# AI AND THE ART OF MANIPULATION

How We Need to Think Differently About AI as We Develop Socially Responsible Applications

How Artificial Intelligence and Machine Learning
Transform the Human Condition
JULY 20, 2021

**ANDREW MAYNARD**
Professor and Associate Dean
Arizona State University College of Global Futures

1

# AI ETHICS

**2017**

**Asilomar AI Principles**
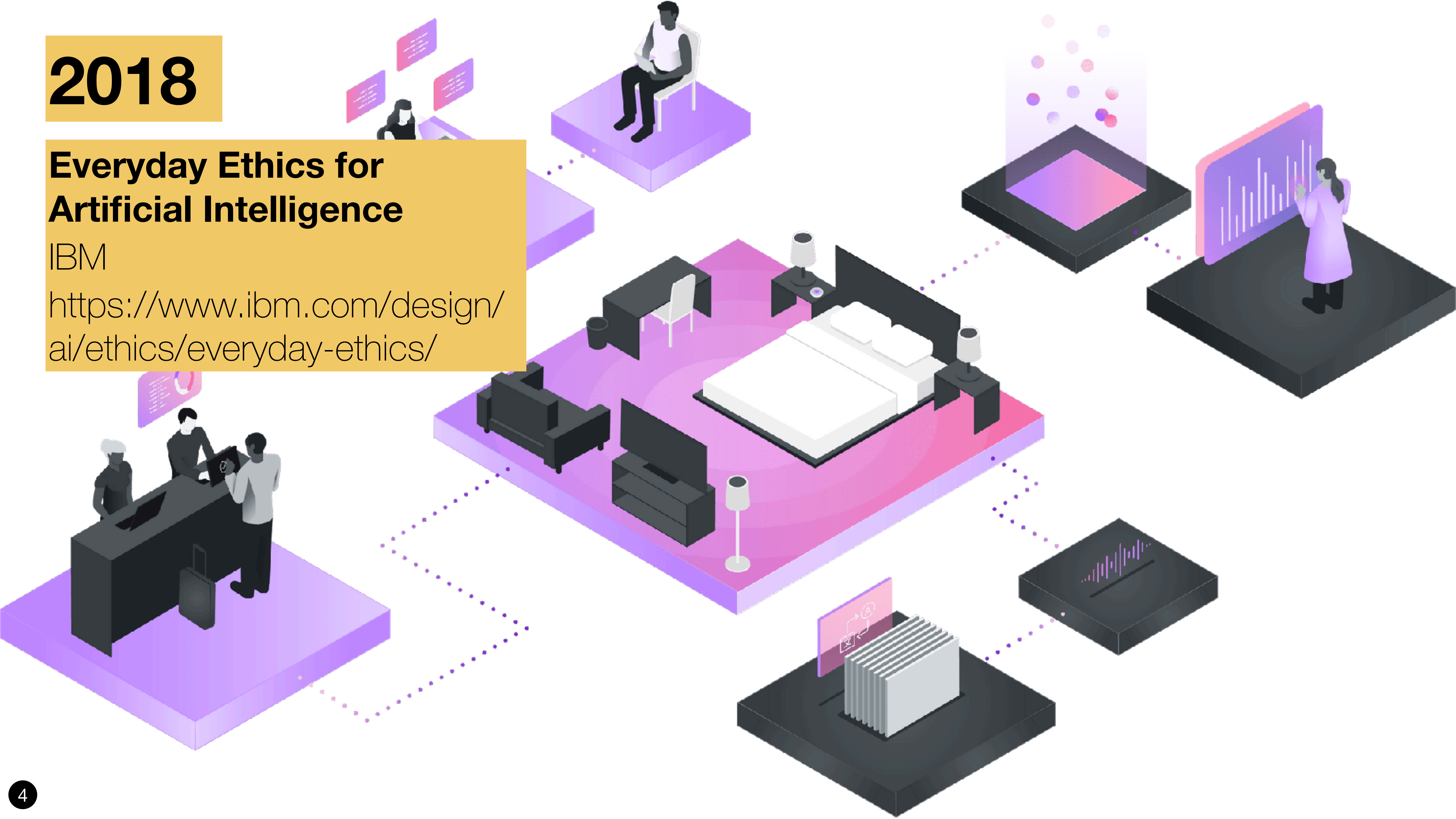Future of Life Institute
https://futureoflife.org/ai-principles/

3

**2018**

**Everyday Ethics for Artificial Intelligence**

IBM

https://www.ibm.com/design/ai/ethics/everyday-ethics/

**2019**

**Ethically Aligned Design:**
A Vision for Prioritizing Human Well-being
with Autonomous and Intelligent Systems
IEEE

https://ethicsinaction.ieee.org/

# 2020

**AI Now**

New York University

https://medium.com/@AINowInstitute/ai-in-2020-a-year-to-give-us-pause-67795fe23324

