

The State of Ethical AI Frameworks

Making sense of worldwide initiatives in Ethical AI

Presented by:

**Kathleen Walch
Ronald Schmelzer
Megan Yarbrough**

About Cognilytica

- Cognilytica is an **AI-focused analyst and advisory firm**
- We are a small business, based in Maryland
- Market research, advisory & guidance on **AI, ML, & Cognitive Technology**
- Provides **role-specific education** on AI, ML, and emerging technology
- Focused on **enterprise and public sector adoption** of AI, ML, and Cognitive Technology
- Kathleen Walch and Ron Schmelzer are **Principal Analysts and Managing Partners** and Megan Yarbrough is a **Research Associate** at Cognilytica.
- Produce the popular **AI Today podcast** and **AI communities including AI in Government and Data for AI**



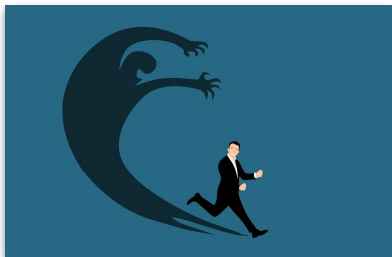
Some Caveats

- We are not ethicists nor experts in technology ethics
- Definitions are imprecise and often conflicting
- No widespread consensus for terminology or definitions
- Boundaries between concepts are fuzzy
- Many concepts are highly culture and context specific



... however, organizations have many of the same challenges, and despite those issues, they are making decisions and creating frameworks...

Tackling the Fears and Concerns of AI (And Big Data)



AI Fears	AI Concerns
<ul style="list-style-type: none">● Intelligent systems will take over the world● Robots will take my job● I'll lose control over my privacy and data● Too much data in too few hands● I'll live in a surveillance state● The robots will kill me	<ul style="list-style-type: none">● Lack of transparency in algorithmic decision-making● Bad actors doing bad things with AI● AI systems vulnerable to tampering and data security issues● Systems are susceptible to bias● Laws not keeping up with technology● How can I really trust these systems?

Ethical AI: “Right vs. Wrong”

AI Ethics: “AI ethics is a set of values, principles, and techniques that employ widely accepted standards of right and wrong to guide moral conduct in the development and use of AI technologies”¹

Bad Machines	Bad People Doing Bad Things	Bad Practices	Bad Visibility
<ul style="list-style-type: none"> ● Threatening lives ● Threatening freedoms ● Threatening control ● Threatening dignity ● Threatening environment ● Not in best interests of humanity 	<ul style="list-style-type: none"> ● Violating our laws ● Violating our trust ● Violating our privacy ● Violating our lives 	<ul style="list-style-type: none"> ● Lack of Safety ● Lack of accountability ● Lack of positive purpose ● Lack of care for workforce disruption 	<ul style="list-style-type: none"> ● Limited visibility into data and processes ● Limited disclosure ● Limited consent ● Limited visibility into algorithmic behavior ● Limited repeatable, consistent processes
			

Some Inspiration from the Past (1942)

First Law

- A robot may not injure a human being or, through inaction, allow a human being to come to harm.

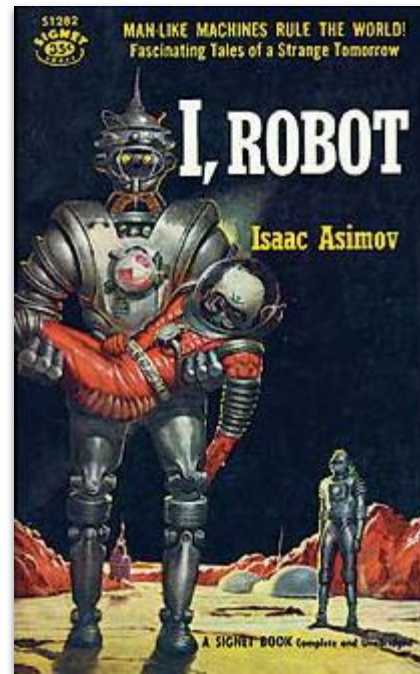
Second Law

- A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.

