

ACM Intelligent User Interfaces Conference, April 13, 2021

https://iui.acm.org/2021/hcai_tutorial.html

Human-Centered AI: Reliable, Safe & Trustworthy Part 2

Ben Shneiderman @benbendc

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

Member, National Academy of Engineering

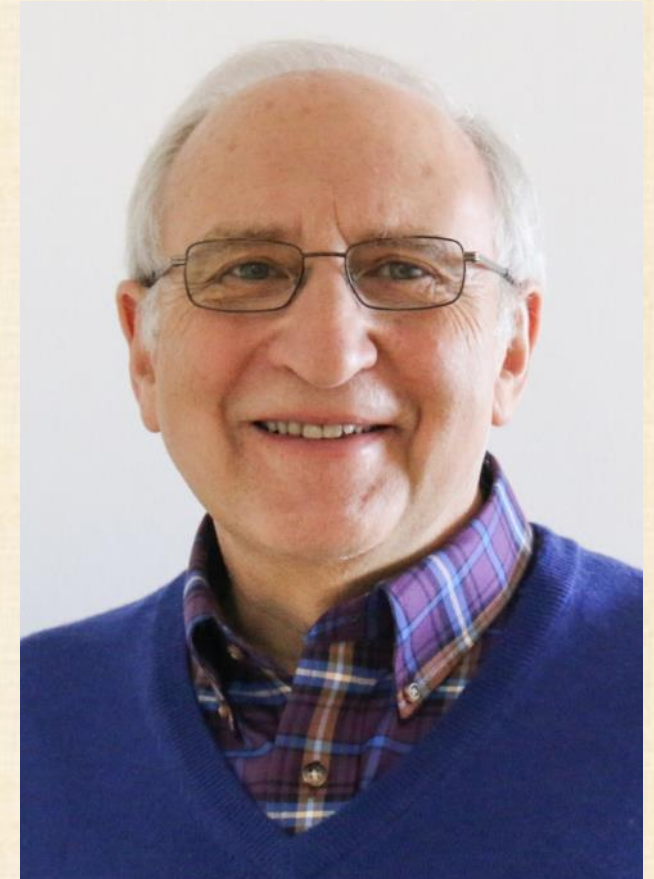


Photo: BK Adams



Interdisciplinary research community

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

hcil.umd.edu
vimeo.com/72440805



Annual Symposium: FREE Virtual
May 27, 2021, Thursday

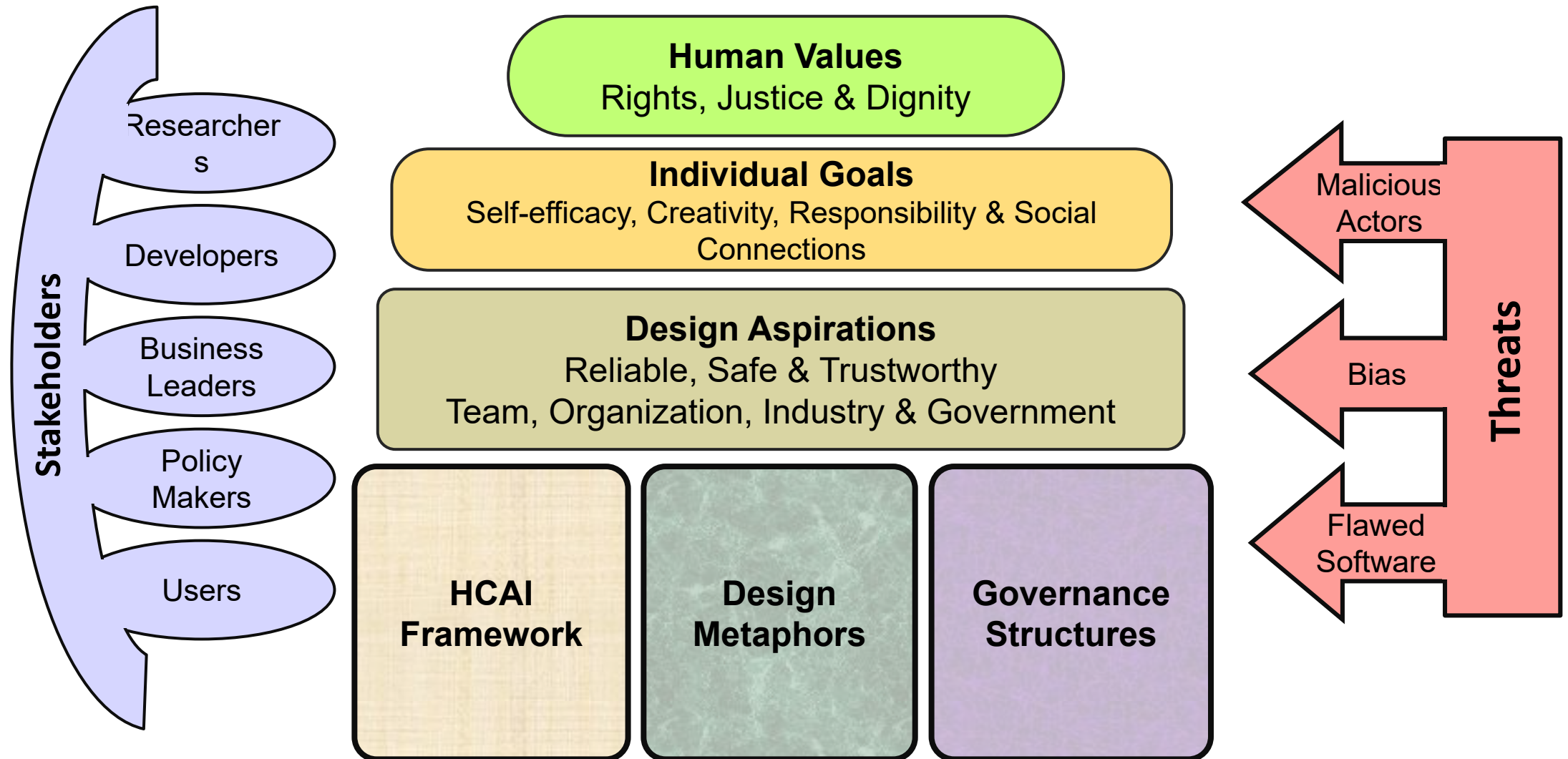
<https://hcil.umd.edu/2021-symposium/>

What is Human-Centered AI?