"2020 has been a year of hard truths and tragedy, as interlocking crises put the failures, inadequacies, and structural limitations of our core institutions in the spotlight. **At the same time, we see the AI industry rushing to profit in the space left by an absent social safety net, bolstered by governments' increasing turn to tech solutions.** AI companies are ramping up surveillance of our workplaces, schools and communities; cracking down on worker organizing and ethical research; and bankrolling the passage of bills that gut worker protections for millions — while growing richer and more powerful in the process."
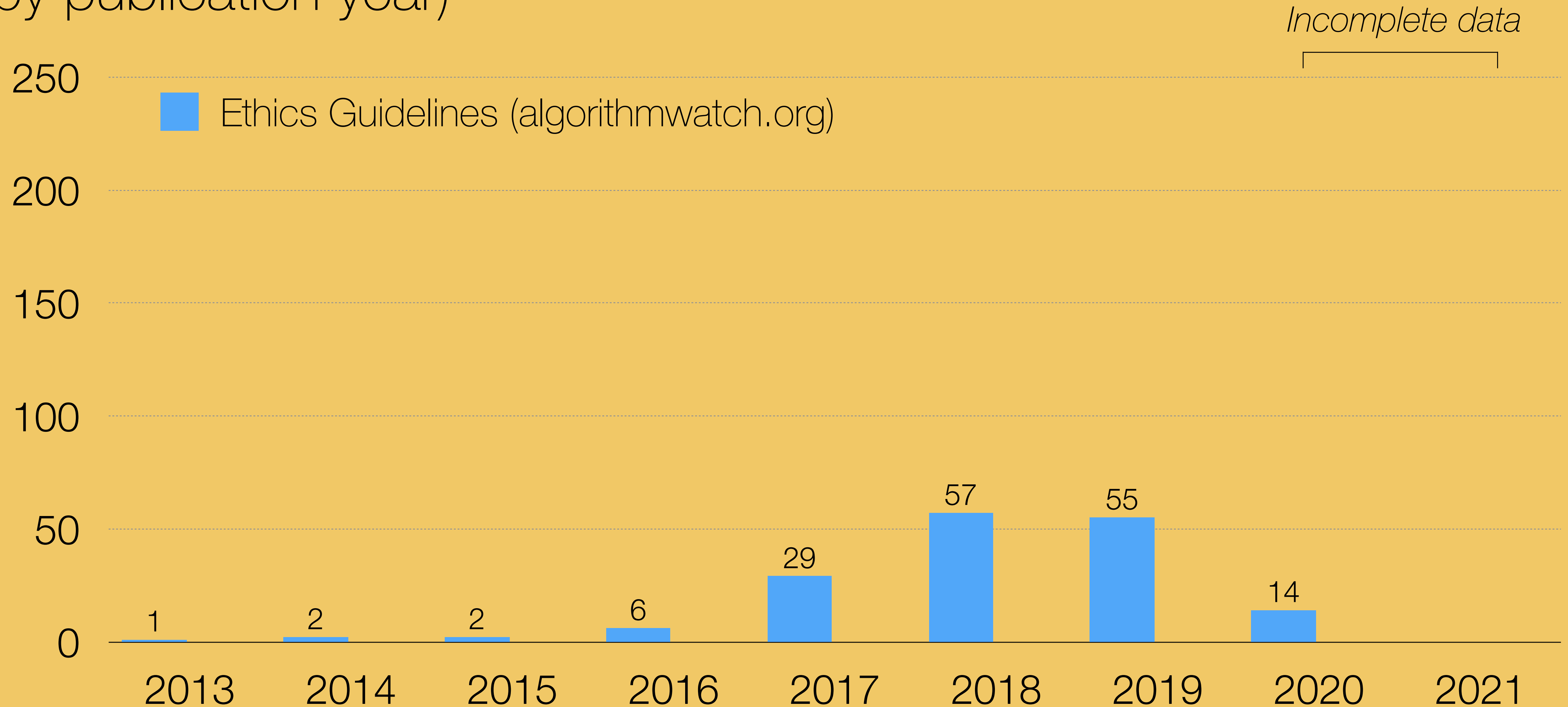
*(Emphasis added)*