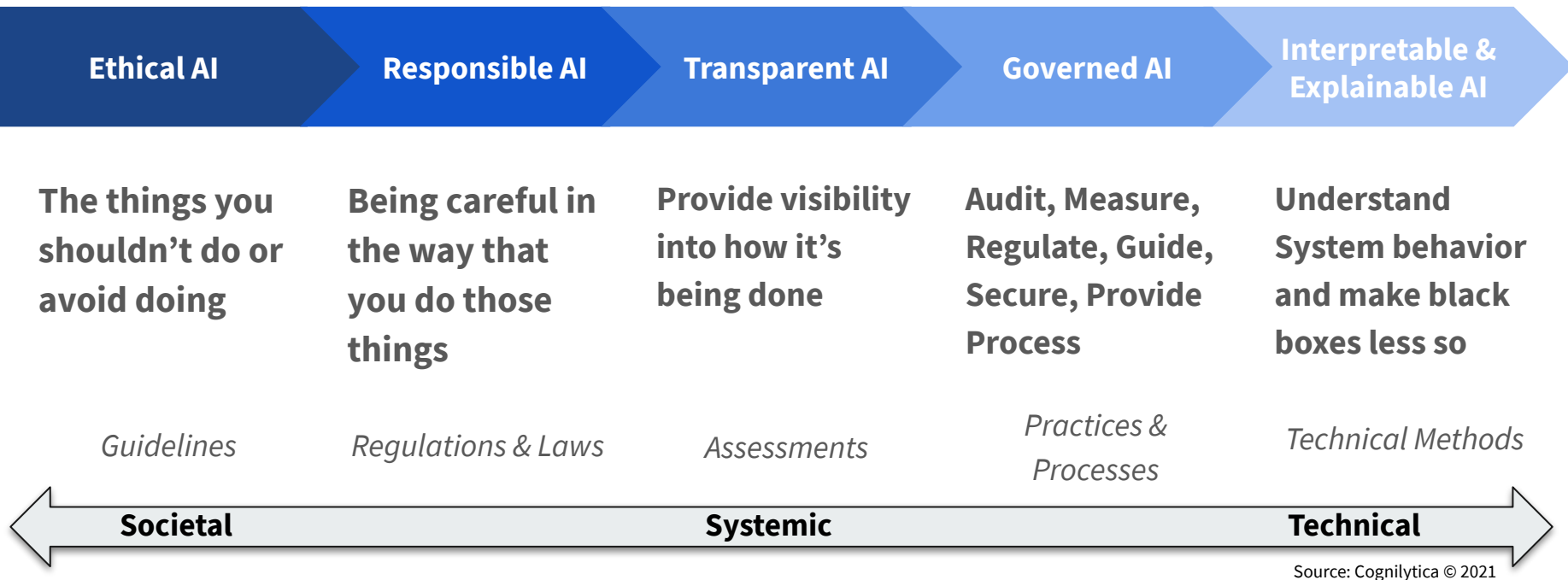
Third Law

- A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.

But this is clearly not enough. And it's also science fiction.



