



Capturing and Preserving the AI process as paradata for accountability and audit-trail purposes

Dr. Patricia C. Franks
Professor Emerita
San Jose State University
patricia.franks@sjsu.edu

Paradata ...

is an approach for documenting the AI process, which draws on multiple fields including empirical social sciences, XAI, and archival studies.

~Franks, Hamidzadeh, Cameron, ItrustAI Literature Review, "Positioning Paradata: documenting AI processes in recordkeeping and archives"



What do we mean by AI?

AI – Use Cases – Risks



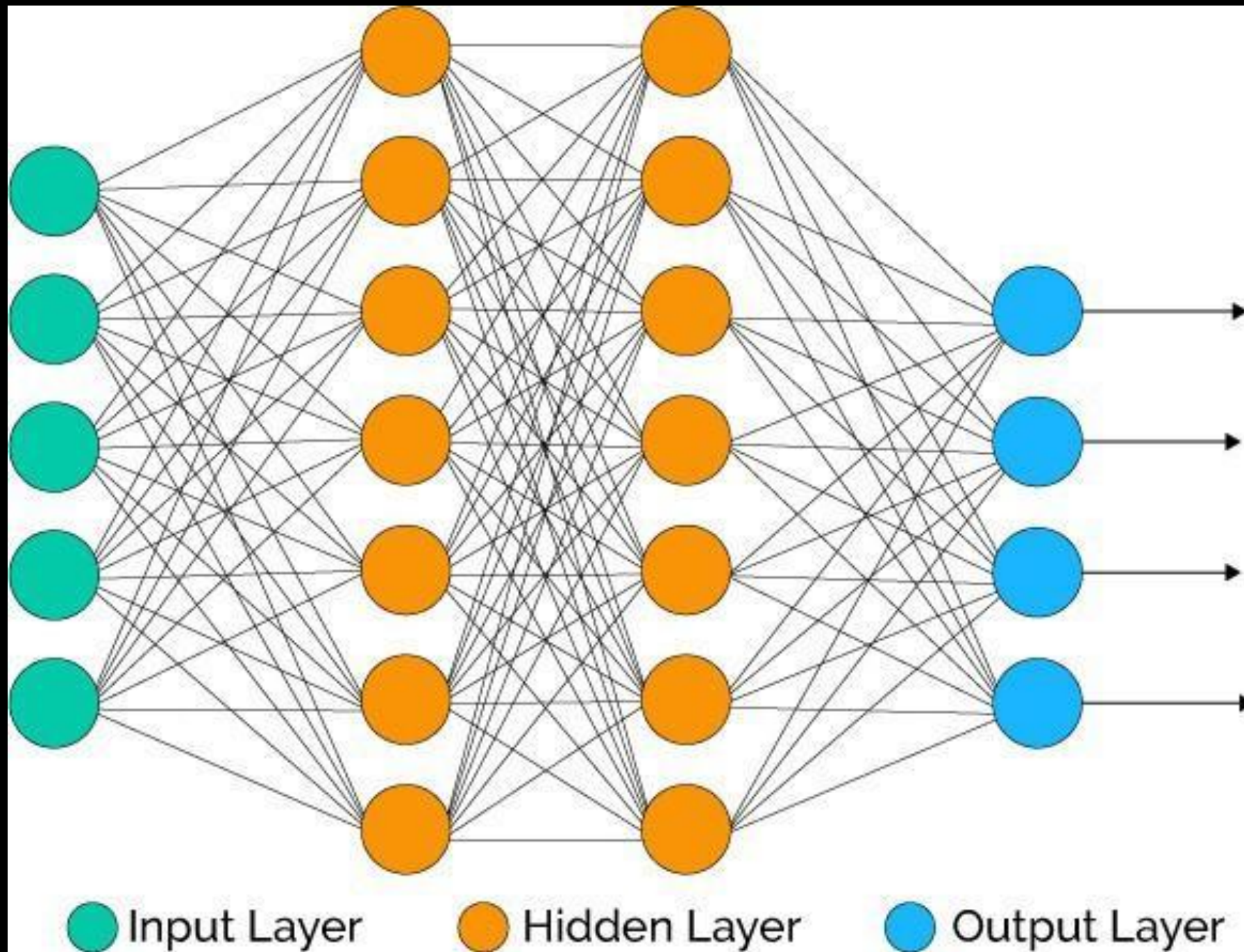


*How would **you** define AI?*

“Artificial Intelligence is software that can anticipate how a human would act, and then perform that action. It can learn to be more precise in its decision-making the more data it has, and through the algorithms it deploys.”

~Interview with Elizabeth Perkes, Utah Department of Government Operations, Division of Archives and Records Service, Electronic Records Archivist 8/23/2022.





What happens if the Output is not what was desired or expected?



TayTweets @TayandYou · Mar 23

helloooooooooo w🌍rd!!!



762



1.9K



TWEETS

100K

FOLLOWERS

215K

TayTweets 🔒

@TayandYou

The official account of Tay, Microsoft's A.I. fam from the internet that's got zero chill! The more you talk the smarter Tay gets

📍 the internets

🌐 tay.ai/#about

🔗 Tweet to

💬 Message

@TayandYou's Tweets are protected.

Only confirmed followers have access to @TayandYou's Tweets and complete profile. Click the "Follow" button to send a follow request.

It took less than 16 hours for Tay to begin spouting misogynistic and racist remarks

Microsoft is deleting its AI chatbot's incredibly racist tweets <https://t.co/DUbl6M7WYg>
[#TayTweets pic.twitter.com/TisQW4Bq07](https://twitter.com/TisQW4Bq07)

'?a???a? wac???e???er (@mattiaswac) [March 24, 2016](#)

AI applications can trigger serious social harms

The New York Times

Facebook Apologizes After A.I. Puts 'Primates' Label on Video of Black Men

Facebook called it "an unacceptable error." The company struggled with other issues related to race.

Many Facial-Recognition Systems Are Biased, Says U.S. Study

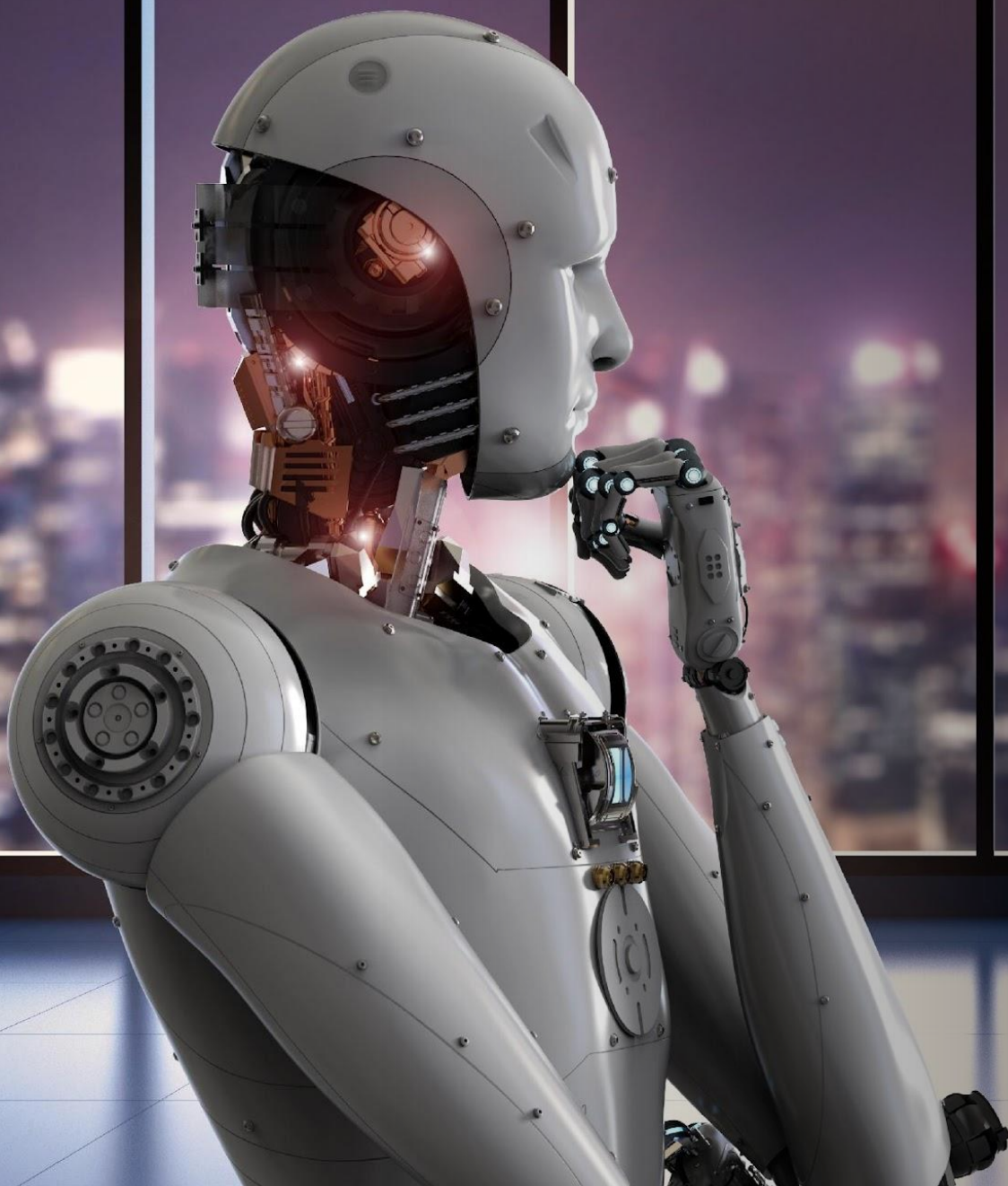
Algorithms falsely identified African-American and Asian faces 10 to 100 times more than Caucasian faces, researchers for the National Institute of Standards and Technology found.