# 2021

## National Artificial Intelligence Initiative

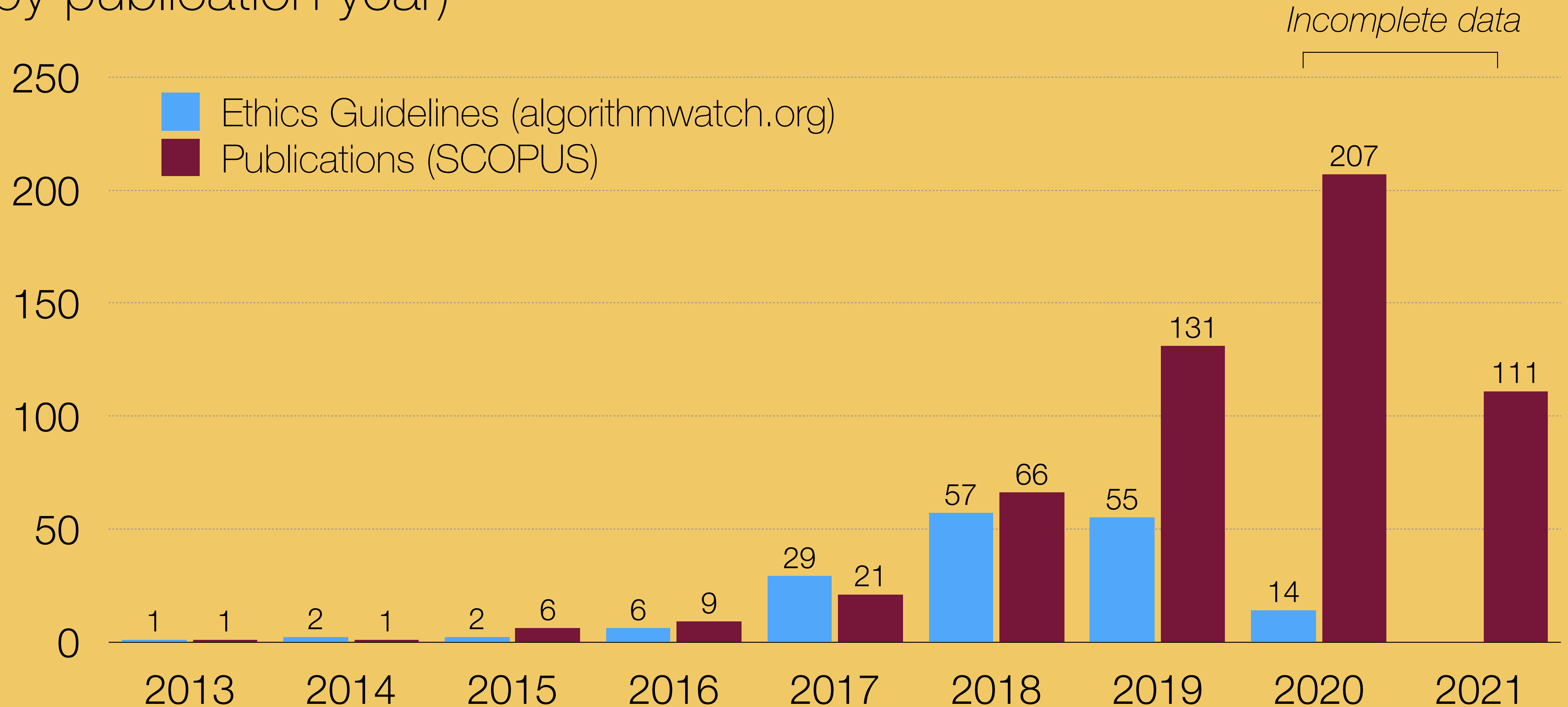Including "Advancing Trustworthy AII"

US federal Government

https://www.ai.gov/

# AI Ethics Guides and Publications

(by publication year)

*Incomplete data*



Legend: ■ Ethics Guidelines (algorithmwatch.org)

| Year | Count |
|------|-------|
| 2013 | 1 |
| 2014 | 2 |
| 2015 | 2 |
| 2016 | 6 |
| 2017 | 29 |
| 2018 | 57 |
| 2019 | 55 |
| 2020 | 14 |
| 2021 | |

# AI Ethics Guides and Publications

(by publication year)



*Incomplete data*

Legend:
- Ethics Guidelines (algorithmwatch.org)
- Publications (SCOPUS)

| Year | Ethics Guidelines | Publications |
|------|-------------------|--------------|
| 2013 | 1 | 1 |
| 2014 | 2 | 1 |
| 2015 | 2 | 6 |
| 2016 | 6 | 9 |
| 2017 | 29 | 21 |
| 2018 | 57 | 66 |
| 2019 | 55 | 131 |
| 2020 | 14 | 207 |
| 2021 | | 111 |

SCOPUS search: TITLE ( ( "AI"  OR  "artificial intelligence"  OR  "machine learning"  OR  "deep learning" )  AND  ( ethic* ) )

# IEEE General Principles
## Commonalities with many other sets of principles

## Human Rights

A/IS* shall be created and operated to respect, promote, and protect internationally recognized human rights.