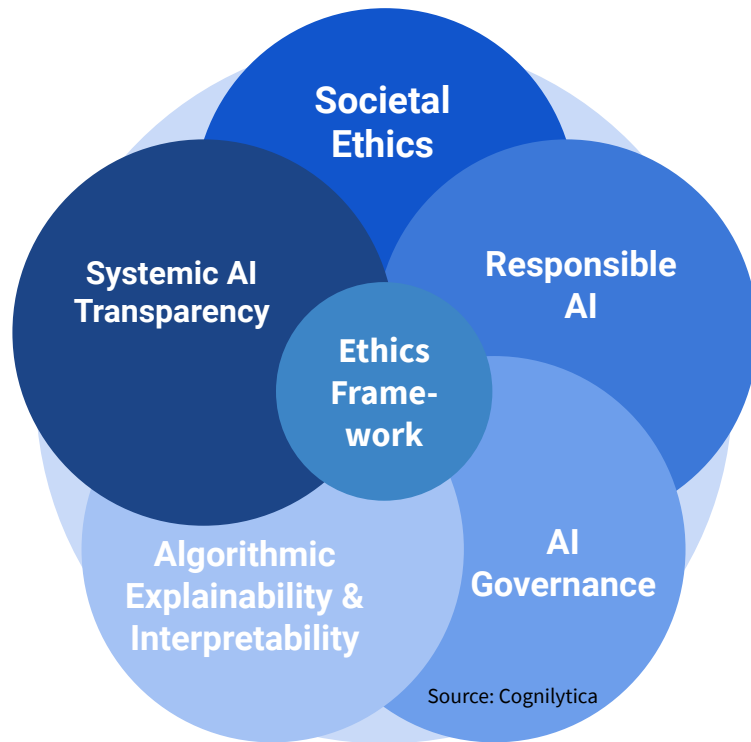
Ethical AI Concepts: A Spectrum



Ethical AI Frameworks: A Bit Of Everything

- Cognilytica analyzed over **60 National, Organizational, and Corporate Ethical AI Frameworks**
- Mostly a ***mish-mash of ethical concepts***
- The words they use often ***don't match the ethical concept***
- There are ***no comprehensive ethical AI frameworks***
- The ***gaps in ethical AI frameworks*** are eye opening
- ***Mix of recommendations*** for technology users and technology creators

We analyzed them all, normalized the terminology and compared them against each other.



Societal Ethics: The Main Elements

★ Human Values

- Machine-based systems should exhibit the same values that we have as humans. Do no harm.

★ Dignity

- AI systems should not treat humans as machines.

★ Fairness

- AI systems should not favor one group over another

★ Diversity & Inclusion

- AI systems should be built for and incorporate data from the breadth of humanity

★ Bias & Discrimination

- AI systems shouldn't further bias or discrimination

★ Freedom & Agency

- AI systems shouldn't limit human choice or freedom of action.

★ Human Benefit

- AI systems should be built for the benefit of the widest group of humanity, and not for the benefit of a few.

★ Human Control

- AI systems should never operate without humans in control

★ Respect of the Environment

- AI systems should take care not to abuse the environment. Do no harm.

Very complete, combine with 2 other Canadian efforts

Not only provides principles,
but methods to ensure their
application

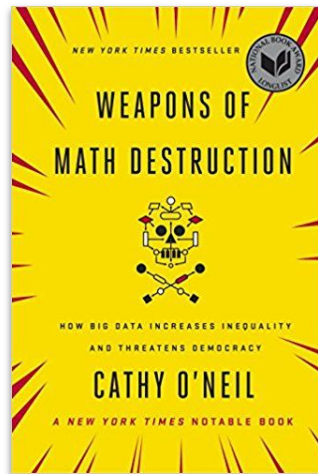
[illegible]

“Harmony & Friendship; Share responsibilities; Tolerance and sharing; Fairness and justice”

Surprisingly weak in societal ethics.
“Equitable (Bias only); Governable (Human Control)”

What is Responsible AI?

- **Responsible AI:**
 - “AI solutions that help maintain individual trust and minimize privacy invasion. Responsible AI places human (e.g., end-users) at the center”²
- Just because you can do it, even ethically, how can you do it the **right** way?
- How can AI be done in a way that doesn’t violate people’s trust?
- How can misusing AI technology be avoided?
 - **“Don’t misuse or abuse AI technology”**
- Examples of Responsible AI Concerns
 - Facial recognition as responsible vs. irresponsible AI
 - Responsible vs. Irresponsible Algorithmic decision-making
 - Development of invasive user profiles
- **What is a Responsible AI Framework?**
 - Set guidelines for the proper use of AI and guardrails to prevent abuse or misuse



Responsible AI: The Main Elements

★ Positive Purpose

- AI systems must be built for some positive purpose

★ Safety & Security

- AI systems should be safe and secure

★ Trust

- AI systems should not violate human trust or cause people to mistrust entities

★ Human Accountability

- Human individuals should be identified who are responsible and accountable for the behavior and operation of the AI systems

★ Privacy

- AI systems should not violate the privacy of humans or impose state-wide surveillance

★ Misuse, Abuse, and Compliance with Laws

- Humans should not misuse AI systems for any criminal or non-law abiding purpose, or to use AI to circumvent laws or regulations.

★ Lethal Autonomous Weapons

- AI systems should not be built to enable lethal autonomous weapons

★ Workforce Disruption

- AI systems should not be built that have as an intentional goal the mass replacement of human workers, disruption to economies

“Ethical Framework for a Good AI Society”

The most complete with respect to Responsible AI principles

Multinational and Org Efforts beat Country Efforts:

Countries are not paying enough attention to Responsible AI uses. Especially around AI applied to weapons.