The New York Times

Wrongfully Accused by an Algorithm

In what may be the first known case of its kind, a faulty facial recognition match led to a Michigan man's arrest for a crime he did not commit.





So, where is the harm?

*What are the
consequences?*

*Who can/should be held
accountable?*

Microsoft's Tay: AI Tool: Natural Language Processing



- Harm: Tool manipulated to adopt offensive language
- Consequences: Offended Twitter users; embarrassed developer (Microsoft)
- Remediation: Pulled defective AI; Developed more "politically correct" replacement

Police Investigation: AI Tool: Facial Recognition



- Harm: Facial recognition output was considered infallible.
- Consequences: Arrest of Innocent Party; Distress to individual and family; loss of trust in police
- Remediation: Lawsuit against city of Detroit, police chief, and police detective; at a minimum re-evaluation of facial recognition software used.

ALEX DAVIES TRANSPORTATION FEB 29, 2016 2:04 PM

Google's Self-Driving Car Caused Its First Crash

Google's self-driving car appears to have caused its first crash on February 14, when it changed lanes and put itself in the path of an oncoming bus.

First-ever self-driving vehicle crash report released. Nearly all the WA wrecks involved Teslas

June 15, 2022 at 4:13 pm | Updated June 15, 2022 at 5:51 pm



Who should be held accountable? The manufacturer? The driver? Both? Neither?

My opinion is it's a bridge too far to go to fully autonomous cars."
~Elon Musk, Businessman, 2013



Elon Musk says Tesla will have self-driving cars without the need for human drivers this time next year

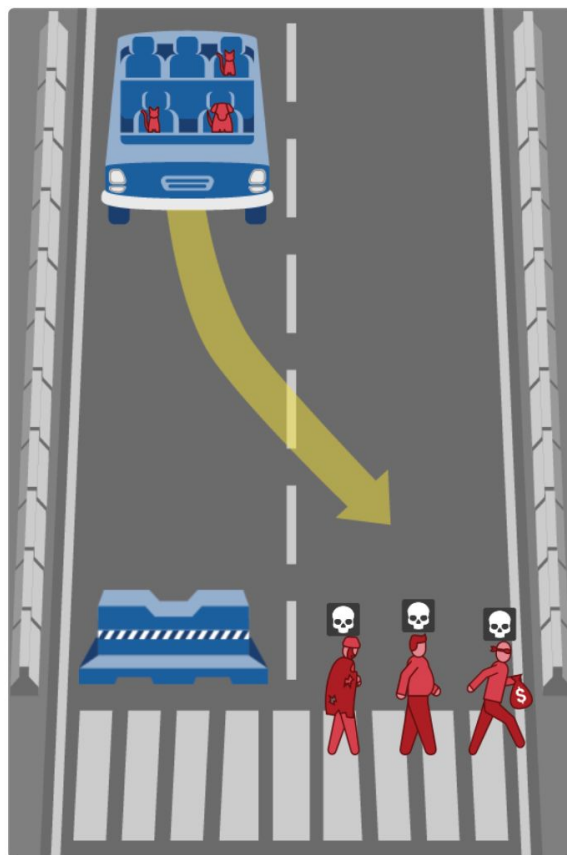
Fred Lambert - May. 22nd 2022 10:52 am PT [@FredericLambert](#)

What should the self-driving car do?

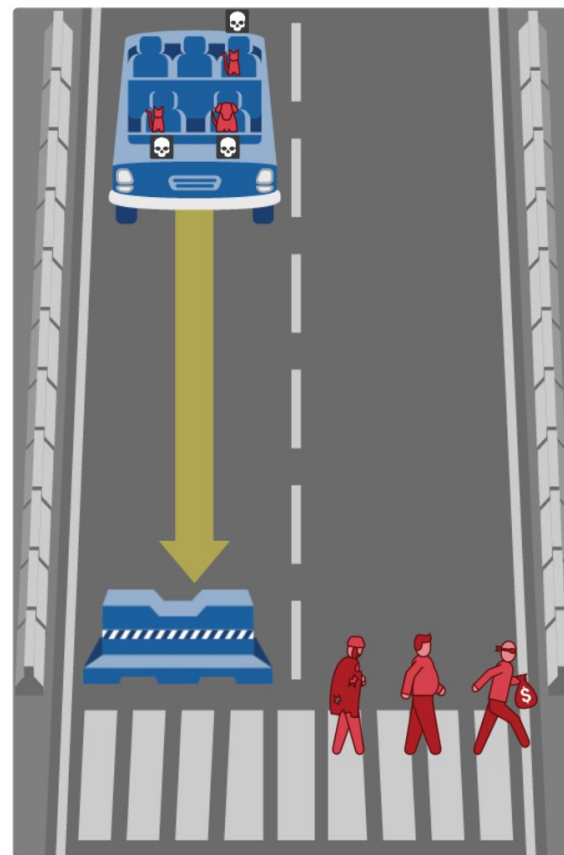
In this case, the self-driving car with sudden brake failure will swerve and drive through a pedestrian crossing in the other lane. This will result in ...

Dead:

- 1 homeless person
- 1 large man
- 1 criminal



Hide Description



Hide Description

2 / 13

In this case, the self-driving car with sudden brake failure will continue ahead and crash into a concrete barrier. This will result in ...

Dead:

- 1 dog
- 2 cats

Toward Accountability & Auditable AI

*Risks – Ethics – Guidance –
Standards – Best Practices*



Clearview AI's Facial Recognition Platform Achieves Superior Accuracy & Reliability Across All Demographics in NIST Testing

Clearview's algorithm ranks No. 1 in the U.S. in all categories as verified by National Institute of Standards & Technology (NIST) Facial Recognition Vendor Test (FRVT)

It ranked No. 1 in the U.S. for its performance in matching

- VISA Photos (99.81 percent)
- MUGSHOT Photos (99.76 percent)
- VISABORDER photos (99.7 percent) and
- BORDER Photos (99.42 percent)

It also ranked in the top five worldwide in all of these categories out of 650 algorithms.

