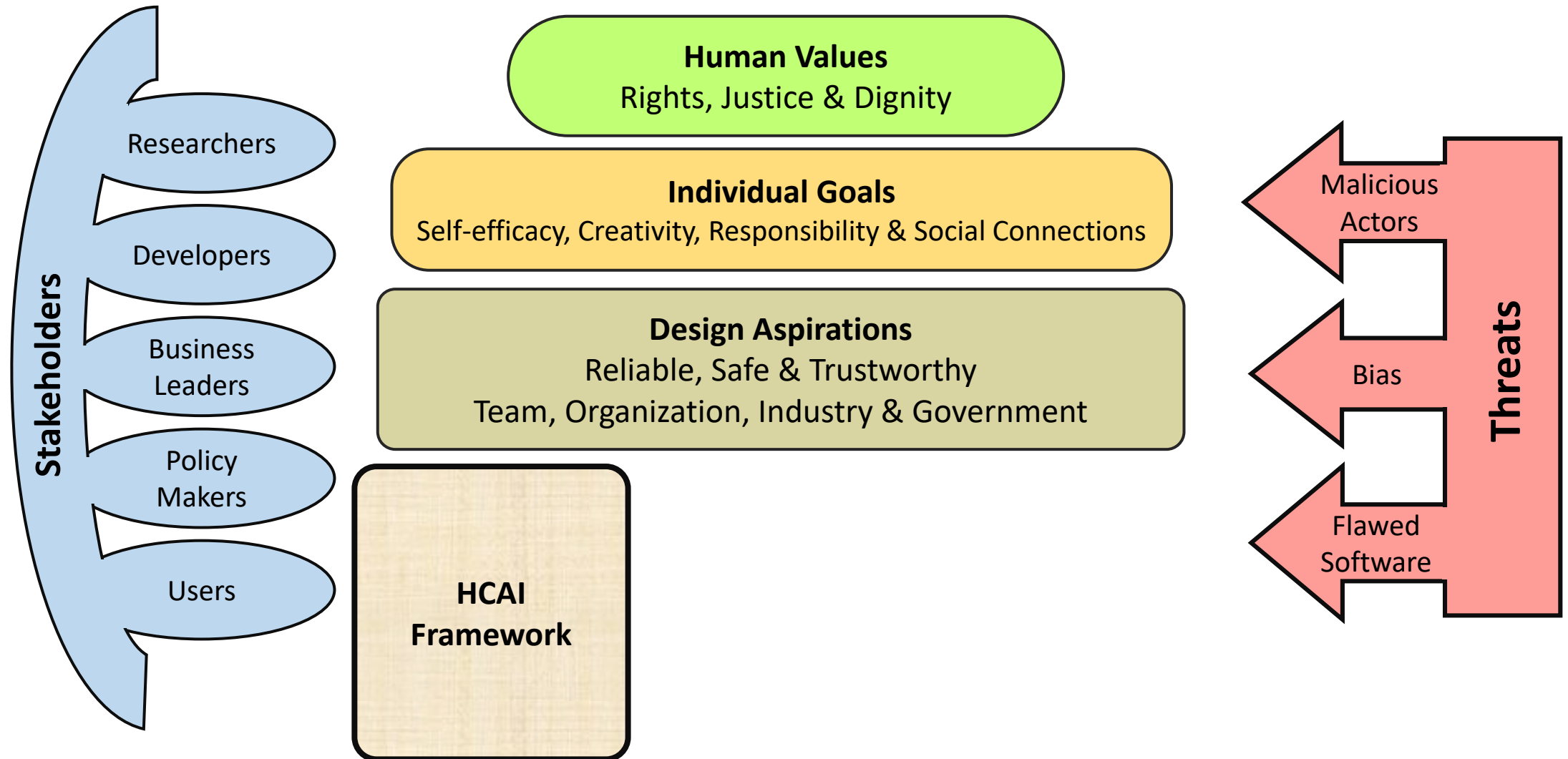**Human Values**
Rights, Justice & Dignity

**Individual Goals**
Self-efficacy, Creativity, Responsibility & Social Connections

**Design Aspirations**
Reliable, Safe & Trustworthy
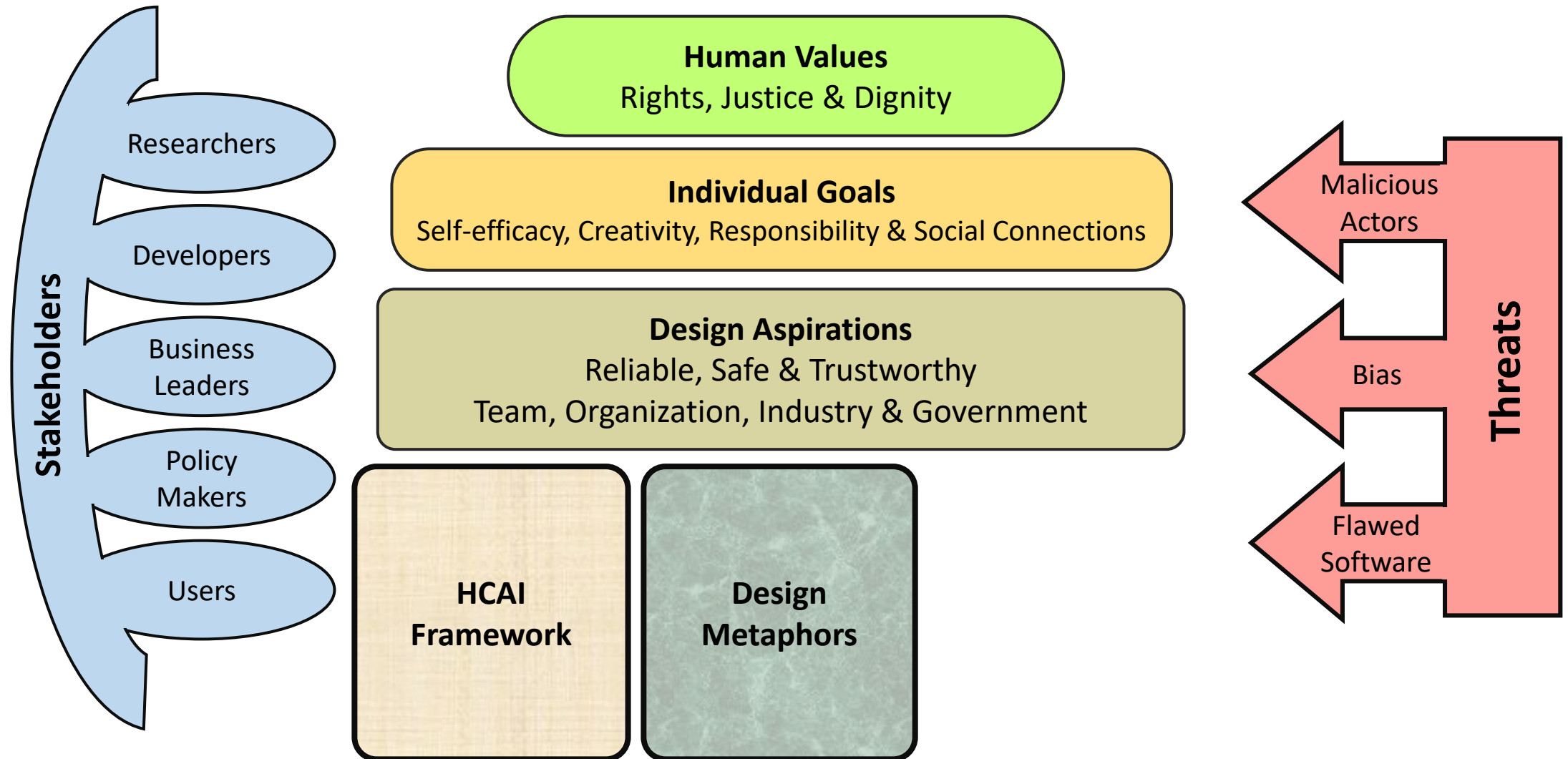Team, Organization, Industry & Government