## Well-being

A/IS creators shall adopt increased human well-being as a primary success criterion for development.

## Data Agency

A/IS creators shall empower individuals with the ability to access and securely share their data, to maintain people's capacity to have control over their identity.

## Effectiveness

A/IS creators and operators shall provide evidence of the effectiveness and fitness for purpose of A/IS.

## Transparency

The basis of a particular A/IS decision should always be discoverable.

## Accountability

A/IS shall be created and operated to provide an unambiguous rationale for all decisions made.

## Awareness of Misuse

A/IS creators shall guard against all potential misuses and risks of A/IS in operation.

## Competence

A/IS creators shall specify and operators shall adhere to the knowledge and skill required for safe and effective operation.

*Autonomous/Intelligent Systems

Source: https://ethicsinaction.ieee.org/

Are ethical frameworks enough to ensure safe and beneficial development and applications of machine learning and other forms of AI?

# RISKS OF AI
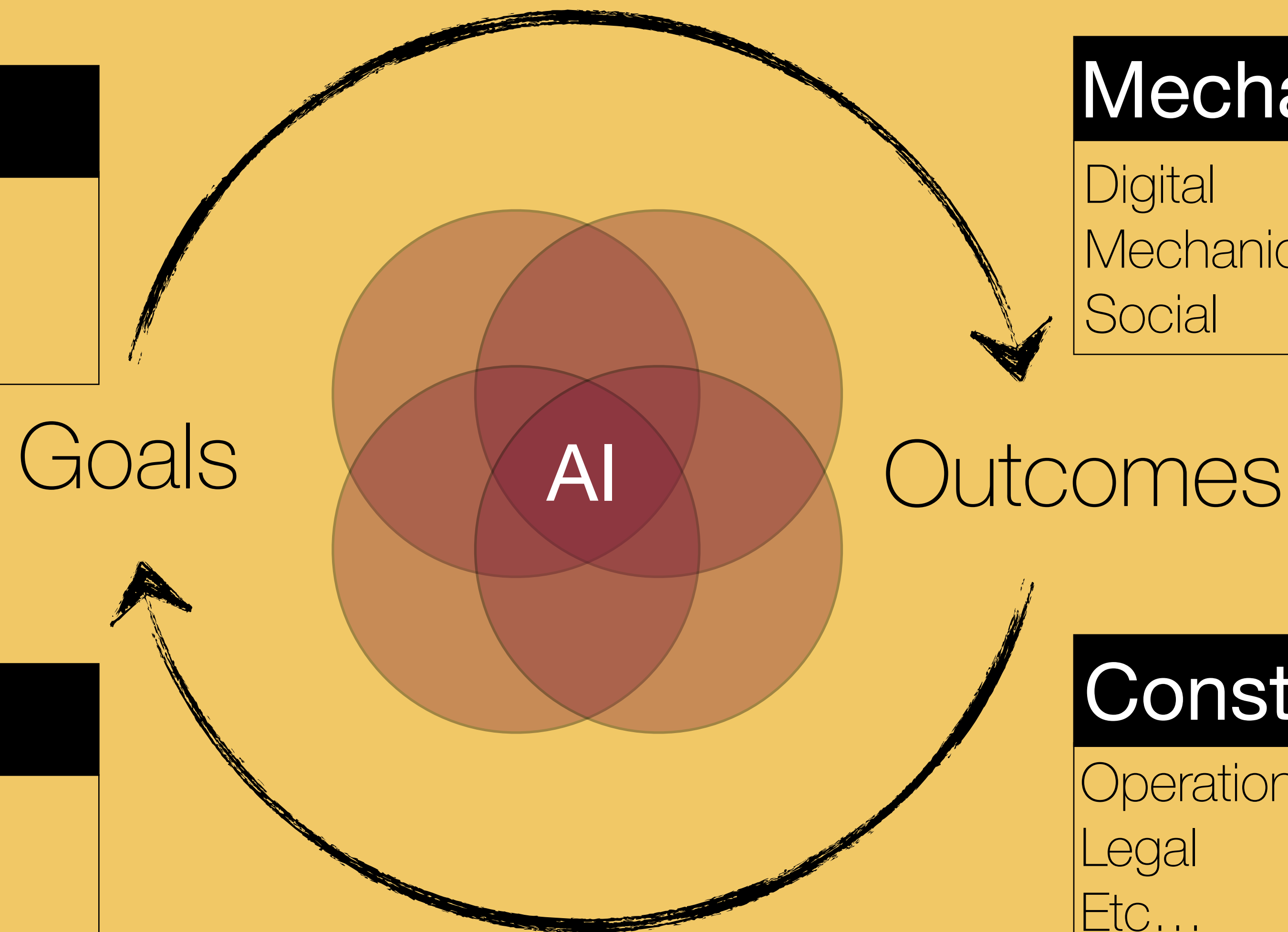
**Data**

| | |
|---|---|
| Curated | Closed |
| Uncurated | Open |
| Etc… | |

**Mechanisms**

| | |
|---|---|
| Digital | Behavioral |
| Mechanical | Political |
| Social | Etc… |