Amplify, Augment, Empower & Enhance People

Human-Centered AI



Governance Structures



ARTIFICIAL INTELLIGENCE AND LIFE IN 2030

Stanford University

SEPTEMBER 2016

NATIONAL (INDIA) STRATEGY FOR ARTIFICIAL INTELLIGENCE #AIFORALL

Pew Research Center 

Artificial Intelligence and the Future of Humans

HUMAN RIGHTS IN THE AGE OF ARTIFICIAL INTELLIGENCE

 accessnow

Artificial Intelligence, Robotics and 'Autonomous' Systems

*European Group on
Ethics in Science and
New Technologies*

House of Commons
Science and Technology
Committee

Algorithms in decision- making

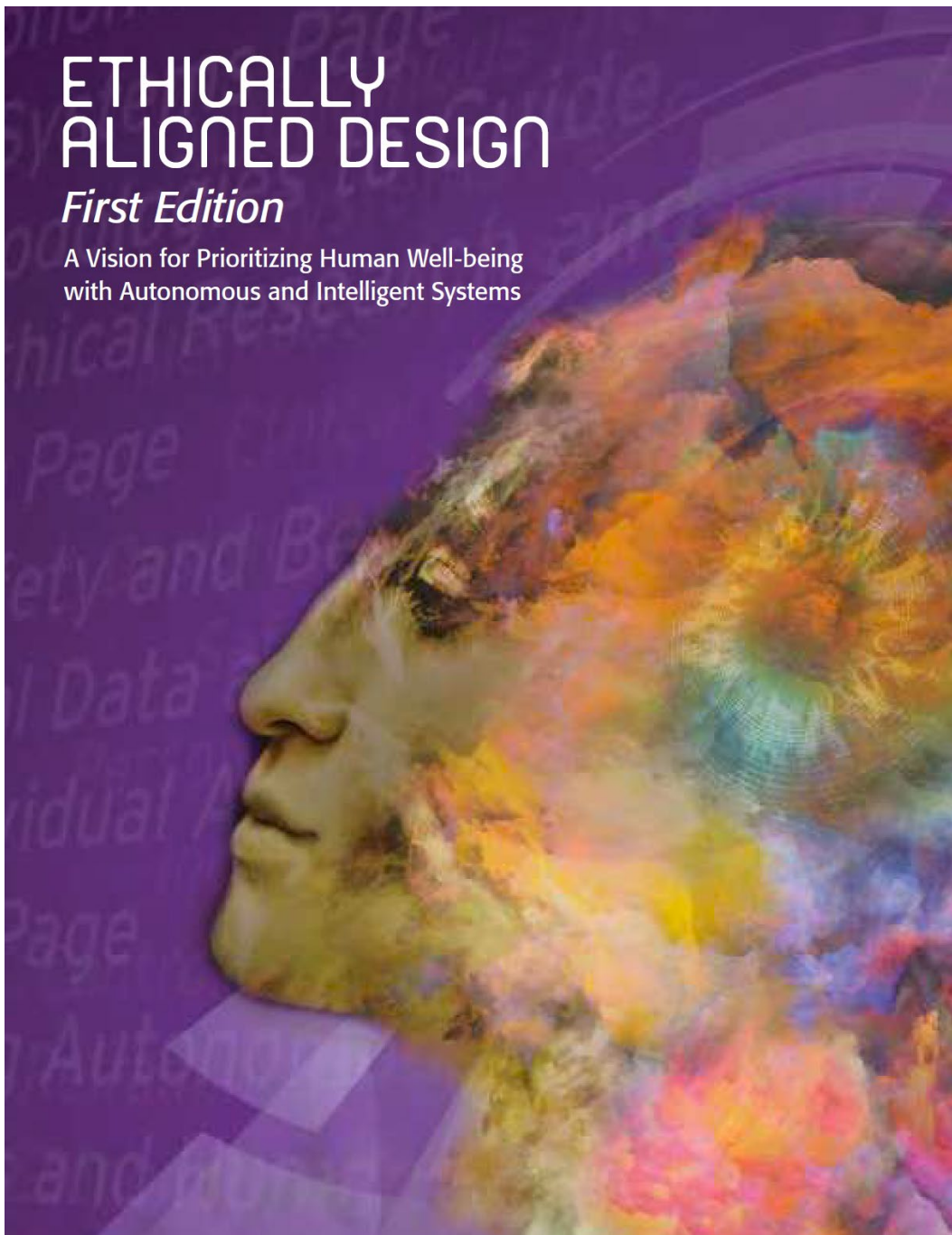
Governing Artificial
Intelligence:
UPHOLDING
HUMAN RIGHTS
& DIGNITY

Mark Latonero
Data&Society

ETHICALLY ALIGNED DESIGN

First Edition

A Vision for Prioritizing Human Well-being
with Autonomous and Intelligent Systems



PRINCIPLED ARTIFICIAL INTELLIGENCE

A Map of Ethical and Rights-Based Approaches to Principles for AI



Ethical AI Principles

Berkman Klein Center

IEEE Ethically Aligned Design

**Close
Match**

Accountability
Transparency & explainability
Promotion of human values
Safety & security

Accountability
Transparency
Human rights
Well-being

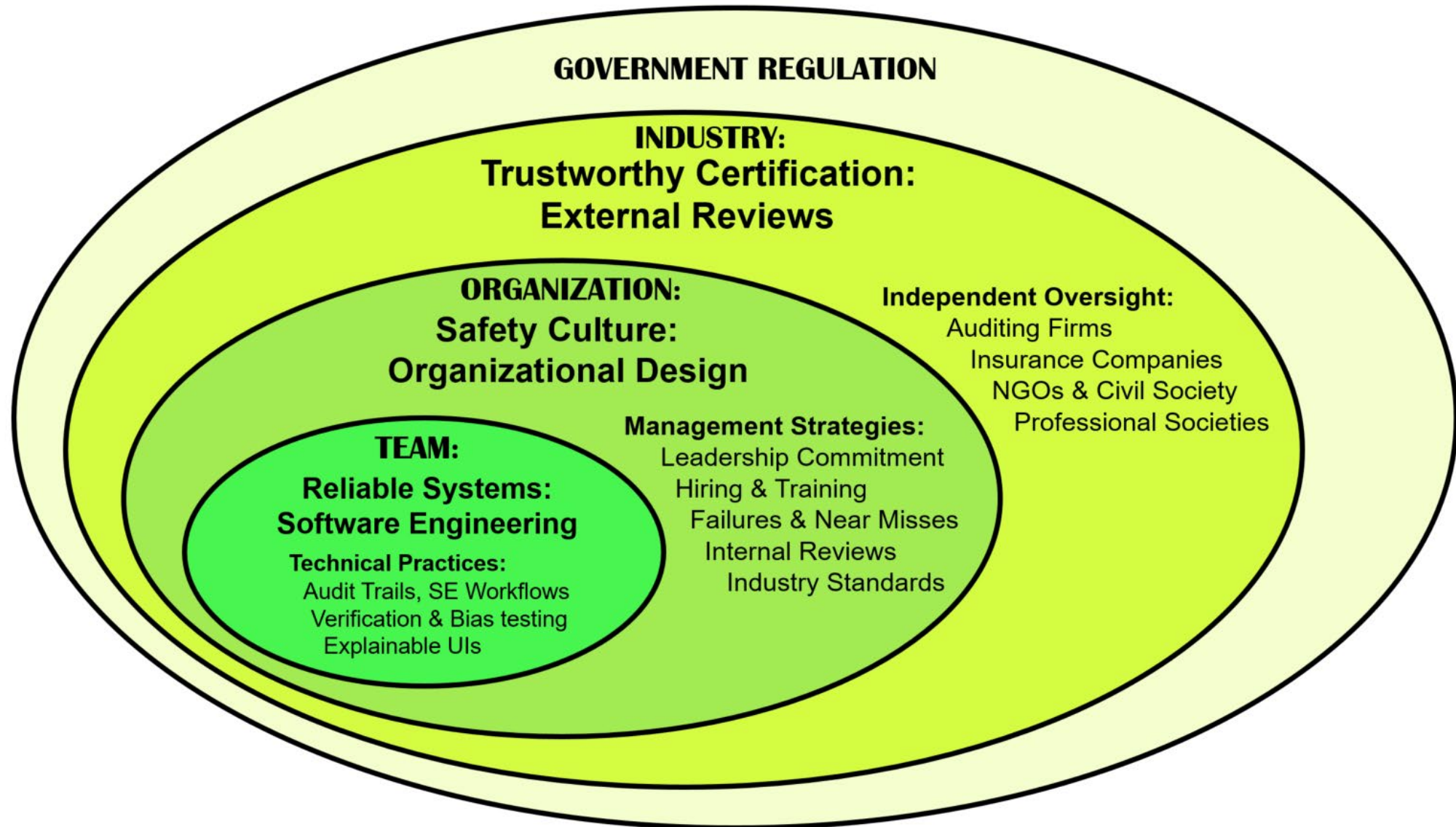
Similar

Human control of technology
Fairness & non-discrimination
Professional responsibility
Privacy