**HCAI Framework**

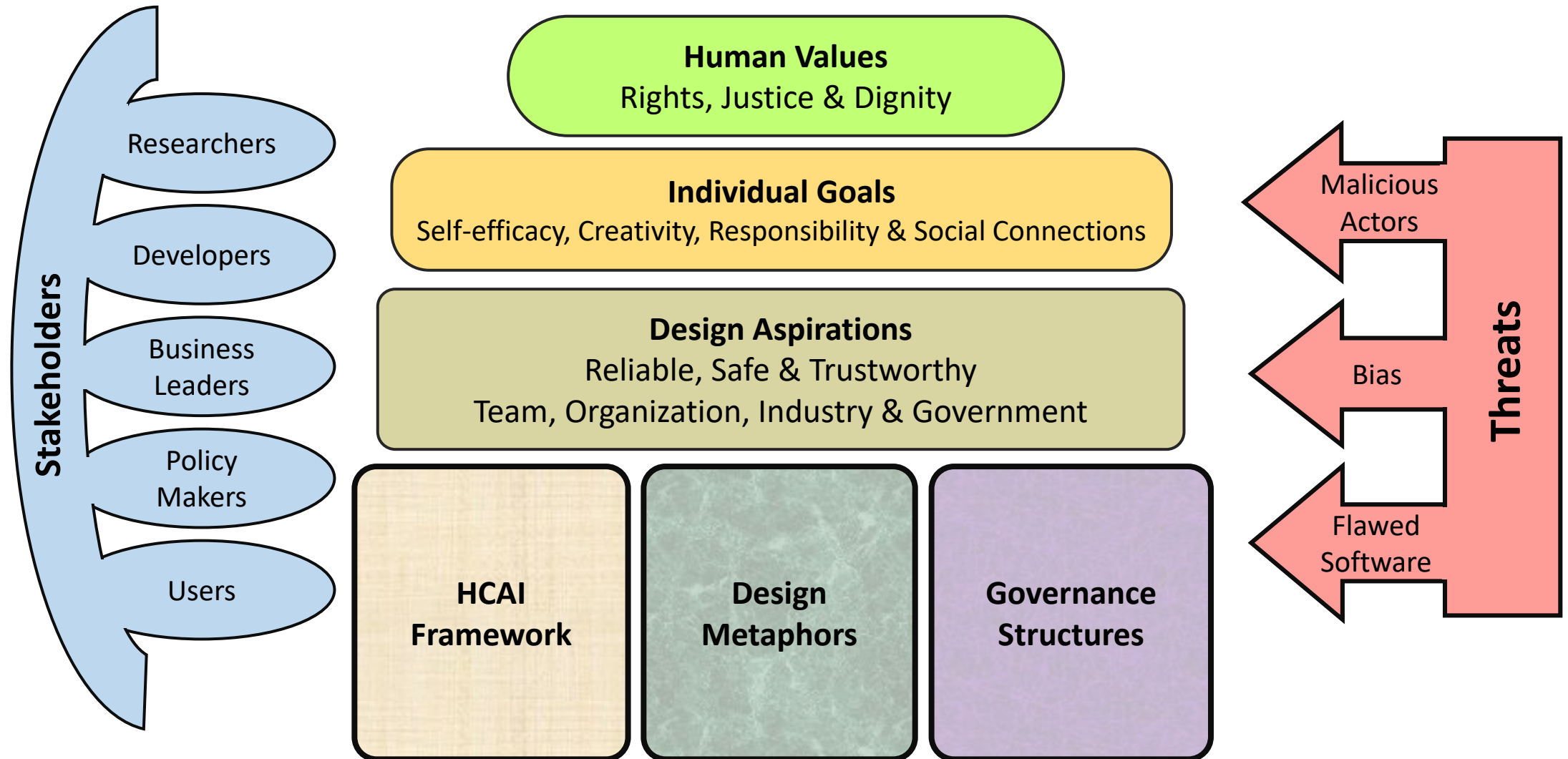**Design Metaphors**

**Governance Structures**

**Threats**
- Malicious Actors
- Bias
- Flawed Software

Oxford University Press (Early 2022)  `https://hcil.umd.edu/human-centered-ai/`

# UN Sustainable Development Goals



(https://sdgs.un.org/goals)

# People are not computers; Computers are not people

# HCAI Attributes

PREDICTABLE
PRIVATE
SECURE
UNBIASED

CONTROLLABLE
FAIR

DIRECTABLE

RELIABLE SAFE

TRUSTWORTHY

USABLE
AGILE

OBSERVABLE

COMPREHENSIBLE
EXPLAINABLE
DEPENDABLE
ROBUST

TRANSPARENT
INTERPRETABLE
RESILIENT

AVAILABLE

TRACEABLE
TRACKABLE

AUDITABLE

## HCAI Attributes that Are Candidates for Assessment

**General virtues of the system itself**
- **Trustworthy**: Can users trust the system to perform correctly?
- **Responsible/Humane**: Has the system been designed, developed, and tested in a responsible way?
- **Ethical Design**: Were stakeholders involved in the design?
- **Ethical Data**: Was the data collected in an ethical manner?
- **Ethical Use**: Will the system's outcome be used in an ethical manner?
- **Well-being/Benevolence**: Does the system support human health, comfort, and values?
- **Secure**: How vulnerable is the system to attack?
- **Private**: Does the system protect a person's identity and data?

**Performs well in practice**
- **Robust/Agile**: Does the system perform well when inputs change?
- **Reliable/Dependable**: Does the system do the right thing?
- **Available**: Is the system running when needed?
- **Resilient/Adaptive**: Can the system recover from disruptions?
- **Testable/Verifiable/Validatable/Certifiable:** Can be tested to verify adherence to requirements?
- **Safe**: Does the system have a history of safe use?

### Clarity to stakeholders

- **Accurate**: Does the system deliver correct results on test cases and real world cases?
- **Fair/Unbiased**: Are the system's biases understood and reported?
- **Accountable/Liable**: Who or what is responsible for the system's outcome?
- **Transparent**: Is it clear to an external observer how the system's outcome was produced?
- **Interpretable/Explainable/Intelligible/Explicable:** Can the system explain the outcome?
- **Usable**: Can a human use it easily?

### Enables independent oversight

- **Auditable**: Can the system be audited by others for retrospective forensic analysis of failures?
- **Trackable:** Does the system display status and next steps so human intervention is possible?
- **Traceable:** Is the system designed to allow tracing back from an outcome to the root cause?
- **Redressable**: Is there a process for those harmed to request review and compensation?
- **Insurable**: Does the design permit insurance companies to offer policies?
- **Recorded**: Does the system record activity for retrospective forensic review?
- **Open**: Is code and data publicly available for others to review?
- **Certifiable**: Can it be certified and approved for use?

### Complies with accepted practices

- **Compliant with standards**: Does the system comply with relevant standards, e.g. IEEE P7000 series?
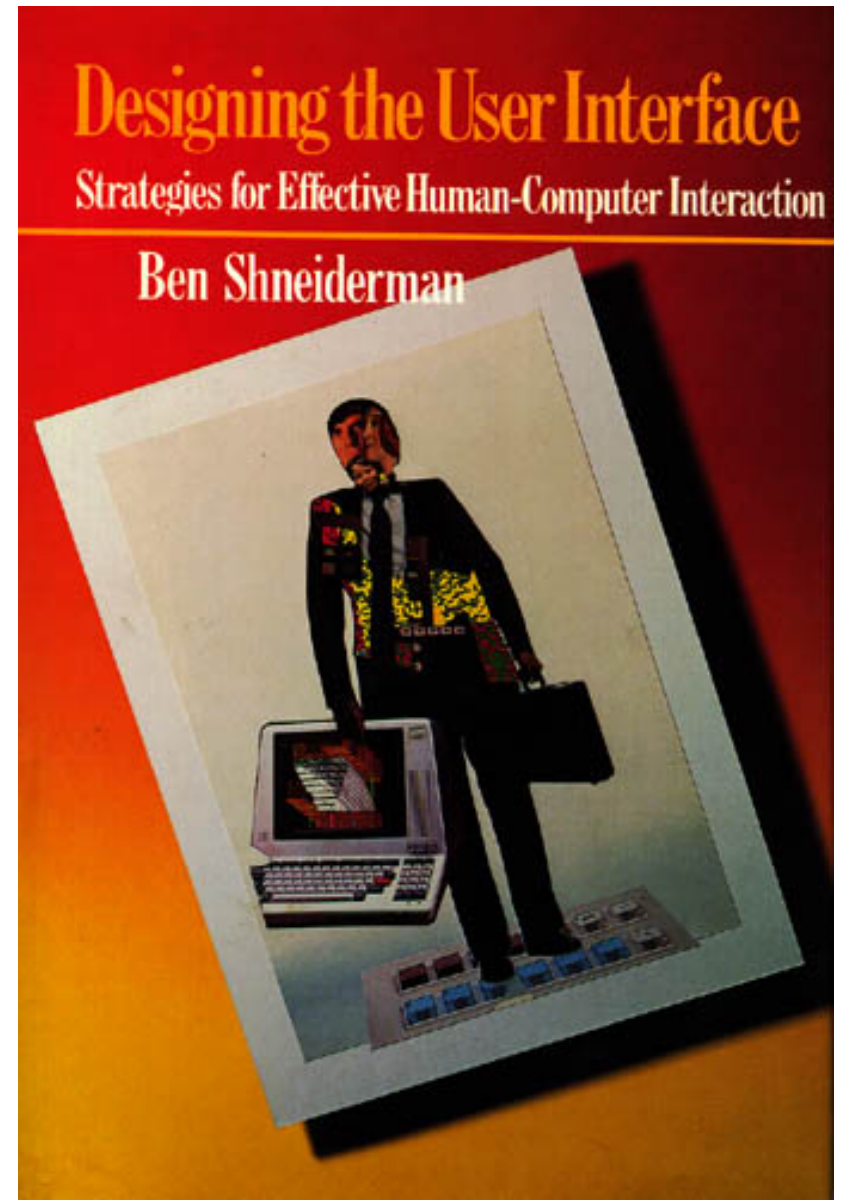- **Compliant with accepted software engineering workflows:** Was a trusted process used?

# HCAI Framework
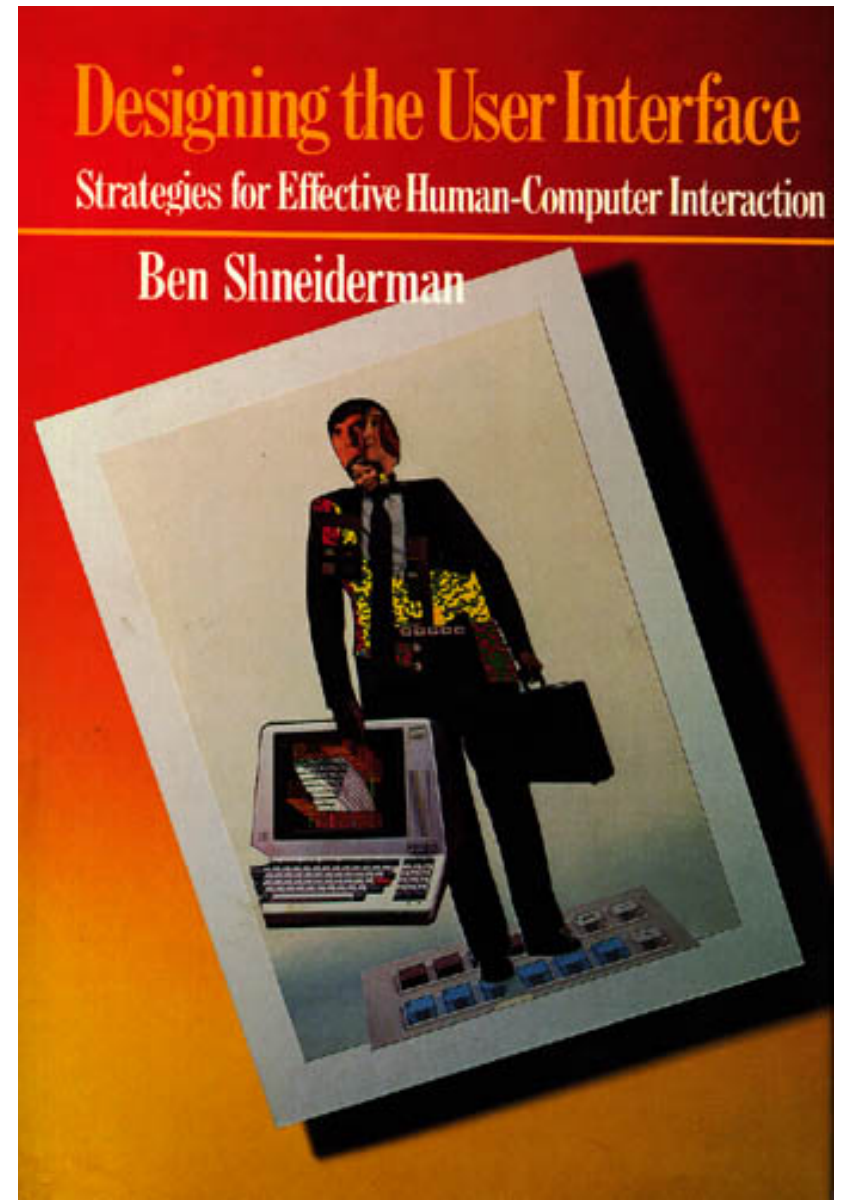
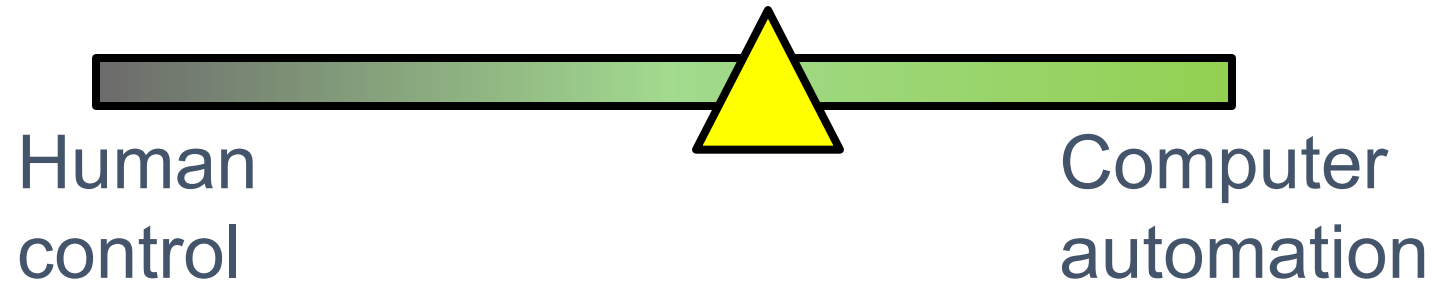# *Designing the User Interface*