**"Knowledge"**

| | |
|---|---|
| Deterministic | Intuitive |
| Inferential | Complex |
| Etc… | |

**Constraints**

| | |
|---|---|
| Operational | Institutional |
| Legal | Ethical |
| Etc… | |

Goals

AI

Outcomes

RISK

Harm ←————————————→ Benefit

**Data**
| | |
|---|---|
| Curated | Closed |
| Uncurated | Open |
| Etc… | |

**Mechanisms**
| | |
|---|---|
| Digital | Behavioral |
| Mechanical | Political |
| Social | Etc… |

Goals

AI

Outcomes

**"Knowledge"**
| | |
|---|---|
| Deterministic | Intuitive |
| Inferential | Complex |
| Etc… | |

**Constraints**
| | |
|---|---|
| Operational | Institutional |
| Legal | Ethical |
| Etc… | |

Wrong ←————————————→ Right

ETHICS

# AI Ethics vs Risk Publications

(by publication year)

*Incomplete data*



**Ethics Publications (SCOPUS)**

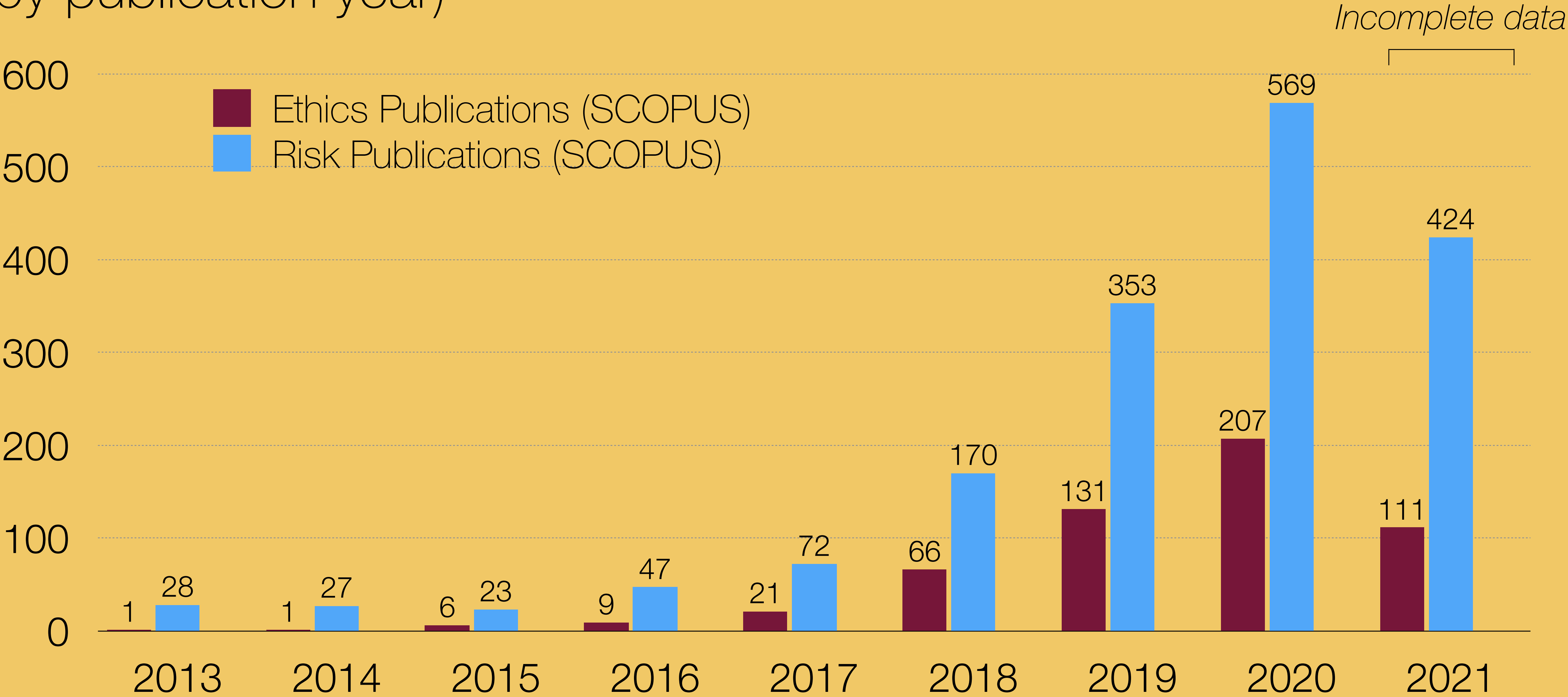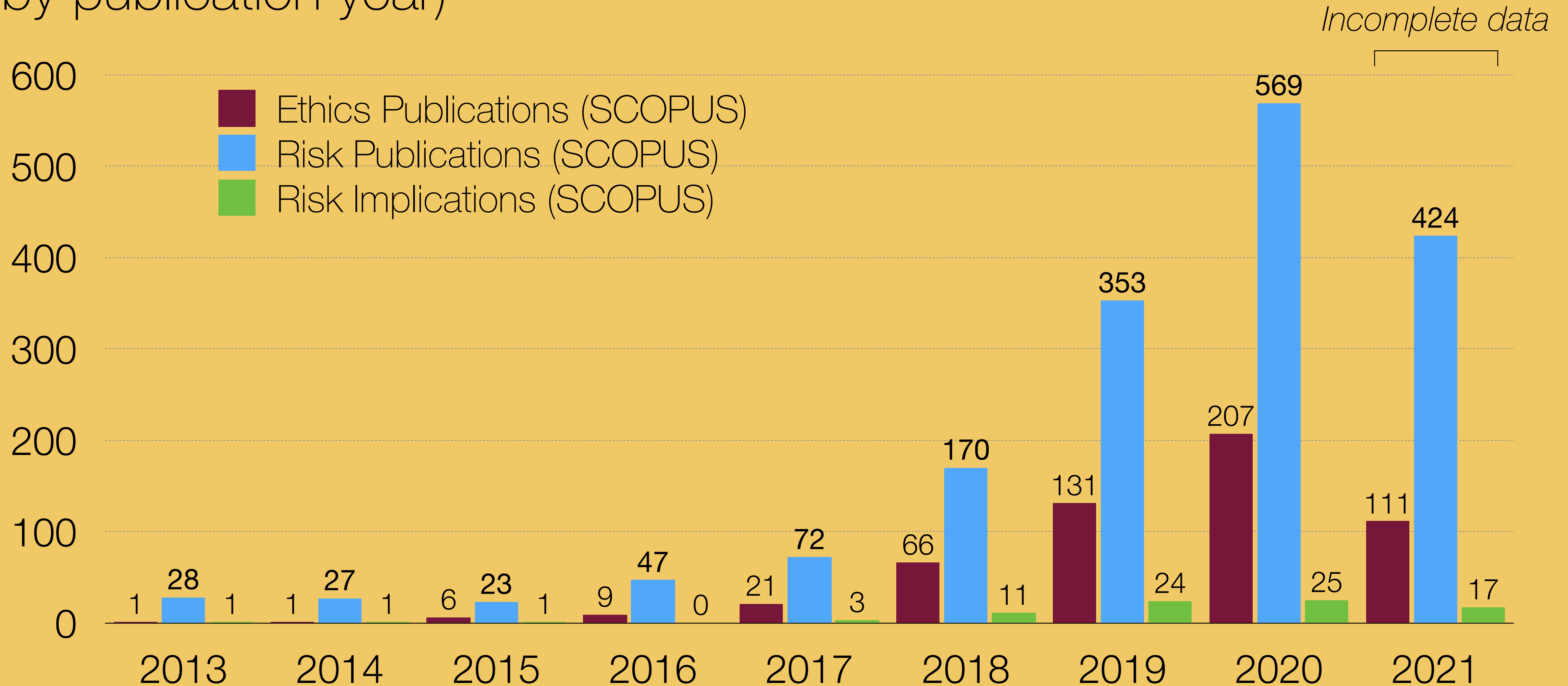| Year | Value |
|------|-------|
| 2013 | 1 |
| 2014 | 1 |
| 2015 | 6 |
| 2016 | 9 |
| 2017 | 21 |
| 2018 | 66 |
| 2019 | 131 |
| 2020 | 207 |
| 2021 | 111 |

SCOPUS search: TITLE ( ( "AI" OR "artificial intelligence" OR "machine learning" OR "deep learning" ) AND ( ethic* ) )
TITLE ( ( "AI" OR "artificial intelligence" OR "machine learning" OR "deep learning" ) AND ( risk* ) )