[illegible]

Workforce Disruption:

Only a few frameworks talk about avoiding workforce disruption as an ethical goal.

China; India; OECD; European and African.

NSF Program on Fairness in Artificial Intelligence in Collaboration with Amazon (FAI):

Details needed on Responsible AI

Not really an ethical AI framework, but it could be...

What is AI System Transparency?

- **AI Transparency:**
 - *Visibility into all the aspects of what went into building an AI system so users can understand the full context of how an AI system is built and used*
- *How do I know what went into making this AI system so I can trust it?*
 - How do I know how the AI models are being built, what data has been used to create the models, what data was being included / excluded, and aspects of bias?
- **Conflicting definitions of transparency**
 - **Transparency as Interpretability / Explainability? “Algorithmic Transparency”**
 - That’s algorithmic transparency, which has to do with how interpretable or explainable the algorithm is
 - Detailed separately.
 - **Transparency as Visibility into how the system was built? “Systemic Transparency”**
 - It’s even more important to know what data went into making that model in the first place.
 - ***Especially if it’s not working how you think it should!***
- *“The ImageNet test set has an estimated label error rate of 5.8%”⁵*
 - A mushroom is labeled a spoon, a frog is labeled a cat, a high note from Ariana Grande is labeled a whistle.
 - Who cares if the model explains why it said "turn right" if the map was wrong
 - ***Without transparency, no visibility!***



You can trust me!

**Error-riddled data sets
are warping our sense of
how good AI really is**



AI Systemic Transparency: The Main Elements

★ AI System Transparency

- AI systems should provide visibility into the data and components of the system with their configuration that is used to generate results
- Human decisions on the operation, versioning, development, and use of the AI system should be disclosed and open.

★ Bias Measurement & Mitigation

- AI systems should provide a means to constantly measure bias of various sources and provide means to mitigate any bias detected.

★ Open Systems

- AI systems as a whole should use open source technology with the mechanism by which the system operates visible to all

★ Disclosure & Consent

- Organizations should disclose when AI systems are being used and when humans are interacting with AI systems
- AI systems should provide a means for humans to consent out of interaction with AI systems, being included in AI models, or otherwise being impacted by the AI system.

Need a way to measure.

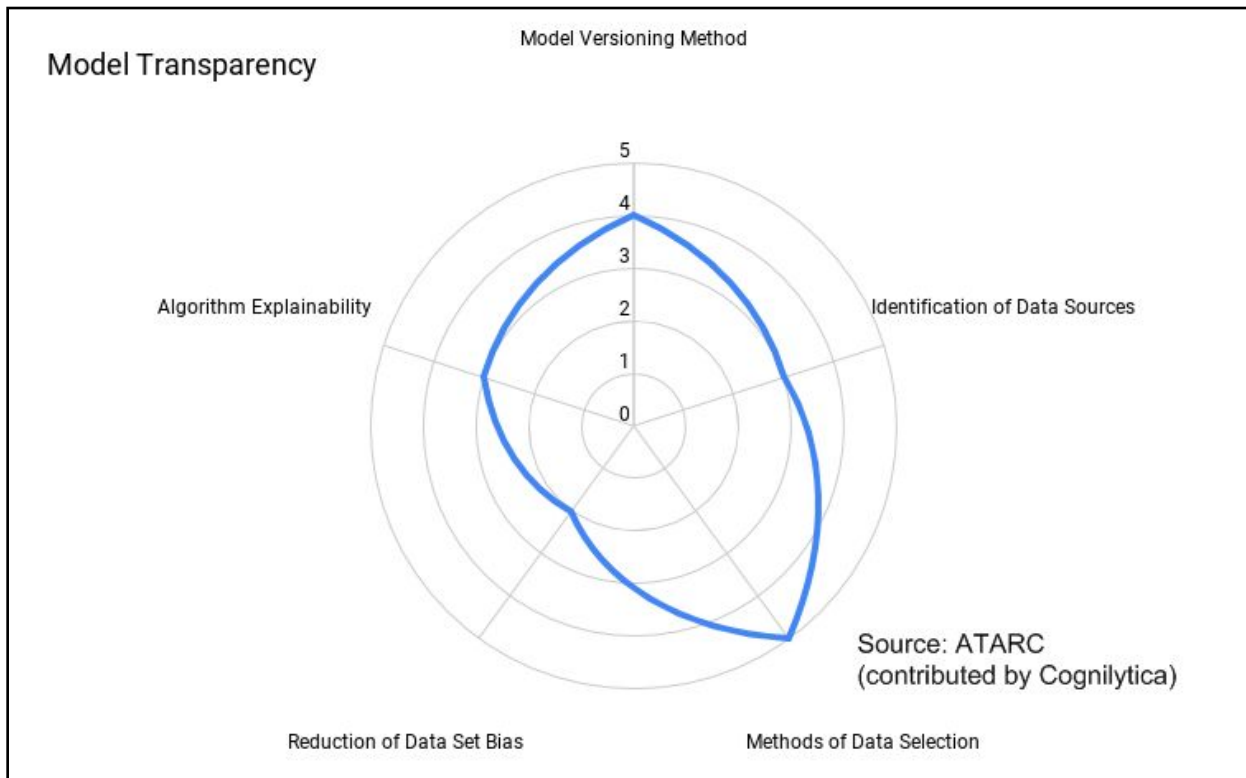
Missing disclosure & consent, visibility into system operation, measurement of bias, and open systems