Balancing automation & human control



**First Edition: 1986**

# *Designing the User Interface*

Balancing automation & human control



Human control

Computer automation



**First Edition: 1986**

# LEVELS OF DRIVING AUTOMATION



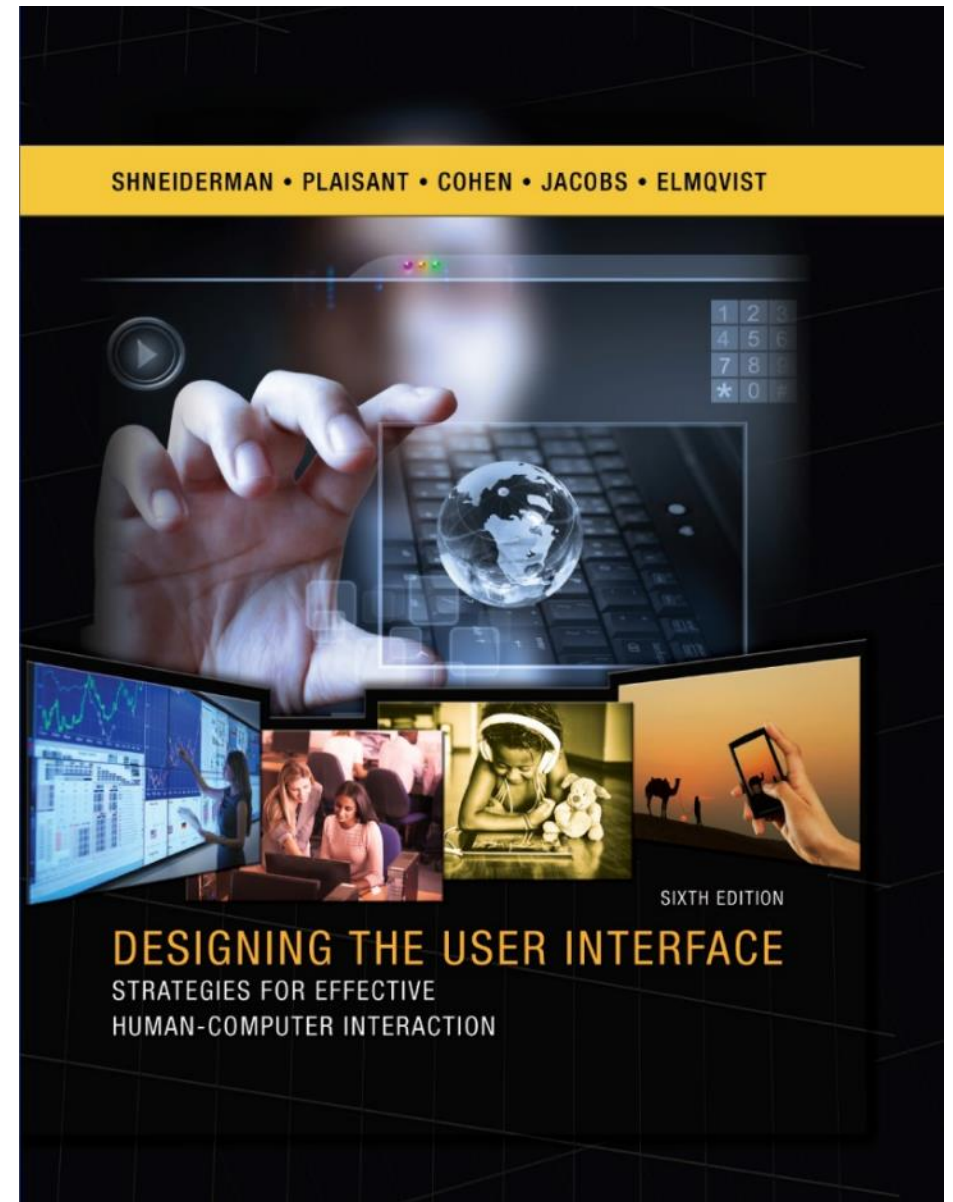| 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| **NO AUTOMATION** | **DRIVER ASSISTANCE** | **PARTIAL AUTOMATION** | **CONDITIONAL AUTOMATION** | **HIGH AUTOMATION** | **FULL AUTOMATION** |
| Manual control. The human performs all driving tasks (steering, acceleration, braking, etc.). | The vehicle features a single automated system (e.g. it monitors speed through cruise control). | ADAS. The vehicle can perform steering and acceleration. The human still monitors all tasks and can take control at any time. | Environmental detection capabilities. The vehicle can perform most driving tasks, but human override is still required. | The vehicle performs all driving tasks under specific circumstances. Geofencing is required. Human override is still an option. | The vehicle performs all driving tasks under all conditions. Zero human attention or interaction is required. |

**THE HUMAN MONITORS THE DRIVING ENVIRONMENT** | **THE AUTOMATED SYSTEM MONITORS THE DRIVING ENVIRONMENT**

(Society of Automotive Engineers, 2016)

# *Designing the User Interface*

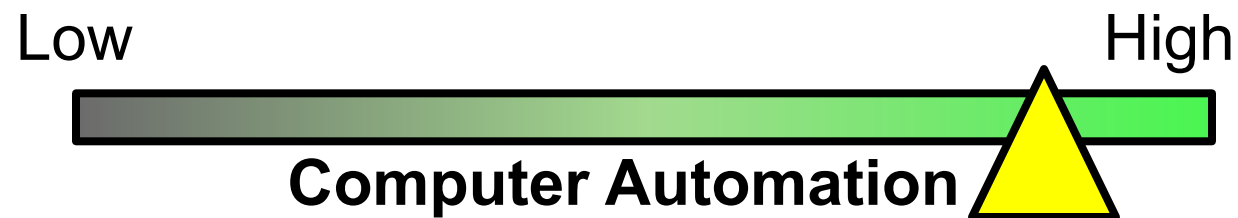Ensuring human control
while increasing automation