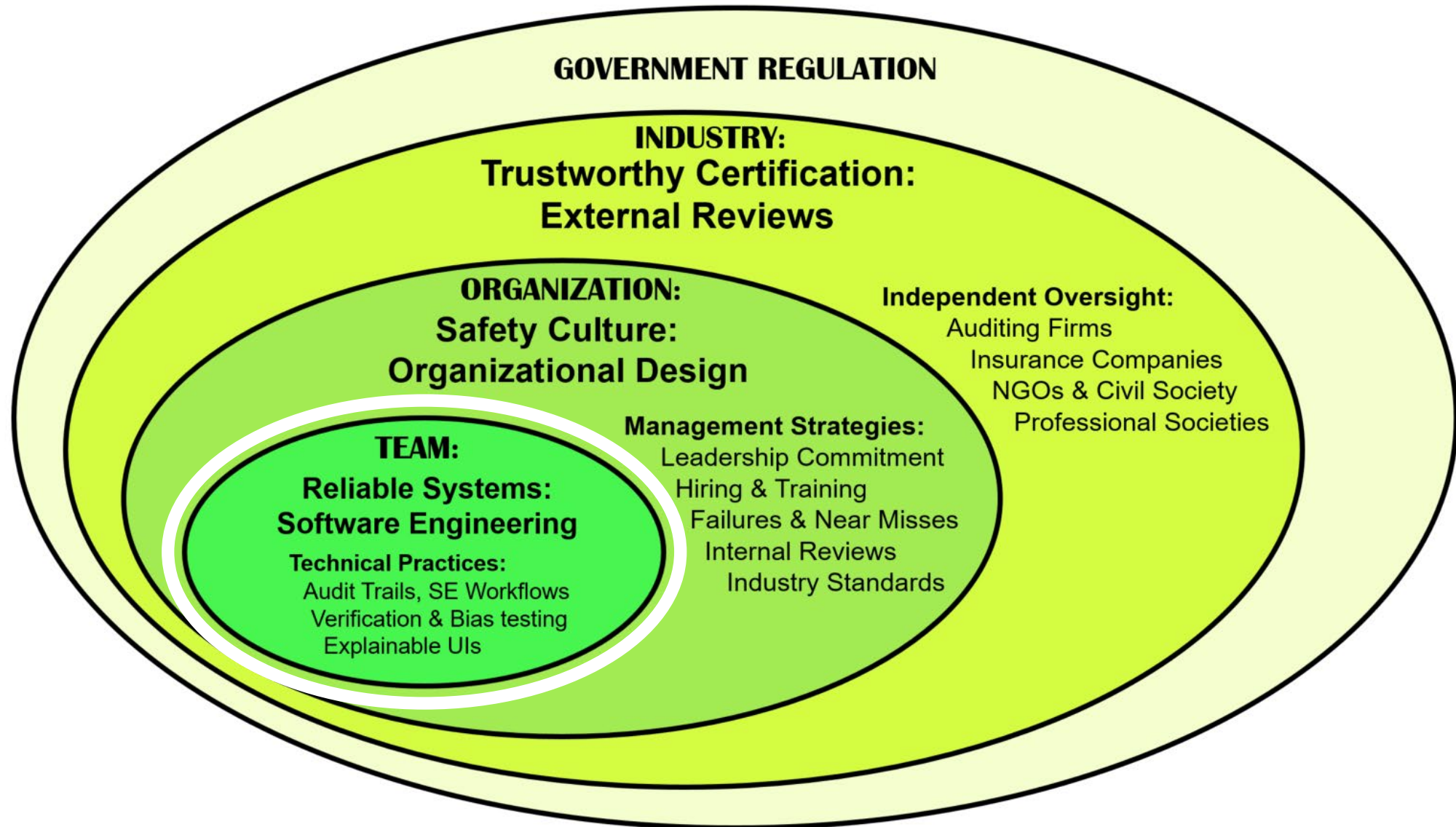
Effectiveness
Awareness of misuse
Competence
Data agency

<https://cyber.harvard.edu/publication/2020/principled-ai>
<https://ethicsinaction.ieee.org/>

Governance Structures for Human-Centered AI



Governance Structures for Human-Centered AI



TEAM

Reliable systems based on software engineering practices

- 1) Audit trails and analysis tools
- 2) Software engineering workflows
- 3) Verification & validation testing
- 4) Bias testing to improve fairness
- 5) Explainable user interfaces

Reliable Systems

Software engineering practices for a TEAM

1) Audit trails and analysis tools

“Flight Data Recorder for Every Robot”

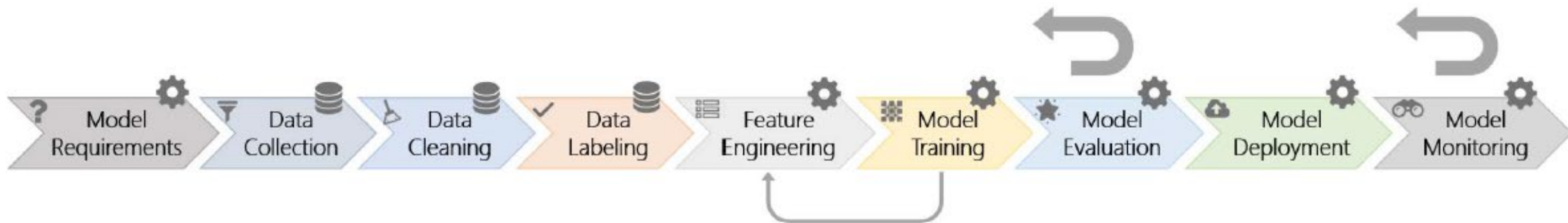
- Retrospective analysis of failures
- Understanding near misses
- Analysis to support preventive maintenance

Reliable Systems

Software engineering practices for a TEAM

2) Software engineering workflows

Microsoft's Workflow for Machine Learning



(Amershi, et al., IEEE/ACM ICSE 2019)

Reliable Systems

Software engineering practices for a TEAM

2) Software engineering workflows

Google Workflow for Algorithmic Auditing

1. Scoping: identify project scope & audit, raise questions of risk
2. Mapping: create stakeholder map & collaborator contact list, conduct interviews & select metrics
3. Artifact Collection: document design process, datasets & machine learning models
4. Testing: conduct adversarial testing to probe edge cases & failure possibilities
5. Reflection: consider risk analysis, failure remediation & record design history

(Raji et al., ACM FAT* 2020)

Reliable Systems

Software engineering practices for a TEAM

3) Verification & validation testing

- Traditional case-based
- User Experience
- Metamorphic
- Red Teams
- Differential
- + Microsoft's Datasheets for Datasets
- + Google's Model Cards
- + IBM FactSheets
- + Track history of bugs, problems, concerns

Reliable Systems

Software engineering practices for a TEAM

4) Bias testing to improve fairness

- Pre-existing, Technical, Emergent Bias

(Friedman & Nissenbaum, 1996; Baeza-Yates, 2018)

- Facial Recognition Intersectional Bias

(Buolamwini & Gebru, 2019; *Coded Bias*, 2021)

- IBM's Fairness 360 toolkit

Reliable Systems

Software engineering practices for a TEAM

5) Explainable user interfaces

- Retrospective explanations (local & global)

New Goal: **Prevent** confusion and surprise

- Prospective user interfaces
- Interactive, visual, exploratory

Mortgage Loan Explanations

Post-hoc Report

Enter amounts to request mortgage:

Mortgage amount requested	<input type="text" value="375000"/>
Household monthly income	<input type="text" value="7000"/>
Liquid assets	<input type="text" value="48000"/>



Mortgage Loan Explanations

Post-hoc Report

Enter amounts to request mortgage:

Mortgage amount requested	<input type="text" value="375000"/>
Household monthly income	<input type="text" value="7000"/>
Liquid assets	<input type="text" value="48000"/>

Enter amounts to request mortgage:

Mortgage amount requested	<input type="text" value="375000"/>
Household monthly income	<input type="text" value="7000"/>
Liquid assets	<input type="text" value="48000"/>

We're sorry, your mortgage loan was not approved. You might be approved if you reduce the Mortgage amount requested, increase your Household monthly income, or increase your Liquid assets.



Mortgage Loan Explanations

Post-hoc Report

Enter amounts to request mortgage:

Mortgage amount requested	<input type="text" value="375000"/>
Household monthly income	<input type="text" value="7000"/>
Liquid assets	<input type="text" value="48000"/>

Enter amounts to request mortgage:

Mortgage amount requested	<input type="text" value="375000"/>
Household monthly income	<input type="text" value="7000"/>
Liquid assets	<input type="text" value="48000"/>

We're sorry, your mortgage loan was not approved. You might be approved if you reduce the Mortgage amount requested, increase your Household monthly income, or increase your Liquid assets.