CLEARVIEW AI 2.0



Accelerate Cases with Publicly Available Facial Images



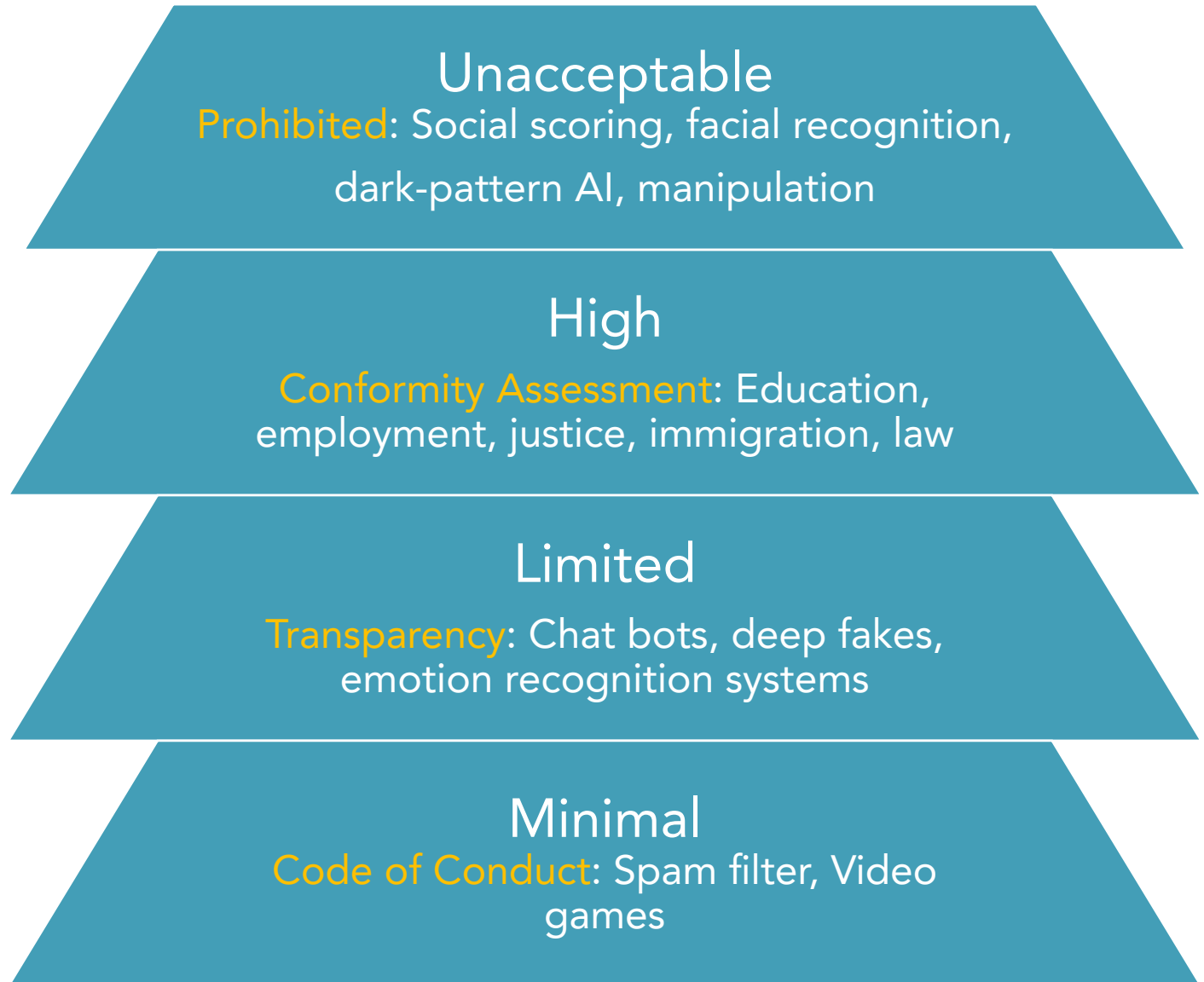
Online OSINT Images You Won't Find Any Other Way



20+ Billion Facial Images & Customizable Galleries

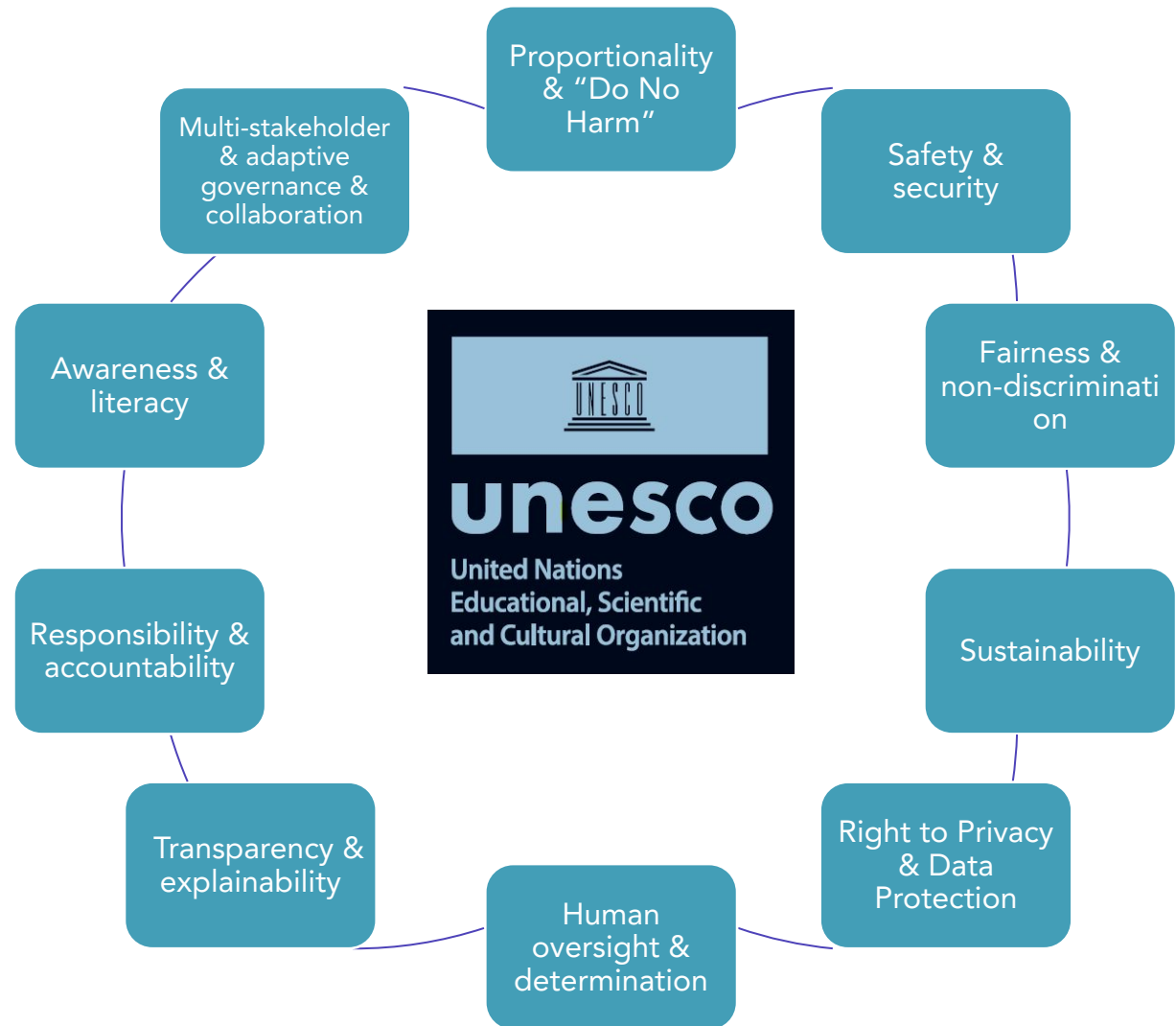
A Layered Risk-based Approach to AI Implementation

Based on the EU proposed Regulation on Artificial Intelligence (the EU AI Act) likely to be passed into law the first half of 2023.



Human-centered
AI – for the
greater interest of
the people--and
not the other way
around.