**Sixth Edition: 2016**
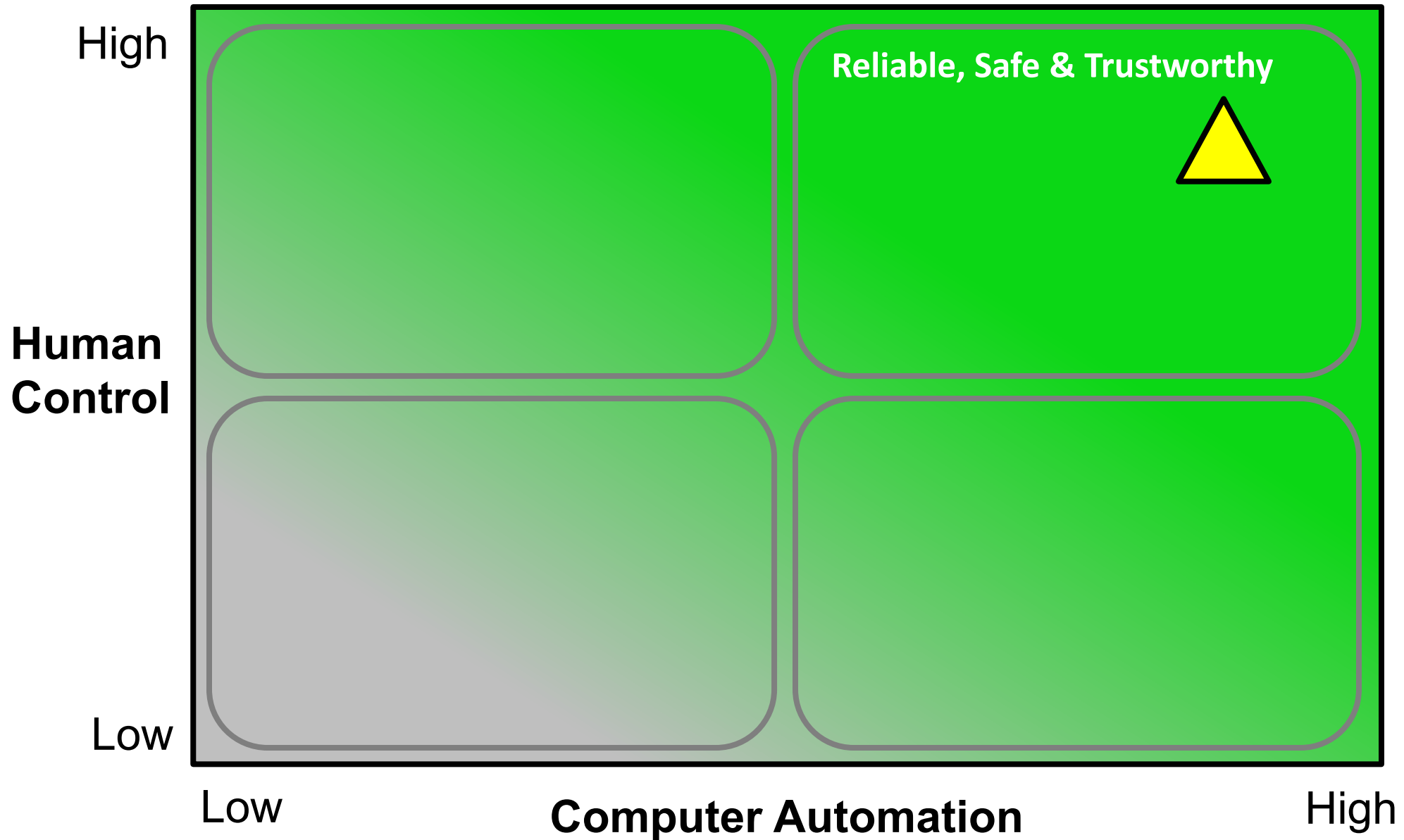
# Designing the User Interface

Ensuring human control
while increasing automation

Low                                    High

**Human Control**

Low                                    High

**Computer Automation**

SHNEIDERMAN • PLAISANT • COHEN • JACOBS • ELMQVIST

SIXTH EDITION

DESIGNING THE USER INTERFACE
STRATEGIES FOR EFFECTIVE
HUMAN-COMPUTER INTERACTION

**Sixth Edition: 2016**

# Human-Centered AI



**Reliable, Safe & Trustworthy**

**Human Control** (vertical axis: Low to High)

**Computer Automation** (horizontal axis: Low to High)

# Human-Centered AI



**Computer Automation**

# Human-Centered AI

| | | |
|---|---|---|
| **Human Mastery**<br>Bicycle<br>Piano | **Reliable, Safe & Trustworthy**<br><br>Elevator<br>Camera | |
| Music box<br>Landmine | Pacemaker<br>Airbag<br>**Computer Control** | |

High

Human Control

Low

Low    **Computer Automation**    High

# Human-Centered AI



Excessive Human Control

Excessive Automation

| | **Human Mastery** | **Reliable, Safe & Trustworthy** |
|---|---|---|
| **High** | Bicycle<br>Piano | Elevator<br>Camera |
| **Human Control** | Music box<br>Landmine | Pacemaker<br>Airbag |
| **Low** | | **Computer Control** |

Low — Computer Automation — High

Pain Control Designs

Excessive Human Control

High — Human Mastery | Reliable, Safe & Trustworthy

Human Control

Morphine drip bag | Computer Control

Low

Low — Computer Automation — High

Excessive Automation

# Pain Control Designs



Excessive Human Control

Excessive Automation

**Human Control**

High

Low

| | |
|---|---|
| **Human Mastery** | **Reliable, Safe & Trustworthy** |
| Morphine drip bag | Automatic dispenser |

**Computer Control**

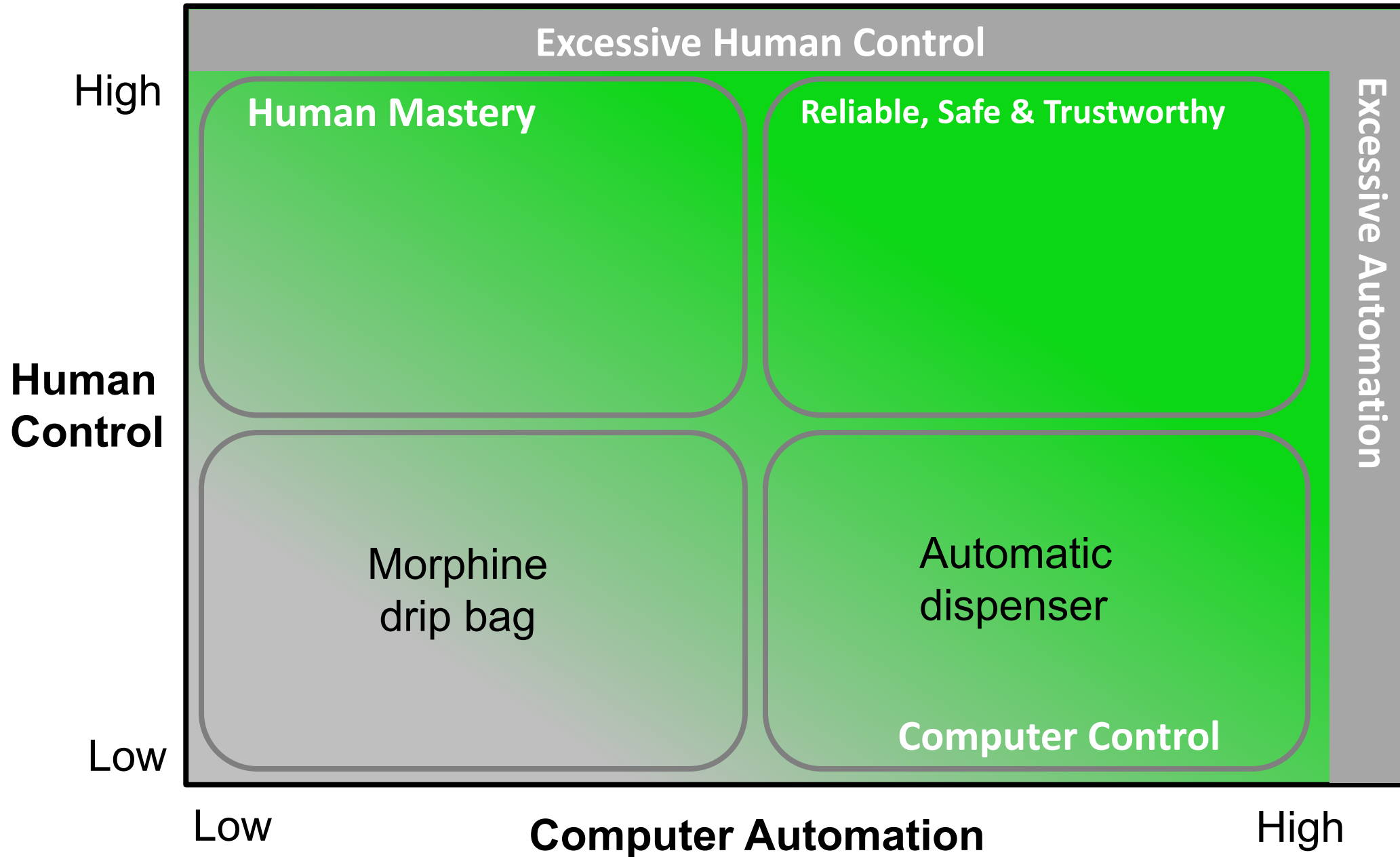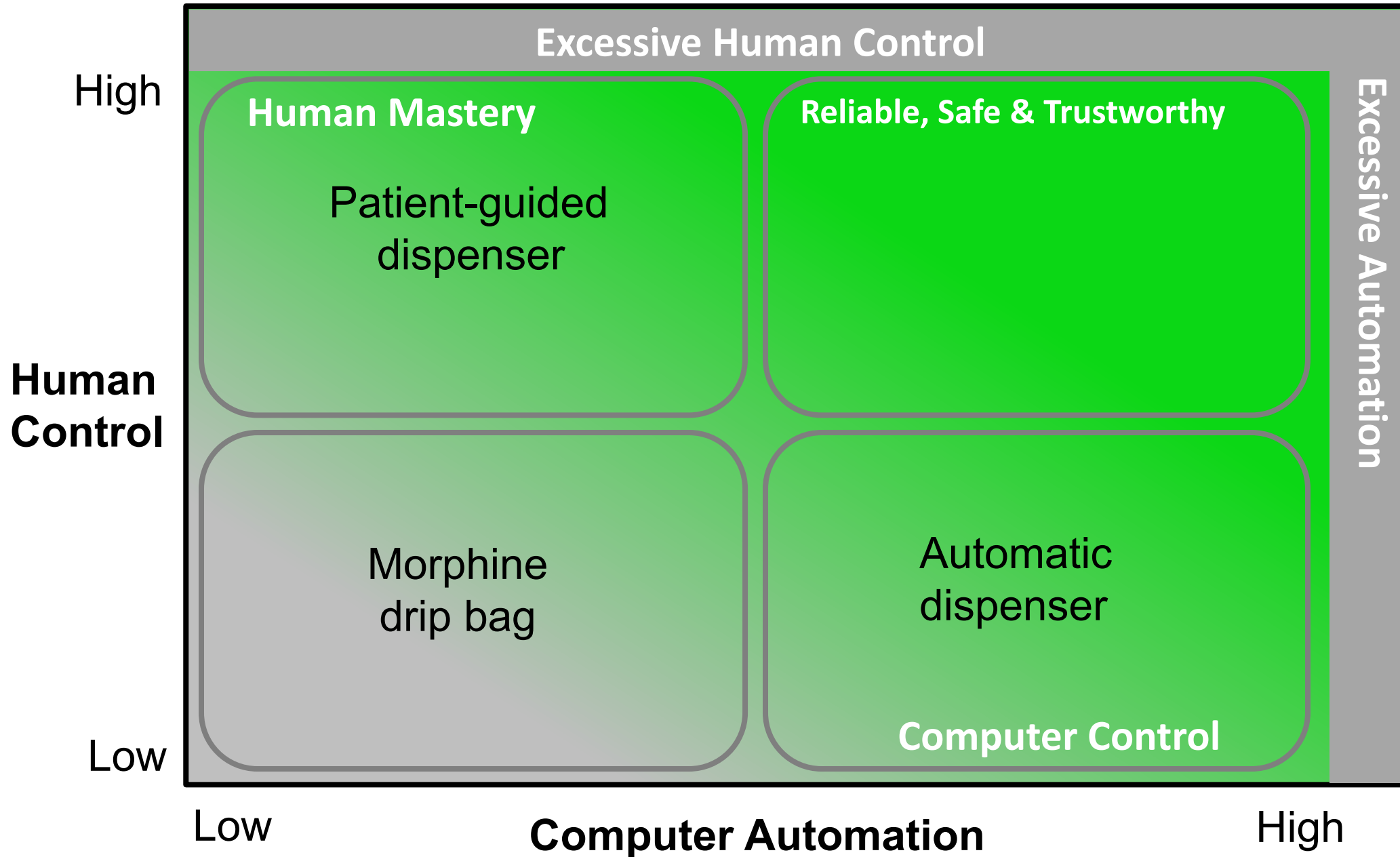Low — Computer Automation — High

# Pain Control Designs

# Pain Control Designs



Pain Control Designs quadrant chart with axes "Computer Automation" (Low to High) on the horizontal axis and "Human Control" (Low to High) on the vertical axis.