Prospective User Interface

Adjust sliders to report your situation:

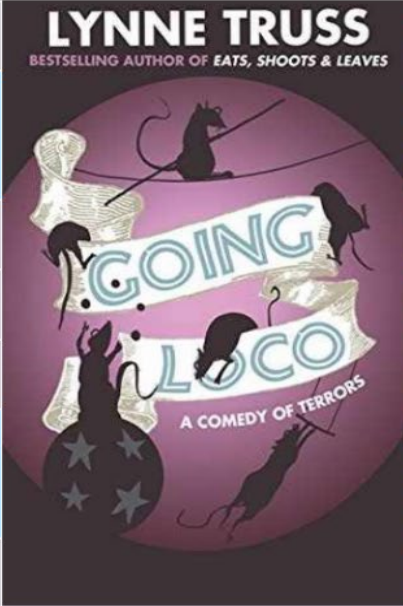
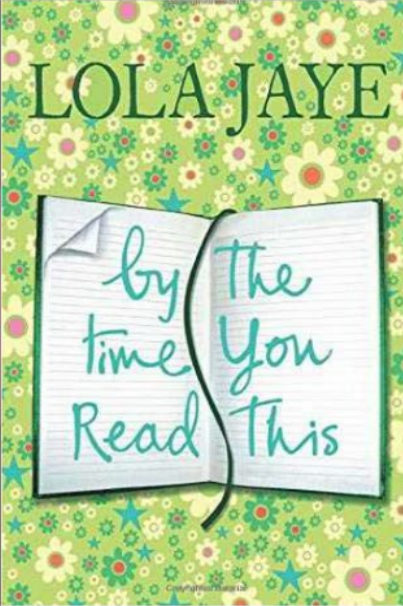
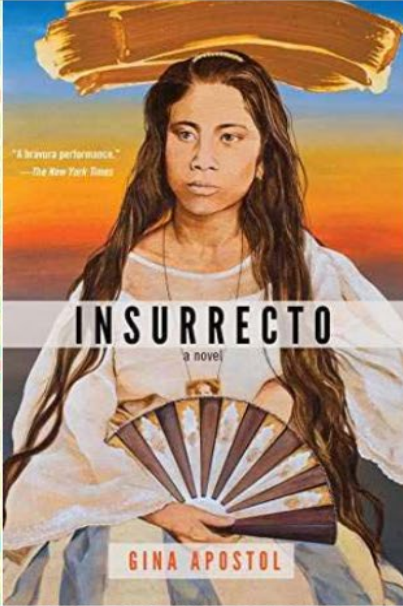
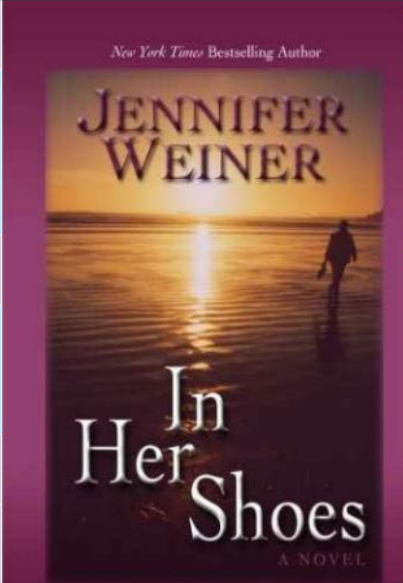
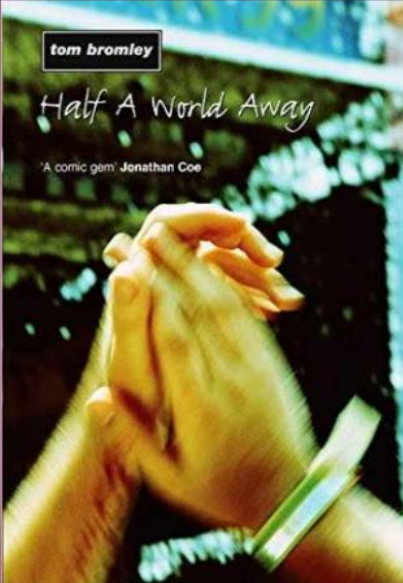

Mortgage amount requested	<input type="range" value="375000"/>	<input type="range" value="0.5"/>
Household monthly income	<input type="range" value="7000"/>	
Liquid assets	<input type="range" value="48000"/>	

Score needed for approval

Your score



Recommenders: Whichbook.net

Happy <input type="checkbox"/> Sad			
Funny <input checked="" type="checkbox"/> Serious			
Safe <input type="checkbox"/> Disturbing			
Expected <input type="checkbox"/> Unpredictable			
Larger than life <input type="checkbox"/> Down to earth			
Beautiful <input checked="" type="checkbox"/> Disgusting			
Gentle <input type="checkbox"/> Violent			
Easy <input type="checkbox"/> Demanding			
No sexual content <input checked="" type="checkbox"/> Explicit sexual content			
Conventional <input type="checkbox"/> Unusual			
Optimistic <input checked="" type="checkbox"/> Bleak			
Short <input type="checkbox"/> Long			
Select up to 4 sliders			

Recommender Control Panels

Modify Attributes

acousticness: 40

instrumentalness: 60

danceability: 80

valence: 60

energy: 40

Slider technique

Recommended Songs

Calculate Recommendations

Task: Make a playlist of songs to listen to during your personal maintenance.

Click or to keep or dismiss a song.
More info:

- Oblivion Grimes
- My December Linkin Park
- I've got that tune Chinese Man
- Got the Life Korn
- Good Riddance (Time of Y... Green Day
- Burn It To The Ground Nickelback

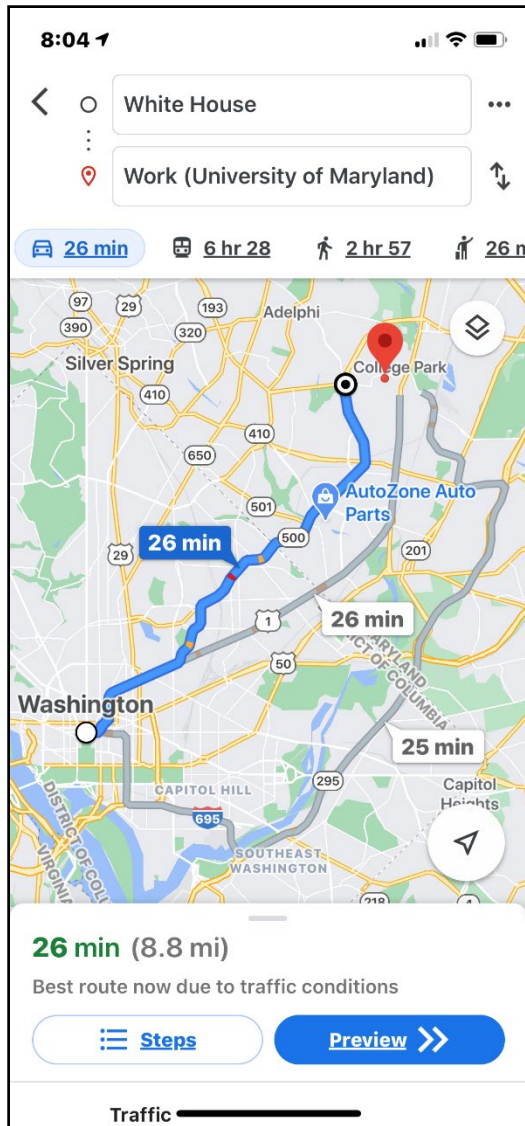
To get more songs, modify the attribute(s) and click "Calculate Recommendations" again.