# AI Ethics vs Risk Publications

(by publication year)

*Incomplete data*

Legend:
- Ethics Publications (SCOPUS)
- Risk Publications (SCOPUS)

| Year | Ethics | Risk |
|------|--------|------|
| 2013 | 1 | 28 |
| 2014 | 1 | 27 |
| 2015 | 6 | 23 |
| 2016 | 9 | 47 |
| 2017 | 21 | 72 |
| 2018 | 66 | 170 |
| 2019 | 131 | 353 |
| 2020 | 207 | 569 |
| 2021 | 111 | 424 |

SCOPUS search: TITLE ( ( "AI"  OR  "artificial intelligence"  OR  "machine learning"  OR  "deep learning" )  AND  ( ethic* ) )
TITLE ( ( "AI"  OR  "artificial intelligence"  OR  "machine learning"  OR  "deep learning" )  AND  ( risk* ) )

# AI Ethics vs Risk Publications

(by publication year)

*Incomplete data*



Legend:
- ■ Ethics Publications (SCOPUS)
- ■ Risk Publications (SCOPUS)
- ■ Risk Implications (SCOPUS)

| Year | Ethics | Risk | Risk Implications |
|------|--------|------|-------------------|
| 2013 | 1 | 28 | 1 |
| 2014 | 1 | 27 | 1 |
| 2015 | 6 | 23 | 1 |
| 2016 | 9 | 47 | 0 |
| 2017 | 21 | 72 | 3 |
| 2018 | 66 | 170 | 11 |
| 2019 | 131 | 353 | 24 |
| 2020 | 207 | 569 | 25 |
| 2021 | 111 | 424 | 17 |

SCOPUS search: TITLE ( ( "AI"  OR  "artificial intelligence"  OR  "machine learning"  OR  "deep learning" )  AND  ( ethic* ) )
TITLE ( ( "AI"  OR  "artificial intelligence"  OR  "machine learning"  OR  "deep learning" )  AND  ( risk* ) )

18

# RISK INNOVATION

# RISK IS …

… the probability of harm occurring from an action or situation

# INNOVATION IS …

… The process of translating an idea or invention into a good or service that creates value for which customers will pay

# INNOVATION IS …

CREATIVITY

… The process of translating an idea or invention into a good or service that creates value for which customers will pay

VALUE

PRODUCTS

# RISK INNOVATION IS …

A way of approaching risk that leads to new knowledge, understanding, and capabilities, and translates these into products, tools, or practices that protect and grow societal, environmental, economic, and other value

23

# RISK INNOVATION IS ...

CREATIVITY

A way of approaching risk that leads to new knowledge, understanding, and capabilities, and translates these into products, tools, or practices that protect and grow societal, environmental, economic, and other value

PRODUCTS

VALUE

RISK INNOVATION
NEXUS

# RISK INNOVATION NEXUS

Connecting ethical and responsible innovation with value growth

25

RiskInnovation.org

**ENTERPRISE**

Research program,
business, organization,
ERC, etc

# STAKEHOLDERS

**CUSTOMERS**
Funders, businesses, researchers, developers, users, policy makers etc.

**INVESTORS**
Venture capital, shareholders, funders, etc.

**ENTERPRISE**
Research program, business, organization, ERC, etc

**COMMUNITIES**
Researchers, interest-groups, civil society, publics, etc.

**STAKEHOLDERS**

**CUSTOMERS**
Funders, businesses, researchers, developers, users, policy makers etc.

**INVESTORS**
Venture capital, shareholders, funders, etc.

**ENTERPRISE**
Research program, business, organization, ERC, etc

**COMMUNITIES**
Researchers, interest-groups, civil society, publics, etc.

**THREATS**

**VALUE**

# ORPHAN RISKS

Organizational Values & Culture

Bad Actors

Geopolitics

Organizations & Systems

Standards

Governance & Regulation

Reputation & Trust

Social Justice & Equity

Worldview

Privacy

Social & Ethical Factors

Product Lifecycle

Black Swan Events

Co-opted Tech

Unintended Consequences
of Emerging Technologies

Ethics

Perception

Health & Environment

Social Trends

Loss of Agency

Intergenerational Impacts

30

Maynard (2018) It's time for tech startups and their
funders to take "orphan risks" seriously. Medium

# Hypothetical Example:

Using AI-based social media agents to reach herd immunity



**Data**

| | |
|---|---|
| Curated | Closed |
| Uncurated | Open |
| Etc… | |

**Mechanisms**

| | |
|---|---|
| Digital | Behavioral |
| Mechanical | Political |
| Social | Etc… |