- **Excessive Human Control** (top border)
- **Excessive Automation** (right border)
- **Human Mastery** (top-left quadrant): Patient-guided dispenser
- **Reliable, Safe & Trustworthy** (top-right quadrant): Patient-guided & clinician-monitored system
- (bottom-left quadrant): Morphine drip bag
- **Computer Control** (bottom-right quadrant): Automatic dispenser
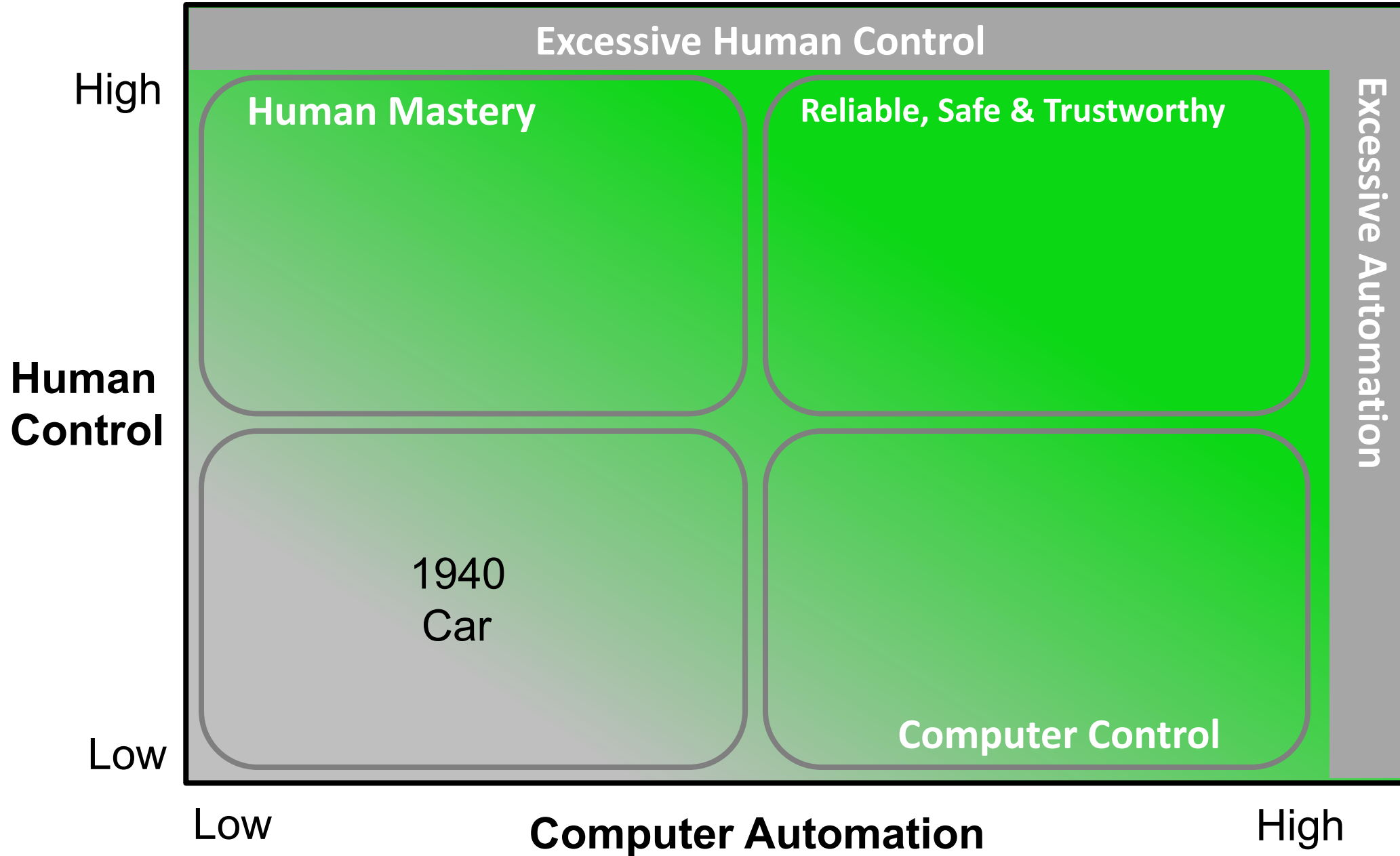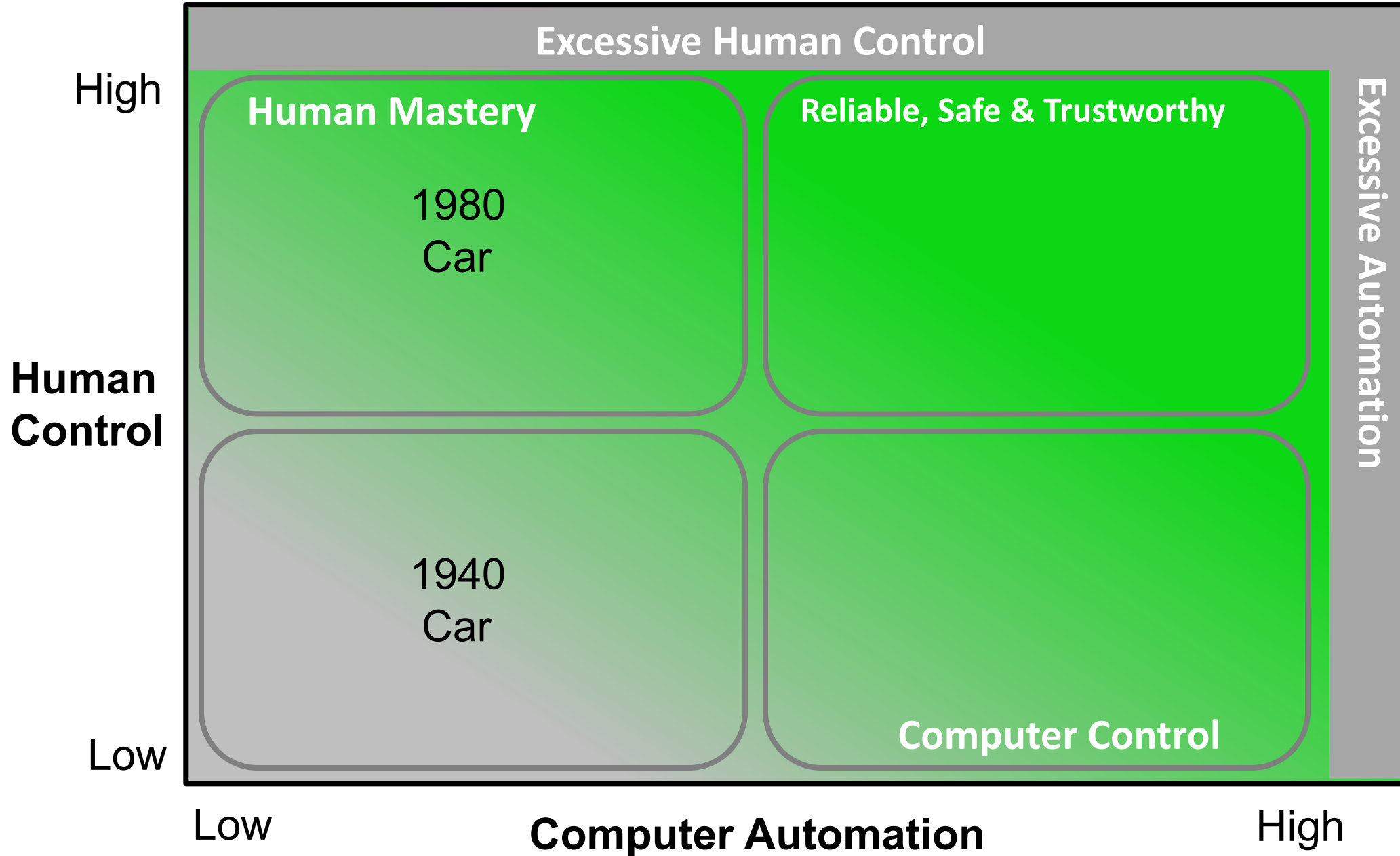
# Johns Hopkins University Hospital Control Center

# Car Control Designs



**Excessive Human Control**

| | |
|---|---|
| **Human Mastery** | **Reliable, Safe & Trustworthy** |
| 1940 Car | **Computer Control** |

High — Low (Human Control)

Low — High (Computer Automation)

**Excessive Automation**

**Human Control**

**Computer Automation**

# Car Control Designs



**Excessive Human Control**

**Excessive Automation**

**Human Control** — High / Low

**Computer Automation** — Low / High

**Human Mastery**

1980 Car

**Reliable, Safe & Trustworthy**

1940 Car

**Computer Control**

# Car Control Designs

# Car Control Designs

# Wheelchair Designs



**Excessive Human Control**

**Excessive Automation**

**High**

**Human Mastery**

Hand-powered
& user guided

**Reliable, Safe & Trustworthy**

Motorized, joystick
controlled, tele-operated
& programmable

**Human Control**

Push chair
(requires caretaker)

Robotic
(navigates to destination)