Create Your Better Life Index

Rate the topics according to their importance to you:

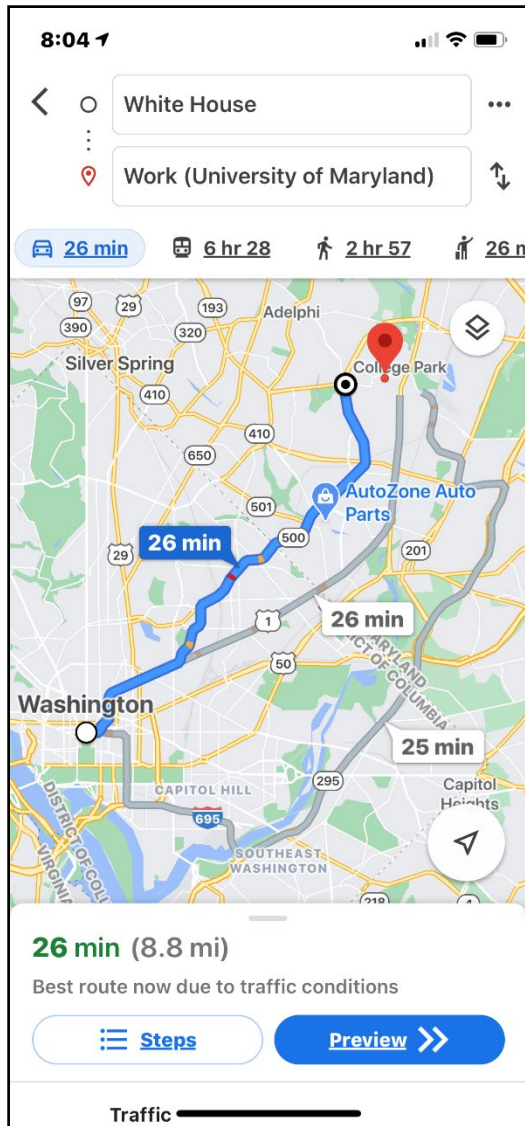
	-	+
Housing		
Income		
Jobs		
Community		
Education		
Environment		
Civic Engagement		
Health		
Life Satisfaction		
Safety		
Work-Life Balance		

Visual Previews



- 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

Visual Previews

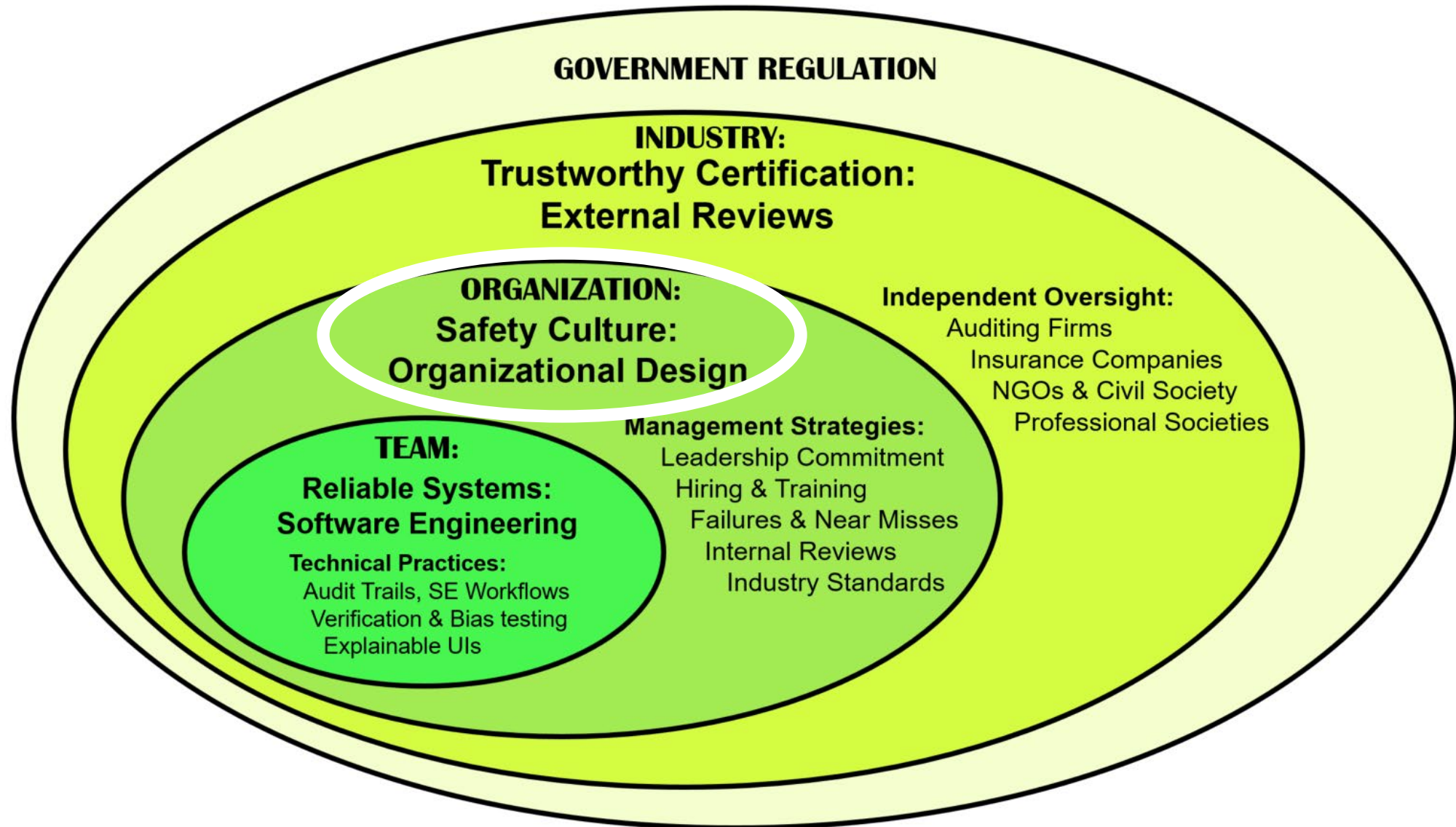


- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand
- Overview first, zoom & filter, then details-on-demand

- Preview first, select & initiate, then show execution
- Preview first, select & initiate, then show execution
- Preview first, select & initiate, then show execution
- Preview first, select & initiate, then show execution



Governance Structures for Human-Centered AI



ORGANIZATION

Safety culture through business management strategies

- 6) Leadership commitment to safety
- 7) Hiring and training oriented to safety
- 8) Extensive reporting of failures and near misses
- 9) Internal review boards for problems and future plans
- 10) Alignment with industry standard practices

Safety Culture

Business management strategies for an ORGANIZATION

6) Leadership commitment to safety

- Normal Accident Theory
- High Reliability Organizations
- Resilience Engineering
- Safety Culture: Repeated public statements, Vision statements, Budget, Openness about failures, Annual reports on safety, Competitive advantage

Safety Culture

Business management strategies for an ORGANIZATION

7) Hiring and training oriented to safety

- Job notices consistently emphasize safety
- Interviews focus on safety
- Training for safety
- Retraining & simulated disasters

Safety Culture

Business management strategies for an ORGANIZATION

8) Extensive reporting of failures and near misses

- Public & internal reporting
- Data analysis & reviews
- US FAA, FDA Adverse event reporting
- Bug & bias bounties
- AI Incident Database (<https://incidentdatabase.ai/>)
- Tesladeaths.com

Safety Culture

Business management strategies for an ORGANIZATION

9) Internal review boards for problems & future plans

- Regular reviews & public reporting
- Oversight group gains experience
- Audit committees