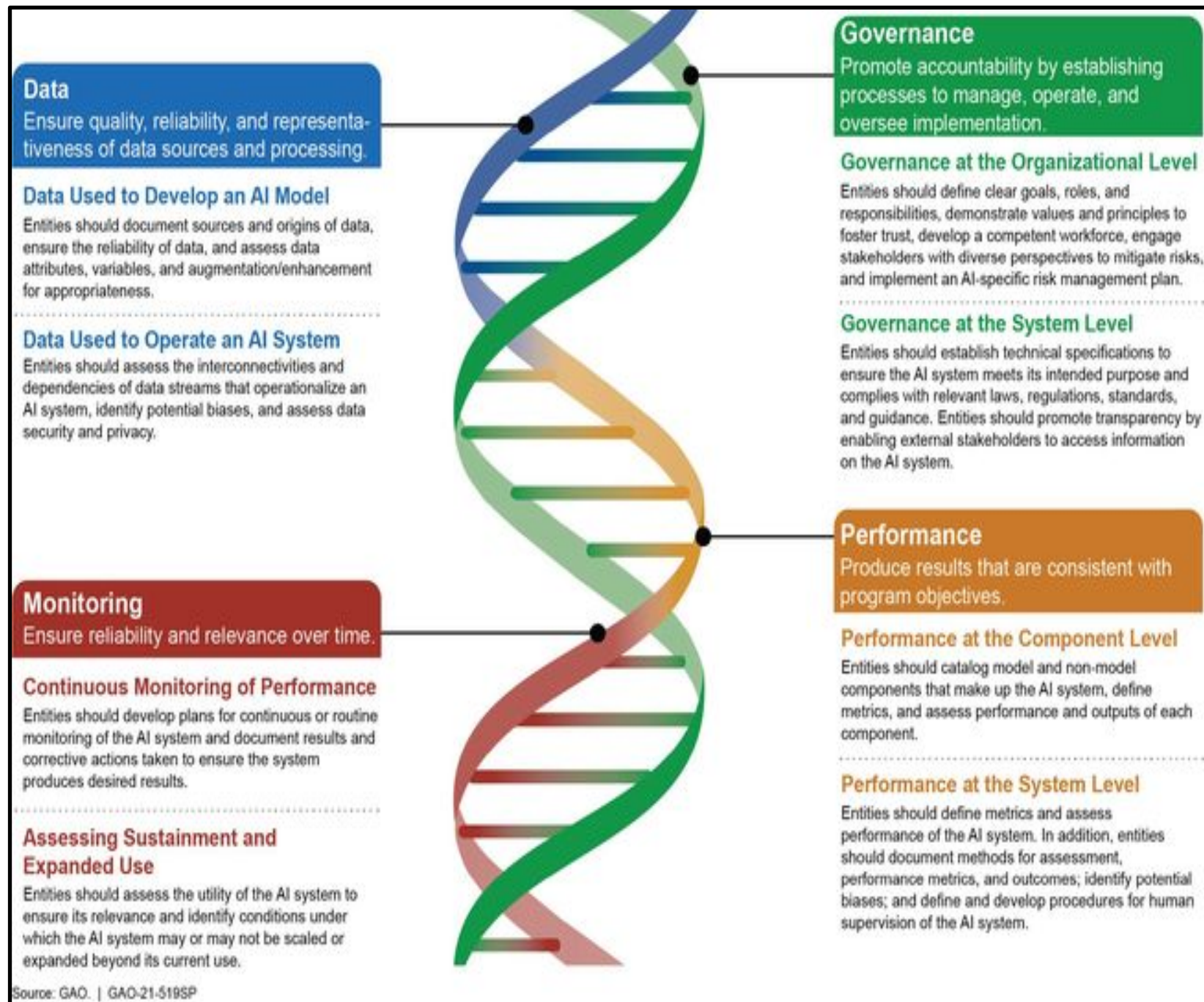
*UNESCO plan for “ethical AI”
was adopted in 2021.*

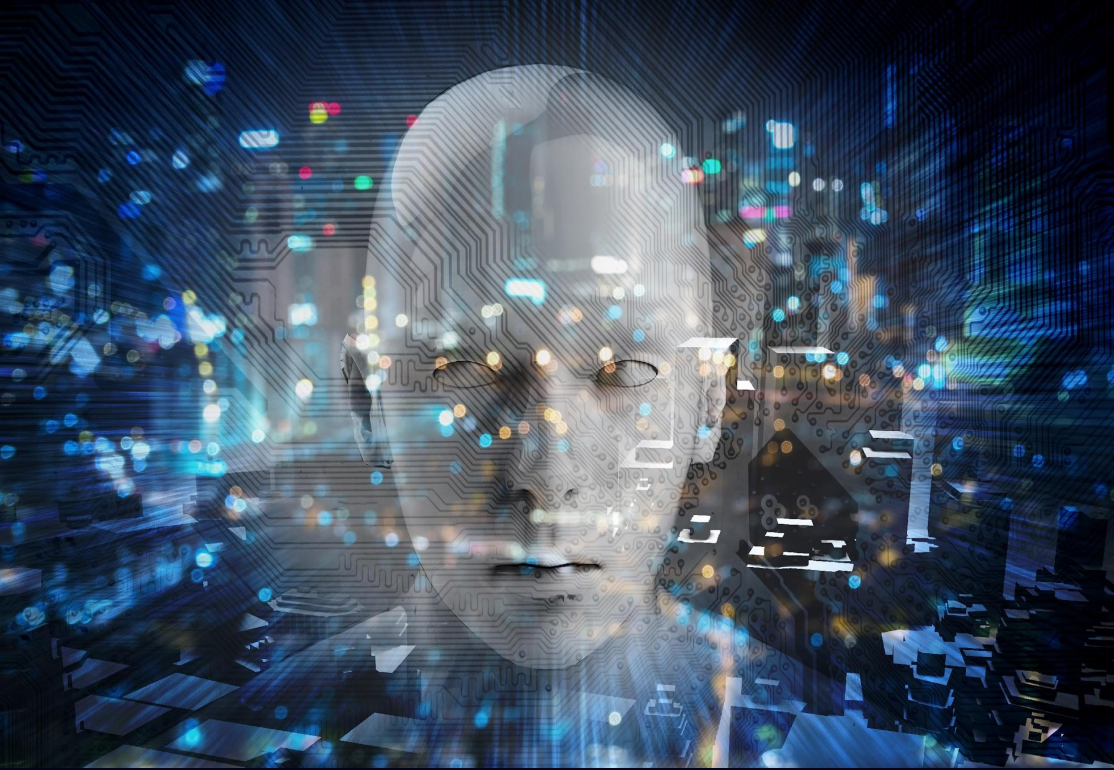


AI Accountability Framework

U.S. Government
Accountability
Office (2021)

<https://www.gao.gov/products/gao-21-519sp>





What do we mean by Audit Trail & Auditable AI?

- An audit trail is the documented flow of a transaction. It is a detailed, chronological record whereby project details are tracked and traced.
- It should include information to establish what events occurred and who (or what) caused them).
- Beyond responsible/accountable AI, is Auditable AI, an audit trail of a company's documented development governance standard during the production of an AI model.



How do we explain the transaction?

What evidence can we produce to support our position?

Documenting the AI Process

*Documentation – Records -
Evidence*

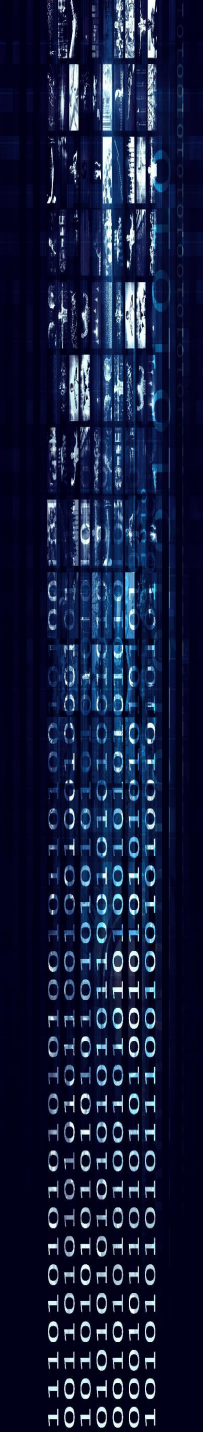




"Defining an AI record and developing methods for capturing AI records is a project the profession should take on."

~Norman Mooradian, Ph.D.





"If business is no longer to be transacted only by human beings, but also by AI agents, or some combination of the two, what will evidence of those transactions look like, what will the record be?"

~Jenny Bunn



Recordkeepers may ask:

What records are created within AI research teams to document their process?

What records are created of the decisions to procure or deploy systems utilising AI?

What records are created of the decisions and impact of such systems?

Are the created records sufficient to meet existing legal provisions?

Do the created records meet the required standards of quality?