**Computer Control**

**Low**

**Low**　　**Computer Automation**　　**High**

# Application Categories

- **Recommenders**
- **Consequential**
- **Life-critical**

# Application Categories

- **Recommenders**
- **Consequential**
- **Life-critical**

- **Rapid performance**
- **Long duration**
- **Remote locations**

# Micro-Structure of Design

- **Automate where**
    - **+ Reliable performance is possible**
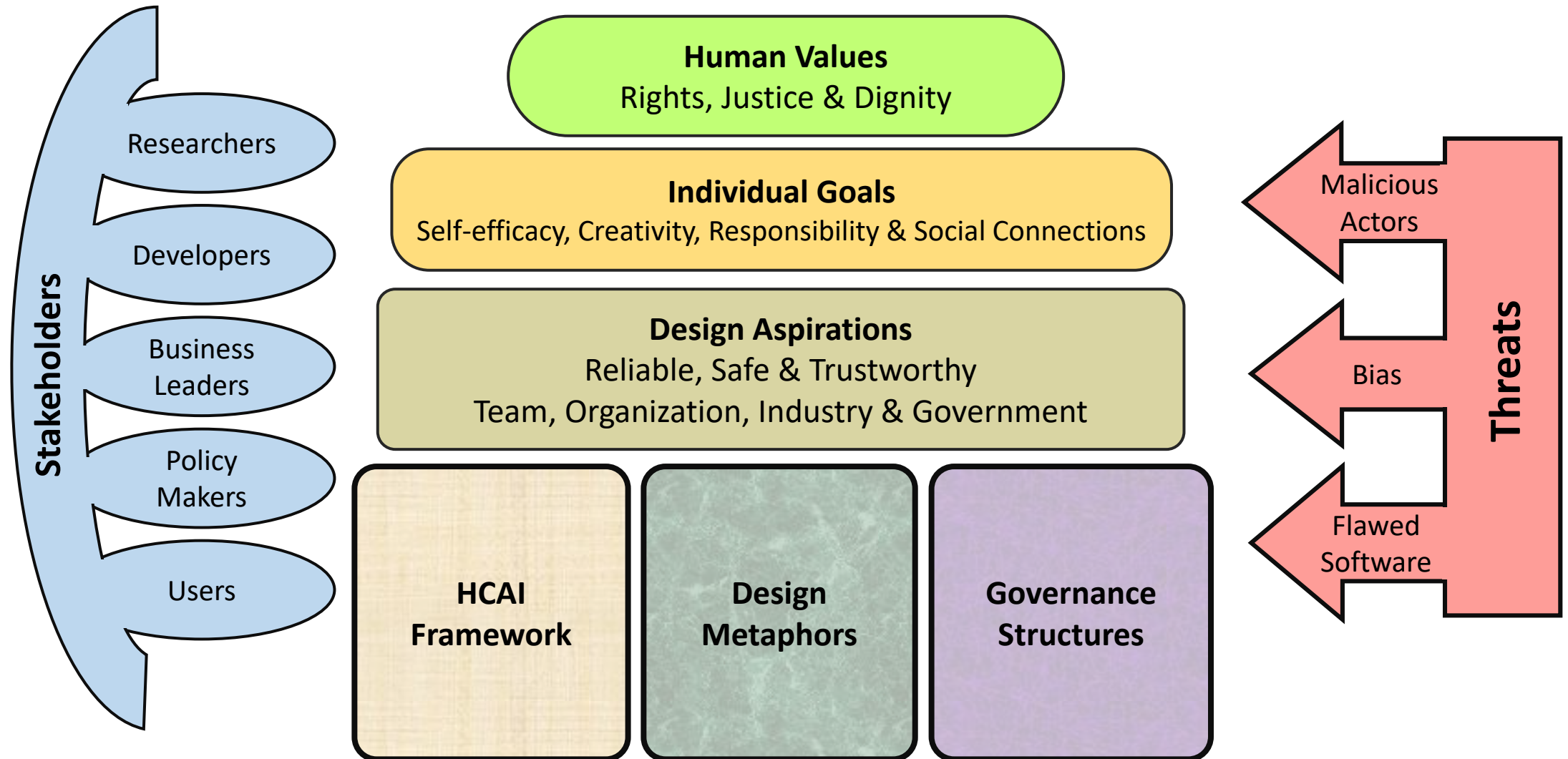    - **- But allow overrides**

# Micro-Structure of Design

- **Automate where**
  - **+ Reliable performance is possible**
  - **- But allow overrides**

- **Give human control where**
  - **+ Desired for creative flexibility**
  - **+ Automation is uncertain**
  - **- But prevent human errors**

# Micro-Structure of Design

- **Automate where**
    - **+ Reliable performance is possible**
    - **- But allow overrides**
- **Give human control where**
    - **+ Desired for creative flexibility**
    - **+ Automation is uncertain**
    - **- But prevent human errors**
- **Design supervisory control**
    - **+ Teleoperate remotely**
    - **+ Collect aggregate data**

# Human-Centered AI

# Design Metaphors

# Design Metaphors



**Science Goal**

**Innovation Goal**

**Combined Designs**

**Intelligent Agents**
Thinking Machine, Cognitive Actor, Artificial Intelligence, Knowledgeable

**Supertools**
Extend Abilities, Empower Users, Enhance Human Performance

**Teammates**
Co-active Collaborator, Colleague, Helpful Partner, Smart Co-worker

**Tele-operated Devices**
Steerable Instrument, Powerful Prosthetic, Boost Human Perceptual & Motor Skills

**Assured Autonomy**
Independent, Self-directed, Goal-setting, Self-monitored

**Supervised Autonomy**
Human Control & Oversight, Situation Awareness, Predictable Actions

**Social Robots**
Anthropomorphic, Humanoid, Android, Bionic, Bio-inspired

**Active Appliances**
Consumer-oriented, Wide Use, Low Cost Comprehensible Control Panels

# Teammate & Supertool

- **Social Teammate:**
  Since many people respond socially to robots,
  -> design robots to be human-like social teammates.

# Teammate & Supertool

- **Social Teammate:**
  Since many people respond socially to robots,
  -> design robots to be human-like social teammates.


- **Human-Centered Conjecture:**
  Since only humans can be responsible &
  computers have distinct capabilities (speed, storage, display...)
  -> design computers to be supertools

# Teammate & Supertool

- **Social Teammate:**
  Since many people respond socially to robots,
  -> design robots to be human-like social teammates.


- **Human-Centered Conjecture:**
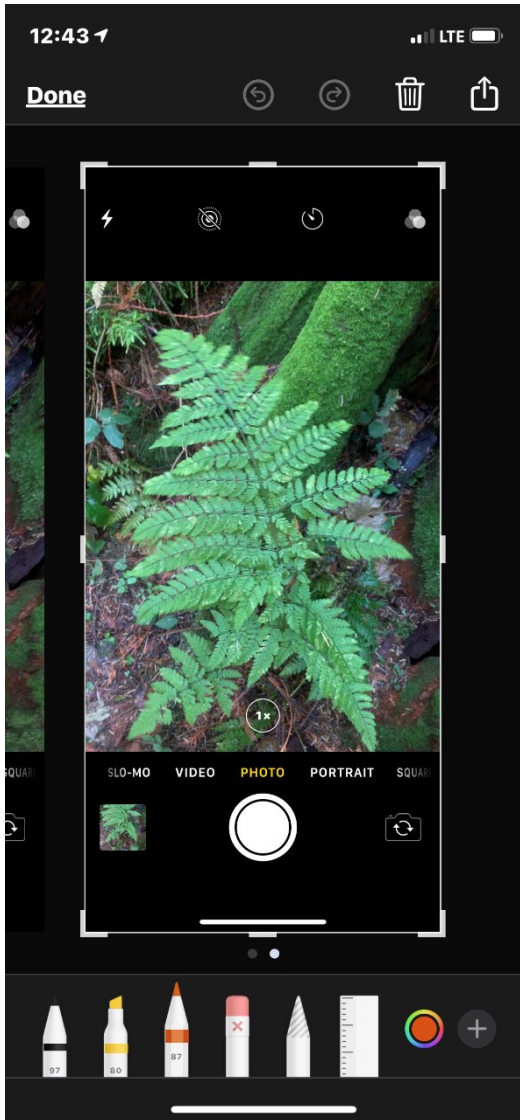  Since only humans can be responsible &
  computers have distinct capabilities (speed, storage, display...)
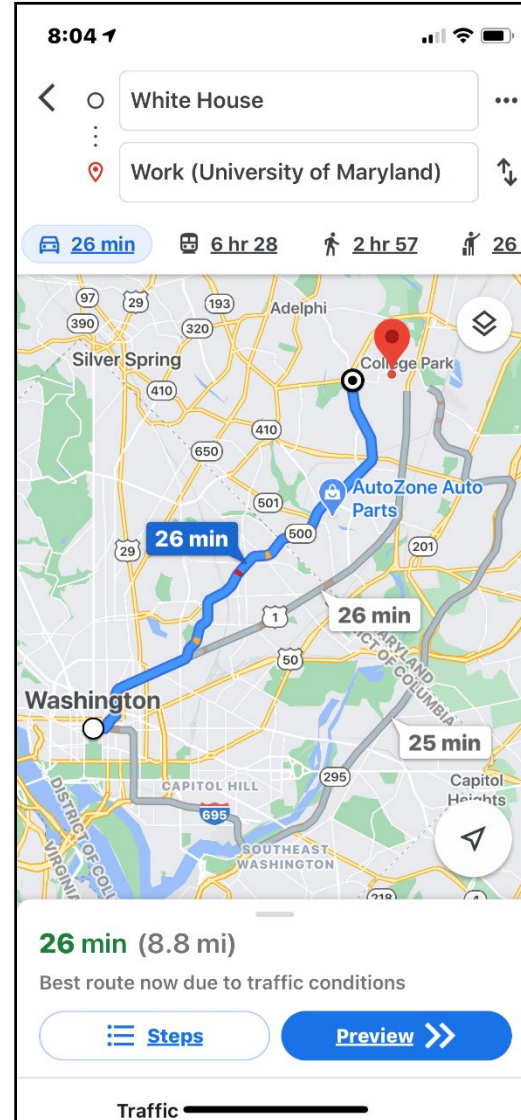  -> design computers to be supertools
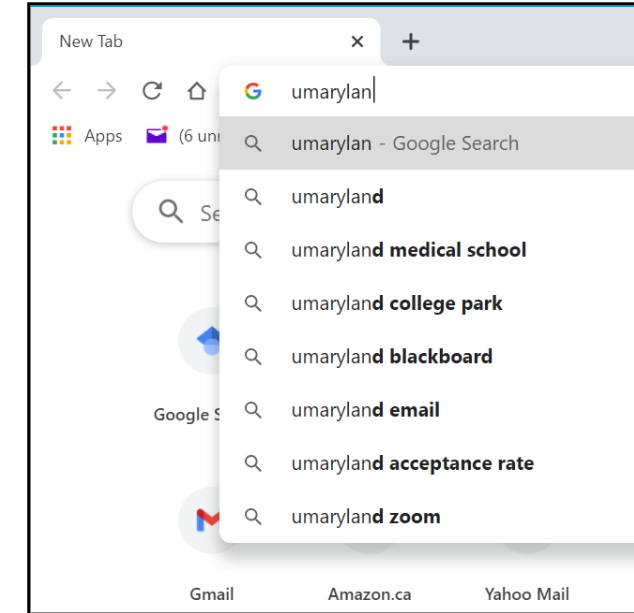  -> invite users to fix, personalize & extend the design
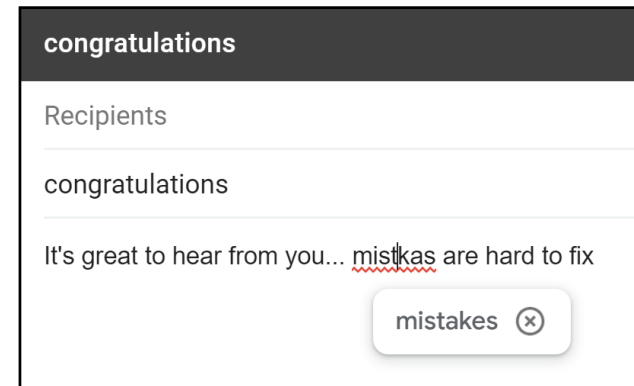
# Supertools