Safety Culture

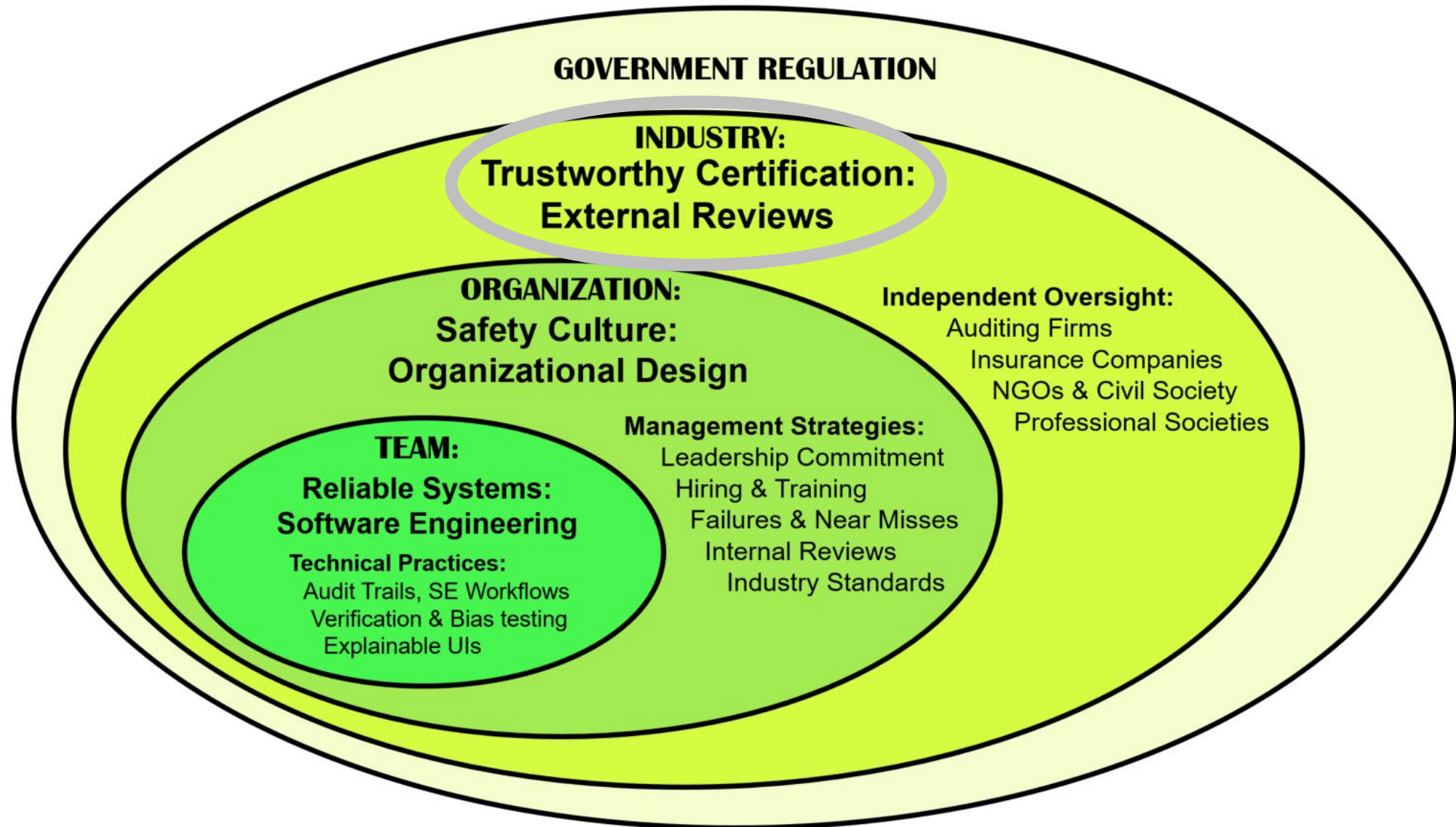
Business management strategies for an ORGANIZATION

10) Alignment with industry standard practices

- ISO Technical Committee on Robotics
- Robotics Industry Association
- Underwriters Laboratory, Consumer Reports
- IEEE P7000 Series: Ethics, Wellbeing, Transparency
- Capability Maturity Models



Governance Structures for Human-Centered AI



INDUSTRY

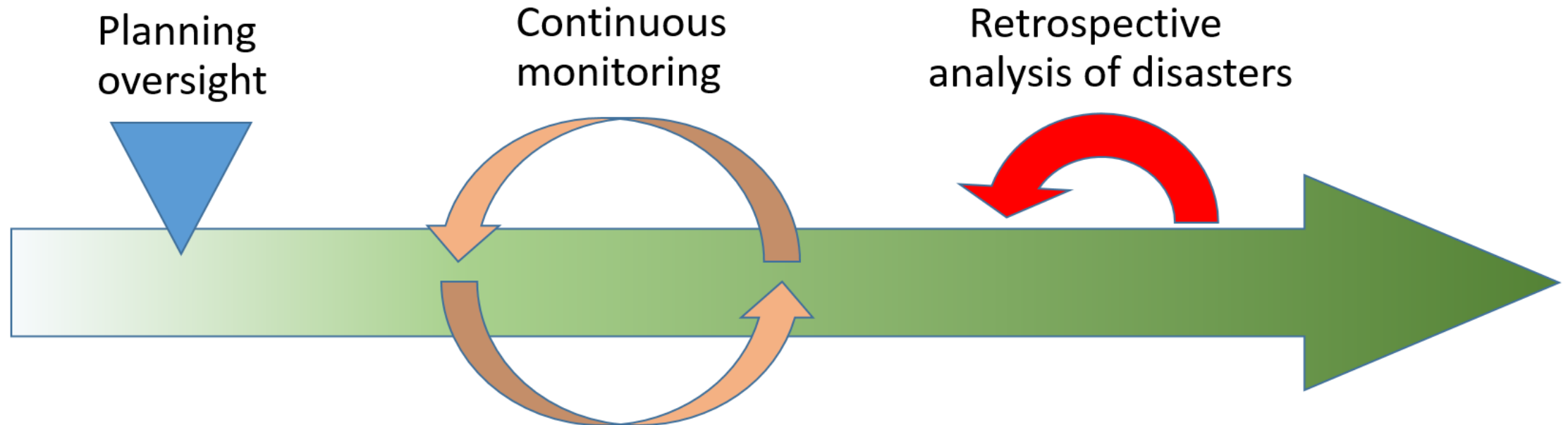
Trustworthy certification by independent oversight

- 11) Accounting firms conduct external audits
- 12) Insurance companies compensate for failures
- 13) Non-governmental and civil society organizations
- 14) Professional organizations and research institutes

- 15) Government interventions and regulation

Trustworthy certification

Independent oversight for an INDUSTRY



- Degree of Independence, subpoena power
- Powers to enforce recommendations

PNAS Opinion: (November 29, 2016) <http://www.pnas.org/content/113/48/13538.full>

To mitigate the dangers of faulty, biased, or malicious algorithms requires independent oversight

Planning oversight



Continuous monitoring



Federal Reserve

Retrospective analysis



NTSB Tweet: 6:25 PM - 7 Jul 2013

Trustworthy Systems

Certification by independent oversight for an INDUSTRY

11) Accounting firms conduct external audits

- Accounting firms can now include AI Audits
Deloitte, Ernst & Young, KPMG, PwC
- Consulting companies may play a role
Accenture, Boston Consulting, McKinsey & Co
- Experience across companies adds value

Trustworthy Systems

Certification by independent oversight for an INDUSTRY

12) Insurance companies compensate for failures

- Success in construction, healthcare, transportation,...
- Building codes for building code (Carl Landwehr)
- Self-driving car insurance premiums guide safety
- Skeptics don't trust insurance companies

Trustworthy Systems

Certification by independent oversight for an INDUSTRY

13) Non-governmental and civil society organizations

- Many NGOs already active
- Algorithmic Justice League success: facial recognition