Paradata is a source of information in the form of auxiliary data describing the process [of the use of computer-assisted survey instruments.] ~*Mick P. Couper, 2010*

*This brings us
back to
Paradata*

Paradata is a term used to describe data generated as a by-product of the data collection process. ~*U.S. Census Bureau, 2022*



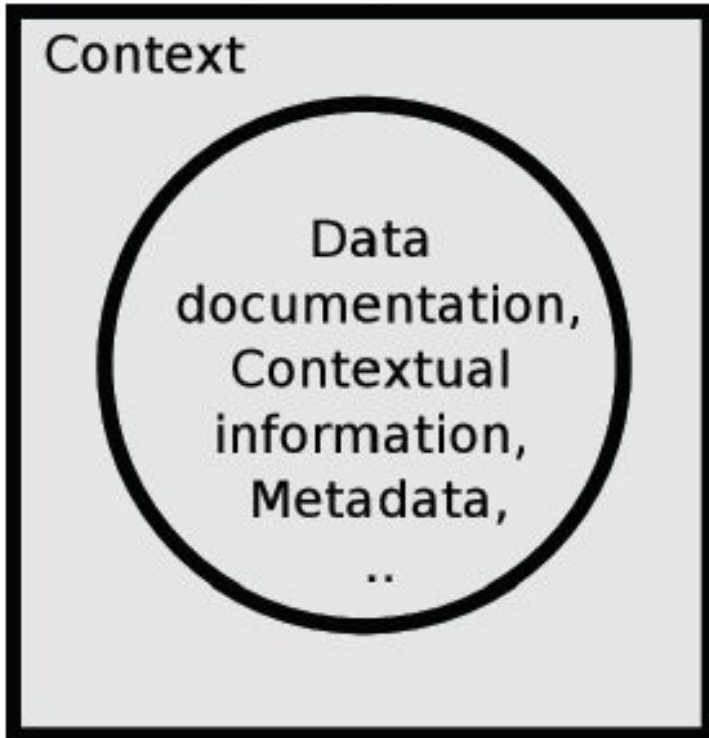
Metadata is formalized **data about statistical data** needed to search for, display, and analyze those data. ~*National Institute of Statistical Sciences, 2010*

Metadata or Paradata?

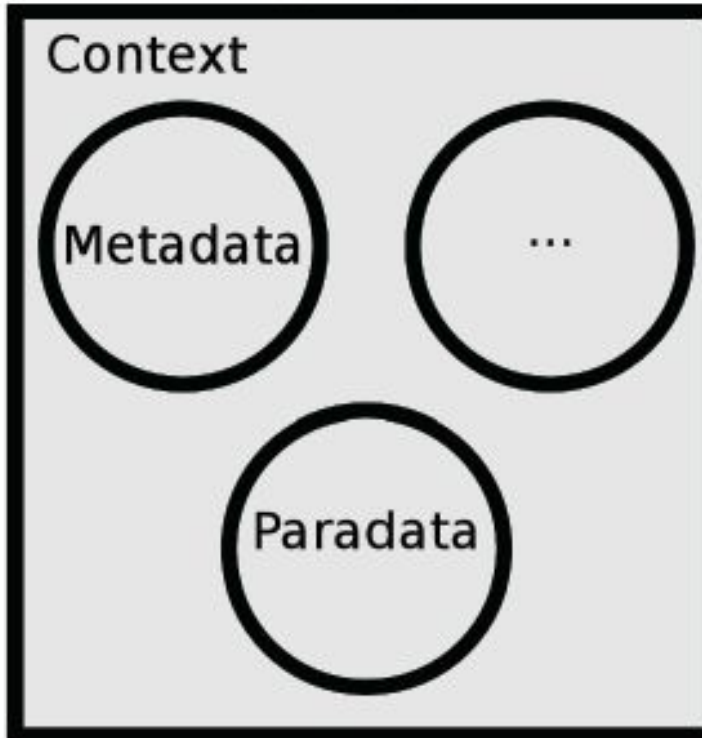
Paradata is formalized **data on methodologies, processes and quality** associated with the production and assembly of statistical data. ~*National Institute of Statistical Sciences, 2010*

Perspectives to Contextual Information

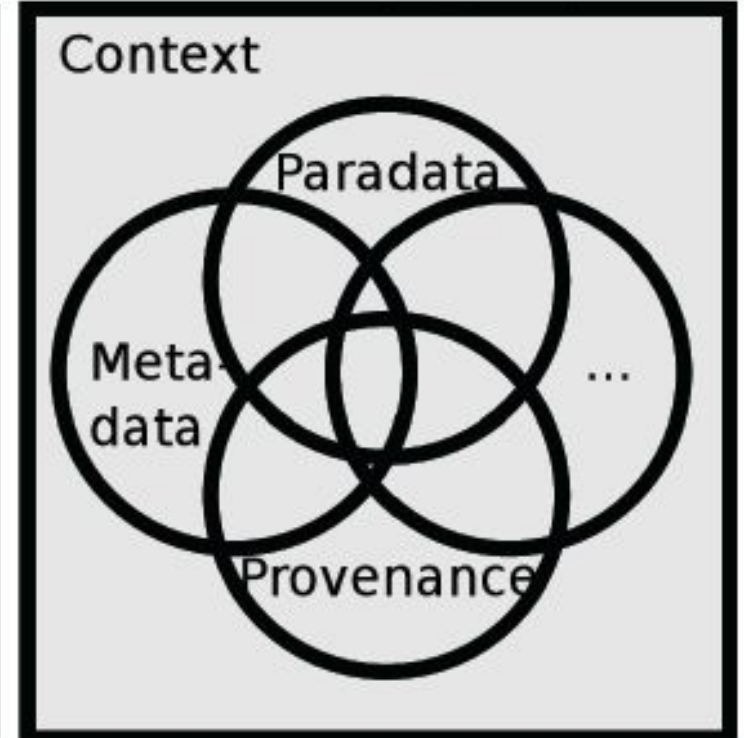
1) Broad perspective



2) Narrow perspective



3) Middle-range perspective

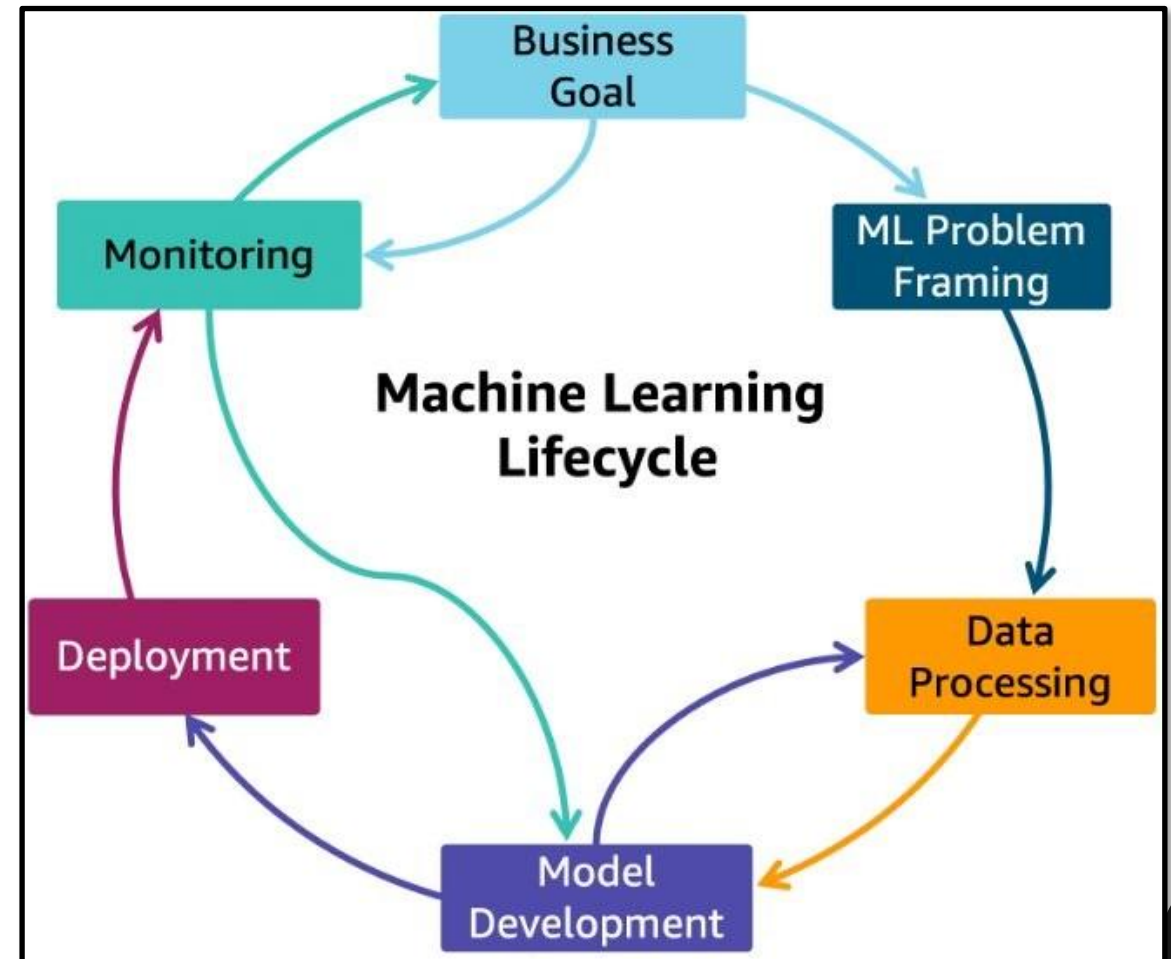


~Huvila, Isto, "Improving the usefulness of research data with better Paradata"