## Digital Camera Controls



## Navigation Choices



## Texting Autocompletion



## Spelling correction

# Active Appliances

## Coffee maker, Rice cooker, Blender

## Dishwasher, Clothes Washer/Dryer



Cuisinart Grind & Brew Coffee Maker



Panasonic Rice Cooker



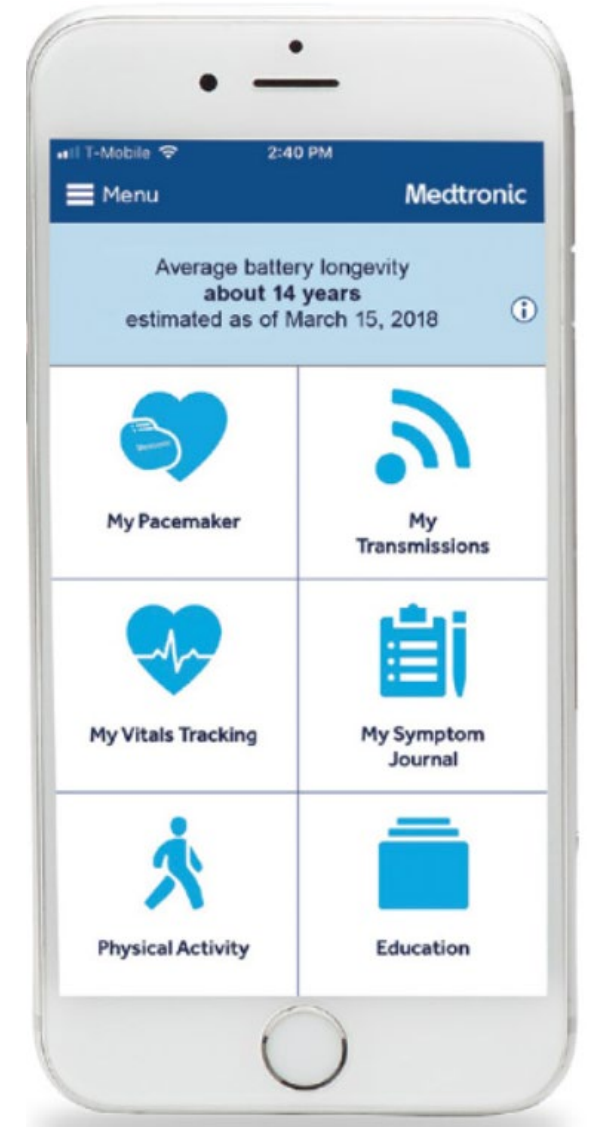Nutri Ninja Blender



Miele Dishwasher



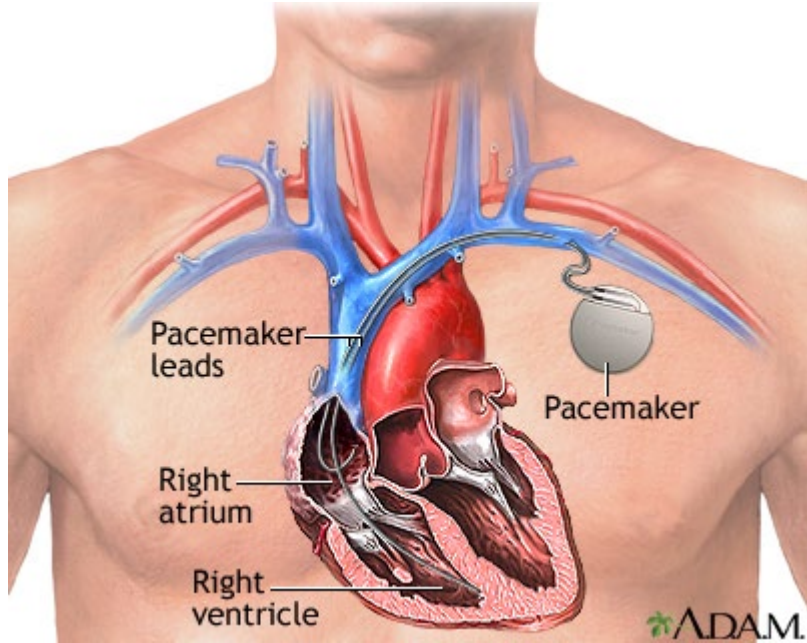General Electric Washer



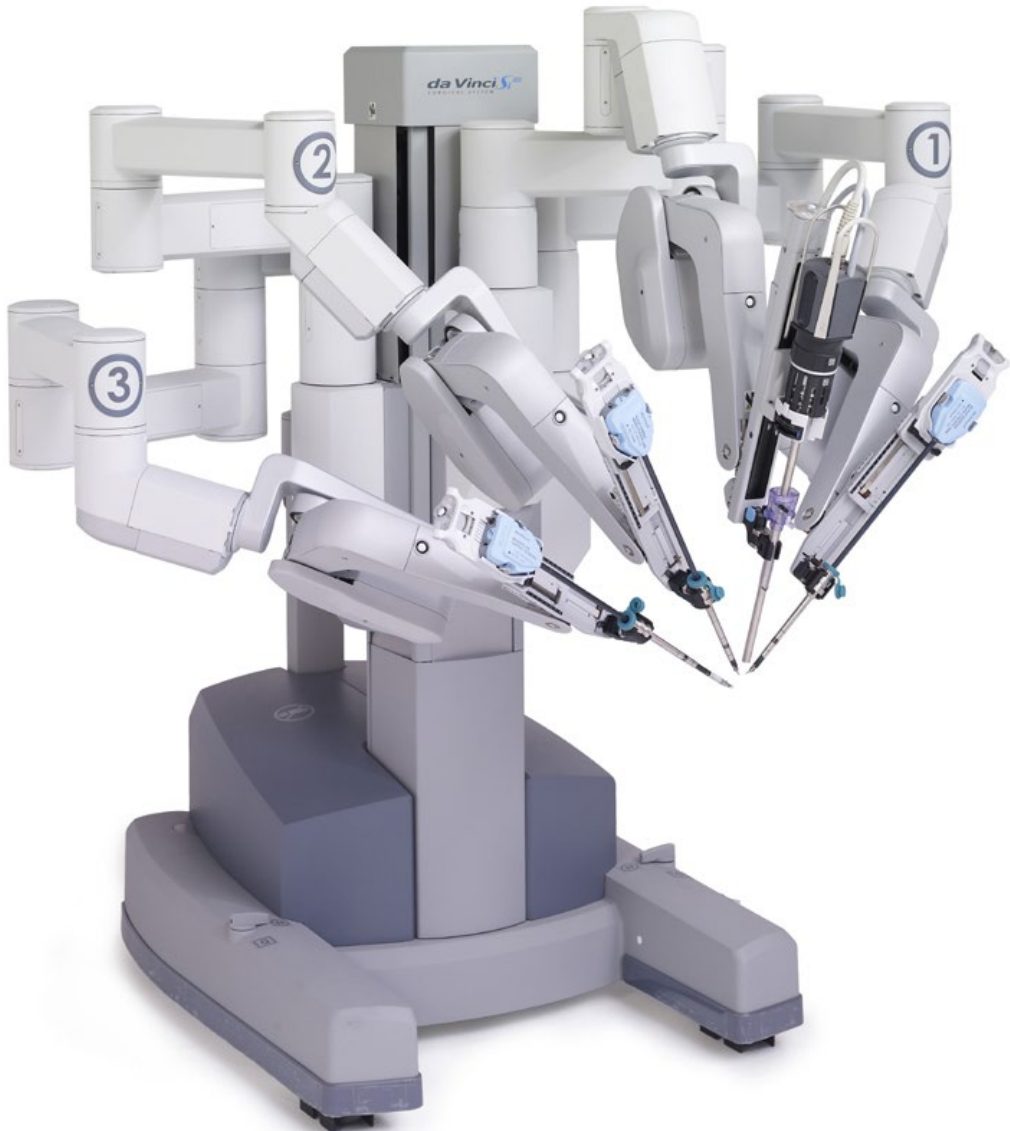General Electric Dryer

# Implanted Cardiac Pacemakers

# NASA Mars Rovers are Tele-Operated

# Da Vinci Tele-Operated Surgery





"Robots don't perform surgery. Your surgeon performs surgery with da Vinci by using instruments that he or she guides via a console."

https://www.davincisurgery.com/

# Bloomberg Terminal

# Hospital Control Center

# Counter Terrorism Center

# Design Guidelines

| Eight Golden Rules |
| --- |
| 1. Strive for consistency |
| 2. Seek universal usability |
| 3. Offer informative feedback |
| 4. Design dialogs to yield closure |
| 5. Prevent errors |
| 6. Permit easy reversal of actions |
| 7. Keep users in control |
| 8. Reduce short-term memory load |

https://www.cs.umd.edu/~ben/goldenrules.html

# Design Guidelines

| Eight Golden Rules |
|---|
| 1. Strive for consistency |
| 2. Seek universal usability |
| 3. Offer informative feedback |
| 4. Design dialogs to yield closure |
| 5. Prevent errors |
| 6. Permit easy reversal of actions |
| 7. Keep users in control |
| 8. Reduce short-term memory load |