Trustworthy Systems

Certification by independent oversight for an INDUSTRY

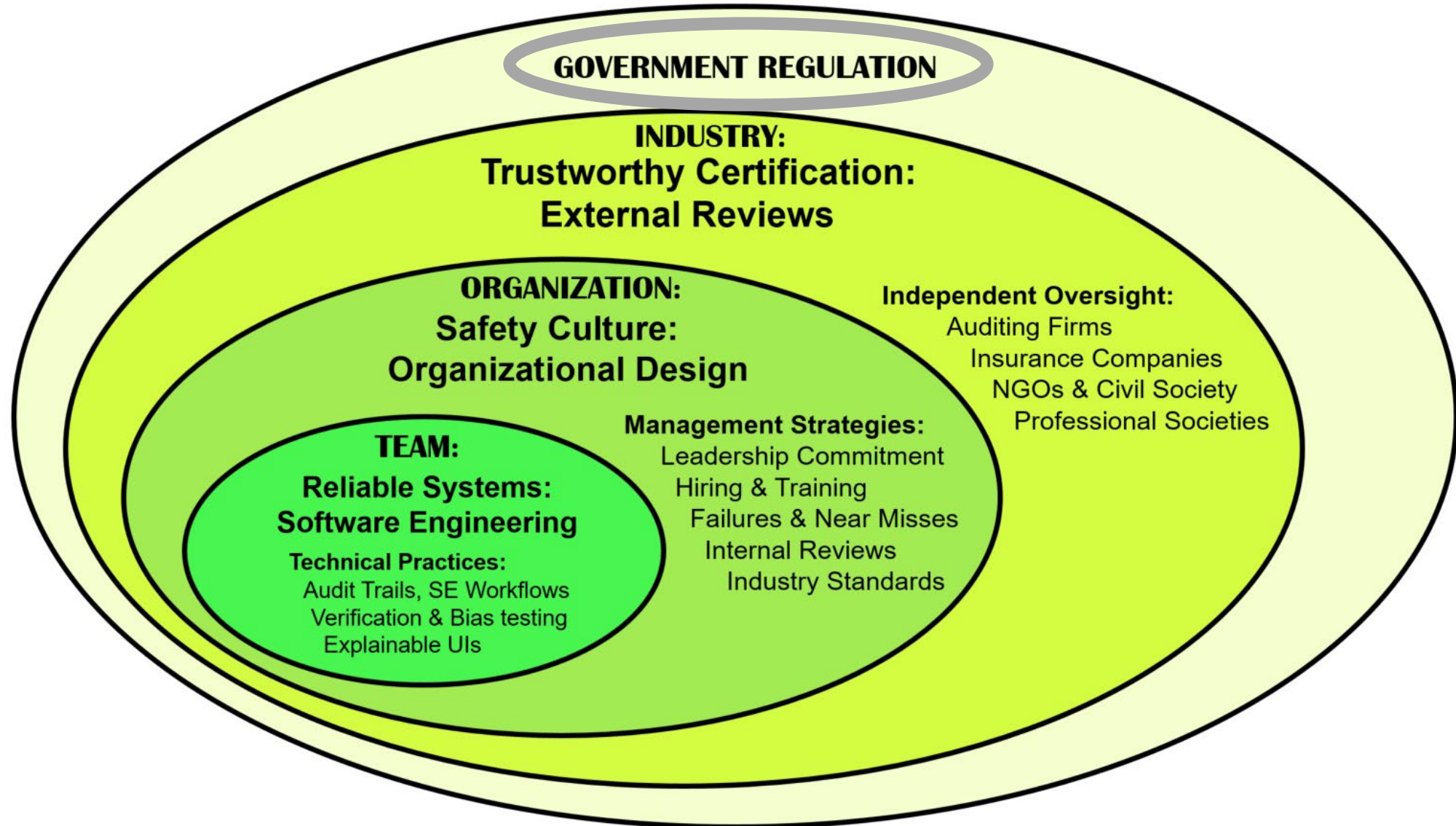
14) Professional organizations and research institutes

- IEEE P7000 series of standards
- IEEE Ethics of Autonomous & Intelligent Systems
- Montreal AI Ethics Institute
- OECD AI Policy Observatory
- ACM Technology Policy Committee
- Partnership on AI

University Research Groups

- Brown Univ Humanity Centered Robotics Initiative
- Columbia Univ Data Science Institute
- Harvard Univ Berkman Klein Center for Internet and Society
- Johns Hopkins Univ Institute for Assured Autonomy
- Monash Univ, Australia Human Centered AI
- New York Univ Center for Responsible AI
- Northwestern Univ Center for Human-Computer Interaction + Design
- Stanford Univ Human-centered AI (HAI) Institute
- Univ of British Columbia, Canada Human-AI Interaction
- Univ of California-Berkeley Center for Human-Compatible AI
- Univ of Cambridge, UK Leverhulme Centre for the Future of Intelligence
- Univ of Canberra, Australia Human Centred Technology Research Centre
- Univ of Chicago Chicago Human+AI Lab (CHAI)
- Univ of Oxford, UK Internet Institute, Future of Humanity Institute
- Univ of Toronto, Canada Ethics of AI Lab
- Utrecht Univ, Netherlands Human-centered AI

Governance Structures for Human-Centered AI



Government Regulation

Certification by independent oversight for an INDUSTRY

15) Government interventions and regulation

- Regulations: Good or Bad?
- U.S. FAA, FDA, FTC, NIST
- EU General Data Protection Regulation
- OECD Principles of HCAI
- UN AI for Good Global Summit

Government Regulation

Certification by independent oversight for an INDUSTRY

15) Government interventions and regulation

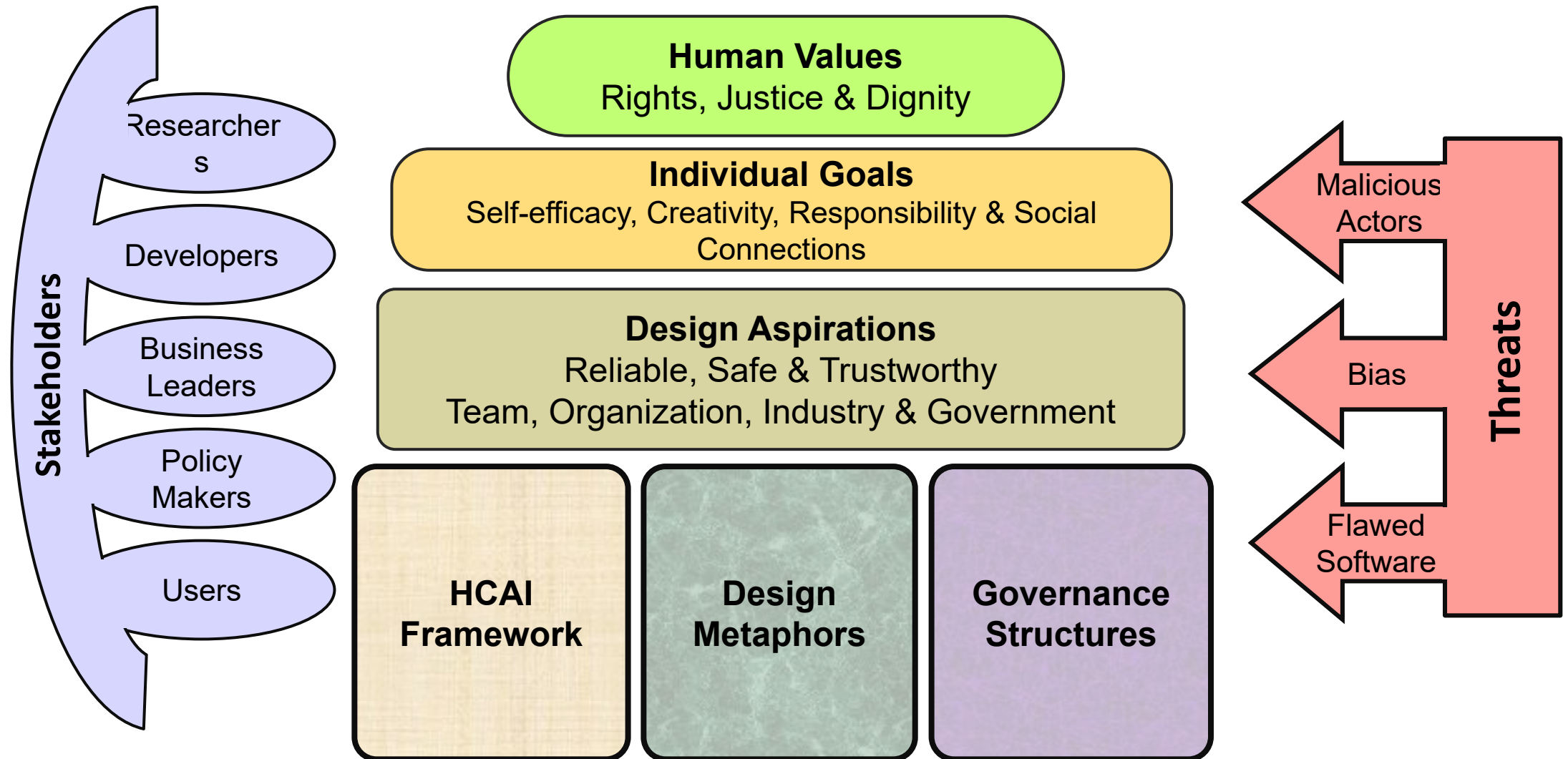
- Regulations: Good or Bad?
- U.S. FAA, FDA, FTC, NIST
- EU General Data Protection Regulation
- OECD Principles of HCAI
- UN AI for Good Global Summit