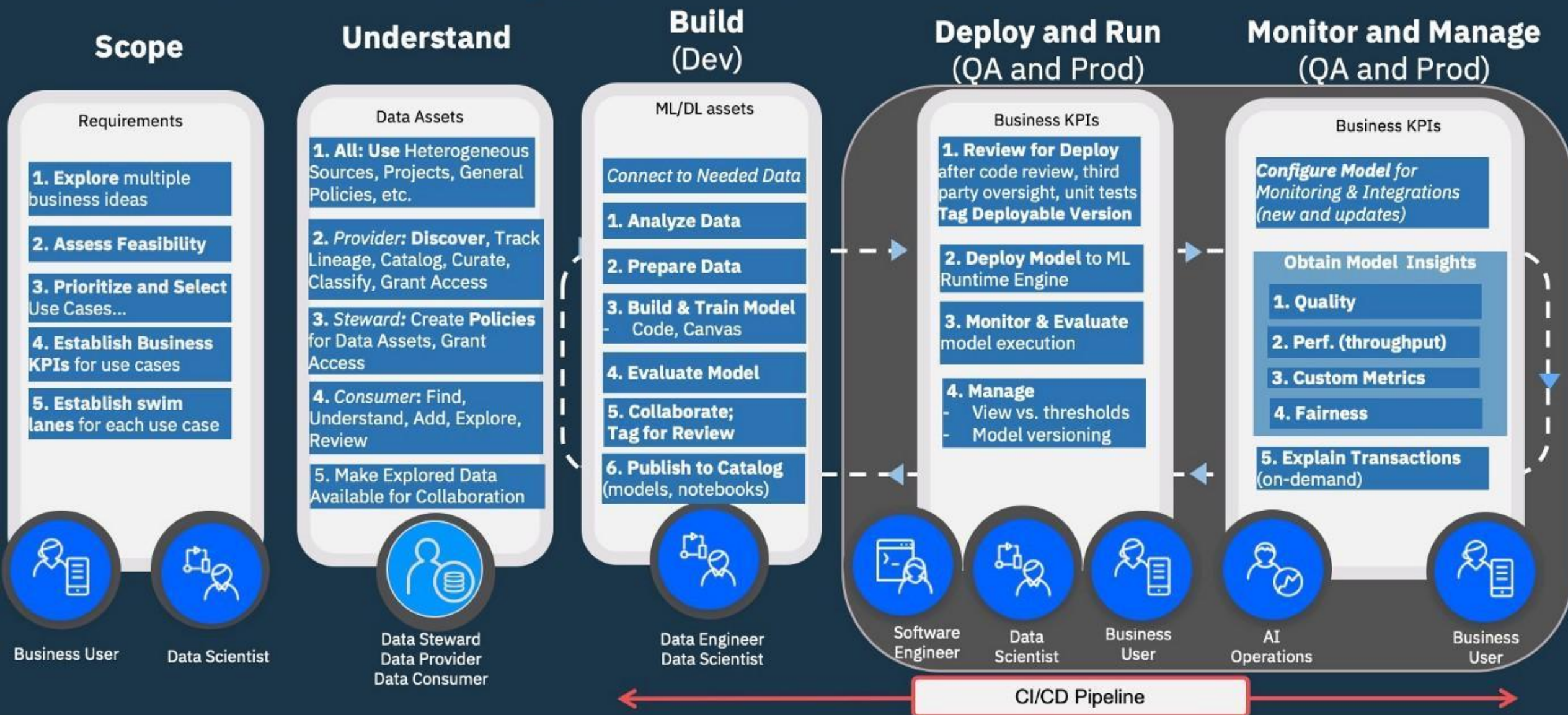
PARADATA & AI Process

Paradata is the formation about the procedure(s) and tools used to create and process information resources, along with information about the persons carrying out those procedures.

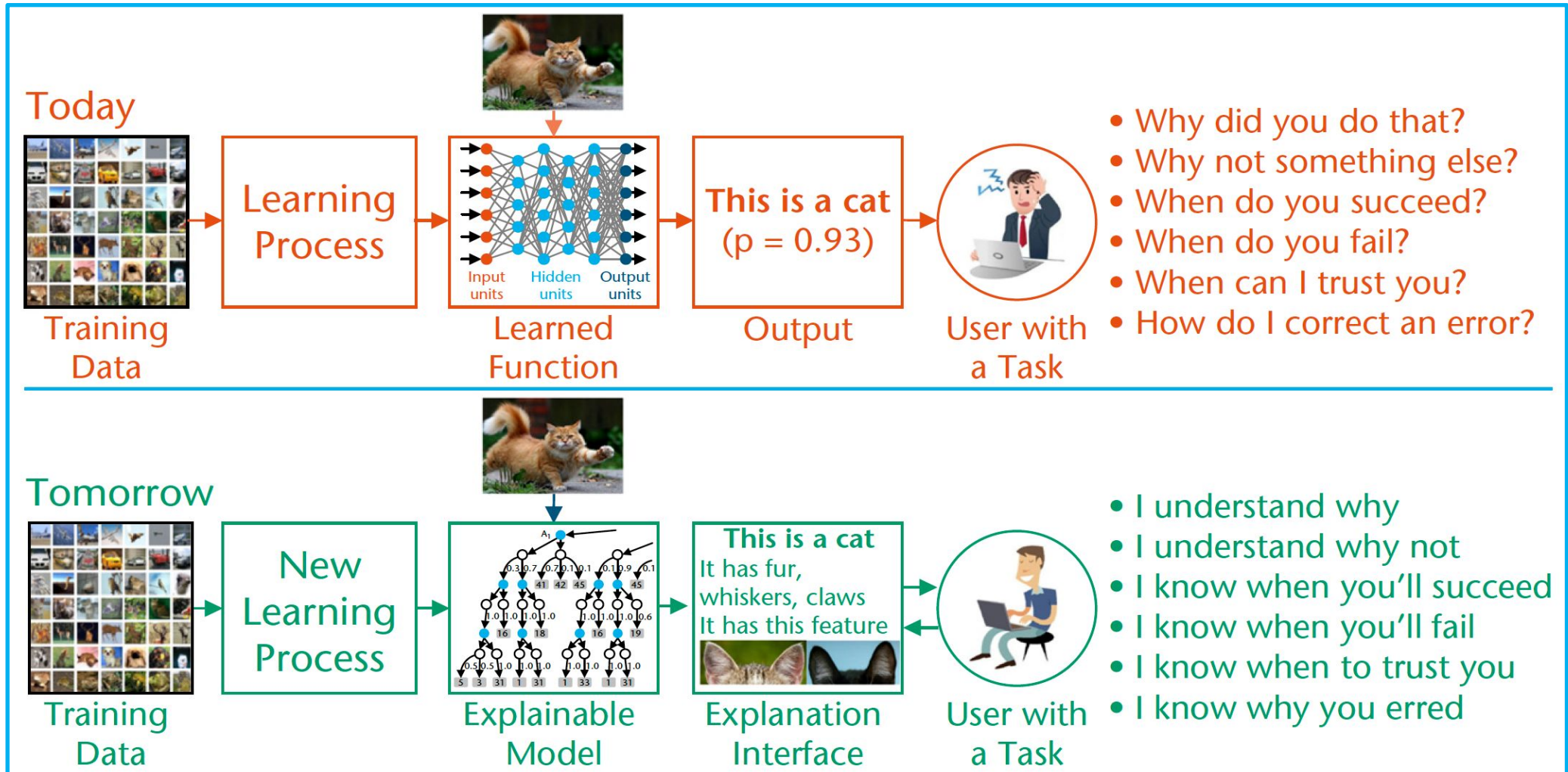
~ITrustAI working definition



Data Science & AI Lifecycle - a General View



XAI Concept



Modes of Explanation

- Causal – How it functions
- Epistemic – How we know it functions
- Justificatory – On what grounds it functions

Justificatory -- Can refer to AI system properties (e.g., datasets and algorithms).