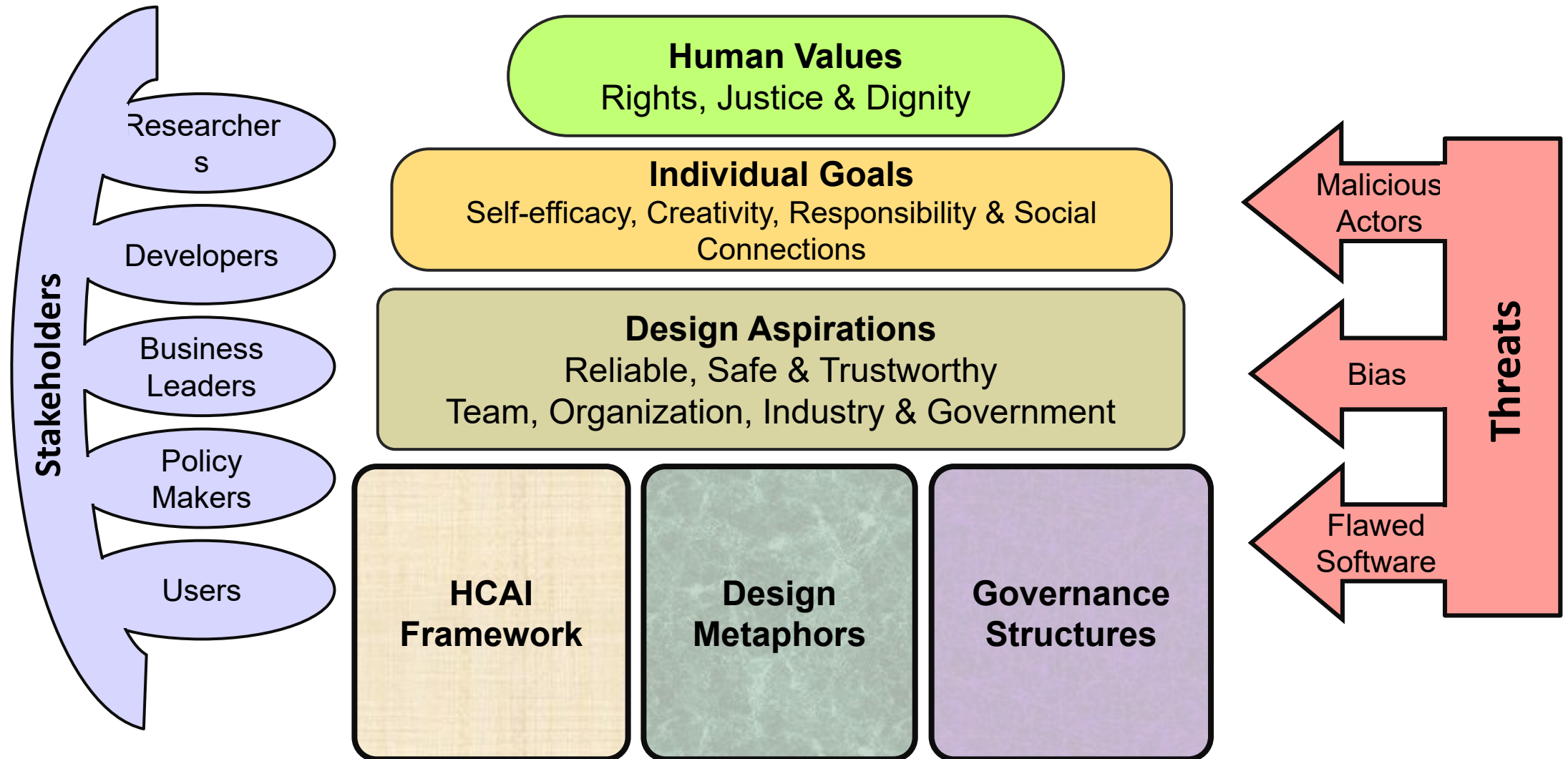| Eight Silver Slogans for HCAI Systems |
|---|
| 1. Store rich data from powerful sensors |
| 2. Design information abundant displays |
| 3. Provide interactive information visualization |
| 4. Make predictive models visual |
| 5. Smooth human-to-human communication |
| 6. Create clear control panels |
| 7. Implement audit trails |
| 8. Develop incident reporting websites |

https://www.cs.umd.edu/~ben/goldenrules.html

# Summary

# Human-Centered AI

**Stakeholders**
- Researchers
- Developers
- Business Leaders
- Policy Makers
- Users

**Human Values**
Rights, Justice & Dignity

**Individual Goals**
Self-efficacy, Creativity, Responsibility & Social Connections

**Design Aspirations**
Reliable, Safe & Trustworthy
Team, Organization, Industry & Government

**HCAI Framework**

**Design Metaphors**

**Governance Structures**

**Threats**
- Malicious Actors
- Bias
- Flawed Software

Oxford University Press (Early 2022)  `https://hcil.umd.edu/human-centered-ai/`

# Human-Centered AI Framework



**Excessive Human Control**

**Human Control** — High / Low

**Computer Automation** — Low / High

**Excessive Automation**

| | |
|---|---|
| **Human Mastery**<br>Bicycle<br>Piano | **Reliable, Safe & Trustworthy**<br>Elevator<br>Camera |
| Music box<br>Landmine | Pacemaker<br>Airbag<br>**Computer Control** |

# Design Metaphors

**Science Goal**

**Innovation Goal**

**Combined Designs**

**Intelligent Agents**
Thinking Machine, Cognitive Actor, Artificial Intelligence, Knowledgeable

**Supertools**
Extend Abilities, Empower Users, Enhance Human Performance

**Teammates**
Co-active Collaborator, Colleague, Helpful Partner, Smart Co-worker

**Tele-operated Devices**
Steerable Instrument, Powerful Prosthetic, Boost Human Perceptual & Motor Skills

**Assured Autonomy**
Independent, Self-directed, Goal-setting, Self-monitored

**Supervised Autonomy**
Human Control & Oversight, Situation Awareness, Predictable Actions

**Social Robots**
Anthropomorphic, Humanoid, Android, Bionic, Bio-inspired

**Active Appliances**
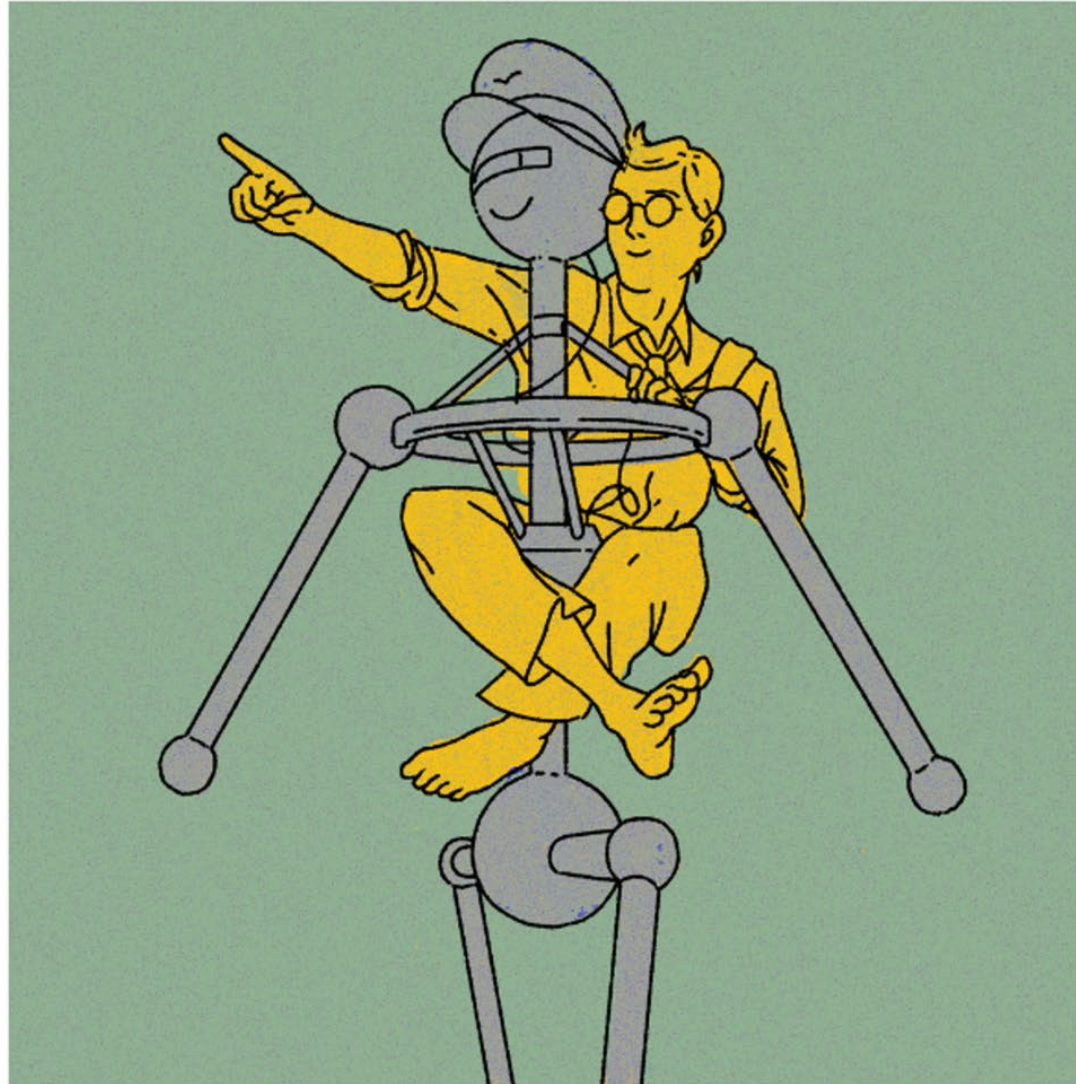Consumer-oriented, Wide Use, Low Cost Comprehensible Control Panels

# Technology

## A Case for Cooperation Between Machines and Humans

A computer scientist argues that the quest for fully automated robots is misguided, perhaps even dangerous. His decades of warnings are gaining more attention.

By **John Markoff**

May 21, 2020    Updated 3:09 p.m. ET

https://www.nytimes.com/2020/05/21/technology/ben-shneiderman-automation-humans.html

Human-Centered Artificial Intelligence: Reliable, safe & trustworthy, *International Journal of Human-Computer Interaction 36,* 6 (March 2020). https://doi.org/10.1080/10447318.2020.1741118

Design lessons from AI's two grand goals: Human emulation and useful applications, *IEEE Transactions on Technology & Society 1,* 2 (June 2020). https://ieeexplore.ieee.org/document/9088114

Bridging the gap between ethics and practice: Guidelines for reliable, safe, and trustworthy Human-Centered AI systems, *ACM Trans. on Interactive Intelligent Systems 10,* 4 (Oct 2020). https://dl.acm.org/doi/10.1145/3419764

Human-Centered Artificial Intelligence: Three fresh ideas, *AIS Trans. on Human-Computer Interaction 12,* 3 (Oct 2020). https://aisel.aisnet.org/thci/vol12/iss3/1/

Human-Centered AI. *NAS ISSUES 37, 2* (Winter 2021). https://issues.org/human-centered-ai/

Summary & resources: https://hcil.umd.edu/human-centered-ai/

# The Future is Human-Centered