[illegible]

Frameworks originating within university and research settings most often mention openness in research, data, shared platforms, and systemic AI transparency

Completely absent of any systemic transparency principles

Assessing AI Transparency: Model Assessment Approach



AI Governance: The Main Elements

★ System Auditability

- AI systems should provide ways to audit all aspects of operation and behavior

★ Contestability

- AI systems need to provide ways to contest or appeal AI decisions for human review.

★ Risk Assessment & Mitigation

- Organizations need established methods to assess ongoing risk to AI systems and identified means to mitigate those risks

★ System Monitoring & Quality

- Systems operating within ethical guidelines need to make sure that AI systems are always operating within acceptable performance, usage, and other parameters

★ Education & Training

- Ethical frameworks should require those who are involved in AI system creation or use to be trained in their proper development and use.

★ Regulation & Certification

- AI systems should comply with requirements of regulatory bodies and third-party certifications with regular third-party audits and certifications of ethical operation

Addresses almost all of the points of AI Governance

Spend the most amount of time talking about regulation, certification, and third-party credentialing or accountability

Policy Frameworks & Guidelines		AI Governance		AI Governance		AI Governance		AI Governance	
Entity/Company	Objective	Initiative/Policy	Document/Link	Year/Ref	Link	Risk Assessment	Impact	Privacy	Security
ACM/EU	Regulation	International EU AI Ethics Approach	Artificial Intelligence Ethics	2019	https://www.acm.edu/ai-ethics-approach	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Amnesty	Organization	N/A	Ethical Framework for a Good AI	2018	https://www.amnesty.org/en/documents/2018/05/24/437667a8-3b9d-4941-b061-386f18c8e000/	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Australia	Government	Large Industry & AI Ethics Framework	AI Ethics Framework	2019	https://www.australiainnovation.gov.au/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Canada	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www150.ca/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
China	Government	Beijing Declaration on AI Ethics Principles	AI Ethics Principles	2019	https://www.bjnews.com.cn/ai-ethics-principles	AI Ethics	AI Ethics	AI Ethics	AI Ethics
China (Hong Kong)	Government	Healthy Development of AI Ethics Framework	AI Ethics Framework	2019	https://www.hkma.gov.hk/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Colombia	Government	President's Council on AI Ethics Framework	AI Ethics Framework	2019	https://www.colombia.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Costa Rica	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.costa-rica.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Croatia	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.croatia.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Denmark	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.denmark.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Euro Union	Multi-Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.euro-union.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
European Union	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.european-union.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
France	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.france.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Germany	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.germany.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Google	Multi-Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.google.com/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Harvard's Berkman Center for Internet & Society	Organization	AI Ethics Framework	AI Ethics Framework	2019	https://www.berkman-center.org/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Healthcare & AI	Organization	AI Ethics Framework	AI Ethics Framework	2019	https://www.healthcare-ai.org/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
India	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.india.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Indonesia	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.indonesia.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Japan	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.japan.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Malaysia	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.malaysia.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Mexico	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.mexico.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Netherlands	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.netherlands.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Norway	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.norway.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
OECD	Multi-Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.oecd.org/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Poland	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.poland.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Portugal	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.portugal.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Russia	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.russia.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Saudi Arabia	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.saudi-arabia.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
South Korea	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.south-korea.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Spain	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.spain.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Sweden	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.sweden.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Switzerland	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.switzerland.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Singapore	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.singapore.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
South Africa	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.south-africa.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
South Korea	Government	AI Ethics Framework	AI Ethics Framework	2019	https://www.south-korea.gov/ai-ethics-framework	AI Ethics	AI Ethics	AI Ethics	AI Ethics
Spain	Government	AI Ethics Framework							