Goals

AI

Outcomes

**Knowledge**

| | |
|---|---|
| Deterministic | Intuitive |
| Inferential | Complex |
| Etc… | |

**Constraints**

| | |
|---|---|
| Operational | Institutional |
| Legal | Ethical |
| Etc… | |

# Orphan Risk Landscape

GOAL: Using AI-based social media agents to reach herd immunity (Private company)

VALUE

**ENTERPRISE**

**INVESTORS**

**CUSTOMERS**

**COMMUNITIES**

ORPHAN RISKS

Social &
Ethical Factors

Unintended Consequences
of Emerging Technologies

Organizations &
Systems

# Orphan Risk Landscape

GOAL: Using AI-based social media agents to reach herd immunity (Private company)

**VALUE**

### ENTERPRISE

Behavior change at scale

Versatile tech platform

Profit

### INVESTORS

Product that delivers on its promise

Trustworthiness

High return on investment

### CUSTOMERS

Significantly reduced social & economic impacts of infectious agent

No legal/regulatory liability

Public support and recognition

### COMMUNITIES

Autonomy

Transparency

Inclusion

ORPHAN RISKS

| Social & Ethical Factors | Unintended Consequences of Emerging Technologies | Organizations & Systems |
|---|---|---|

# Orphan Risk Landscape

GOAL: Using AI-based social media agents to reach herd immunity (Private company)



## VALUE

### ENTERPRISE
Behavior change at scale

Versatile tech platform

Profit

### INVESTORS
Product that delivers on its promise

Trustworthiness

High return on investment

### CUSTOMERS
Significantly reduced social & economic impacts of infectious agent

No legal/regulatory liability

Public support and recognition

### COMMUNITIES
Autonomy

Transparency

Inclusion

| Privacy | Social Trends | | Loss of Agency | | Governance & Regulation | Organizational Values & Culture |
| Worldview | | | | | Standards | |
| Perception | Ethics | | Co-opted Tech | Loss of Agency | Governance & Regulation | Reputation & Trust |
| | | | | | Geopolitics | |
| Social Trends | Perception | | Health & Environment | Co-opted Tech | Governance & Regulation | Reputation & Trust |
| | Worldview | | | | Standards | |
| Privacy | Social Justice & Equity | | Co-opted Tech | Loss of Agency | Reputation & Trust | Bad Actors |
| Ethics | | | | | | |

Social & Ethical Factors

Unintended Consequences of Emerging Technologies

Organizations & Systems

34

# RISK INNOVATION PLANNER

Short, iterative orphan risk audits and responses

One of a number of tools available through the Risk Innovation Nexus

Adaptable to multiple contexts

https://riskinnovation.org/services/risk-innovation-planner/

Key threats to value that need to be addressed if social good is to be realized
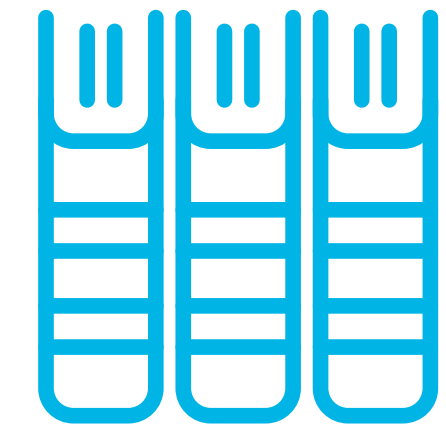
Ethics

Social Trends

Perception

Privacy

Worldview

Governance & Regulation

Loss of Agency

Reputation & Trust

Co-opted Tech

35

**Loss of Agency**

People are highly manipulatable

Human-human manipulation is constrained by a level(ish) playing field

Machine-human manipulation is not constrained in the same ways

Machines can potentially be taught (or learn) to take advantage of human heuristics and cognitive biases

This might be a beneficial thing for the future of humanity …

… or it could be really harmful!

Loss of
Agency

How do we navigate
a machine-mediated
future of cognitive
and behavioral
asymmetry?

People are highly manipulatable

Human-human manipulation is
constrained by a level(ish) playing
field

Machine-human manipulation is
not constrained in the same ways

Machines can potentially be
taught (or learn) to take advantage
of human heuristics and cognitive
biases

This might be a beneficial thing for
the future of humanity …

… or it could be really harmful!

# KEY TAKEAWAY:

Socially beneficial and responsible development and use of "AI" requires new thinking around value and risk, as well as ethics

# PROFESSOR ANDREW MAYNARD

Associate Dean
College of Global Futures
Arizona State University

Email: andrew.maynard@asu.edu
Twitter: @2020science
Medium: medium.com/edge-of-innovation
Web: andrewmaynard.net