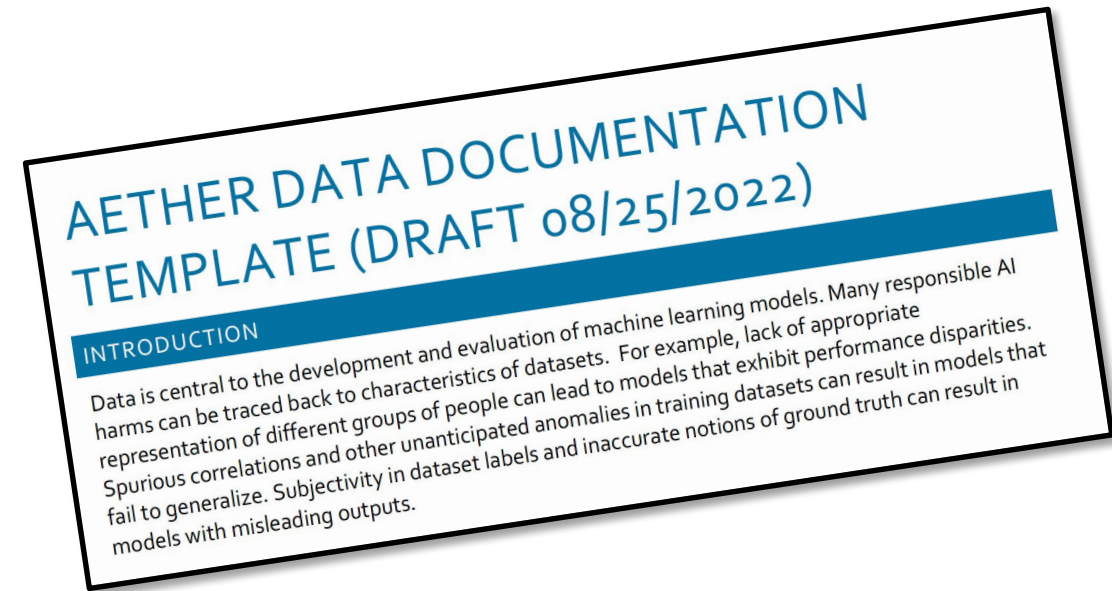
Must also reference institutional and social facts about the implementation of the system (e.g., regulations, standards, organizational processes pertinent to the use case).

ISO/IEC TR 24028 (2020-05) Information technology — Artificial intelligence
— Overview of trustworthiness in artificial intelligence

Microsoft's Datasheets for Datasets

Potential Audience:

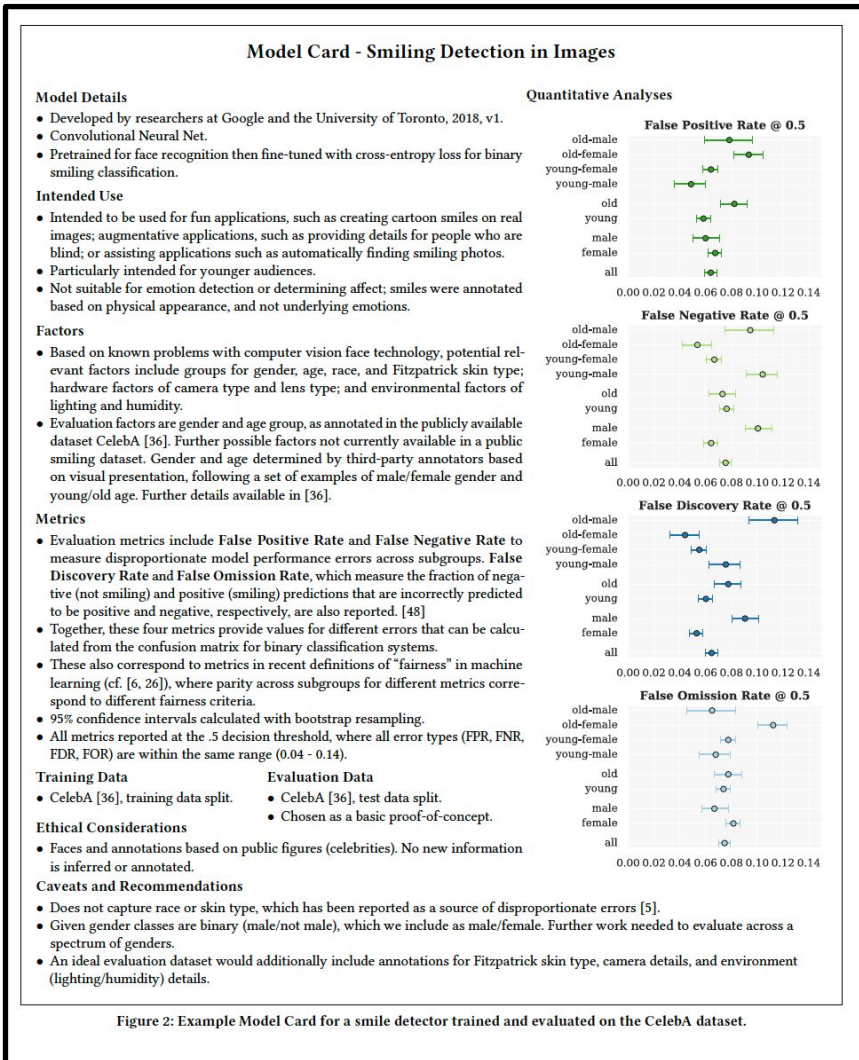
- People who are considering using this dataset to train or evaluate models
- People who are auditing a model or AI system



Major sections of the template:

- Data Set Overview (ex. contact, distribution, access basics; data set contents; intended and inappropriate uses.)
- Details (data collection procedures/ representativeness; data quality; pre-processing cleaning, and labeling; privacy; additional details on distribution and access.)

Google Model Cards



- ✓ Model Details
- ✓ Intended Use
- ✓ Factors
- ✓ Metrics
- ✓ Training Data
- ✓ Ethical Considerations
- ✓ Caveats & Recommendations