The topic of education and training of the workforce on AI comes up frequently in ethical AI frameworks

A lot of talk about bias, but little governance effort
Only a few ethical AI frameworks require organizations to measure the extent to which AI systems exhibit bias

What is Interpretable AI and Explainable AI?

- Discomfort with “Black Boxes”
- **Explainable AI:**
 - **Can the model tell you how it arrived at a particular prediction / conclusion?**
 - “An explainable model is a ... method/technique to be able to peer into the black-box and understand how the model works”⁵
- **Interpretable AI:**
 - **If not, can you provide ways to understand the components of the decision-making or the key factors contributing to a result?**
 - “Interpretability is the degree to which a human can understand the cause of a decision” or “The degree to which a human can consistently predict the model’s result”³
- **“Root Cause Analysis” / “Failure Analysis”**
- Can we understand how algorithmic decisions are being made?
 - **“Algorithmic Transparency”**
- How do I know what the AI system is doing?
- Do I have visibility into how the systems are succeeding or failing and why?
- Is it possible to “debug” AI systems?



AI Explainability & Interpretability: The Main Elements

★ Understandability / Root Cause Explanations

- When AI systems fail to provide expected results, AI systems should always provide a human understandable means to understand the root cause of any failures. Explanations without necessarily algorithmic explanation.

★ Algorithmic Interpretability

- AI systems should provide a means to interpret AI results so that cause and effect can be understood, even with limited algorithmic explainability

★ Algorithmic Explainability

- AI systems should use algorithms that provide a direct means to explain how outcomes were arrived from input data

“Common sense explanations”

Focus is on understandability for the purposes of giving users disclosure and ability to contest results