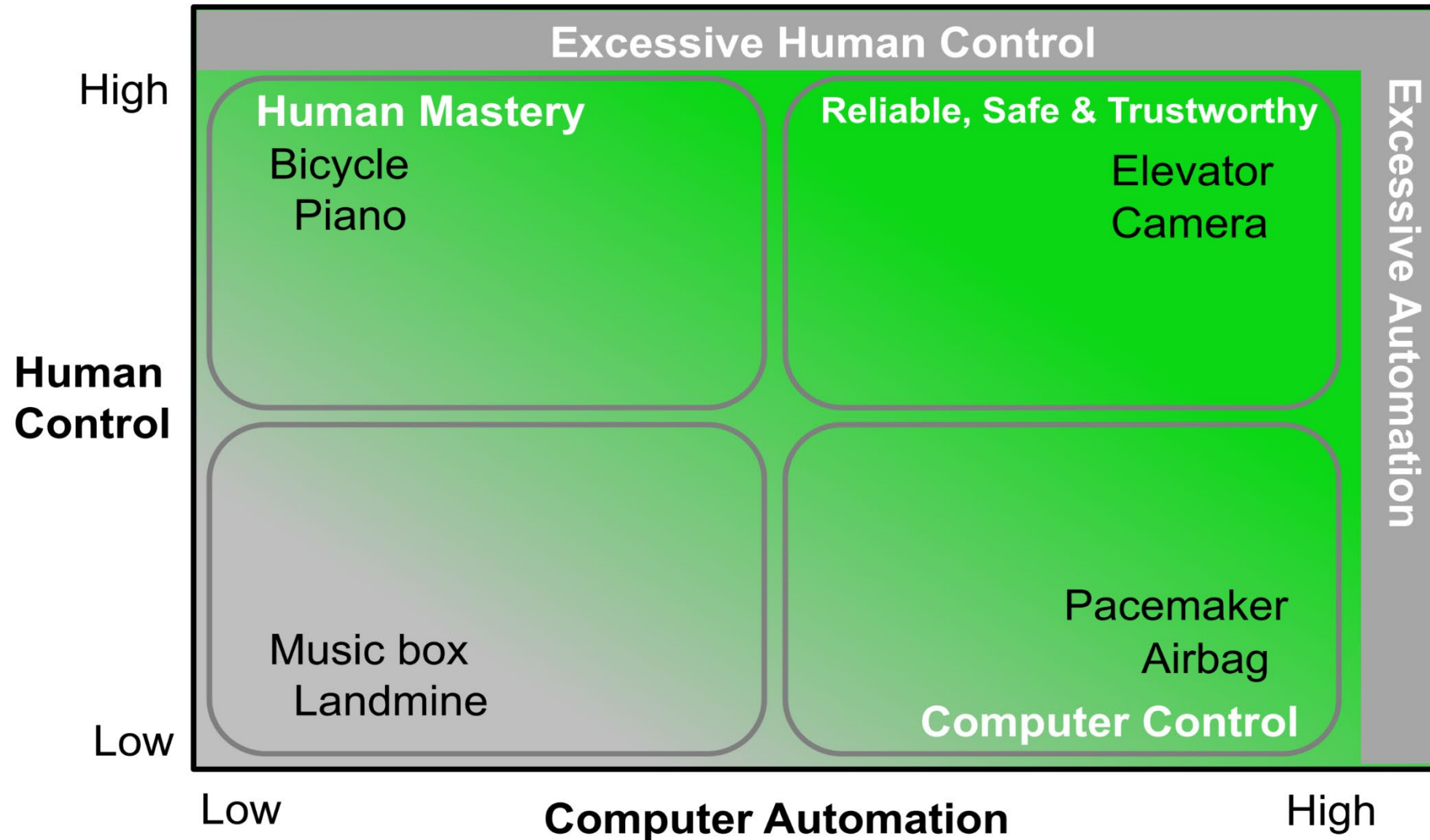
Summary



Human-Centered AI



HCAI Framework



Design Metaphors

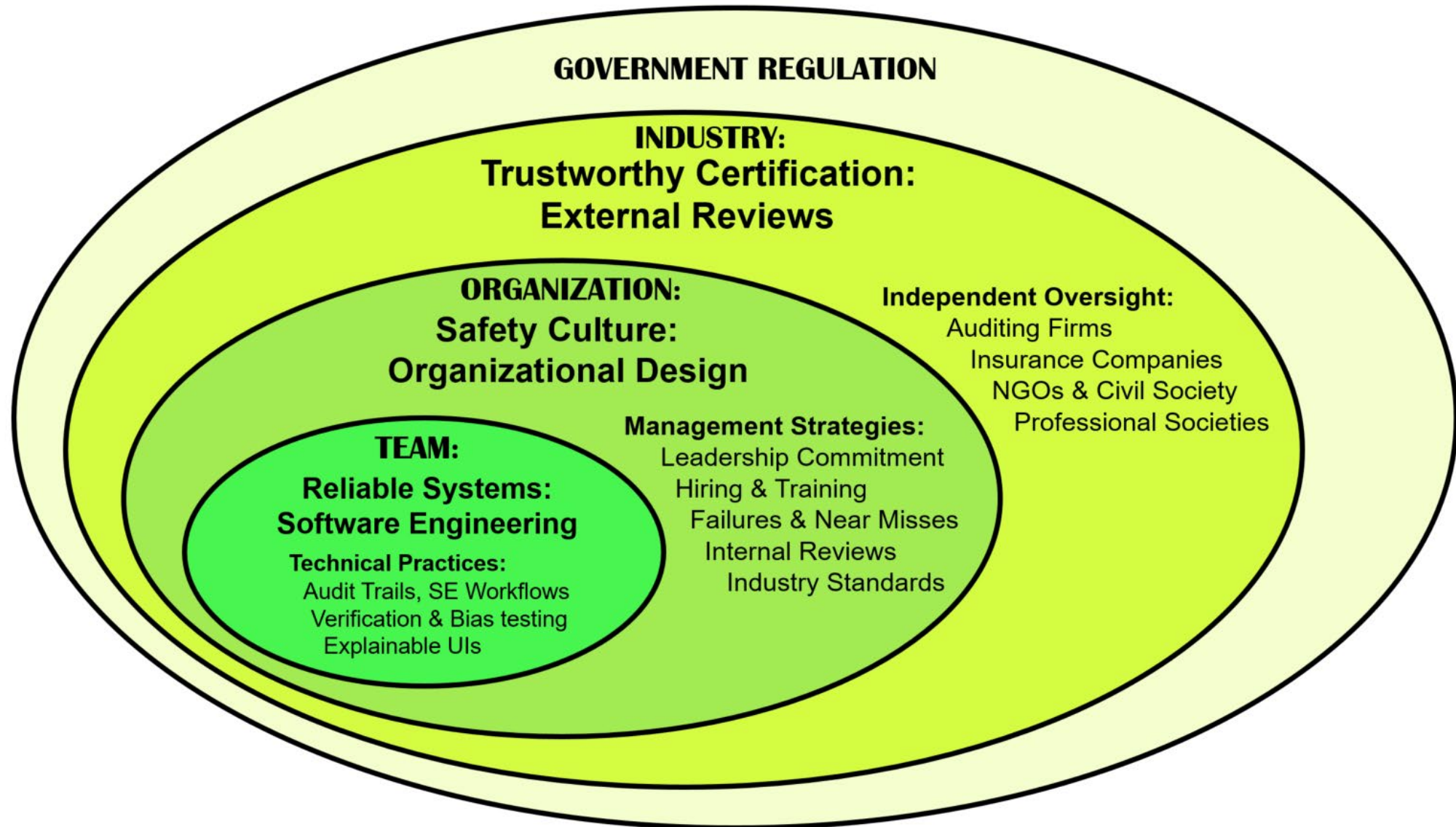
**Science
Goal**

**Innovation
Goal**

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

Combined Designs

Governance Structures for Human-Centered AI



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/>

Human-Centered AI: Google Group

<https://groups.google.com/g/human-centered-ai>



Annual Symposium: FREE Virtual
May 27, 2021, Thursday

<https://hcil.umd.edu/2021-symposium/>