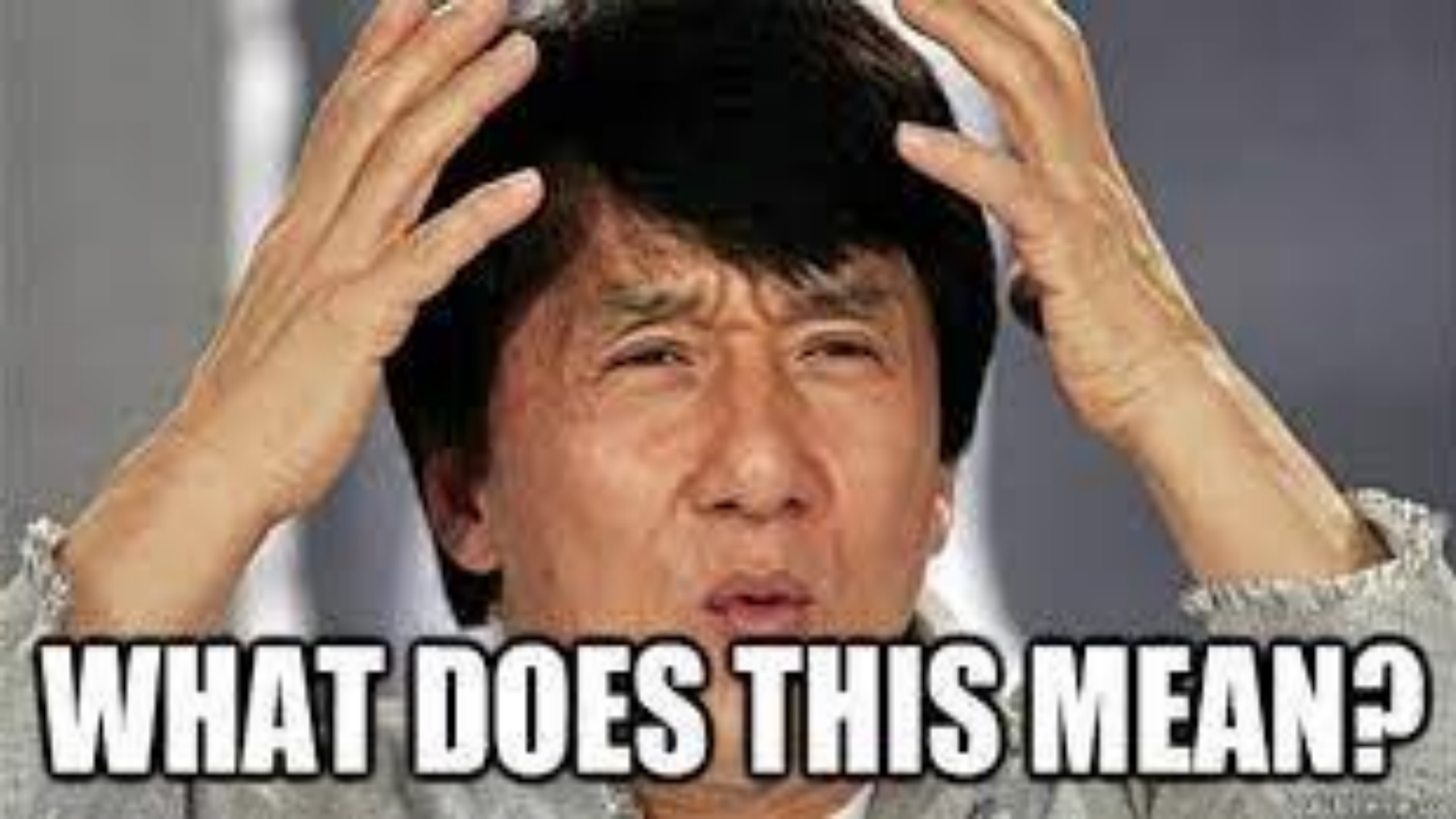
Updated model cards at <https://modelcards.withgoogle.com/about>

IBM's AI FactSheets 360

5 views – all facts and based on roles

The screenshot displays the 'Mortgage Evaluator Governance FactSheet' interface. At the top, it states 'Created to demonstrate how development and deployment facts of a mortgage evaluation model can be recorded and viewed'. Below this, there are five view options: 'All Facts View' (Every fact collected from concept to deployment), 'Business Owner's View' (Filtered to show just business relevant facts), 'Data Scientist's View' (Primarily data and model metrics), 'Model Validator's View' (Compares challenge model metrics), and 'AI Ops Engineer's View' (Compares deployment metrics). The 'Business Owner's View' is selected, indicated by a green arrow. Below the view selection, the 'Mortgage Evaluator' section shows a 'Business Request' table with columns 'Purpose' and 'Risk Level'. The 'Purpose' is 'Predict mortgage approval' and the 'Risk Level' is 'High'. The 'Model Policy' section lists seven items: 1. Datasets must be approved and in data catalog, 2. Race, ethnicity, and gender of applicant cannot be used in models used to make mortgage related decisions, 3. Model predictive performance metrics must minimally include accuracy, balanced_accuracy and AUC score, 4. Models must be checked for bias using Disparate Impact, 5. Models must be checked for faithfulness of explanations, 6. Models must be checked for robustness to Adversarial attacks using Empirical Robustness metric, 7. Models must be checked for robustness to dataset shift.

This screenshot shows the same 'Mortgage Evaluator Governance FactSheet' interface, but with the 'Business Owner's View' selected and the 'Risk Level' highlighted in a green box. The 'Business Request' table shows 'Purpose' as 'Predict mortgage approval' and 'Risk Level' as 'High'. The 'Model Policy' section lists seven items: 1. Datasets must be approved and in data catalog, 2. Race, ethnicity, and gender of applicant cannot be used in models used to make mortgage related decisions, 3. Model predictive performance metrics must minimally include accuracy, balanced_accuracy and AUC score, 4. Models must be checked for bias using Disparate Impact, 5. Models must be checked for faithfulness of explanations, 6. Models must be checked for robustness to Adversarial attacks using Empirical Robustness metric, 7. Models must be checked for robustness to dataset shift.



WHAT DOES THIS MEAN?

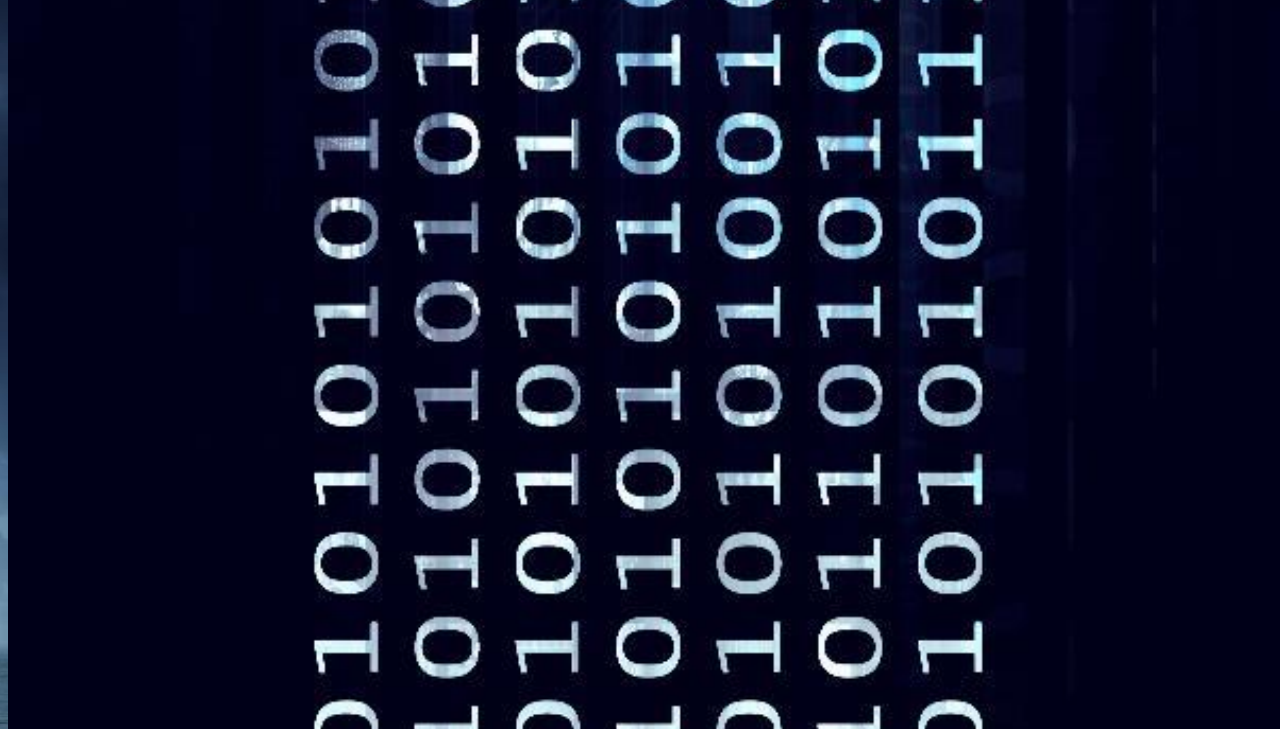
Interpretable implying some sense of understanding how the technology works;

Explainable, implying that a wider range of users can understand why or how a conclusion was reached;

Transparent implying some level of accessibility to the data or algorithm;

Justifiable implying there is an understanding of the case in support of a particular outcome

Contestable implying users have the information they need to argue against a decision or classification.



Paradata Can Help!

Continue research into...

- the nature of Paradata to document the AI process,
- the relationship between and potential overlap across Metadata and Paradata,
- the actions that take place along the AI lifecycle that require documentation,
- a recommended risk-management approach when determining the extent of documentation needed,
- the best form of representation, method of capture, and preservation,
- And finally, the identification and development of AI tools and techniques that can be employed to aid us in this task.



Thank you!

Patricia C. Franks, PhD
Professor Emerita
ItrustAI Researcher
San Jose State University
patricia.franks@sjsu.edu