Policy Frameworks & Guidelines			Algorithm Interpretability / Explainability		
Country/Region	Doc Type / Organization	Substance / Org	Framework / Initiative	Doc Title / Link	Relevant Concepts / Topics
ACTA	Organization	EMERGING TECH	Ethical Framework of Artificial Intel	2020	https://www.acta.org/
Australia	Government	N/A	Ethical Framework for a Good AI	2019	https://www.australiagov.au/ai-ethics-framework
Austria	Company	N/A	Responsible AI	2019	https://www.austrianai.at/en/ai-ethics
Canada	Government	University of the	The Moral Declaration for Big	2017	https://www.universityoftoronto.ca/ai-ethics
Canada	Government	N/A	Rules on Responsible use of AI	2019	https://www.canada.ca/en/innovation-and-science/ai-ethics
Canada	Government	N/A	Direction on Automated Decision	2019	https://www.canada.ca/en/innovation-and-science/ai-ethics
China	Government	Office of the	Regulatory Framework for AI	2020	https://www.china-ai-ethics.com/
China	Government	Beijing Academy	Big Data Ethics Principles	2019	https://www.bjia.ac.cn/ai-ethics
China	Government	Ministry of Science	Development of Responsible AI	2019	https://www.most.gov.cn/ai-ethics
China (Hong Kong)	Government	Office of the	Regulatory Framework for AI	2020	https://www.china-ai-ethics.com/
Colombia	Government	Presidential Council	Ethical AI Framework Draft	2020	https://www.colombia-ai-ethics.com/
Deep Mind	Company	N/A	Deep Mind Ethics	2017	https://www.deepmind.com/ai-ethics
Denmark	Government	Committee on	Use of AI for the Benefit of the People	2019	https://www.denmark-ai-ethics.com/
EU	Multi-Group	High-Level Expert	Group of Experts	2020	https://www.eu-ai-ethics.com/
European Union	Government	Commission on	AI Ethics Framework of ethical	2020	https://www.eu-ai-ethics.com/
European Union	Government	European Group	Statement on ethical	2019	https://www.eu-ai-ethics.com/
France	Government	High Council	AI Ethics Framework	2019	https://www.france-ai-ethics.com/
Future of Life Institute	Organization	Future of Life	Institute's AI Ethics	2019	https://www.fli-ai-ethics.com/
Germany	Government	German Ethics	Commission of the	2017	https://www.german-ai-ethics.com/
Global Partnership	Multi-Group	N/A	Principles for responsible	2019	https://www.gp-ai-ethics.com/
Google	Company	N/A	Principles for responsible	2019	https://www.google.com/ai-ethics
Hanover's Research	Organization	N/A	Principles of Artificial Intelligence	2019	https://www.hanover-ai-ethics.com/
Hamsterdam Data	Organization	N/A	A Framework for the Ethical Use	2019	https://www.hamsterdam-ai-ethics.com/
HAJAC	Organization	Consultation on	AI - Key Themes	2019	https://www.hajac-ai-ethics.com/
India	Government	N/A	AI Ethics Framework	2019	https://www.india-ai-ethics.com/
INEC	Organization	The IEEE Global	ETHICAL & AI-RELATED	2019	https://www.ieee-ai-ethics.com/
INEC	Government	N/A	AI Ethics and Human Rights	2019	https://www.inec-ai-ethics.com/
Japan	Government	High Council	AI Ethics Framework	2019	https://www.japan-ai-ethics.com/
Machin Intelligence	Organization	Japanese Society	for Artificial Intelligence	2019	https://www.machin-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government	Malta's AI Ethics	Framework	2019	https://www.malta-ai-ethics.com/
Malta	Government				

Checks all the boxes for understandability, interpretability, and explainability

Does a good job of identifying needs for explainability in an ethical AI Framework

“Agencies shall ensure that the operations and outcomes of their AI applications are sufficiently understandable by subject matter experts, users, and others, as appropriate.”

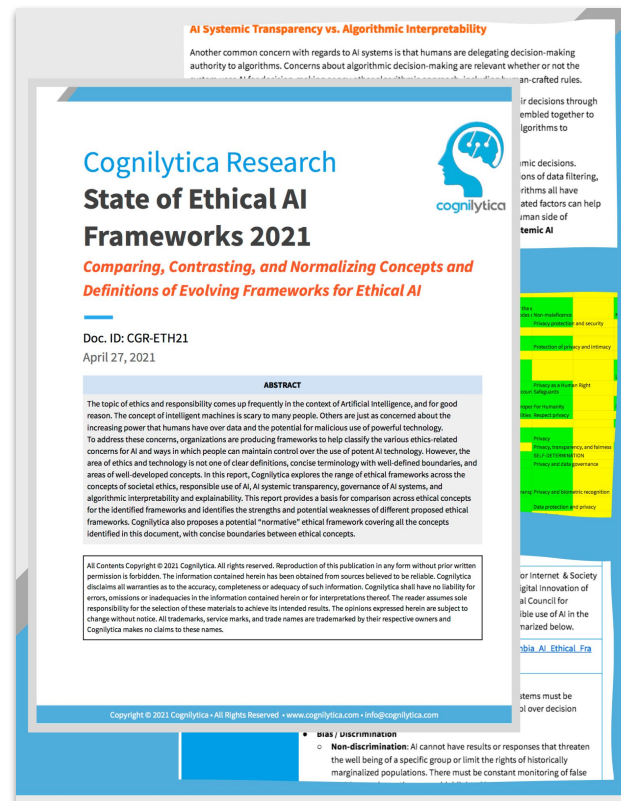
In Search of a Comprehensive Ethical AI Framework

... wait... we just presented it!

Societal Ethics	Responsible AI	Systemic AI Transparency	AI Governance	AI Explainability & Interpretability
<ul style="list-style-type: none"> ★ Human Values ★ AI for Human Benefit ★ Dignity ★ Fairness ★ Respect Diversity ★ Avoid Bias ★ Ensure Freedom & Human Agency ★ Keep the Human in Control ★ Respect the Environment 	<ul style="list-style-type: none"> ★ Positive Purpose ★ Law Abiding use of AI ★ Ensure Safety ★ Maintain Trust ★ Provide Human Accountability ★ Ensure Privacy ★ No Lethal Autonomous Weapons ★ Avoid Workforce Disruption 	<ul style="list-style-type: none"> ★ Visibility into AI Data ★ Bias Measurement & Mitigation ★ Open Systems & Human Decisions ★ Provide Disclosure and Respect Consent 	<ul style="list-style-type: none"> ★ Risk Assessment & Mitigation ★ System Auditability ★ Contestability ★ Continuous System Quality ★ External Regulatory Bodies & Third-party Certifications ★ Training 	<ul style="list-style-type: none"> ★ Algorithmic Explainability ★ Algorithmic Interpretability ★ Root Cause Explanations

Cognilytica's State of Ethical AI Frameworks 2021 Report

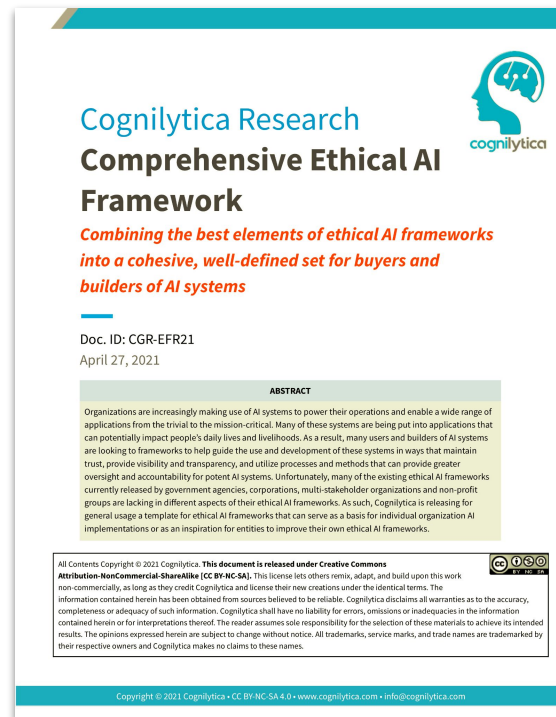
- **State of Ethical AI Frameworks 2021** report published April 2021
- Thorough evaluation of over **60 ethical AI Frameworks**
- Provides **normalized definitions and means for comparison**
- **Guided Questions** for buyers and builders of AI systems
- This report is 150+ pages with comprehensive overviews of each AI framework. It's available for Cognilytica research subscribers to download or for individual purchase.



Cognilytica Comprehensive Ethical AI Framework

- In addition, we produced *Cognilytica's Comprehensive Ethical AI Framework* as a free framework that serves as a guide for others.

Get a copy of Cognilytica's Comprehensive Ethical AI Framework



Thank You

Presenters:

Kathleen Walch, Ronald Schmelzer, Megan Yarbrough

Cognilytica - <http://www.cognilytica.com>

info@cognilytica.com

Appendix / Sources

Sources

1. Leslie, D. (2019). Understanding artificial intelligence ethics and safety: A guide for the responsible design and implementation of AI systems in the public sector. The Alan Turing Institute.
<https://doi.org/10.5281/zenodo.3240529>
2. Wang, Yichuan & Xiong, Mengran & Olya, Hossein. (2019). Toward an Understanding of Responsible Artificial Intelligence Practices. 10.24251/HICSS.2020.610.
3. Carvalho DV, Pereira EM, Cardoso JS. Machine Learning Interpretability: A Survey on Methods and Metrics. Electronics. 2019; 8(8):832. <https://doi.org/10.3390/electronics8080832>
4. <https://towardsdatascience.com/interperable-vs-explainable-machine-learning-1fa525e12f48>
5. “Error-riddled data sets are warping our sense of how good AI really is” -
<https://www.technologyreview.com/2021/04/01/1021619/ai-data-errors-warp-machine-learning-progress/>